Folien zur Vorlesung am 17.06.2025
3D Computer Vision
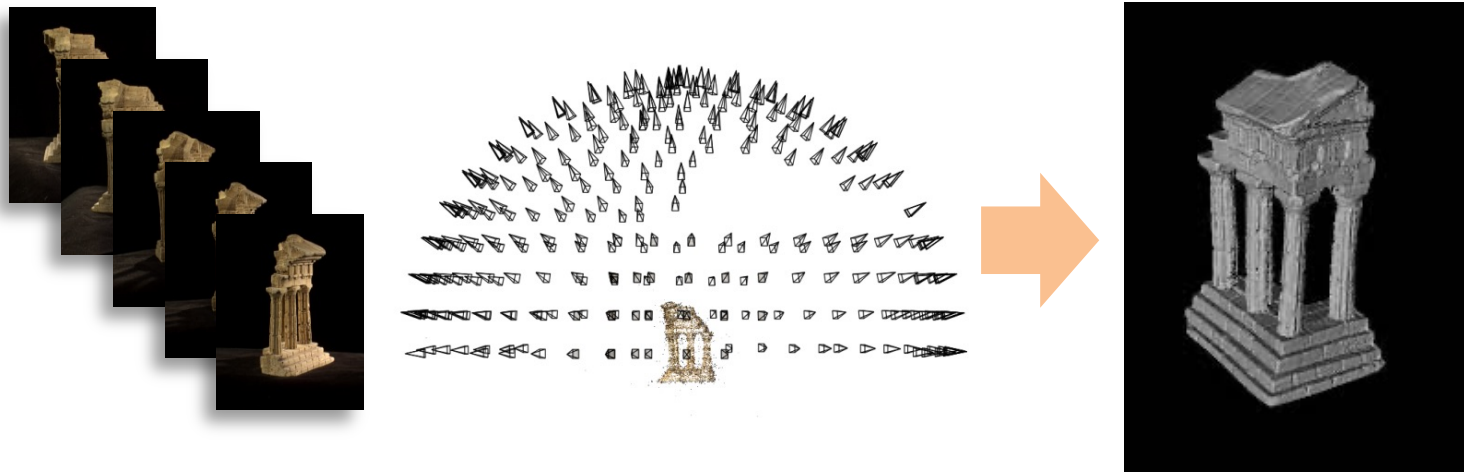
# STRUCTURE FROM MOTION

# References

- Book:
  - Computer Vision: Algorithms and Applications, Richard Szeliski → In the 2nd Edition Chapter 11, available at http://szeliski.org/book
- Slides:
  - Introduction to Computer Vision, CS5670, Spring 2022, Cornell Tech → Noah Snavely and others: Yung-Yu Chuang, Fredo Durand, Alexei Efros, William Freeman, James Hays, Svetlana Lazebnik, Andrej Karpathy, Fei-Fei Li, Srinivasa Narasimhan, Silvio Savarese, Steve Seitz, Richard Szeliski, and Li Zhang.
    - Most slides are copied from there.
- Videos:
  - Shree Nayar: https://fpcv.cs.columbia.edu/ → Today: a summary of the Reconstruction II → Structure from Motion series

# Reminder: multi-view stereo

**Problem formulation:** given several images of the same object or scene, compute a representation of its 3D shape

# Structure from motion

- Multi-view stereo assumes that cameras are calibrated
    - Extrinsics and intrinsics are known for all views

- How do we compute calibration if we don't know it? In general, this is called *structure from motion*

# Large-scale structure from motion



Dubrovnik, Croatia. 4,619 images (out of an initial 57,845).
Total reconstruction time: 23 hours
Number of cores: 352

From „Building Rome in a day"

# Two views



- Solve for Fundamental matrix / Essential matrix
- Factorize into intrinsics, rotation, and translation

# What about more than two views?

- The geometry of three views is described by a 3 x 3 x 3 tensor called the *trifocal tensor*

- The geometry of four views is described by a 3 x 3 x 3 x 3 tensor called the *quadrifocal tensor*

- After this it starts to get complicated...
  - Instead, we explicitly solve for camera poses and scene geometry
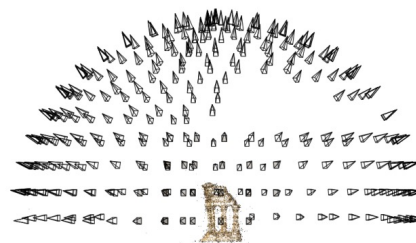
# Structure from motion

- Given many images, how can we
    - a) figure out where they were all taken from?
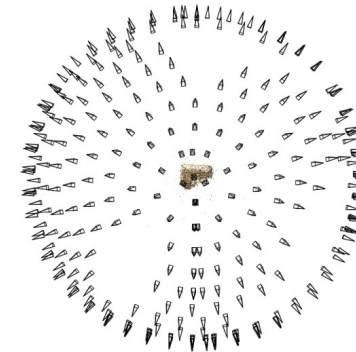    - b) build a 3D model of the scene?



This is the **structure from motion (SfM)** problem
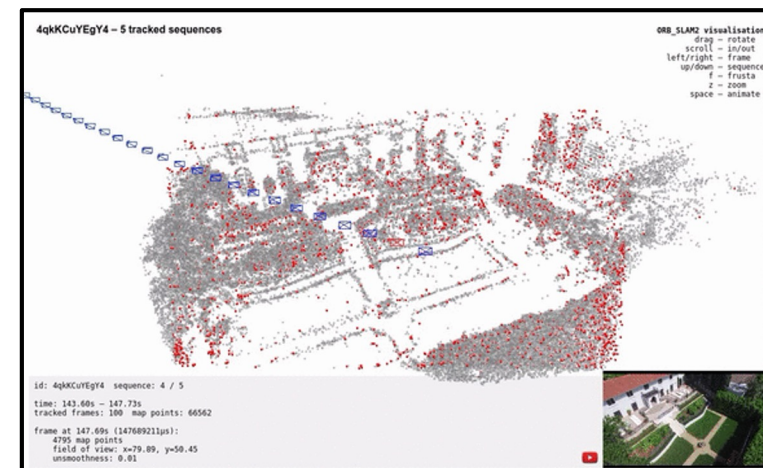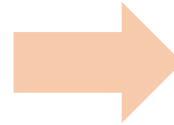
# Structure from motion



Reconstruction (side)

(top)

- Input: images with 2D points in correspondence
$$p_{ij} = (u_{ij}, v_{ij})$$

- Output
  - structure: 3D location $\mathbf{x}_i$ for each point $p_i$
  - motion: camera parameters $\mathbf{R}_j$, $\mathbf{t}_j$ & possibly $\mathbf{K}_j$

- Objective function: minimize *reprojection error*
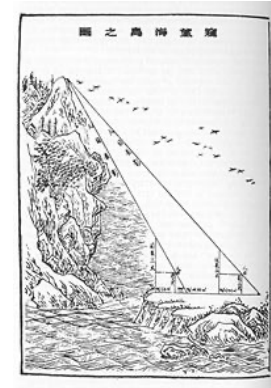
# Also doable from video

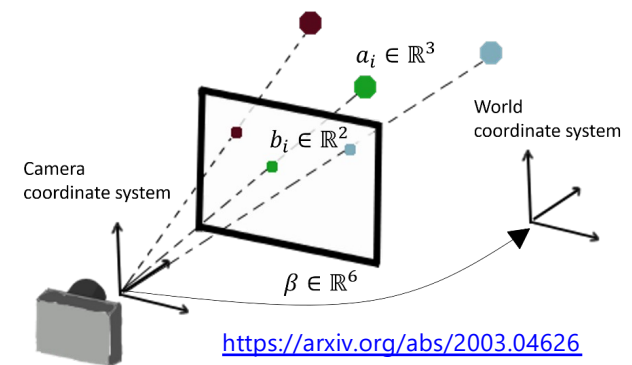# What we've seen so far...

- 2D transformations between images
  - Translations, affine transformations, homographies...

- Fundamental matrices
  - Represent relationships between 2D images in the form of corresponding 2D lines

- **What's new:** Explicitly representing 3D geometry of cameras *and points*

# Triangulation & camera calibration

- Suppose we have known camera parameters, each of which observes a point
  - How can we compute the 3D location of that point?
  - This is called *triangulation* (known since ancient times)

- On the other hand: Suppose we have known 3D points
  - And have matches between these points and an image
  - How can we compute the camera parameters?
  - This is called *camera calibration* (or *camera resectioning*)

Liu Hui (c. 263), How to measure the height of a sea island.

$a_i \in \mathbb{R}^3$

World coordinate system

$b_i \in \mathbb{R}^2$

Camera coordinate system

$\beta \in \mathbb{R}^6$
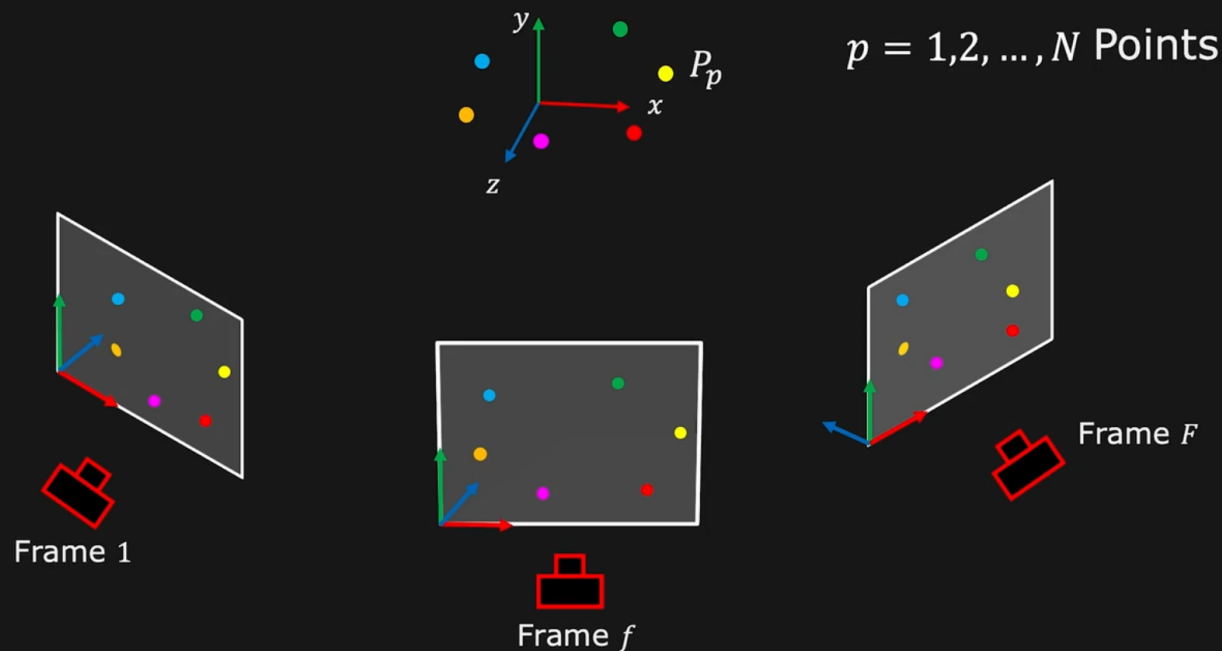
https://arxiv.org/abs/2003.04626

# Structure from motion

- What if we don't know 3D points or camera parameters?
- SfM solves both of these problems *at once*
- A kind of chicken-and-egg problem
  - (but solvable)

Next: Classic approach using an assumption that simplifies the problem.

Afterwards: Extended approaches...

# Orthographic Structure from Motion



Given sets of corresponding image points (2D): $(u_{f,p}, v_{f,p})$

Find scene points (3D) $P_p$, assuming orthographic camera.

[Tomasi 1992]

Image from Shree Nayar FPCV videos

# Orthographic Structure from Motion



$p = 1,2,\ldots,N$ Points

Given sets of corresponding image points (2D): $(u_{f,p}, v_{f,p})$

Find scene points (3D) $P_p$, assuming orthographic camera.

[Tomasi 1992]

Image from Shree Nayar FPCV videos
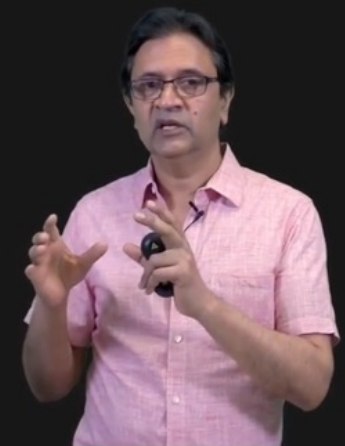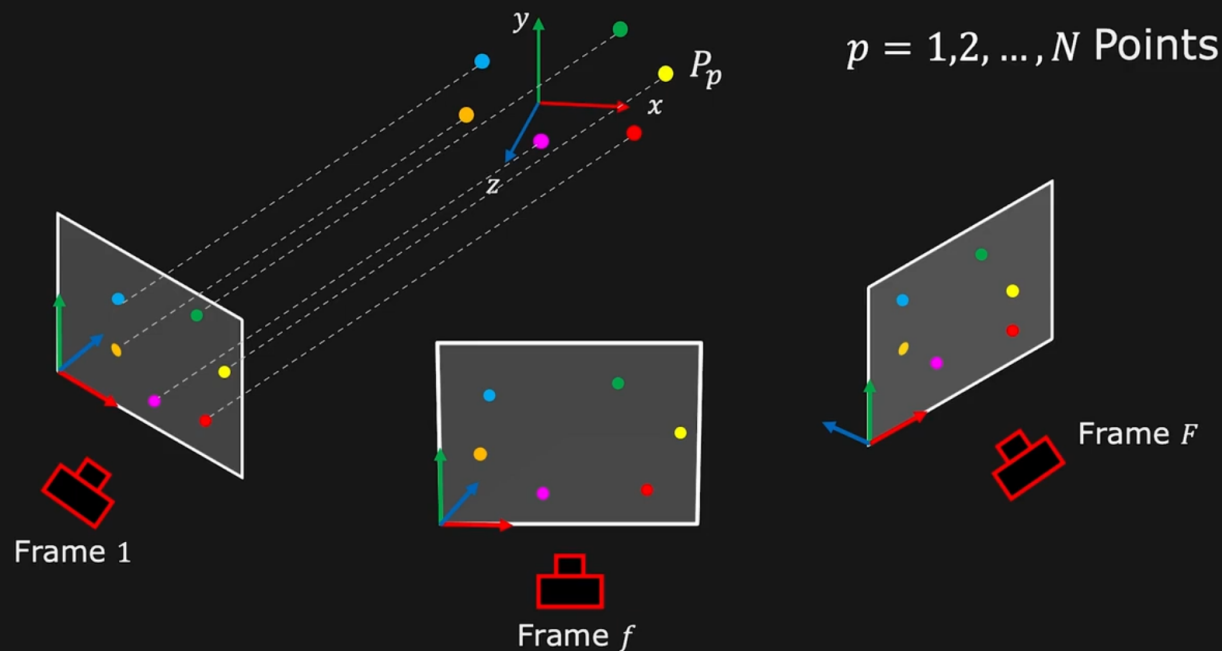
# From 3D to 2D: Orthographic Projection
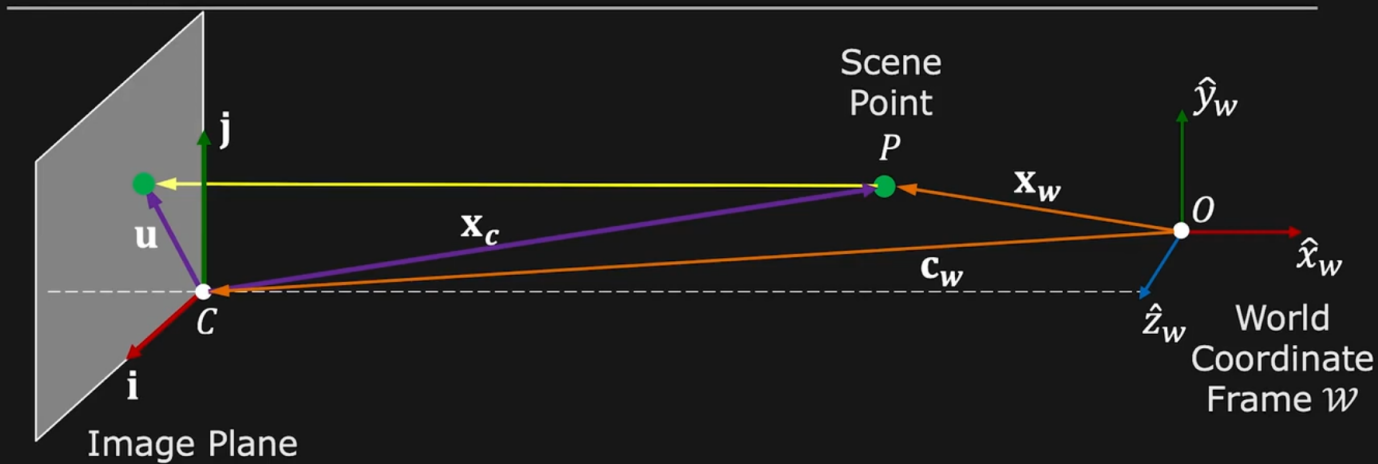
Scene Point $P$

$\hat{y}_w$

$\mathbf{x}_w$

$O$

$\hat{x}_w$

$\mathbf{x}_c$

$\mathbf{c}_w$

$\hat{z}_w$ World Coordinate Frame $\mathcal{W}$
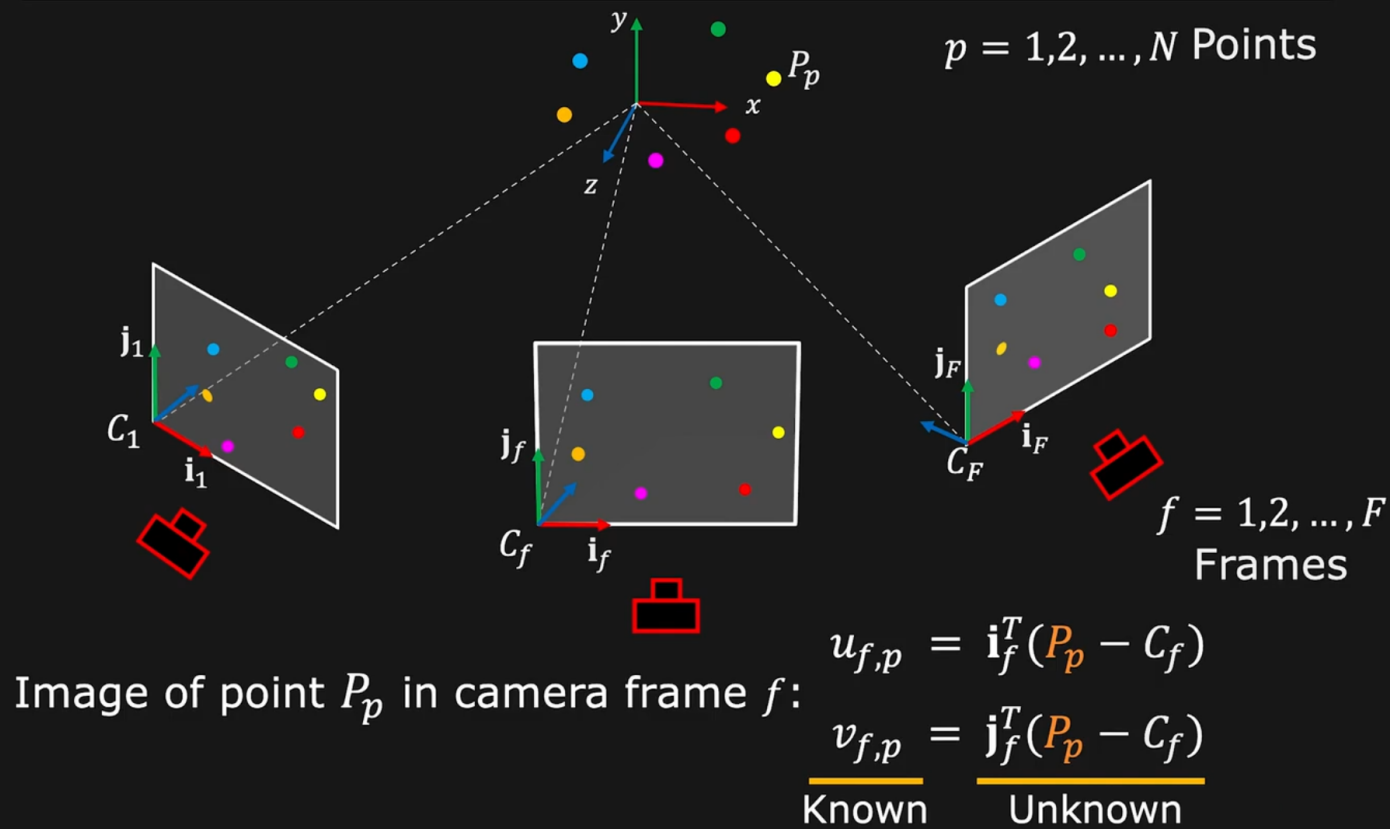
$\mathbf{j}$

$\mathbf{u}$

$C$

$\mathbf{i}$

Image Plane

$$u = \mathbf{i}^T \mathbf{x}_c = \mathbf{i}^T(\mathbf{x}_w - \mathbf{c}_w) = \mathbf{i}^T(P - C)$$

$$v = \mathbf{j}^T \mathbf{x}_c = \mathbf{j}^T(\mathbf{x}_w - \mathbf{c}_w) = \mathbf{j}^T(P - C)$$

$$u = \mathbf{i}^T(P - C)$$
$$v = \mathbf{j}^T(P - C)$$

Image from Shree Nayar FPCV videos

# Orthographic SFM



$p = 1, 2, \ldots, N$ Points

$f = 1, 2, \ldots, F$ Frames

Image of point $P_p$ in camera frame $f$:

$$u_{f,p} = \mathbf{i}_f^T (P_p - C_f)$$

$$v_{f,p} = \mathbf{j}_f^T (P_p - C_f)$$

Known ___ Unknown ___

We can remove $C_f$ from equations to simply SFM problem.
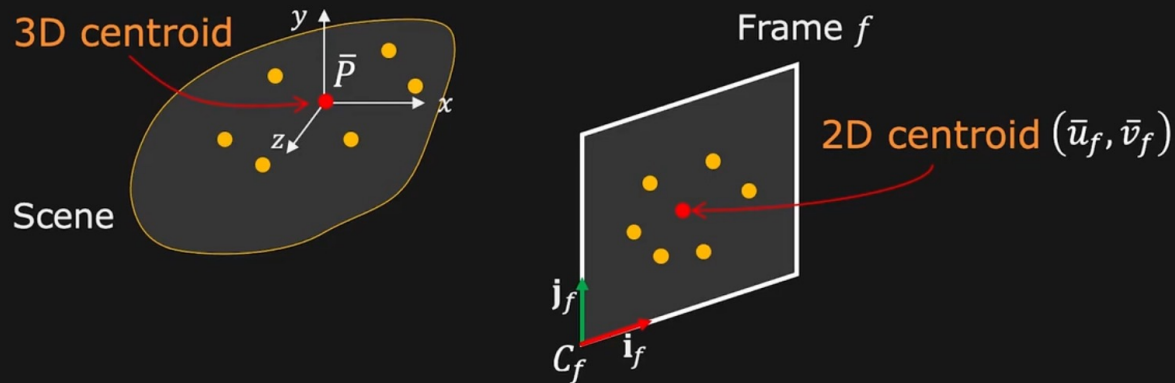
Image from Shree Nayar FPCV videos

# Centering Trick



Assume origin of world at centroid of scene points:

$$\frac{1}{N}\sum_{p=1}^{N} P_p = \bar{P} = \mathbf{0}$$

We will compute scene points w.r.t their centroid!

Image from Shree Nayar FPCV videos

# Centering Trick



3D centroid — $\bar{P}$ — Scene

Frame $f$ — 2D centroid $(\bar{u}_f, \bar{v}_f)$

$\mathbf{j}_f$  $C_f$  $\mathbf{i}_f$

Centroid $(\bar{u}_f, \bar{v}_f)$ of the image points in frame $f$:

$$\bar{u}_f = \frac{1}{N}\sum_{p=1}^{N} u_{f,p} = \frac{1}{N}\sum_{p=1}^{N} \mathbf{i}_f^T(P_p - C_f)$$
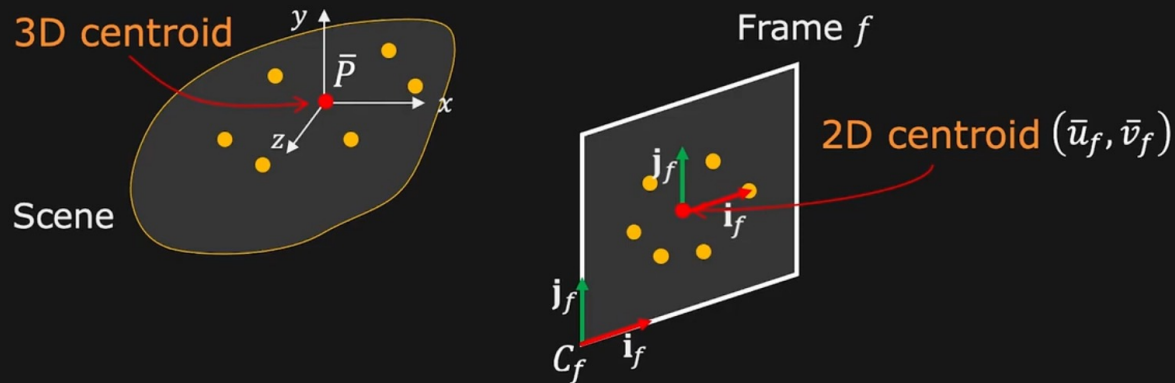
$$\bar{v}_f = \frac{1}{N}\sum_{p=1}^{N} v_{f,p} = \frac{1}{N}\sum_{p=1}^{N} \mathbf{j}_f^T(P_p - C_f)$$

$$\bar{u}_f = \frac{1}{N}\mathbf{i}_f^T\sum_{p=1}^{N} P_p - \frac{1}{N}\sum_{p=1}^{N} \mathbf{i}_f^T C_f$$

$$\bar{v}_f = \frac{1}{N}\mathbf{j}_f^T\sum_{p=1}^{N} P_p - \frac{1}{N}\sum_{p=1}^{N} \mathbf{j}_f^T C_f$$

$$\boxed{\bar{u}_f = -\mathbf{i}_f^T C_f}$$

$$\boxed{\bar{v}_f = -\mathbf{j}_f^T C_f}$$

Image from Shree Nayar FPCV videos

# Centering Trick



3D centroid $\bar{P}$

Scene

Frame $f$

2D centroid $(\bar{u}_f, \bar{v}_f)$

Shift camera origin to the centroid $(\bar{u}_f, \bar{v}_f)$.

Image points w.r.t. $(\bar{u}_f, \bar{v}_f)$:

$$\tilde{u}_{f,p} = u_{f,p} - \bar{u}_f \qquad\qquad \tilde{v}_{f,p} = v_{f,p} - \bar{v}_f$$

$$= \mathbf{i}_f^T(P_p - C_f) - \mathbf{i}_f^T C_f \qquad\qquad = \mathbf{j}_f^T(P_p - C_f) - \mathbf{j}_f^T C_f$$

$$\boxed{\tilde{u}_{f,p} = \mathbf{i}_f^T P_p} \qquad\qquad \boxed{\tilde{v}_{f,p} = \mathbf{j}_f^T P_p}$$

Camera locations $C_f$ now removed from equations.

Image from Shree Nayar FPCV videos

# Observation Matrix $W$

$$\tilde{u}_{f,p} = \mathbf{i}_f^T P_p$$

$$\tilde{v}_{f,p} = \mathbf{j}_f^T P_p$$

$\Rightarrow$

$$\begin{bmatrix} \tilde{u}_{f,p} \\ \tilde{v}_{f,p} \end{bmatrix} = \begin{bmatrix} \mathbf{i}_f^T \\ \mathbf{j}_f^T \end{bmatrix} P_p$$

Image from Shree Nayar FPCV videos

# Observation Matrix $W$

$$\tilde{u}_{f,p} = \mathbf{i}_f^T P_p$$
$$\tilde{v}_{f,p} = \mathbf{j}_f^T P_p$$

$\Rightarrow$

$$\begin{bmatrix} \tilde{u}_{f,p} \\ \tilde{v}_{f,p} \end{bmatrix} = \begin{bmatrix} \mathbf{i}_f^T \\ \mathbf{j}_f^T \end{bmatrix} P_p$$

$$
\begin{array}{c}
\begin{array}{cccc} \text{Point 1} & \text{Point 2} & & \text{Point N} \end{array} \\
\begin{array}{c} \text{Image 1} \\ \text{Image 2} \\ \\ \text{Image F} \\ \text{Image 1} \\ \text{Image 2} \\ \\ \text{Image F} \end{array}
\begin{bmatrix}
\tilde{u}_{1,1} & \tilde{u}_{1,2} & \cdots & \tilde{u}_{1,N} \\
\tilde{u}_{2,1} & \tilde{u}_{2,2} & \cdots & \tilde{u}_{2,N} \\
\vdots & \vdots & \vdots & \vdots \\
\tilde{u}_{F,1} & \tilde{u}_{F,2} & \cdots & \tilde{u}_{F,N} \\
\tilde{v}_{1,1} & \tilde{v}_{1,2} & \cdots & \tilde{v}_{1,N} \\
\tilde{u}_{2,1} & \tilde{u}_{2,2} & \cdots & \tilde{v}_{2,N} \\
\vdots & \vdots & \vdots & \vdots \\
\tilde{v}_{F,1} & \tilde{v}_{F,2} & \cdots & \tilde{v}_{F,N}
\end{bmatrix}
\end{array}
=
\begin{bmatrix}
\mathbf{i}_1^T \\ \mathbf{i}_2^T \\ \vdots \\ \mathbf{i}_F^T \\ \mathbf{j}_1^T \\ \mathbf{j}_2^T \\ \vdots \\ \mathbf{j}_F^T
\end{bmatrix}
\begin{bmatrix} P_1 & P_2 & \cdots & P_N \end{bmatrix}
$$

$W_{2F \times N}$

**Centroid-Subtracted Feature Points** (Known)

$M_{2F \times 3}$

**Camera Motion** (Unknown)

Point 1  Point 2  Point N

$S_{3 \times N}$

**Scene Structure** (Unknown)

Image from [Shree Nayar FPCV videos](#)

# Observation Matrix $W$

$$
\underbrace{
\begin{array}{c}
\text{Image 1} \\ \text{Image 2} \\ \\ \text{Image F} \\ \text{Image 1} \\ \text{Image 2} \\ \\ \text{Image F}
\end{array}
\overset{\text{Point 1 \quad Point 2 \qquad\qquad Point N}}{
\begin{bmatrix}
\tilde{u}_{1,1} & \tilde{u}_{1,2} & \dots & \tilde{u}_{1,N} \\
\tilde{u}_{2,1} & \tilde{u}_{2,2} & \dots & \tilde{u}_{2,N} \\
\vdots & \vdots & \vdots & \vdots \\
\tilde{u}_{F,1} & \tilde{u}_{F,2} & \dots & \tilde{u}_{F,N} \\
\tilde{v}_{1,1} & \tilde{v}_{1,2} & \dots & \tilde{v}_{1,N} \\
\tilde{u}_{2,1} & \tilde{u}_{2,2} & \dots & \tilde{v}_{2,N} \\
\vdots & \vdots & \vdots & \vdots \\
\tilde{v}_{F,1} & \tilde{v}_{F,2} & \dots & \tilde{v}_{F,N}
\end{bmatrix}}
}_{\substack{W_{2F\times N} \\ \text{Centroid-Subtracted} \\ \text{Feature Points (Known)}}}
=
\underbrace{
\begin{bmatrix}
\mathbf{i}_1^T \\ \mathbf{i}_2^T \\ \vdots \\ \mathbf{i}_F^T \\ \mathbf{j}_1^T \\ \mathbf{j}_2^T \\ \vdots \\ \mathbf{j}_F^T
\end{bmatrix}}_{\substack{M_{2F\times 3} \\ \text{Camera Motion} \\ \text{(Unknown)}}}
\underbrace{
\overset{\text{Point 1 \quad Point 2 \qquad Point N}}{
\begin{bmatrix} P_1 & P_2 & \dots & P_N \end{bmatrix}}
}_{\substack{S_{3\times N} \\ \text{Scene Structure} \\ \text{(Unknown)}}}
$$

Can we find $M$ and $S$ from $W$?

Image from Shree Nayar FPCV videos

# Important Properties of Matrix Rank

- $Rank(A^T) = Rank(A)$

- $Rank\left(A_{m\times n}B_{n\times p}\right) = \min\left(Rank(A_{m\times n}), Rank\left(B_{n\times p}\right)\right)$
$$\leq \min(m, n, p)$$

- $Rank(AA^T) = Rank(A^TA) = Rank(A^T) = Rank(A)$

- $A_{m\times m}$ is invertible iff $Rank(A_{m\times m}) = m$

MATH PRIMER

Image from Shree Nayar FPCV videos

# Rank of Observation Matrix

$$W = M \times S$$

$$\underset{2F \times N}{} \quad \underset{2F \times 3}{} \quad \underset{3 \times N}{}$$

We know:

$$Rank(MS) \leq Rank(M) \qquad Rank(MS) \leq Rank(S)$$

⟹ $$Rank(MS) \leq \min(3, 2F) \qquad Rank(MS) \leq \min(3, N)$$

⟹ $$Rank(W) = Rank(MS) \leq \min(3, N, 2F)$$

Rank Theorem: $Rank(W) \leq 3$

[Tomasi 1992]

Image from Shree Nayar FPCV videos

# Singular Value Decomposition (SVD)

For any matrix $A$ there exists a factorization:

$$A_{M \times N} = U_{M \times M} \cdot \Sigma_{M \times N} \cdot V^T_{N \times N}$$

where $U$ and $V^T$ are orthonormal and $\Sigma$ is diagonal.

MATLAB: `[U,S,V] = svd(A)`

$$\Sigma_{M \times N} = \begin{pmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & 0 & \dots & 0 \\ 0 & 0 & 0 & \sigma_4 & \dots & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 0 & \dots & \sigma_N \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{pmatrix}$$

$\sigma_1, \dots, \sigma_N$: Singular Values

If $Rank(A) = r$ then $A$ has $r$ non-zero singular values.



$$M = U \cdot \Sigma \cdot V^*$$

© 2021 MATH PRIMER

Image from Shree Nayar FPCV videos

# Enforcing Rank Constraint

Using SVD: $\qquad W = U \, \Sigma \, V^T$

$$= \begin{bmatrix} & & \\ & U & \\ & & \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} \begin{bmatrix} & & \\ & V^T & \\ & & \end{bmatrix}$$

$\qquad\quad 2F \times 2F \qquad\qquad\qquad 2F \times N \qquad\qquad\qquad N \times N$

Since $Rank(W) \le 3$, $Rank(\Sigma) \le 3$.

All except first 3 diagonal elements of $\Sigma$ must be 0.

Image from Shree Nayar FPCV videos

# Enforcing Rank Constraint

Using SVD: $\qquad W = U \Sigma V^T$

$$= \begin{bmatrix} U_1 & U_2 \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{bmatrix} \begin{bmatrix} V_1^T \\ V_2^T \end{bmatrix}$$

$\qquad\quad 3 \qquad 2F-3 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad 3, \; N-3$

$\qquad\quad 2F \times 2F \qquad\qquad\qquad\qquad 2F \times N \qquad\qquad\qquad N \times N$

Since $Rank(W) \leq 3$, $Rank(\Sigma) \leq 3$.

Submatrices $U_2$ and $V_2^T$ do not contribute to $W$.

Image from [Shree Nayar FPCV videos](Shree Nayar FPCV videos)

# Enforcing Rank Constraint

Using SVD: $\qquad W = U \Sigma V^T$

$$= \begin{pmatrix} U_1 & U_2 \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2 & 0 & 0 & \dots & 0 \\ 0 & 0 & \sigma_3 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \dots & \vdots \end{pmatrix} \begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix}$$

$\qquad\qquad$ 3 $\qquad$ $2F-3$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ 3

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $N-3$

$\qquad\qquad$ $2F \times 2F$ $\qquad\qquad\qquad\qquad\qquad$ $2F \times N$ $\qquad\qquad\qquad\qquad$ $N \times N$

$$\boxed{W = U_1 \; \Sigma_1 \; V_1^T}$$

$$(2F \times 3)\,(3 \times 3)\,(3 \times P)$$

Image from Shree Nayar FPCV videos

# Factorization (Finding $M, S$)

$$W \;=\; \boxed{U_1 \, (\Sigma_1)^{1/2}} \; \boxed{(\Sigma_1)^{1/2} V_1^T}$$

$\qquad\qquad\quad (2F \times 3) \qquad\qquad (3 \times N)$

$\qquad\qquad\quad = M? \qquad\qquad = S?$

Image from Shree Nayar FPCV videos

# Factorization (Finding $M, S$)

$$W = \underbrace{U_1\,(\Sigma_1)^{1/2}}_{(2F\times 3)}\ \underbrace{(\Sigma_1)^{1/2}V_1^T}_{(3\times N)}$$

$$= M? \qquad = S?$$

Not so fast. Decomposition not unique!

For any 3x3 non-singular matrix $Q$:

$$W = \underbrace{U_1\,(\Sigma_1)^{1/2}Q}_{(2F\times 3)}\ \underbrace{Q^{-1}(\Sigma_1)^{1/2}V_1^T}_{(3\times N)} \quad \text{is also valid.}$$

$$= M \qquad\qquad = S \ \text{... for some } Q.$$

How to find the matrix $Q$ ?

Image from Shree Nayar FPCV videos

# Orthonormality of $M$

The Motion Matrix $M$:

$$M = \begin{bmatrix} \mathbf{i}_1^T \\ \vdots \\ \mathbf{i}_F^T \\ \mathbf{j}_1^T \\ \vdots \\ \mathbf{j}_F^T \end{bmatrix} = U_1 (\Sigma_1)^{1/2} Q = \begin{bmatrix} \hat{\mathbf{i}}_1^T \\ \vdots \\ \hat{\mathbf{i}}_F^T \\ \hat{\mathbf{j}}_1^T \\ \vdots \\ \hat{\mathbf{j}}_F^T \end{bmatrix} Q = \begin{bmatrix} \hat{\mathbf{i}}_1^T Q \\ \vdots \\ \hat{\mathbf{i}}_F^T Q \\ \hat{\mathbf{j}}_1^T Q \\ \vdots \\ \hat{\mathbf{j}}_F^T Q \end{bmatrix}$$

Computed

Computed

Orthonormality Constraints:

$$\mathbf{i}_f \cdot \mathbf{i}_f = \mathbf{i}_f^T \mathbf{i}_f = 1$$

$$\mathbf{j}_f \cdot \mathbf{j}_f = \mathbf{j}_f^T \mathbf{j}_f = 1$$

$$\mathbf{i}_f \cdot \mathbf{j}_f = \mathbf{i}_f^T \mathbf{j}_f = 0$$

$\Rightarrow$

$$\hat{\mathbf{i}}_f^T Q Q^T \hat{\mathbf{i}}_f = 1$$

$$\hat{\mathbf{j}}_f^T Q Q^T \hat{\mathbf{j}}_f = 1$$

$$\hat{\mathbf{i}}_f^T Q Q^T \hat{\mathbf{j}}_f = 0$$

Image from Shree Nayar FPCV videos

# Orthonormality of $M$

- We have computed $(\hat{\mathbf{i}}_f^T, \hat{\mathbf{j}}_f^T)$ for $f = 1, \dots, F$.

$$\left.\begin{array}{l} \hat{\mathbf{i}}_f^T Q Q^T \hat{\mathbf{i}}_f = 1 \\[10pt] \hat{\mathbf{j}}_f^T Q Q^T \hat{\mathbf{j}}_f = 1 \\[10pt] \hat{\mathbf{i}}_f^T Q Q^T \hat{\mathbf{j}}_f = 0 \end{array}\right\} \quad Q \text{ is unknown.}$$

- $Q$ is 3×3 matrix, 9 variables, $3F$ quadratic equations.

- $Q$ can be solved with 3 or more images ($F \geq 3$) using Newton's method.

Final Solution:

$$\boxed{M = U_1\,(\Sigma_1)^{1/2} Q} \qquad \boxed{S = Q^{-1}(\Sigma_1)^{1/2} V_1^T}$$

Camera Motion          Scene Structure

Image from Shree Nayar FPCV videos

# Summary: Orthographic SFM

1. Detect and track feature points.

2. Create the centroid subtracted matrix $W$ of corresponding feature points.

3. Compute SVD of $W$ and enforce rank constraint.

$$W = U \Sigma V^T = U_1 \Sigma_1 V_1^T$$

$$(2F \times 3) \quad (3 \times 3) \quad (3 \times P)$$

4. Set $M = U_1 (\Sigma_1)^{1/2} Q$ and $S = Q^{-1} (\Sigma_1)^{1/2} V_1^T$.

5. Find $Q$ by enforcing the orthonormality constraint.

[Tomasi 1992]

Image from Shree Nayar FPCV videos

# Results from 1992

- Input



Fig. 2a. The "Hotel" stream: four of the 150 frames.

Image from [Tomasi1992]

# Results from 1992

- Resulting 3D reconstruction and new viewpoint rendering



(a)                                                                (b)

*Fig. 4.* Qualitative shape results for the "Hotel" stream: top view of the (a) computed and (b) actual shape.

Image from [Tomasi1992]

# More modern results: Photo Tourism

# First step: how to get correspondence?

- Feature detection and matching

# Feature detection

Detect features using SIFT [Lowe, IJCV 2004]

# Feature detection

Detect features using [SIFT](#) [Lowe, IJCV 2004]

# Video: Feature Detection

https://www.youtube.com/watch?v=4AvTMVD9ig0



Histograms of the gradients

Courtesy of Gil Levi

See also:
https://gilscvblog.com/2013/08/18/a-short-introduction-to-descriptors/

**SIFT - 5 Minutes with Cyrill**

# Feature matching

Match features between each pair of images

# Feature distance

How to measure the difference between $f_1$ and $f_2$, $f_3$ or $f_4$?

- A simple approach: $L_2$-distance, $\| \vec{x}_{f_1} - \vec{x}_{f_i} \|$ (aka SSD)*

*Attention: It is about the distance of the feature vectors (so for SIFT a 128 dimensional space)



$$\vec{x}_{f_1} = (x_1, x_2, \ldots, x_n)$$



$$\vec{x}_{f_2} = (x_1', x_2', \ldots, x_n')$$
$$\vec{x}_{f_3} = (x_1'', x_2'', \ldots, x_n'')$$
$$\vec{x}_{f_4} = (x_1''', x_2''', \ldots, x_n''')$$

# Feature distance

How to measure the difference between $f_1$ and $f_2$?

- A simple approach: $L_2$-distance, $\| \vec{x}_{f_1} - \vec{x}_{f_i} \|$ (aka SSD)



$L_2$-distance: $\left\| \vec{x}_{f_1} - \vec{x}_{f_2} \right\| = \sqrt{(x_1 - x_1')^2 + (x_2 - x_2')^2 + \cdots + (x_n - x_n')^2}$

# Feature distance

How to measure the difference between $f_1$ and $f_2$, $f_3$ or $f_4$?

- A simple approach: $L_2$-distance, $\| \vec{x}_{f_1} - \vec{x}_{f_i} \|$ (aka SSD)



Which $L_2$-distance is the smallest?

$\times$ $L_2(f_1, f_2) = \left\| \vec{x}_{f_1} - \vec{x}_{f_2} \right\|$

$\times$ $L_2(f_1, f_3) = \left\| \vec{x}_{f_1} - \vec{x}_{f_3} \right\|$

$L_2(f_1, f_4) = \left\| \vec{x}_{f_1} - \vec{x}_{f_4} \right\|$ ✓

# Feature distance

Attention: this can result in small distances for ambiguous (false) matches



$I_1$                                                $I_2$

# Feature distance

- Better approach:  ratio distance = $\| f1 - f2 \| / \| f1 - f2' \|$
  - f2 is the best SSD match to f1 in I2
  - f2' is the 2nd best SSD match to f1 in I2
  - gives large values for ambiguous matches



$I_1$                  $I_2$

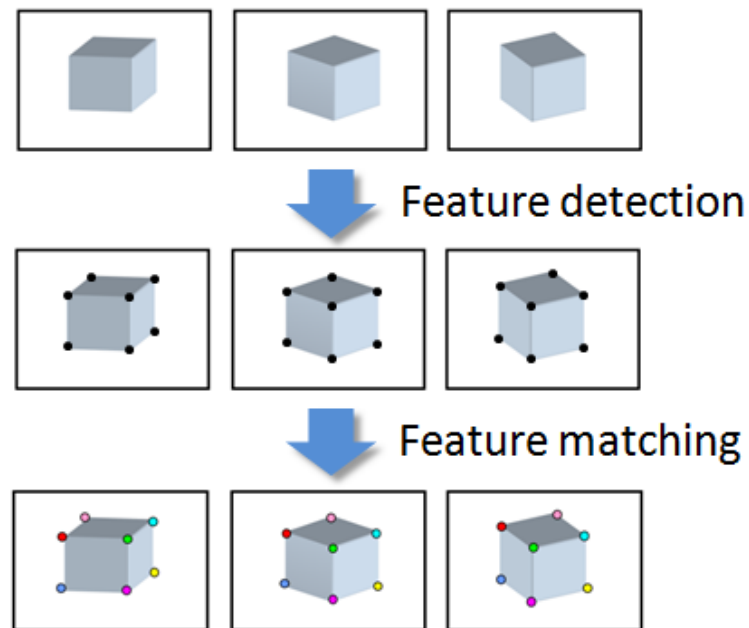# Feature matching

Refine matching using [RANSAC](#) to estimate fundamental matrix between each pair

# Correspondence estimation

- Link up pairwise matches to form connected components of matches across several images



| Image 1 | Image 2 | Image 3 | Image 4 |

# Input to Structure from Motion

# Structure from motion



minimize

$$g(\mathbf{R}, \mathbf{T}, \mathbf{X})$$

*non-linear least squares*

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

Camera 1
$R_1, t_1$

Camera 2
$R_2, t_2$

Camera 3
$R_3, t_3$

# Problem size

- What are the variables?
  - Cameras and points
- How many variables per camera?
  - 6 (if calibrated), more if uncalibrated
- How many variables per point?
  - 3

- Trevi Fountain collection
    466 input photos
  + > 100,000 3D points
    = very large optimization problem

# Structure from motion

- Minimize sum of squared reprojection errors:

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

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

*predicted* image location

*observed* image location

- Minimizing this function is called *bundle adjustment*
  - Optimized using non-linear least squares, e.g. Levenberg-Marquardt algorithm

# Is SfM always uniquely solvable?

- No…

# Incremental structure from motion

# Complete SfM pipeline



[images from Hays2018, Szeliski2020 and Schönberger2016]

# Incremental structure from motion
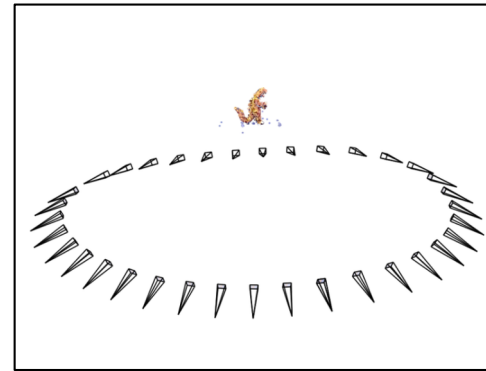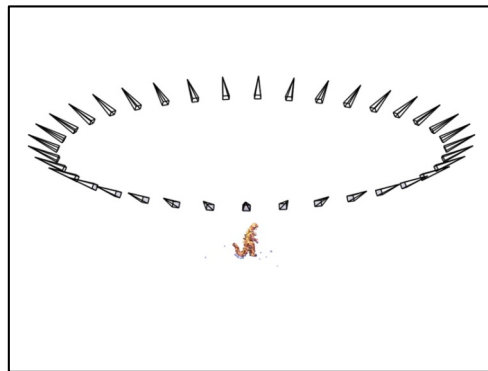


Video from Noah Snavely: [3D Old City of Dubrovnik](#)

# BigSFM: Reconstructing the World from Internet Photos



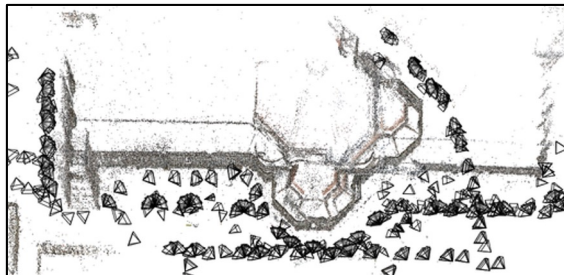Images from https://www.cs.cornell.edu/projects/bigsfm/ check the site for more results.
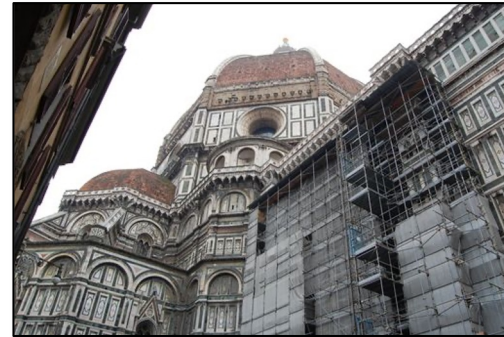
# SfM – Failure cases

- Necker reversal

# Structure from Motion – Failure cases

- Repetitive structures: Symmetries in man-made scenes
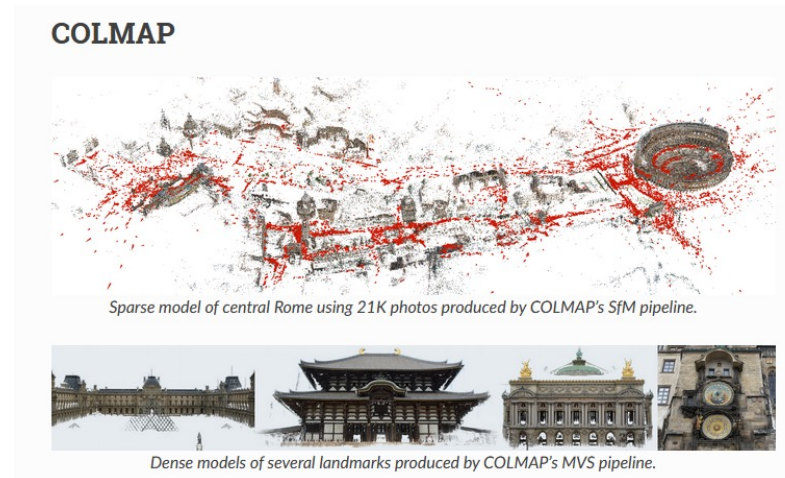
# SfM applications

- 3D modeling

- Surveying

- Robot navigation and mapmaking

- Virtual and augmented reality

- Visual effects ("Match moving")
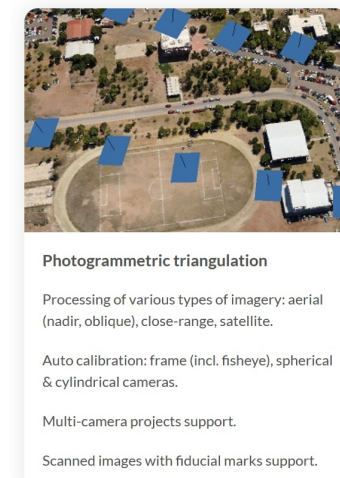  - https://www.youtube.com/watch?v=RdYWp70P_kY

# SfM implementations

- The scientific standard (foss):
  https://colmap.github.io/



**COLMAP**

Sparse model of central Rome using 21K photos produced by COLMAP's SfM pipeline.

Dense models of several landmarks produced by COLMAP's MVS pipeline.

- Agisoft Metashape:
  https://www.agisoft.com/features/professional-edition/
  – Licenses available at HFU



**Photogrammetric triangulation**

Processing of various types of imagery: aerial (nadir, oblique), close-range, satellite.

Auto calibration: frame (incl. fisheye), spherical & cylindrical cameras.

Multi-camera projects support.

Scanned images with fiducial marks support.

- see more at https://tu-dresden.de/ Photogrammetrie

# Applications: Virtual Reality & Augmented Reality



Oculus
https://www.youtube.com/watch?v=KOG7yTz1iTA
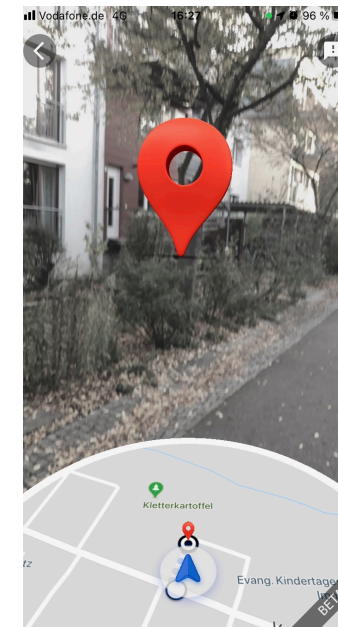
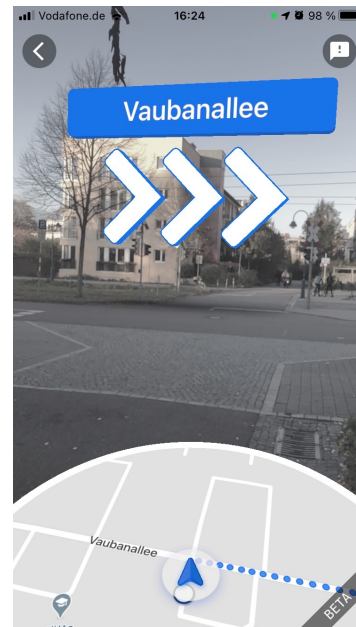

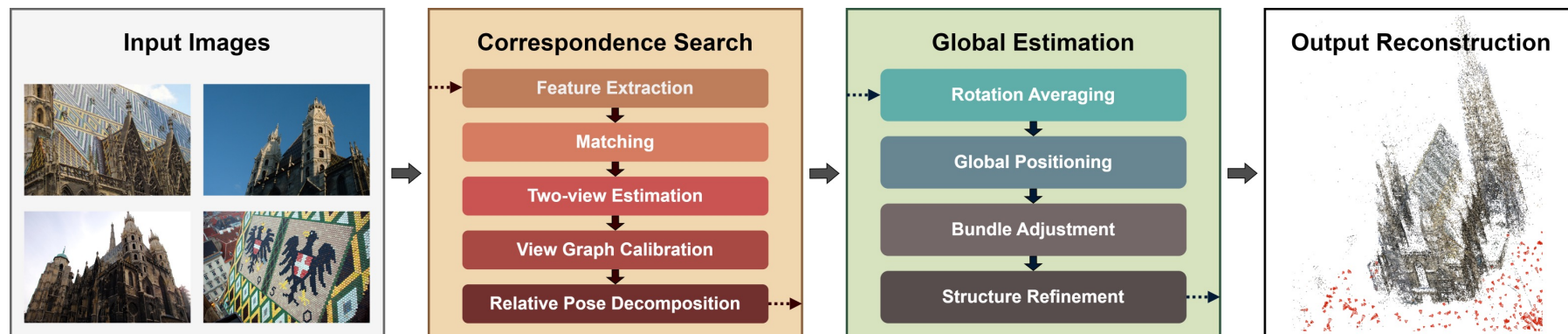Hololens
https://www.youtube.com/watch?v=FMtvrTGnP04

# Application: AR walking directions

- Feature based localization → ARCore 2022

# New trends

- [GLOMAP](#)
  - Improvement of COLMAP

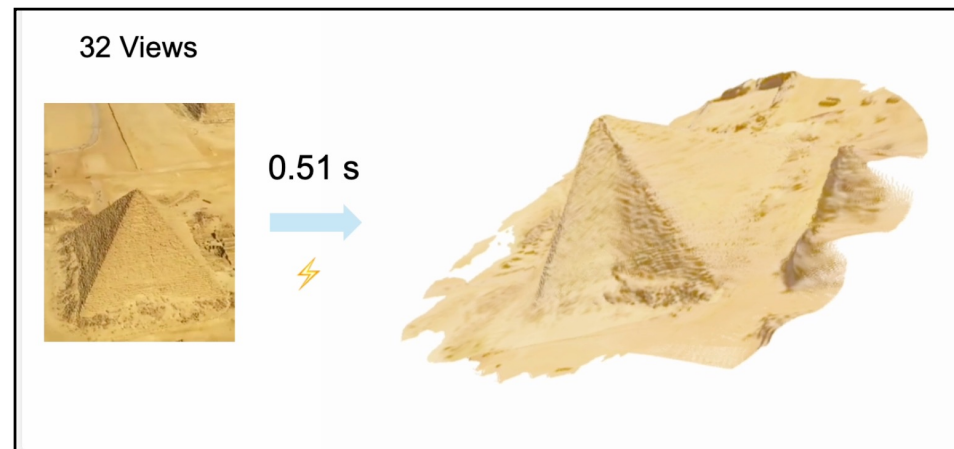# New trends

Dust3r and Mast3r (Naver Labs)
- the idea: do not start from scratch for each scene
- using large visual foundation models (similar to DepthAnything)

VGGT (Oxford + Meta AI)



- More to be explored in the next workshop.