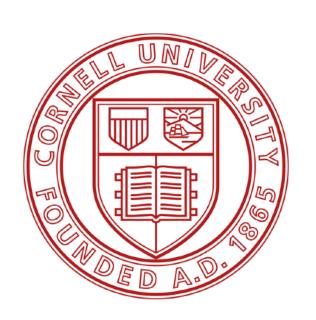
Lecture 21: Structure from motion

CS 5670: Introduction to Computer Vision



Announcements

PS5 out tonight (panorama stitching)

Triangulation

Given projection $\mathbf{p_i}$ of unknown 3D point \mathbf{X} in two or more images (with known cameras $\mathbf{P_i}$), find \mathbf{X}



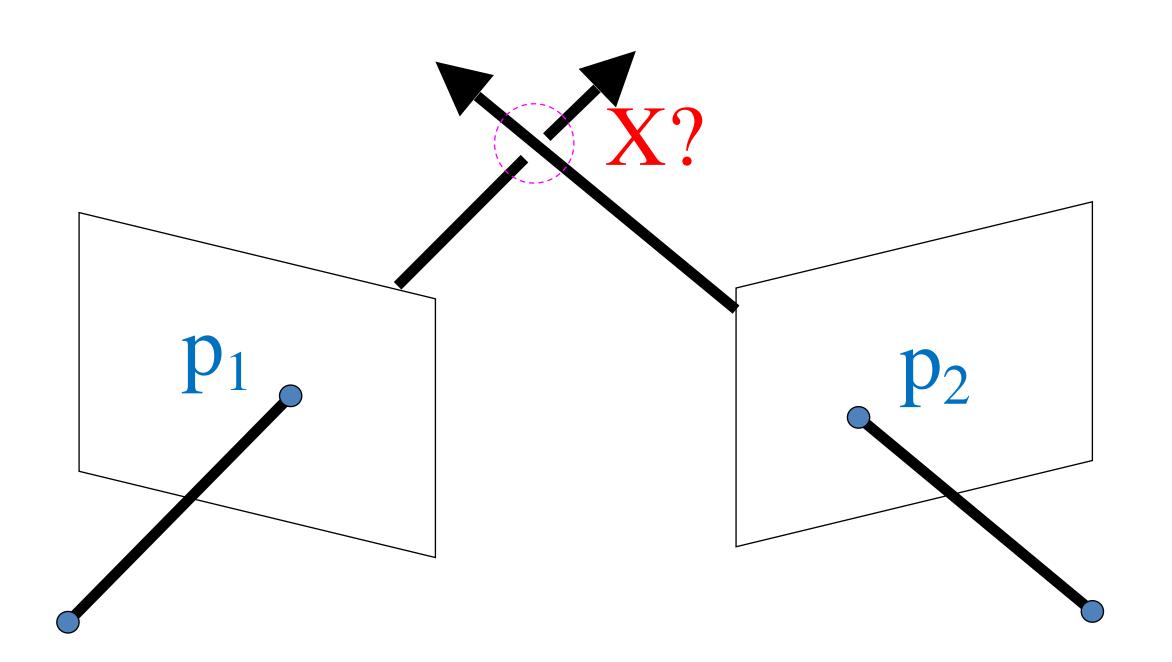




Source: D. Fouhey

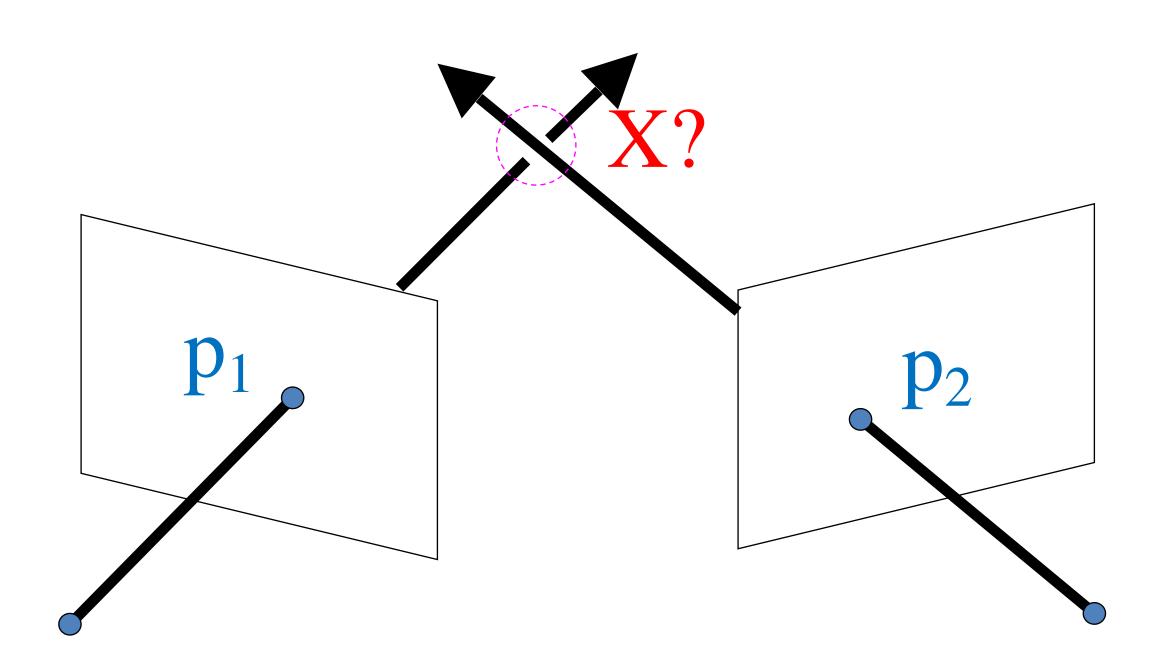
Triangulation

Given projection $\mathbf{p_i}$ of unknown 3D point \mathbf{X} in two or more images (with known camera projection matrices $\mathbf{P_i}$), find \mathbf{X}



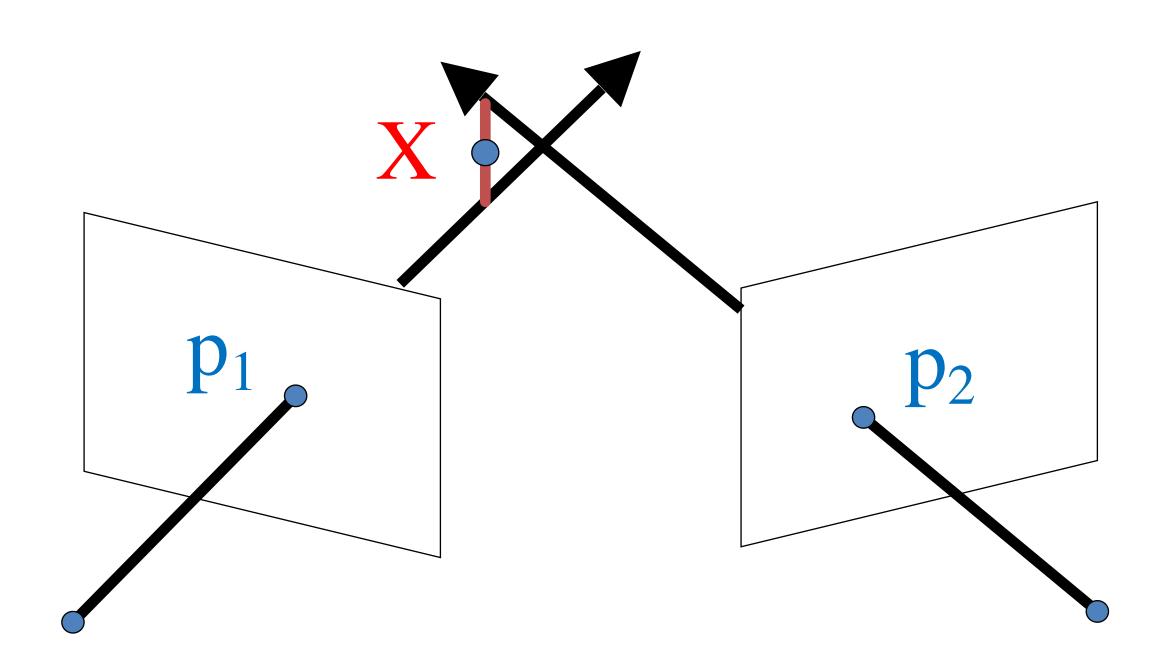
Triangulation

Rays in principle should intersect, but in practice usually don't exactly due to noise, numerical errors.



Triangulation: Geometry

Find shortest segment between viewing rays, set X to be the midpoint of the segment.

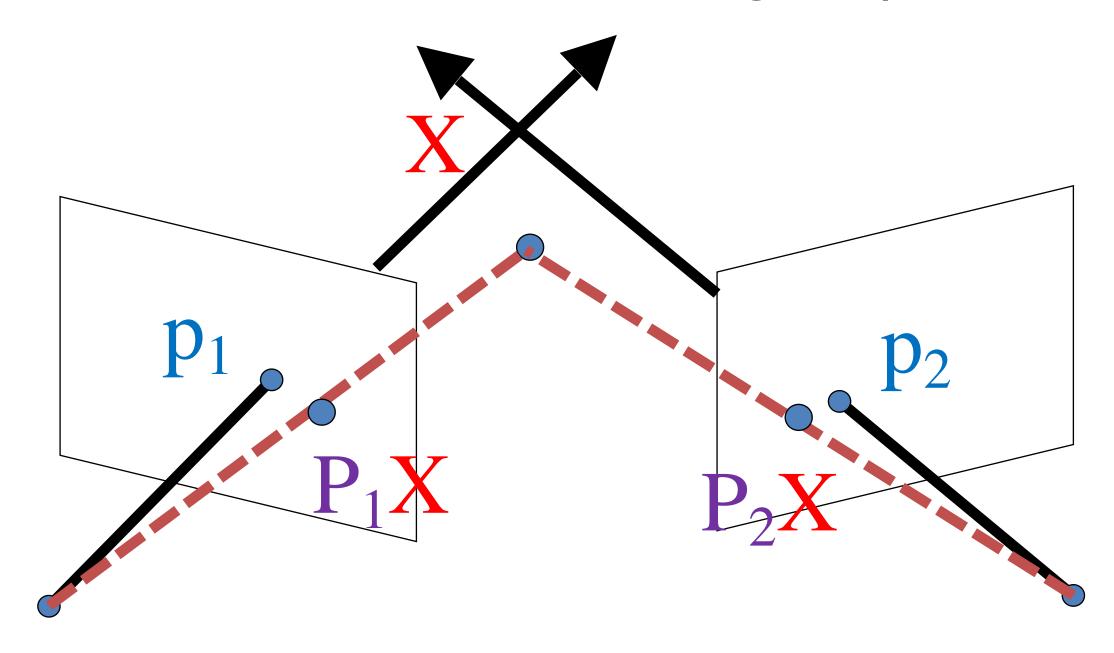


Triangulation: Non-linear Optim.

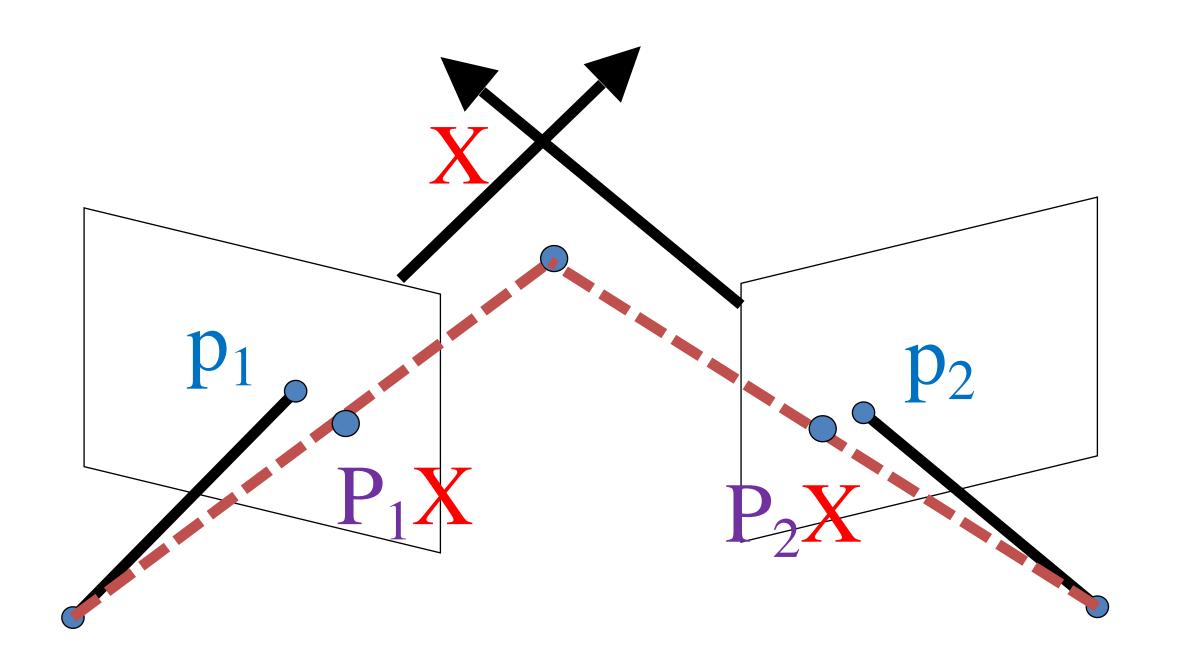
Use nonlinear least squares. Find X that minimizes:

$$d(p_1, P_1X)^2 + d(p_2, P_2X)^2$$

where d is distance in image space.



Triangulation: Linear method



First: A better way to handle homogeneous coordinates in linear optimization

Projection in homogeneous coordinates.

$$p_i \equiv PX_i$$

i.e., PX; & p; are proportional/scaled copies of each other

$$p_i = \lambda P X_i, \ \lambda \neq 0$$

This implies their cross product is **0**, since

$$a \times b = ||a|| ||b|| \sin(\theta).$$

$$p_i \times PX_i = 0$$

Handles the "divide by 0" issue when solving.

Triangulation: Linear method

$$p_1 \equiv P_1 X$$
 $p_1 \times P_1 X = 0$ $[p_{1x}]P_1 X = 0$ $p_2 \equiv P_2 X$ $p_2 \times P_2 X = 0$ $[p_{2x}]P_2 X = 0$

Cross product as matrix

$$\mathbf{a} \times \mathbf{b} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = \begin{bmatrix} \mathbf{a}_x \end{bmatrix} \mathbf{b}$$

$$[p_{1x}]P_1X = 0$$

$$[p_{2x}]P_2X = 0$$

$$([p_{1x}]P_1)X = 0$$

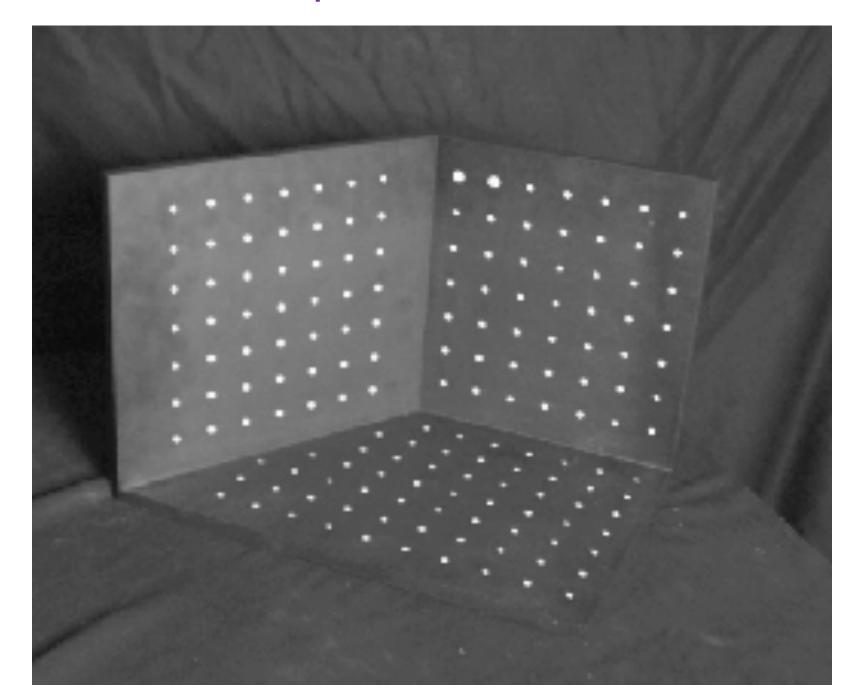
$$([p_{2x}]P_2)X = 0$$
Two equations per camera for 3 unknown in X

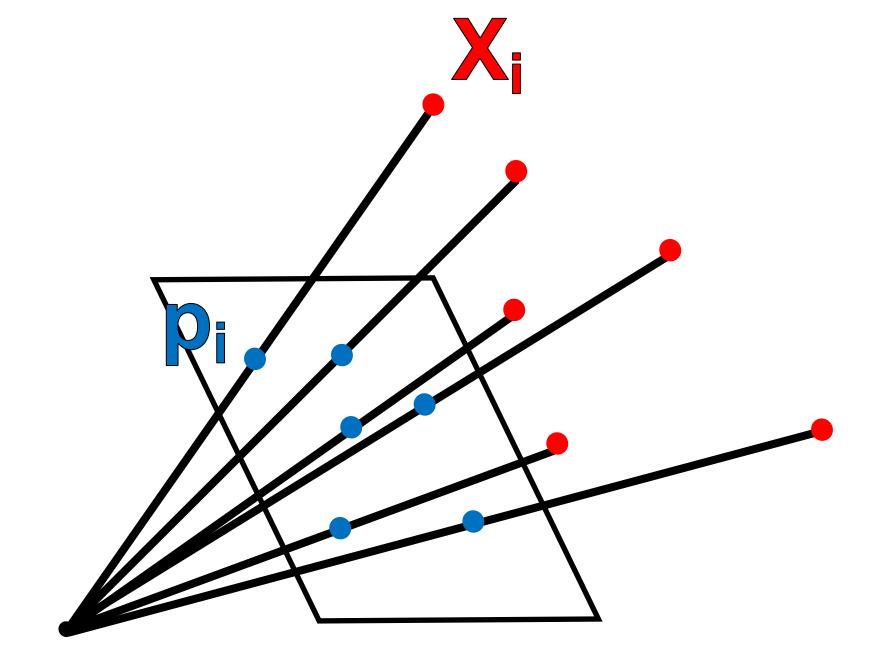
Camera calibration

• Can we estimate camera pose from points? Known as camera calibration or camera resectioning.

• Given n points with known 3D coordinates X_i and known image projections p_i ,

estimate the camera parameters.





Camera Calibration: Linear Method

$$p_i \equiv PX_i$$

Remember (from geometry): this implies $\mathbf{MX_i}$ & $\mathbf{p_i}$ are proportional/scaled copies of each other $p_i = \lambda PX_i, \ \lambda \neq 0$

Recall that this implies their cross product is 0

$$p_i \times PX_i = 0$$

Camera Calibration: Linear Method

$$p_i \times PX_i = 0$$

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix} \times \begin{bmatrix} P_1 P_i \\ P_2 X_i \\ P_3 X_i \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Using these constraints, you can derive:

$$\begin{bmatrix} \mathbf{0}^T & -\mathbf{X}_i^T & \mathbf{v}_i \mathbf{X}_i^T \\ \mathbf{X}_i^T & \mathbf{0}^T & -\mathbf{u}_i \mathbf{X}_i^T \\ -\mathbf{v}_i \mathbf{X}_i^T & \mathbf{u}_i \mathbf{X}_i^T & \mathbf{0}^T \end{bmatrix} \begin{bmatrix} P_1^T \\ P_2^T \\ P_3^T \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

where P_i is row i of P

Camera calibration

- The linear solution does not optimize the right objective function.
- Optimize using nonlinear least squares:

$$\sum \left\| \operatorname{proj}(PX_i) - \left[u_i, v_i\right]^T \right\|_2^2$$

- Can initialize using the linear solution.
- Other advantages: can also add radial distortion,
 not optimize over known variables, add constraints

Structure from motion

- We can estimate points from cameras and cameras from points. Can we do both at once?
- Given many images, how can we...
 - 1. Figure out where they were all taken from?
 - 2. Build a 3D model of the scene?

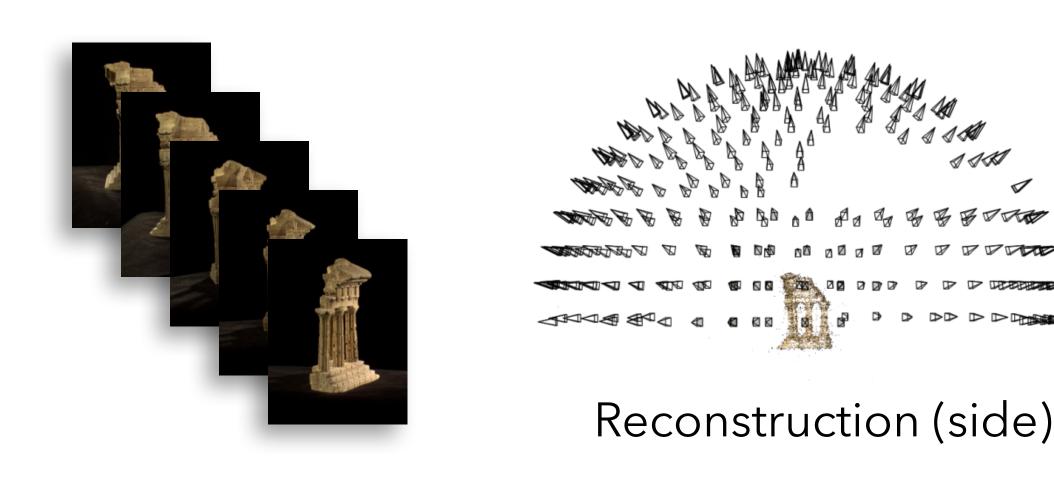


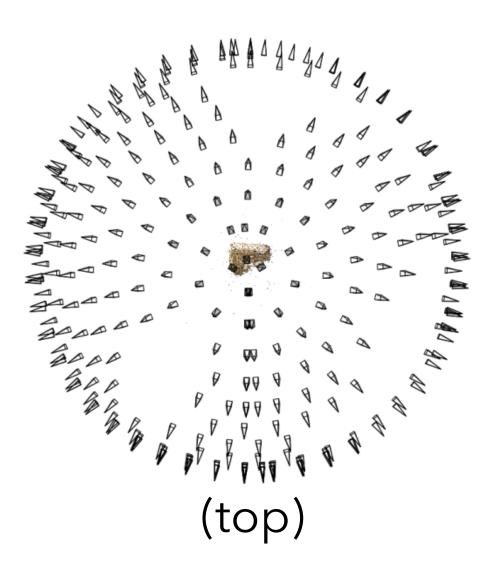
This is the **structure from motion** problem.

Today

- Structure from motion
- Multi-view stereo
- Radiance fields

Structure from motion





 $p_{i,j} = (u_{i,j}, v_{i,j})$

- Input: images with pixels in correspondence
- 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

Camera calibration & triangulation

- Suppose we know 3D points
 - And have matches between these points and an image
 - Computing camera parameters similar to homography estimation
- Suppose we have know camera parameters, each of which observes a point
 - We can solve for the 3D location
- Seems like a chicken-and-egg problem, but in SfM we can solve both at once



Source: N. Snavely

Example: Photo Tourism



Feature detection

- Same process as with homography estimation
- Detect features using SIFT



























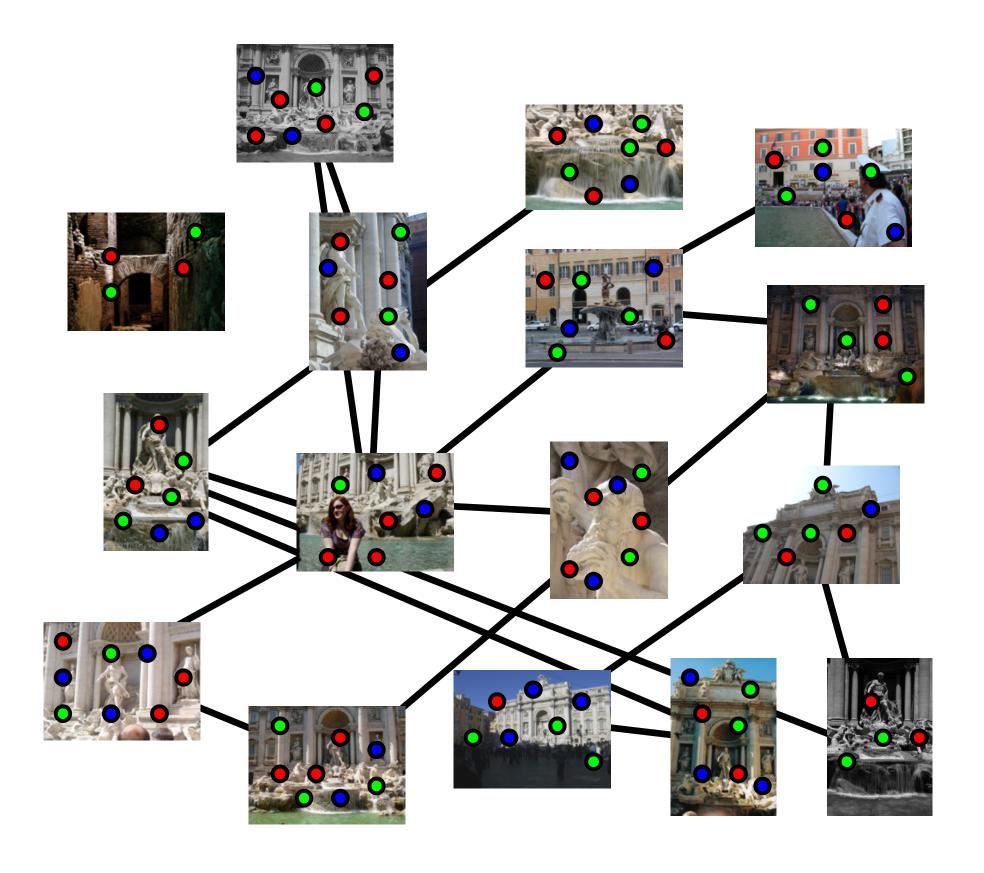






Feature matching

Match features between each pair of images



Feature matching

- Remove bad matches using ratio test.
- Other tricks: throw out matches that aren't on epipolar lines.

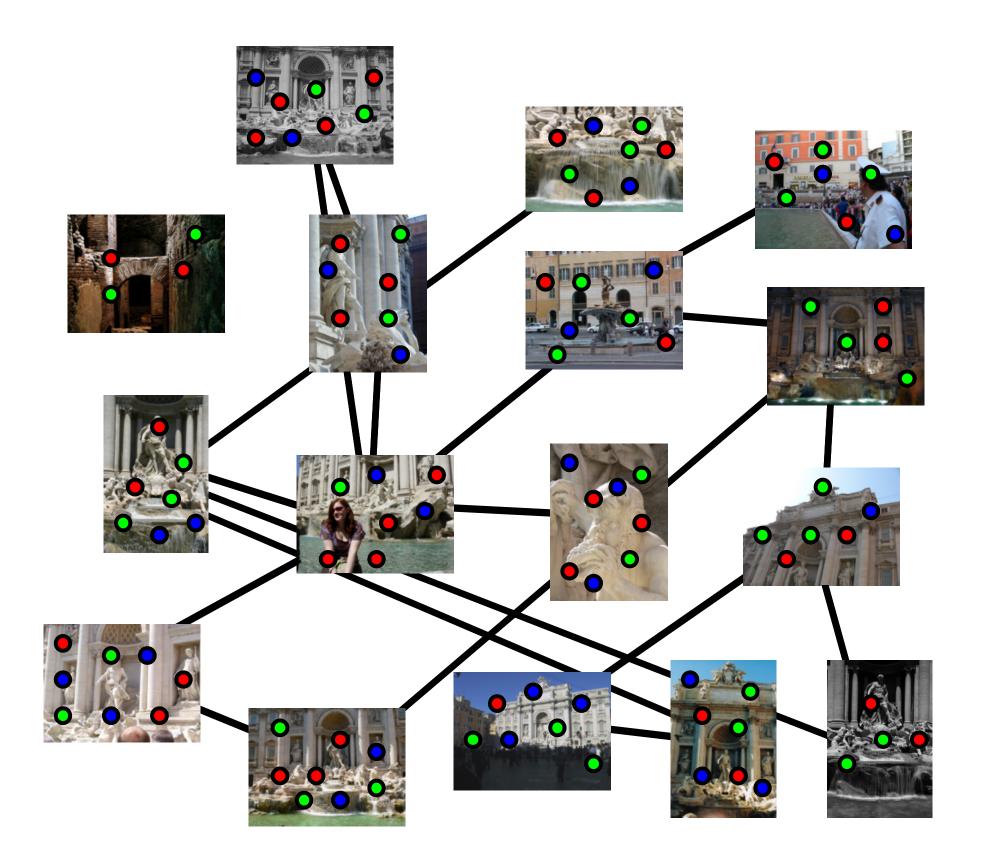
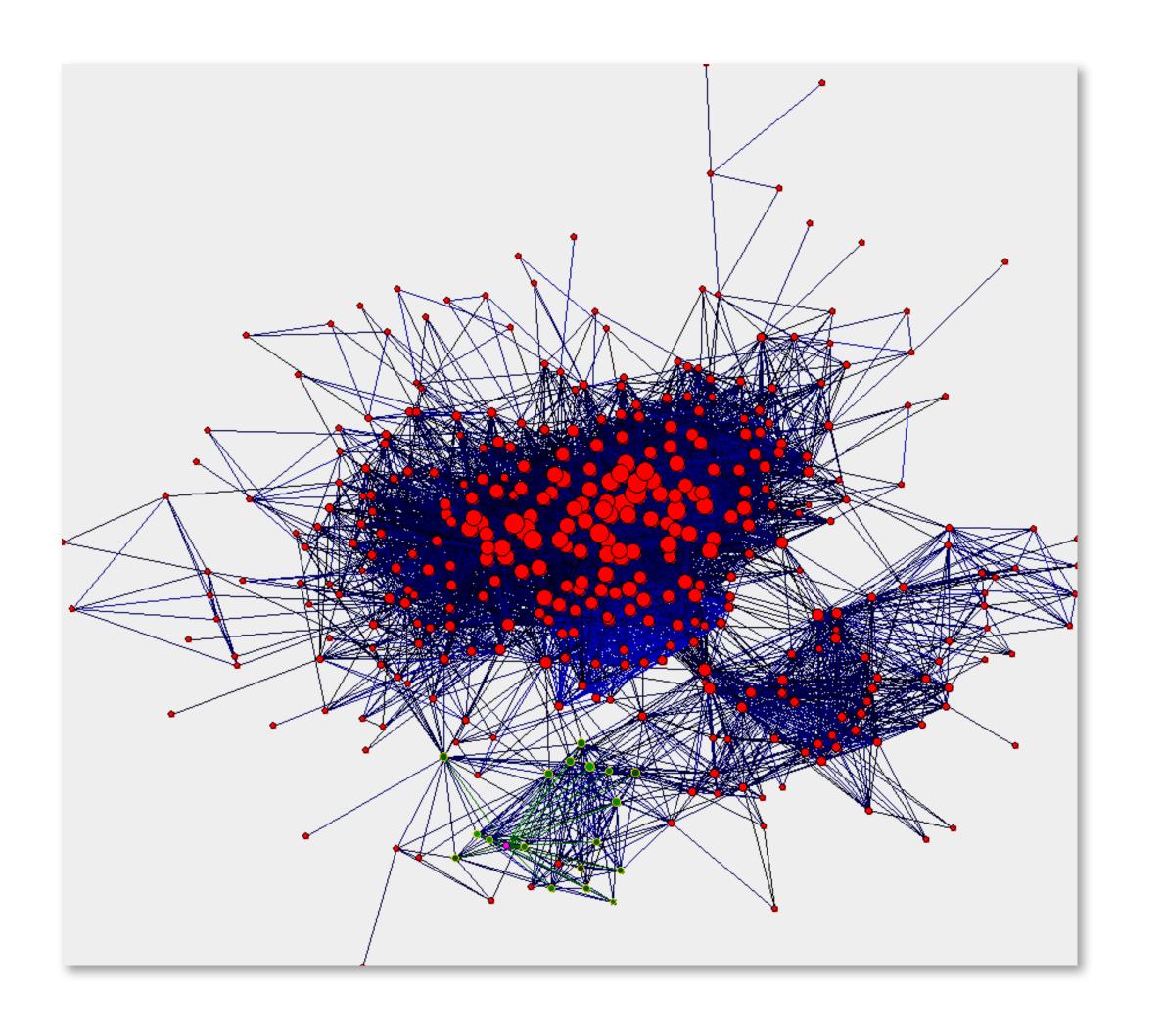
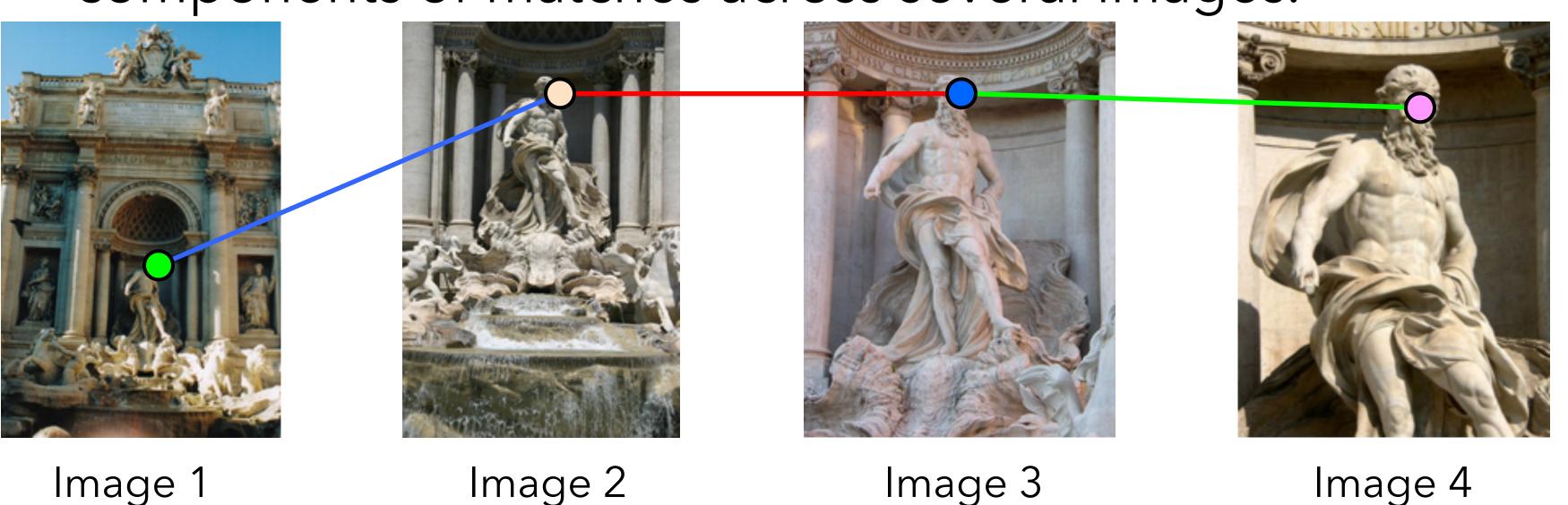


Image connectivity graph



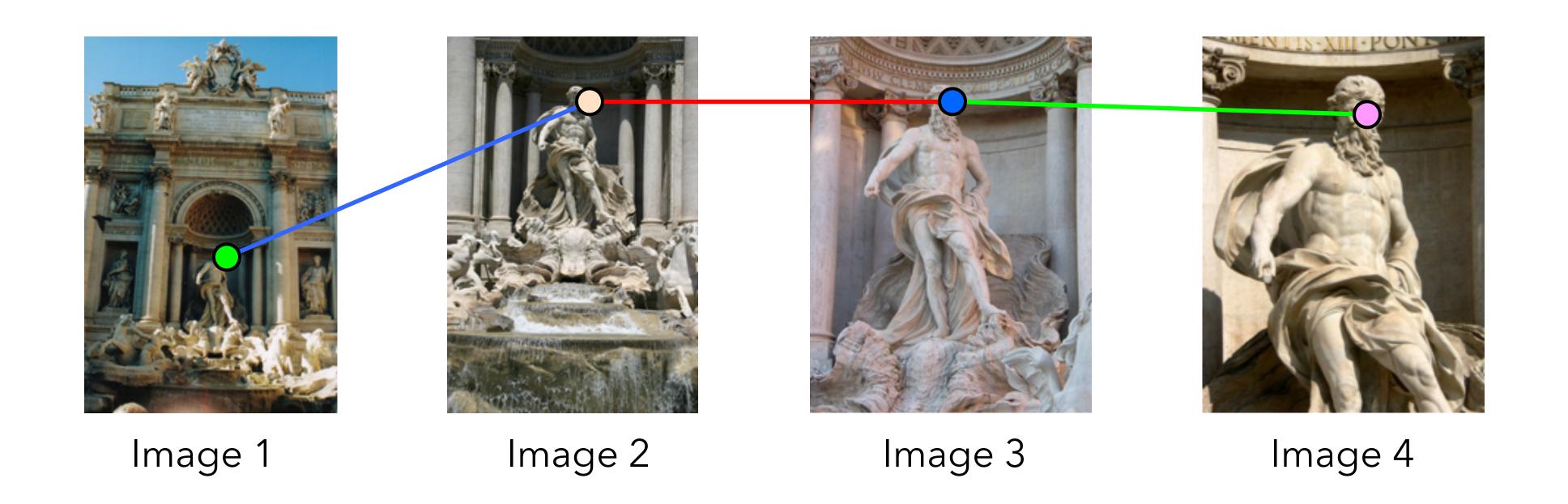
Correspondence estimation

- Track each feature across the dataset.
- Link up pairwise matches to form connected components of matches across several images.

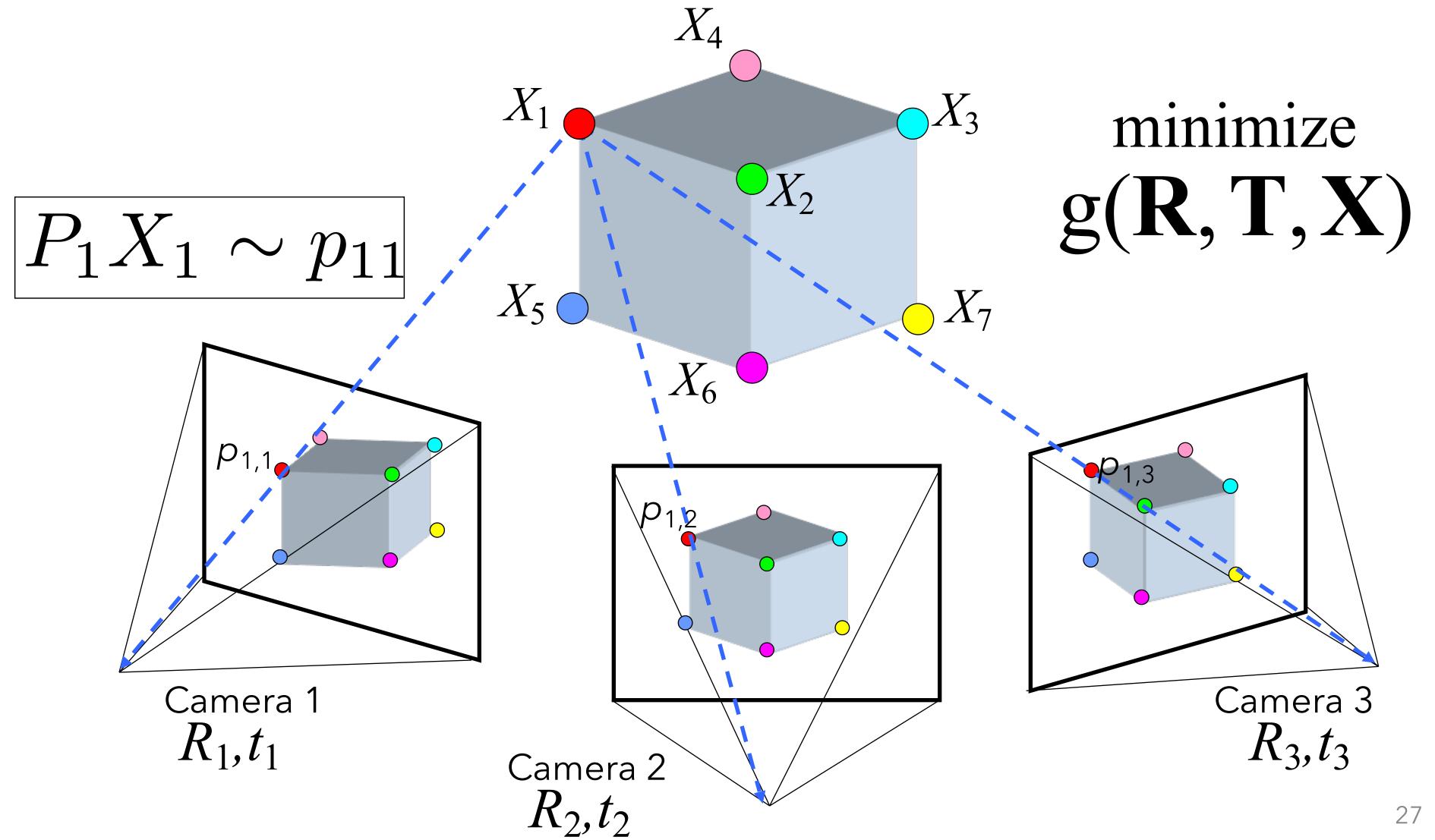


Correspondence estimation

A point track: the same 3D point projects to all 4 image positions.



Structure from motion



Structure from motion

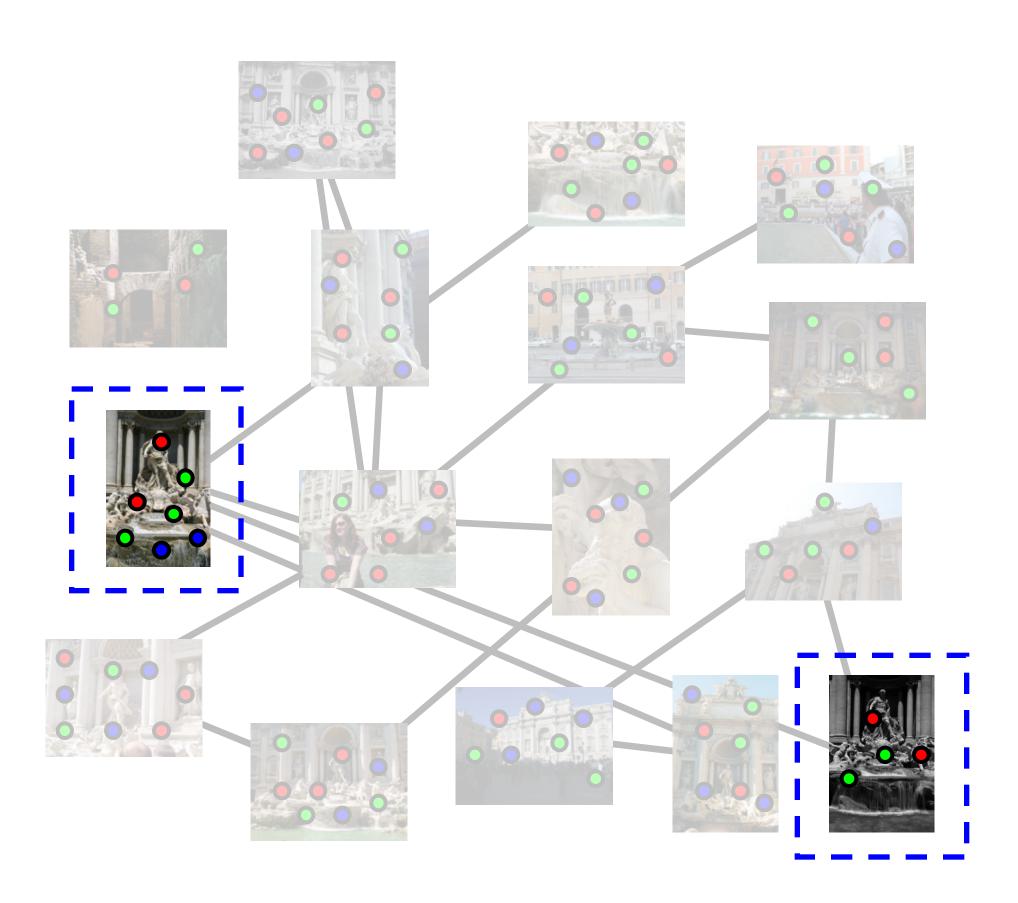
• Minimize sum of squared reprojection errors:

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

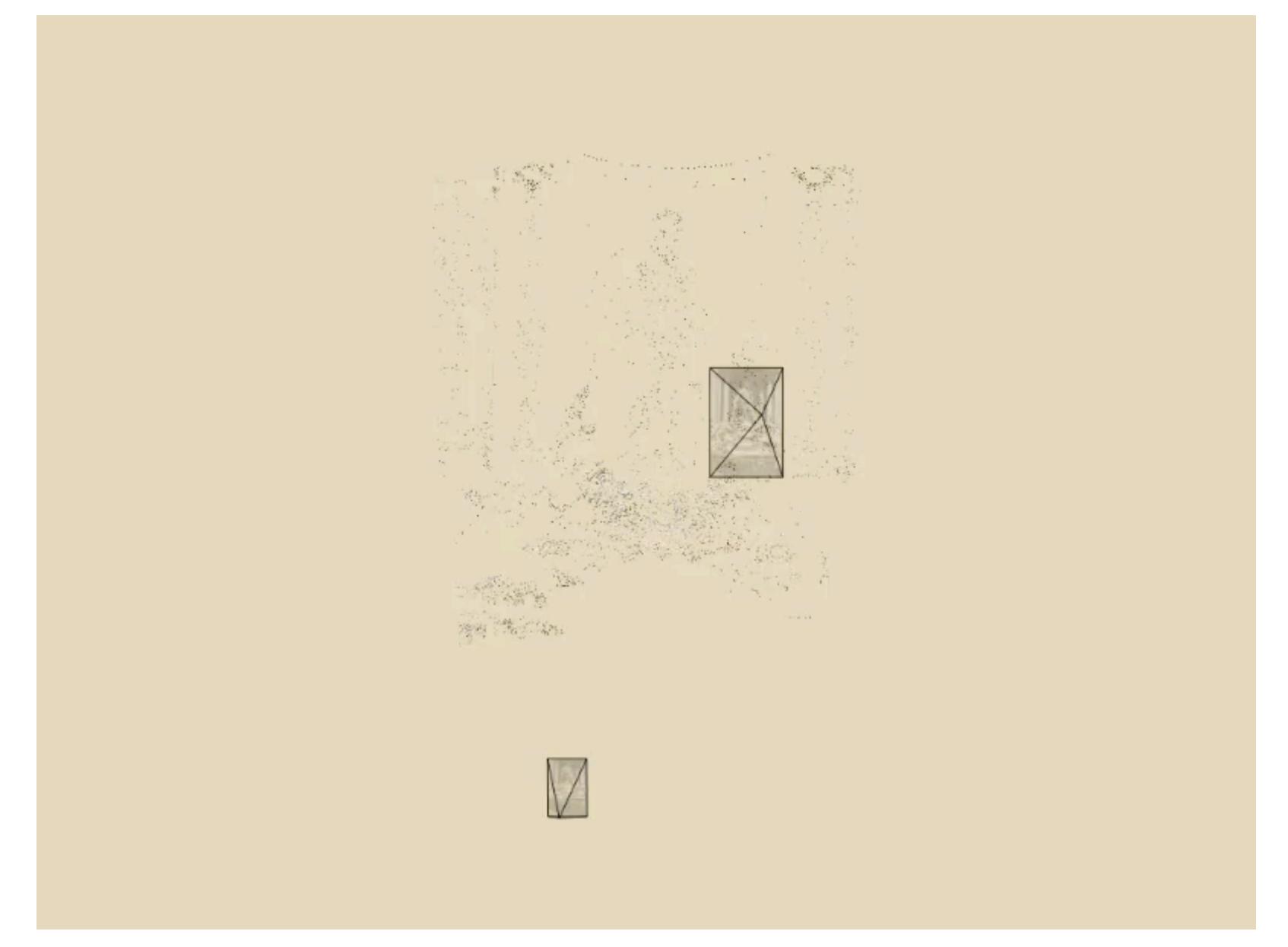
$$\downarrow \qquad \qquad \qquad predicted \qquad observed \qquad image location \qquad indicator variable: \qquad is point i visible in image j?$$

- Minimizing this function is called bundle adjustment.
 - Optimized using non-linear least squares
- Lots of outliers: use robust loss functions (e.g., Huber) and solve incrementally

Incremental structure from motion

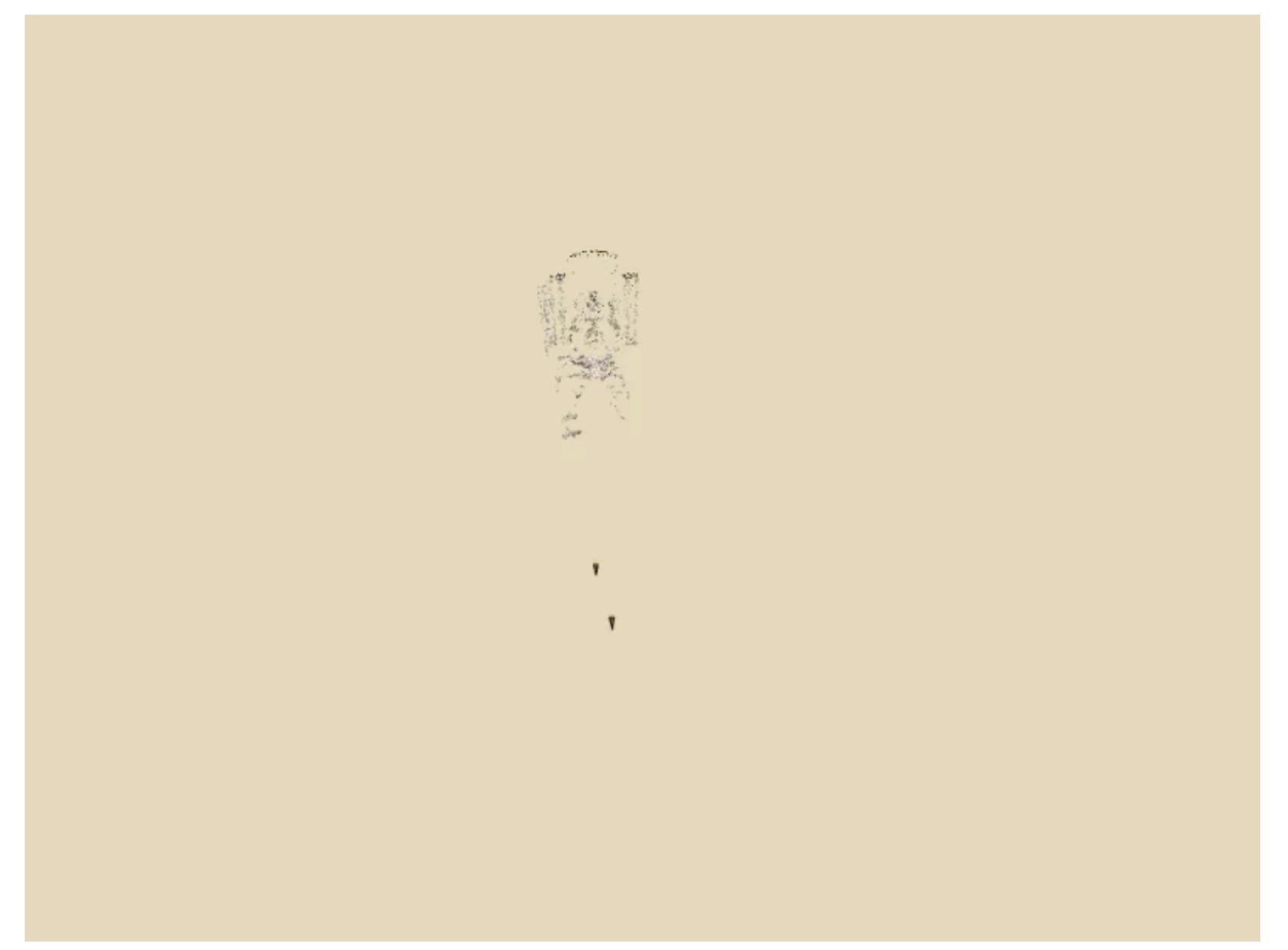


Incremental structure from motion



Source: N. Snavely

Incremental structure from motion



Source: N. Snavely

Photo Tourism Exploring photo collections in 3D

Noah Snavely Steven M. Seitz Richard Szeliski

University of Washington Microsoft Research

SIGGRAPH 2006



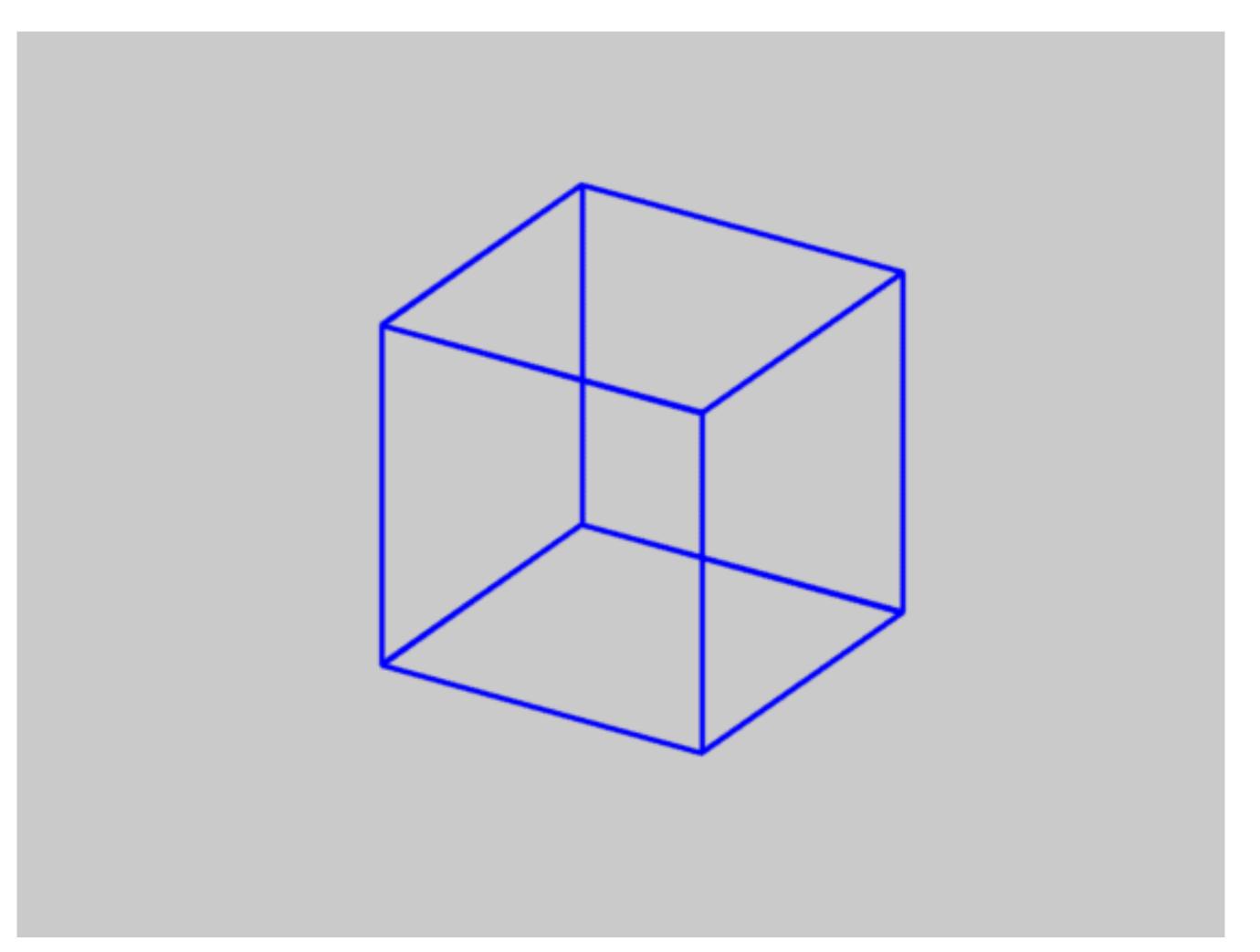




Dubrovnik, Croatia. 4,619 images (out of an initial 57,845). Total reconstruction time: 23 hours on 352 cores

Is SfM always uniquely solvable?

No. Consider the Necker cube:



Source: N. Snavely

Is SfM always uniquely solvable?

Two interpretations:

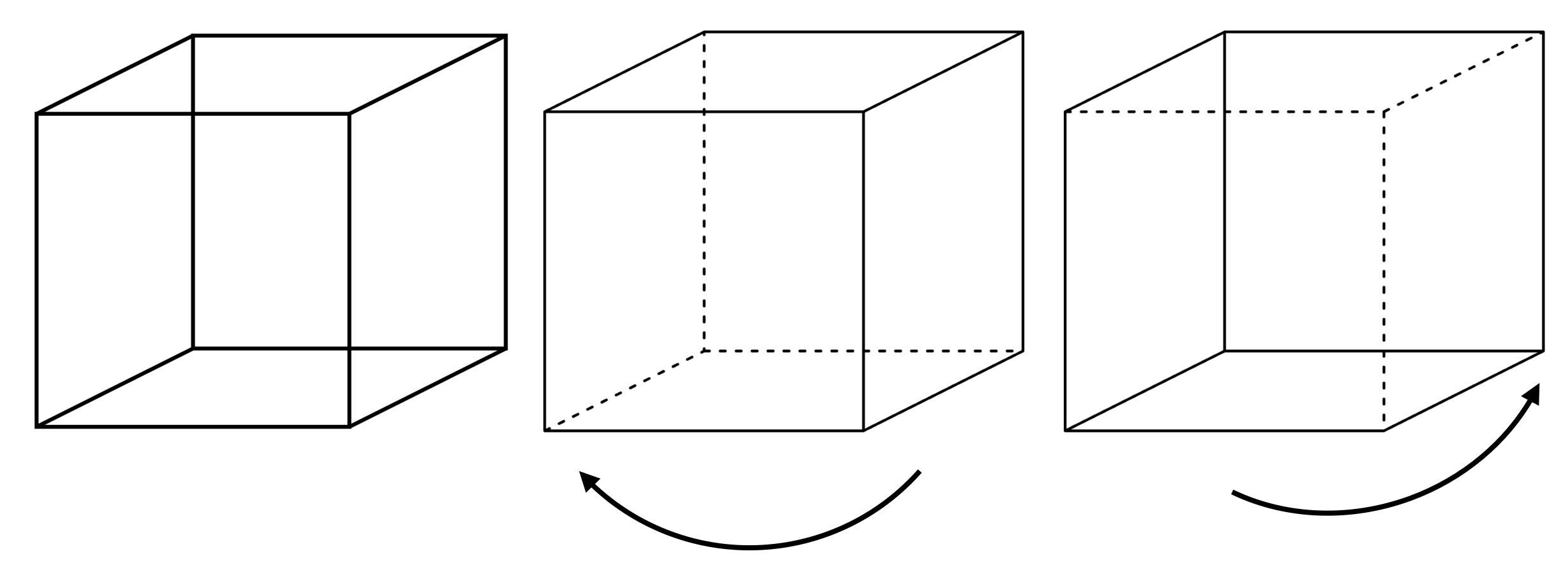
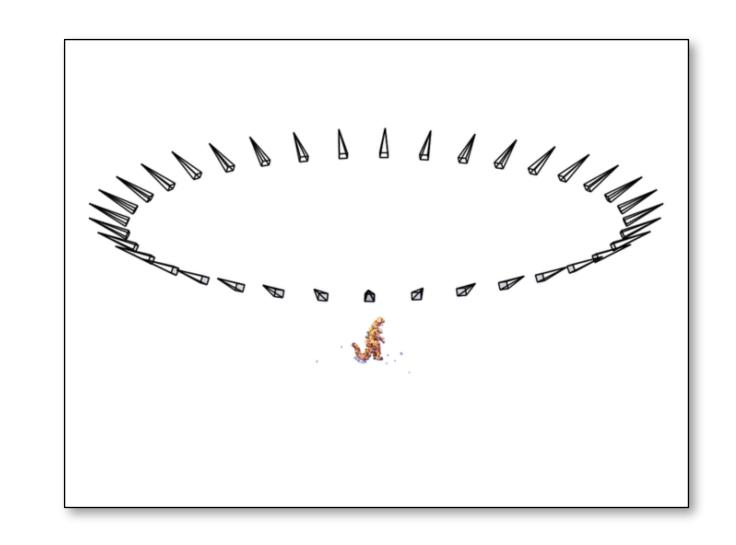


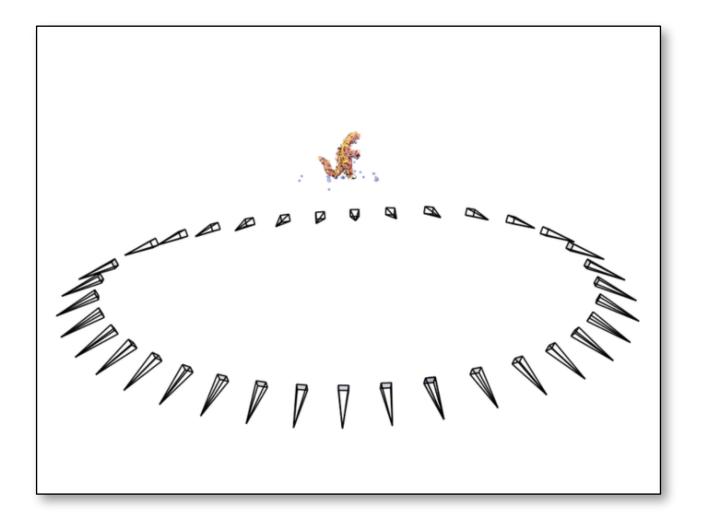
Image source: Wikipedia

Failure case: Necker reversal

- Under orthographic camera, object rotation by θ produces same image as mirror image rotated by $-\theta$.
- Can occur in perspective cameras, e.g., when objects are far from the camera







Ambiguity up to similarity transformation

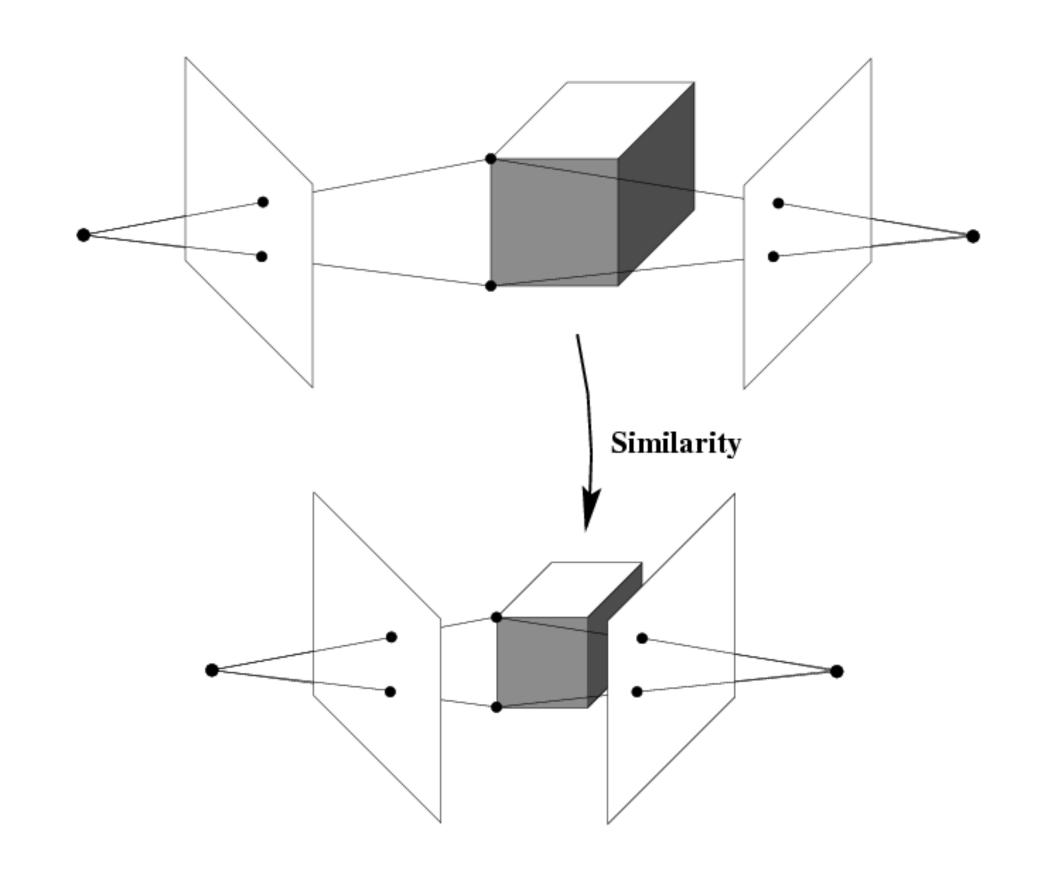
$$x \cong PX = (PQ_S^{-1})(Q_SX)$$

$$3 \times 3$$

$$\text{rotation}$$

$$\text{matrix}$$

$$Q_S = \begin{bmatrix} sR & t \\ oT & 1 \end{bmatrix}$$



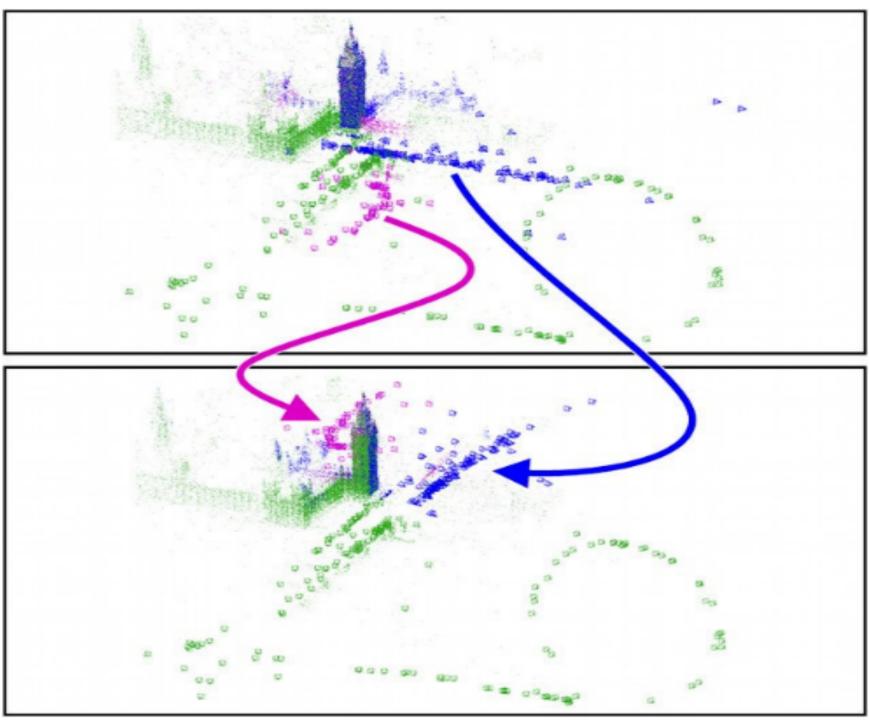
Special case: scale ambiguity

Repetitive structures cause catastrophic failures



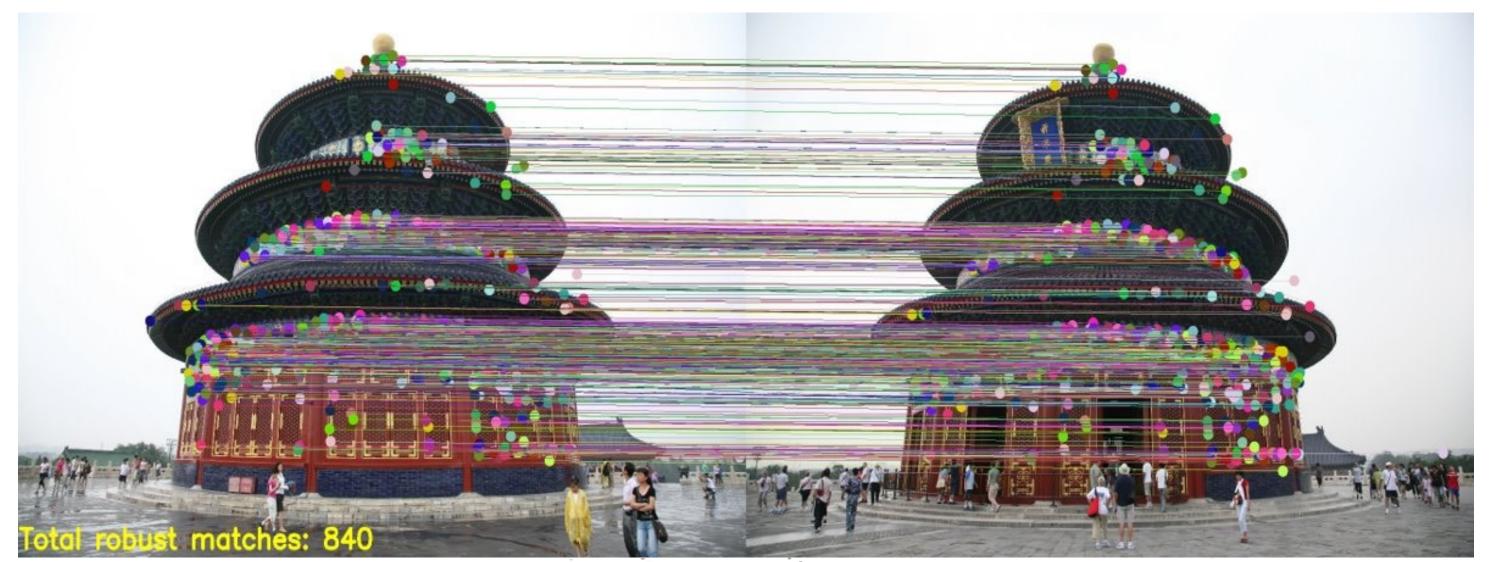


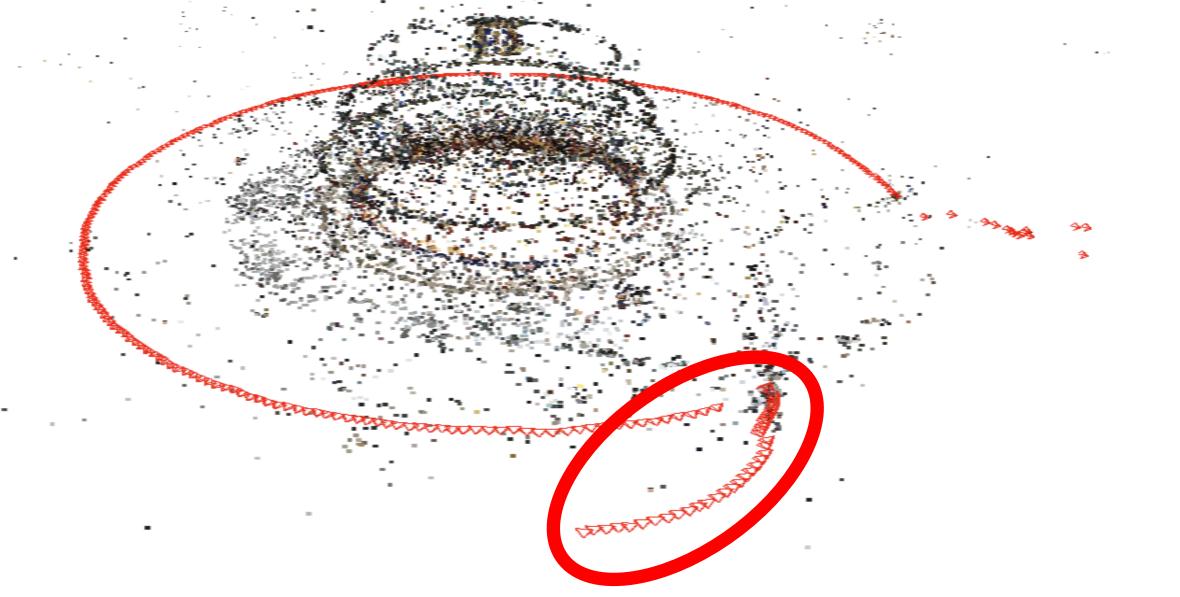




https://demuc.de/tutorials/cvpr2017/sparse-modeling.pdf

Repetitive structures cause catastrophic failures

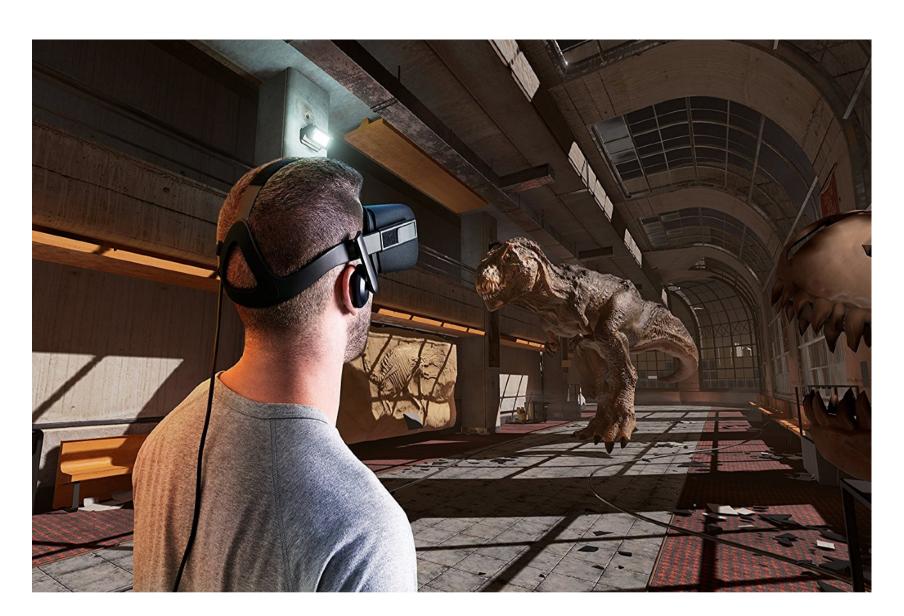




Can also reconstruct from video



Applications: Visual Reality & Augmented Reality





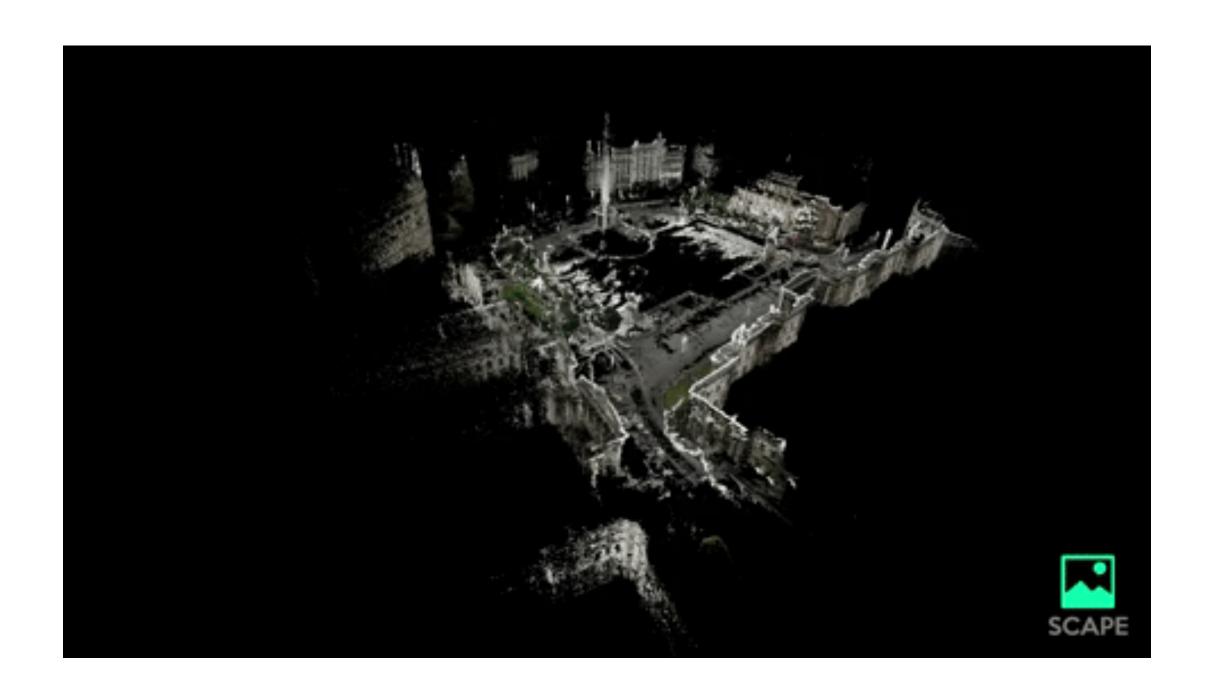
Oculus

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

Hololens

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

Application: Simultaneous localization and mapping (SLAM)



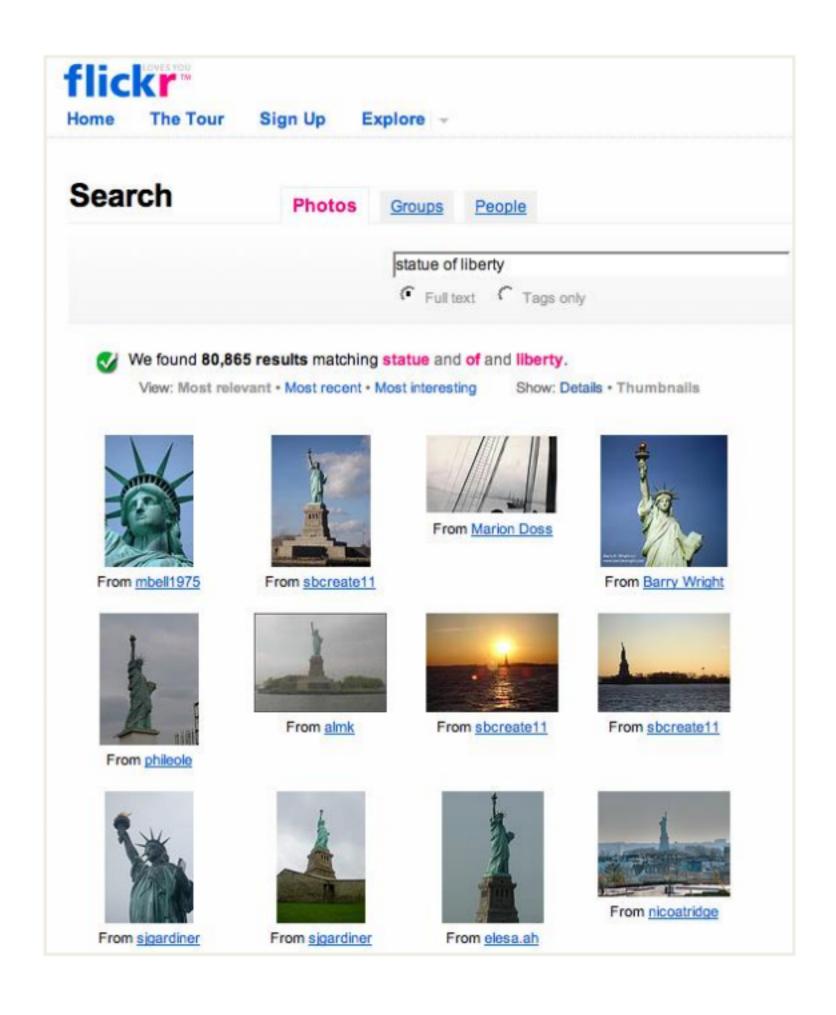
Scape: Building the 'AR Cloud': Part Three –3D Maps, the Digital Scaffolding of the 21st Century

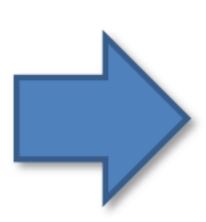
https://medium.com/scape-technologies/building-the-ar-cloud-part-three-3d-maps-the-digital-scaffolding-of-the-21st-century-465fa55782dd

Source: N. Snavely

Today

- Structure from motion
- Multi-view stereo
- Stereo matching algorithms

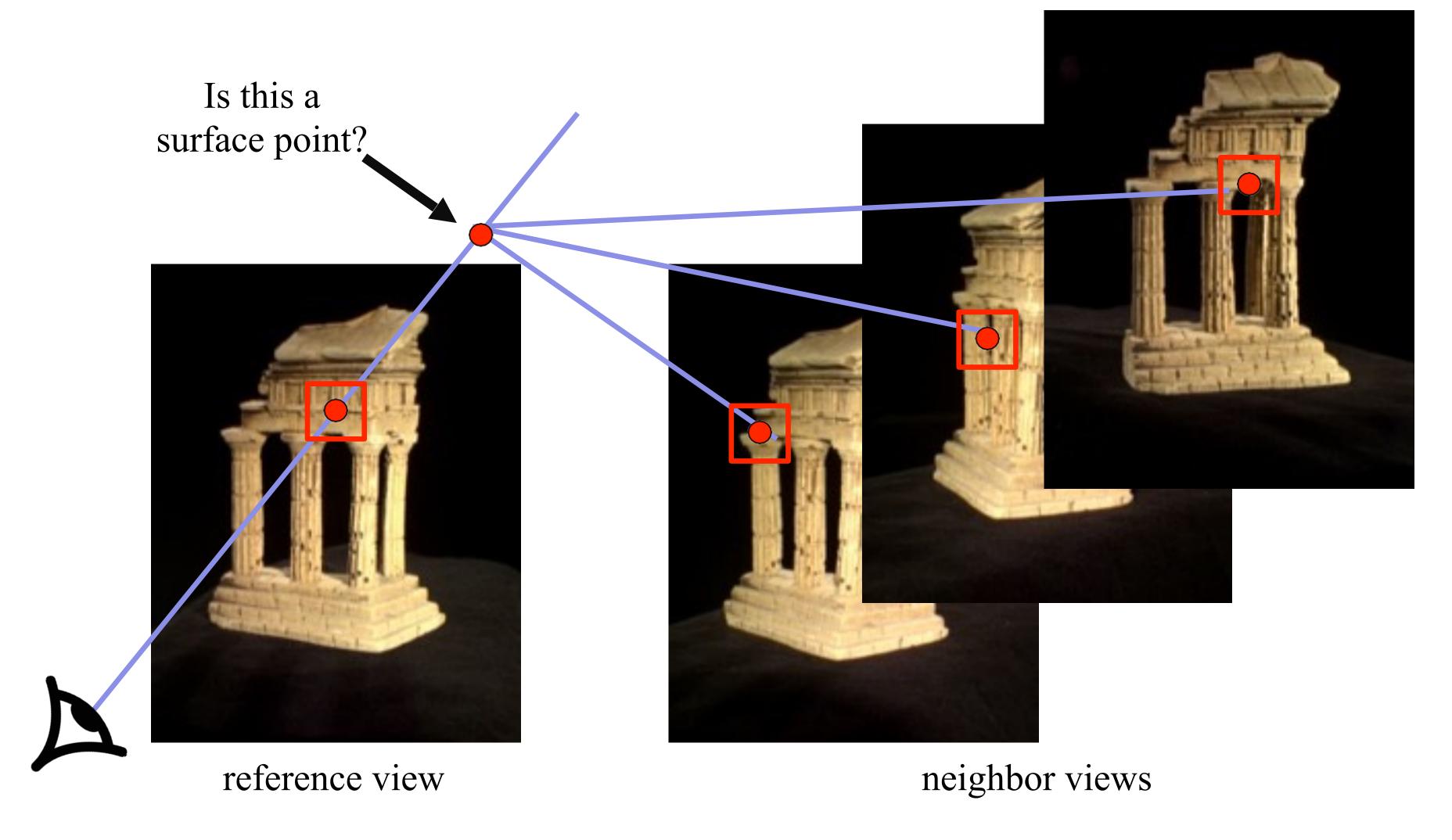






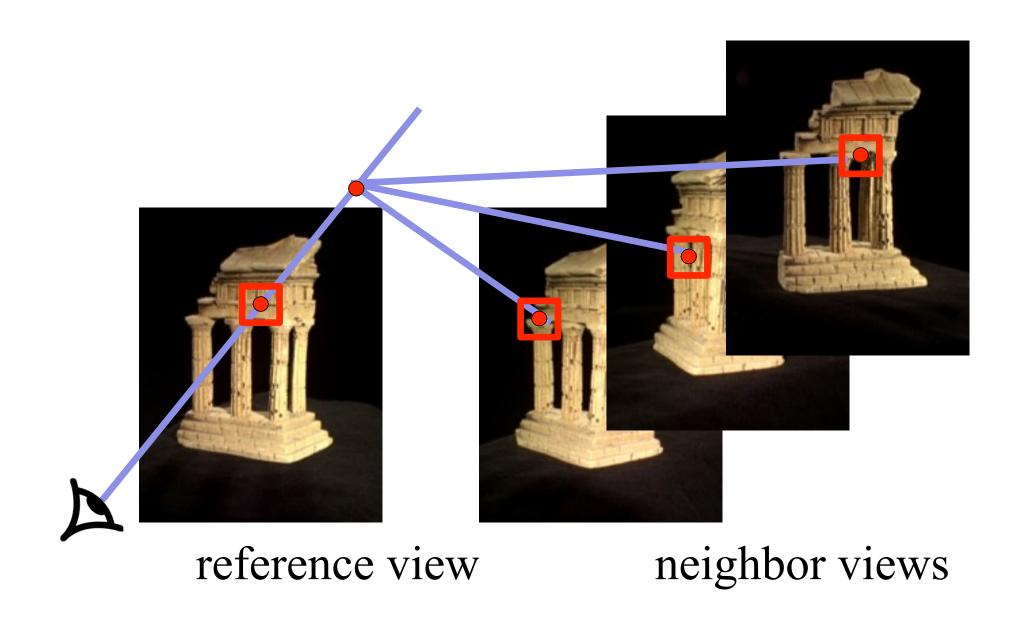
Can we estimate depth, now that we have pose?

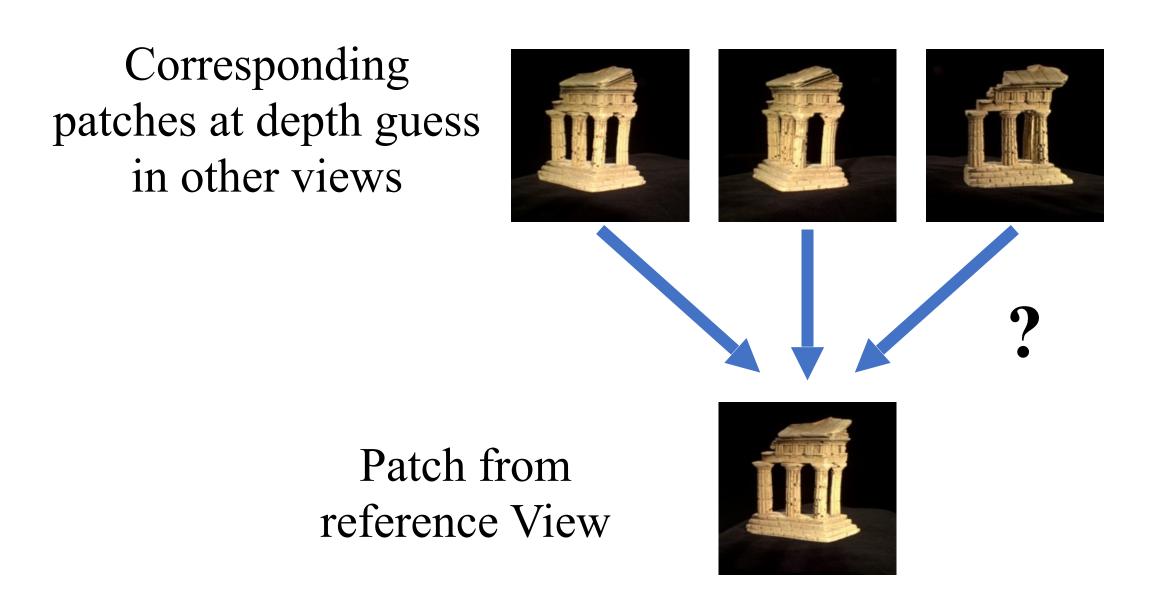
Source: N. Snavely

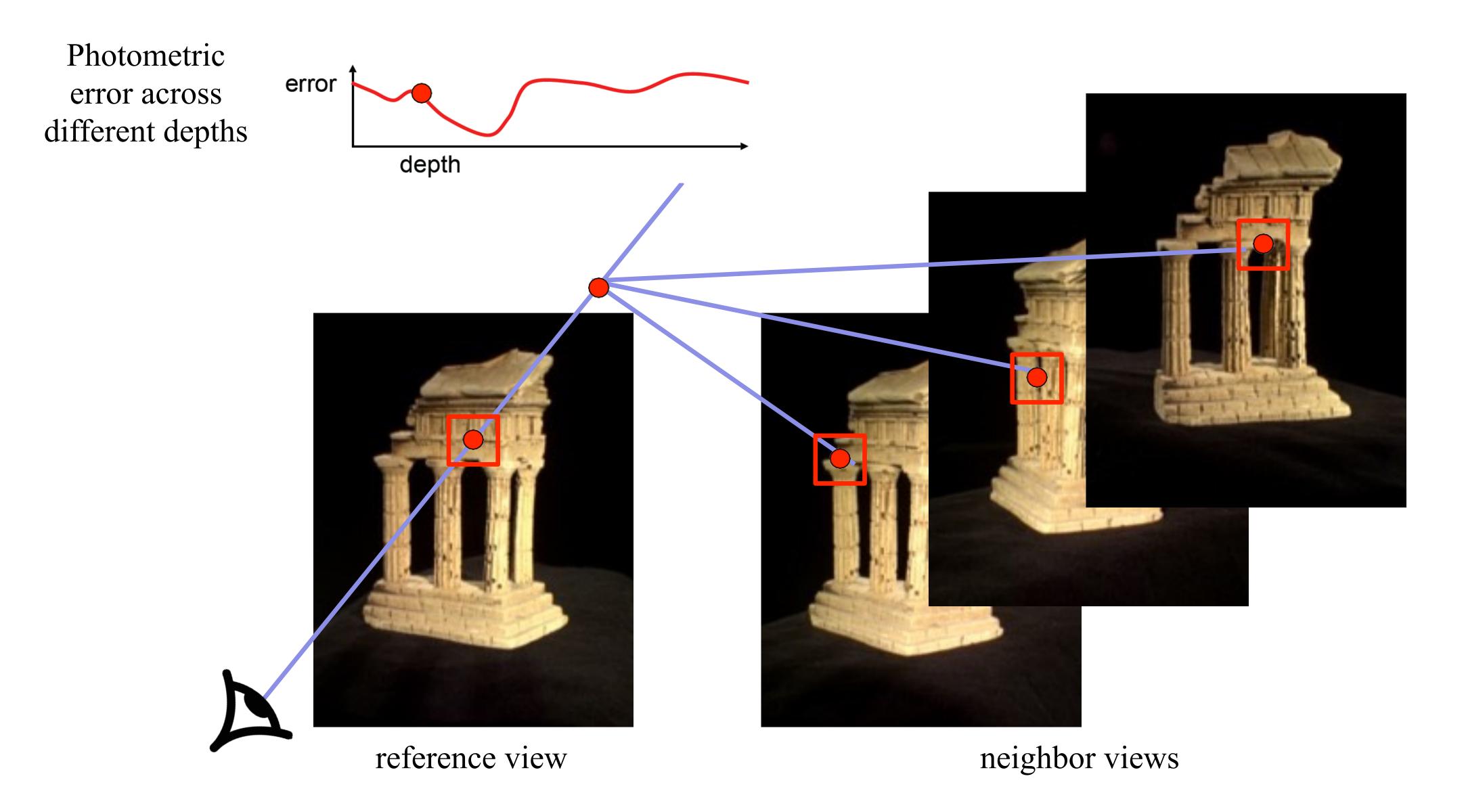


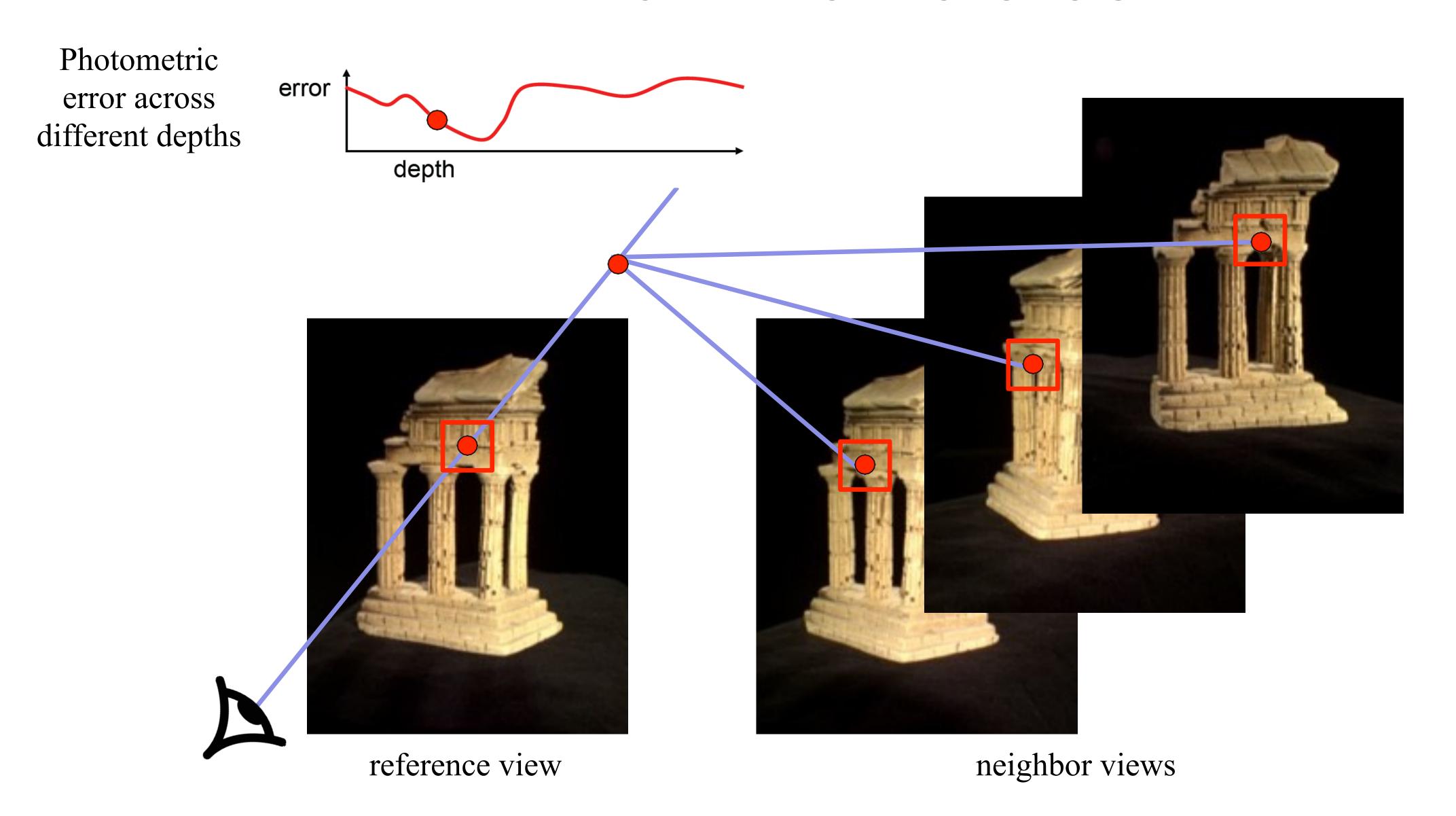
Source: Y. Furukawa

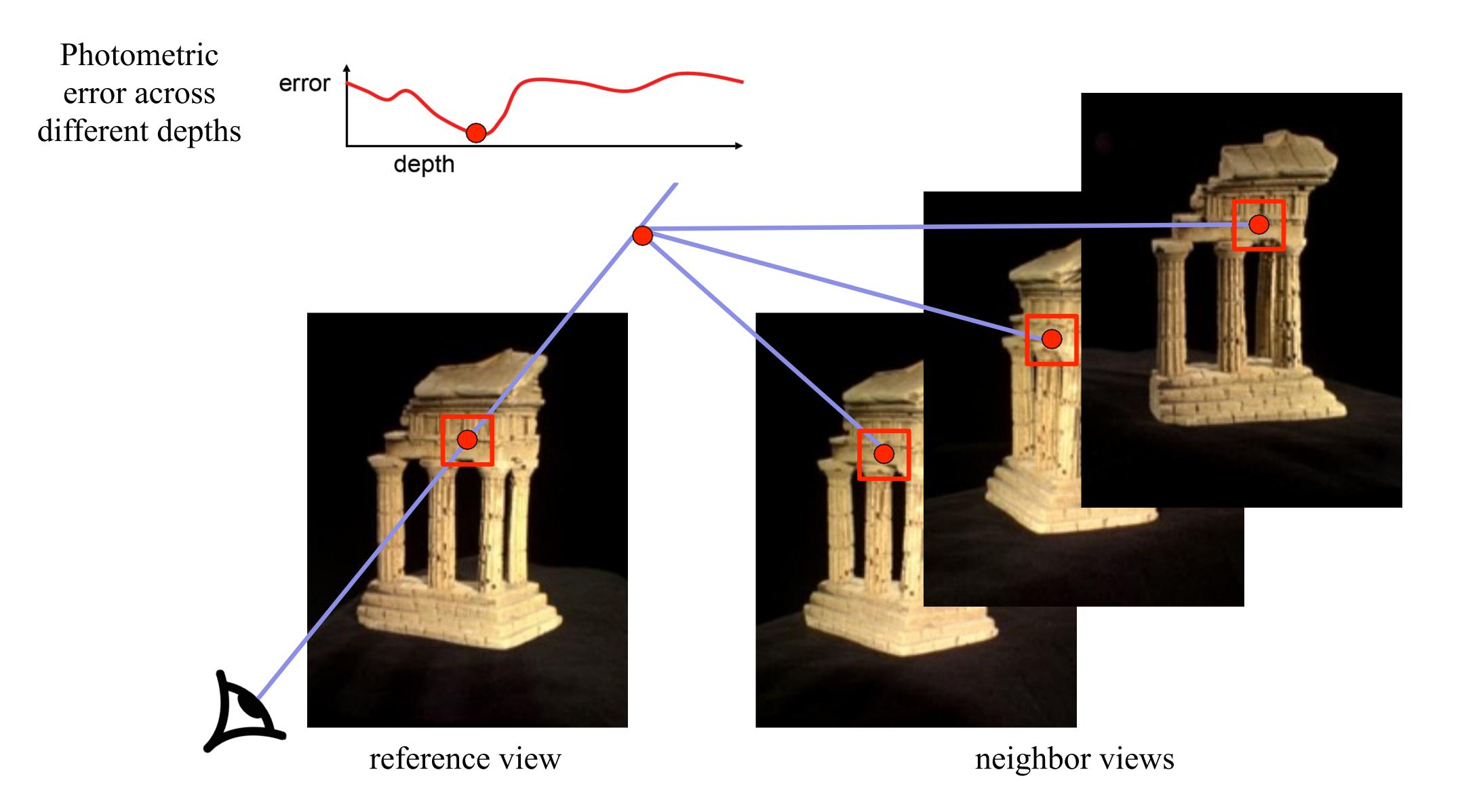
Evaluate the likelihood of a particular depth for a particular reference patch:



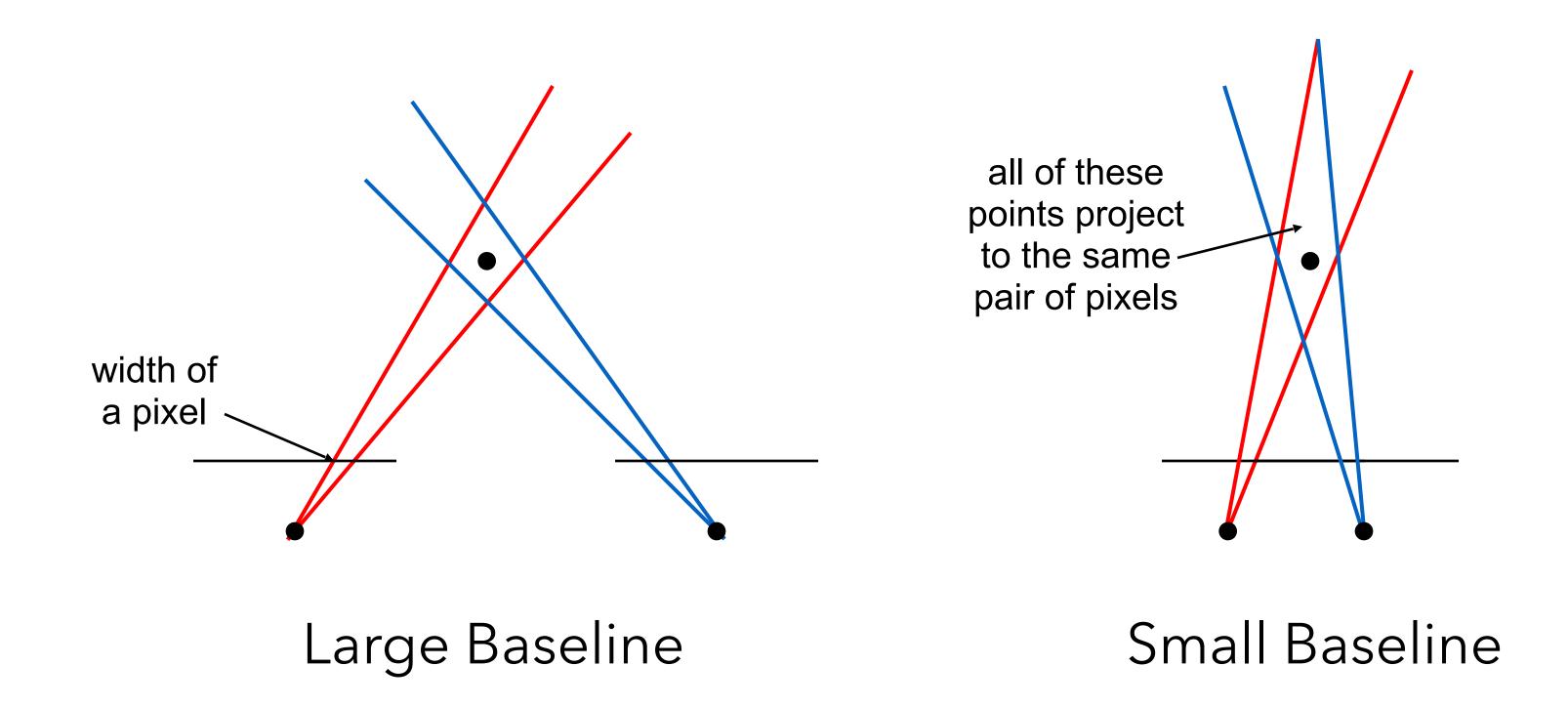








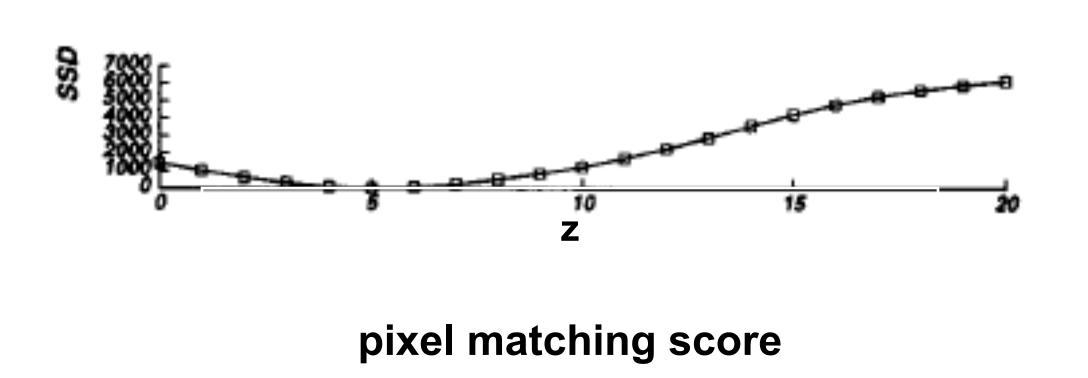
Multiple-baseline stereo



What's the optimal baseline?

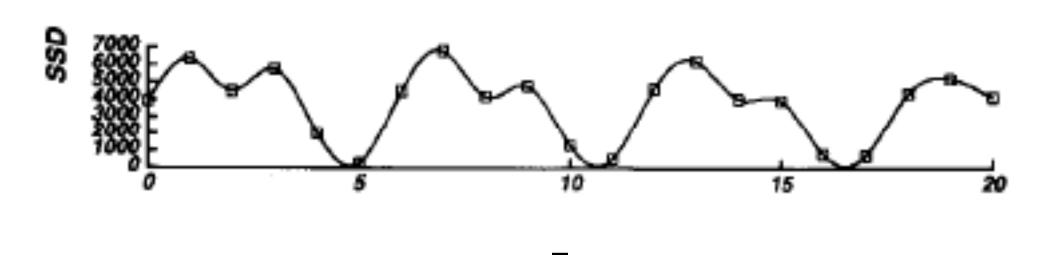
- Too small: large depth error
- Too large: difficult search problem

Multiple-baseline stereo



width of a pixel

 For short baselines, estimated depth will be less precise due to narrow triangulation



width of a pixel

For larger baselines, must search larger area in second image

Next class: neural fields