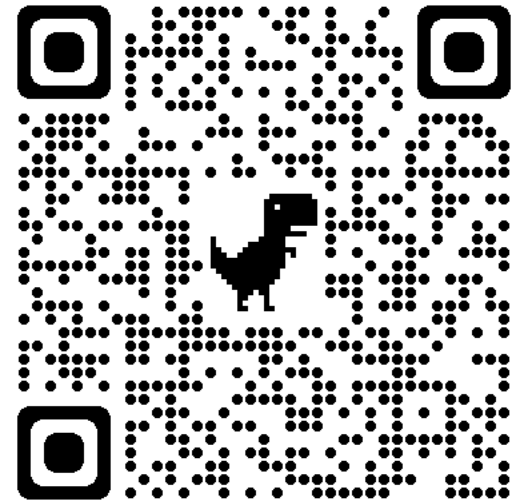


Lecture 16: Structure from Motion

COMP 590/776: Computer Vision

Instructor: Soumyadip (Roni) Sengupta

TA: Mykhailo (Misha) Shvets



Course Website:
Scan Me!

Recap

Geometry: How do we represent shape of an object?

2.5D representation:

1) Depth & Normal map

Easy to predict with 2D neural networks, efficient but do not give full 3D information.

Explicit representation:

2) Mesh

Hard for neural network but most Graphics pipeline use it. Very efficient with memory.

3) Voxels

Easy for neural network but high memory consumption

4) Point Cloud

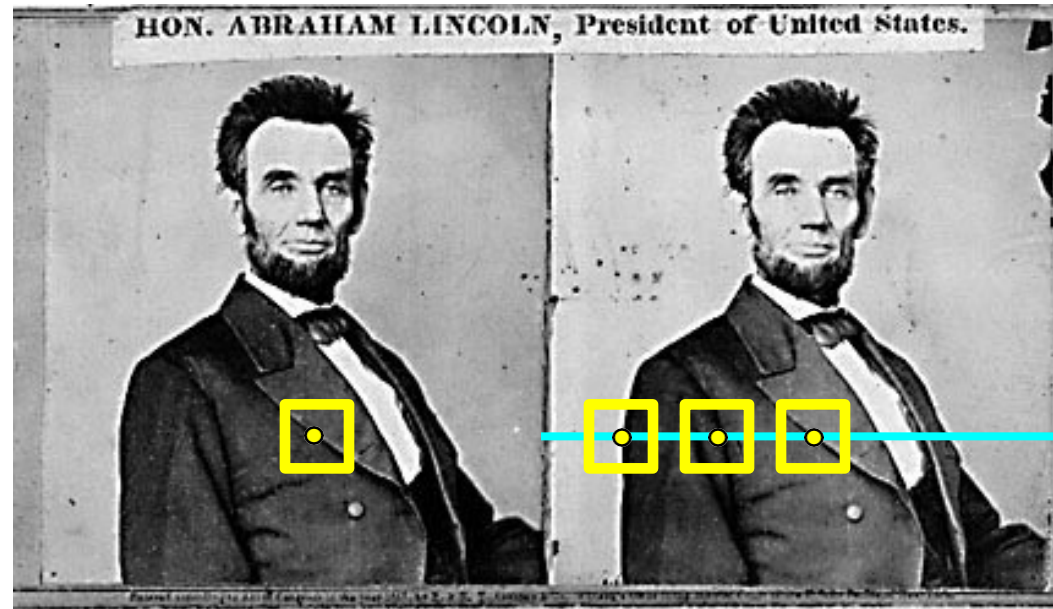
Output of many RGBD sensors or RGB algorithms

Implicit representation:

5) Surface Representation (SDF)

Memory efficient and deep networks can predict it. But need to convert it to mesh/voxel to be usable in Graphics engines.

Stereo

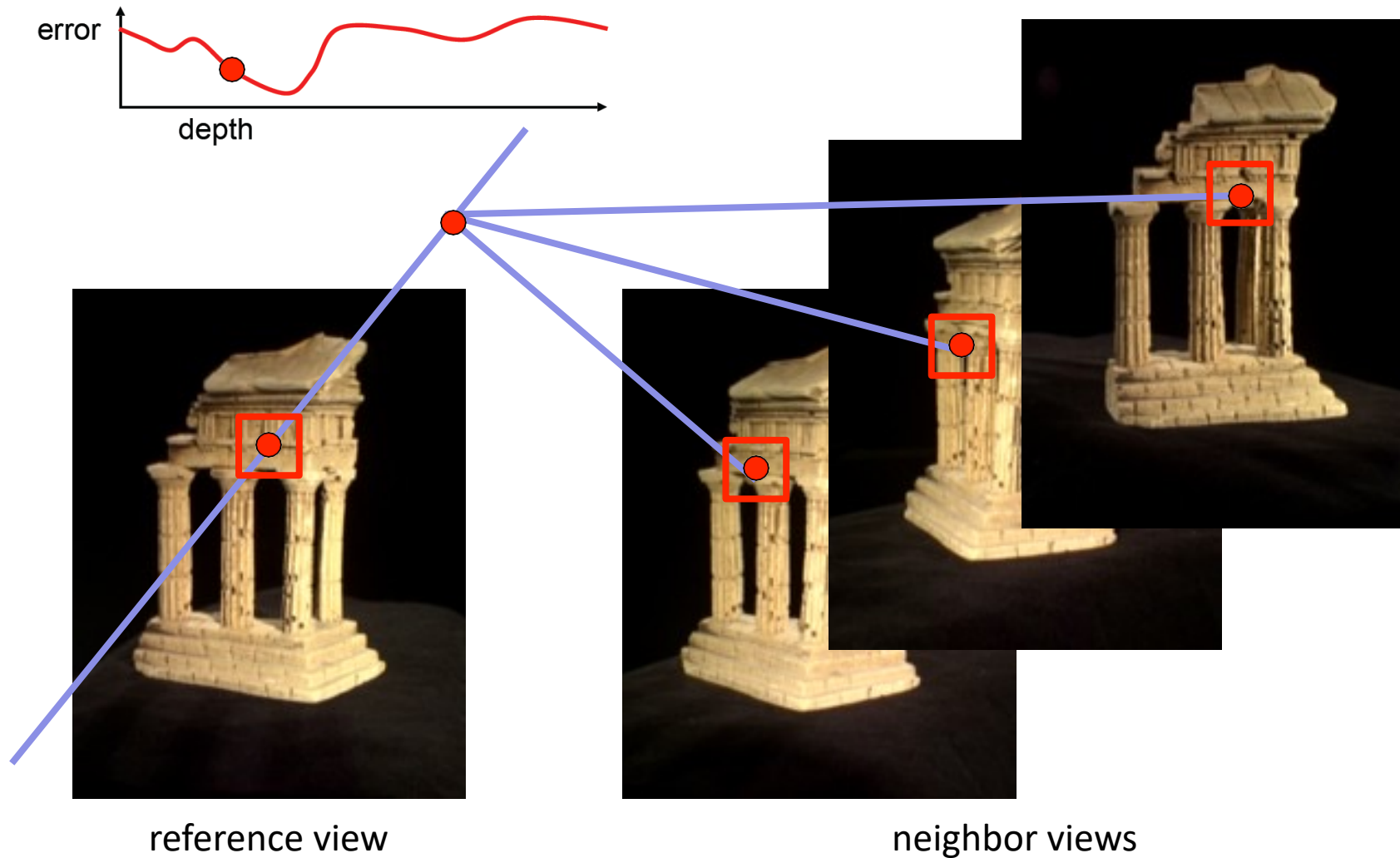


1. Rectify images
(make epipolar lines horizontal)
2. For each pixel
 - a. Find epipolar line
 - b. Scan line for best match
 - c. Compute depth from disparity

$$Z = \frac{bf}{d}$$

How can you make the epipolar lines horizontal?

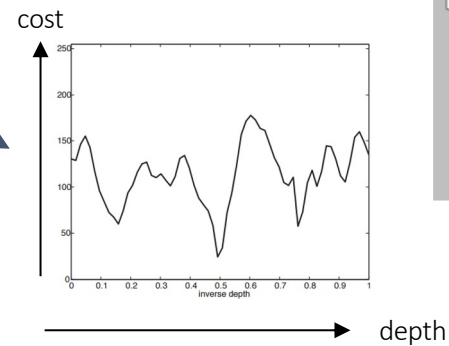
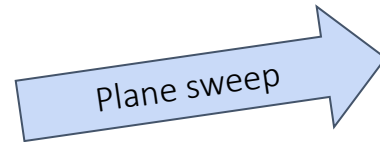
Multi-view stereo: Basic idea



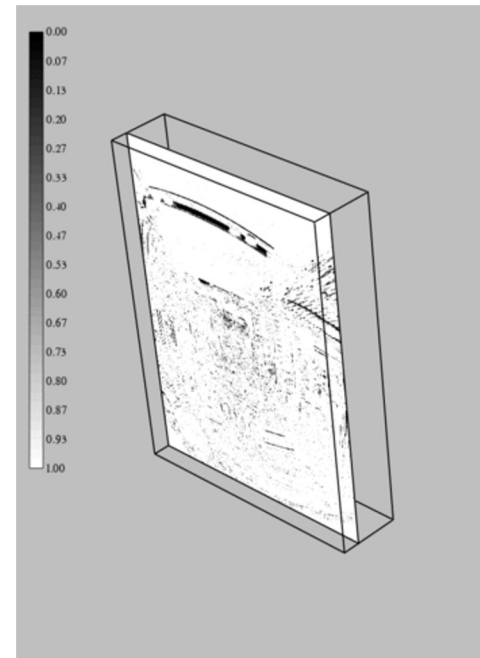
Plane Sweep Stereo: Cost Volumes -> Depth Maps



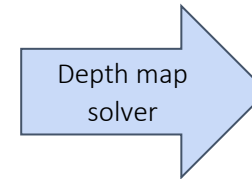
Reference image



Single pixel's cost profile



Full cost volume



(Belief propagation, graph cuts, etc.)



Plane-Sweep Stereo

- The family of depth planes in the coordinate frame of the reference view

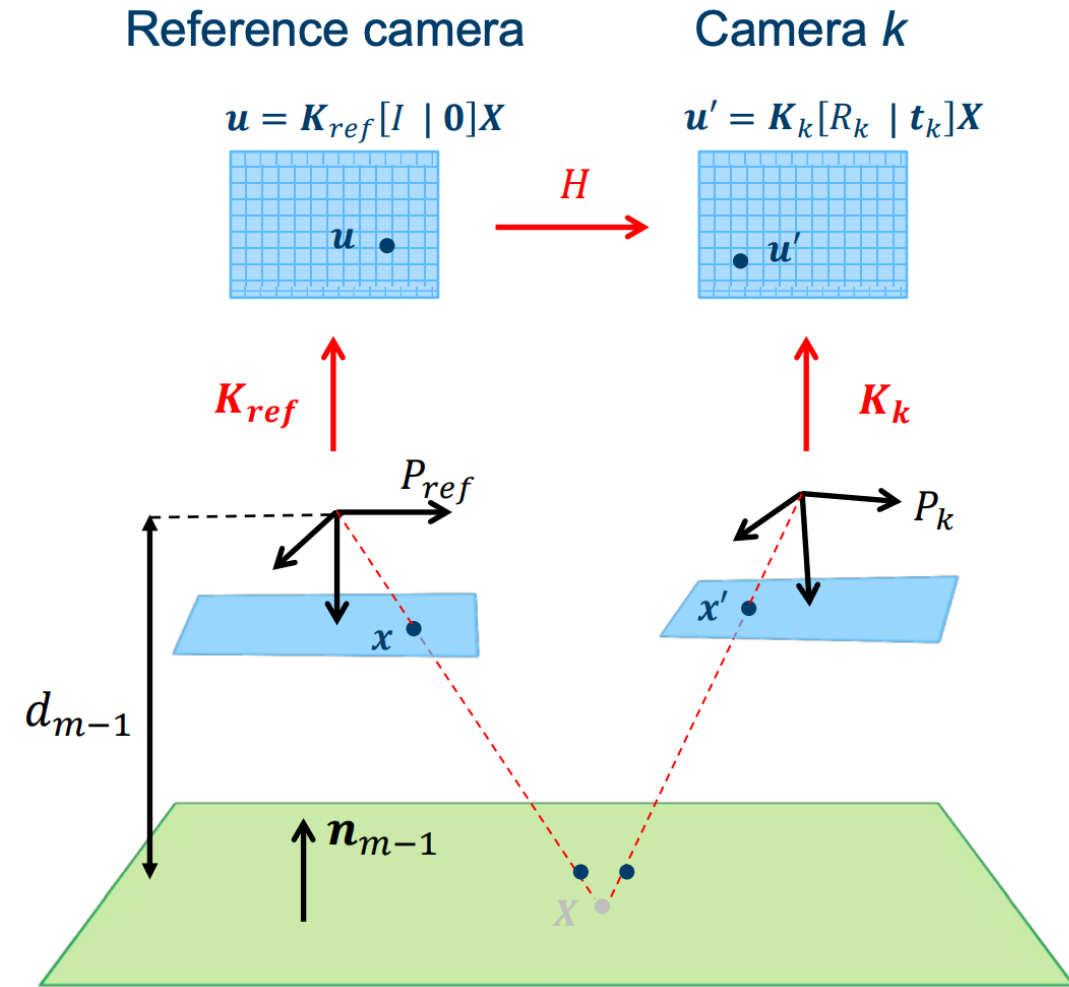
$$\Pi_m = \begin{bmatrix} \mathbf{n}_m^T & -d_m \end{bmatrix}$$

- The mapping from the reference camera P_{ref} onto the plane Π_m and back to camera P_k is described by the homography induced by the plane Π_m

$$H_{\Pi_m, P_k} = K_k \left(R_k - \mathbf{t}_k \mathbf{n}_m^T / d_m \right) K_{ref}^{-1}$$

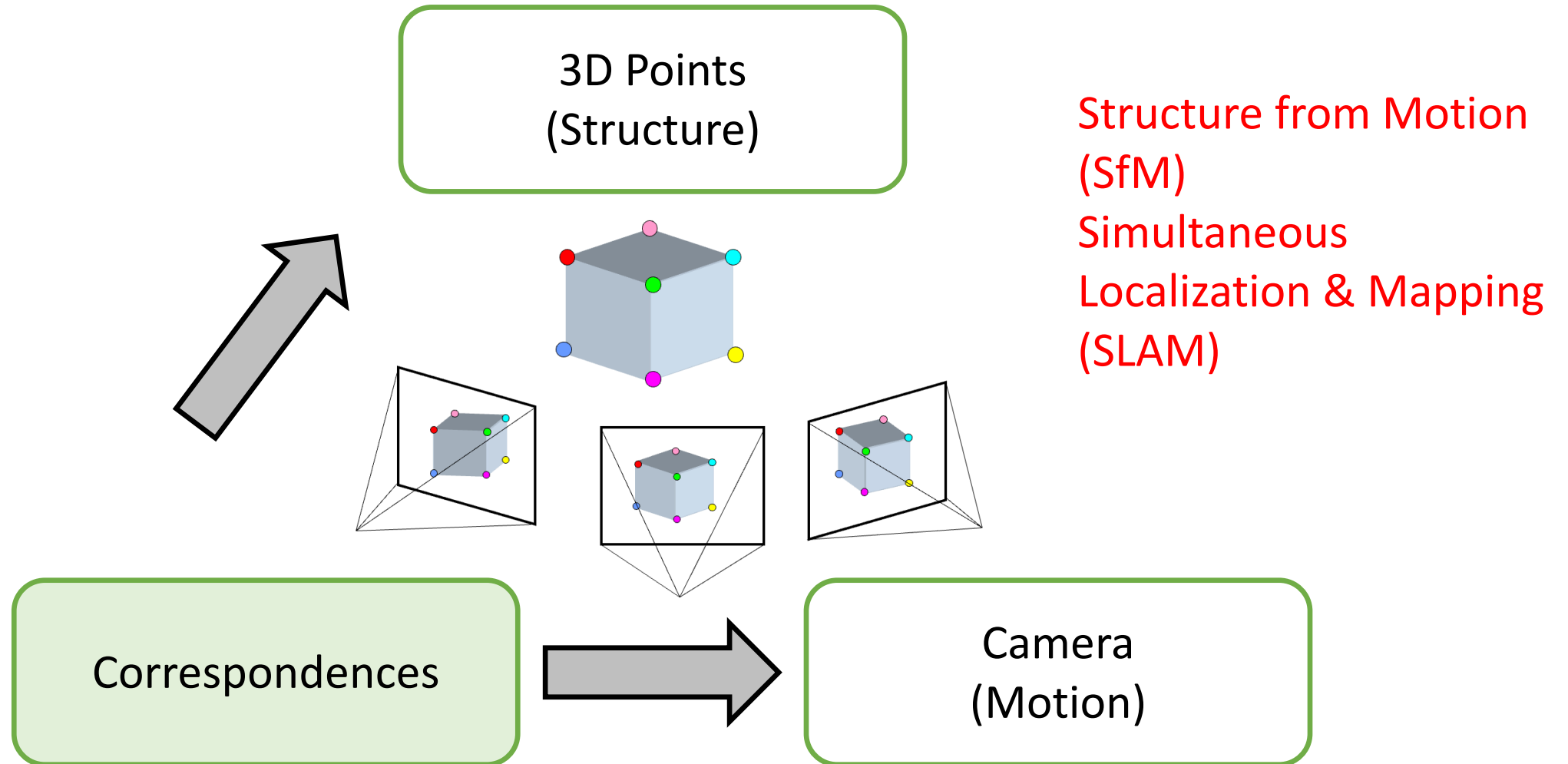
Try the proof in HW!

- The mapping from P_k to P_{ref} induced by Π_m is the inverse homography H_{Π_m, P_k}^{-1}



Slight abuse of notation. In equation (x,y) are image co-ordinates, in figure u is image co-ordinate.

Big picture: 3 key components in 3D



Structure from motion

- SfM solves both of these problems *at once*
- A kind of chicken-and-egg problem
 - (but solvable)

Structure from Motion (SfM)

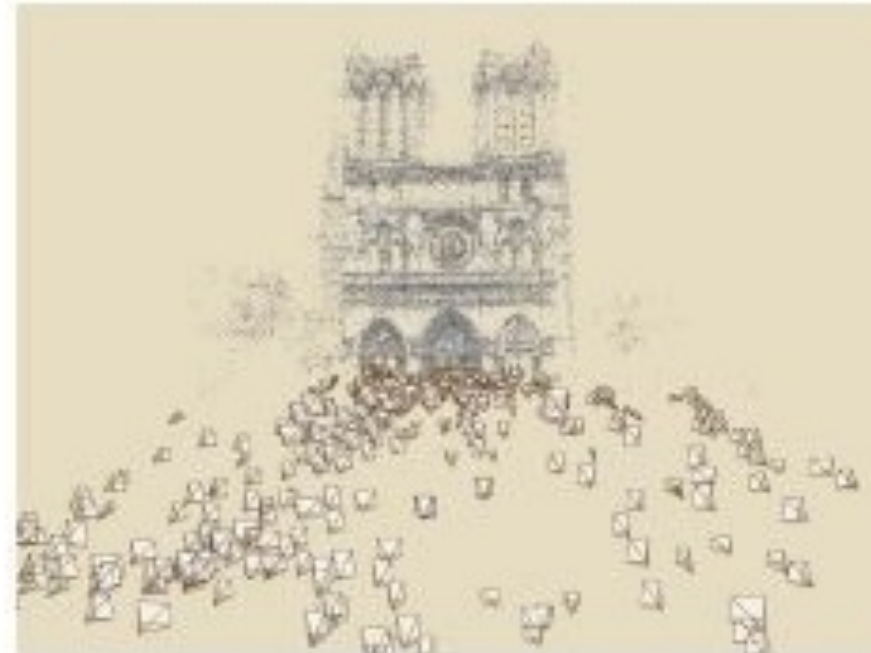
- Given many images, how can we
 - a) figure out where they were all taken from?
 - b) build a 3D model of the scene?



This is (roughly) the **structure from motion** problem

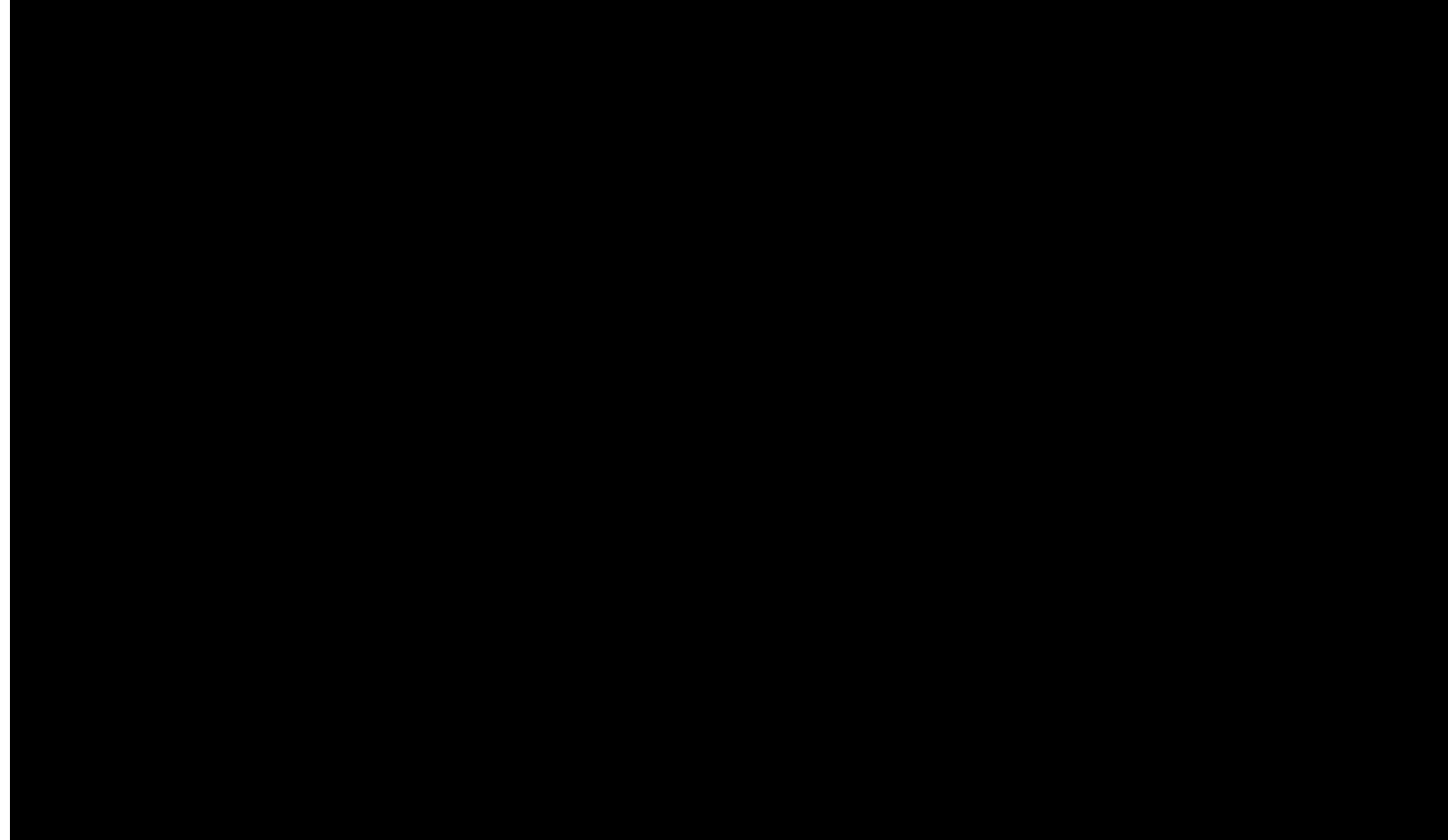
Photo Tourism

Noah Snavely, Steven M. Seitz, Richard Szeliski, "[Photo tourism: Exploring photo collections in 3D](#)," SIGGRAPH 2006



<https://youtu.be/mTBPGuPLI5Y>

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

Large-scale structure from motion



Rome's Colosseum

Reconstructing the World in Six Days,

Jared Heinly, Johannes L. Schönberger, Enrique Dunn, Jan-Michael Frahm, CVPR 2015.

Work done at UNC CS!



St. Peter's Basilica, Vatican City

Yahoo Flickr Creative Commons 100M Dataset

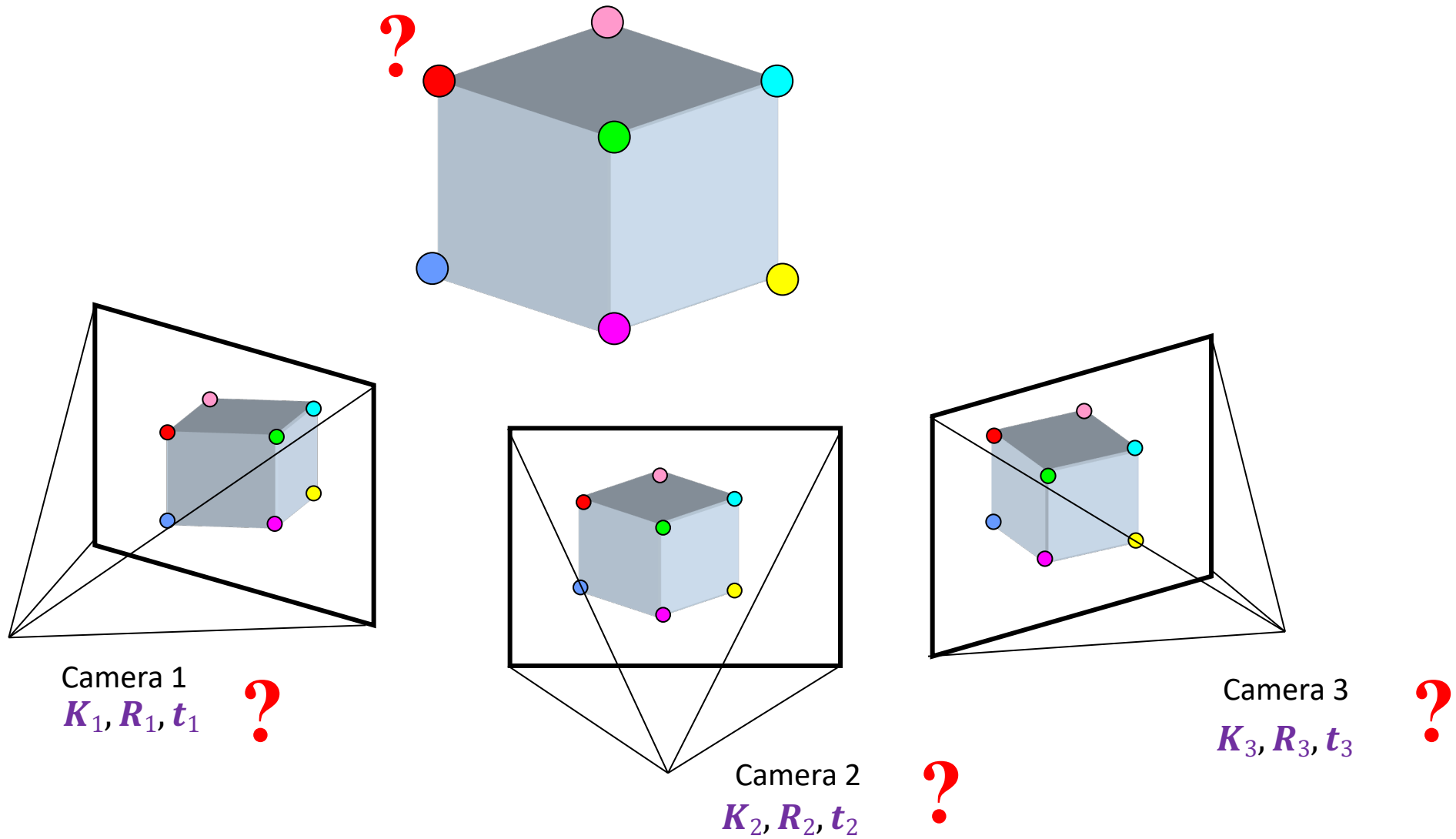
Today's Class

- Ambiguities in SfM
- Affine SfM
- Projective SfM
 - Global SfM
 - Incremental SfM
- Challenges and Applications

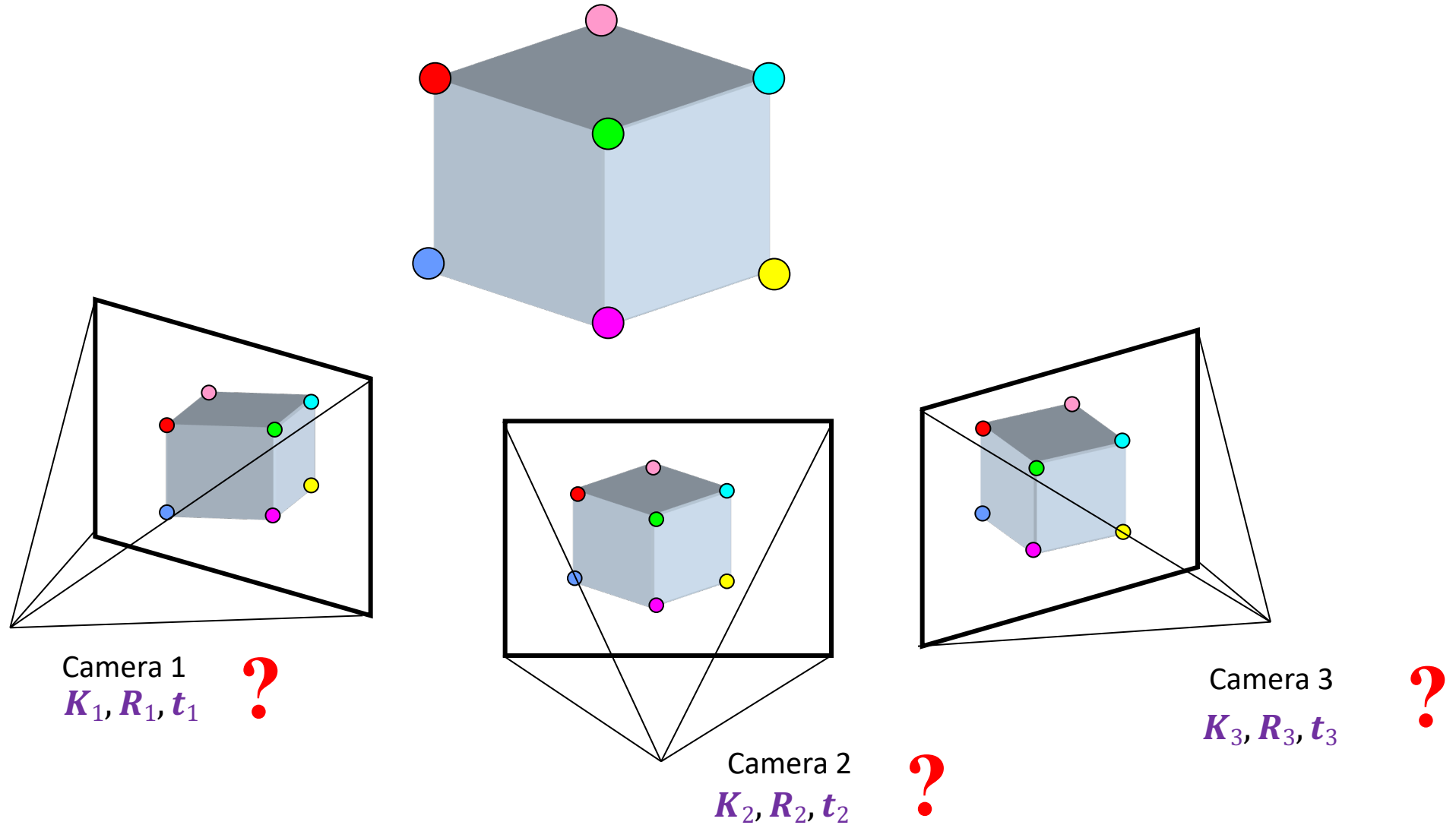
Today's Class

- Ambiguities in SfM
- Affine SfM
- Projective SfM
 - Global SfM
 - Incremental SfM
- Challenges and Applications

Structure from motion

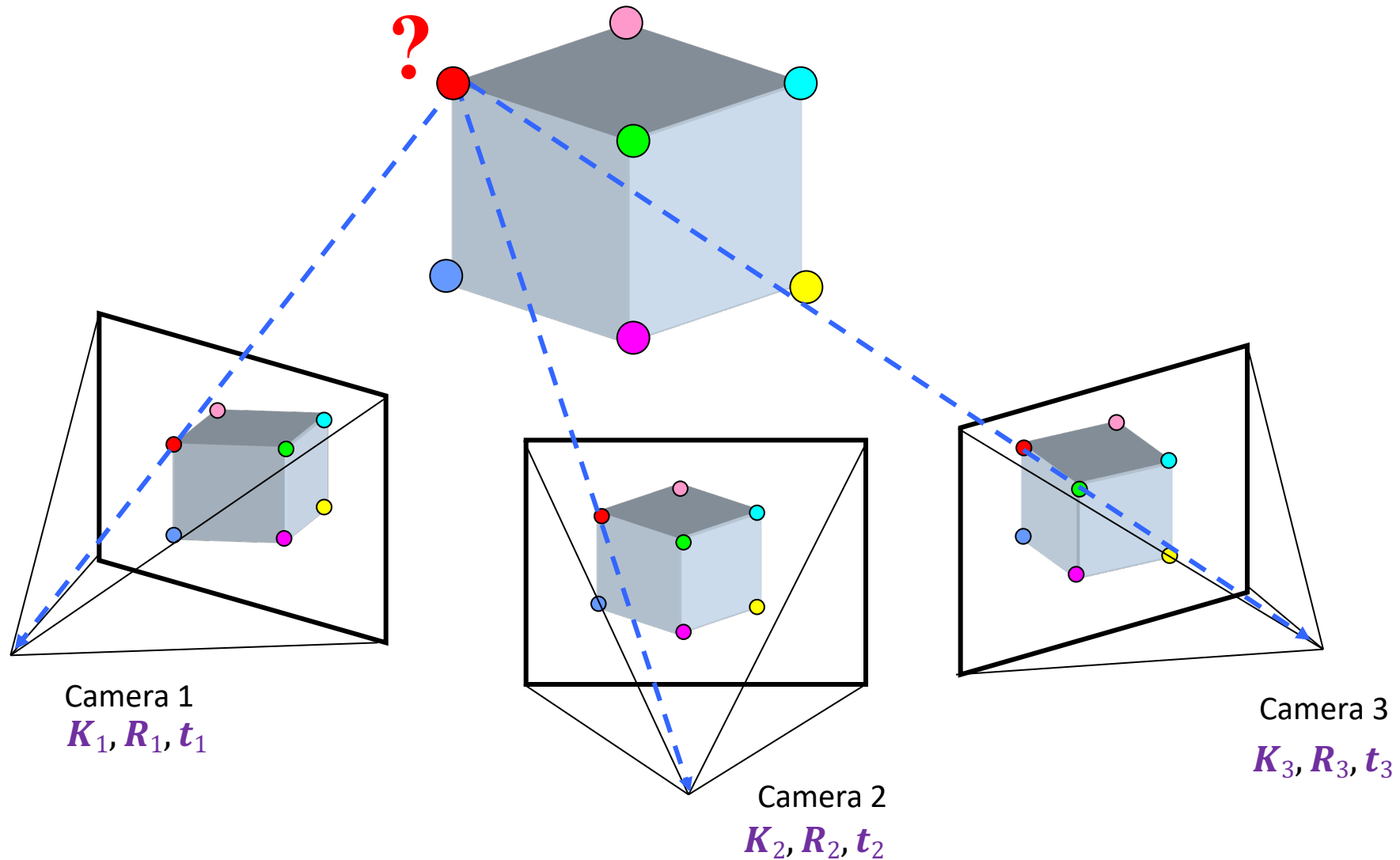


Recall: Calibration



- Given a set of *known* 3D points seen by a camera, compute the camera parameters

Recall: Triangulation



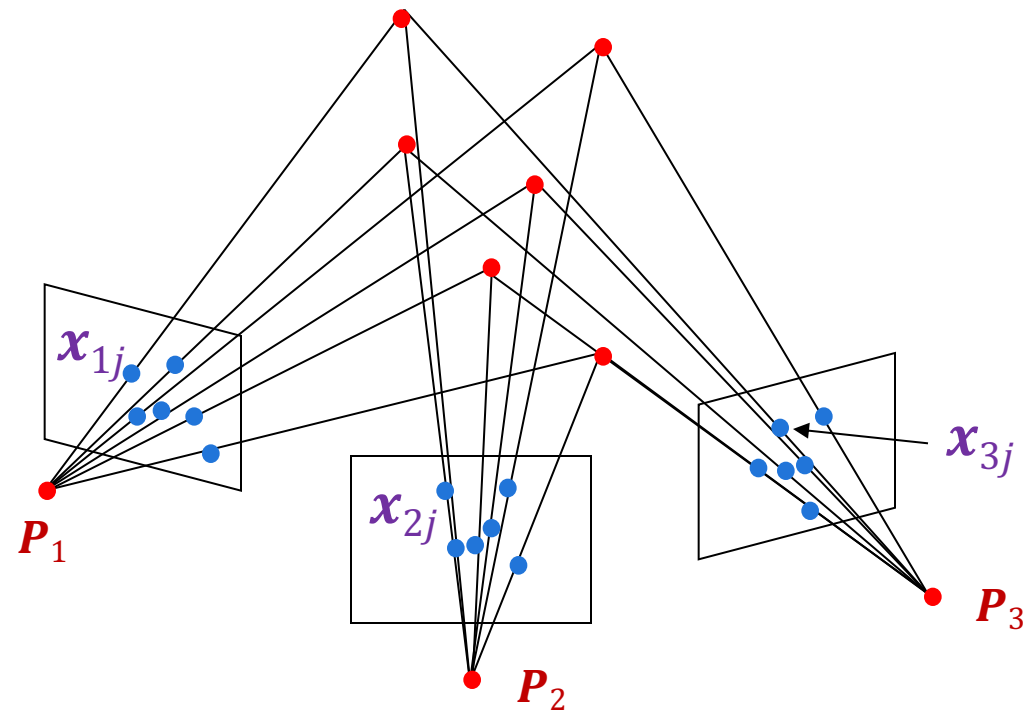
- Given *known cameras* and projections of the same 3D point in two or more images, compute the 3D coordinates of that point

Structure from motion: Problem formulation

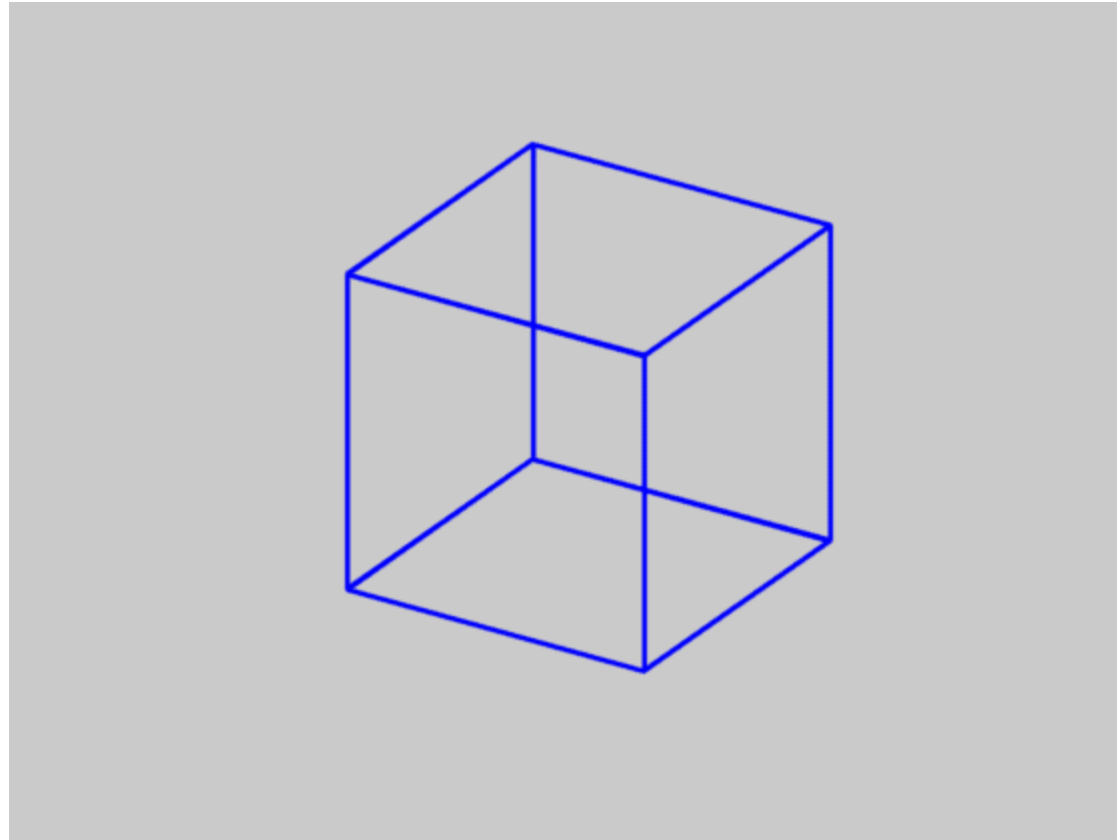
- Given: m images of n fixed 3D points such that (ignoring visibility)

$$\bullet \mathbf{x}_{ij} \cong \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- Problem: estimate m projection matrices \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}



Is SFM always uniquely solvable?



- Necker cube

Structure from motion ambiguity

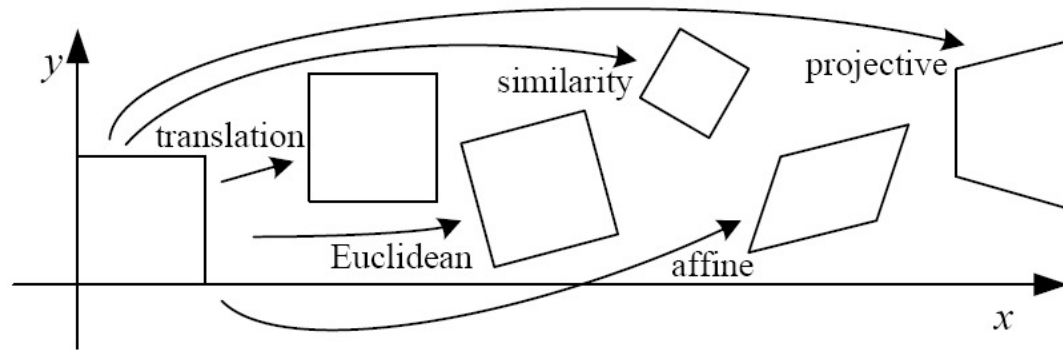
- If we scale the entire scene by some factor k and, at the same time, scale the camera matrices by the factor of $1/k$, the projections of the scene points remain exactly the same:

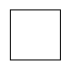
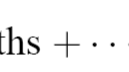



$$\bullet \mathbf{x} \cong \mathbf{P}\mathbf{X} = \left(\frac{1}{k}\mathbf{P}\right)(k\mathbf{X})$$

- Without a reference measurement, it is impossible to recover the absolute scale of the scene!
- In general, if we transform the scene using a transformation Q and apply the inverse transformation to the camera matrices, then the image observations do not change:

$$\bullet \mathbf{x} \cong \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}^{-1})(\mathbf{Q}\mathbf{X})$$

Recall: 2D image transformations



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

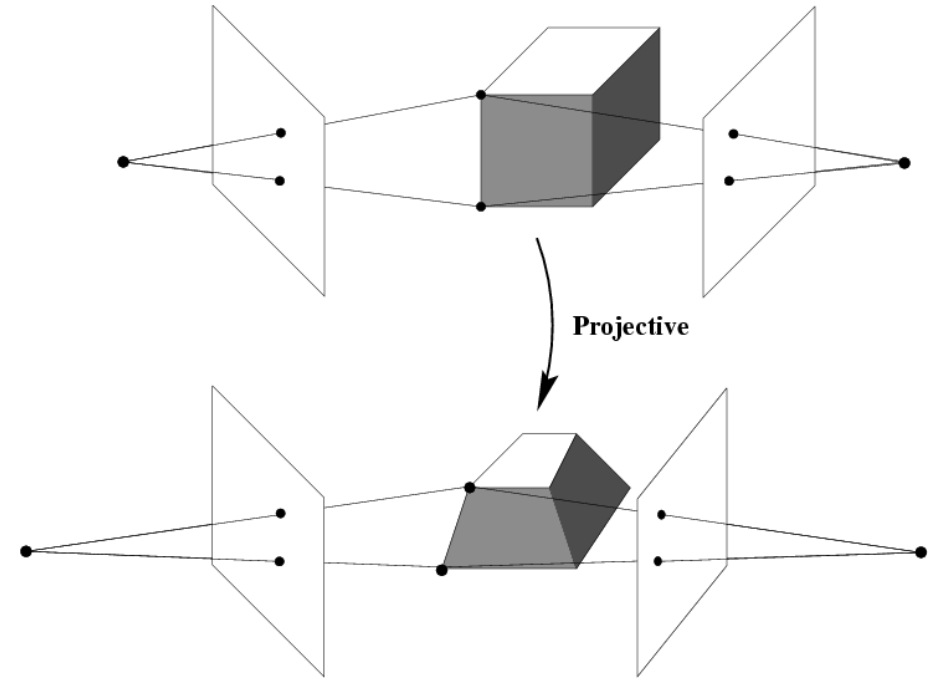
Now, lets extend this to 3D.

Projective ambiguity

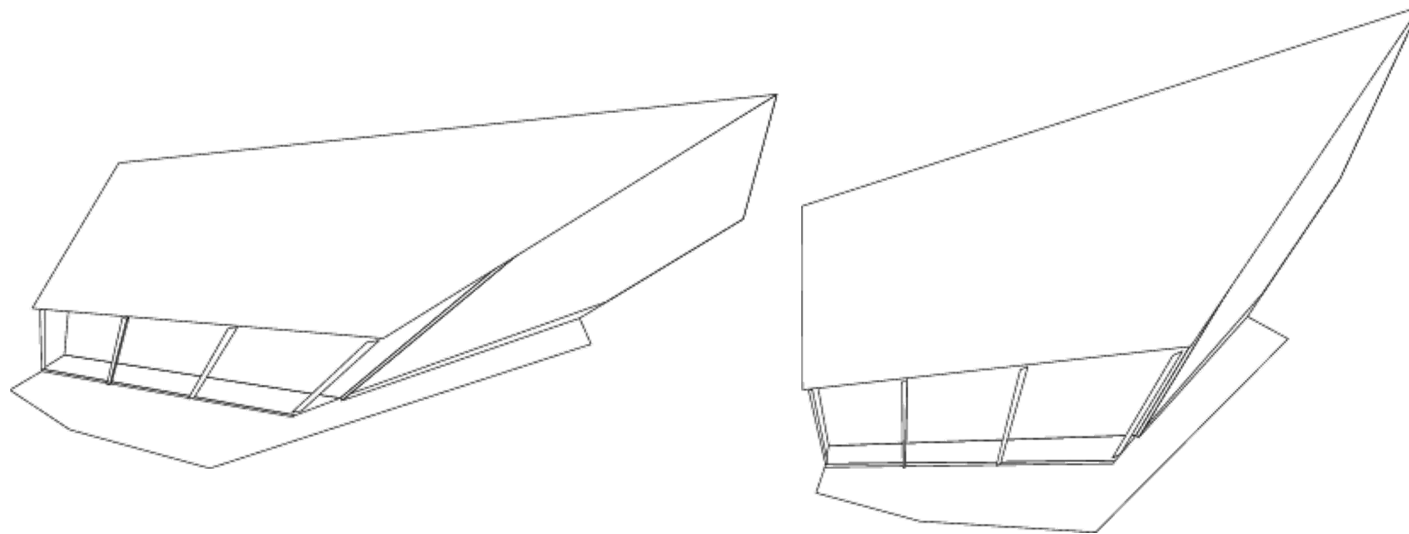
- With no constraints on the camera calibration matrices or on the scene, we can reconstruct up to a *projective* ambiguity:

$$\mathbf{x} \cong \mathbf{P}\mathbf{X} = (\mathbf{P}\mathbf{Q}^{-1})(\mathbf{Q}\mathbf{X})$$

\mathbf{Q} is a general full-rank 4×4 matrix



Projective ambiguity



Affine ambiguity

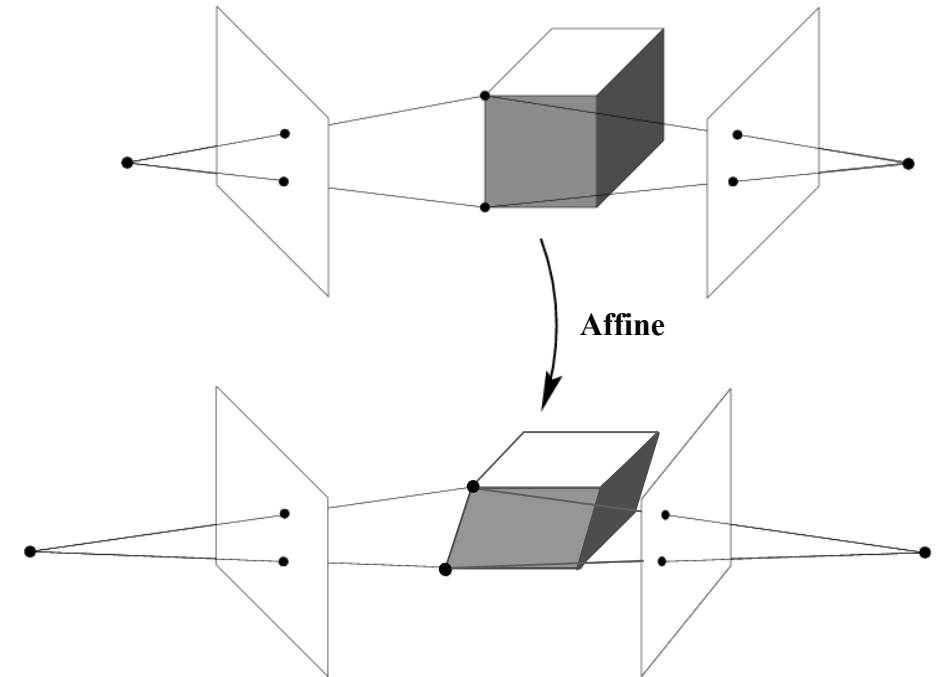
- If we impose parallelism constraints, we can get a reconstruction up to an *affine* ambiguity:

$$x \cong PX = (PQ_A^{-1})(Q_A X)$$

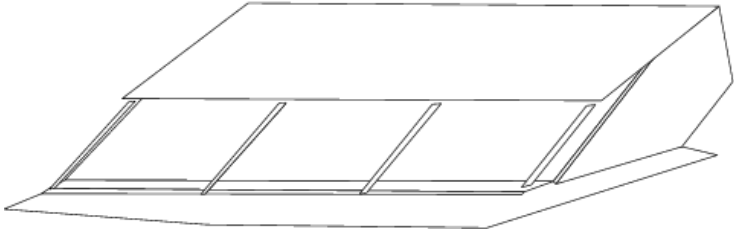
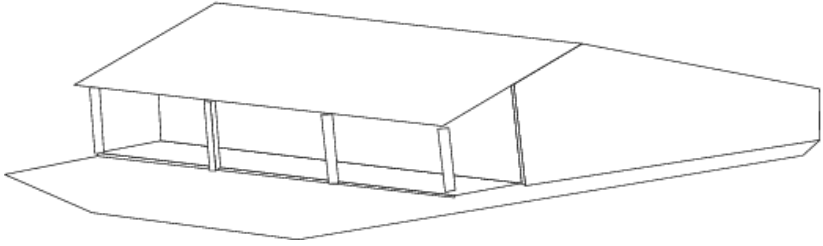
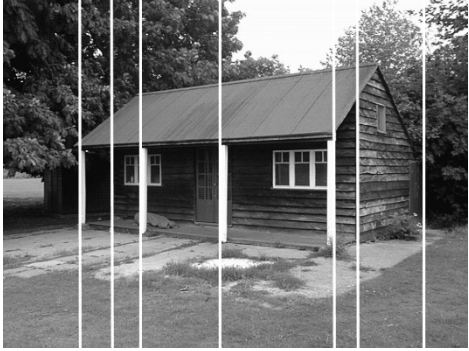
3×3 full-rank matrix

3×1 translation vector

$$Q_A = \begin{bmatrix} A & t \\ \mathbf{0}^T & 1 \end{bmatrix}$$



Affine ambiguity



Similarity ambiguity

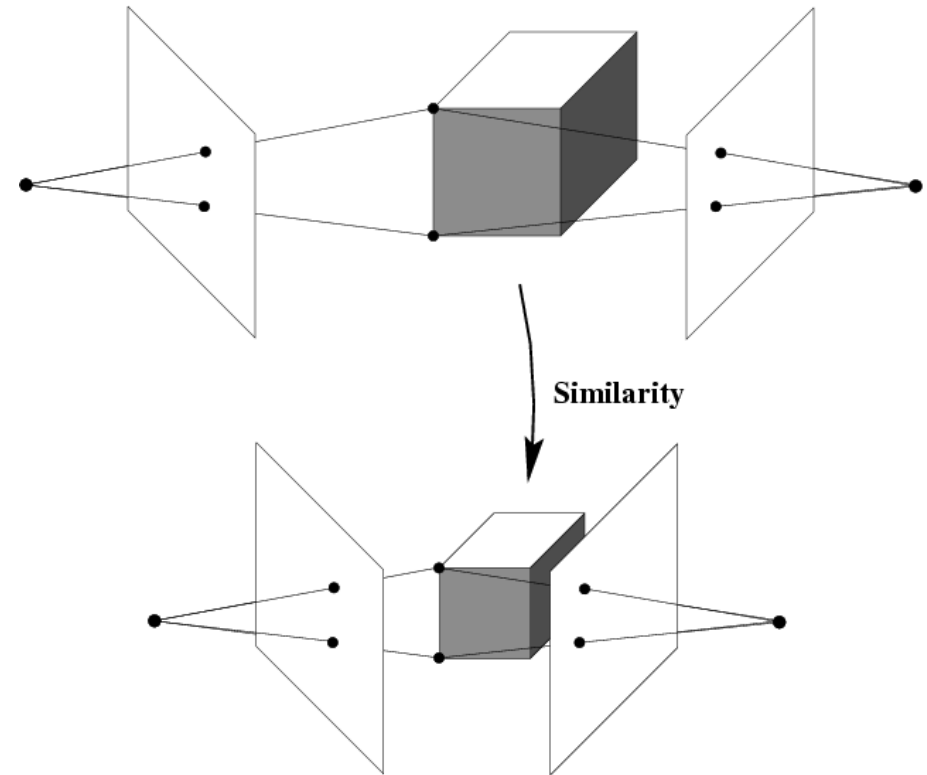
- A reconstruction that obeys orthogonality constraints on camera parameters and/or scene

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

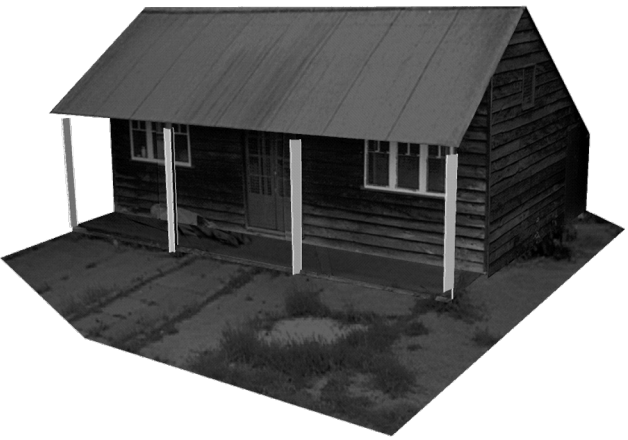
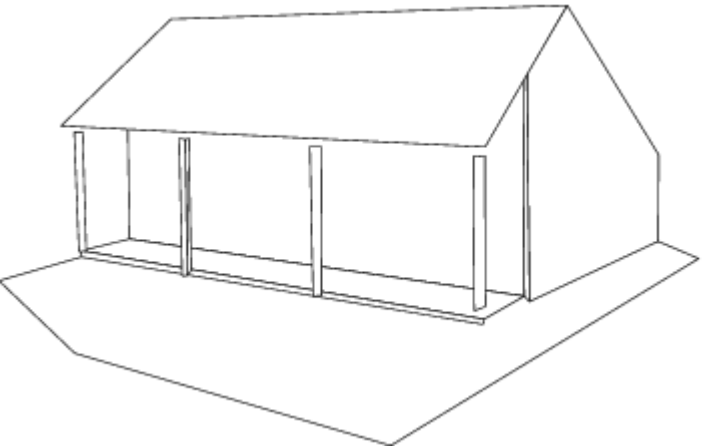
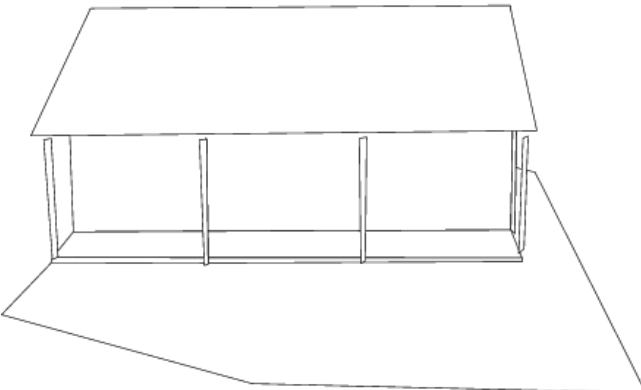
3×3
rotation
matrix

3×1 translation
vector

$$Q_S = \begin{bmatrix} sR & t \\ \mathbf{0}^T & 1 \end{bmatrix}$$



Similarity ambiguity



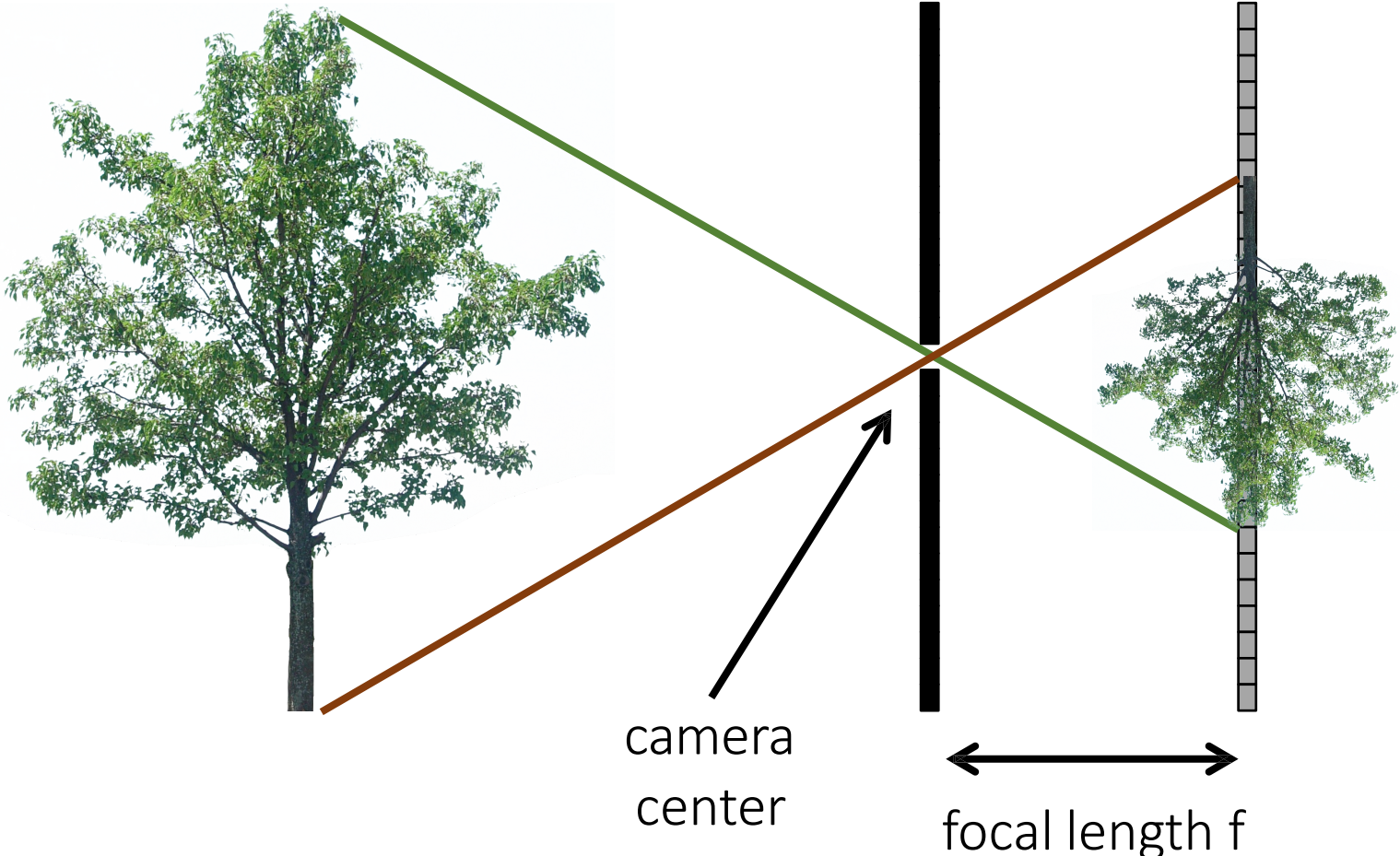
Today's Class

- Ambiguities in SfM
- **Affine SfM**
- Projective SfM
 - Global SfM
 - Incremental SfM
- Challenges and Applications

The pinhole camera

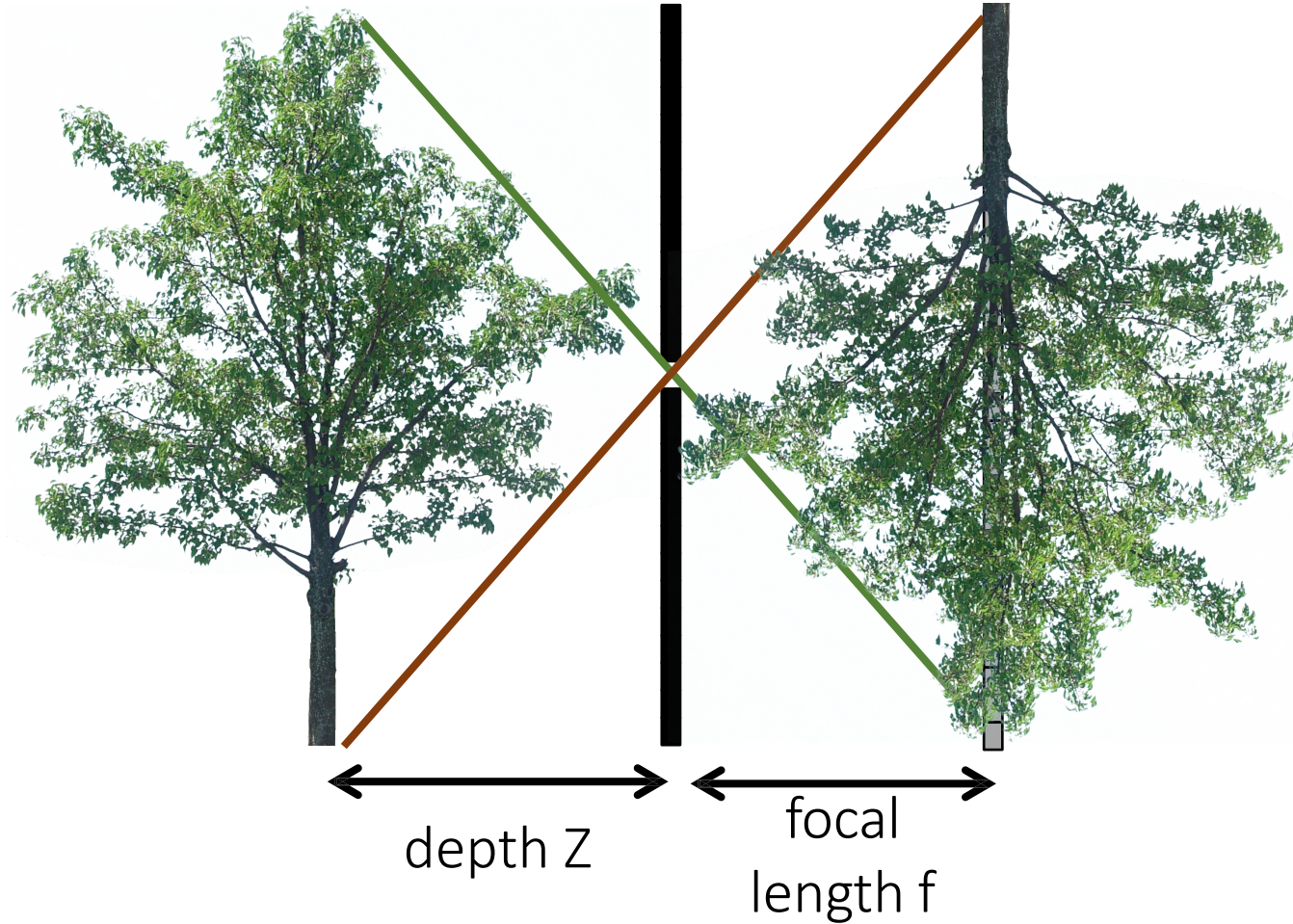
$$(x, y, z) \rightarrow \left(f \frac{x}{z}, f \frac{y}{z}\right)$$

real-world
object



What if...

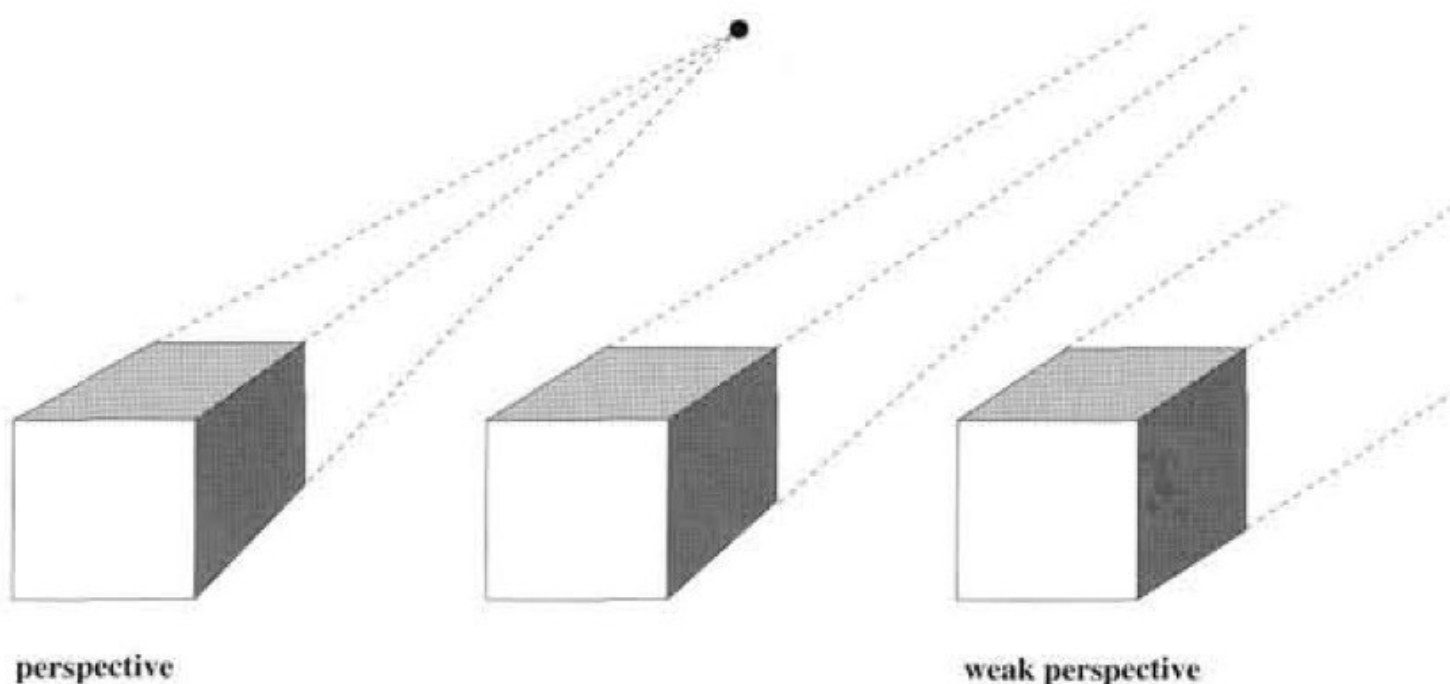
real-world
object



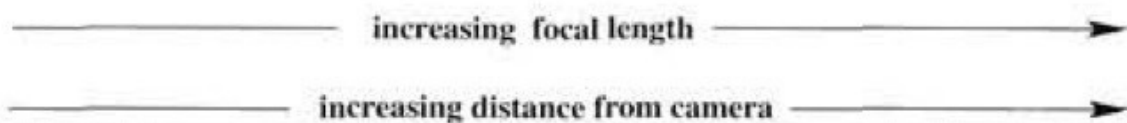
... we continue increasing Z
and f while maintaining
same magnification?

$$f \rightarrow \infty \text{ and } \frac{f}{Z} = \text{constant}$$

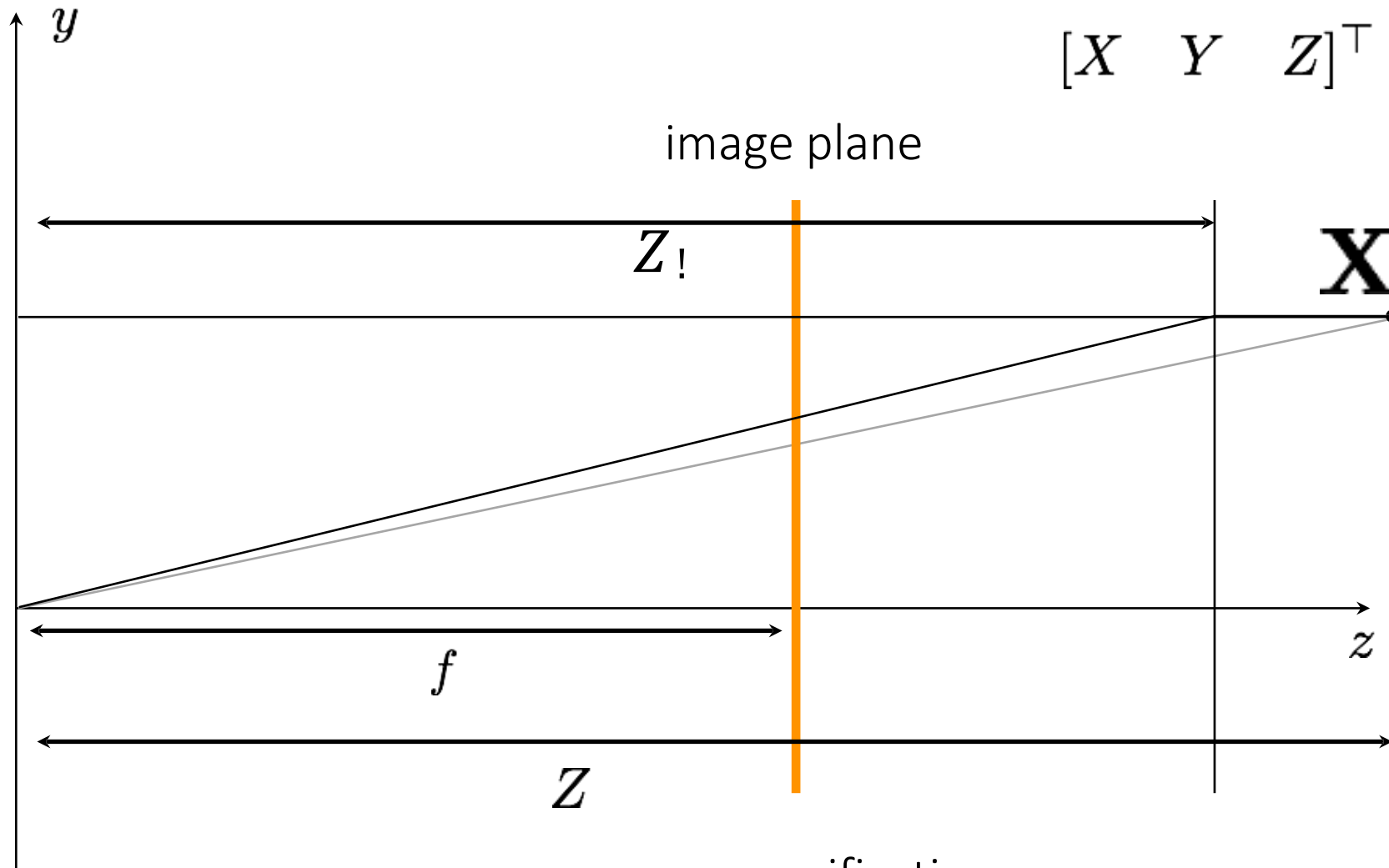
camera is *close*
to object and has
small focal length



camera is *far* from
object and has
large focal length



Weak perspective vs perspective camera



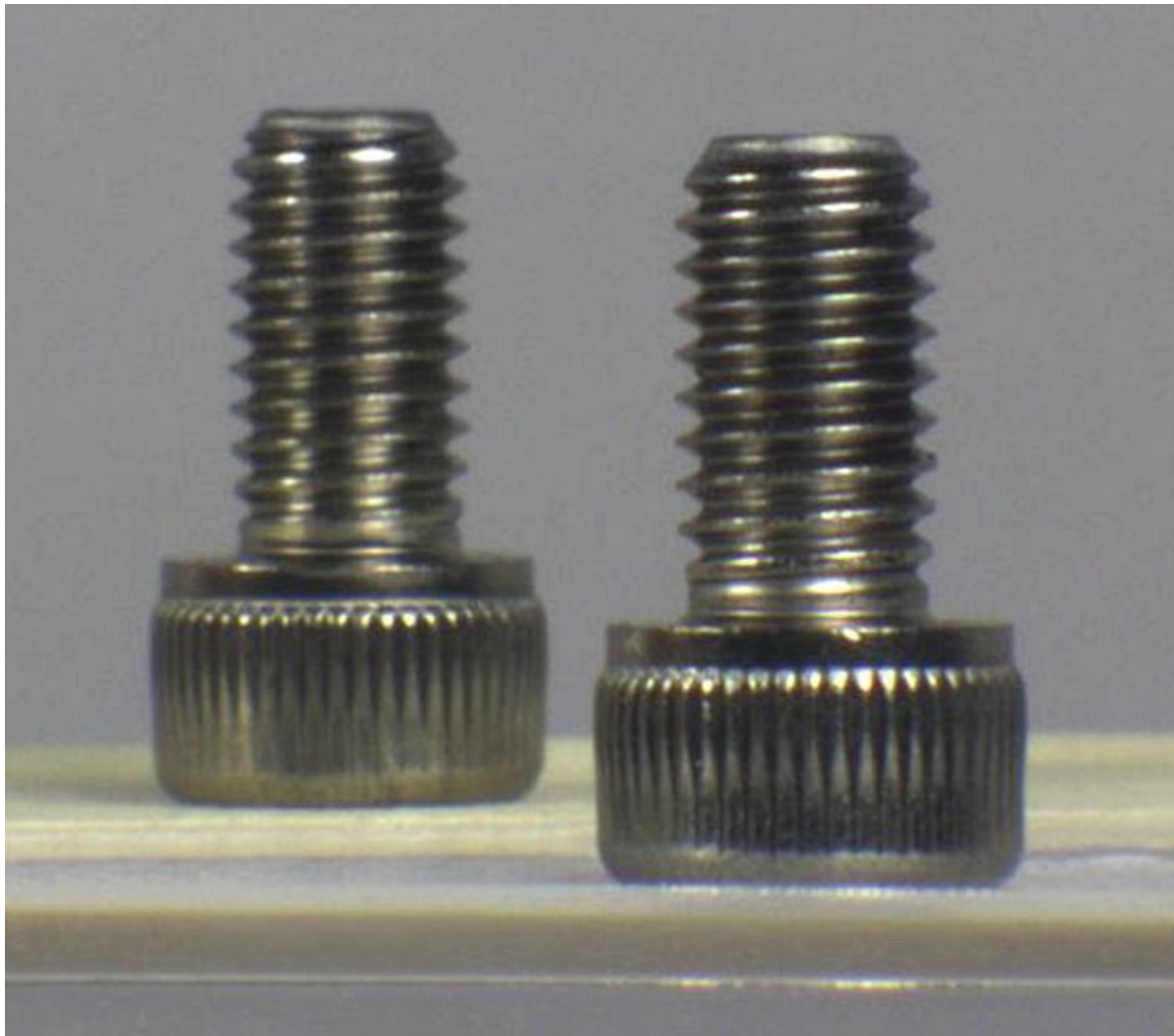
$$[X \ Y \ Z]^T \mapsto [fX/Z_0 \ fY/Z_0]^T$$

- magnification does not change with depth
- *constant* magnification depending on f and Z_0

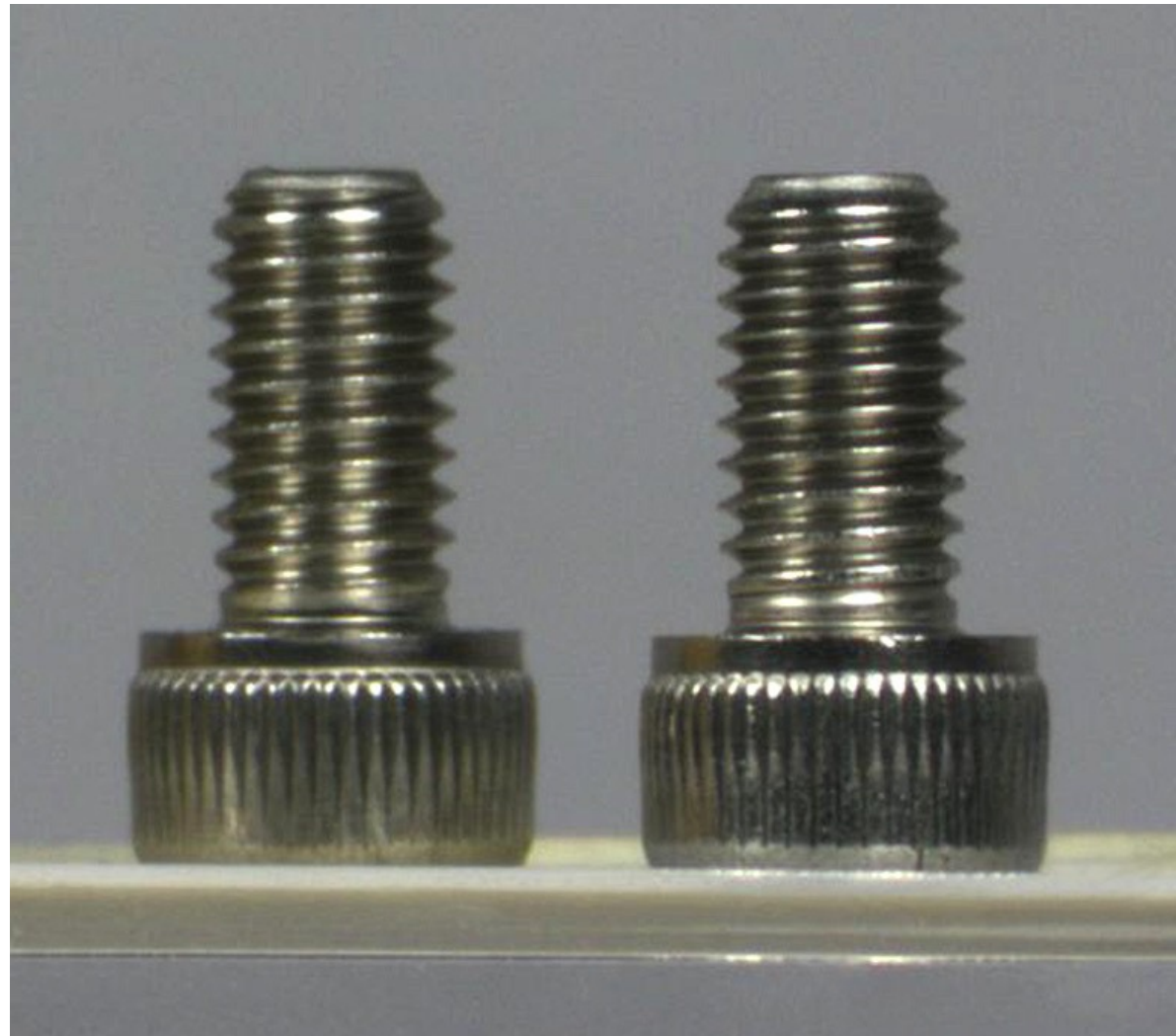
magnification
changes with depth

$$[X \ Y \ Z]^T \mapsto [fX/Z \ fY/Z]^T$$

Different cameras



perspective camera



weak perspective camera

When can we assume a weak perspective camera?

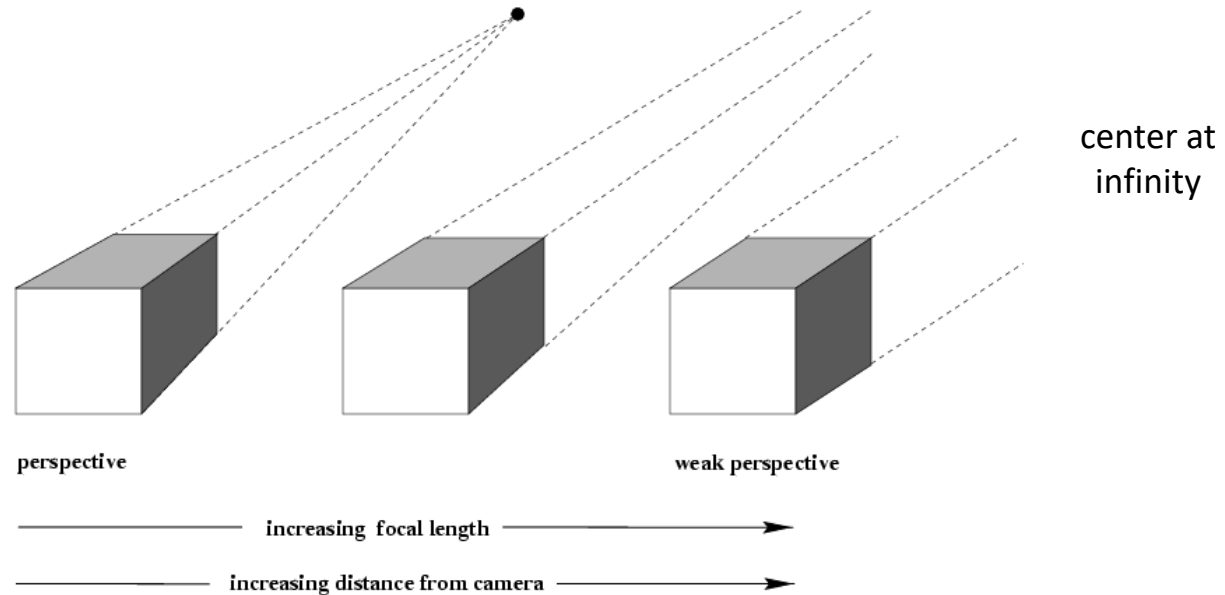
When the scene (or parts of it) is very far away.



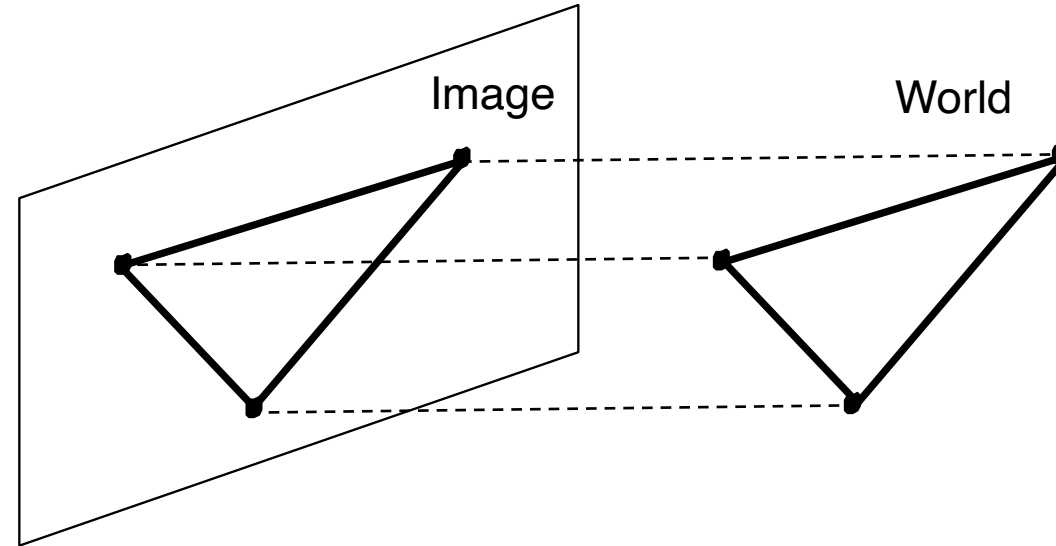
Weak perspective projection applies to the mountains.

Affine structure from motion

- Let's start with *affine* or *weak perspective* cameras



Orthographic projection



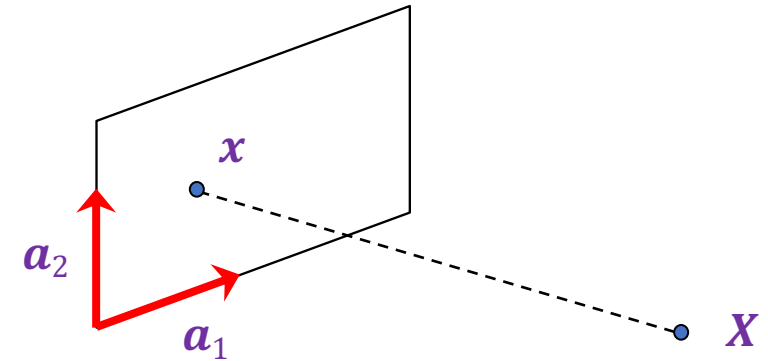
Just drop the z coordinate!

$$\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix}$$

General affine projection

- A general affine projection is a 3D-to-2D linear mapping plus translation:

$$P = \begin{bmatrix} a_{11} & a_{12} & a_{13} & t_1 \\ a_{21} & a_{22} & a_{23} & t_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} A & t \\ \mathbf{0}^T & 1 \end{bmatrix}$$



- In non-homogeneous coordinates:

$$\begin{pmatrix} x \\ y \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} t_1 \\ t_2 \end{pmatrix} = AX + t$$

a_1, a_2 : rows of projection matrix

Projection of
world origin

Affine structure from motion

- **Given:** m images of n fixed 3D points such that
 - $\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{X}_j + \mathbf{t}_i, \quad i = 1, \dots, m, j = 1, \dots, n$
- **Problem:** use the mn correspondences \mathbf{x}_{ij} to estimate m projection matrices \mathbf{A}_i and translation vectors \mathbf{t}_i , and n points \mathbf{X}_j
- The reconstruction is defined up to an arbitrary *affine* transformation \mathbf{Q} (12 degrees of freedom):

$$\begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{A} & \mathbf{t} \\ \mathbf{0}^T & 1 \end{bmatrix} \mathbf{Q}^{-1}, \quad \begin{pmatrix} \mathbf{X}_j \\ 1 \end{pmatrix} \rightarrow \mathbf{Q} \begin{pmatrix} \mathbf{X}_j \\ 1 \end{pmatrix}$$

- How many knowns and unknowns for m images and n points?
 - $2mn$ knowns and $8m + 3n$ unknowns
 - To be able to solve this problem, we must have $2mn \geq 8m + 3n - 12$ (affine ambiguity takes away 12 dof)
 - E.g., for **two** views, we need **four** point correspondences

Affine structure from motion

- First, center the data by subtracting the centroid of the image points in each view:

$$\begin{aligned}\hat{\mathbf{x}}_{ij} &= \mathbf{x}_{ij} - \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik} \\ &= \mathbf{A}_i \mathbf{X}_j + \mathbf{t}_i - \frac{1}{n} \sum_{k=1}^n (\mathbf{A}_i \mathbf{X}_k + \mathbf{t}_i) \\ &= \mathbf{A}_i \left(\mathbf{X}_j - \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \right) \\ &= \mathbf{A}_i \hat{\mathbf{X}}_j\end{aligned}$$

Affine structure from motion

- After centering, each normalized 2D point \hat{x}_{ij} is related to the 3D point by

$$\bullet \hat{x}_{ij} = A_i \hat{X}_j$$

- We can get rid of the need to center the 3D data (and the translation ambiguity) by defining the origin of the world coordinate system as the centroid of the 3D points

Affine structure from motion

- Let's create a $2m \times n$ data (measurement) matrix:

$$\bullet D = \begin{bmatrix} \hat{x}_{11} & \hat{x}_{12} & \cdots & \hat{x}_{1n} \\ \hat{x}_{21} & \hat{x}_{22} & \cdots & \hat{x}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{x}_{m1} & \hat{x}_{m2} & \cdots & \hat{x}_{mn} \end{bmatrix}$$

points (n)

cameras
($2m$)

$$\hat{x}_{ij} = A_i X_j$$

Affine structure from motion

- Let's create a $2m \times n$ data (measurement) matrix:

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}$$

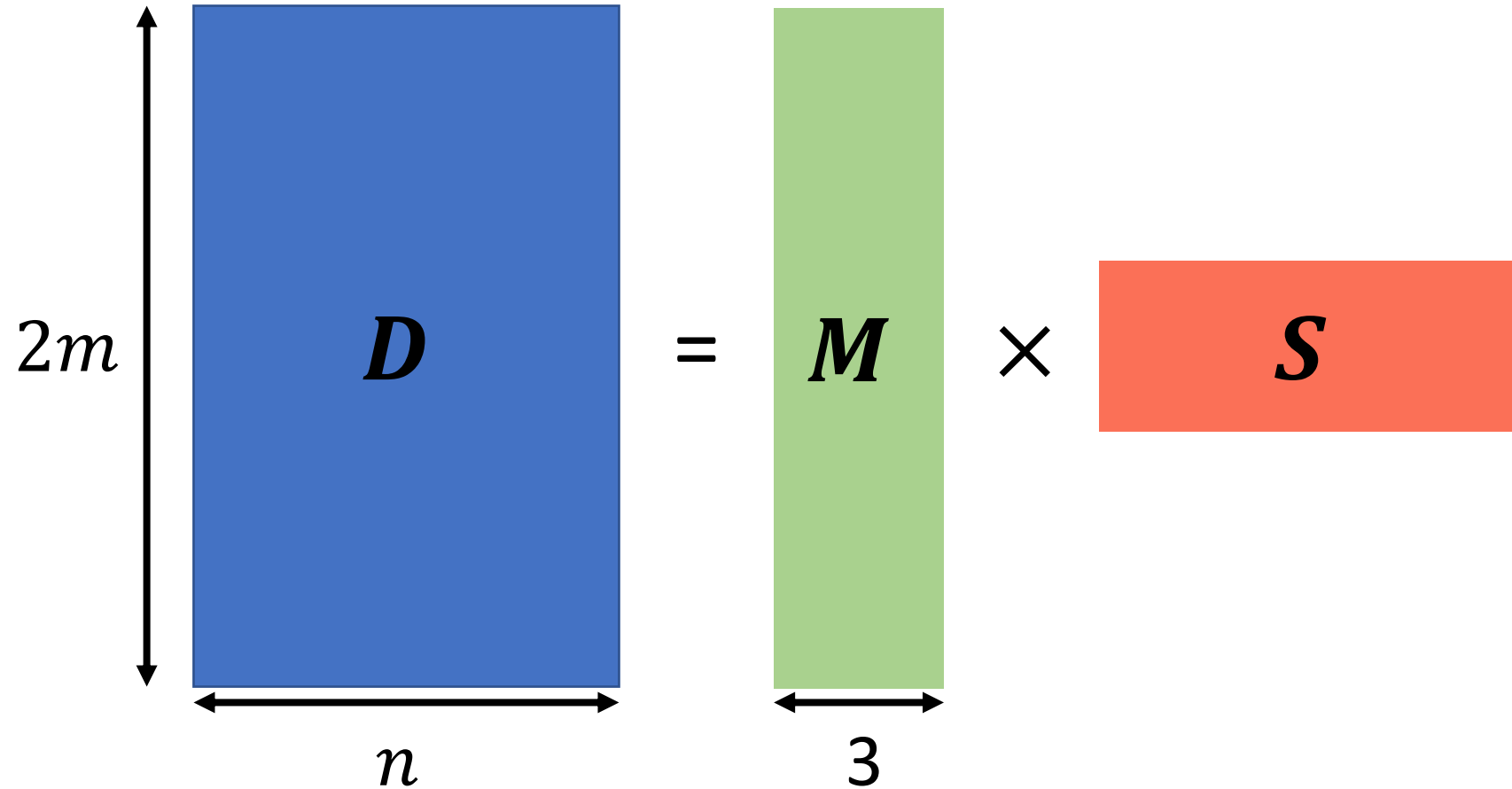
\mathbf{M}
cameras
($2m \times 3$)

\mathbf{S}
points ($3 \times n$)

- What must be the rank of the measurement matrix $\mathbf{D} = \mathbf{MS}$?

Factorizing the measurement matrix

- We want:



Factorizing the measurement matrix

- Perform SVD of D :

The diagram illustrates the SVD factorization of matrix D . It shows a blue rectangle representing matrix D with dimensions $2m \times n$. This is followed by an equals sign, then a green rectangle representing matrix U with dimensions $2m \times n$. This is followed by a multiplication sign, then a red square representing matrix Σ with dimensions $n \times n$. This is followed by another multiplication sign, then another red square representing matrix V^T with dimensions $n \times n$.

$$\begin{matrix} D \\ 2m \times n \end{matrix} = \begin{matrix} U \\ 2m \times n \end{matrix} \times \begin{matrix} \Sigma \\ n \times n \end{matrix} \times \begin{matrix} V^T \\ n \times n \end{matrix}$$

Factorizing the measurement matrix

- Keep top 3 singular values:

- This is the closest approximation of D with a rank-3 matrix in terms of Frobenius norm

The diagram illustrates the factorization of a matrix D into three matrices: U_3 , Σ_3 , and V_3^T . On the left is a blue rectangle representing matrix D with dimensions $2m \times n$. This is followed by an equals sign. To the right of the equals sign is a green rectangle representing matrix U_3 with dimensions $2m \times 3$. This is followed by a multiplication sign. To the right of the multiplication sign is a white rectangle representing matrix Σ_3 with dimensions 3×3 . This is followed by another multiplication sign. To the right of the second multiplication sign is a white rectangle representing matrix V_3^T with dimensions $3 \times n$. The top-left corner of the Σ_3 and V_3^T rectangles is shaded red.

- What to do about Σ_3 ?

- One solution: $M = U_3 \Sigma_3^{\frac{1}{2}}$, $S = \Sigma_3^{\frac{1}{2}} V_3^T$

Factorizing the measurement matrix

- One possible solution:

$$\begin{matrix} \mathbf{D} \\ 2m \times n \end{matrix} = \begin{matrix} \mathbf{M} \\ 2m \times 3 \end{matrix} \times \begin{matrix} \mathbf{S} \\ 3 \times n \end{matrix}$$
$$\mathbf{M} = \mathbf{U}_3 \boldsymbol{\Sigma}_3^{\frac{1}{2}}$$
$$\mathbf{S} = \boldsymbol{\Sigma}_3^{\frac{1}{2}} \mathbf{V}_3^T$$

Factorizing the measurement matrix

- Other possible solutions (Ambiguity in Reconstruction)

The diagram illustrates the factorization of the measurement matrix D into a product of four matrices: M , Q , Q^{-1} , and S . The matrix D is represented by a blue rectangle with dimensions $2m \times n$. It is equal to the product of a green rectangle M with dimensions $2m \times 3$, a pink square Q with dimensions 3×3 , another pink square Q^{-1} with dimensions 3×3 , and an orange rectangle S with dimensions $3 \times n$. The matrices are arranged in a sequence from left to right, connected by multiplication symbols.

$$D_{2m \times n} = M_{2m \times 3} \times Q_{3 \times 3} \times Q^{-1}_{3 \times 3} \times S_{3 \times n}$$

How to eliminate ambiguity?

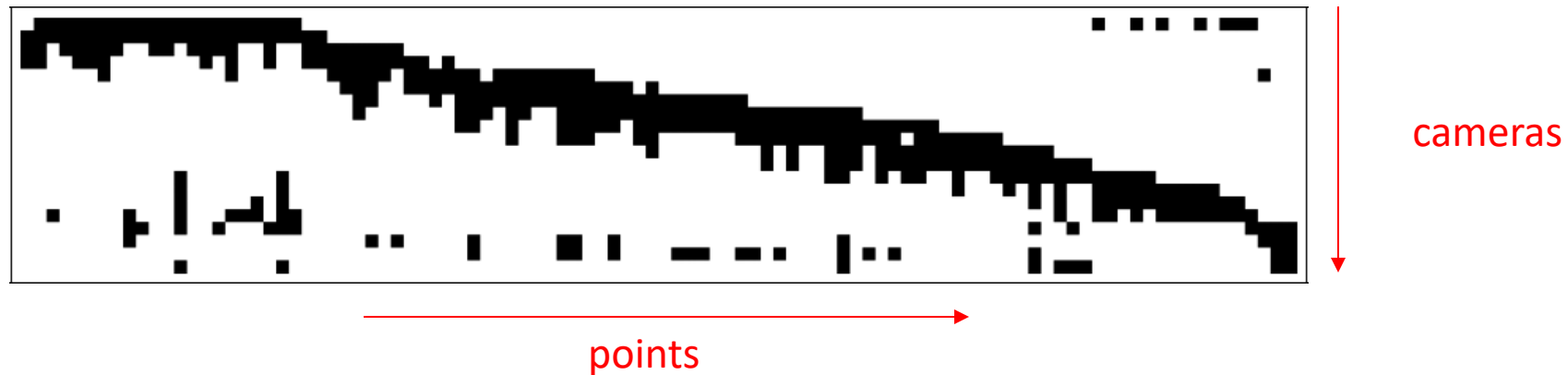
Assume certain special structure of the projection matrix.

Assume certain conditions about the 3D structure.

We can estimate Q to give the camera matrices in M desirable properties, like orthographic projection

Dealing with missing data

- So far, we have assumed that all points are visible in all views
- In reality, the measurement matrix typically looks something like this:



- These kind of problems are called Low-rank Matrix Completion problems (aka the Netflix Problem). Solved with convex/non-convex optimizations.
- Very popular before deep learning era!

Today's Class

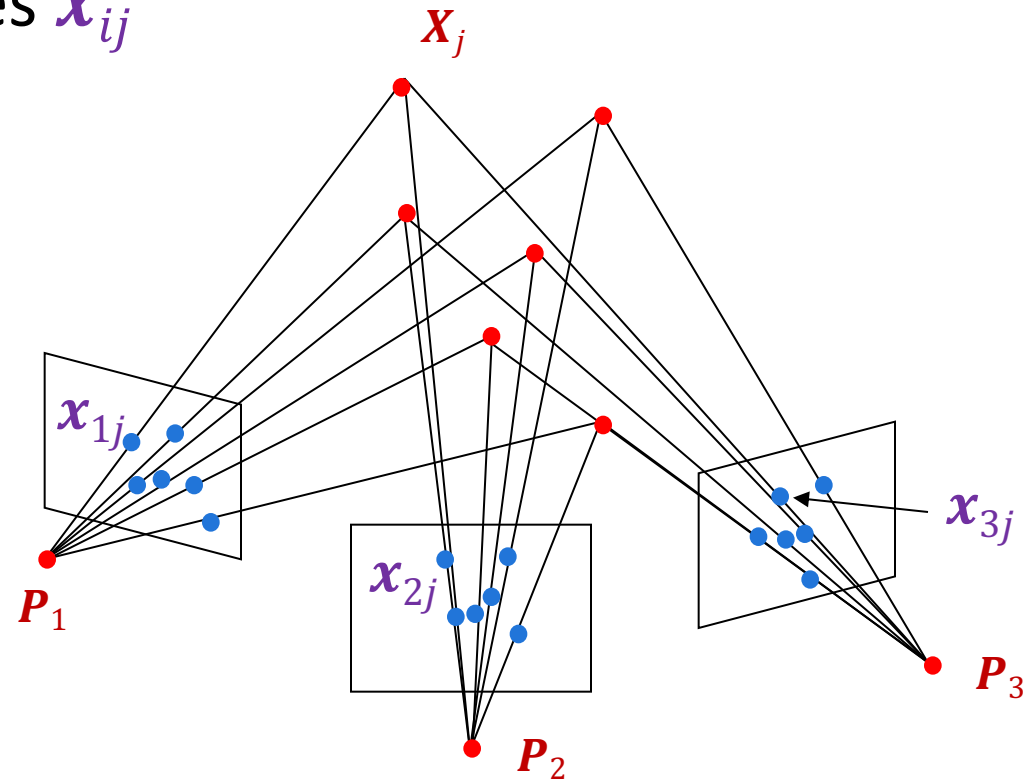
- Ambiguities in SfM
- Affine SfM
- **Projective SfM**
 - Global SfM
 - Incremental SfM
- Challenges and Applications

Projective structure from motion

- **Given:** m images of n fixed 3D points such that (ignoring visibility):

$$\bullet \mathbf{x}_{ij} \cong \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- **Problem:** estimate m projection matrices \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}



Projective structure from motion

- **Given:** m images of n fixed 3D points such that (ignoring visibility):

$$\bullet \mathbf{x}_{ij} \cong \mathbf{P}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

- **Problem:** estimate m projection matrices \mathbf{P}_i and n 3D points \mathbf{X}_j from the mn correspondences \mathbf{x}_{ij}
- With no calibration info, cameras and points can only be recovered up to a 4×4 projective transformation \mathbf{Q} :

$$\bullet \mathbf{X} \rightarrow \mathbf{QX}, \mathbf{P} \rightarrow \mathbf{PQ}^{-1}$$

- We can solve for structure and motion when $2mn \geq 11m + 3n - 15$
- For two cameras, at least **7 points** are needed
- You can solve it similar to Affine SfM with matrix factorization.
- Algebraic methods are good for initializing a non-linear optimization problem.

Bundle adjustment

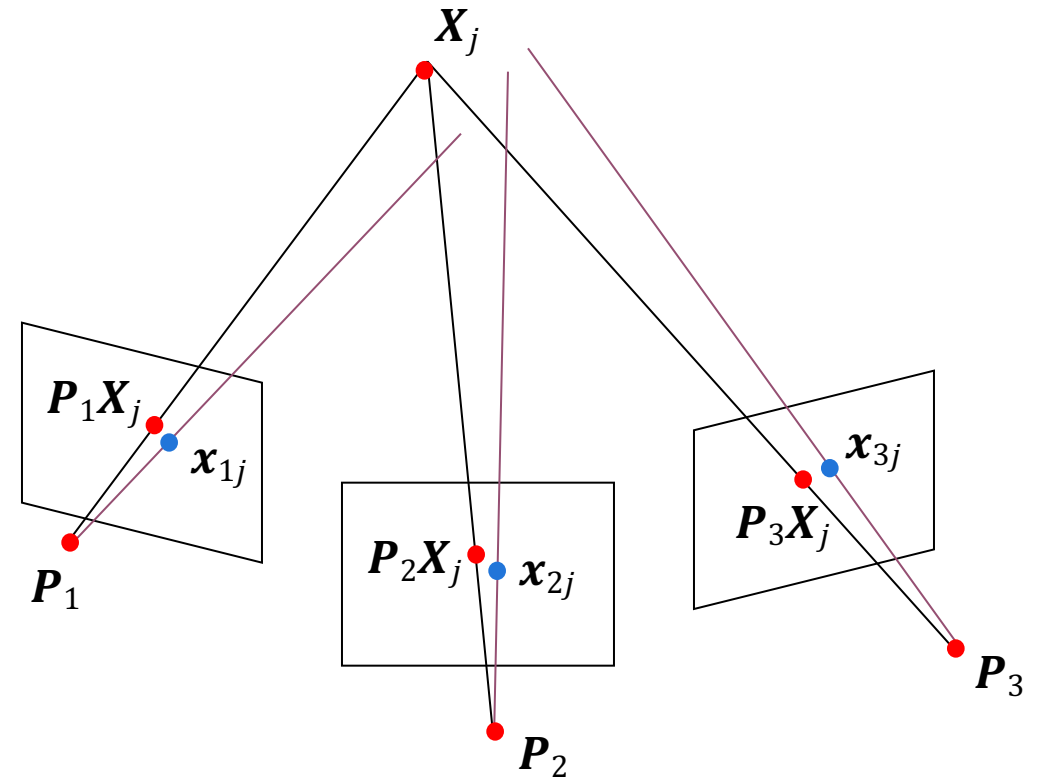
- Non-linear method for refining structure and motion
- Minimize reprojection error (with lots of bells and whistles):

$$\bullet \sum_{i=1}^m \sum_{j=1}^n w_{ij} d \left(\mathbf{x}_{ij} - \text{proj}(\mathbf{P}_i \mathbf{X}_j) \right)^2$$



visibility flag:
is point j visible in view i ?

Factorization based SfM works well for very small scenes with limited number of images, even then it produces poor result for most practical purposes.



Global Structure from Motion

SfM for large scale scenes

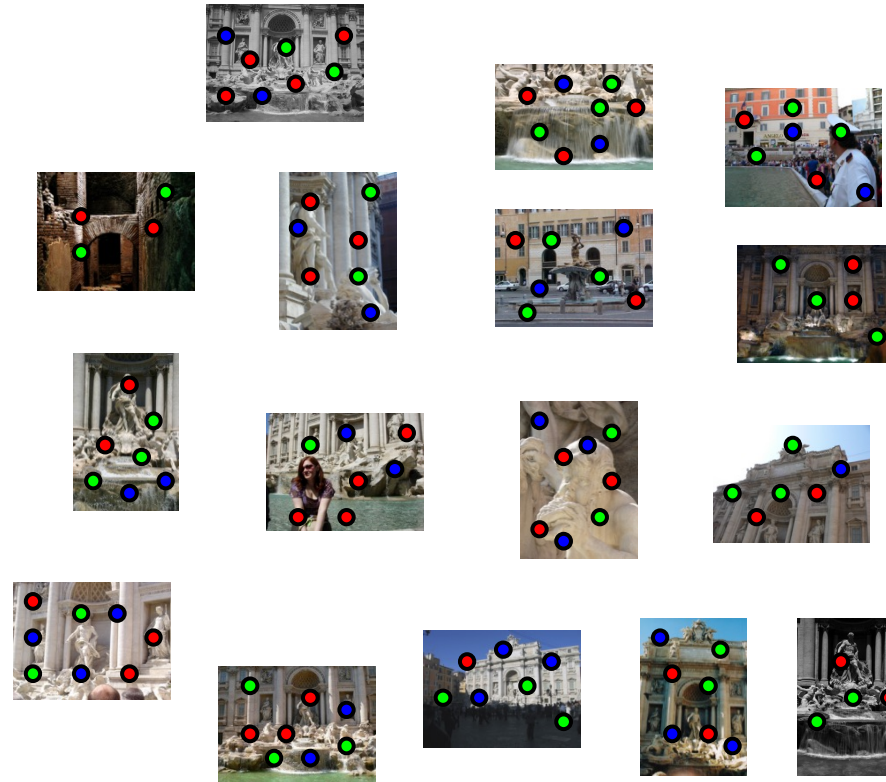
Feature detection

Detect features using SIFT [Lowe, IJCV 2004]



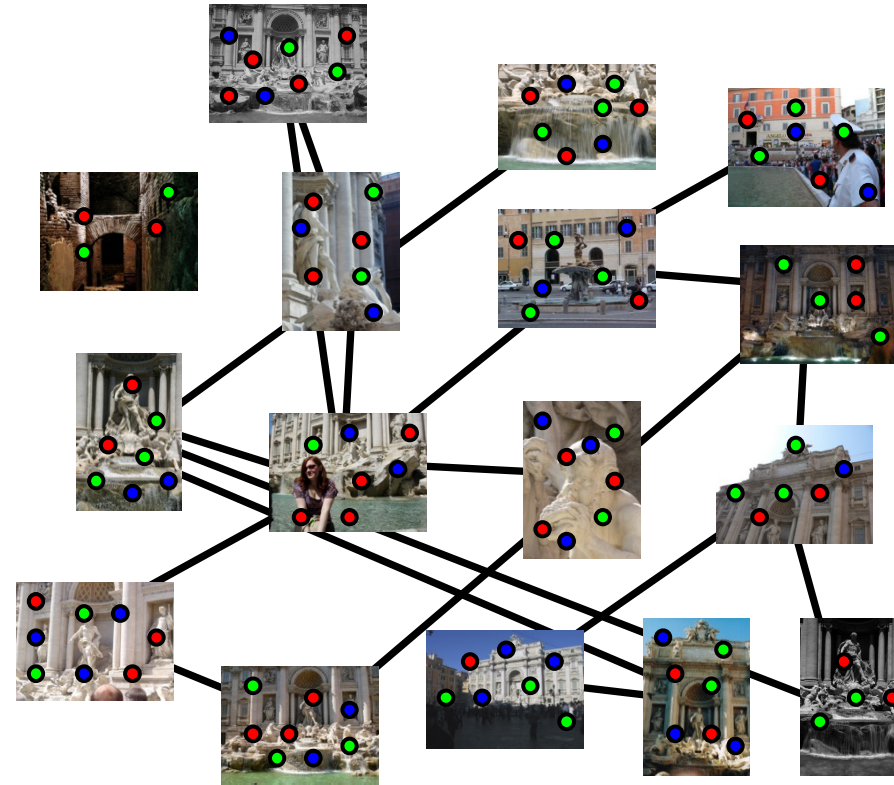
Feature detection

Detect features using SIFT [Lowe, IJCV 2004]



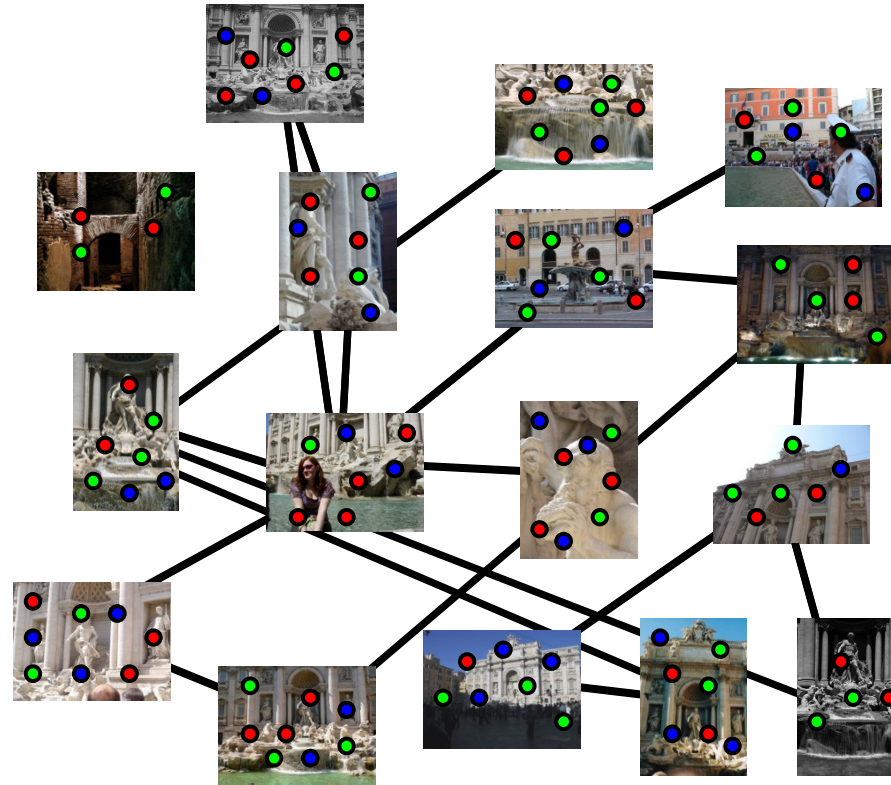
Feature matching

Match features between each pair of images



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

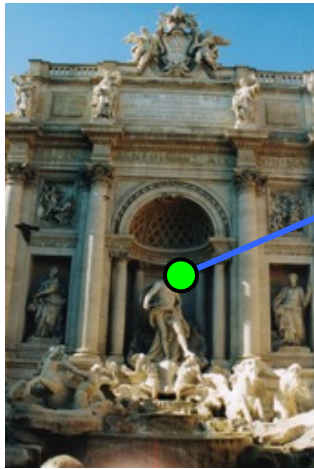


Image 1



Image 2

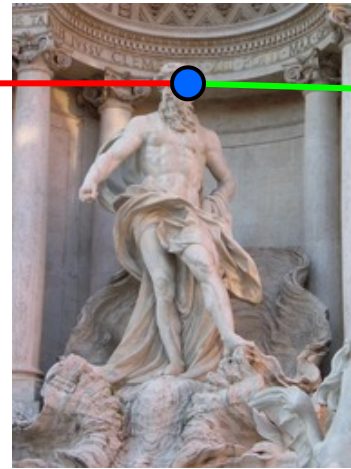


Image 3

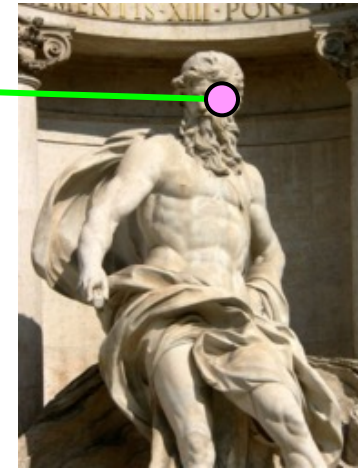
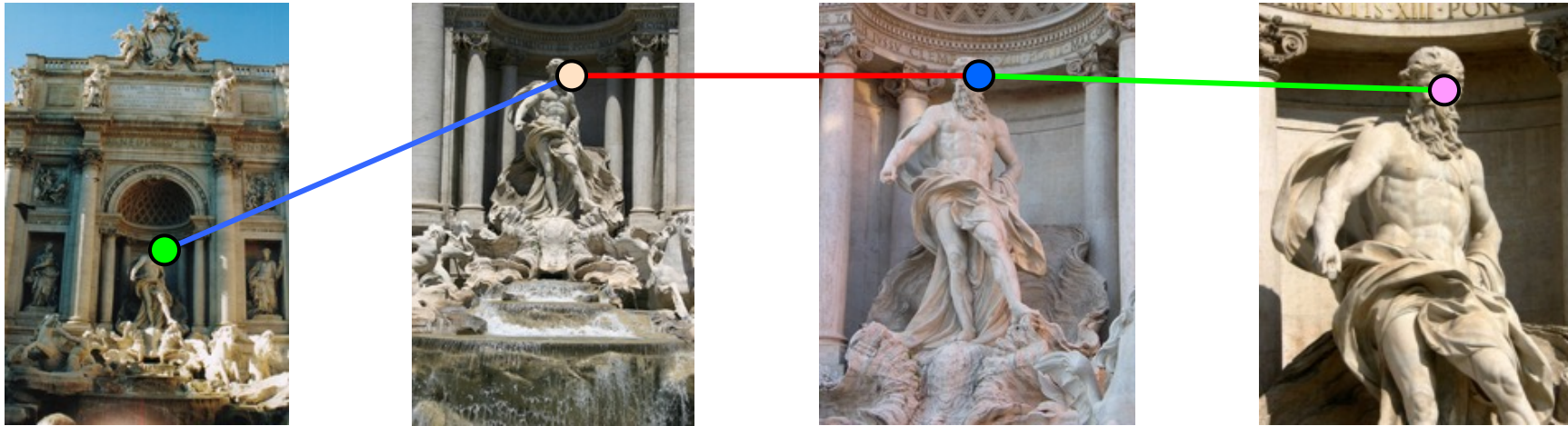


Image 4

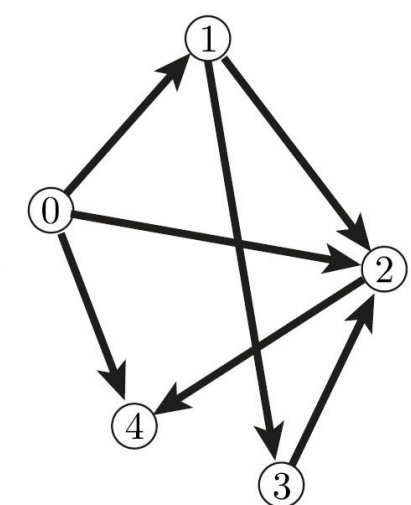
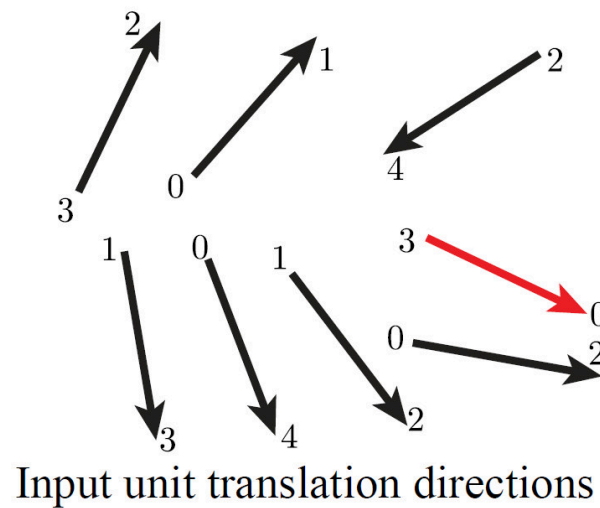
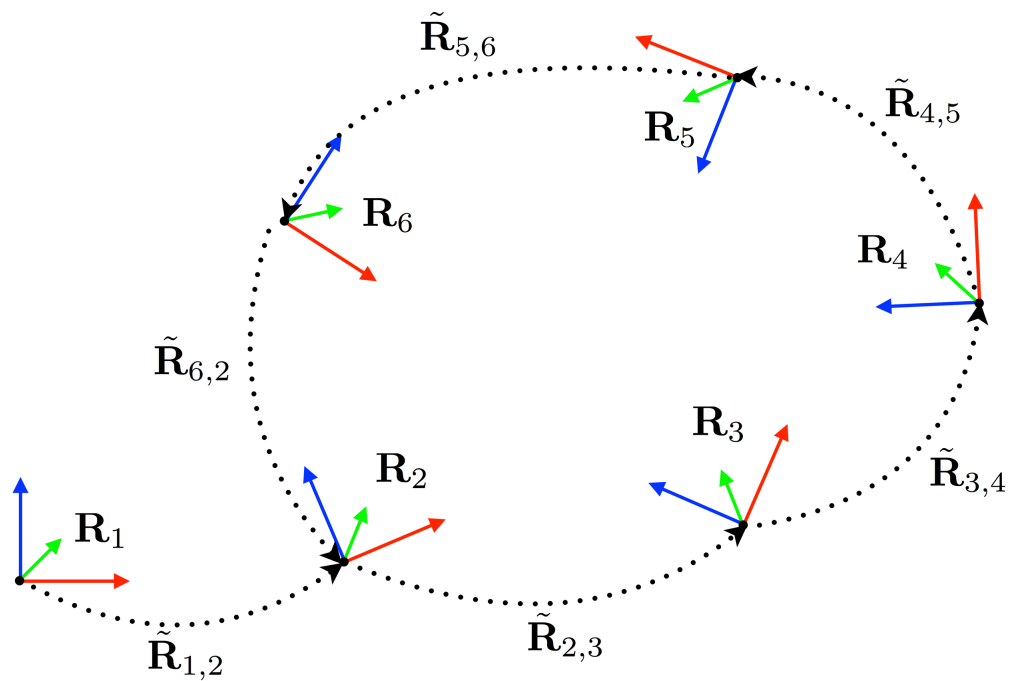


Global SfM

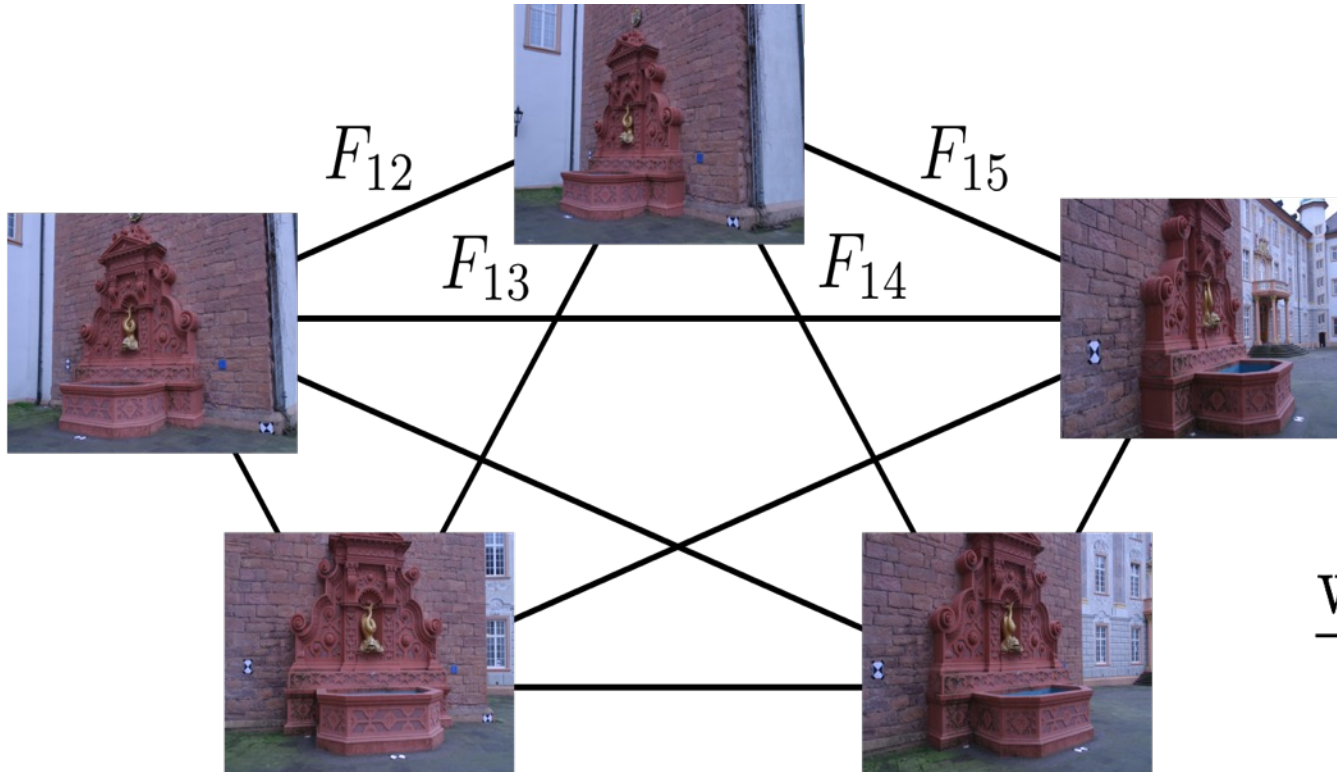
- Given N images, there are ${}^N C_2$ pairs. Many of these pairs will have no overlaps in views and/or Fundamental/Essential matrix between them can not be reliably estimated using RANSAC.
 - Consider we have N_0 ($N_0 < {}^N C_2$) pairs of images with fundamental matrix estimated
- For each N_0 pairs of images decompose essential matrix into relative rotation and translation between two cameras: R_{ij} and t_{ij} .
- Can we solve for global (world coordinate) rotation and translation of the cameras, given pairwise measurements, i.e.
 - Given R_{ij} and t_{ij} for N_0 pairs, find R_k & T_k for N cameras.
- Once we have the cameras we can better initialize the Bundle Adjustment problem.

Rotation & Translation Averaging

Given R_{ij} and t_{ij} for N_0 pairs, find R_k & T_k for N cameras



Camera Pose estimation as matrix completion over Fundamental matrices



$$F = \begin{bmatrix} 0 & F_{12} & F_{13} & F_{14} & F_{15} \\ F_{21} & 0 & F_{23} & F_{24} & F_{25} \\ F_{31} & F_{32} & 0 & F_{34} & F_{35} \\ F_{41} & F_{42} & F_{43} & 0 & F_{45} \\ F_{51} & F_{52} & F_{53} & F_{54} & 0 \end{bmatrix}$$

with $F = A + A^T$ and $rank(A) = 3$.

- Proves a low-rank property of all the cameras capturing different images of a scene.
- Solves a low-rank camera pose recovery algorithm from Structure from Motion.

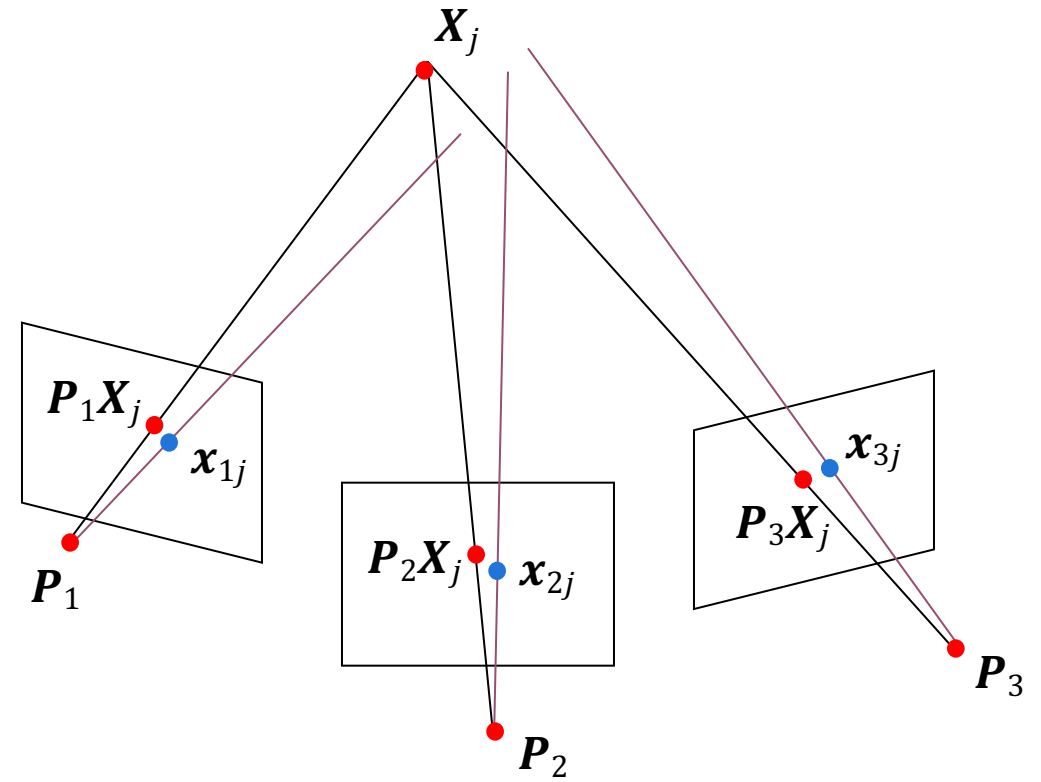
Bundle adjustment

- Non-linear method for refining structure and motion
- Minimize reprojection error (with lots of bells and whistles):
-

$$\bullet \sum_{i=1}^m \sum_{j=1}^n w_{ij} d \left(\mathbf{x}_{ij} - \text{proj}(\mathbf{P}_i \mathbf{X}_j) \right)^2$$

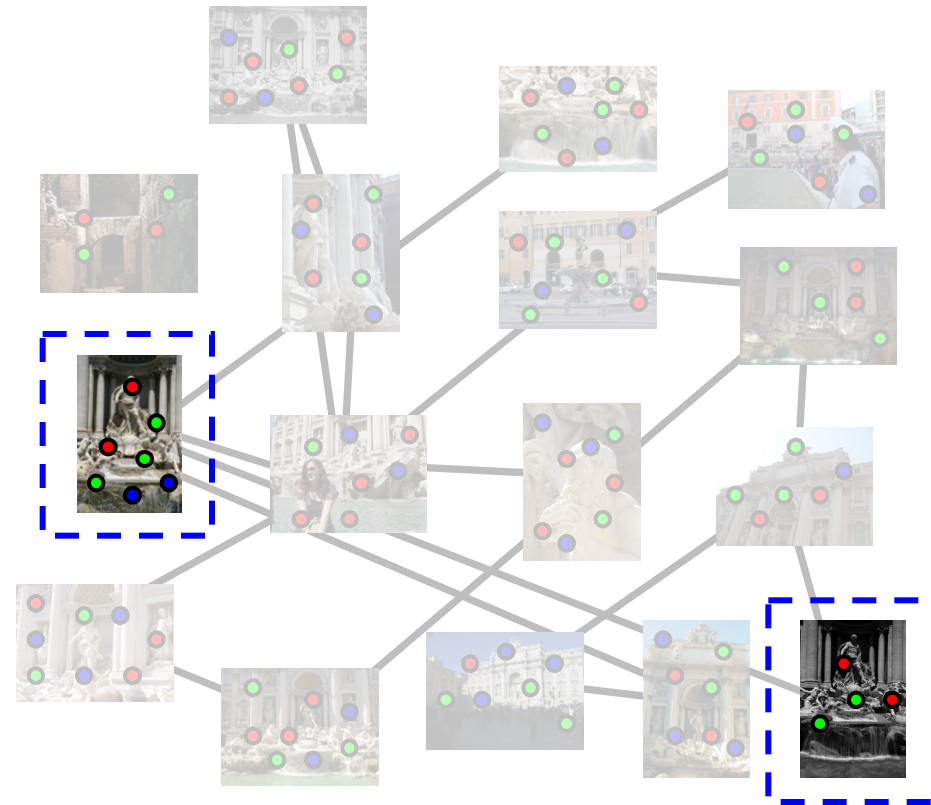
↑
visibility flag: is
point j visible in
view i ?

- **Initialize \mathbf{P}_i 's by solving global SfM**
 - **Rotation Averaging**
 - **Translation Averaging**



Incremental SfM

Can handle large scale scene, more than Global SfM



- Automatically select an initial pair of images

1. Picking the initial pair

- We want a pair with many matches, but which has as large a baseline as possible



✔ lots of matches
✘ small baseline



✔ large baseline
✘ very few matches



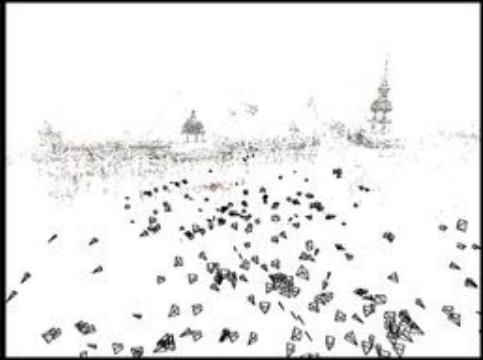
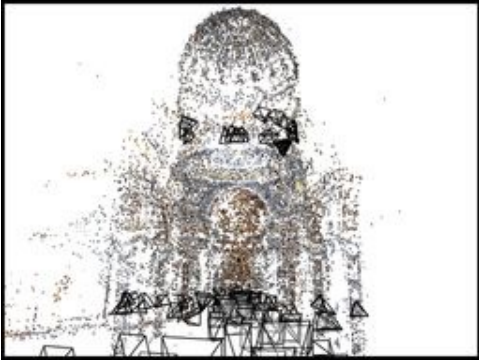
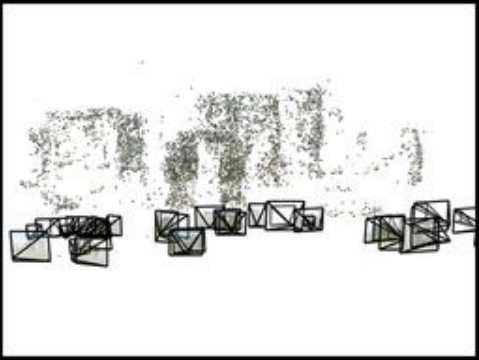
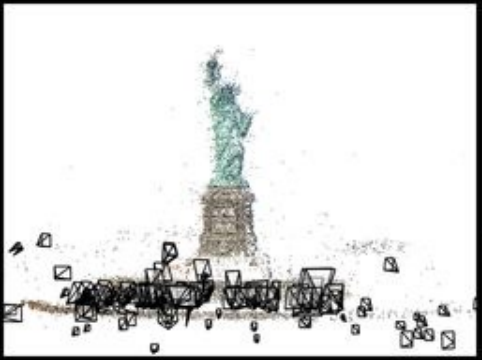
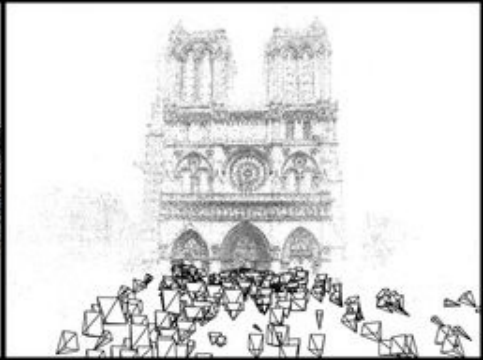
✔ large baseline
✔ lots of matches

Incremental SFM

- Pick a pair of images with lots of inliers (and preferably, good EXIF data)
 - Initialize intrinsic parameters (focal length, principal point) from EXIF
 - Estimate extrinsic parameters (R and t) using [five-point algorithm](#) (similar to 8-pt algorithm but for essential matrix)
 - Use triangulation to initialize model points
- While remaining images exist
 - Find an image with many feature matches with images in the model
 - Run RANSAC on feature matches to register new image to 3D model points
 - Triangulate new points
 - Perform bundle adjustment to re-optimize everything
 - Optionally, align with GPS from EXIF data or ground control points

Next Best View Problem

- Choice of next view impacts reconstruction quality
 - almost identical view => high uncertainty in triangulation
 - very different view => low overlap and high camera uncertainty
 - single bad choice may impact the whole reconstruction
- Popular next best view methods:
 - choose view with seeing the most triangulated points
 - minimize reconstruction uncertainty
 - depends on number of observations
 - distribution in the image



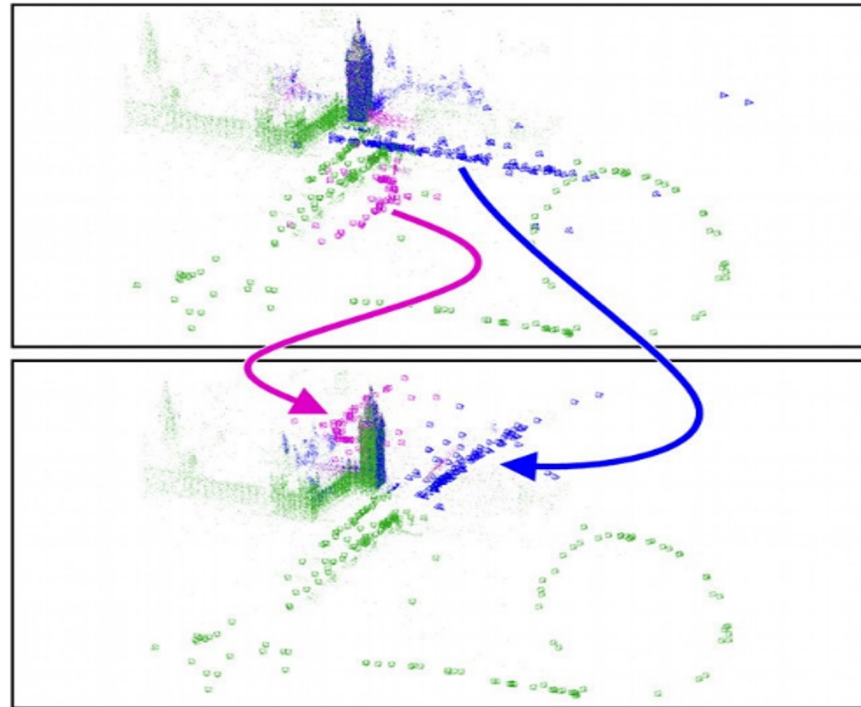
Today's Class

- Ambiguities in SfM
- Affine SfM
- Projective SfM
 - Global SfM
 - Incremental SfM
- **Challenges and Applications**

The devil is in the details

- Handling degenerate configurations (e.g., homographies)
- Filtering out incorrect matches
- Dealing with repetitions and symmetries

Repetitive structures cause catastrophic failures



Repetitive structures cause catastrophic failures



The devil is in the details

- Handling degenerate configurations (e.g., homographies)
- Filtering out incorrect matches
- Dealing with repetitions and symmetries
- Reducing error accumulation and closing loops

Loop Detection/Closure

- Problem:
 - Structure from motion is an incremental process
 - Drift accumulates
- Mitigation:
 - Retrieval of long range connections

Reducing error accumulation and closing loops



seattle1

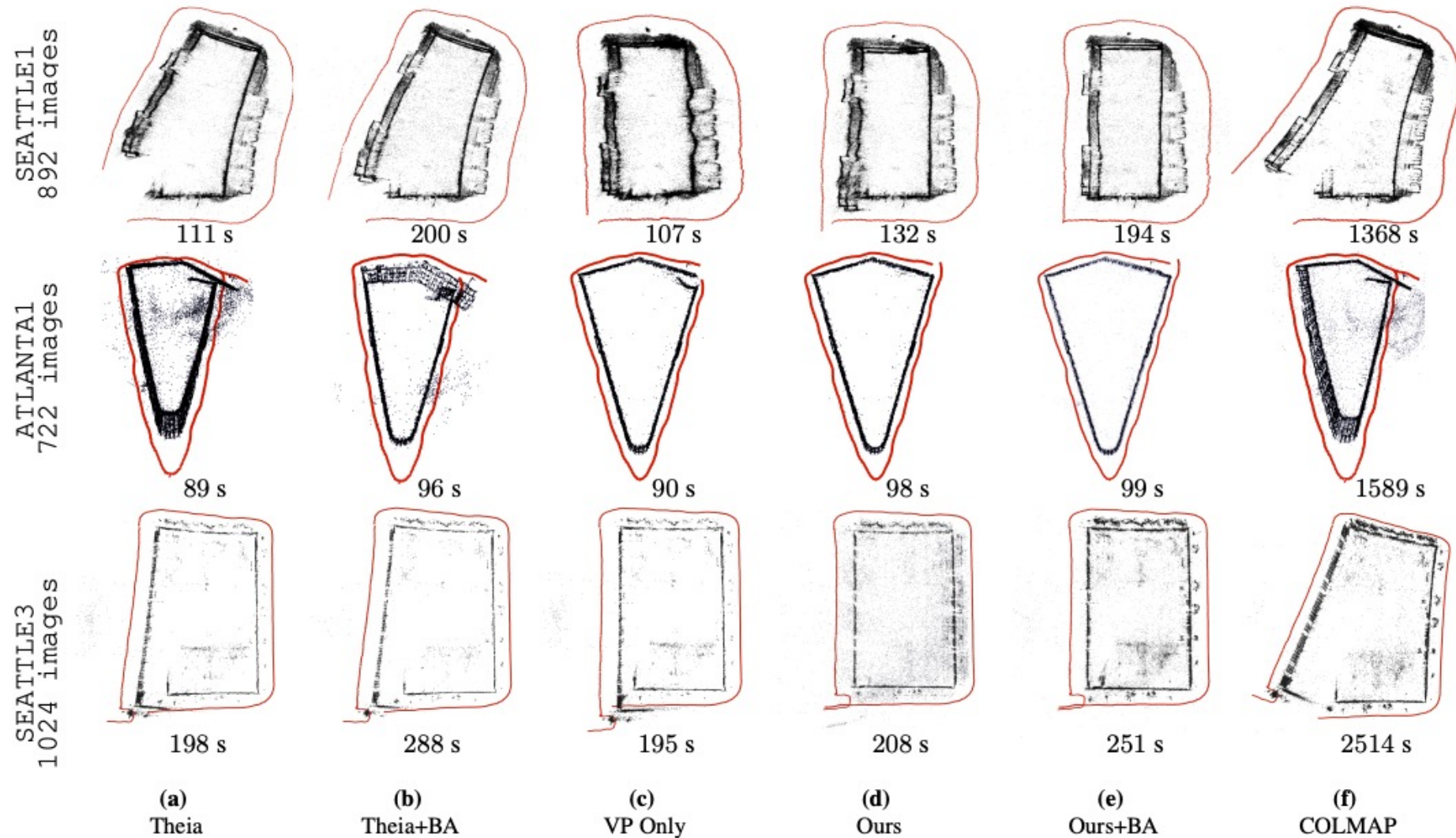
more_half

seattle2

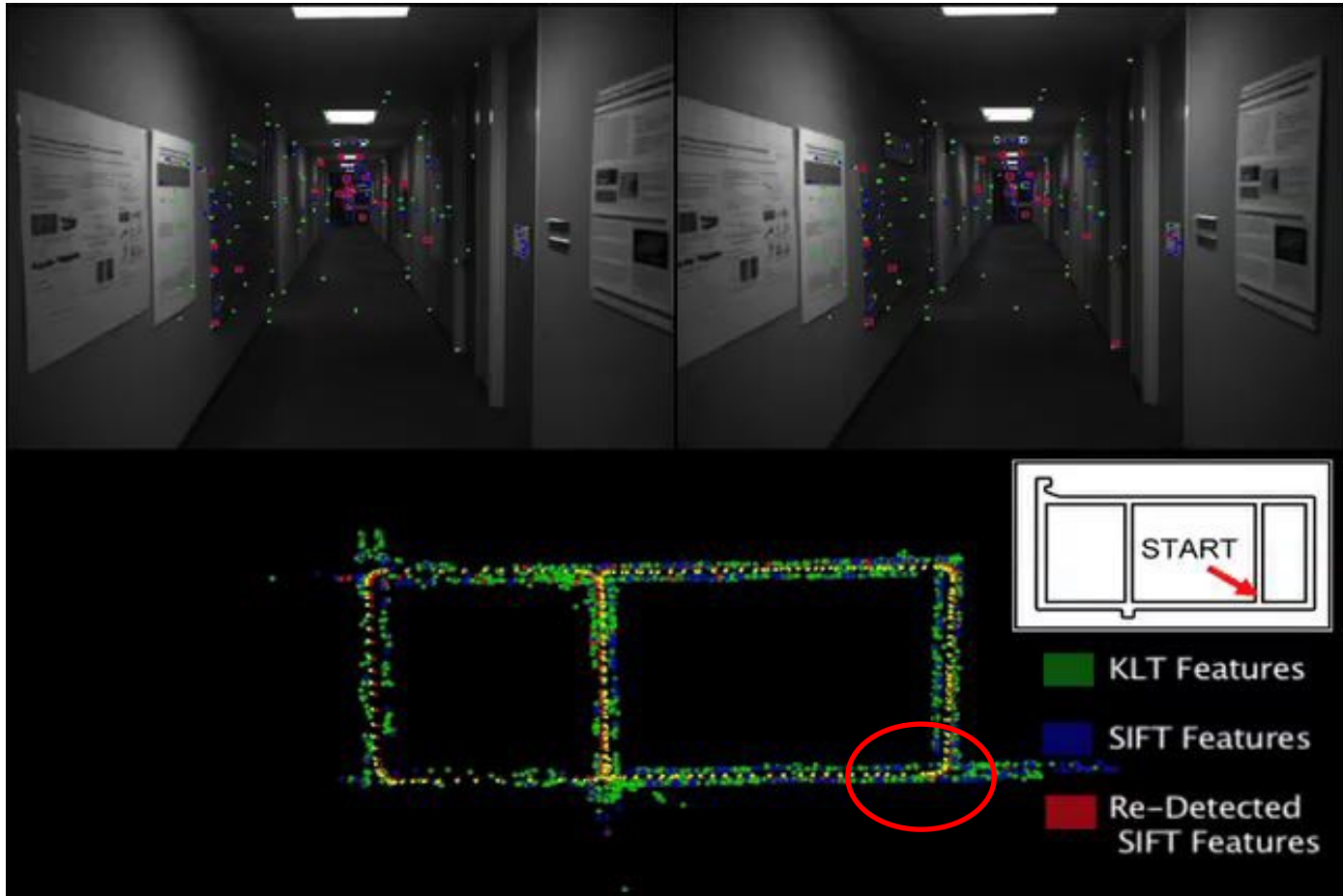
atlanta1

seattle3

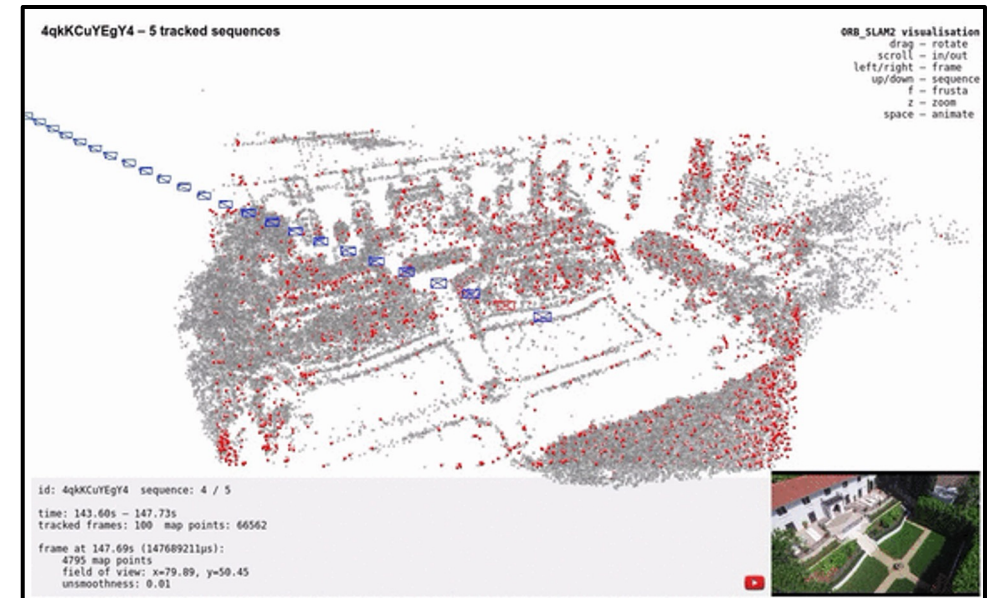
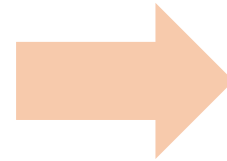
Reducing error accumulation and closing loops



Loop Closure

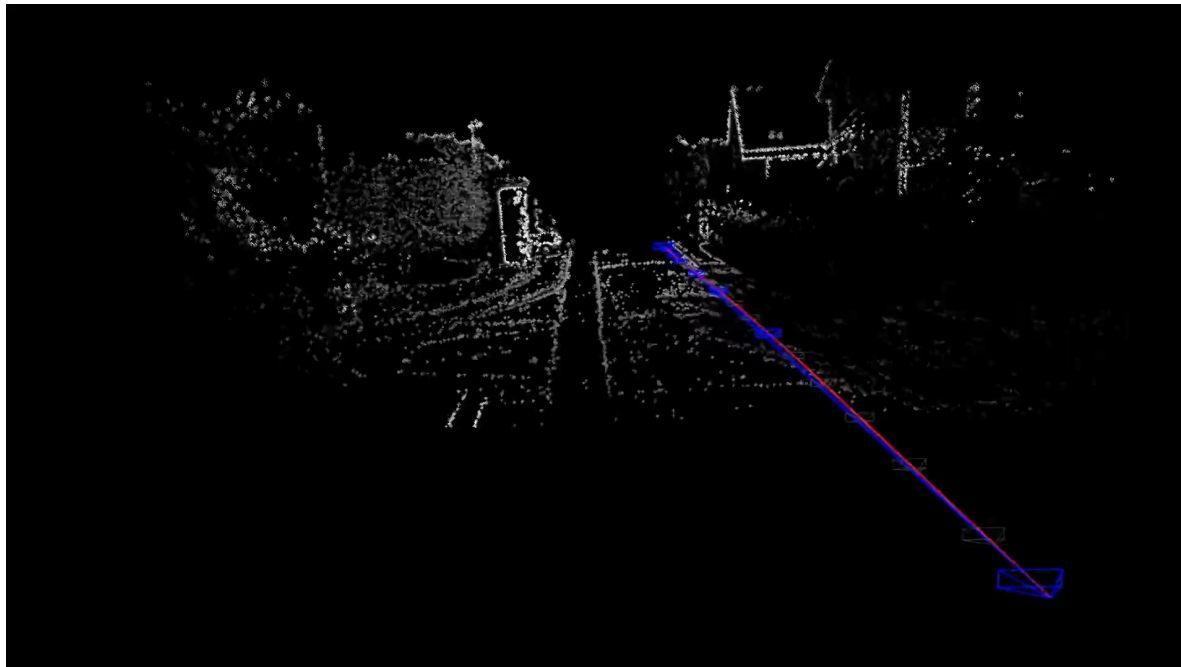


Can also compute camera poses from video (often called Visual SLAM)



Visual Simultaneous Localization and Mapping (V-SLAM)

- Main differences with SfM:
 - Continuous visual input from sensor(s) over time
 - Gives rise to problems such as loop closure
 - Often the goal is to be online / real-time



SFM software

- [Bundler](#)
- [OpenSfM](#)
- [OpenMVG](#)
- [VisualSFM](#)
- [COLMAP](#) ([Structure-from-motion revisited](#), JL Schonberger, JM Frahm, CVPR 2016, from UNC!)

- See also [Wikipedia's list of toolboxes](#)

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

Applications: Match Moving

Or Motion tracking, solving for camera trajectory
Integral for visual effects (VFX)

Why?



Applications: Visual Reality & Augmented Reality



Oculus

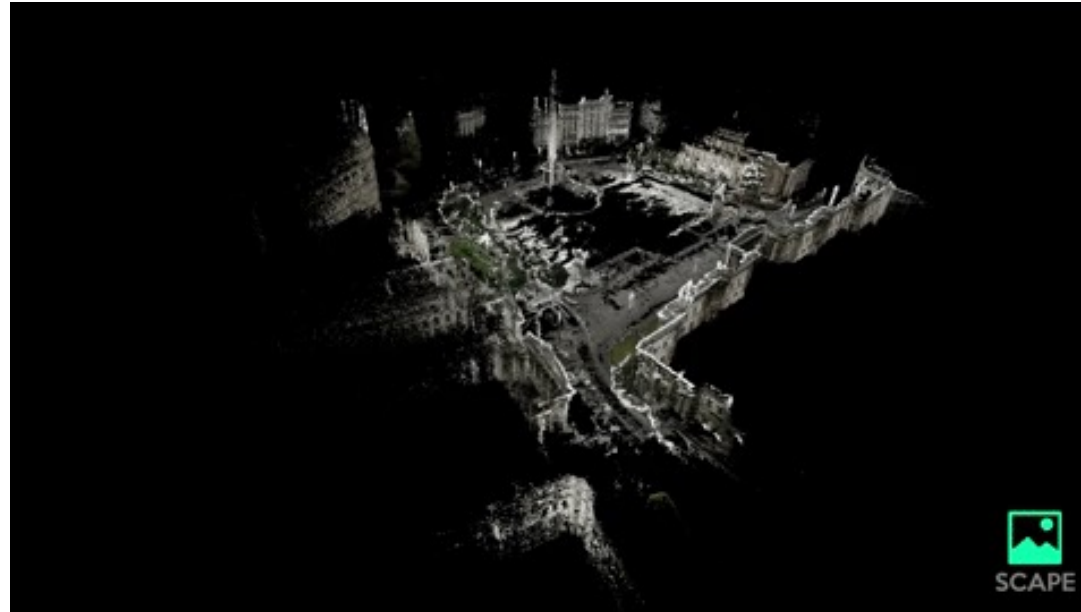
<https://www.youtube.com/watch?v=KOG7yTz1iTA>



Hololens

<https://www.youtube.com/watch?v=FMtvrTGnPO4>

Applications: Visual Reality & Augmented Reality



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

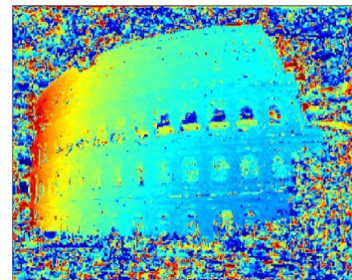
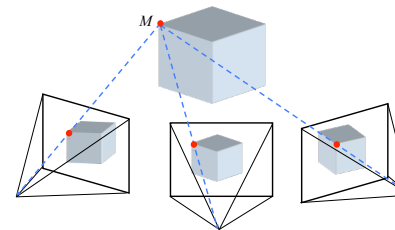
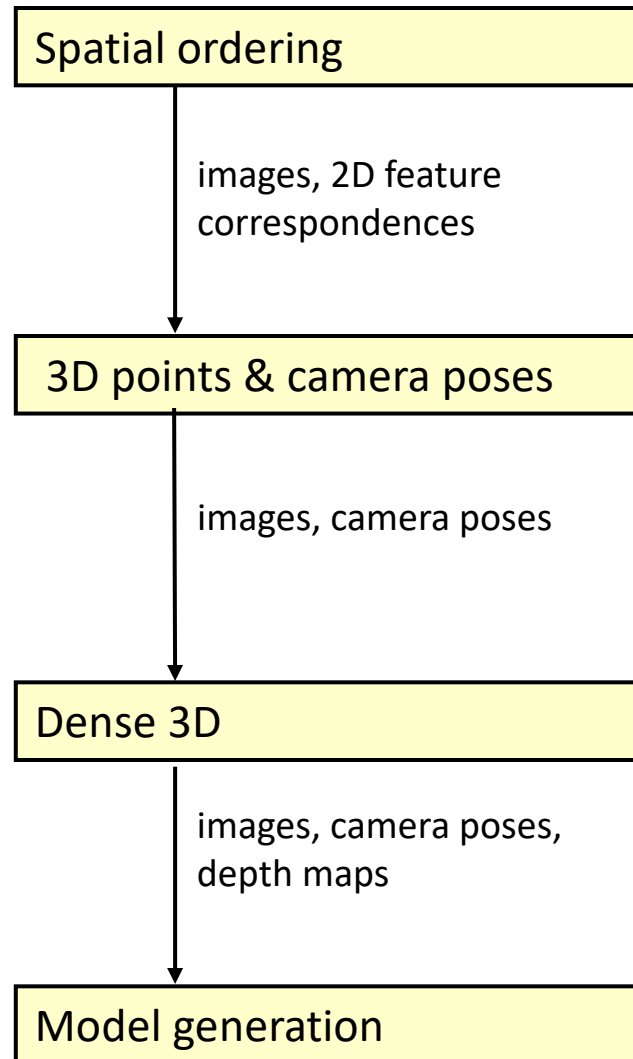
<https://medium.com/scape-technologies/building-the-ar-cloud-part-three-3d-maps-the-digital-scaffolding-of-the-21st-century-465fa55782dd>

Application: AR walking directions



<https://www.theverge.com/2019/8/8/20776247/google-maps-live-view-ar-walking-directions-ios-android-feature>

3D model from video



Summary: 3D geometric vision

- Fundamentals:
 - Camera Models: Intrinsic & Extrinsic
 - 3D to 2D projections, perspective distortions
 - Vanishing Points & Lines
 - Epipolar Geometry
 - Essential & Fundamental Matrices
- Core problems:
 - Camera calibration: single camera + two camera (estimate E/F matrix)
 - Stereo: depth from two calibrated cameras
- Reconstruction Techniques:
 - Active Stereo
 - Multi-view Stereo
 - Structure from Motion
 - Photometric Stereo (next class)

Slide Credits

- [CS5670, Introduction to Computer Vision](#), **Cornell Tech**, by **Noah Snavely**.
- [CS 194-26/294-26: Intro to Computer Vision and Computational Photography](#), **UC Berkeley**, by **Angjoo Kanazawa**.
- **CS 543** [Computer Vision](#), by **Stevlana Lazebnik**, **UIUC**.
- **COMP 776**, by **Jan-Michael Frahm**, **UNC**