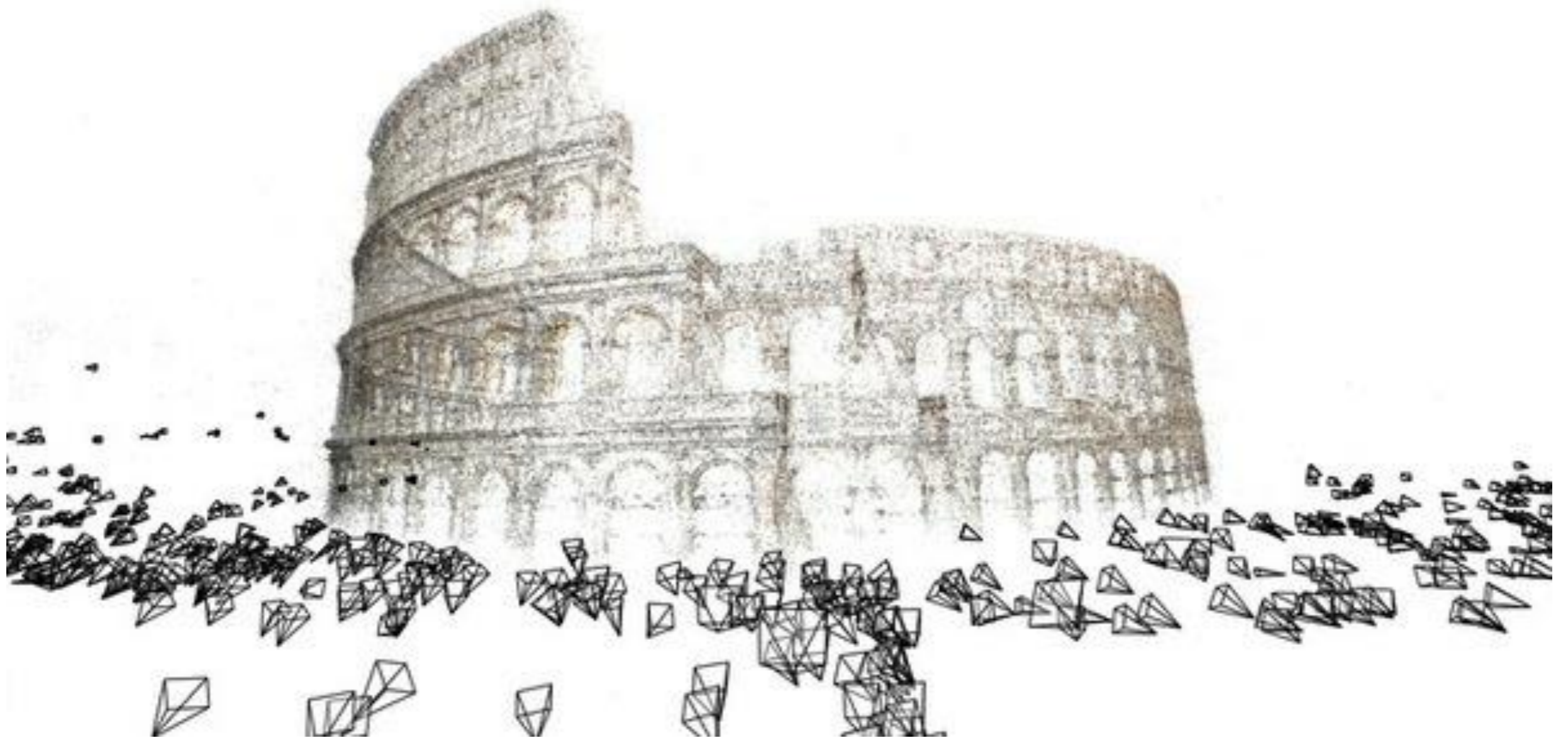# Structure from motion

# Course announcements

- Homework 3 has been posted and is due on March 10th.
  - Yes, this is during the spring break *per popular demand.*
  - No, you don't *have* to work during spring break:
    - -- This is the same homework that was originally planned for March 8th.
    - -- You can finish the homework by March 8th.
    - -- Shifting the deadline to March 10th means that everyone gets two extra late days *for free.*
  - Any questions about the homework?
  - How many of you have looked at/started/finished homework 3?

- Grades for homework 1 will be posted tonight.

- Grades for homework 2 will be posted before the mid-semester grades are due.

- Yannis will have extra office hours **Tuesday 3-5 pm.**

# Overview of today's lecture

Leftover from lecture 11:

- Template matching.

- Structured light.

New in lecture 12:

- A note on normalization.

- Two-view structure from motion.

- Ambiguities in structure from motion.

- Affine structure from motion.

- Multi-view structure from motion.

- Large-scale structure from motion.

# Slide credits

Many of these slides were adapted from:

- Kris Kitani (16-385, Spring 2017).

- Noah Snavely (Cornell University).

- Rob Fergus (New York University).

# A note on normalization

# Estimating F – 8-point algorithm

- The fundamental matrix F is defined by

$$\mathbf{x'}^{\mathrm{T}} \mathbf{Fx} = 0$$

  for any pair of matches x and x' in two images.

- Let x=$(u,v,1)^{\mathrm{T}}$ and x'=$(u',v',1)^{\mathrm{T}}$,   $\mathbf{F} = \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}$

  each match gives a linear equation

$$uu'f_{11} + vu'f_{12} + u'f_{13} + uv'f_{21} + vv'f_{22} + v'f_{23} + uf_{31} + vf_{32} + f_{33} = 0$$

# Problem with 8-point algorithm

$$\begin{bmatrix} u_1 u_1' & v_1 u_1' & u_1' & u_1 v_1' & v_1 v_1' & v_1' & u_1 & v_1 & 1 \\ u_2 u_2' & v_2 u_2' & u_2' & u_2 v_2' & v_2 v_2' & v_2' & u_2 & v_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u_n u_n' & v_n u_n' & u_n' & u_n v_n' & v_n v_n' & v_n' & u_n & v_n & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = 0$$

~10000    ~10000    ~100    ~10000    ~10000    ~100    ~100    ~100    1
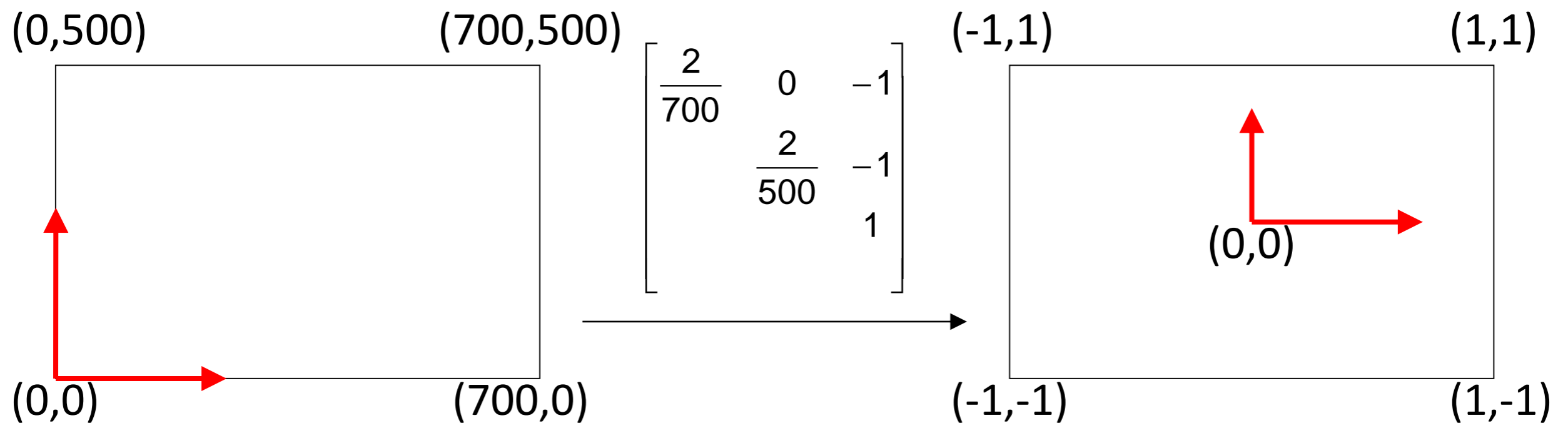
⚠ Orders of magnitude difference
between column of data matrix
→ least-squares yields poor results

# Normalized 8-point algorithm

normalized least squares yields good results

Transform image to ~[-1,1]x[-1,1]

# Normalized 8-point algorithm

1. Transform input by $\hat{\mathbf{x}}_i = \mathbf{T}\mathbf{x}_i$, $\hat{\mathbf{x}}_i' = \mathbf{T}\mathbf{x}_i'$

2. Call 8-point on $\hat{\mathbf{x}}_i$, $\hat{\mathbf{x}}_i'$ to obtain $\hat{\mathbf{F}}$

3. $\mathbf{F} = \mathbf{T}'^{\mathrm{T}}\hat{\mathbf{F}}\mathbf{T}$

$$\mathbf{x}'^{\mathrm{T}}\mathbf{F}\mathbf{x} = 0$$

$$\underbrace{\hat{\mathbf{x}}'^{\mathrm{T}}\mathbf{T}'^{-\mathrm{T}}\,\mathbf{F}\,\mathbf{T}^{-1}\hat{\mathbf{x}}}_{\hat{\mathbf{F}}} = 0$$

# Normalized 8-point algorithm

```
[x1, T1] = normalise2dpts(x1);
[x2, T2] = normalise2dpts(x2);

A = [x2(1,:)'.*x1(1,:)'   x2(1,:)'.*x1(2,:)'  x2(1,:)' ...
      x2(2,:)'.*x1(1,:)'   x2(2,:)'.*x1(2,:)'  x2(2,:)' ...
      x1(1,:)'             x1(2,:)'            ones(npts,1) ];

[U,D,V] = svd(A);

F = reshape(V(:,9),3,3)';

[U,D,V] = svd(F);
F = U*diag([D(1,1) D(2,2) 0])*V';

% Denormalise
F = T2'*F*T1;
```
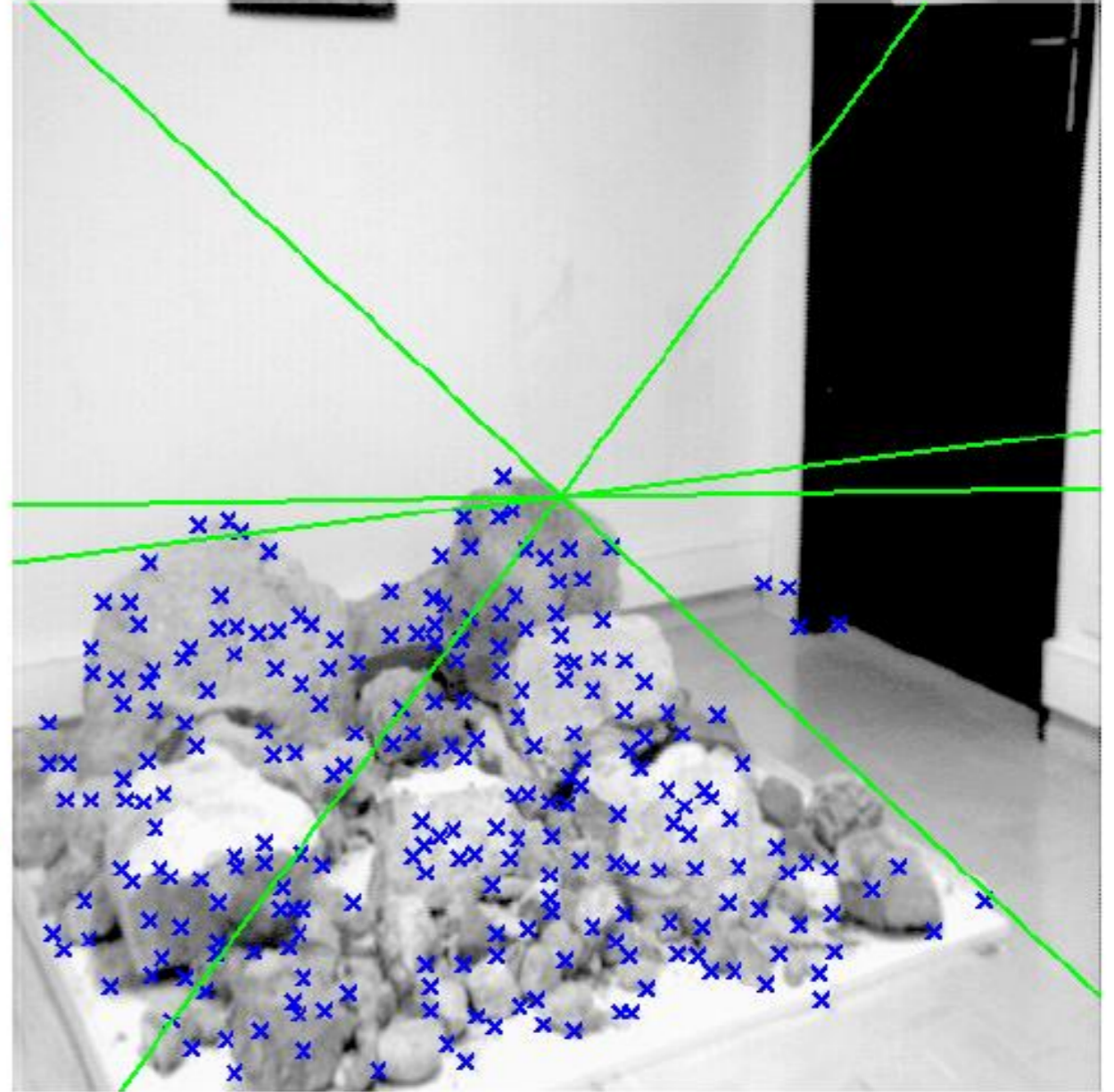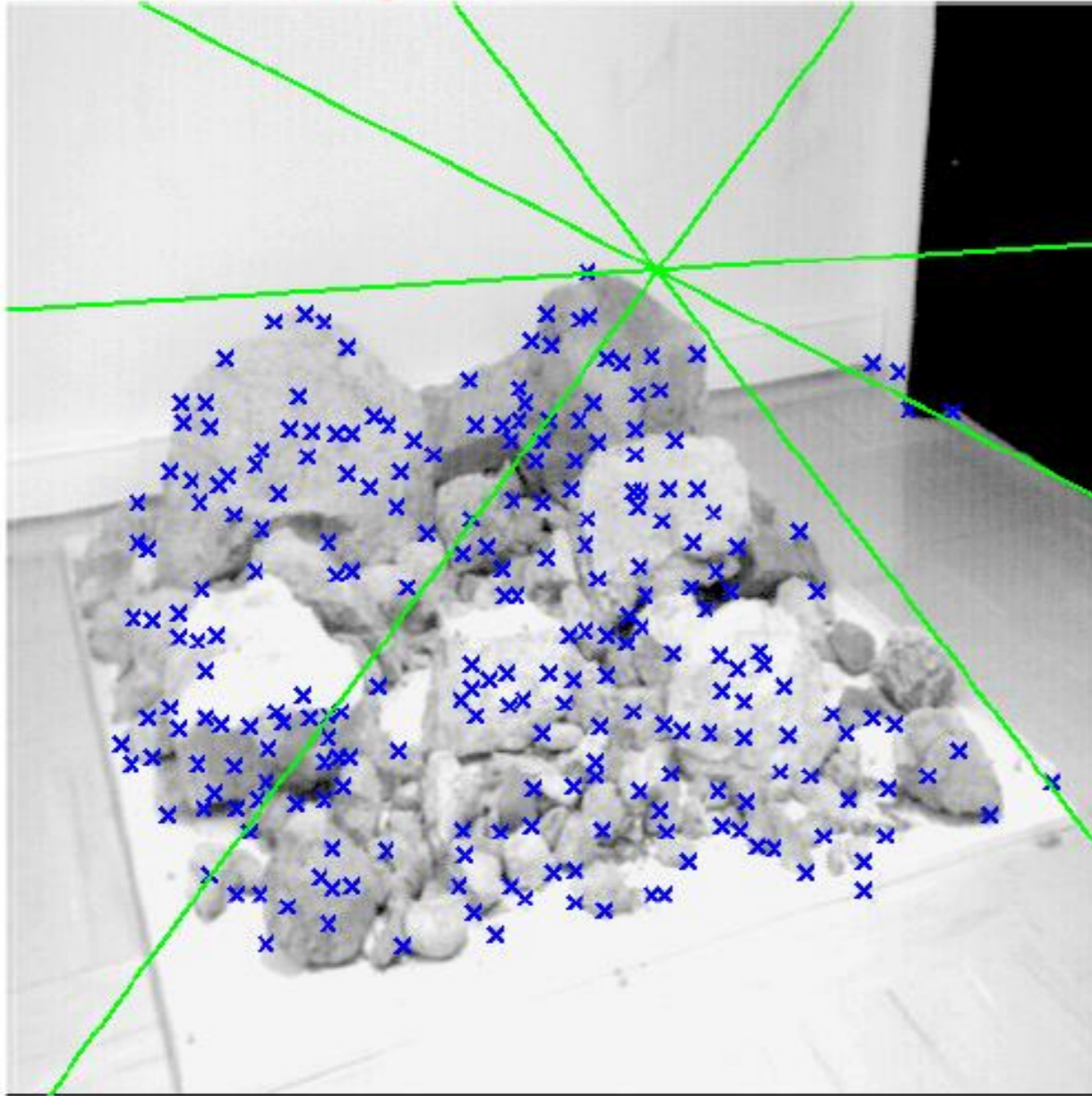
# Results (ground truth)



■ **Ground truth**   with standard stereo calibration
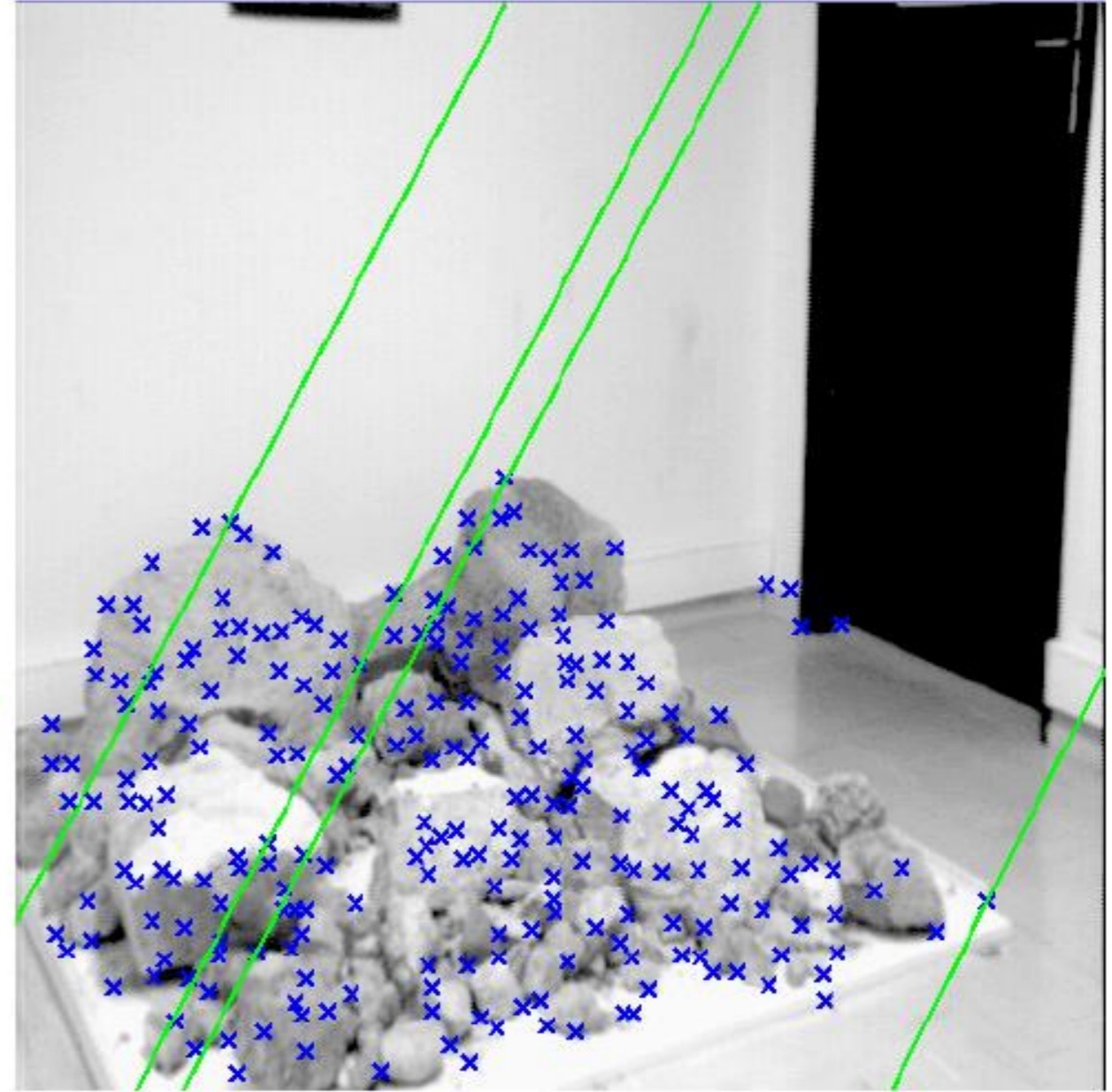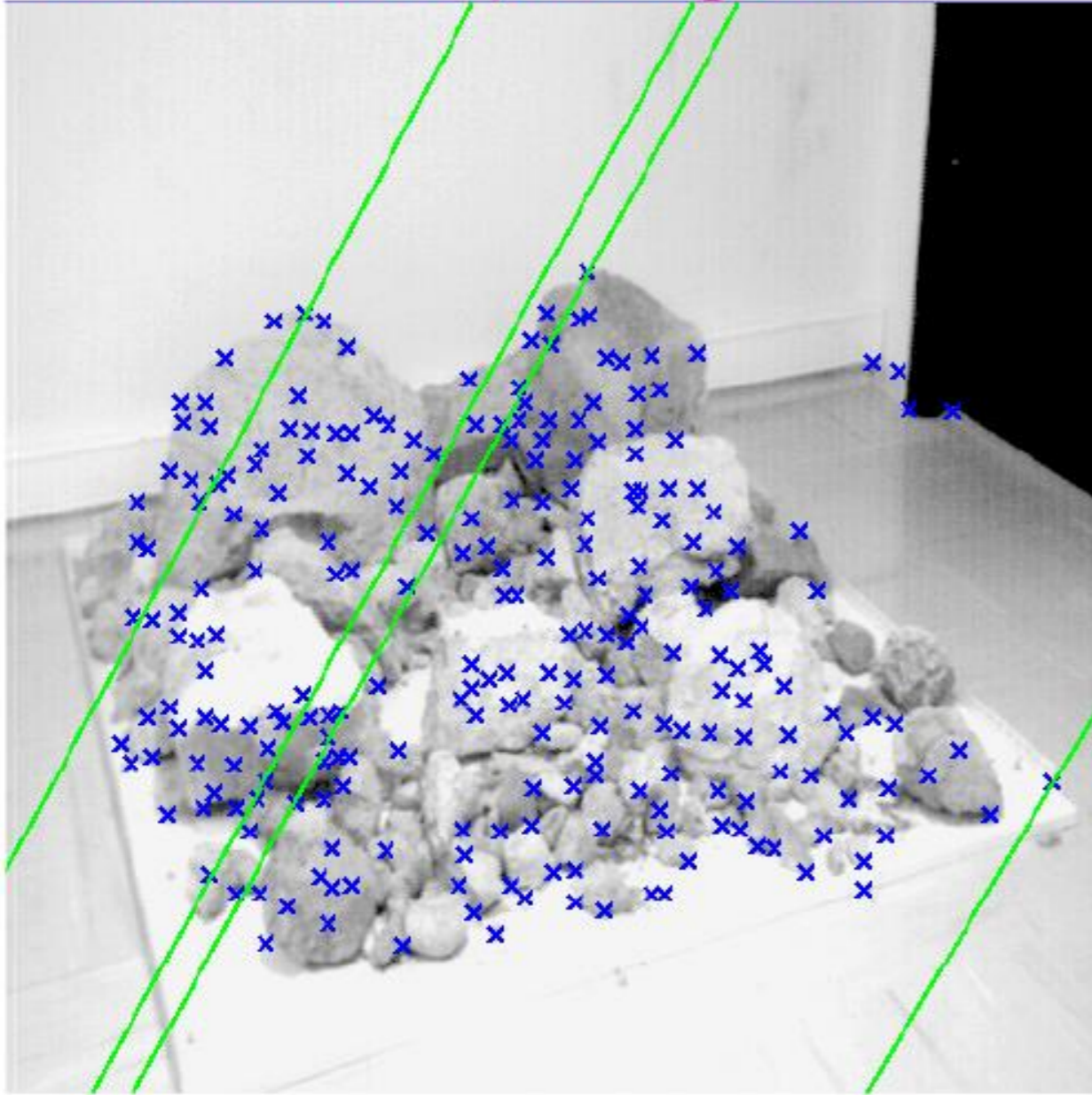
# Results (8-point algorithm)

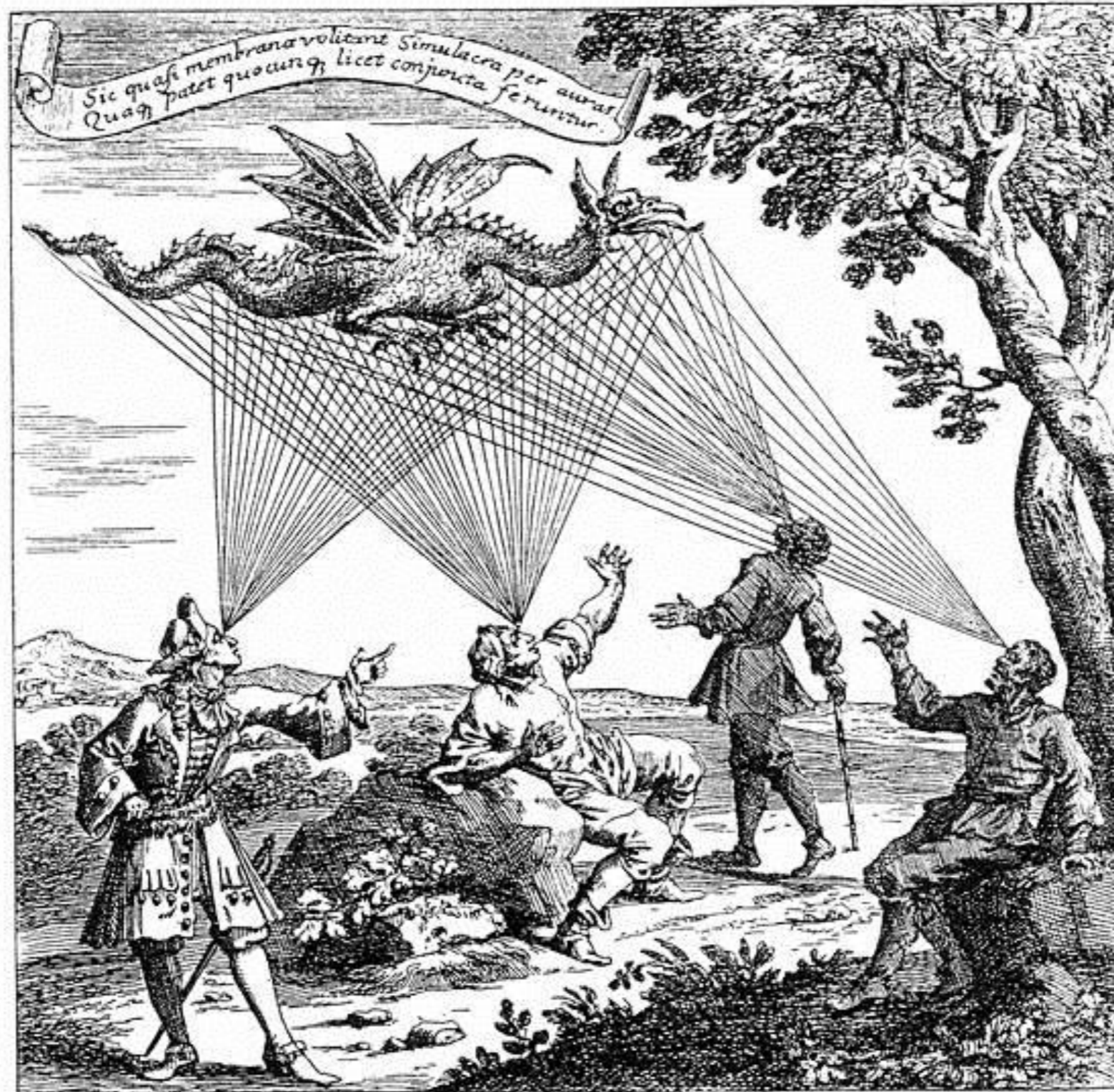# Results (normalized 8-point algorithm)



Normalized 8-point algorithm

# Two-view structure from motion

|  | Structure (scene geometry) | Motion (camera geometry) | Measurements |
|---|---|---|---|
| Pose Estimation | known | estimate | 3D to 2D correspondences |
| Triangulation | estimate | known | 2D to 2D cooorespondences |
| Reconstruction | estimate | estimate | 2D to 2D cooorespondences |

# Structure from motion



*Sic quasi membrana volitant Simulacra per auras*
*Quaq; patet quoaunq; licet conjecta feruntur.*

Драконъ, видимый подъ различными углами зрѣнія

По гравюрѣ на мѣди изъ „Oculus artificialis teledioptricus" Цана. 1702 года.

# Camera calibration & triangulation

- Suppose we know 3D points

  - And have matches between these points and an image

  - How can we compute the camera parameters?

- Suppose we have know camera parameters, each of which observes a point

  - How can we compute the 3D location of that point?

# Structure from motion

- SfM solves both of these problems *at once*

- A kind of chicken-and-egg problem

  - (but solvable)

# Reconstruction

(2 view structure from motion)

Given a set of matched points

$$\{x_i, x_i'\}$$

Estimate the camera matrices

$$\mathbf{P}, \mathbf{P}'$$

Estimate the 3D point

$$\mathbf{X}$$

# Reconstruction

(2 view structure from motion)

Given a set of matched points

$$\{x_i, x_i'\}$$

Estimate the camera matrices

$$\mathbf{P}, \mathbf{P'}$$ ← 'motion'
(of the cameras)

Estimate the 3D point

$$\mathbf{X}$$ ← 'structure'

# Two-view SfM

1. Compute the Fundamental Matrix **F** from points correspondences

$$x'^{\top}_m \mathbf{F} x_m = 0$$

# Two-view SfM

1. Compute the Fundamental Matrix **F** from points correspondences
   **8-point algorithm**

$$x'^{\top}_m \mathbf{F} x_m = 0$$

# Two-view SfM

1. Compute the Fundamental Matrix **F** from points correspondences
   **8-point algorithm**

2. Compute the camera matrices **P** from the Fundamental matrix
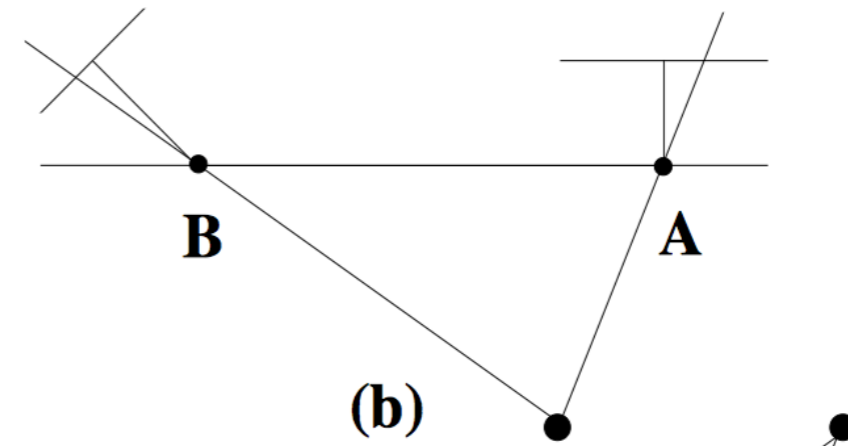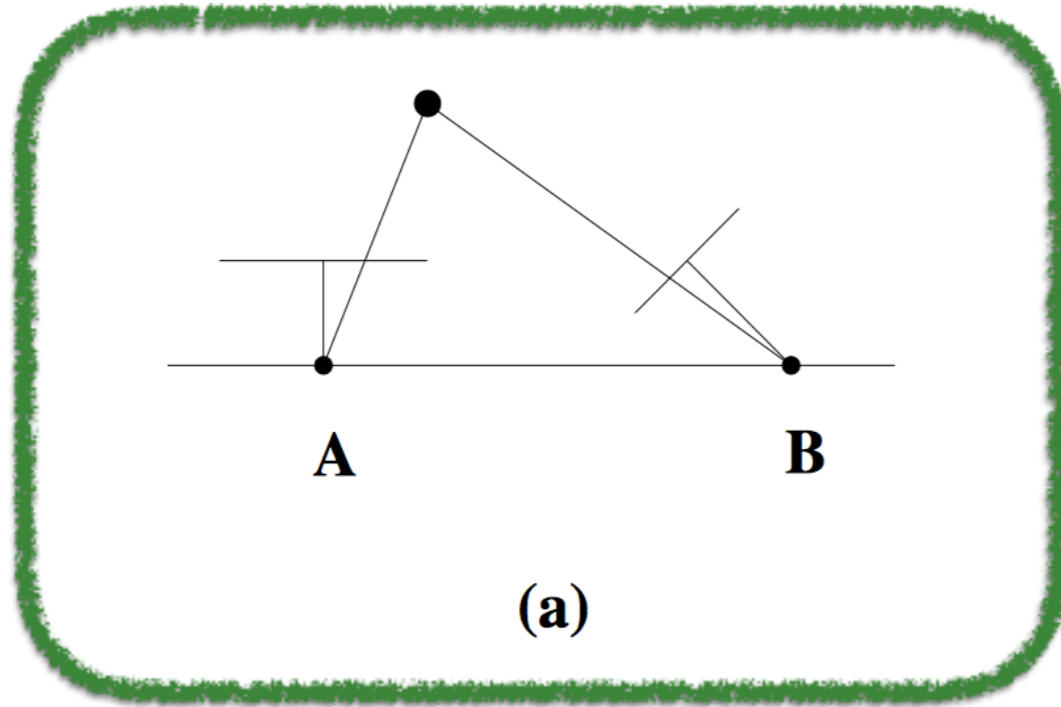   **P** = [ **I** | **0** ]  and  **P'** = [ [e$_x$]**F** | **e'** ]

Camera matrices corresponding to the fundamental matrix **F** may be chosen as

$$\mathbf{P} = [\mathbf{I}|\mathbf{0}] \quad \mathbf{P}' = [[e_\times]\mathbf{F}|e']$$

(See Hartley and Zisserman C.9 for proof)

# Find the configuration where the points is in front of both cameras



**(a)**

**(b)**

**(c)**

**(d)**

# Two-view SfM

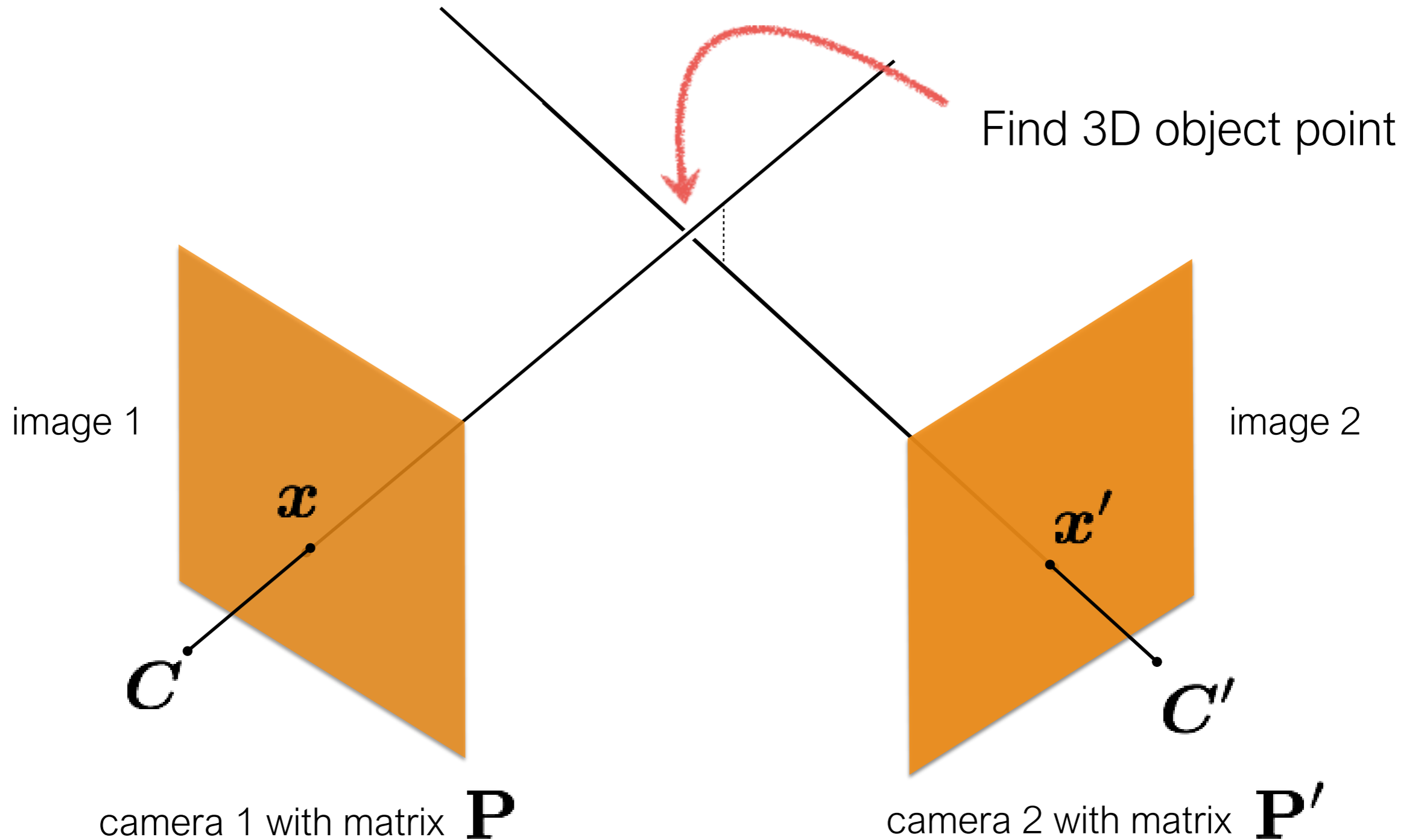1. Compute the Fundamental Matrix **F** from points correspondences
   **8-point algorithm**

2. Compute the camera matrices **P** from the Fundamental matrix
   $\mathbf{P} = [\ \mathbf{I}\ |\ \mathbf{0}\ ]$ and $\mathbf{P'} = [\ [\mathbf{e'}_x]\mathbf{F}\ |\ \mathbf{e'}\ ]$

3. For each point correspondence, compute the point **X** in 3D space (triangularization)
   **DLT** with $x = \mathbf{P}\ X$ and $x' = \mathbf{P'}\ X$

# Triangulation



Find 3D object point

image 1

$x$

$C$

camera 1 with matrix $\mathbf{P}$

image 2

$x'$

$C'$

camera 2 with matrix $\mathbf{P}'$

# Two-view SfM

1. Compute the Fundamental Matrix **F** from points correspondences
   **8-point algorithm**

2. Compute the camera matrices **P** from the Fundamental matrix
   $\mathbf{P} = [\ \mathbf{I}\ |\ \mathbf{0}\ ]$  and  $\mathbf{P'} = [\ [\mathbf{e'}_x]\mathbf{F}\ |\ \mathbf{e'}\ ]$

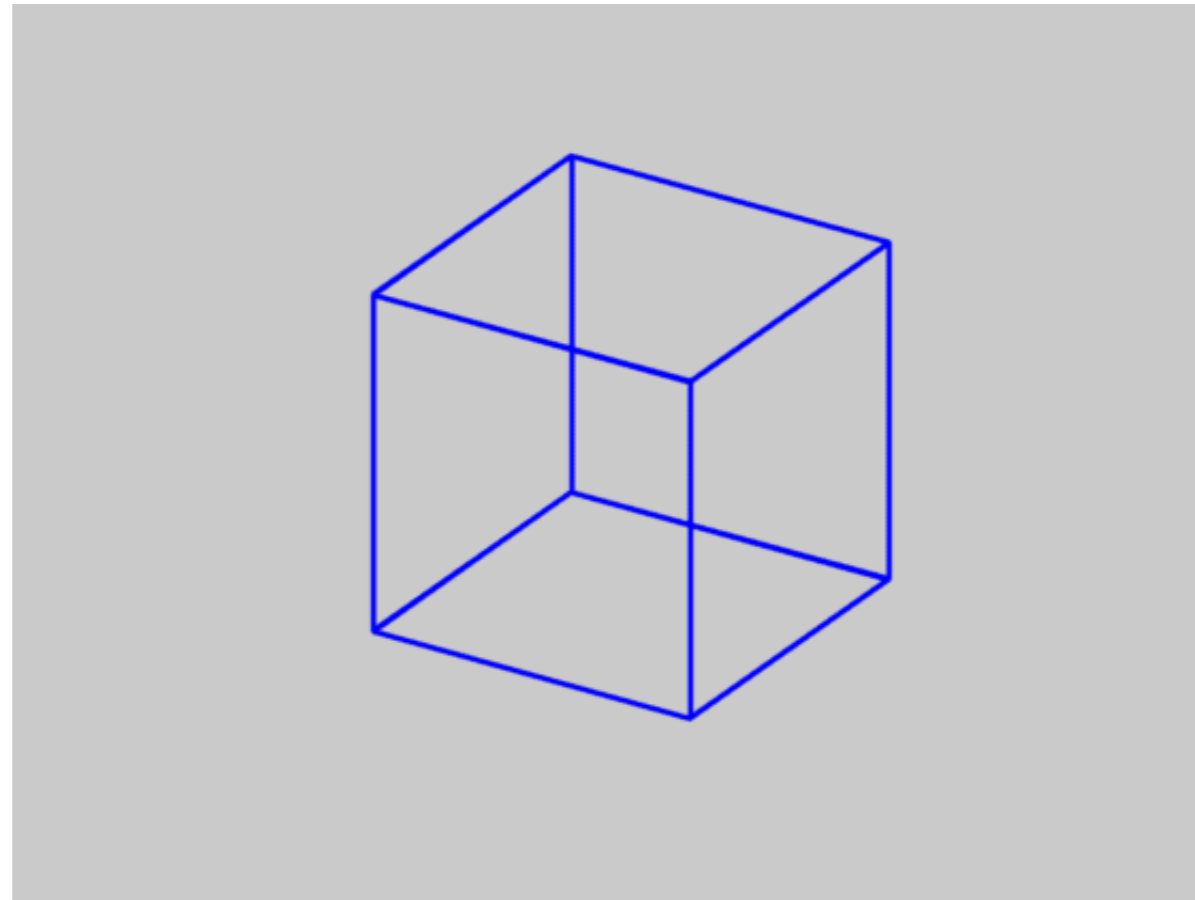3. For each point correspondence, compute the point **X** in 3D space (triangularization)
   **DLT** with $x = \mathbf{P}\,X$ and $x' = \mathbf{P'}\,X$

# Is SfM always uniquely solvable?
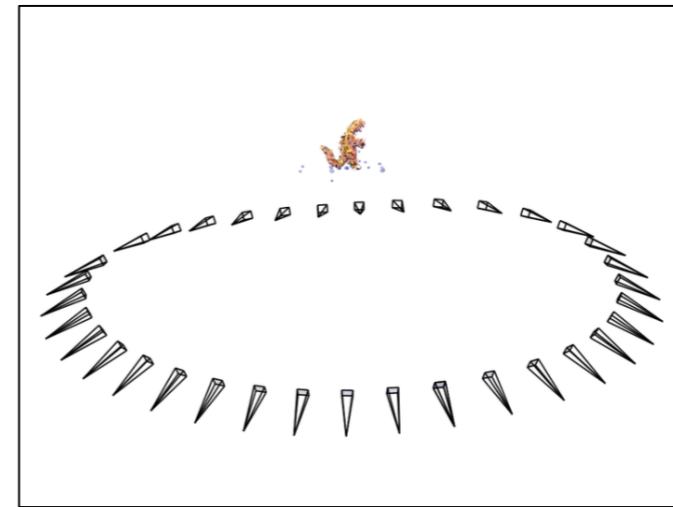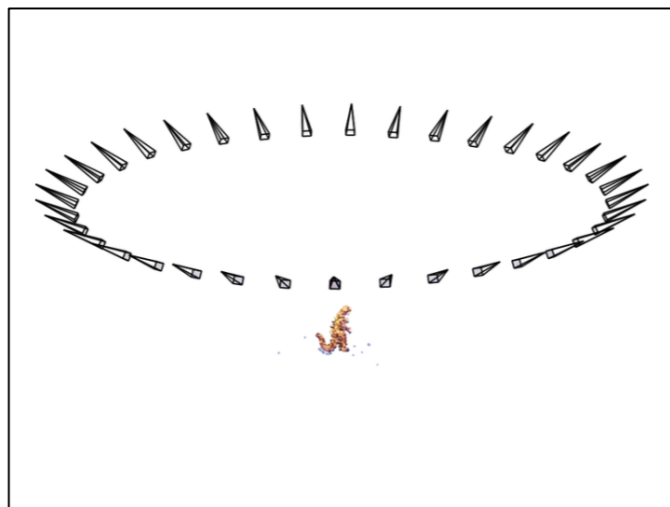
# Ambiguities in structure from motion

# Is SfM always uniquely solvable?

- No…

# SfM – Failure cases

- Necker reversal

# Projective Ambiguity

- Reconstruction is ambiguous by an arbitrary 3D projective transformation without prior knowledge of camera parameters
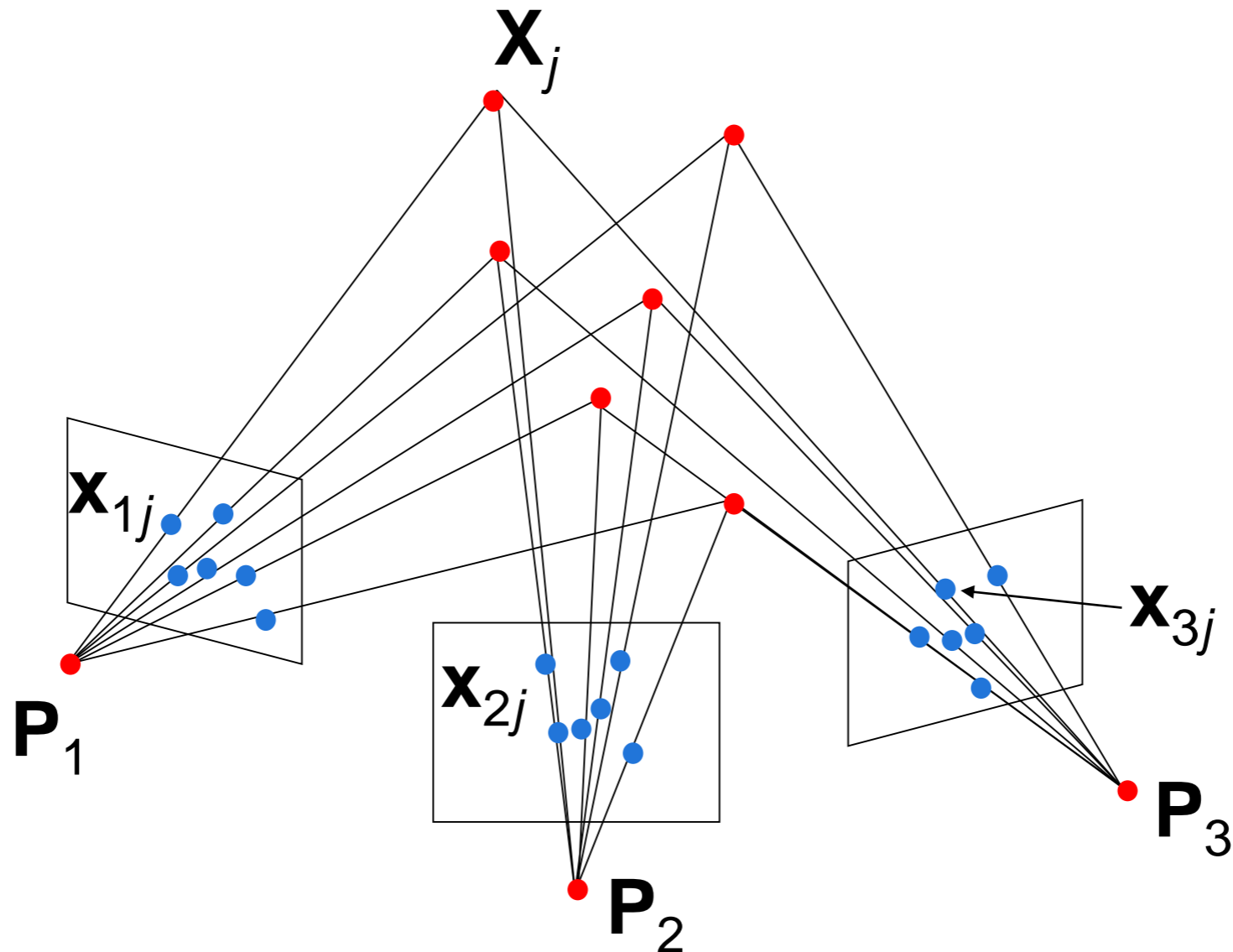
# Structure from motion

- Given: $m$ images of $n$ fixed 3D points

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

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

# 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 in the image remain exactly the same:

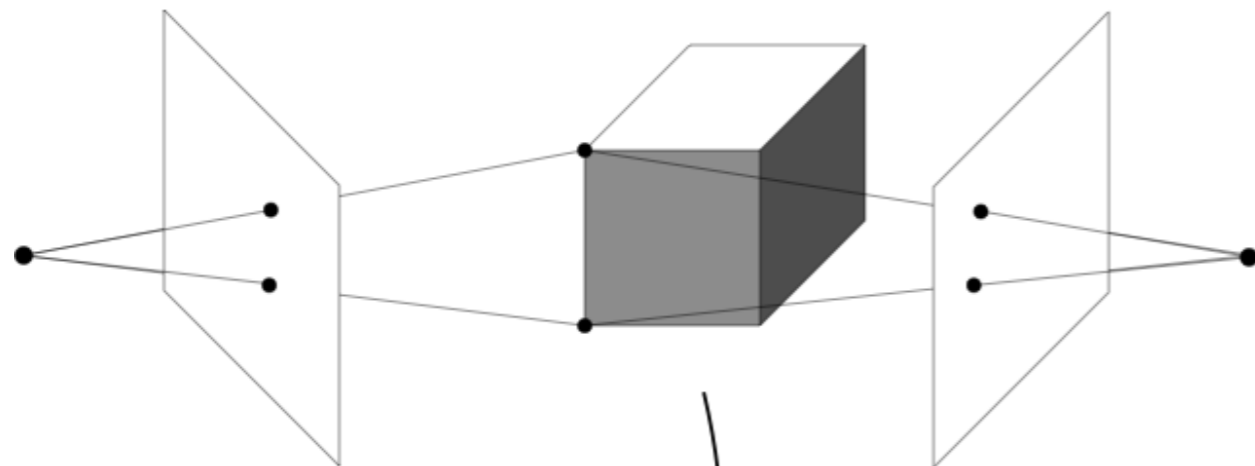$$\mathbf{x} = \mathbf{PX} = \left(\frac{1}{k}\mathbf{P}\right)(k\,\mathbf{X})$$

It is impossible to recover the absolute scale of the scene!

# Structure from motion ambiguity
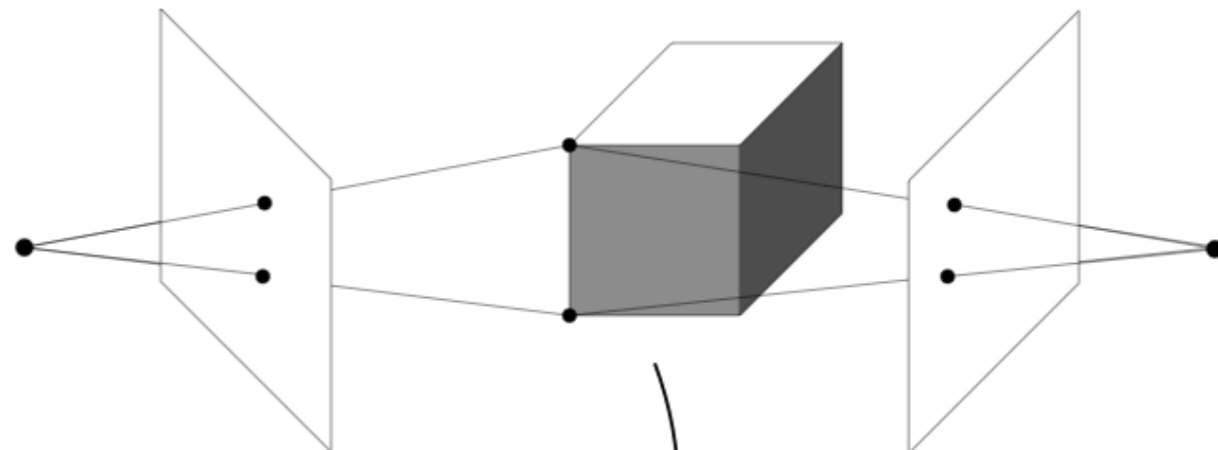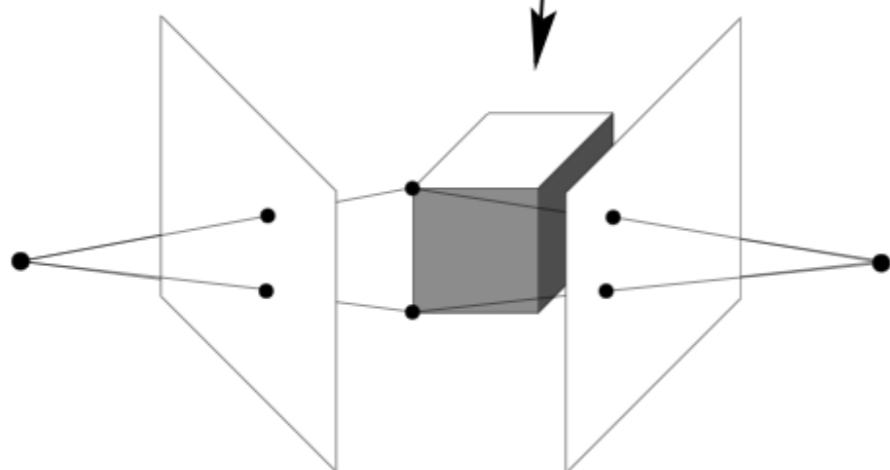
- If we scale the entire scene by some factor $k$ and, at the same time, scale the camera matrices by the factor of $1/k$, the projections of the scene points in the image remain exactly the same

- More generally: if we transform the scene using a transformation $\mathbf{Q}$ and apply the inverse transformation to the camera matrices, then the images do not change

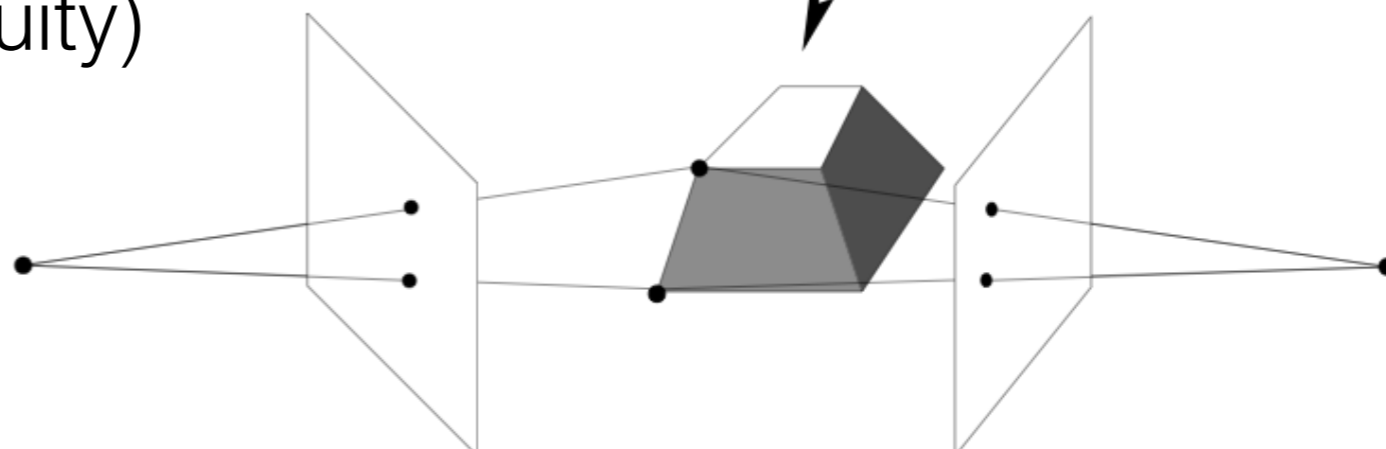$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ^{-1}}\right)\left(\mathbf{QX}\right)$$

**Calibrated cameras**
(similarity projection ambiguity)

Similarity

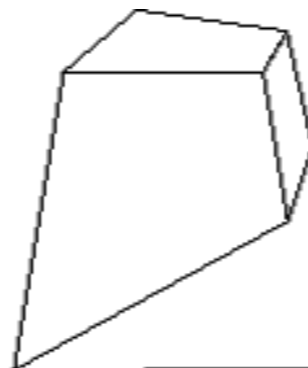**Uncalibrated cameras**
(projective projection ambiguity)

Projective

# Types of ambiguity

Projective
15dof
$$\begin{bmatrix} A & t \\ v^\top & v \end{bmatrix}$$
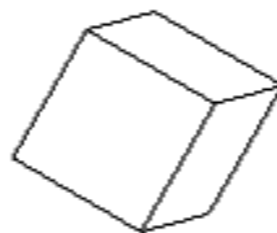Preserves intersection and tangency

Affine
12dof
$$\begin{bmatrix} A & t \\ 0^\top & 1 \end{bmatrix}$$
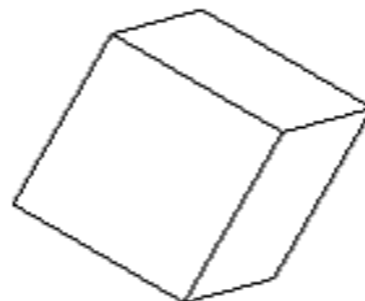Preserves parallellism, volume ratios

Similarity
7dof
$$\begin{bmatrix} s\,R & t \\ 0^\top & 1 \end{bmatrix}$$
Preserves angles, ratios of length

Euclidean
6dof
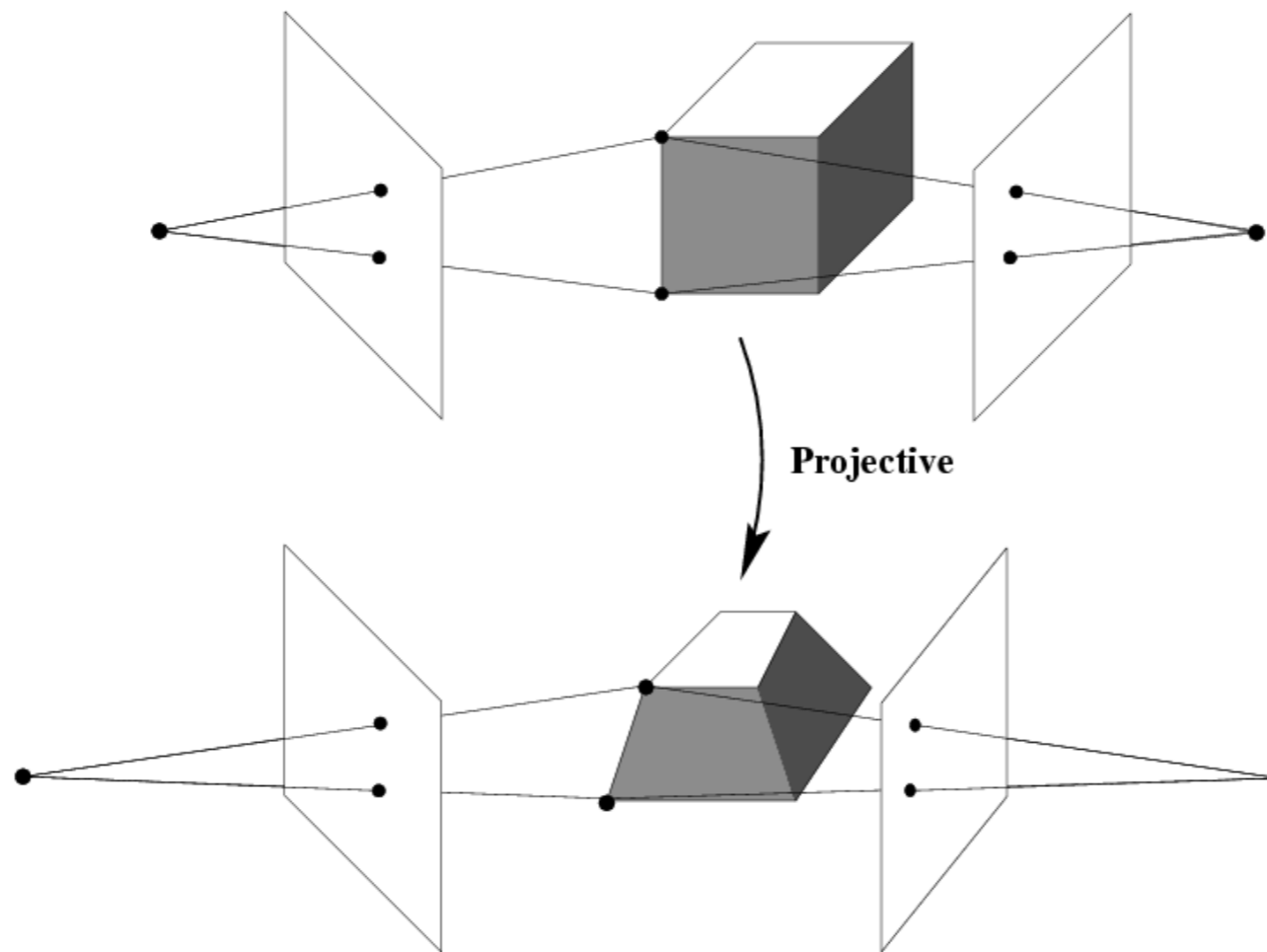$$\begin{bmatrix} R & t \\ 0^\top & 1 \end{bmatrix}$$
Preserves angles, lengths

- With no constraints on the camera calibration matrix or on the scene, we get a *projective* reconstruction
- Need additional information to *upgrade* the reconstruction to affine, similarity, or Euclidean
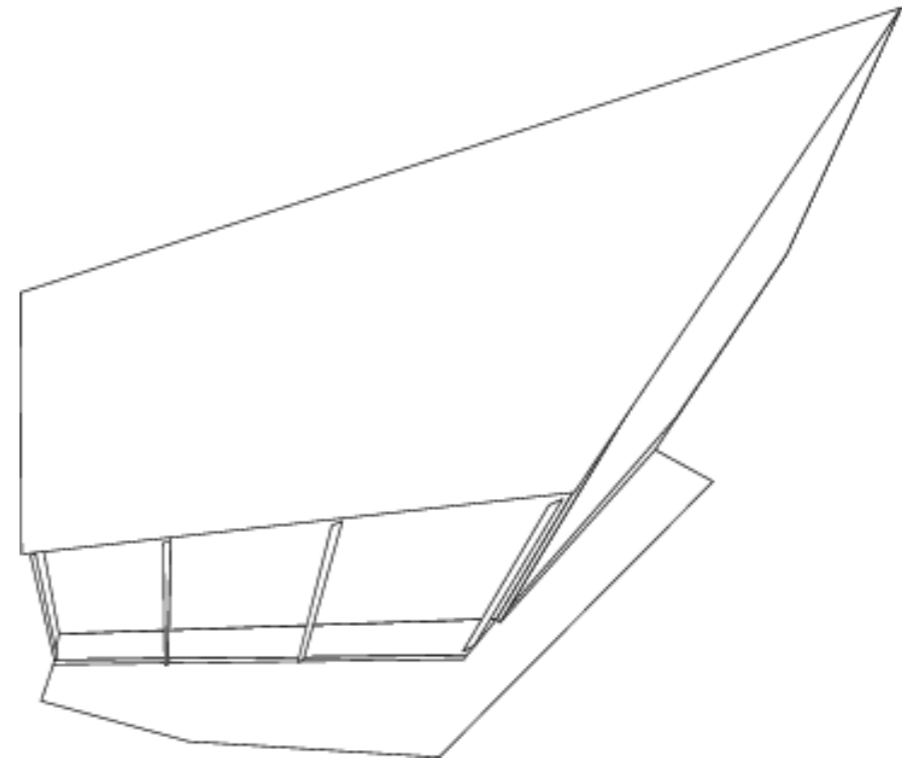
# Projective ambiguity



Projective

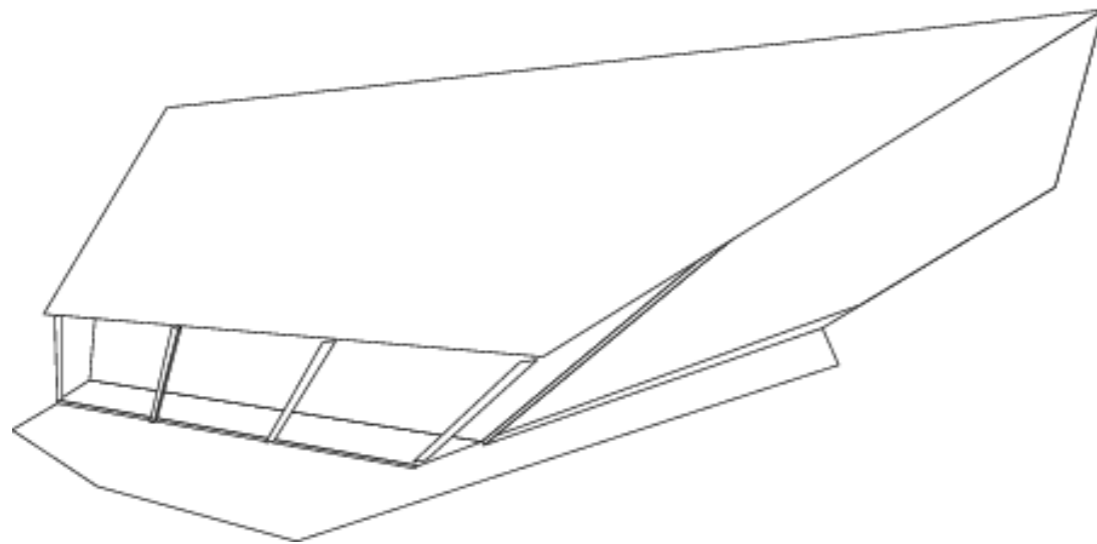$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ_P^{-1}}\right)\left(\mathbf{Q_P\,X}\right)$$

# Projective ambiguity

# Affine ambiguity



$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ_A^{-1}}\right)\left(\mathbf{Q_A\, X}\right)$$

# Affine ambiguity

# Similarity ambiguity



$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ}_S^{-1}\right)\left(\mathbf{Q}_S\mathbf{X}\right)$$

# Similarity ambiguity

# What can we do to remove ambiguities?

# Affine structure from motion

# Structure from motion

- Let's start with *affine cameras* (the math is easier)



perspective            weak perspective

center at infinity

increasing focal length →

increasing distance from camera →

# Recall: Orthographic Projection

## Special case of perspective projection

- Distance from center of projection to image plane is infinite

Image                    World

- Projection matrix:

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \Rightarrow (x, y)$$

# Affine cameras

Orthographic Projection

Parallel Projection

# Affine cameras

- A general affine camera combines the effects of an affine transformation of the 3D space, orthographic projection, and an affine transformation of the image:

$$\mathbf{P} = [3 \times 3 \text{ affine}] \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} [4 \times 4 \text{ affine}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & \mathbf{1} \end{bmatrix}$$
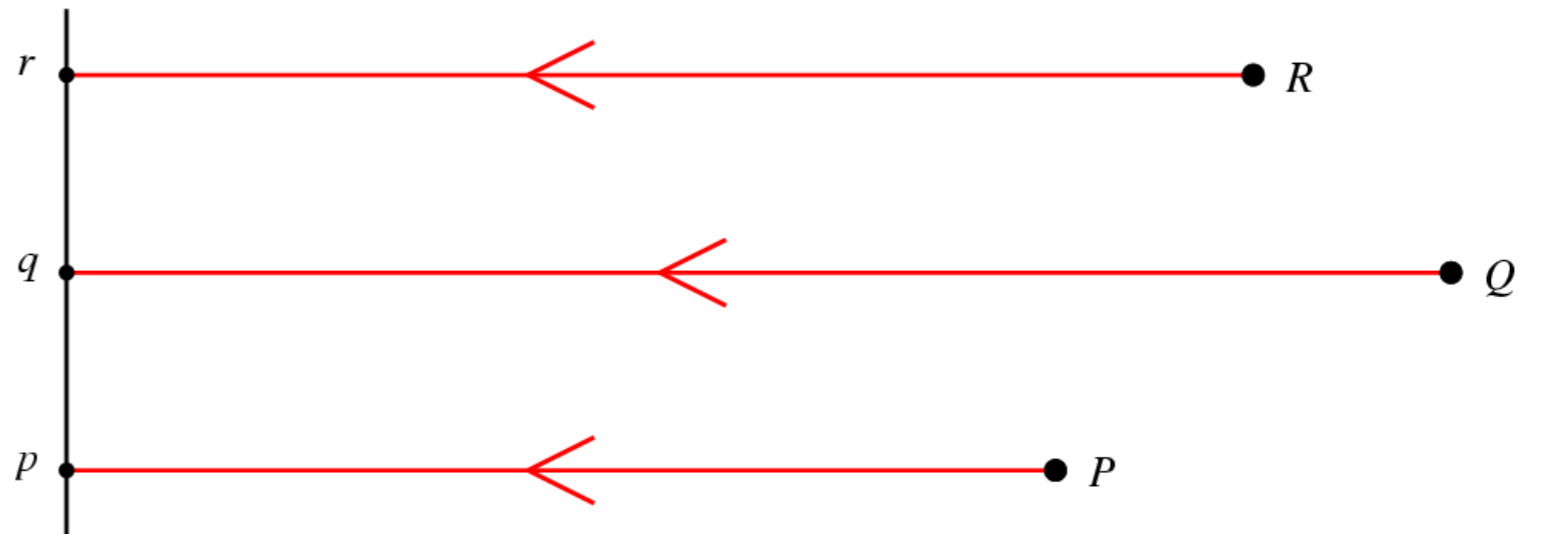
- Affine projection is a linear mapping + translation in inhomogeneous coordinates



$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = \mathbf{AX} + \mathbf{b}$$

$\mathbf{a}_2$

$\mathbf{a}_1$

$\mathbf{x}$

$\mathbf{X}$

Projection of world origin

# Affine structure from motion

- Given: $m$ images of $n$ fixed 3D points:

$$\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i, \quad i = 1, \dots, m, \ j = 1, \dots, n$$

- Problem: use the $mn$ correspondences $\mathbf{x}_{ij}$ to estimate $m$ projection matrices $\mathbf{A}_i$ and translation vectors $\mathbf{b}_i$, and $n$ points $\mathbf{X}_j$

- The reconstruction is defined up to an arbitrary *affine* transformation $\mathbf{Q}$ (12 degrees of freedom):

$$\begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \mathbf{Q}^{-1}, \qquad \begin{pmatrix} \mathbf{X} \\ \mathbf{1} \end{pmatrix} \rightarrow \mathbf{Q} \begin{pmatrix} \mathbf{X} \\ \mathbf{1} \end{pmatrix}$$

- We have $2mn$ knowns and $8m + 3n$ unknowns (minus 12 dof for affine ambiguity)

- Thus, we must have $2mn >= 8m + 3n - 12$

- For two views, we need four point correspondences

# Affine structure from motion

- Centering: subtract the centroid of the image points

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

$$= \mathbf{A}_i\left(\mathbf{X}_j - \frac{1}{n}\sum_{k=1}^{n}\mathbf{X}_k\right) = \mathbf{A}_i\hat{\mathbf{X}}_j$$

- For simplicity, assume that the origin of the world coordinate system is at the centroid of the 3D points

- After centering, each normalized point $\mathbf{x}_{ij}$ is related to the 3D point $\mathbf{X}_i$ by

$$\hat{\mathbf{x}}_{ij} = \mathbf{A}_i\mathbf{X}_j$$

# Affine structure from motion

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

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

cameras ($2m$)

points ($n$)

C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization method. *IJCV*, 9(2):137-154, November 1992.

# Affine structure from motion

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

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ & & \ddots & \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}$$
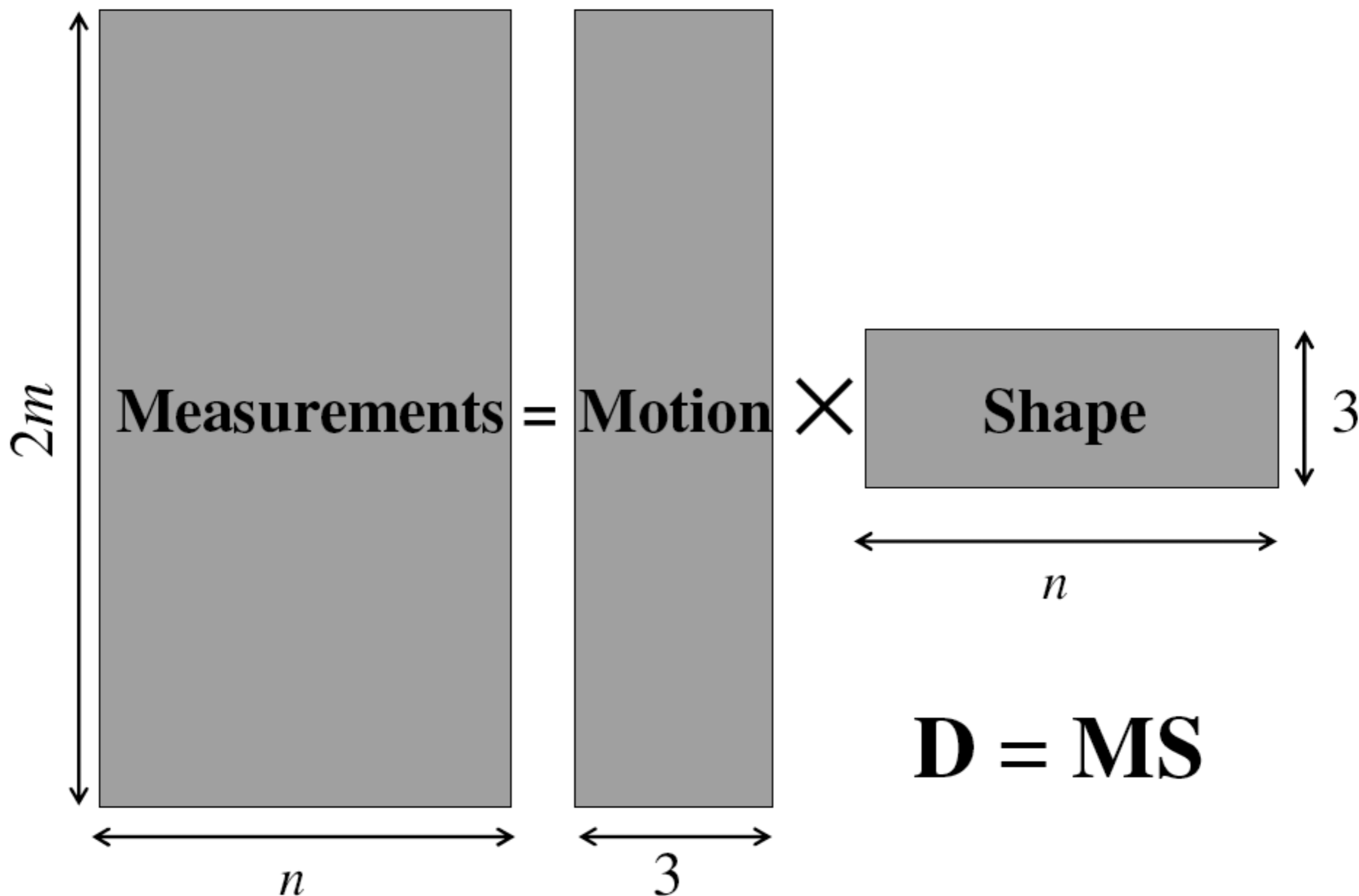
<span style="color:red">points ($3 \times n$)</span>

<span style="color:red">cameras ($2m \times 3$)</span>

## The measurement matrix $\mathbf{D} = \mathbf{MS}$ must have rank 3!

C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization method. *IJCV*, 9(2):137-154, November 1992.

# Factorizing the measurement matrix



$$\textbf{Measurements} = \textbf{Motion} \times \textbf{Shape}$$

$2m$, $n$, $3$, $n$, $3$

$$\textbf{D} = \textbf{MS}$$

# Factorizing the measurement matrix

- Singular value decomposition of D:
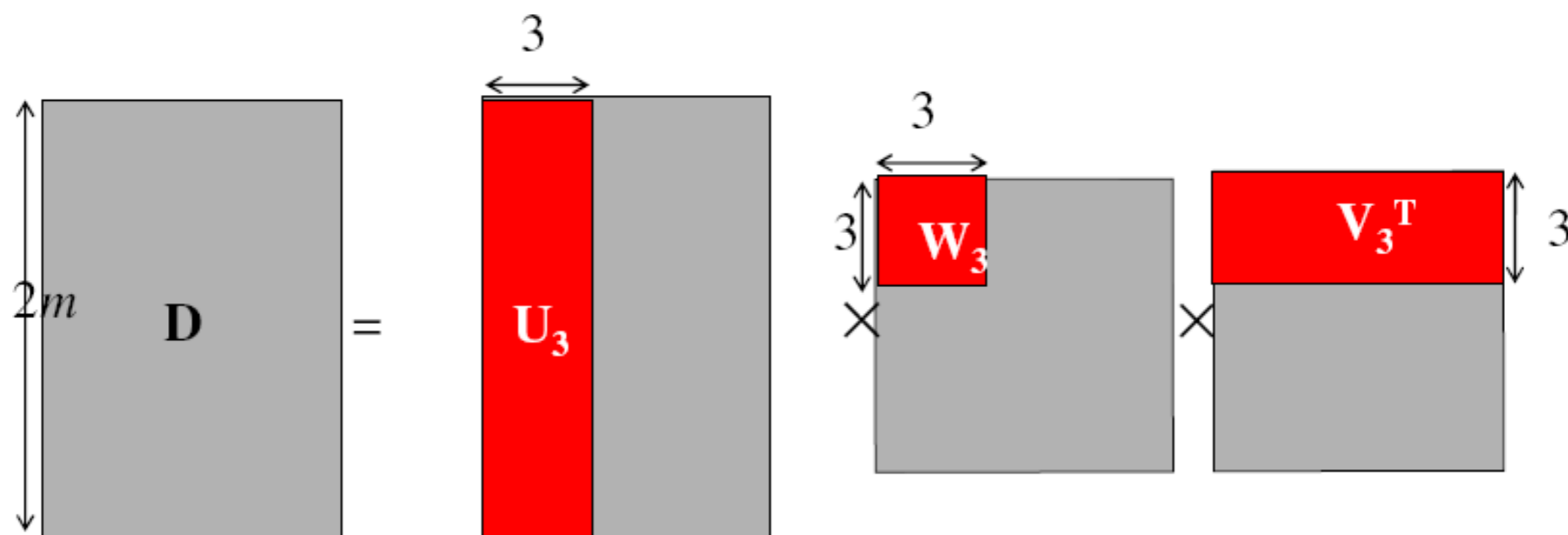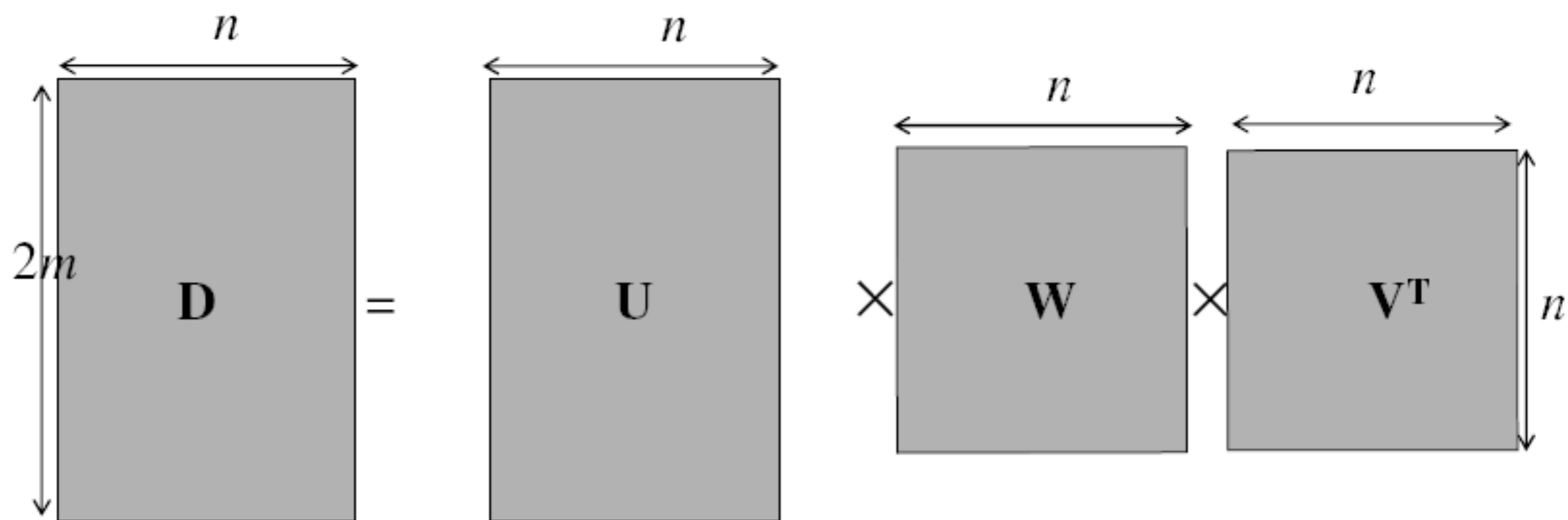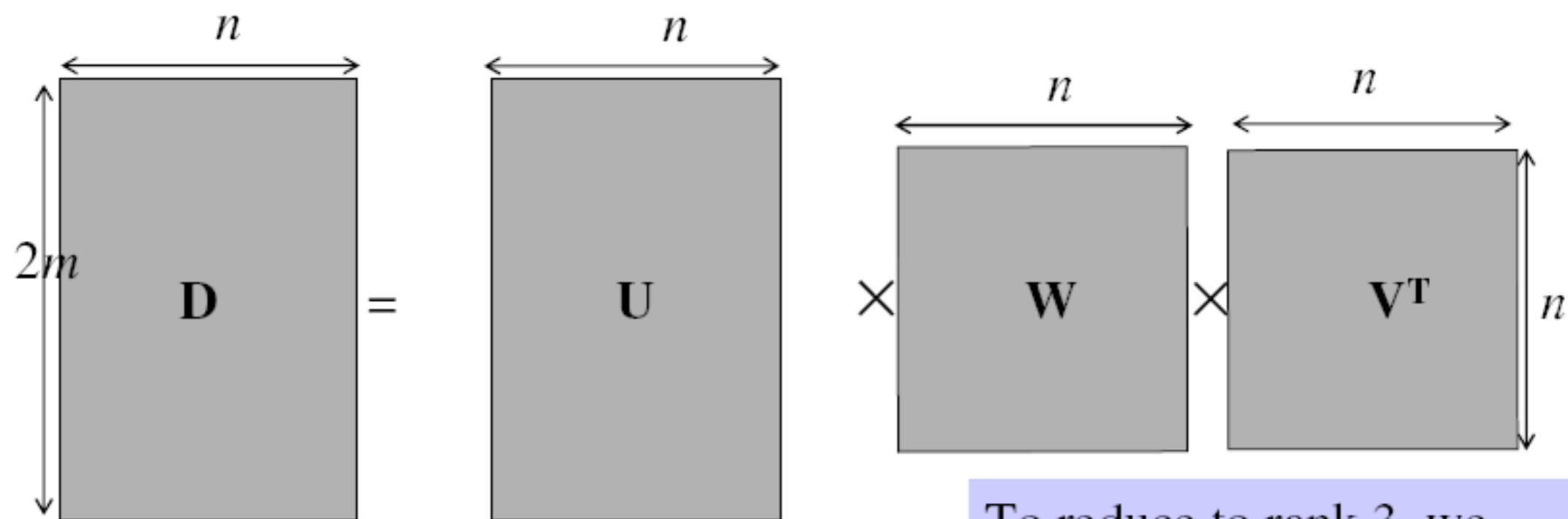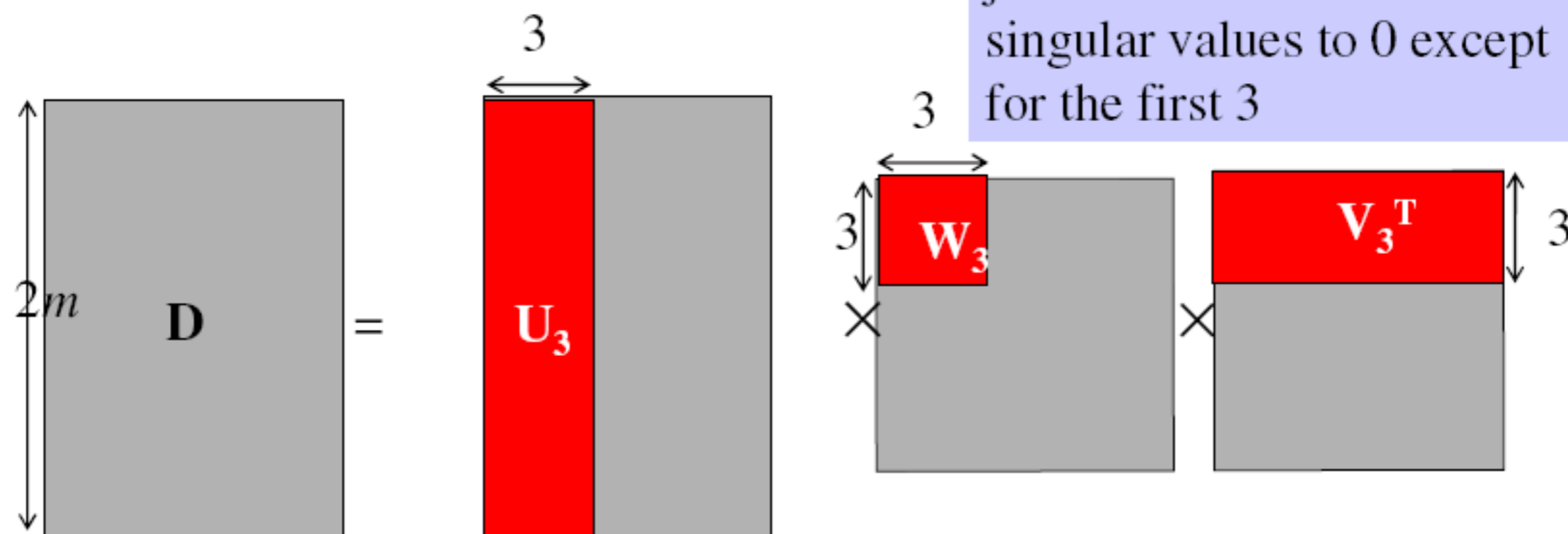
# Factorizing the measurement matrix

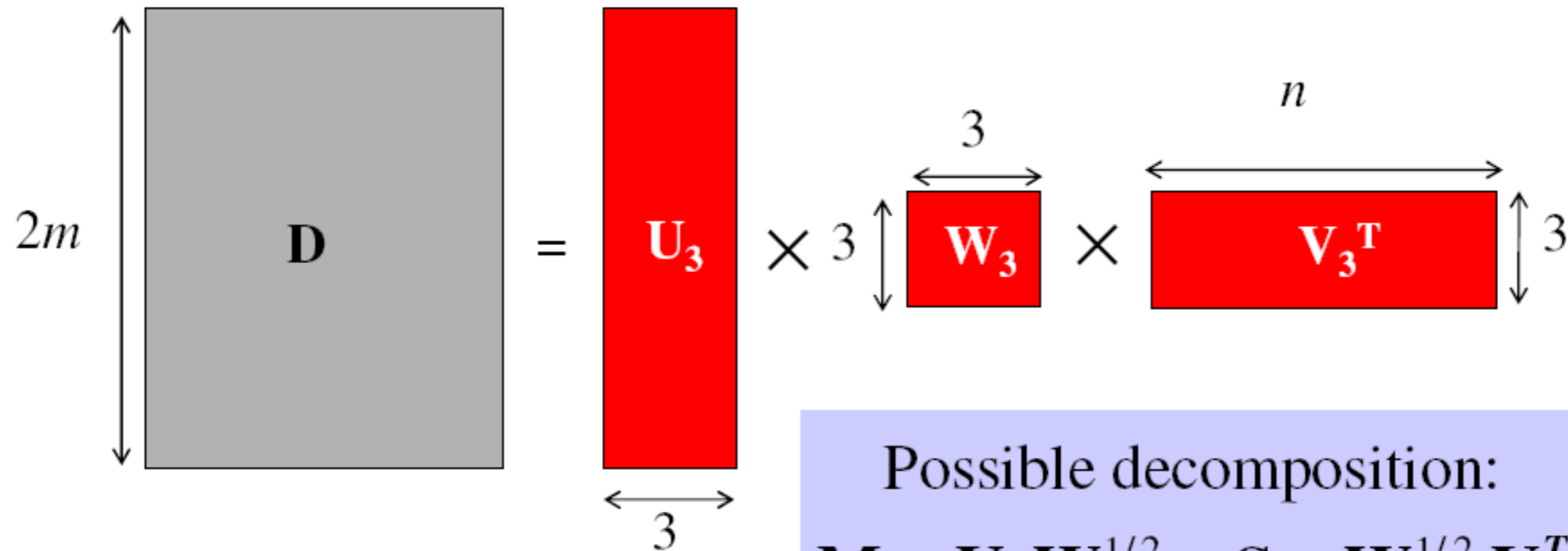- Singular value decomposition of D:



To reduce to rank 3, we just need to set all the singular values to 0 except for the first 3
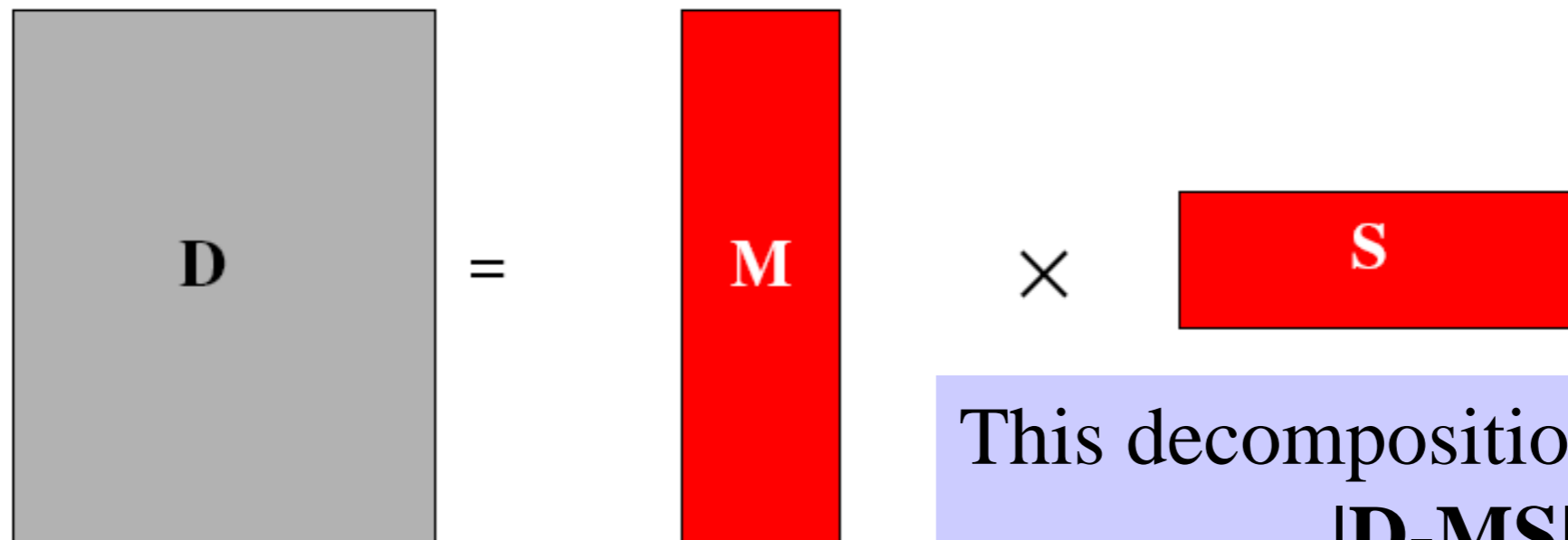
# Factorizing the measurement matrix

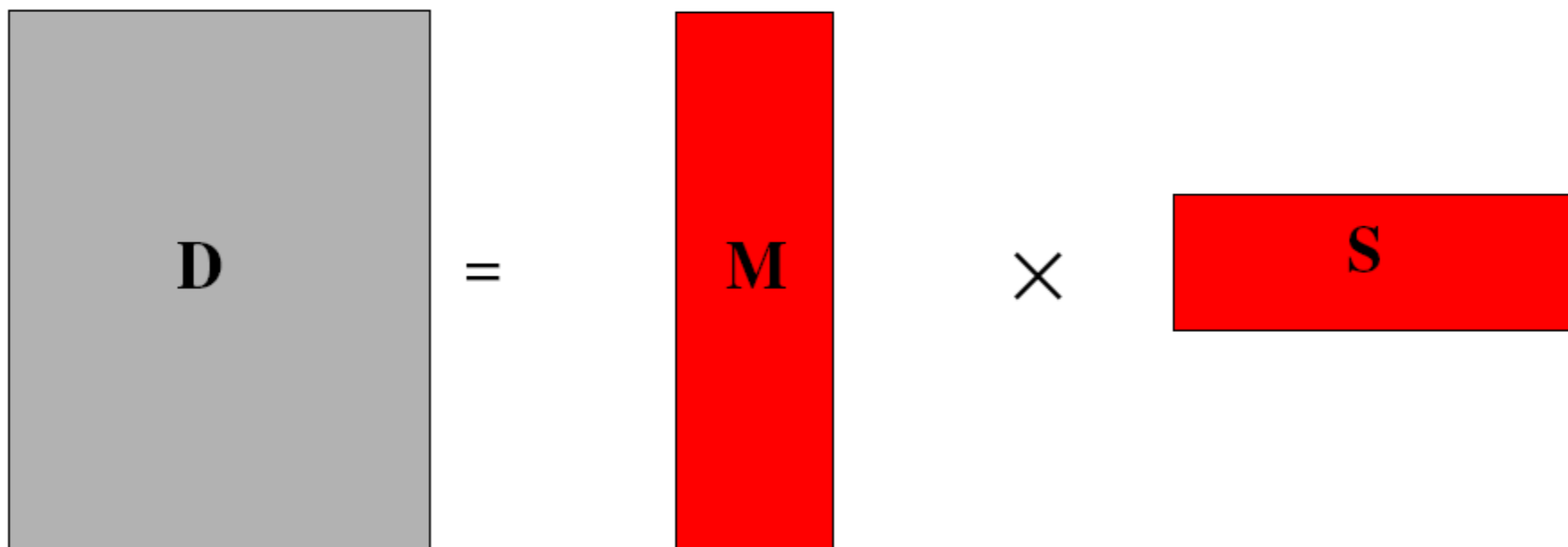- Obtaining a factorization from SVD:



Possible decomposition:

$$\mathbf{M} = \mathbf{U}_3 \mathbf{W}_3^{1/2} \quad \mathbf{S} = \mathbf{W}_3^{1/2} \mathbf{V}_3^{T}$$

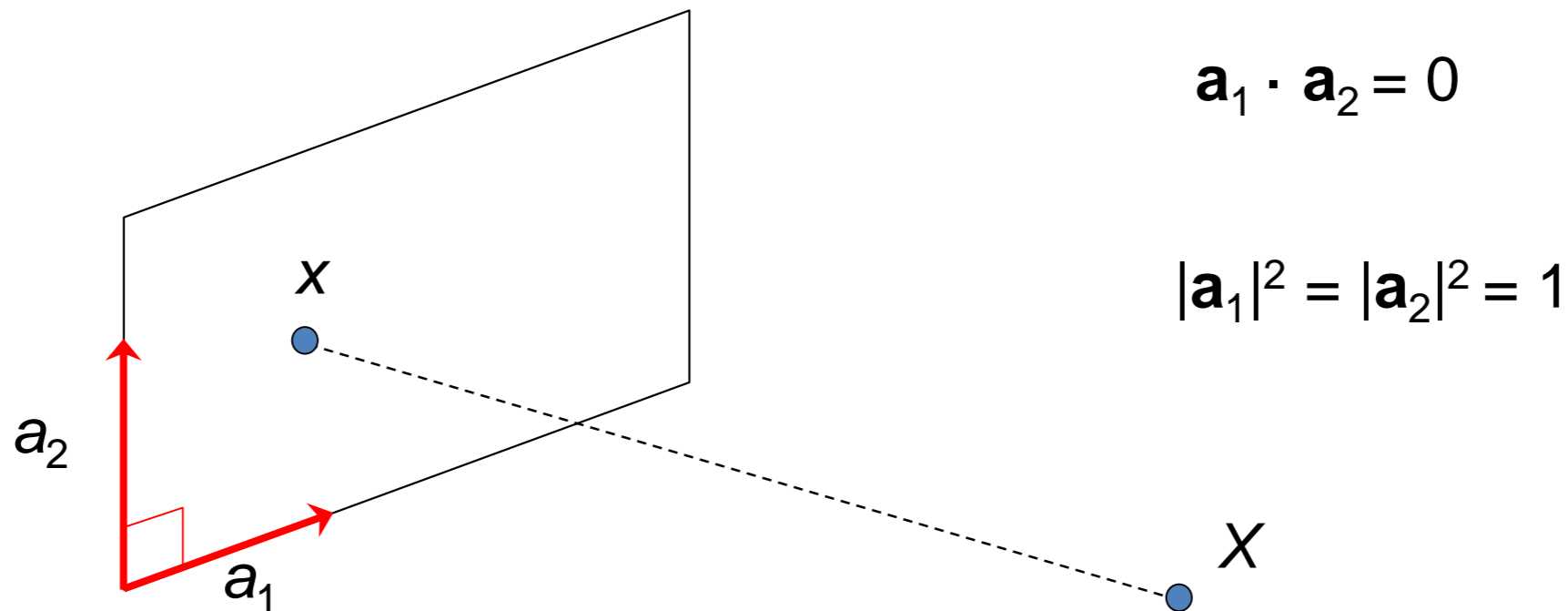This decomposition minimizes $|\mathbf{D}-\mathbf{MS}|^2$

# Affine ambiguity



$$\mathbf{D} = \mathbf{M} \times \mathbf{S}$$

- The decomposition is not unique. We get the same **D** by using any 3×3 matrix **C** and applying the transformations $\mathbf{M} \rightarrow \mathbf{MC}$, $\mathbf{S} \rightarrow \mathbf{C^{-1}S}$

- That is because we have only an affine transformation and we have not enforced any Euclidean constraints (like forcing the image axes to be perpendicular, for example)

# Eliminating the affine ambiguity

- Orthographic: image axes are perpendicular and of unit length

$$\mathbf{a}_1 \cdot \mathbf{a}_2 = 0$$

$$|\mathbf{a}_1|^2 = |\mathbf{a}_2|^2 = 1$$

*x*

*a*$_2$

*a*$_1$

*X*

# Solve for orthographic constraints

Three equations for each image i

$$\tilde{\mathbf{a}}_{i1}^T \mathbf{C} \mathbf{C}^T \tilde{\mathbf{a}}_{i1}^T = 1$$

$$\tilde{\mathbf{a}}_{i2}^T \mathbf{C} \mathbf{C}^T \tilde{\mathbf{a}}_{i2}^T = 1 \quad \text{where} \quad \tilde{\mathbf{A}}_i = \begin{bmatrix} \tilde{\mathbf{a}}_{i1}^T \\ \tilde{\mathbf{a}}_{i2}^T \end{bmatrix}$$

$$\tilde{\mathbf{a}}_{i1}^T \mathbf{C} \mathbf{C}^T \tilde{\mathbf{a}}_{i2}^T = 0$$

- Solve for $\mathbf{L} = \mathbf{CC}^T$

- Recover $\mathbf{C}$ from $\mathbf{L}$ by Cholesky decomposition: $\mathbf{L} = \mathbf{CC}^T$

- Update $\mathbf{A}$ and $\mathbf{X}$: $\mathbf{A} = \tilde{\mathbf{A}}\mathbf{C}, \mathbf{X} = \mathbf{C}^{-1}\tilde{\mathbf{X}}$

# Algorithm summary

- Given: $m$ images and $n$ features $\mathbf{x}_{ij}$
- For each image $i$, center the feature coordinates
- Construct a $2m \times n$ measurement matrix $\mathbf{D}$:
  - Column $j$ contains the projection of point $j$ in all views
  - Row $i$ contains one coordinate of the projections of all the $n$ points in image $i$
- Factorize $\mathbf{D}$:
  - Compute SVD: $\mathbf{D} = \mathbf{U} \, \mathbf{W} \, \mathbf{V}^{\mathsf{T}}$
  - Create $\mathbf{U}_3$ by taking the first 3 columns of $\mathbf{U}$
  - Create $\mathbf{V}_3$ by taking the first 3 columns of $\mathbf{V}$
  - Create $\mathbf{W}_3$ by taking the upper left $3 \times 3$ block of $\mathbf{W}$
- Create the motion and shape matrices:
  - $\mathbf{M} = \mathbf{U}_3 \mathbf{W}_3^{\frac{1}{2}}$ and $\mathbf{S} = \mathbf{W}_3^{\frac{1}{2}} \mathbf{V}_3^{\mathsf{T}}$ (**or** $\mathbf{M} = \mathbf{U}_3$ and $\mathbf{S} = \mathbf{W}_3 \mathbf{V}_3^{\mathsf{T}}$)
- Eliminate affine ambiguity

# Reconstruction results



C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization method. *IJCV*, 9(2):137-154, November 1992.
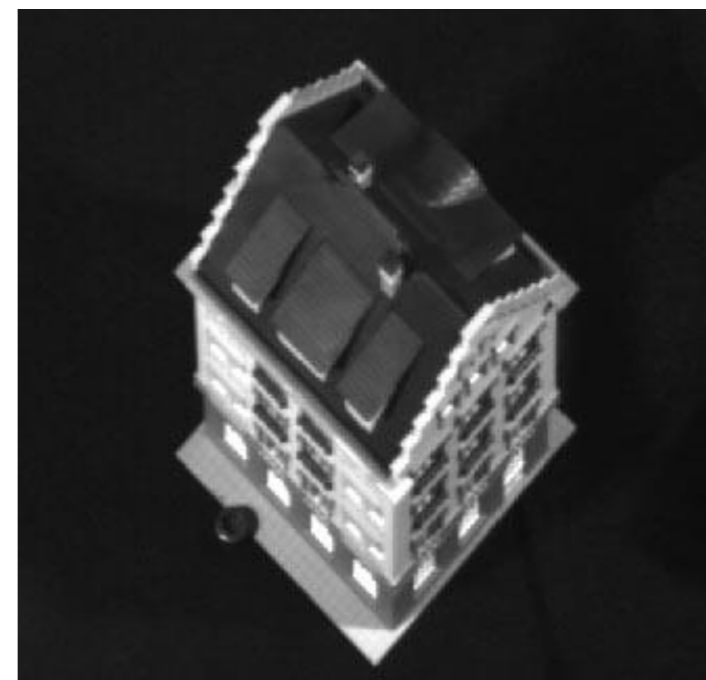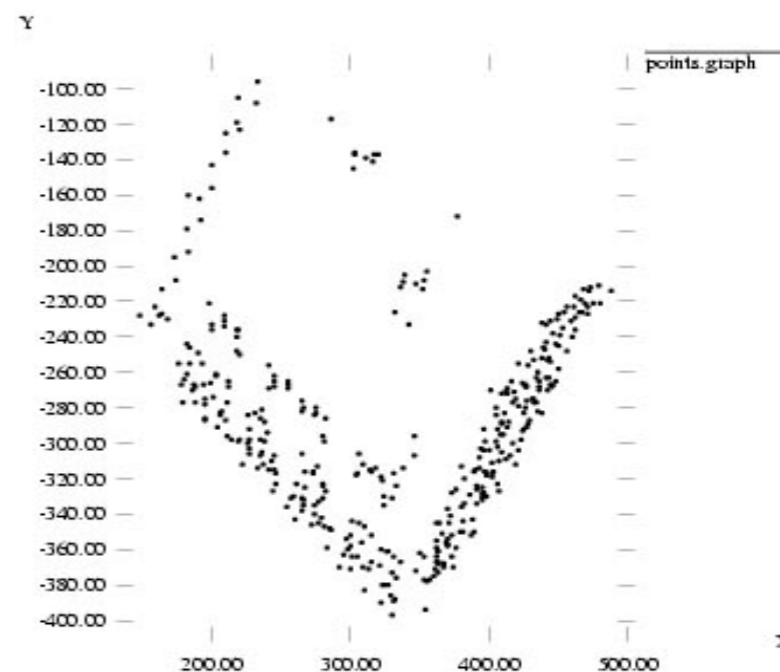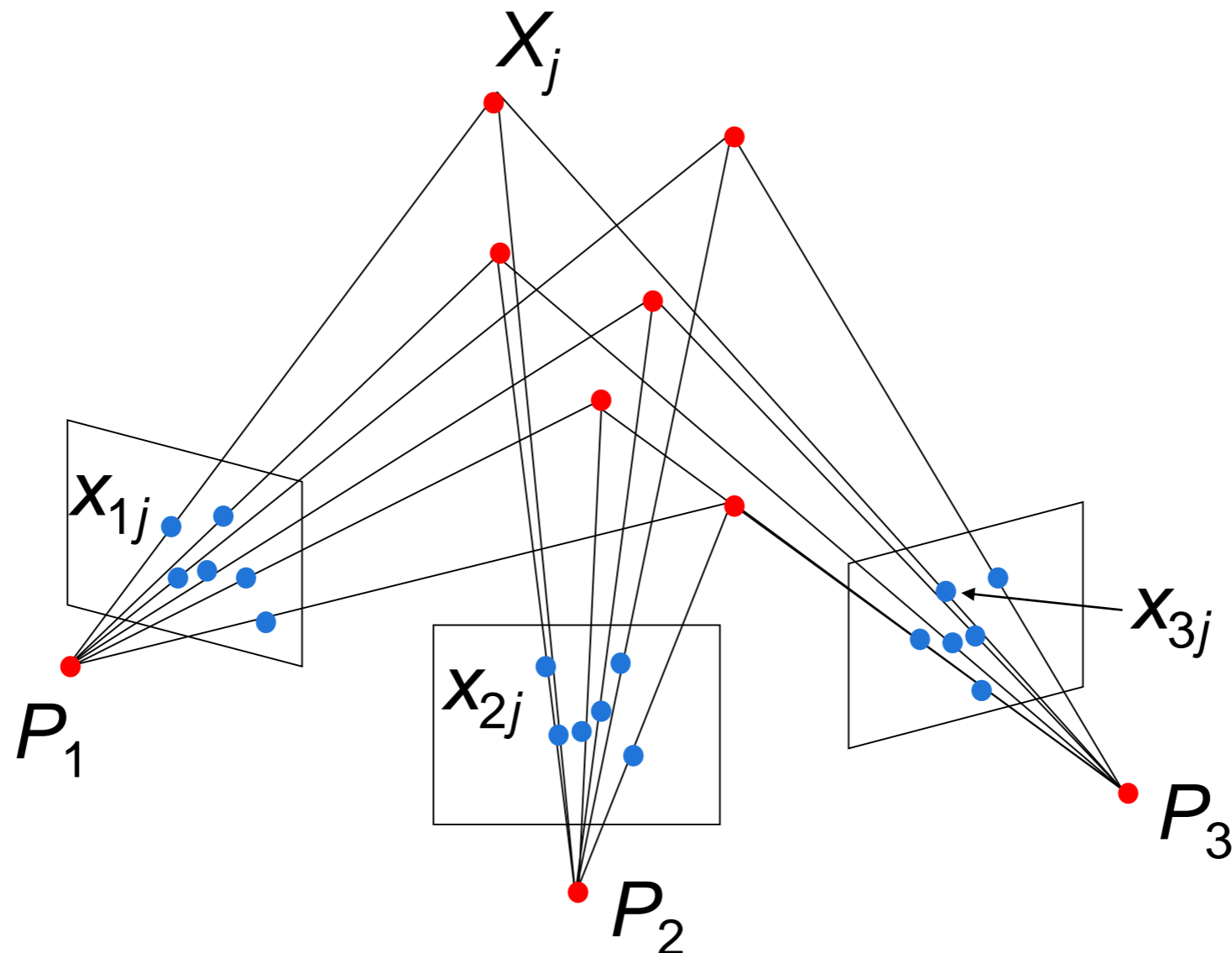
# Multi-view projective structure from motion

# Projective structure from motion

- Given: $m$ images of $n$ fixed 3D points

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

- 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

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

- 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 4x4 projective transformation $\mathbf{Q}$:

$$\mathbf{X} \rightarrow \mathbf{Q}\mathbf{X}, \ \mathbf{P} \rightarrow \mathbf{P}\mathbf{Q}^{-1}$$

- We can solve for structure and motion when

$$2mn >= 11m + 3n - 15$$

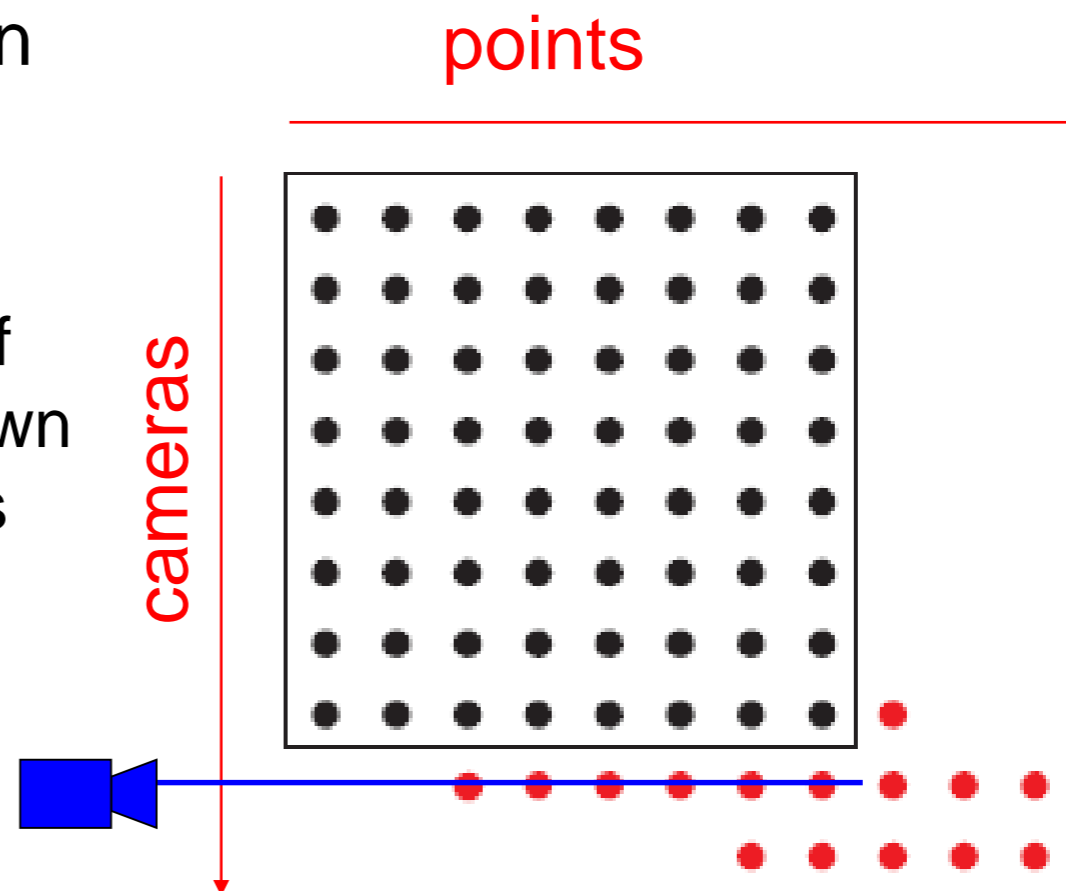- For two cameras, at least 7 points are needed

# Projective SFM: Two-camera case

- Compute fundamental matrix $\mathbf{F}$ between the two views
- First camera matrix: $[\mathbf{I}|\mathbf{0}]$
- Second camera matrix: $[\mathbf{A}|\mathbf{b}]$
- Then $\mathbf{b}$ is the epipole $(\mathbf{F}^T\mathbf{b} = 0)$, $\mathbf{A} = -[\mathbf{b}_\times]\mathbf{F}$
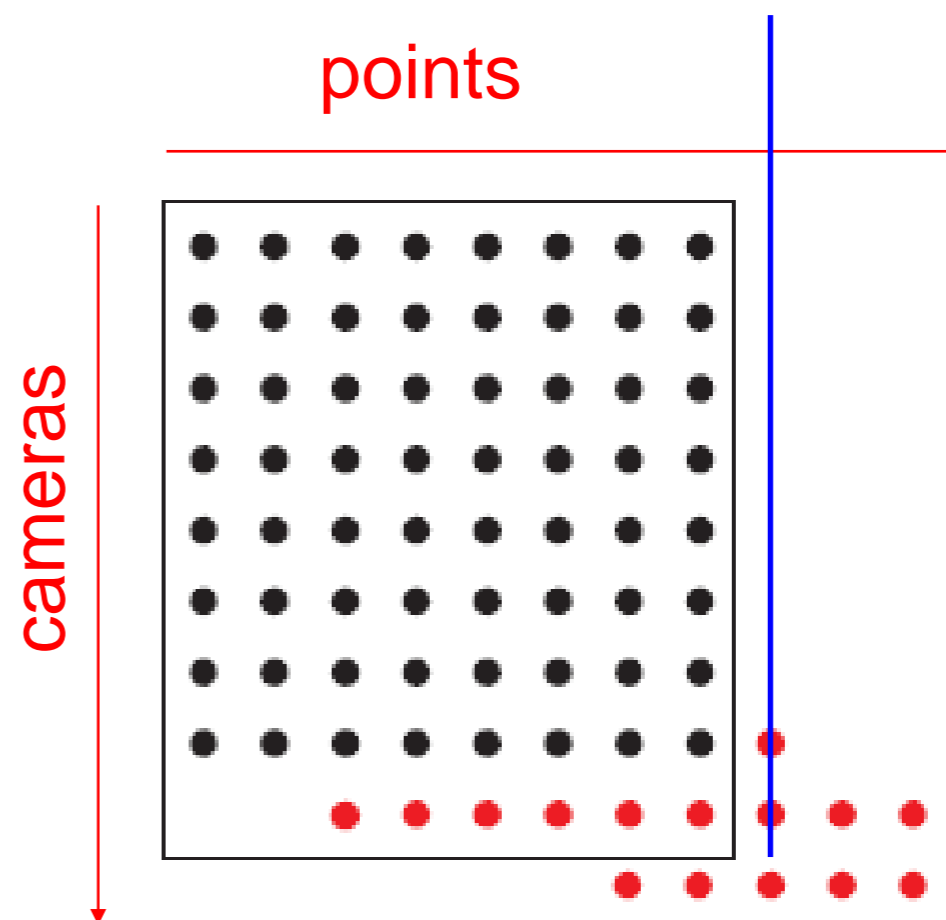
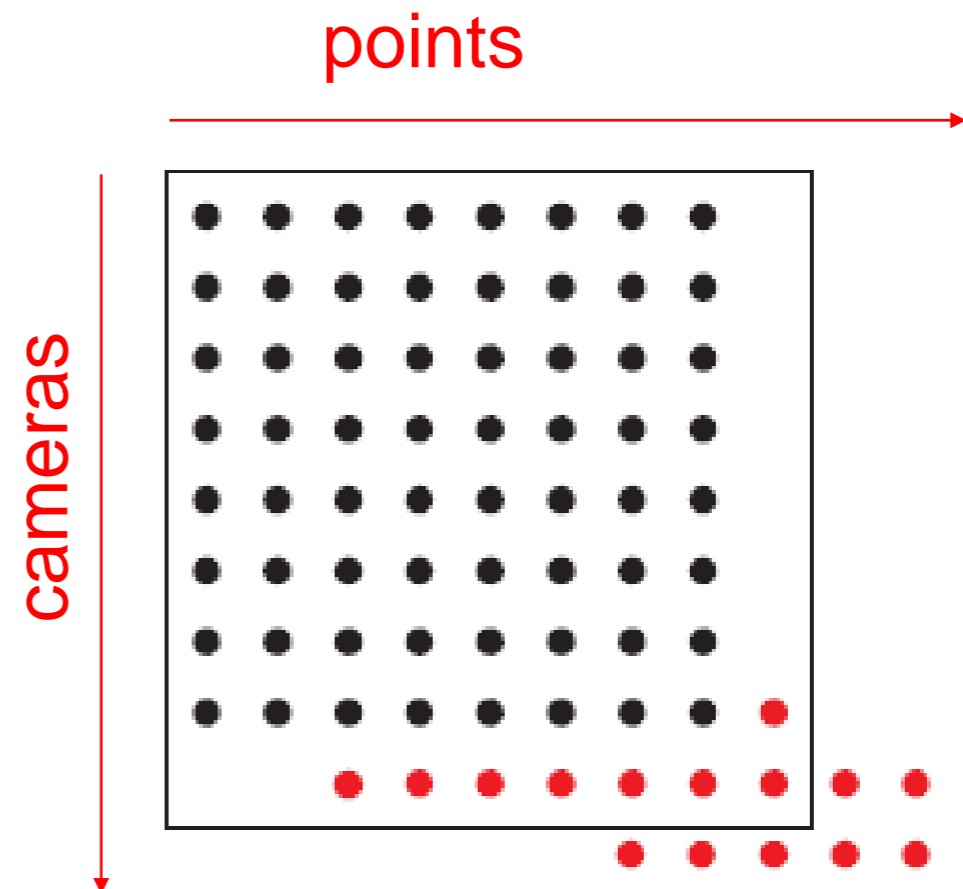# Sequential structure from motion

- Initialize motion from two images using fundamental matrix

- Initialize structure by triangulation

- For each additional view:
  - Determine projection matrix of new camera using all the known 3D points that are visible in its image – *calibration*

points

cameras

# Sequential structure from motion

- Initialize motion from two images using fundamental matrix

- Initialize structure by triangulation

- For each additional view:

  - Determine projection matrix of new camera using all the known 3D points that are visible in its image – *calibration*

  - Refine and extend structure: compute new 3D points, re-optimize existing points that are also seen by this camera – *triangulation*

points

cameras

# Sequential structure from motion

- Initialize motion from two images using fundamental matrix

- Initialize structure by triangulation

- For each additional view:
  - Determine projection matrix of new camera using all the known 3D points that are visible in its image – *calibration*

  - Refine and extend structure: compute new 3D points, re-optimize existing points that are also seen by this camera – *triangulation*
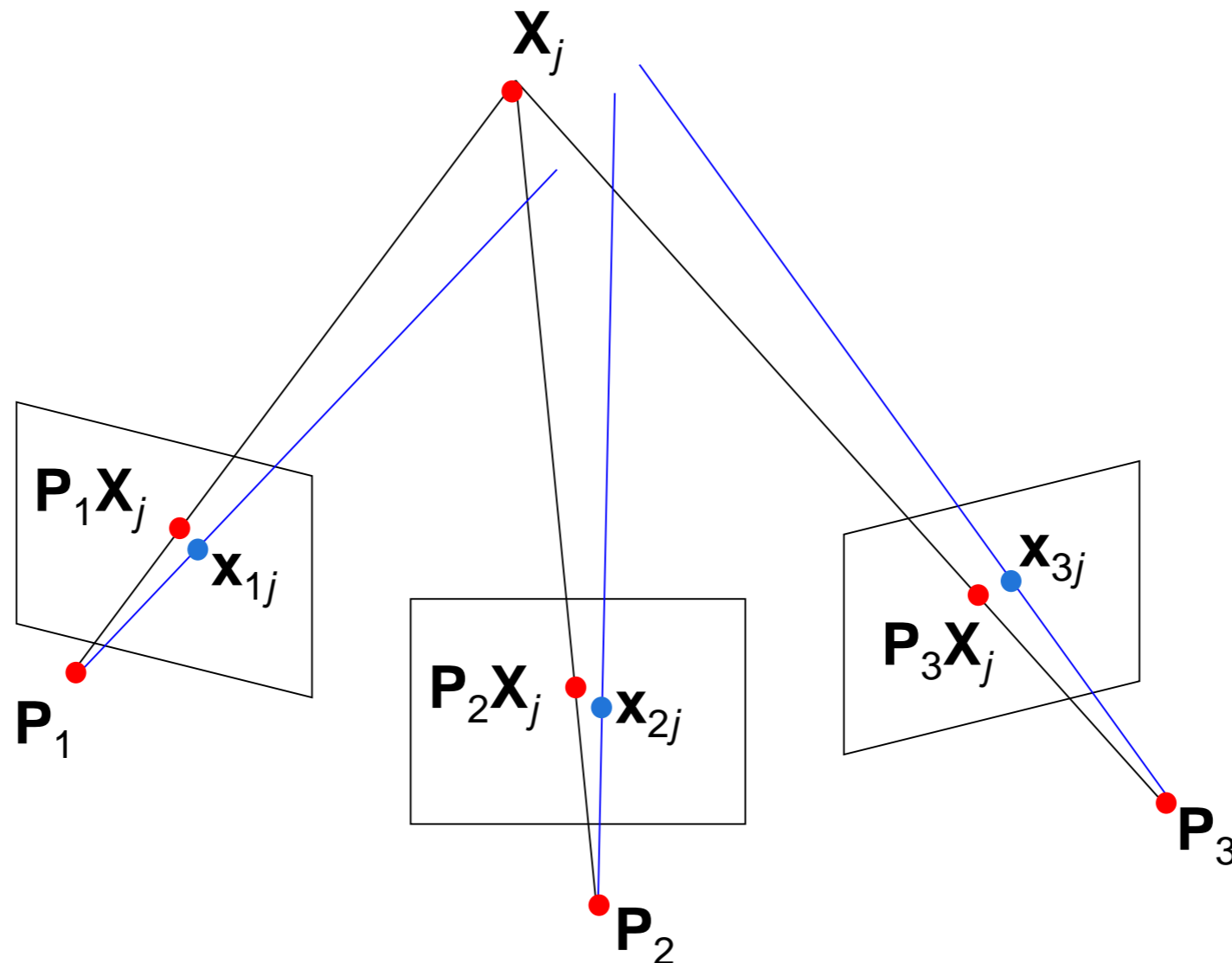
- Refine structure and motion: bundle adjustment

points

cameras

# Bundle adjustment

- Non-linear method for refining structure and motion

- Minimizing reprojection error

$$E(\mathbf{P}, \mathbf{X}) = \sum_{i=1}^{m} \sum_{j=1}^{n} D\left(\mathbf{x}_{ij}, \mathbf{P}_i \mathbf{X}_j\right)^2$$
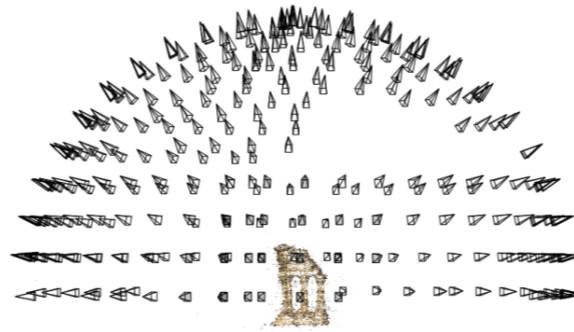
# Review: Structure from motion

- Ambiguity

- Affine structure from motion

  - Factorization

- Dealing with missing data

  - Incremental structure from motion

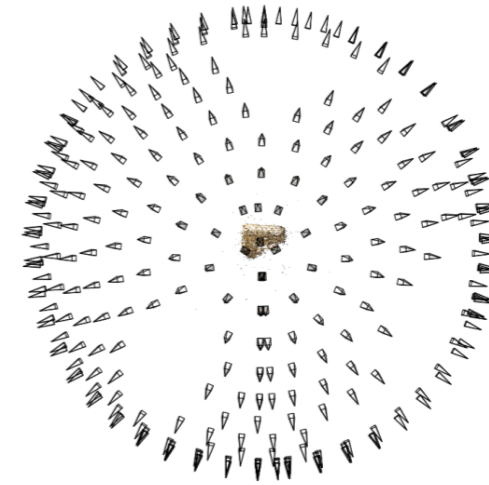- Projective structure from motion

  - Bundle adjustment

|  | Structure (scene geometry) | Motion (camera geometry) | Measurements |
|---|---|---|---|
| **Pose Estimation** | known | **estimate** | 3D to 2D correspondences |
| **Triangulation** | **estimate** | known | 2D to 2D cooorespondences |
| **Reconstruction** | **estimate** | **estimate** | 2D to 2D cooorespondences |

# Large-scale structure from motion

# Structure from motion



Reconstruction (side)

(top)

- Input: images with points in correspondence
  $p_{i,j} = (u_{i,j}, v_{i,j})$

- Output
  - structure: 3D location $\mathbf{x}_i$ for each point $p_i$
  - motion: camera parameters $\mathbf{R}_j$, $\mathbf{t}_j$ possibly $\mathbf{K}_j$
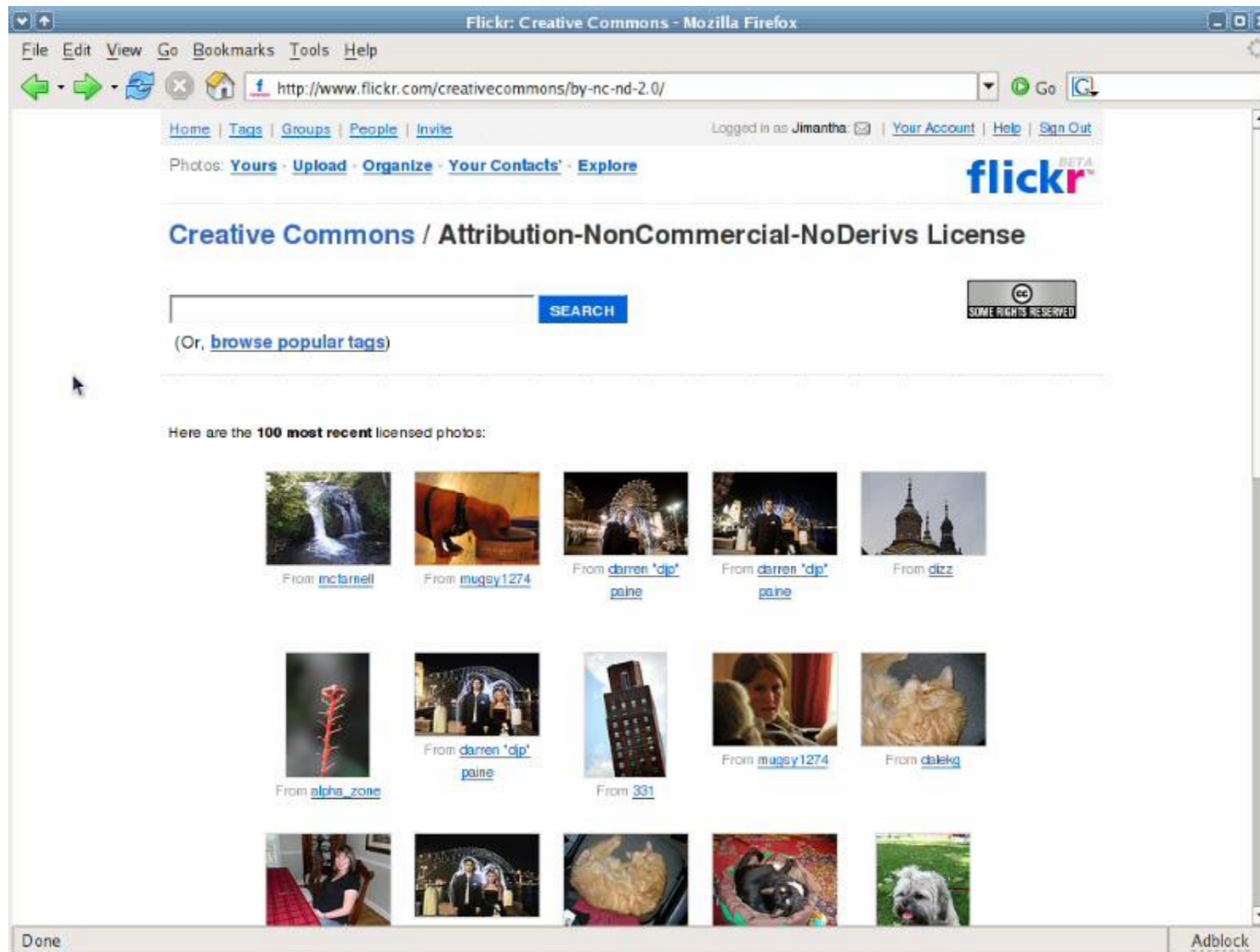
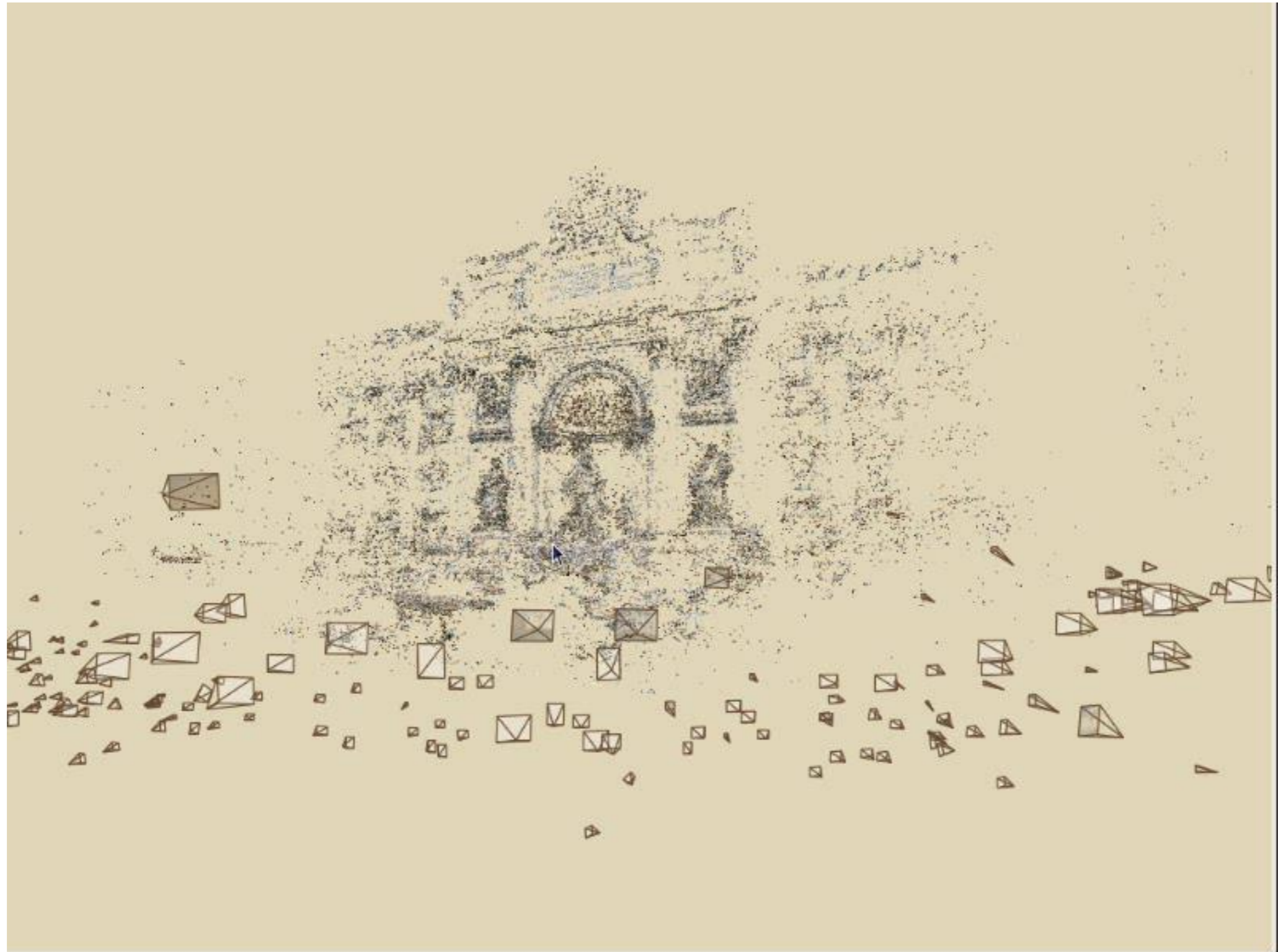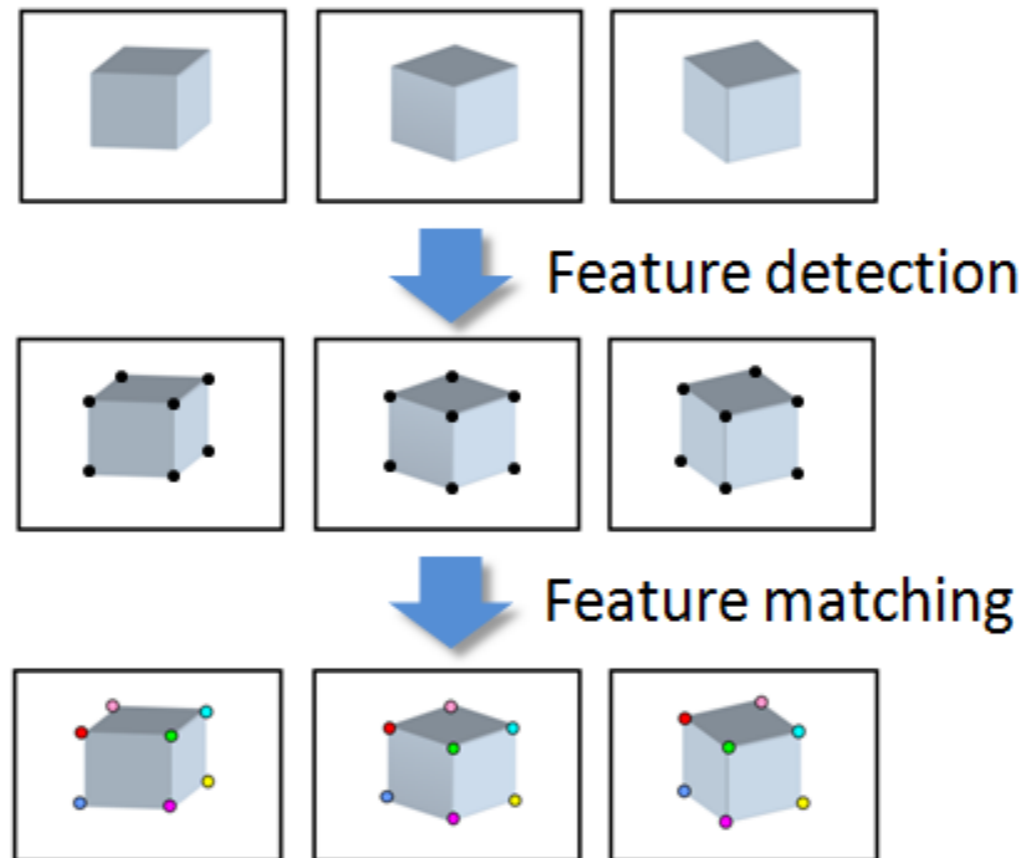- Objective function: minimize *reprojection error*

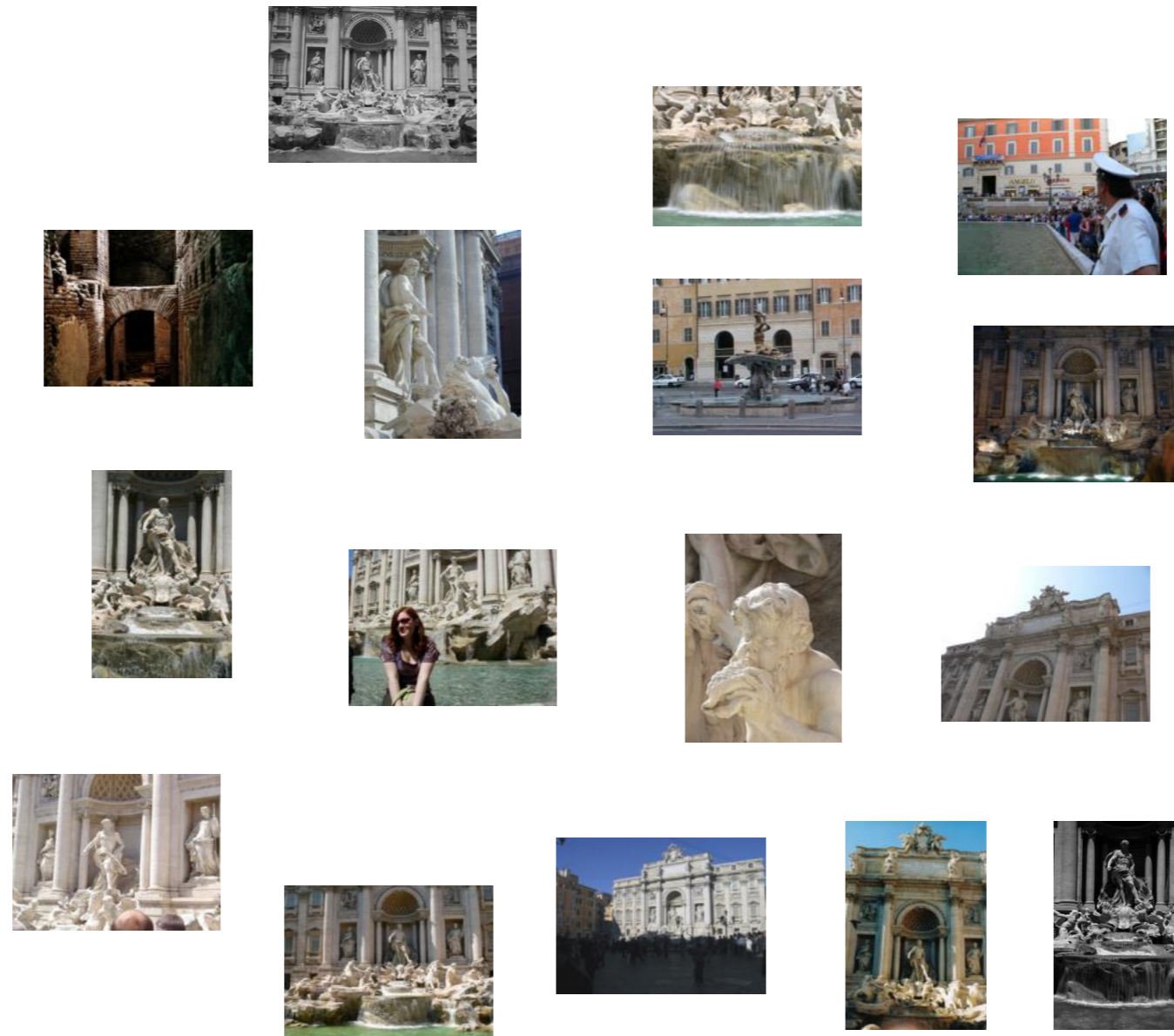15,464

37,383

76,389

# Standard way to view photos

# Photo Tourism

# Input: Point correspondences

# Feature detection

Detect features using SIFT [Lowe, IJCV 2004]

# Feature description

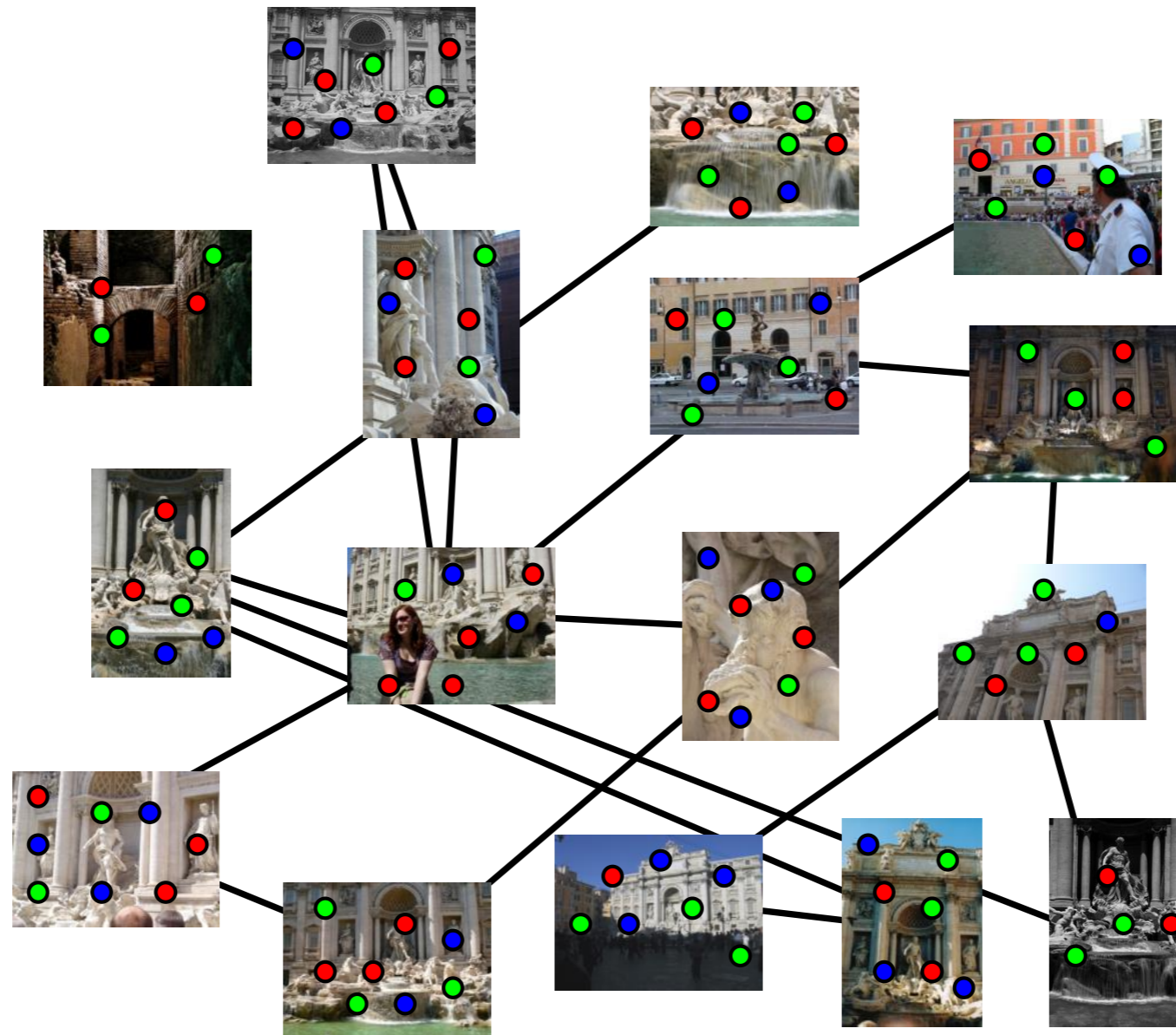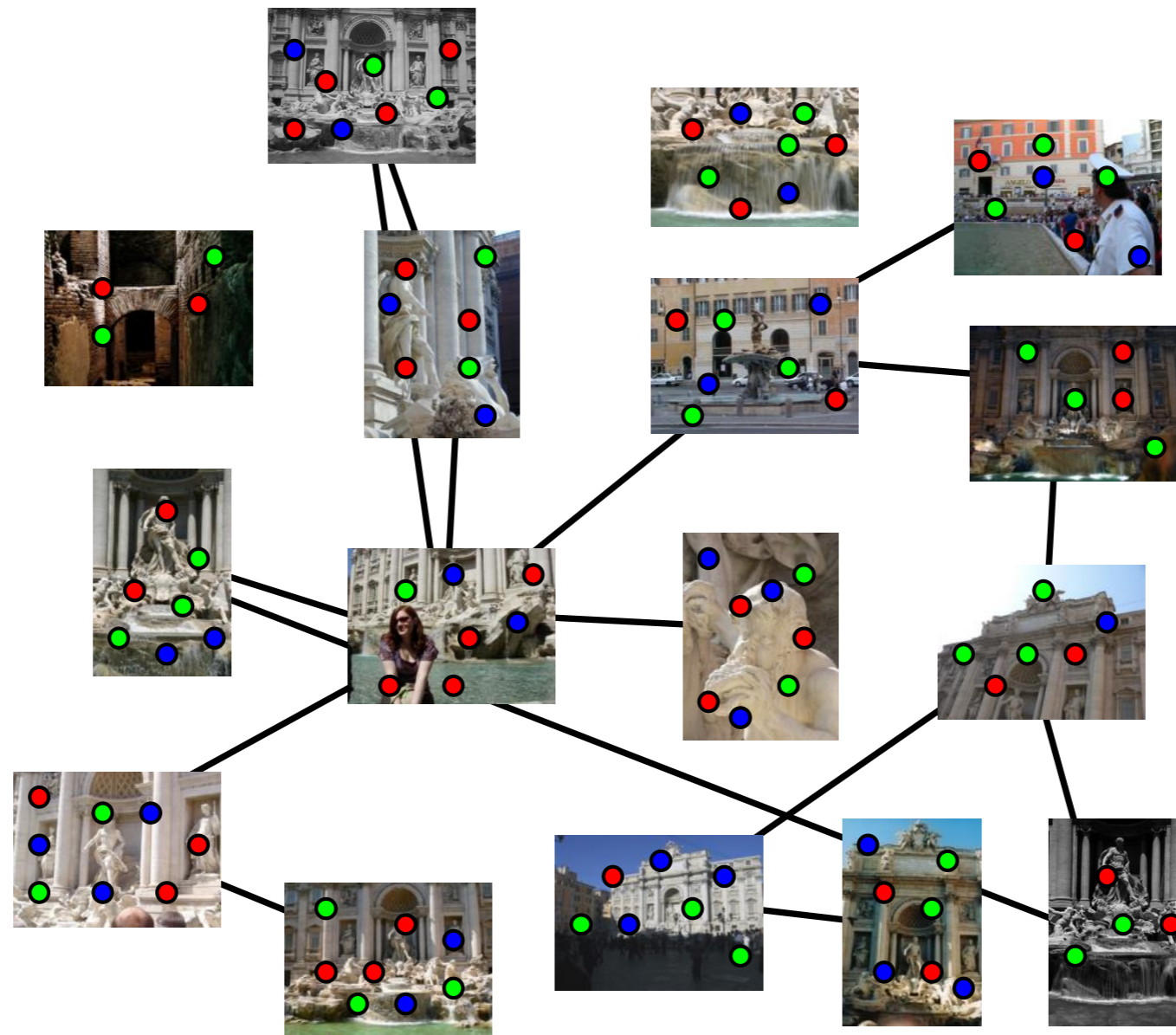Describe features using SIFT [Lowe, IJCV 2004]

# Feature matching

Match features between each pair of images

# Feature matching

Refine matching using RANSAC to estimate fundamental matrix between each pair

# Correspondence estimation

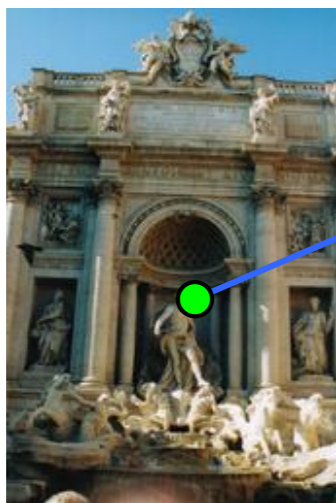- Link up pairwise matches to form connected components of matches across several images
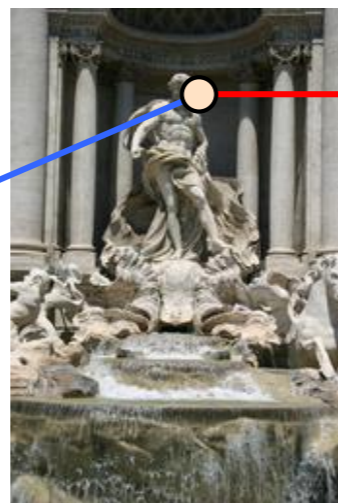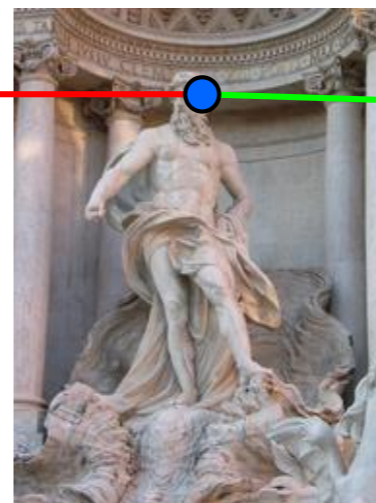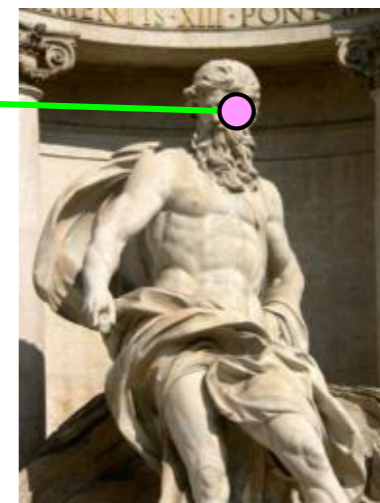


Image 1  Image 2  Image 3  Image 4

# Image connectivity graph



(graph layout produced using the Graphviz toolkit: http://www.graphviz.org/)

# Structure from motion



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

minimize
$$\mathrm{g}(\mathbf{R}, \mathbf{T}, \mathbf{X})$$
*non-linear least squares*

$X_4$

$X_1$     $X_3$

$X_2$

$X_5$    $X_7$

$X_6$

$p_{1,1}$

$p_{1,2}$

$p_{1,3}$

Camera 1
$R_1, t_1$

Camera 2
$R_2, t_2$

Camera 3
$R_3, t_3$

# Global structure from motion

- Minimize sum of squared reprojection errors:

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

*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
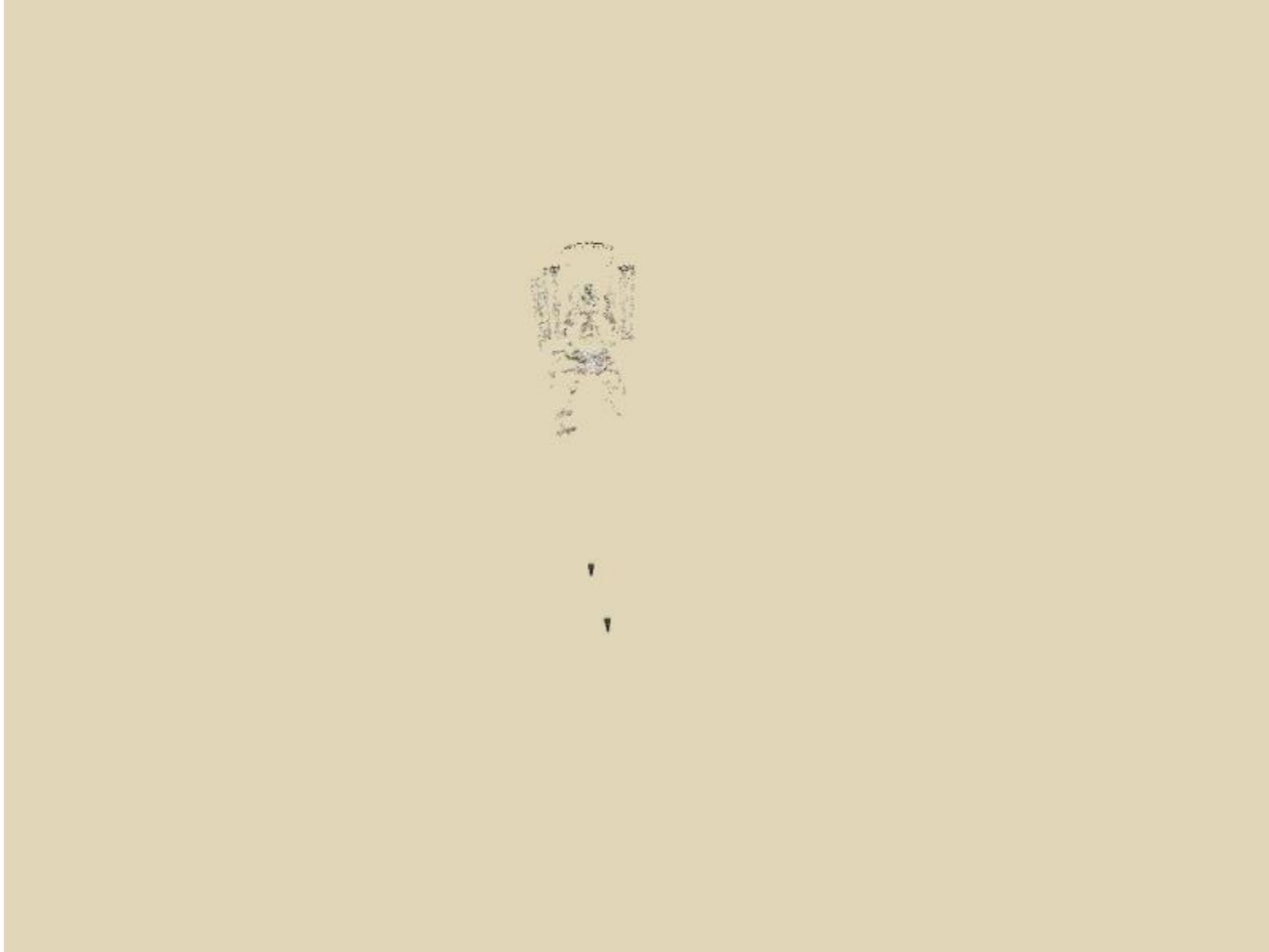
# Problem size

- What are the variables?

- How many variables per camera?

- How many variables per point?


- Trevi Fountain collection

    466 input photos

  + > 100,000 3D points

      = very large optimization problem

# Doing bundle adjustment

- Minimizing *g* is difficult
  - *g* is non-linear due to rotations, perspective division
  - lots of parameters: 3 for each 3D point, 6 for each camera
  - difficult to initialize
  - gauge ambiguity: error is invariant to a similarity transform (translation, rotation, uniform scale)

- Many techniques use non-linear least-squares (NLLS) optimization (*bundle adjustment*)
  - Levenberg-Marquardt is one common algorithm for NLLS
  - Lourakis, **The Design and Implementation of a Generic Sparse Bundle Adjustment Software Package Based on the Levenberg-Marquardt Algorithm**, http://www.ics.forth.gr/~lourakis/sba/
  - http://en.wikipedia.org/wiki/Levenberg-Marquardt_algorithm
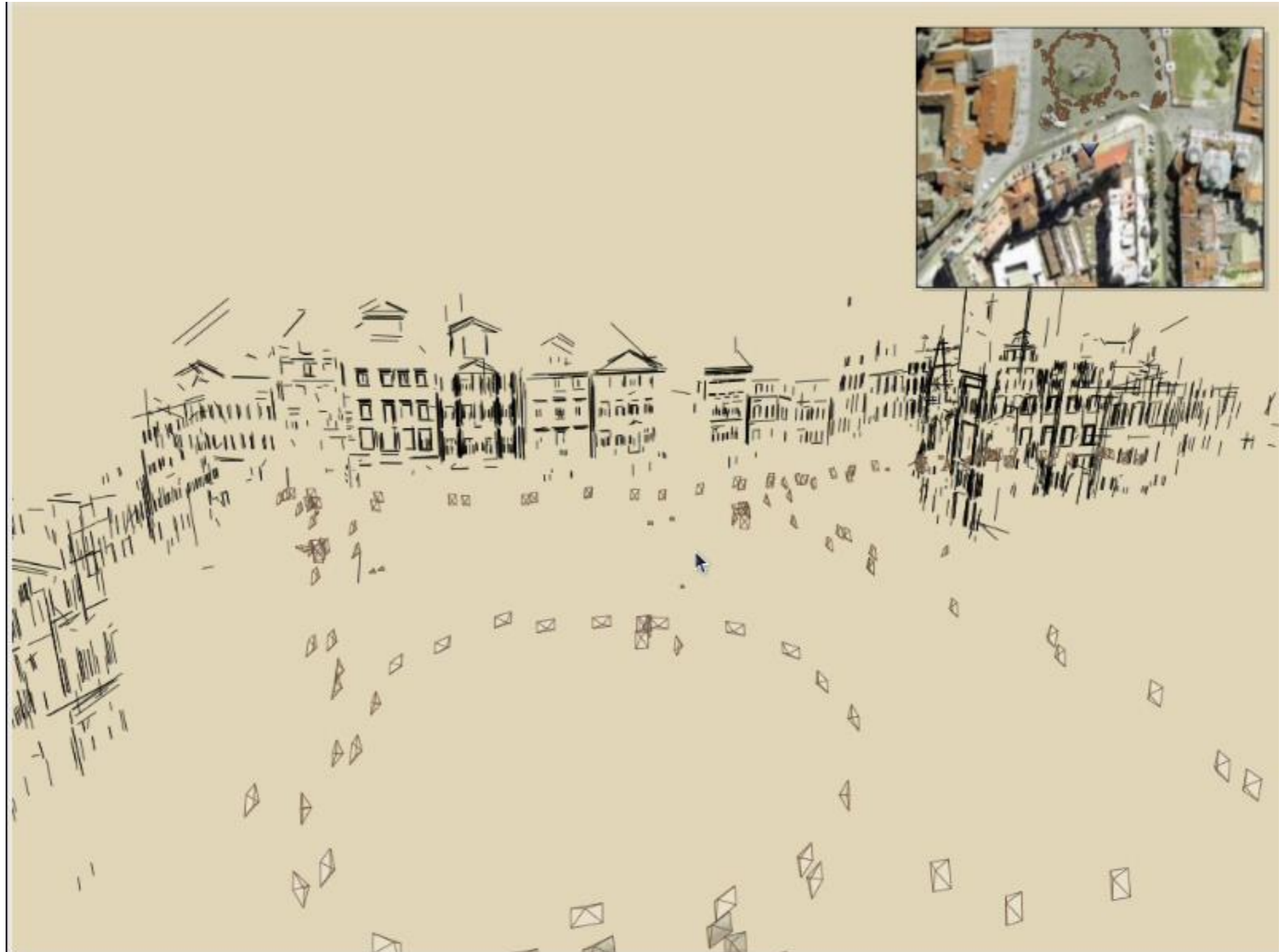
# Incremental structure from motion
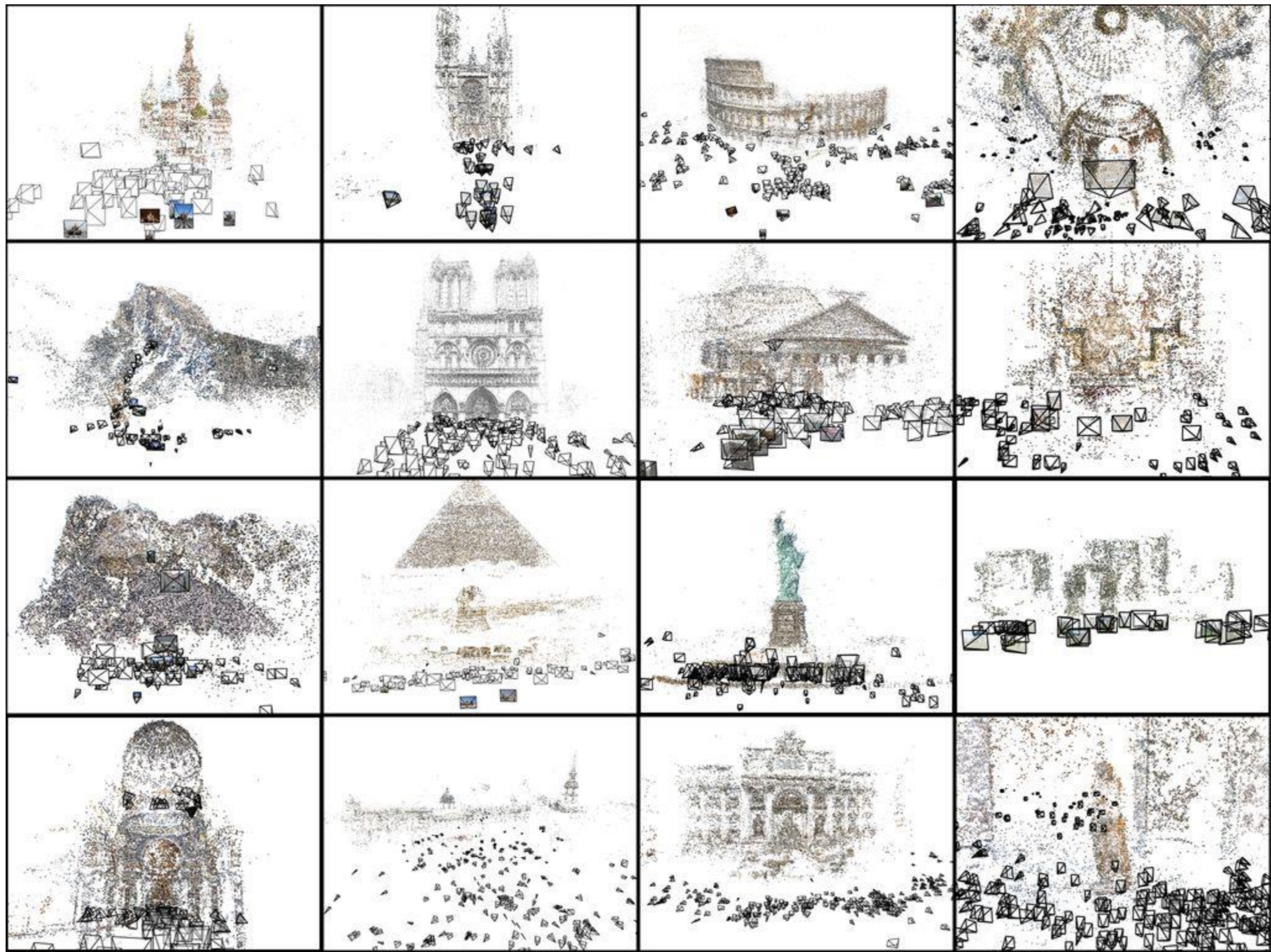
# Final reconstruction

# More examples

# More examples

# More examples

# Even larger scale SfM

City-scale structure from motion

- "Building Rome in a day"

http://grail.cs.washington.edu/projects/rome/

# SfM applications

- 3D modeling
- Surveying
- Robot navigation and mapmaking
- Visual effects ("Match moving")
  - https://www.youtube.com/watch?v=RdYWp70P_kY

# Applications – Photosynth

# Applications – Hyperlapse



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

# Summary: 3D geometric vision

- ## Single-view geometry
  - ### The pinhole camera model
    - Variation: orthographic projection
  - ### The perspective projection matrix
  - ### Intrinsic parameters
  - ### Extrinsic parameters
  - ### Calibration

- ## Multiple-view geometry
  - ### Triangulation
  - ### The epipolar constraint
    - Essential matrix and fundamental matrix
  - ### Stereo
    - Binocular, multi-view
  - ### Structure from motion
    - Reconstruction ambiguity
    - Affine SFM
    - Projective SFM

# References

Basic reading:
- Szeliski textbook, Chapter 7.
- Hartley and Zisserman, Chapter 18.