# Structure from motion



Драконъ, видимый подъ различными углами зрѣнія
По гравюрѣ на мѣди изъ „Oculus artificialis teledioptricus" Цана. 1702 года.

# Multiple-view geometry questions

- **Scene geometry (structure):** Given 2D point matches in two or more images, where are the corresponding points in 3D?

- **Correspondence (stereo matching):** Given a point in just one image, how does it constrain the position of the corresponding point in another image?

- **Camera geometry (motion):** Given a set of corresponding points in two or more images, what are the camera matrices for these views?
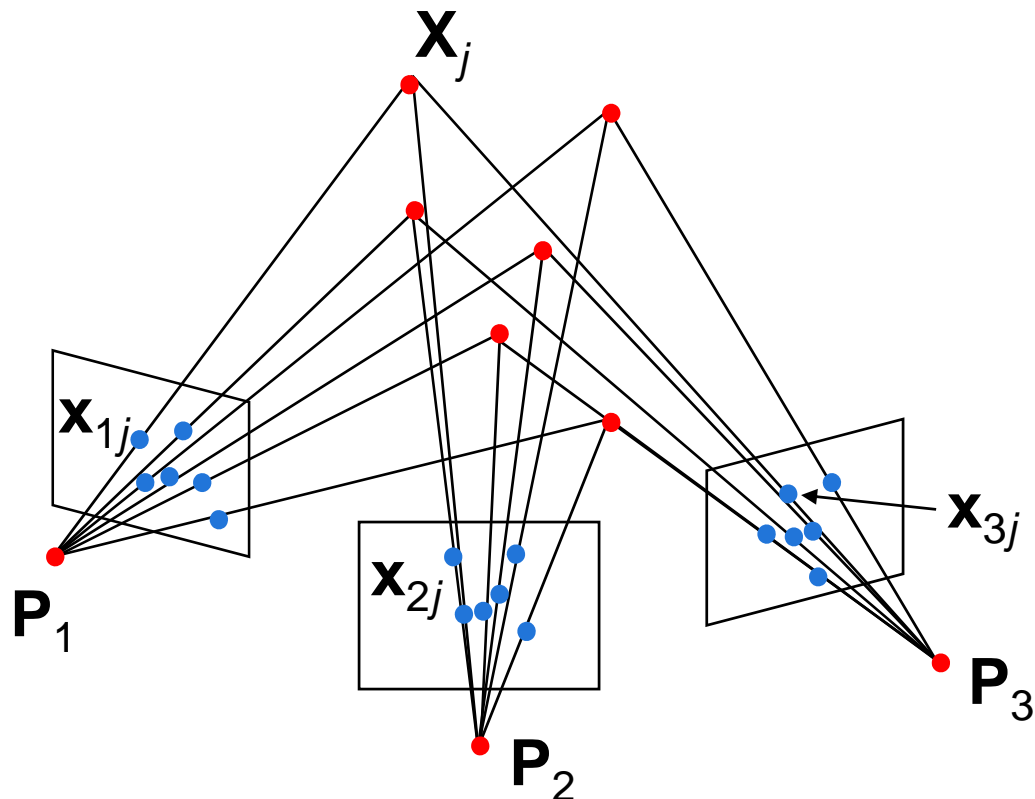
# Structure from motion

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

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

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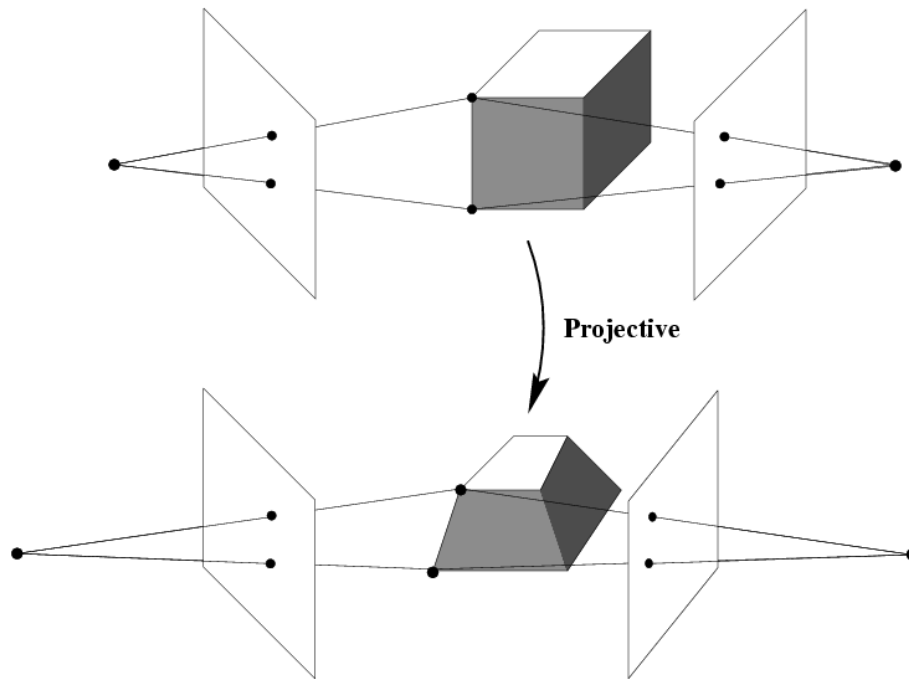
# Structure from motion ambiguity

- 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:

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

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 **Q** and apply the inverse transformation to the camera matrices, then the images do not change

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

# Projective ambiguity



Projective

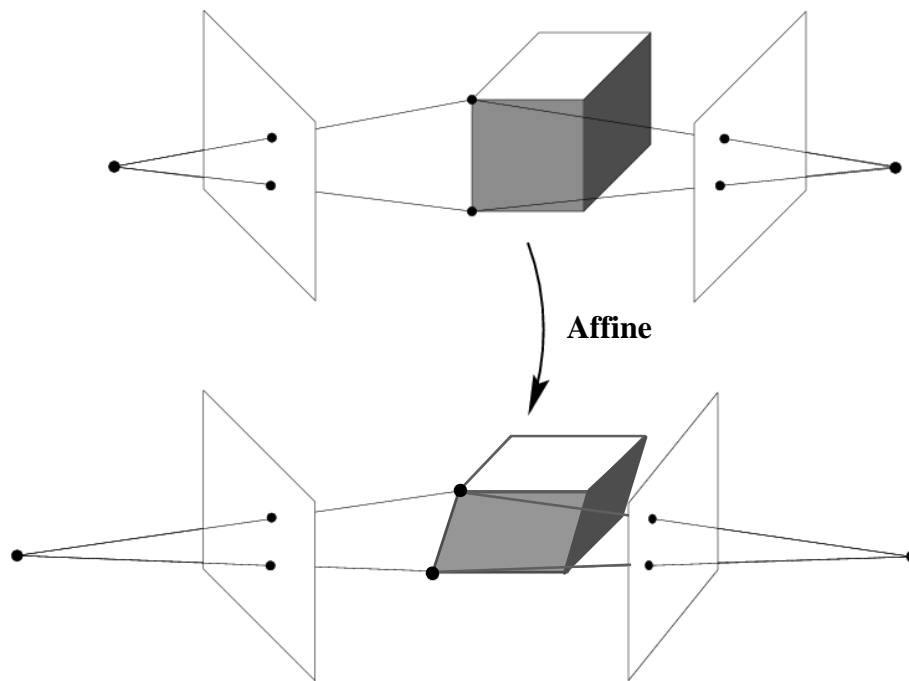$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ_P^{-1}}\right)\left(\mathbf{Q_P\,X}\right)$$
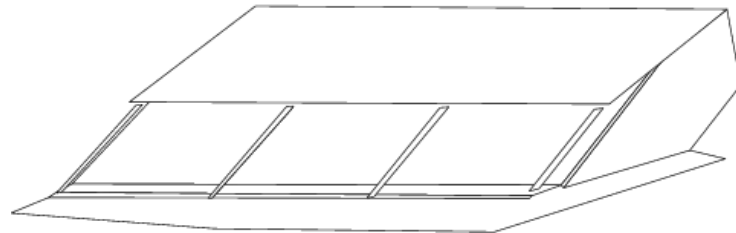
# Projective ambiguity

# Affine ambiguity



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

# Affine ambiguity

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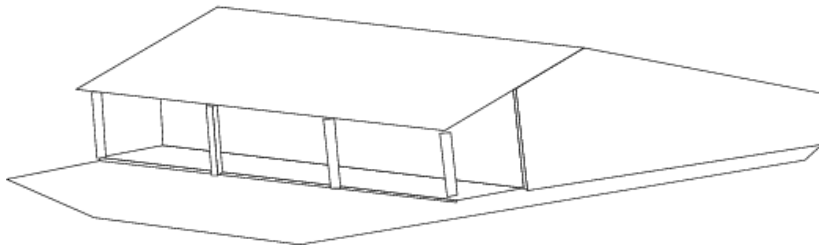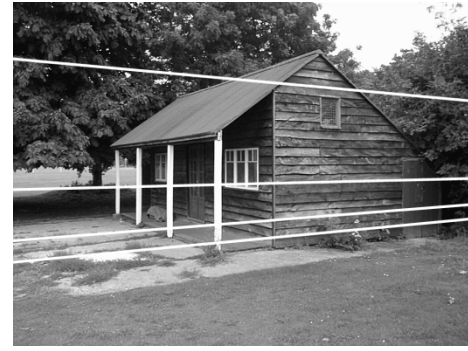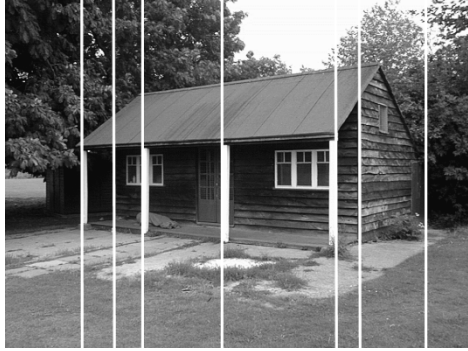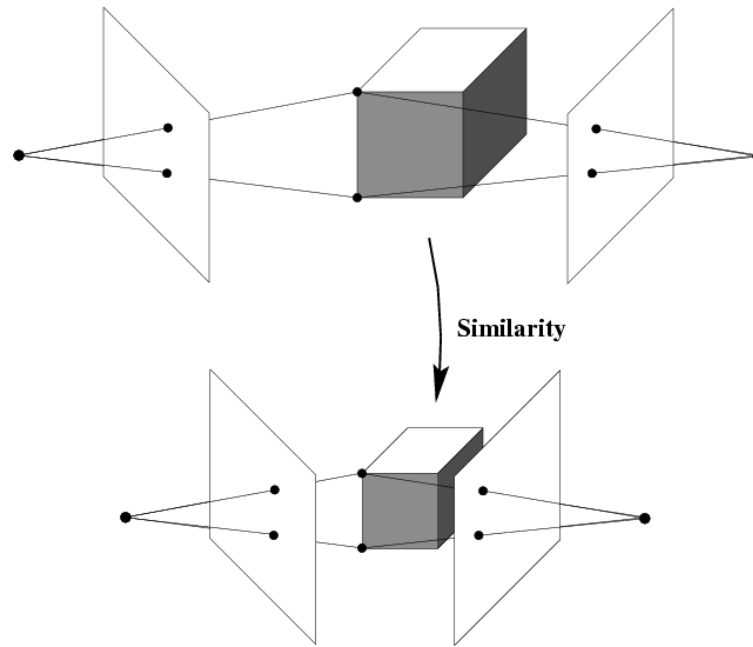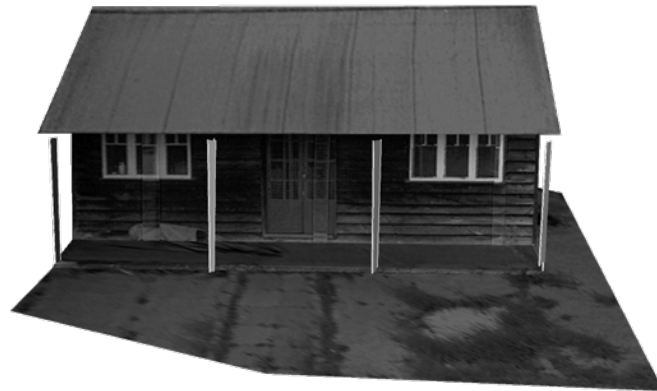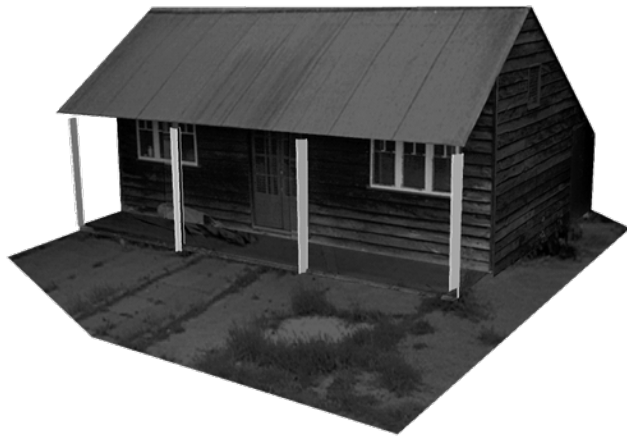
$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ}_S^{-1}\right)\left(\mathbf{Q}_S\mathbf{X}\right)$$

# Similarity ambiguity

# Hierarchy of 3D transformations

Projective
15dof

$$\begin{bmatrix} A & t \\ v^T & v \end{bmatrix}$$

Preserves intersection and tangency

Affine
12dof

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

Preserves parallellism, volume ratios

Similarity
7dof

$$\begin{bmatrix} s\,R & t \\ 0^T & 1 \end{bmatrix}$$

Preserves angles, ratios of length

Euclidean
6dof

$$\begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix}$$

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

# Structure from motion

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



center at
infinity

perspective

weak perspective

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



- 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:

$$\mathbf{P} = [3 \times 3\,\text{affine}] \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} [4 \times 4\,\text{affine}] = \begin{bmatrix} a_{11} & a_{12} & a_{13} & b_1 \\ a_{21} & a_{22} & a_{23} & b_2 \\ 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{0} & 1 \end{bmatrix}$$

- Affine projection is a linear mapping + translation in inhomogeneous coordinates
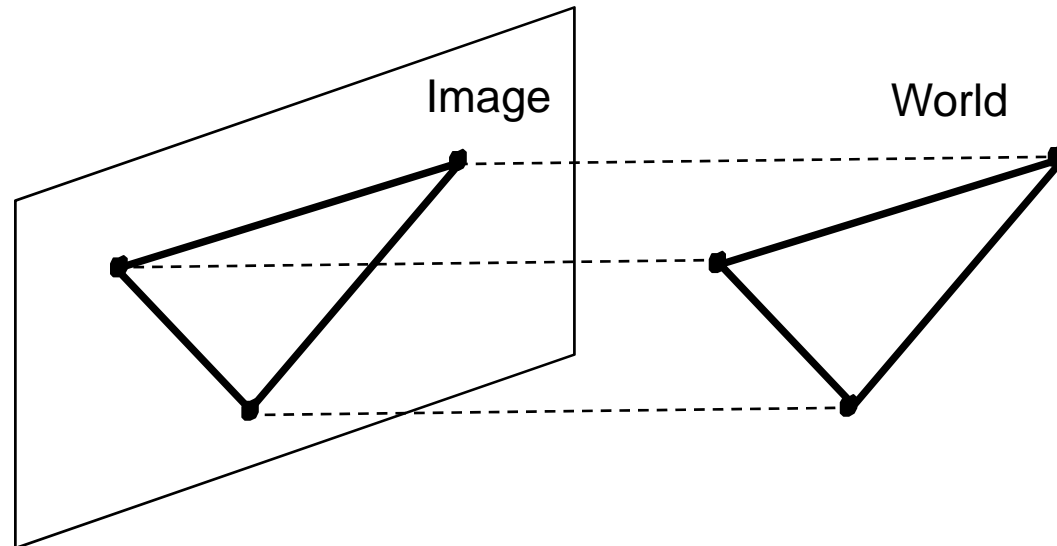
$$\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \begin{pmatrix} X \\ Y \\ Z \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = \mathbf{A}\mathbf{X} + \mathbf{b}$$

$\mathbf{x}$

$\mathbf{a}_2$

$\mathbf{a}_1$

$\mathbf{X}$

Projection of world origin

# Affine structure from motion

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

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

- 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 $2mn \geq 8m + 3n - 12$
- 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>
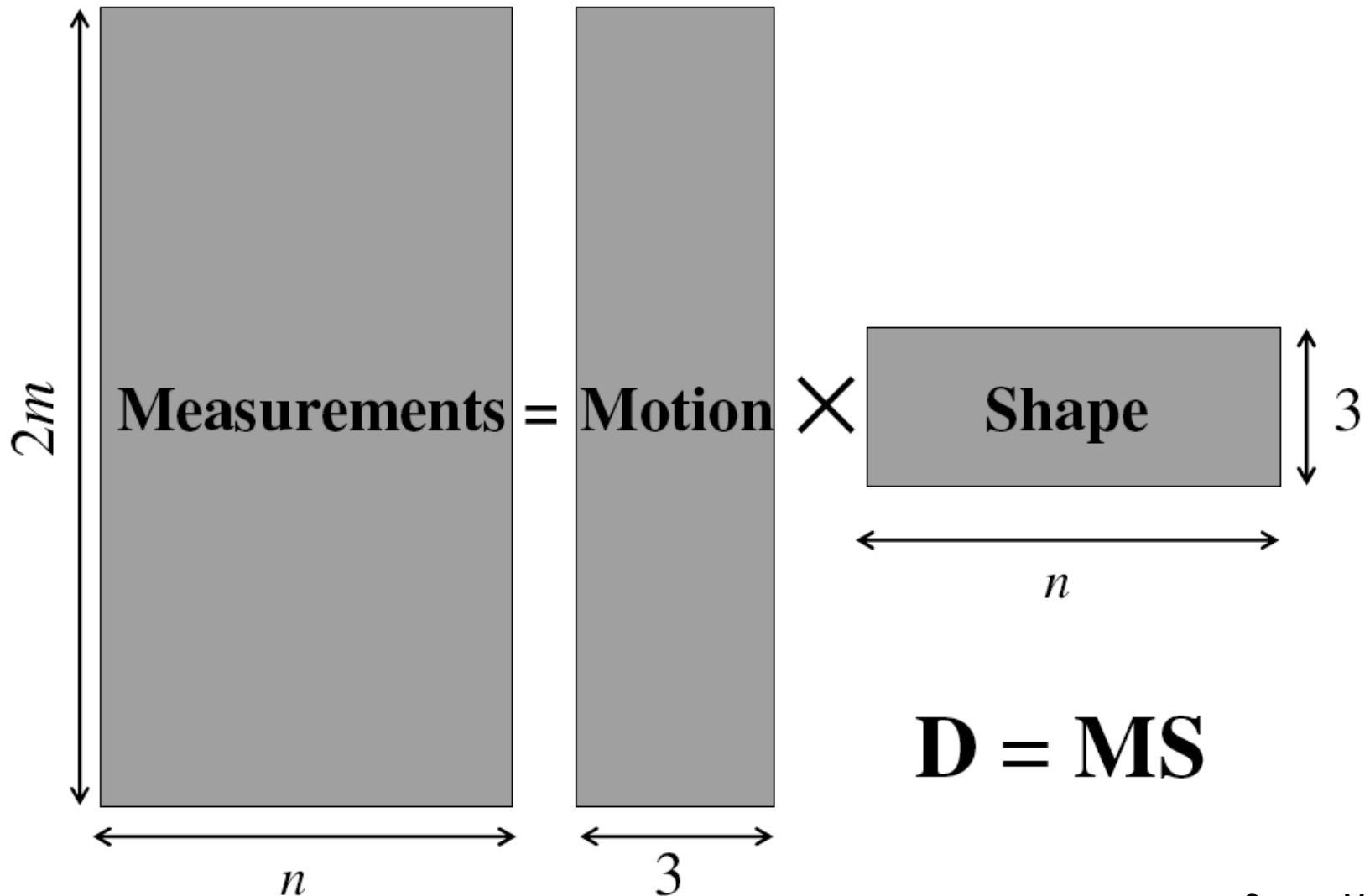
<span style="color:red">cameras ($2m \times 3$)</span>

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



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

$$D = MS$$
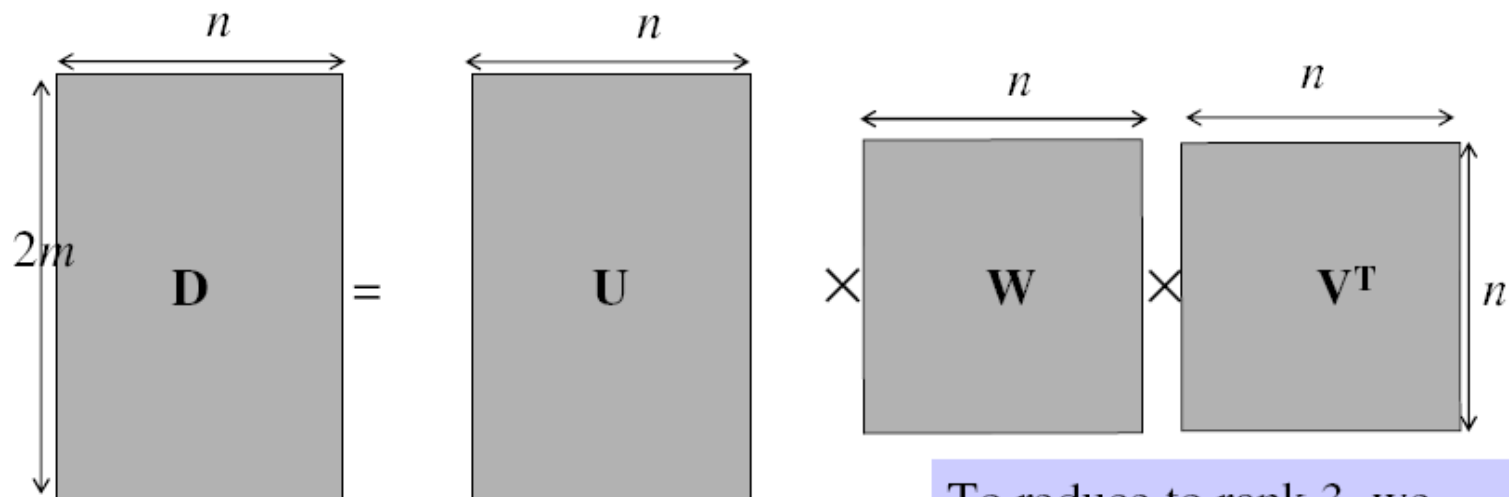
# Factorizing the measurement matrix

- Singular value decomposition of D:



Source: M. Hebert

# 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

Source: M. Hebert

# Factorizing the measurement matrix

- Obtaining a factorization from SVD:

# Factorizing the measurement matrix

- Obtaining a factorization from SVD:

$$D = U_3 \times W_3 \times V_3^T$$

with dimensions $2m$, $3$, $3$, $n$.

Possible decomposition:

$$M = U_3 W_3^{1/2} \qquad S = W_3^{1/2} V_3^T$$

$$D = M \times S$$

This decomposition minimizes $|D-MS|^2$

# Affine ambiguity

$$D = M \times S$$

- The decomposition is not unique. We get the same **D** by using any 3×3 matrix **C** and applying the transformations **M** → **MC, S** →**C⁻¹S**

- 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 scale is 1

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

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

- This translates into $3m$ equations in $\mathbf{L} = \mathbf{CC^T}$ :

$$\mathbf{A}_i \, \mathbf{L} \, \mathbf{A}_i^T = \mathbf{Id}, \qquad i = 1, \ldots, m$$

  - Solve for **L**
  - Recover **C** from **L** by Cholesky decomposition: $\mathbf{L} = \mathbf{CC^T}$
  - Update **M** and **S**: $\mathbf{M = MC}, \mathbf{S = C^{-1}S}$

# 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: $\mathbf{D} = \mathbf{U\ W\ V^T}$
  - 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:
  - $\mathbf{M} = \mathbf{U}_3\mathbf{W}_3^{½}$ and $\mathbf{S} = \mathbf{W}_3^{½}\ \mathbf{V}_3^T$ (**or** $\mathbf{M} = \mathbf{U}_3$ and $\mathbf{S} = \mathbf{W}_3\mathbf{V}_3^T$)
- Eliminate affine ambiguity

# Reconstruction results



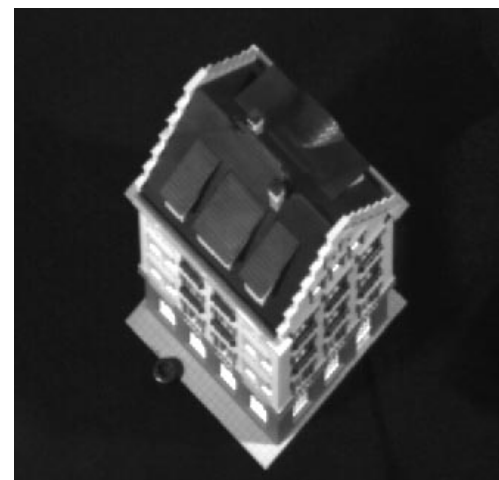C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization method. *IJCV*, 9(2):137-154, November 1992.

# 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:



cameras

points

# Dealing with missing data

- Possible solution: decompose matrix into dense sub-blocks, factorize each sub-block, and fuse the results
  - Finding dense maximal sub-blocks of the matrix is NP-complete (equivalent to finding maximal cliques in a graph)

- Incremental bilinear refinement



(1) Perform factorization on a dense sub-block

(2) Solve for a new 3D point visible by at least two known cameras (linear least squares)

(3) Solve for a new camera that sees at least three known 3D points (linear least squares)

F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce. Segmenting, Modeling, and Matching Video Clips Containing Multiple Moving Objects. PAMI 2007.

# 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

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

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

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

- 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 $z\mathbf{x} = [\mathbf{I} \mid \mathbf{0}]\mathbf{X}, \quad z'\mathbf{x}' = [\mathbf{A} \mid \mathbf{b}]\mathbf{X}$

$$z'\mathbf{x}' = \mathbf{A}[\mathbf{I} \mid \mathbf{0}]\mathbf{X} + \mathbf{b} = z\mathbf{A}\mathbf{x} + \mathbf{b}$$

$$z'\mathbf{x}' \times \mathbf{b} = z\mathbf{A}\mathbf{x} \times \mathbf{b}$$

$$(z'\mathbf{x}' \times \mathbf{b}) \cdot \mathbf{x}' = (z\mathbf{A}\mathbf{x} \times \mathbf{b}) \cdot \mathbf{x}'$$

$$\mathbf{x}'^{\mathrm{T}}[\mathbf{b}_\times]\mathbf{A}\mathbf{x} = 0$$

$$\mathbf{F} = [\mathbf{b}_\times]\mathbf{A} \quad \mathbf{b}: \text{epipole } (\mathbf{F}^{\mathrm{T}}\mathbf{b} = 0), \quad \mathbf{A} = -[\mathbf{b}_\times]\mathbf{F}$$

# Projective factorization

$$\mathbf{D} = \begin{bmatrix} z_{11}\mathbf{X}_{11} & z_{12}\mathbf{X}_{12} & \cdots & z_{1n}\mathbf{X}_{1n} \\ z_{21}\mathbf{X}_{21} & z_{22}\mathbf{X}_{22} & \cdots & z_{2n}\mathbf{X}_{2n} \\ & & \ddots & \\ z_{m1}\mathbf{X}_{m1} & z_{m2}\mathbf{X}_{m2} & \cdots & z_{mn}\mathbf{X}_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{P}_1 \\ \mathbf{P}_2 \\ \vdots \\ \mathbf{P}_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}$$
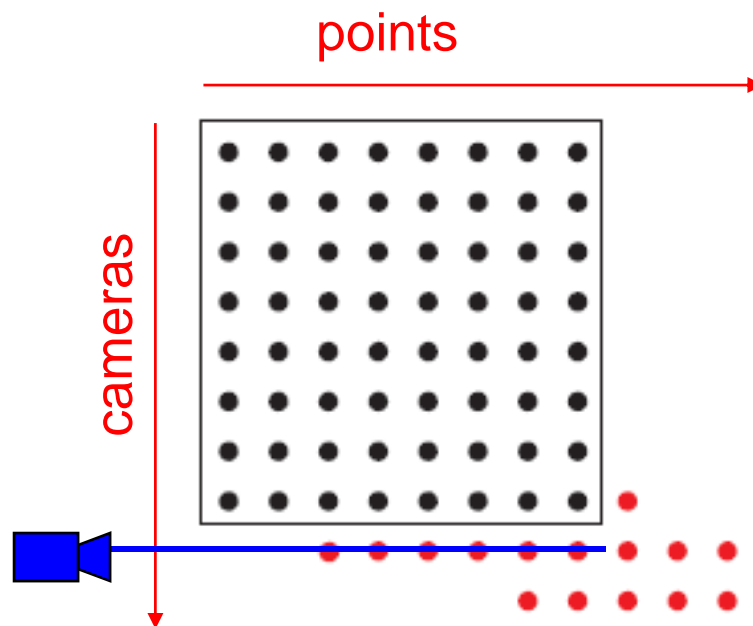
points ($4 \times n$)

cameras
($3m \times 4$)

$$\mathbf{D} = \mathbf{MS} \text{ has rank } 4$$

- If we knew the depths $z$, we could factorize $\mathbf{D}$ to estimate $\mathbf{M}$ and $\mathbf{S}$

- If we knew $\mathbf{M}$ and $\mathbf{S}$, we could solve for $z$

- Solution: iterative approach (alternate between above two steps)

# Sequential structure from motion

- Initialize motion from two images using fundamental matrix

- Initialize structure

- 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

- 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*

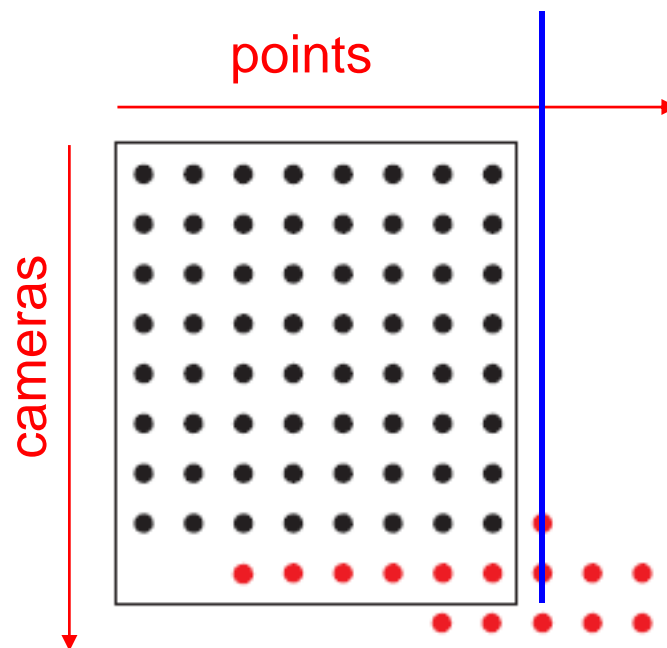points

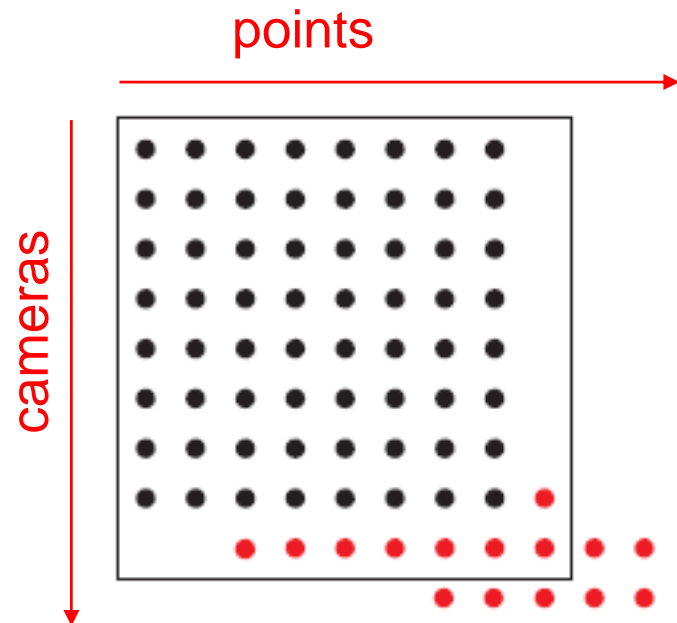cameras

# Sequential structure from motion

- Initialize motion from two images using fundamental matrix

- Initialize structure

- 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
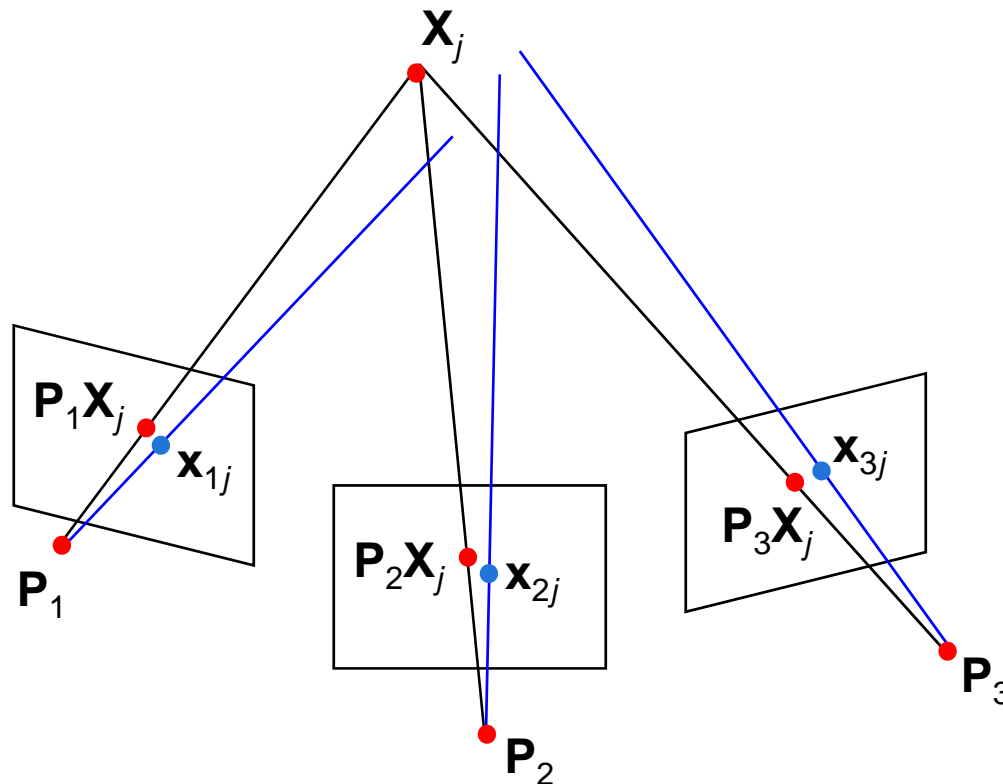
points

cameras

# Bundle adjustment

- Non-linear method for refining structure and motion
- Minimizing reprojection error

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

# Self-calibration

- Self-calibration (auto-calibration) is the process of determining intrinsic camera parameters directly from uncalibrated images

- For example, when the images are acquired by a single moving camera, we can use the constraint that the intrinsic parameter matrix remains fixed for all the images
  - Compute initial projective reconstruction and find 3D projective transformation matrix $\mathbf{Q}$ such that all camera matrices are in the form $\mathbf{P}_i = \mathbf{K} [\mathbf{R}_i \mid \mathbf{t}_i]$

- Can use constraints on the form of the calibration matrix: zero skew

# Summary: Structure from motion

- Ambiguity
- Affine structure from motion: factorization
- Dealing with missing data
- Projective structure from motion: two views
- Projective structure from motion: iterative factorization
- Bundle adjustment
- Self-calibration

# 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