# C280, Computer Vision

Prof. Trevor Darrell
[trevor@eecs.berkeley.edu](mailto:trevor@eecs.berkeley.edu)

Lecture 11: Structure from Motion

# Roadmap

- Previous: Image formation, filtering, local features, (Texture)…
- Tues: Feature-based Alignment
  - Stitching images together
  - Homographies, RANSAC, Warping, Blending
  - Global alignment of planar models
- Today: Dense Motion Models
  - Local motion / feature displacement
  - Parametric optic flow
- No classes next week: ICCV conference
- Oct 6th: Stereo / 'Multi-view': Estimating depth with known inter-camera pose
- **Oct 8th: 'Structure-from-motion': Estimation of pose and 3D structure**
  - **Factorization approaches**
  - **Global alignment with 3D point models**

# Last time: Stereo

- Human stereopsis & stereograms

- Epipolar geometry and the epipolar constraint

  - Case example with parallel optical axes

  - General case with calibrated cameras

- Correspondence search

- The Essential  and the Fundamental Matrix

- Multi-view stereo

# Today: SFM

- SFM problem statement
- Factorization
- Projective SFM
- Bundle Adjustment
- Photo Tourism
- "Rome in a day:

# Structure from motion



Дракопъ, видимый подъ различными углами зрѣнія
По гравюрѣ на мѣди изъ „Oculus artificialis teledioptricus" Цана. 1702 года.
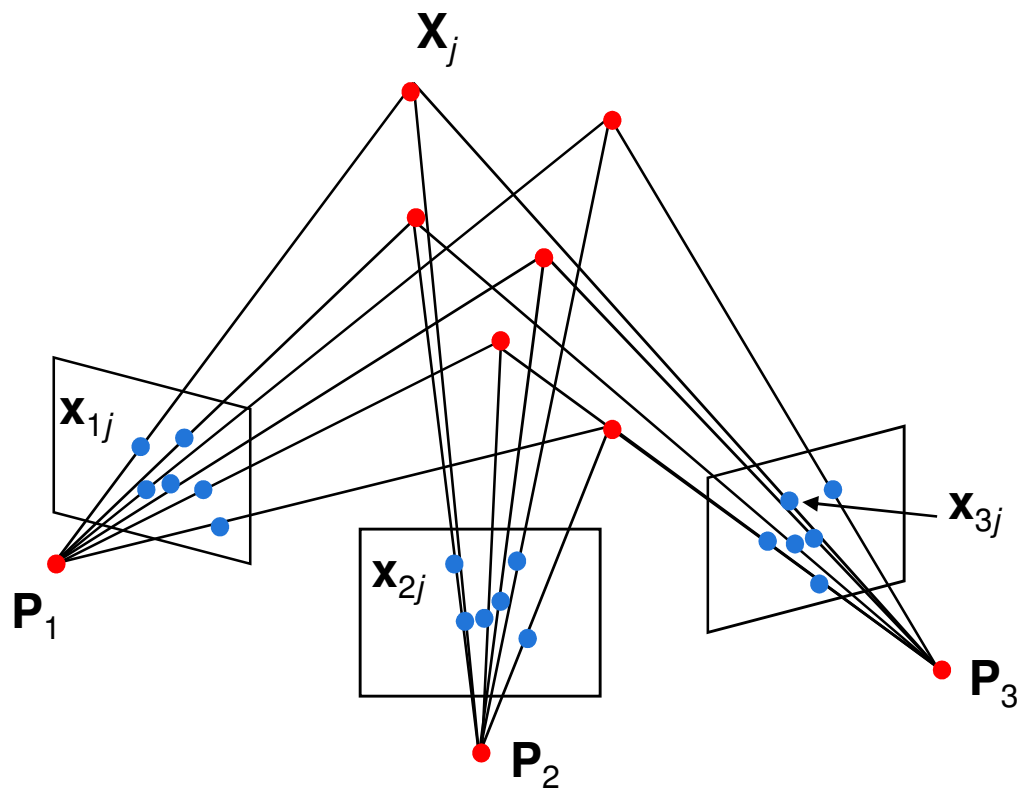
# Multiple-view geometry questions

- **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

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

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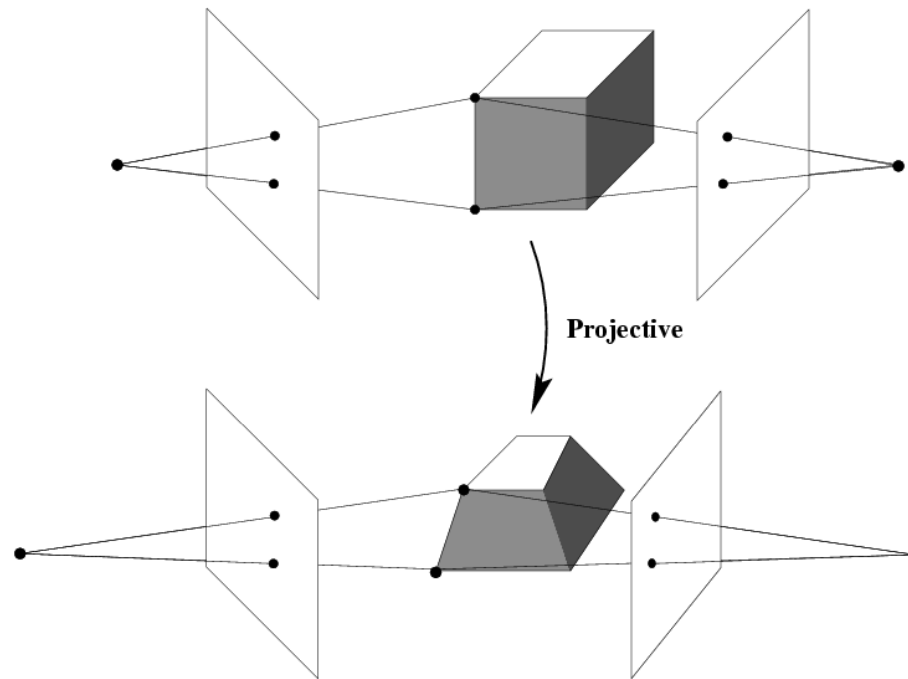
# Structure from motion ambiguity

- 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 $\mathbf{Q}$ and apply the inverse transformation to the camera matrices, then the images do not change
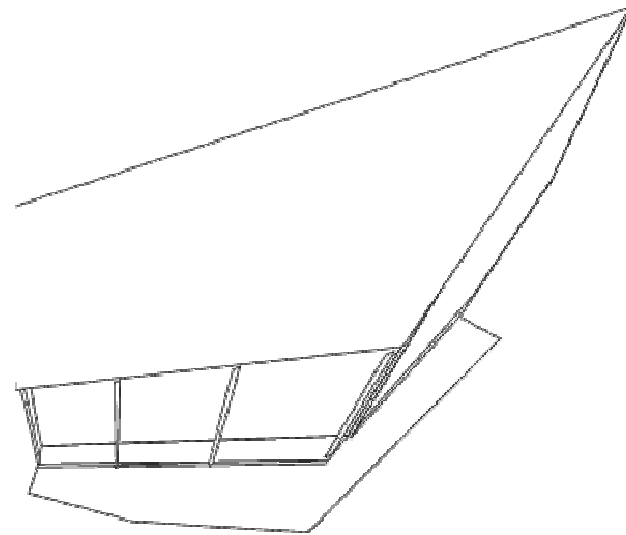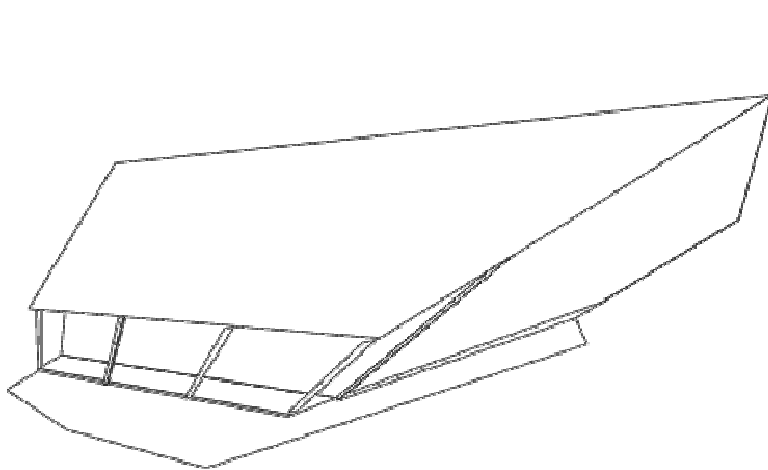
$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ^{-1}}\right)\left(\mathbf{QX}\right)$$

# Projective ambiguity



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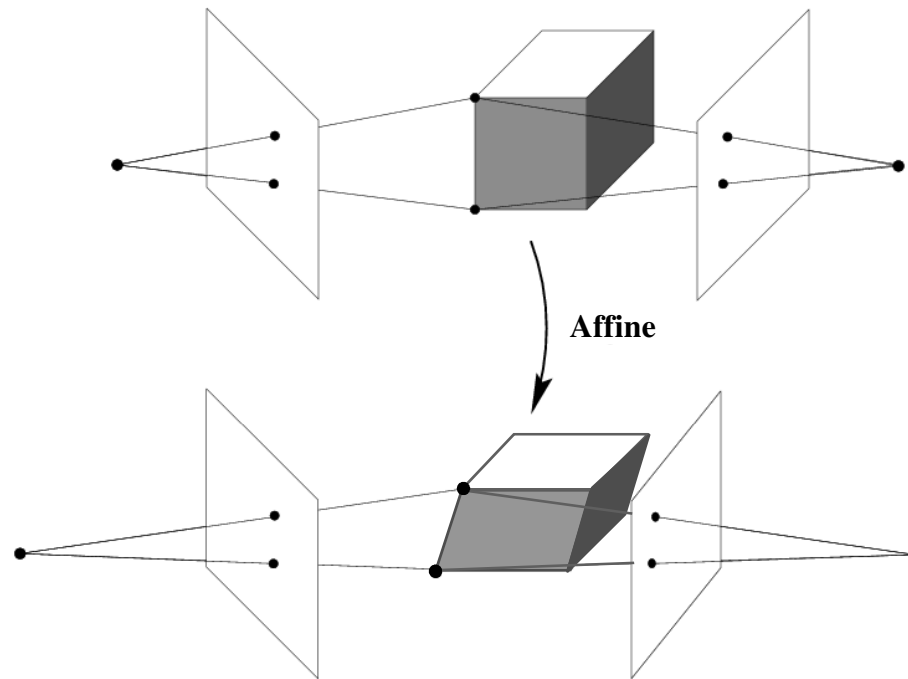
$$\mathbf{x} = \mathbf{PX} = \left(\mathbf{PQ_P^{-1}}\right)\left(\mathbf{Q_P}\,\mathbf{X}\right)$$

# Projective ambiguity

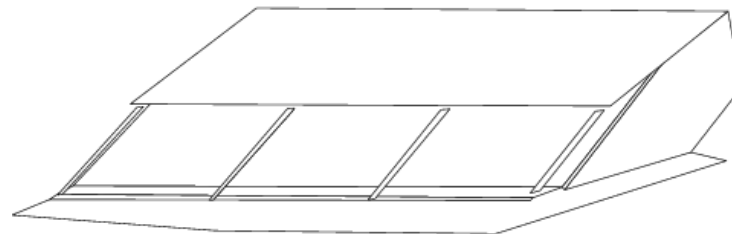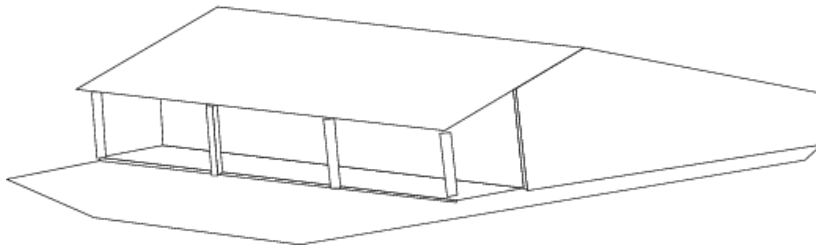# Affine ambiguity



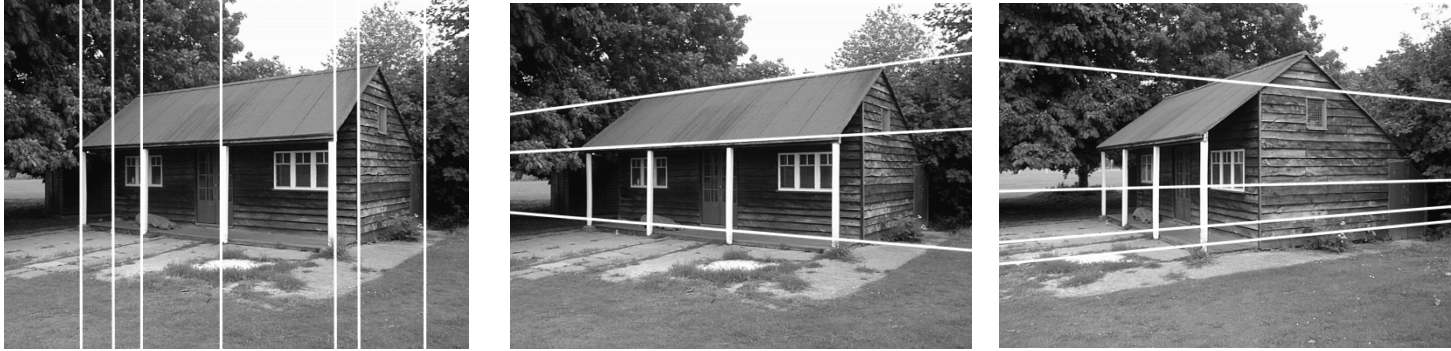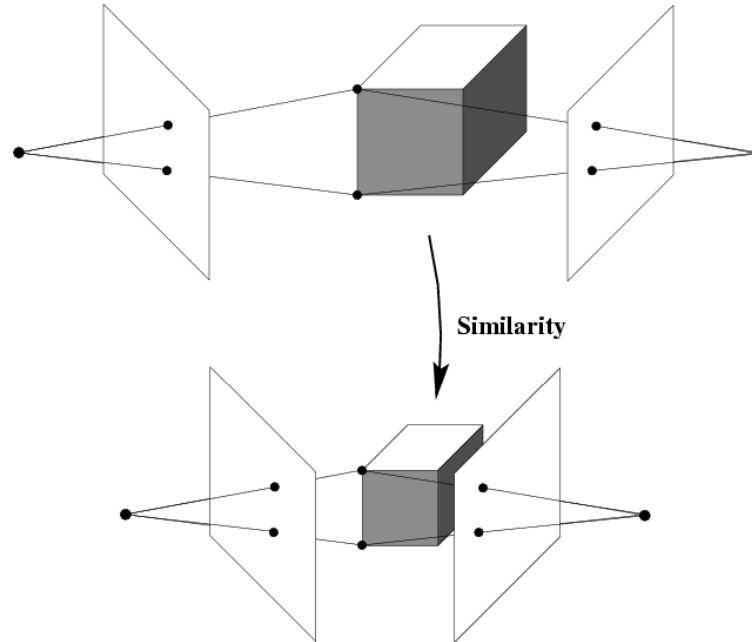Affine

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

# Affine ambiguity

# Similarity ambiguity



Similarity

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

# Similarity ambiguity

# Hierarchy of 3D transformations

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

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

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

Euclidean
6dof
$$\begin{bmatrix} R & t \\ 0^\mathsf{T} & 1 \end{bmatrix}$$
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

Lazebnik

# 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:

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

Lazebnik

# 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

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

$\mathbf{x}$

$\mathbf{a}_2$

$\mathbf{a}_1$

$\mathbf{X}$

Projection of world origin

Lazebnik

# 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 $2mn >= 8m + 3n - 12$
- For two views, we need four point correspondences

Lazebnik

# Affine structure from motion

- Centering: subtract the centroid of the image points

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

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

- 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:

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

cameras
($2\,m$)

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:

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

points ($3 \times n$)

cameras
($2m \times 3$)

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



$$D = MS$$

# Factorizing the measurement matrix

- Singular value decomposition of D:



$2m$ $D$ $=$ $U$ $\times$ $W$ $\times$ $V^T$   $n$

$2m$ $D$ $=$ $U_3$ $\times$ $W_3$ $\times$ $V_3^T$

Source: M. Hebert

Lazebnik

# 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

Lazebnik

Source: M. Hebert

# Factorizing the measurement matrix

- Obtaining a factorization from SVD:

$$2m \begin{vmatrix} \\ \mathbf{D} \\ \\ \end{vmatrix} = \begin{vmatrix} \mathbf{U_3} \end{vmatrix} \times \begin{vmatrix} \mathbf{W_3} \end{vmatrix} \times \begin{vmatrix} \mathbf{V_3^T} \end{vmatrix}$$



Lazebnik

# 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 $|D\text{-}MS|^2$

Lazebnik

# 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 $\mathbf{L} = \mathbf{C}\mathbf{C}^\mathsf{T}$ :

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

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

Lazebnik

# 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

Source: M. Hebert

# 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.

# Further Factorization work

Factorization with uncertainty

<span style="color:orange">(Irani & Anandan, IJCV'02)</span>

Factorization for indep. moving objects (now)

<span style="color:orange">(Costeira and Kanade '94)</span>

Factorization for articulated objects (now)

<span style="color:orange">(Yan and Pollefeys '05)</span>

Factorization for dynamic objects (now)

<span style="color:orange">(Bregler et al. 2000, Brand 2001)</span>

Perspective factorization (next week)

<span style="color:orange">(Sturm & Triggs 1996, …)</span>

Factorization with outliers and missing pts.

<span style="color:orange">(Jacobs '97 (affine), Martinek & Pajdla'01 Aanaes'02 (perspective))</span>

Pollefeys

# Structure from motion of multiple moving objects

$$\mathcal{D}^{(i)} \stackrel{\text{def}}{=} \begin{pmatrix} \boldsymbol{p}_{11}^{(i)} & \cdots & \boldsymbol{p}_{1n_i}^{(i)} \\ \cdots & \cdots & \cdots \\ \boldsymbol{p}_{m1}^{(i)} & \cdots & \boldsymbol{p}_{mn_i}^{(i)} \end{pmatrix} \qquad i = 1, \ldots, k,$$

$$\mathcal{D}^{(i)} = \mathcal{M}^{(i)} \mathcal{P}^{(i)} \qquad \mathcal{M}^{(i)} \stackrel{\text{def}}{=} \begin{pmatrix} \mathcal{M}_1^{(i)} & \boldsymbol{o}_1^{(i)} \\ \cdots & \cdots \\ \mathcal{M}_m^{(i)} & \boldsymbol{o}_m^{(i)} \end{pmatrix} \qquad \mathcal{P}^{(i)} \stackrel{\text{def}}{=} \begin{pmatrix} \boldsymbol{P}_1^{(i)} & \cdots & \boldsymbol{P}_{n_i}^{(i)} \\ 1 & \cdots & 1 \end{pmatrix}$$

$$\mathcal{D} \stackrel{\text{def}}{=} (\mathcal{D}^{(1)} \mathcal{D}^{(2)} \ldots \mathcal{D}^{(k)})$$

$$\mathcal{D} = \mathcal{M} \mathcal{P} \qquad \mathcal{M} \stackrel{\text{def}}{=} (\mathcal{M}^{(1)} \mathcal{M}^{(2)} \ldots \mathcal{M}^{(k)}) \qquad \mathcal{P} \stackrel{\text{def}}{=} \begin{pmatrix} \mathcal{P}^{(1)} & 0 & \cdots & 0 & 0 \\ 0 & \mathcal{P}^{(2)} & \cdots & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & 0 & \mathcal{P}^{(k)} \end{pmatrix}$$

$\mathcal{D}$ has rank $4k$ (or less)

**The Shape Interaction Matrix**

$$\mathcal{Q} \stackrel{\text{def}}{=} \mathcal{V}_r \mathcal{V}_r^T = \mathcal{P}^T (B^T B) \mathcal{P} \qquad\qquad C \stackrel{\text{def}}{=} (\mathcal{P}\mathcal{P}^T)^{-1}$$

$$\mathcal{Q} = \begin{pmatrix} \mathcal{P}^{(1)T} C \mathcal{P}^{(1)} & 0 & \cdots & 0 \\ 0 & \mathcal{P}^{(2)T} C \mathcal{P}^{(2)} & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & \cdots & \mathcal{P}^{(k)T} C \mathcal{P}^{(k)} \end{pmatrix}$$

# Structure from motion of multiple moving objects



Pollefeys

# Shape interaction matrix

Shape interaction matrix for articulated objects looses block diagonal structure

$$\mathbf{Q} = \mathbf{V}\mathbf{V}^\top =$$



Costeira and Kanade's approach is not usable for articulated bodies (assumes independent motions)

# Articulated motion subspaces

Motion subspaces for articulated bodies intersect

Joint (1D intersection)

(Yan and Pollefeys, CVPR'05)
(Tresadern and Reid, CVPR'05)

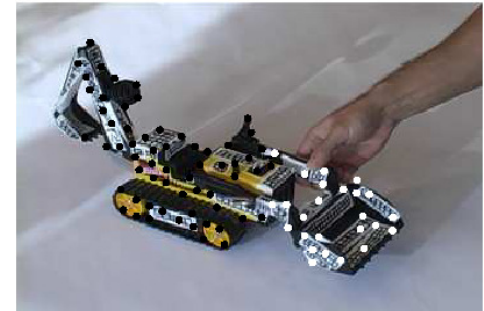$$\mathbf{W} = \begin{bmatrix} \mathbf{R}_1|\mathbf{t}_1 & \mathbf{R}_1\mathbf{R}_1'|\mathbf{t}_1 \\ \vdots & \vdots \\ \mathbf{R}_m|\mathbf{t}_m & \mathbf{R}_m\mathbf{R}_m'|\mathbf{t}_m \end{bmatrix} \begin{bmatrix} \mathbf{S}_1 & 0 \\ 0 & \mathbf{S}_2 \end{bmatrix} \quad \text{(joint=origin)}$$

(rank=8-1)

Hinge (2D intersection)

$$\mathbf{W} = \begin{bmatrix} \mathbf{R}_1|\mathbf{t}_1 & \mathbf{R}_1\mathbf{R}_1'|\mathbf{t}_1 \\ \vdots & \vdots \\ \mathbf{R}_m|\mathbf{t}_m & \mathbf{R}_m\mathbf{R}_m'|\mathbf{t}_m \end{bmatrix} \begin{bmatrix} \mathbf{S}_1 & 0 \\ 0 & \mathbf{S}_2 \end{bmatrix} \quad \text{(hinge=z-axis)}$$

(rank=8-2)

$$\mathbf{R}_i' = \begin{bmatrix} \cos\theta_i & \sin\theta_i & 0 \\ -\sin\theta_i & \cos\theta_i & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Exploit rank constraint to obtain better estimate

Also for non-rigid parts if $\sum_{i=1}^{K} c_i^f = 1$ (Yan & Pollefeys, 06?)

# Results

## Toy truck

### Segmentation



### Intersection

$(\mathbf{W}_1 \cap \mathbf{W}_2)$



## Student

### Segmentation



### Intersection

$(\mathbf{W}_1 \cap \mathbf{W}_2)$



Pollefeys

# Articulated shape and motion factorization

Automated kinematic chain building for articulated & non-rigid obj.

- Estimate principal angles between subspaces
- Compute affinities based on principal angles
- Compute minimum spanning tree



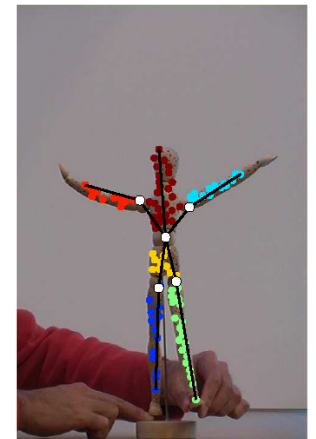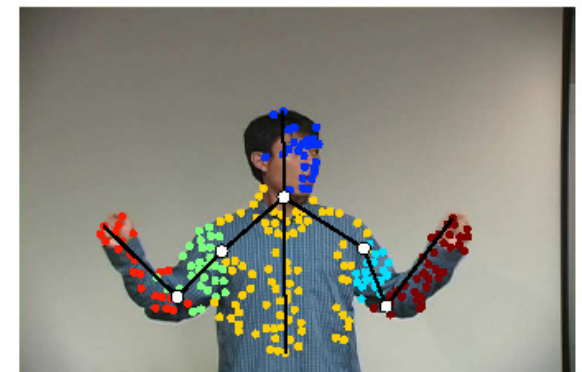|      | larm   | lleg   | hip        | rarm   | body       |
|------|--------|--------|------------|--------|------------|
| rleg | 0.0111 | 0.0007 | **0.0002** | 0.0126 | 0.0006     |
| larm |        | 0.0110 | 0.0060     | 0.0250 | **0.0008** |
| lleg |        |        | **0.0002** | 0.0170 | 0.0006     |
| hip  |        |        |            | 0.0175 | **0.0005** |
| rarm |        |        |            |        | **0.0003** |



|       | luarm  | ruarm  | body       | rlarm      | llarm      |
|-------|--------|--------|------------|------------|------------|
| head  | 0.0015 | 0.0033 | **0.0011** | 0.0035     | 0.0065     |
| luarm |        | 0.0036 | **0.0008** | 0.0058     | **0.0009** |
| ruarm |        |        | **0.0008** | **0.0003** | 0.0145     |
| body  |        |        |            | 0.0018     | 0.0033     |
| rlarm |        |        |            |            | 0.0103     |

Pollefeys

# Structure from motion of deforming objects

(Bregler et al '00;
Brand '01)

Extend factorization approaches to deal with dynamic shapes



Pollefeys

# Representing dynamic shapes
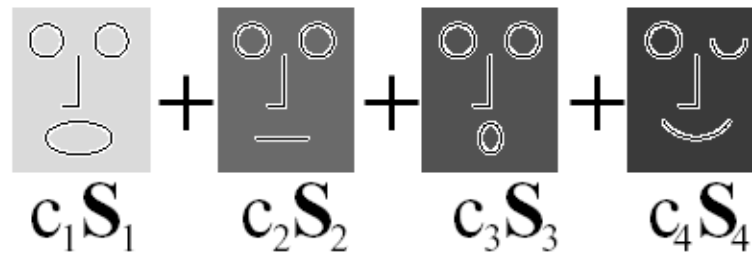
$$S(t) = \sum_{k} c_k(t) S_k$$



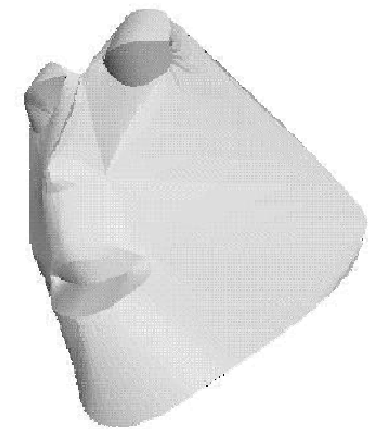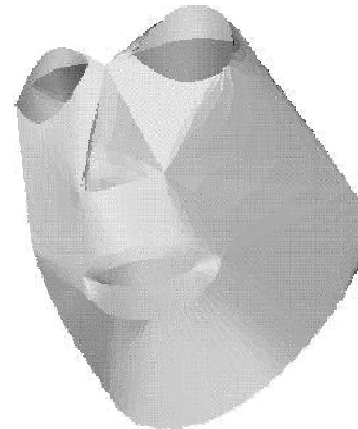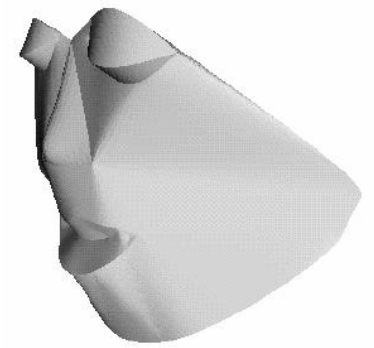$$c_1 S_1 \quad + \quad c_2 S_2 \quad + \quad c_3 S_3 \quad + \quad c_4 S_4$$

(fig. M.Brand)

represent dynamic shape as
varying linear combination of basis shapes

# Results

(Bregler et al '00)

# Dynamic SfM factorization

$$\mathbf{M} \quad \times \quad \mathbf{S} \quad \oplus \quad \mathbf{T} \quad \mathbf{P} \quad =$$

$$\mathbf{c}^\mathsf{T} \otimes \mathbf{R}$$

constraints to be satisfied for M

$$\forall_{\hat{\mathbf{M}}_f \in \hat{\mathbf{M}}} \quad (\underset{3}{\mathrm{vec}}\,\hat{\mathbf{M}}_f^\top)^\top (\underset{3}{\mathrm{vec}}\,\hat{\mathbf{M}}_f^\top) = \tfrac{1}{D}\mathbf{I}_D \otimes ((\underset{3D}{\mathrm{vec}}\,\hat{\mathbf{M}}_f)^\top(\underset{3D}{\mathrm{vec}}\,\hat{\mathbf{M}}_f))$$

constraints to be satisfied for M, use to compute J

$$\hat{\mathbf{M}} = \tilde{\mathbf{M}}\mathbf{J}^{-1} \quad \text{hard!}$$

(different methods are possible,
 not so simple and also not optimal)

# Non-rigid 3D subspace flow

Same is also possible using optical flow in stead of features, also
takes uncertainty into account

$$\begin{bmatrix} \sum_w I_x * I_x & \sum_w I_x * I_y \\ \sum_w I_x * I_y & \sum_w I_x * I_x \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = -\begin{bmatrix} \sum_w I_x * I_t \\ \sum_w I_y * I_t \end{bmatrix}$$

# Results

(Brand '01)



Motion parameters from preschooler frames 0000−0489

Pollefeys

# Results

t=00280

Pollefeys

# 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 $\mathbf{x}_{ij}$



Lazebnik

# 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 **Q**:

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

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

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

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

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

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

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

# Projective factorization

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

<span style="color:red">points (4 × $n$)</span>

<span style="color:red">cameras (3$m$ × 4)</span>

$\mathbf{D} = \mathbf{MS}$ 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)

Lazebnik

# 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*
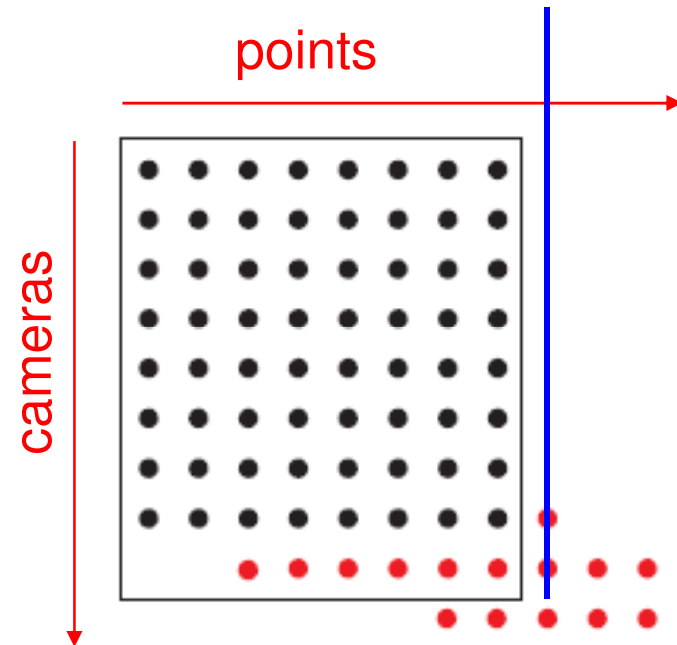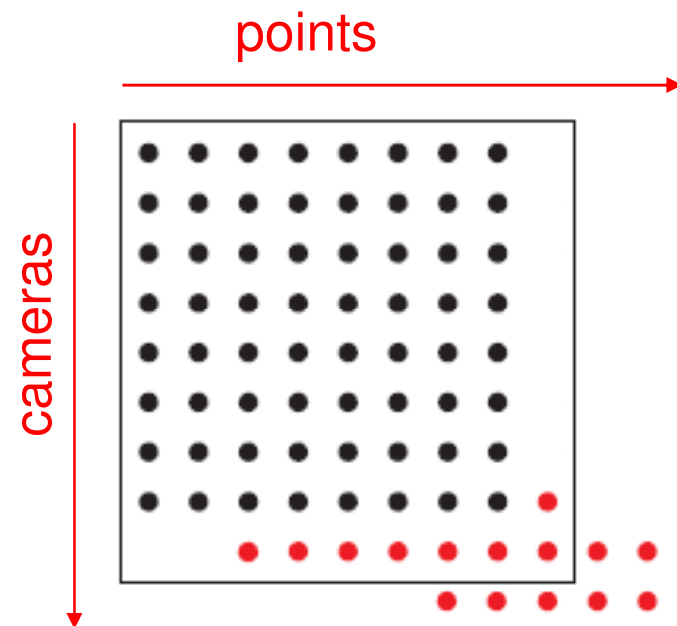


Lazebnik

# 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*

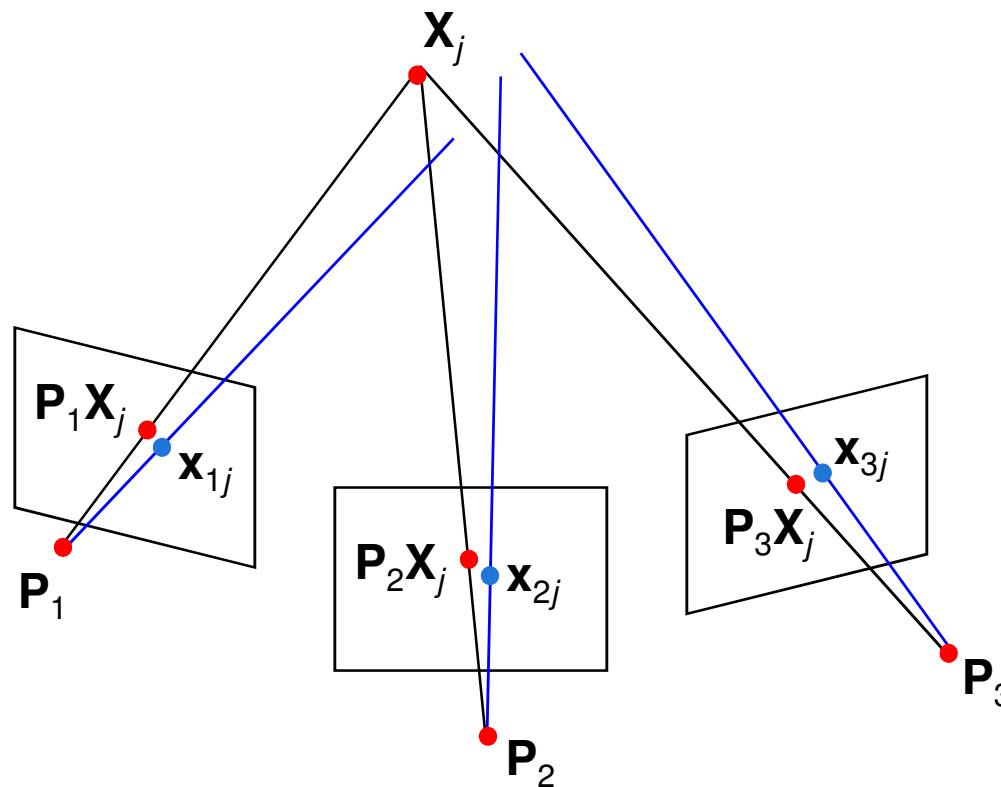- Refine structure and motion: bundle adjustment



points

cameras

Lazebnik

# Bundle adjustment

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

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



Lazebnik

# 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 $\mathbf{P}_i = \mathbf{K} [\mathbf{R}_i | \mathbf{t}_i]$

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

Lazebnik

# 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

# Photo Tourism:
## Exploring Photo Collections in 3D

Noah Snavely

Steven M. Seitz

*University of Washington*

Richard Szeliski

*Microsoft Research*

# Building Rome in a Day

Sameer Agarwal[1,*] Noah Snavely[2] Ian Simon[1] Steven M. Seitz[1] Richard Szeliski[3]

[1]University of Washington [2]Cornell University [3]Microsoft Research

*We present a system that can match and reconstruct 3D scenes from extremely large collections of photographs such as those found by searching for a given city (e.g., Rome) on Internet photo sharing sites. Our system uses a collection of novel parallel distributed matching and reconstruction algorithms, designed to maximize parallelism at each stage in the pipeline and minimize serialization bottlenecks. It is designed to scale gracefully with both the size of the problem and the amount of available computation. We have experimented with a variety of alternative algorithms at each stage of the pipeline and report on which ones work best in a parallel computing environment. Our experimental results demonstrate that it is now possible to reconstruct cities consisting of 150K images in less than a day on a cluster with 500 compute cores.*

| Data set | Images | Cores | Registered | Pairs verified | Pairs found | Time (hrs) | | |
|----------|--------|-------|------------|----------------|-------------|------------|--------------|----------------|
| | | | | | | Matching | Skeletal sets | Reconstruction |
| Dubrovnik | 57,845 | 352 | 11,868 | 2,658,264 | 498,982 | 5 | 1 | 16.5 |
| Rome | 150,000 | 496 | 36,658 | 8,825,256 | 2,712,301 | 13 | 1 | 7 |
| Venice | 250,000 | 496 | 47,925 | 35,465,029 | 6,119,207 | 27 | 21.5 | 16.5 |

Table 1. Matching and reconstruction statisics for the three data sets.

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

# "Rome in a day": Coliseum video

# "Rome in a day": Trevi video

# "Rome in a day": St. Peters video



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

# Slide Credits

- Svetlana Lazebnik
- Marc Pollefeys
- Noah Snaveley & co-authors

# Today: SFM

- SFM problem statement
- Factorization
- Projective SFM
- Bundle Adjustment
- Photo Tourism
- "Rome in a day:

# Roadmap

- Previous: Image formation, filtering, local features, (Texture)…
- Tues: Feature-based Alignment
  - Stitching images together
  - Homographies, RANSAC, Warping, Blending
  - Global alignment of planar models
- Today: Dense Motion Models
  - Local motion / feature displacement
  - Parametric optic flow
- No classes next week: ICCV conference
- Oct 6[th]: Stereo / 'Multi-view': Estimating depth with known inter-camera pose
- **Oct 8[th]: 'Structure-from-motion': Estimation of pose and 3D structure**
  - **Factorization approaches**
  - **Global alignment with 3D point models**