

Stereo matching Calibration

Thursday, Oct 23

Kristen Grauman
UT-Austin

Today

- Some Pset 2 results
- Correspondences, matching for stereo
 - A couple stereo applications
- Camera calibration
- Weak calibration
 - Fundamental matrix
 - 8-point algorithm

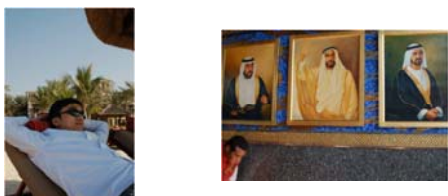
Ming-Jun Chen



Ming-Jun Chen



Ming-Jun Chen



Ming-Jun Chen



Andrew Harp



Wei-Cheng Su



Wei-Cheng Su



Wei-Cheng Su



Wei-Cheng Su



Andy Luong



Andy Luong



Andy Luong

Chia-Sheng Tsai



Chia-Sheng Tsai



Chia-Sheng Tsai



Chia-Sheng Tsai



Amirshahed Mehrtash



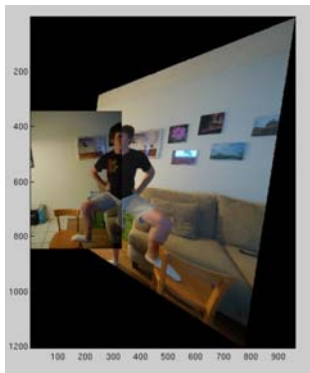
Amirshahed Mehrdash



Amirshahed Mehrdash



Cameron Davison



Jeffrey Dang

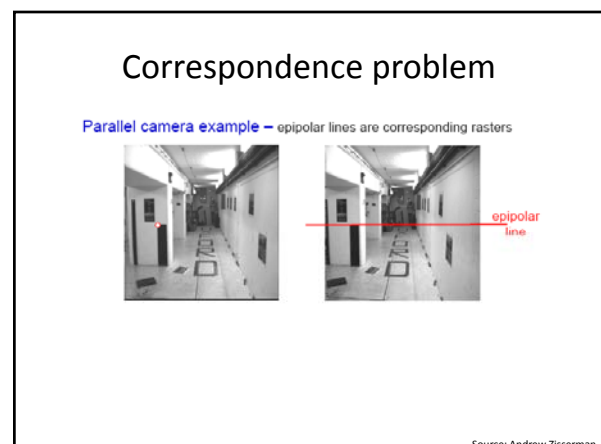
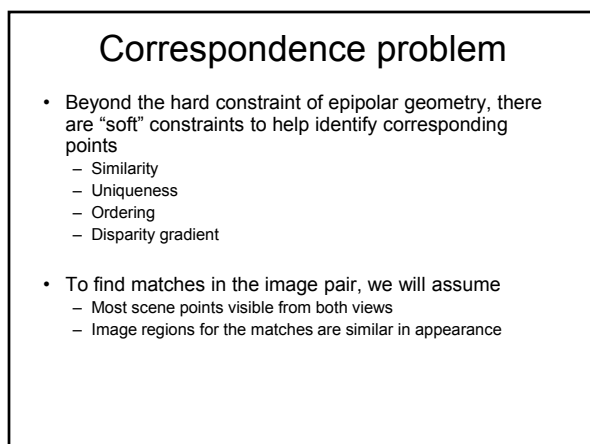
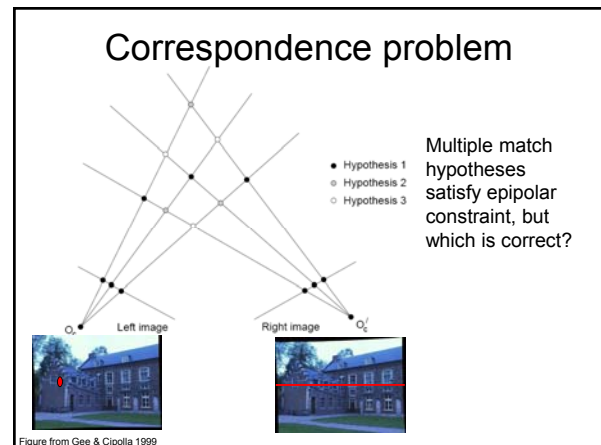
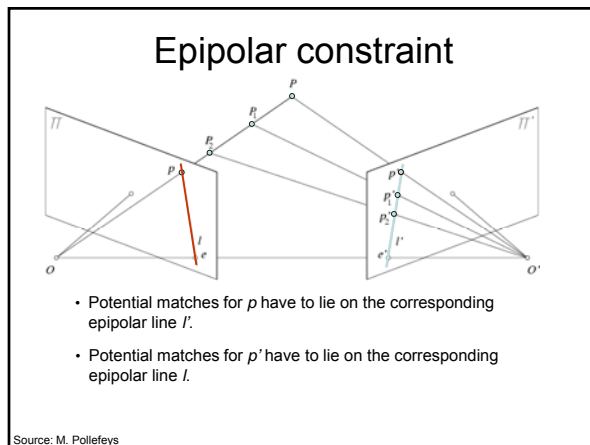
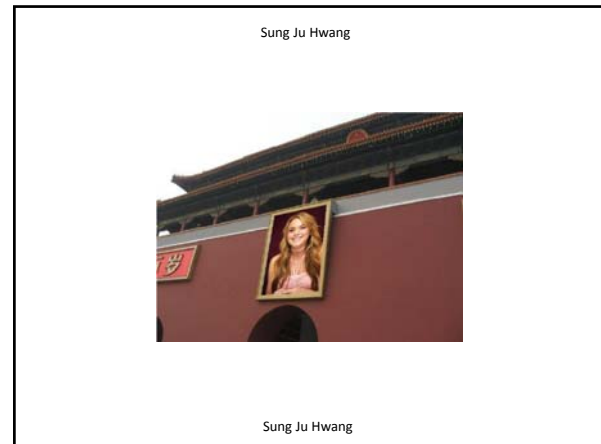


Fei Li

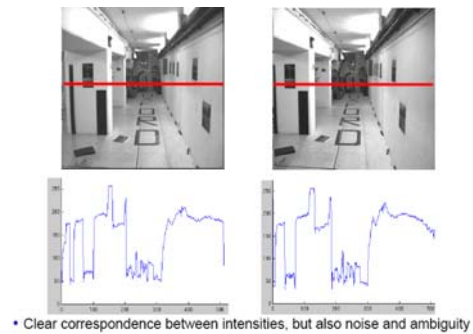


Sung Ju Hwang



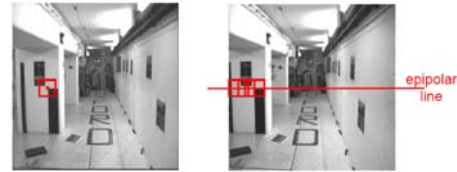


Intensity profiles



Source: Andrew Zisserman

Correspondence problem



Neighborhood of corresponding points are similar in intensity patterns.

Source: Andrew Zisserman

Normalized cross correlation

subtract mean: $A \leftarrow A - \langle A \rangle, B \leftarrow B - \langle B \rangle$

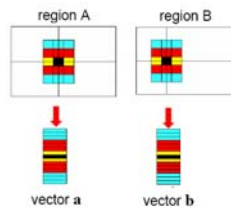
$$NCC = \frac{\sum_i \sum_j A(i,j) B(i,j)}{\sqrt{\sum_i \sum_j A(i,j)^2} \sqrt{\sum_i \sum_j B(i,j)^2}}$$

Write regions as vectors

$A \rightarrow a, B \rightarrow b$

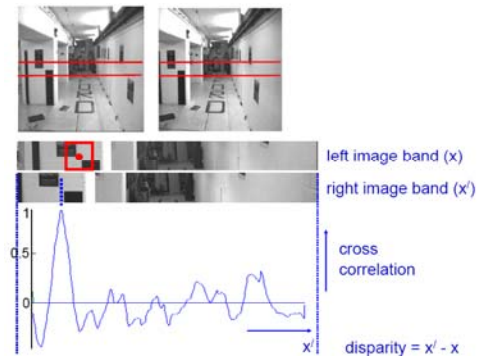
$$NCC = \frac{a \cdot b}{|a| |b|}$$

$$-1 \leq NCC \leq 1$$



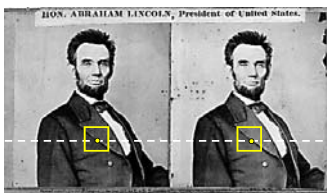
Source: Andrew Zisserman

Correlation-based window matching



Source: Andrew Zisserman

Dense correspondence search



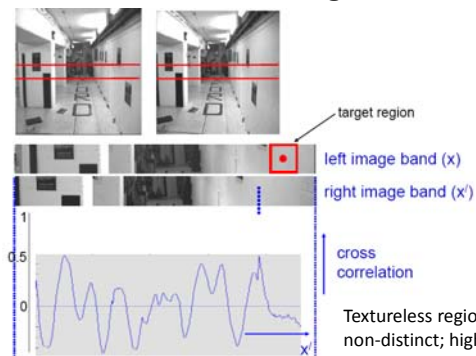
For each epipolar line

For each pixel / window in the left image

- compare with every pixel / window on same epipolar line in right image
- pick position with minimum match cost (e.g., SSD, correlation)

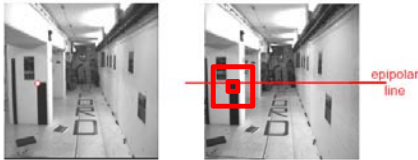
Adapted from Li Zhang

Textureless regions



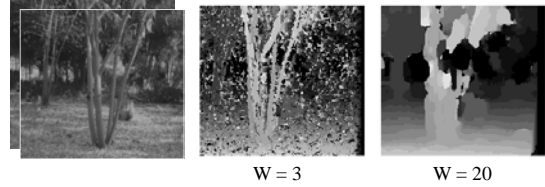
Source: Andrew Zisserman

Effect of window size



Source: Andrew Zisserman

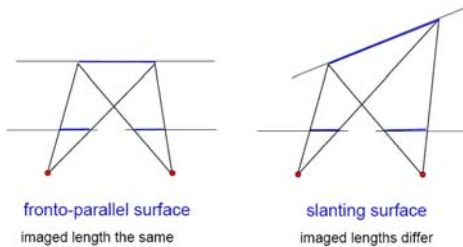
Effect of window size



Want window large enough to have sufficient intensity variation, yet small enough to contain only pixels with about the same disparity.

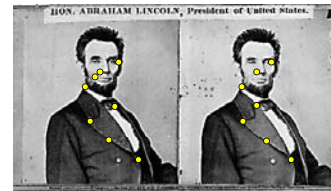
Figures from Li Zhang

Foreshortening effects



Source: Andrew Zisserman

Sparse correspondence search



- Restrict search to sparse set of detected features
- Rather than pixel values (or lists of pixel values) use *feature descriptor* and an associated *feature distance*
- Still narrow search further by epipolar geometry

Correspondence problem

- Beyond the hard constraint of epipolar geometry, there are "soft" constraints to help identify corresponding points
 - Similarity
 - **Uniqueness**
 - Ordering
 - Disparity gradient

Uniqueness

- For opaque objects, up to one match in right image for every point in left image

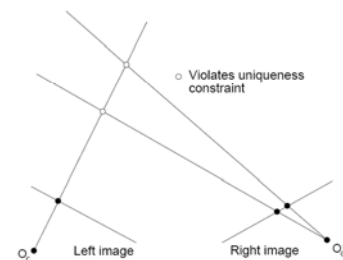


Figure from Gee & Cipolla 1999

Ordering constraint

- Points on **same surface** (opaque object) will be in same order in both views

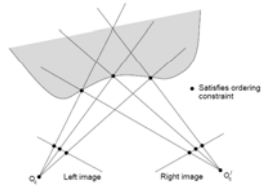
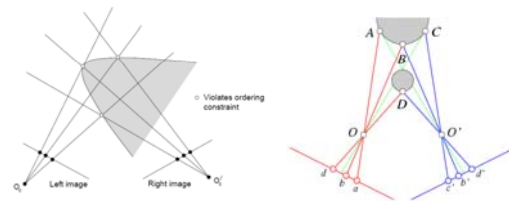


Figure from Gee & Cipolla 1999

Ordering constraint

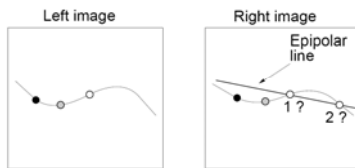
- Won't always hold, e.g. consider transparent object, or an occluding surface



Figures from Forsyth & Ponce

Disparity gradient

- Assume piecewise continuous surface, so want disparity estimates to be locally smooth

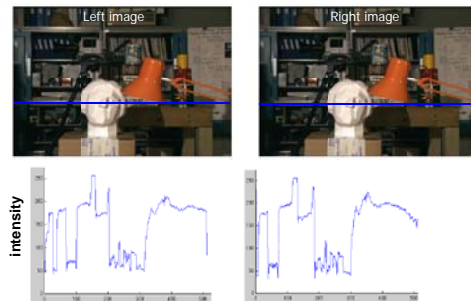


Given matches \bullet and \circ , point \circ in the left image must match point 1 in the right image. Point 2 would exceed the disparity gradient limit.

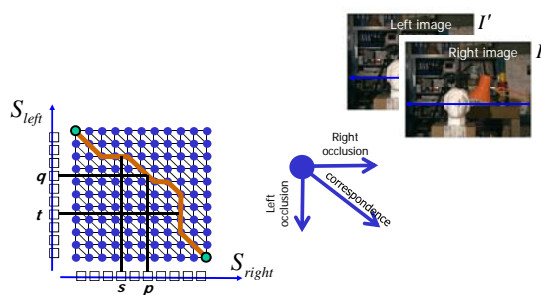
Figure from Gee & Cipolla 1999

Scanline stereo

- Try to coherently match pixels on the entire scanline
- Different scanlines are still optimized independently



"Shortest paths" for scan-line stereo

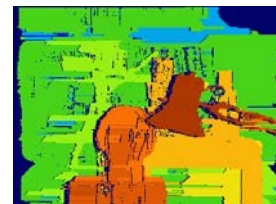


Can be implemented with dynamic programming
Ohta & Kanade '85, Cox et al. '96

Slide credit: Y. Boykov

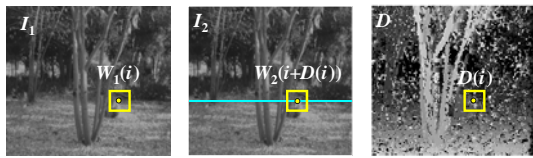
Coherent stereo on 2D grid

- Scanline stereo generates streaking artifacts



- Can't use dynamic programming to find spatially coherent disparities/ correspondences on a 2D grid

Stereo matching as energy minimization

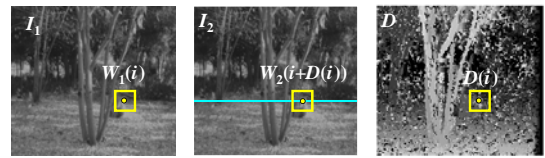


$$E = \alpha E_{\text{data}}(I_1, I_2, D) + \beta E_{\text{smooth}}(D)$$

$$E_{\text{data}} = \sum_i (W_1(i) - W_2(i + D(i)))^2$$

$$E_{\text{smooth}} = \sum_{\text{neighbors } i, j} \rho(D(i) - D(j))$$

Stereo matching as energy minimization



$$E = \alpha E_{\text{data}}(I_1, I_2, D) + \beta E_{\text{smooth}}(D)$$

$$E_{\text{data}} = \sum_i (W_1(i) - W_2(i + D(i)))^2$$

$$E_{\text{smooth}} = \sum_{\text{neighbors } i, j} \rho(D(i) - D(j))$$

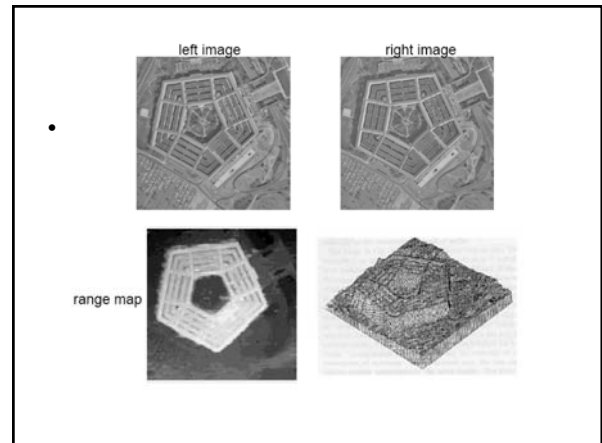
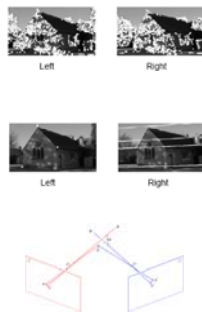
- Energy functions of this form can be minimized using *graph cuts*

Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001

Source: Steve Seitz

Recap: stereo with calibrated cameras

- Image pair
- Detect some features
- Compute **E** from given **R** and **T**
- Match features using the epipolar and other constraints
- Triangulate for 3d structure



Z-keying for virtual reality

- Merge synthetic and real images given depth maps

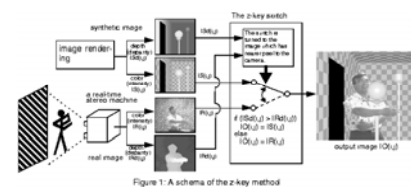
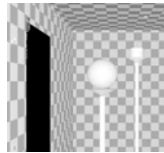


Figure 1: A schema of the z-key method

Kanade et al., CMU, 1995

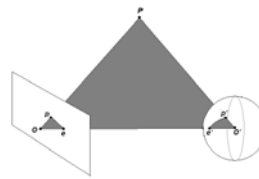
Z-keying for virtual reality



Kanade et al., CMU, 1995

<http://www.cs.cmu.edu/afs/cs/project/stereo-machine/www/z-key.html>

An audio camera & epipolar geometry



Spherical microphone array

Adam O' Donovan, [Ramani Duraiswami](#) and [Jan Neumann](#)
Microphone Arrays as Generalized Cameras for Integrated Audio
Visual Processing, IEEE Conference on Computer Vision and
Pattern Recognition (CVPR), Minneapolis, 2007

First without beamforming

- Adam O' Donovan, [Ramani Duraiswami](#) and [Jan Neumann](#).
Microphone Arrays as Generalized Cameras for Integrated Audio
Visual Processing, IEEE Conference on Computer Vision and
Pattern Recognition (CVPR), Minneapolis, 2007

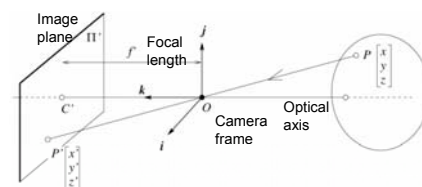
Uncalibrated case

- What if we don't know the camera parameters?

Today

- Some Pset 2 results
- Correspondences, matching for stereo
 - A couple stereo applications
- Camera calibration
- Weak calibration
 - Fundamental matrix
 - 8-point algorithm

Perspective projection



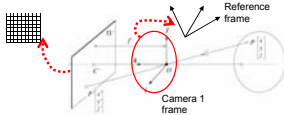
$$(x, y, z) \rightarrow \left(f \frac{x}{z}, f \frac{y}{z}\right)$$

Scene point \rightarrow Image coordinates

Thus far, in camera's reference frame only.

Camera parameters

- **Extrinsic:** location and orientation of camera frame with respect to reference frame
- **Intrinsic:** how to map pixel coordinates to image plane coordinates



Extrinsic camera parameters

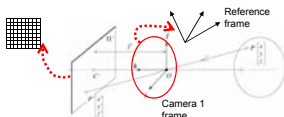
$$\mathbf{P}_c = \mathbf{R}(\mathbf{P}_w - \mathbf{T})$$

\uparrow Camera reference frame \uparrow World reference frame

$$\mathbf{P}_c = (X, Y, Z)^T$$

Camera parameters

- **Extrinsic:** location and orientation of camera frame with respect to reference frame
- **Intrinsic:** how to map pixel coordinates to image plane coordinates



Intrinsic camera parameters

- Ignoring any geometric distortions from optics, we can describe them by:

$$x = -(x_{im} - o_x)s_x$$

$$y = -(y_{im} - o_y)s_y$$

Coordinates of
projected point in
camera reference
frame

Coordinates of
image point in
pixel units

Coordinates of
image center in
pixel units

Effective size of a
pixel (mm)

Camera parameters

- We know that in terms of camera reference frame:

$$x = f \frac{X}{Z} \quad y = f \frac{Y}{Z} \quad \text{and} \quad \mathbf{P}_c = \mathbf{R}(\mathbf{P}_w - \mathbf{T})$$

$$\mathbf{P}_c = (X, Y, Z)^T$$

- Substituting previous eqns describing intrinsic and extrinsic parameters, can relate *pixels coordinates* to *world points*:

$$-(x_{im} - o_x)s_x = f \frac{\mathbf{R}_1 \cdot (\mathbf{P}_w - \mathbf{T})}{\mathbf{R}_3 \cdot (\mathbf{P}_w - \mathbf{T})}$$

\mathbf{R}_i = Row i of rotation matrix

$$-(y_{im} - o_y)s_y = f \frac{\mathbf{R}_2 \cdot (\mathbf{P}_w - \mathbf{T})}{\mathbf{R}_3 \cdot (\mathbf{P}_w - \mathbf{T})}$$

Projection matrix

- This can be rewritten as a matrix product using homogeneous coordinates:

$$\begin{bmatrix} wx_{im} \\ wy_{im} \\ w \\ 1 \end{bmatrix} = \mathbf{M}_{int} \mathbf{M}_{ext} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$x_{im} = wx_{im} / w$
 $y_{im} = wy_{im} / w$

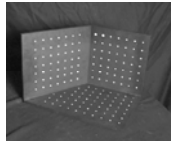
$$\mathbf{M}_{int} = \begin{bmatrix} -f/s_x & 0 & o_x \\ 0 & -f/s_y & o_y \\ 0 & 0 & 1 \end{bmatrix} \quad \mathbf{M}_{ext} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & -\mathbf{R}_1^T \mathbf{T} \\ r_{21} & r_{22} & r_{23} & -\mathbf{R}_2^T \mathbf{T} \\ r_{31} & r_{32} & r_{33} & -\mathbf{R}_3^T \mathbf{T} \end{bmatrix}$$

Calibrating a camera

- Compute intrinsic and extrinsic parameters using observed camera data

Main idea

- Place "calibration object" with known geometry in the scene
- Get correspondences
- Solve for mapping from scene to image: estimate $\mathbf{M} = \mathbf{M}_{\text{int}} \mathbf{M}_{\text{ext}}$



The Opti-CAL Calibration Target Image

Projection matrix

- This can be rewritten as a matrix product using homogeneous coordinates:

$$\begin{pmatrix} wx_{im} \\ wy_{im} \\ w \end{pmatrix} = \underbrace{\mathbf{M}_{\text{int}} \mathbf{M}_{\text{ext}}}_{\mathbf{M}} \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$

\mathbf{P}_w in homog.

$$x_{im} = \frac{\mathbf{M}_1 \cdot \mathbf{P}_w}{\mathbf{M}_3 \cdot \mathbf{P}_w}$$

$$y_{im} = \frac{\mathbf{M}_2 \cdot \mathbf{P}_w}{\mathbf{M}_3 \cdot \mathbf{P}_w}$$

product \mathbf{M} is single **projection matrix** encoding both extrinsic and intrinsic parameters

Let \mathbf{M}_i be row i of matrix \mathbf{M}

Estimating the projection matrix

For a given feature point

$$x_{im} = \frac{\mathbf{M}_1 \cdot \mathbf{P}_w}{\mathbf{M}_3 \cdot \mathbf{P}_w} \longrightarrow 0 = (\mathbf{M}_1 - x_{im} \mathbf{M}_3) \cdot \mathbf{P}_w$$

$$y_{im} = \frac{\mathbf{M}_2 \cdot \mathbf{P}_w}{\mathbf{M}_3 \cdot \mathbf{P}_w} \longrightarrow 0 = (\mathbf{M}_2 - y_{im} \mathbf{M}_3) \cdot \mathbf{P}_w$$

Estimating the projection matrix

$$\begin{aligned} 0 &= (\mathbf{M}_1 - x_{im} \mathbf{M}_3) \cdot \mathbf{P}_w \\ 0 &= (\mathbf{M}_2 - y_{im} \mathbf{M}_3) \cdot \mathbf{P}_w \end{aligned}$$

Expanding this first equation, we have:

$$[m_{11} \ m_{12} \ m_{13} \ m_{14}] - x_{im} [m_{31} \ m_{32} \ m_{33} \ m_{34}] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} = 0$$

$$(X_w m_{11} - x_{im} m_{31}) + (Y_w m_{12} - x_{im} m_{32}) + \dots \\ \dots + (Z_w m_{13} - x_{im} m_{33}) + (m_{14} - x_{im} m_{34}) = 0$$

Estimating the projection matrix

$$\begin{aligned} 0 &= (\mathbf{M}_1 - x_{im} \mathbf{M}_3) \cdot \mathbf{P}_w \\ 0 &= (\mathbf{M}_2 - y_{im} \mathbf{M}_3) \cdot \mathbf{P}_w \end{aligned}$$

$$\begin{pmatrix} X_w & Y_w & Z_w & 1 & 0 & 0 & 0 & 0 & -x_{im} X_w & -x_{im} Y_w & -x_{im} Z_w & -x_{im} \\ 0 & 0 & 0 & 0 & X_w & Y_w & Z_w & 1 & -y_{im} X_w & -y_{im} Y_w & -y_{im} Z_w & -y_{im} \end{pmatrix} \begin{pmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

Estimating the projection matrix

This is true for every feature point, so we can stack up n observed image features and their associated 3d points in single equation:

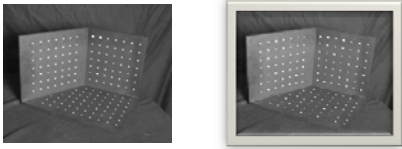
$$\begin{pmatrix} X_w^{(1)} & Y_w^{(1)} & Z_w^{(1)} & 1 & 0 & 0 & 0 & 0 & -x_{im}^{(1)} X_w^{(1)} & -x_{im}^{(1)} Y_w^{(1)} & -x_{im}^{(1)} Z_w^{(1)} & -x_{im}^{(1)} \\ 0 & 0 & 0 & 0 & X_w^{(1)} & Y_w^{(1)} & Z_w^{(1)} & 1 & -y_{im}^{(1)} X_w^{(1)} & -y_{im}^{(1)} Y_w^{(1)} & -y_{im}^{(1)} Z_w^{(1)} & -y_{im}^{(1)} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ X_w^{(n)} & Y_w^{(n)} & Z_w^{(n)} & 1 & 0 & 0 & 0 & 0 & -x_{im}^{(n)} X_w^{(n)} & -x_{im}^{(n)} Y_w^{(n)} & -x_{im}^{(n)} Z_w^{(n)} & -x_{im}^{(n)} \\ 0 & 0 & 0 & 0 & X_w^{(n)} & Y_w^{(n)} & Z_w^{(n)} & 1 & -y_{im}^{(n)} X_w^{(n)} & -y_{im}^{(n)} Y_w^{(n)} & -y_{im}^{(n)} Z_w^{(n)} & -y_{im}^{(n)} \end{pmatrix} \begin{pmatrix} m_{11} \\ m_{12} \\ m_{13} \\ m_{14} \\ m_{21} \\ m_{22} \\ m_{23} \\ m_{24} \\ m_{31} \\ m_{32} \\ m_{33} \\ m_{34} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \dots \\ 0 \\ 0 \end{pmatrix}$$

$\mathbf{P} \mathbf{m} = \mathbf{0}$

Solve for m_i 's (the calibration information)
[F&P Section 3.1]

Summary: camera calibration

- Associate image points with scene points on object with known geometry
- Use together with perspective projection relationship to estimate projection matrix
- (Can also solve for explicit parameters themselves)



When would we calibrate this way?

- Makes sense when geometry of system is not going to change over time
- ...When would it change?

Weak calibration

- Want to estimate world geometry without requiring calibrated cameras
 - Archival videos
 - Photos from multiple unrelated users
 - Dynamic camera system
- Main idea:
 - Estimate epipolar geometry from a (redundant) set of point correspondences between two uncalibrated cameras

Uncalibrated case

For a given camera: $\bar{\mathbf{p}} = \mathbf{M}_{\text{int}} \mathbf{p}$ ← Camera coordinates

So, for two cameras (left and right):

$$\begin{aligned} \mathbf{p}_{(\text{left})} &= \mathbf{M}_{\text{left,int}}^{-1} \bar{\mathbf{p}}_{(\text{left})} \\ \mathbf{p}_{(\text{right})} &= \mathbf{M}_{\text{right,int}}^{-1} \bar{\mathbf{p}}_{(\text{right})} \end{aligned}$$

← Camera coordinates ← Image pixel coordinates

Internal calibration matrices, one per camera

$\mathbf{p}_{(\text{left})} = \mathbf{M}_{\text{left,int}}^{-1} \bar{\mathbf{p}}_{(\text{left})}$
 $\mathbf{p}_{(\text{right})} = \mathbf{M}_{\text{right,int}}^{-1} \bar{\mathbf{p}}_{(\text{right})}$

Uncalibrated case: fundamental matrix

$$\mathbf{p}_{(\text{right})}^T \mathbf{E} \mathbf{p}_{(\text{left})} = 0 \quad \text{From before, the essential matrix } \mathbf{E}.$$

$$(\mathbf{M}_{\text{right,int}}^{-1} \bar{\mathbf{p}}_{\text{right}})^T \mathbf{E} (\mathbf{M}_{\text{left,int}}^{-1} \bar{\mathbf{p}}_{\text{left}}) = 0$$

$$\bar{\mathbf{p}}_{\text{right}}^T (\mathbf{M}_{\text{right,int}}^{-T} \mathbf{E} \mathbf{M}_{\text{left,int}}^{-1}) \bar{\mathbf{p}}_{\text{left}} = 0$$

$$\bar{\mathbf{p}}_{\text{right}}^T \mathbf{F} \bar{\mathbf{p}}_{\text{left}} = 0$$

Fundamental matrix

Fundamental matrix

- Relates pixel coordinates in the two views
- More general form than essential matrix: we remove need to know intrinsic parameters
- If we estimate fundamental matrix from correspondences in pixel coordinates, can reconstruct epipolar geometry without intrinsic or extrinsic parameters

Computing F from correspondences

$$\mathbf{F} = \left(\mathbf{M}_{right,int}^{-T} \mathbf{E} \mathbf{M}_{left,int}^{-1} \right)$$

$$\bar{\mathbf{p}}_{right}^T \mathbf{F} \bar{\mathbf{p}}_{left} = 0$$

- Cameras are uncalibrated: we don't know \mathbf{E} or left or right \mathbf{M}_{int} matrices
- Estimate \mathbf{F} from 8+ point correspondences.

Computing F from correspondences

Each point correspondence generates one constraint on \mathbf{F}

$$\begin{bmatrix} u' & v' & 1 \end{bmatrix} \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

Collect n of these constraints

$$\begin{bmatrix} u'_1 u_1 & u'_1 v_1 & u'_1 & u'_1 v_1 & v'_1 v_1 & v'_1 u_1 & v'_1 & u_1 & v_1 & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = \mathbf{0}$$

Solve for \mathbf{f} , vector of parameters.

Stereo pipeline with weak calibration

- So, where to start with uncalibrated cameras?
 - Need to find fundamental matrix \mathbf{F} and the correspondences (pairs of points $(u',v') \leftrightarrow (u,v)$).

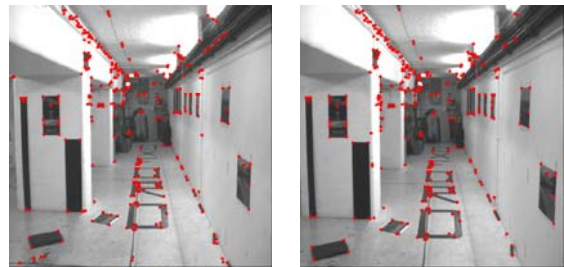


- 1) Find interest points in image (more on this later)
- 2) Compute correspondences
- 3) Compute epipolar geometry
- 4) Refine

Example from Andrew Zisserman

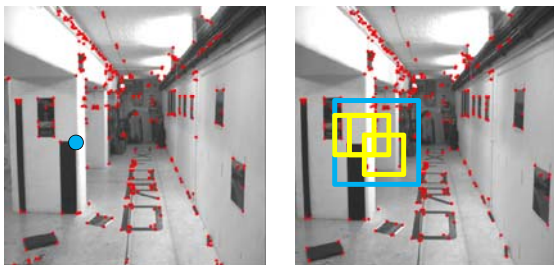
Stereo pipeline with weak calibration

- Find interest points (next week)

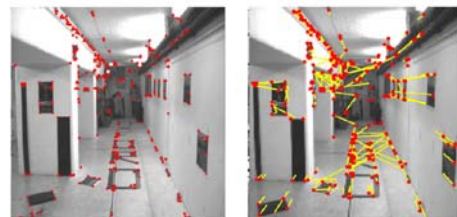


Stereo pipeline with weak calibration

- Match points only using proximity



Putative matches based on correlation search



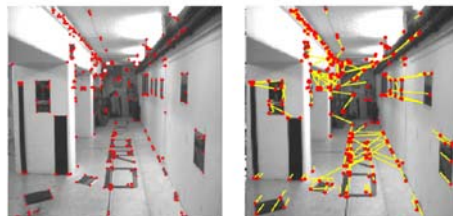
- Many wrong matches (10-50%), but enough to compute \mathbf{F}

RANSAC for robust estimation of the fundamental matrix

- Select random sample of correspondences
- Compute F using them
 - This determines epipolar constraint
- Evaluate amount of support – inliers within threshold distance of epipolar line
- Choose F with most support (inliers)



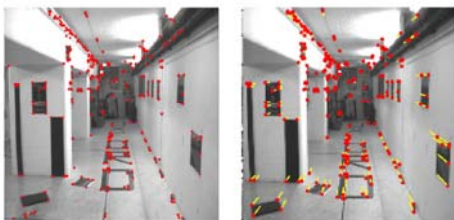
Putative matches based on correlation search



• Many wrong matches (10-50%), but enough to compute F

Pruned matches

- Correspondences consistent with epipolar geometry



- Resulting epipolar geometry



Next

- How to find interest points?
- How to describe local neighborhoods more robustly than with a list of pixel intensities?