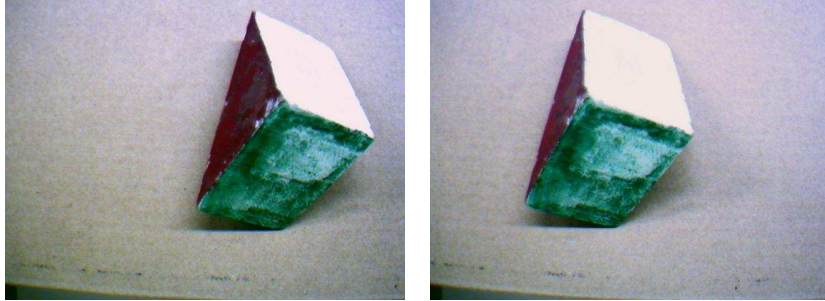


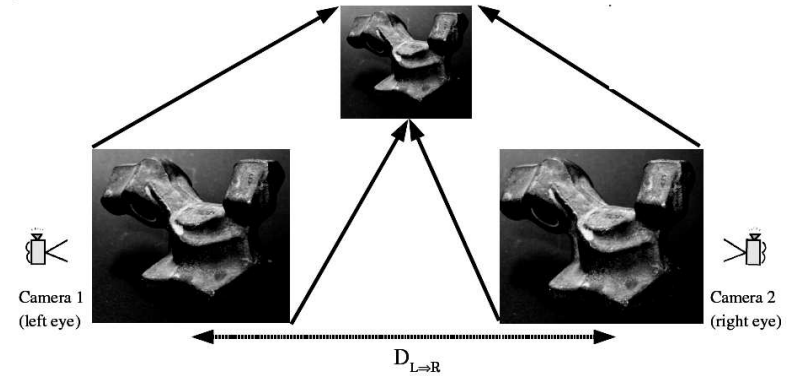
System 6 Introduction

Is there a Wedge in this 3D scene?



Data a stereo pair of images!

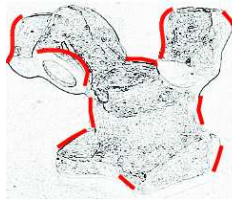
Binocular Stereo Vision



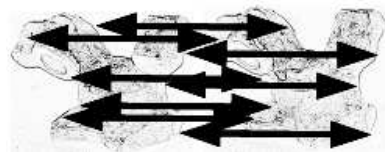
Given two 2D images of an object, how can we reconstruct 3D awareness of it?

Stereo vision - a solution

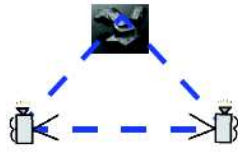
1) **Feature extraction**



2) **Feature matching:**



3) **Triangulation:**

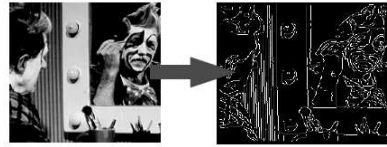


System 6 Overview

3D part recognition using geometric stereo

1. **Feature extraction:**
 - Canny edge detector
 - RANSAC straight line finding
2. **Feature matching:**
 - Stereo correspondence matching lines
3. **Triangulation:**
 - 3D feature position estimation
4. **3D Object recognition:**
 - 3D geometric model
 - Model-data matching using lines
 - 3D pose estimation using lines

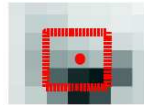
Possible image features



- 1) Edge fragments:
- 2) Edge structures (lines):



- 3) General interest points:



Edge Detector Introduction

- Edge detection: find pixels at large changes in intensity
- Much historical work on this topic in computer vision (Roberts, Sobel)
- Canny edge detector first modern edge detector and still commonly used today
- Edge detection never very accurate process: image noise, areas of low contrast, a question of scale. Humans see edges where none exist.

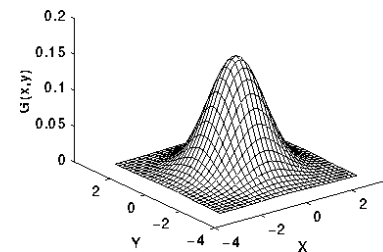
Canny Edge Detector

Four stages:

1. Gaussian smoothing: to reduce noise and smooth away small edges
2. Gradient calculation: to locate potential edge areas
3. Non-maximal suppression: to locate “best” edge positions
4. Hysteresis edge tracking: to locate reliable, but weak edges

Canny: Gaussian Smoothing

Convolve with a 2D Gaussian



$$\frac{1}{273}$$

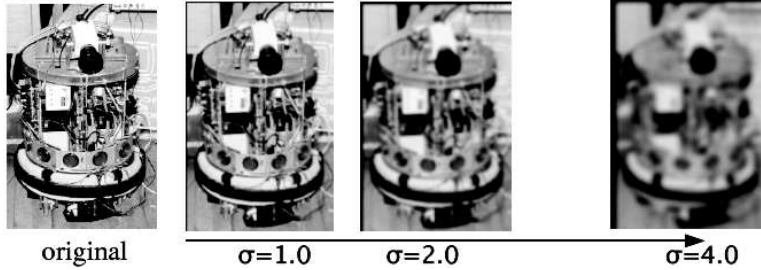
1	4	7	4	1
4	16	26	16	4
7	26	41	26	7
4	16	26	16	4
1	4	7	4	1

Averages pixels with preference near center:
smooths noise without too much blurring of edges

σ of Gaussian controls smoothing explicitly

$$\text{convolution mask}(r, c) = \frac{1}{2\pi\sigma} e^{-\frac{r^2+c^2}{2\sigma^2}}$$

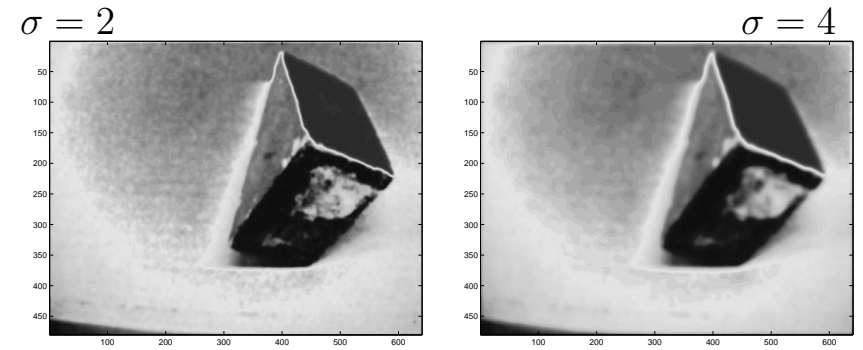
Larger σ gives more smoothing - low pass filter



AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 9

Gaussian Smoothing Examples



AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 10

Conservative Smoothing

Gaussian smoothing inappropriate for salt&pepper/spot noise



Noisy image

Gauss smooth

Conservative

AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 11

Canny: Gradient Magnitude Calculation

$G(r, c)$ is smoothed image

Compute local derivatives in the r and c directions as $G_r(r, c)$, $G_c(r, c)$:

Edge gradient: $\nabla G(r, c) = (G_r(r, c), G_c(r, c))$

AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 12

Gradient magnitude:

$$H(r, c) = \sqrt{G_r(r, c)^2 + G_c(r, c)^2}$$

$$\doteq |G_r(r, c)| + |G_c(r, c)|$$

Gradient direction

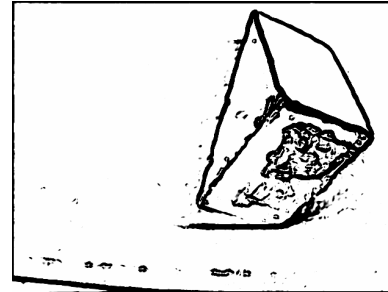
$$\theta(r, c) = \arctan(G_r(r, c), G_c(r, c))$$

$$G_r(r, c) = \frac{\partial G}{\partial r} = \lim_{h \rightarrow 0} \frac{G(r+h, c) - G(r, c)}{h}$$

$$\doteq G(r+1, c) - G(r, c)$$

Gradient Magnitude Examples

$\sigma = 2$



$\sigma = 4$



σ controls amount of smoothing
Smaller σ gives more detail & noise
Larger σ gives less detail & noise

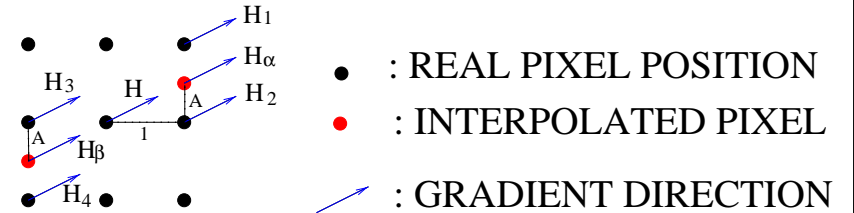
Canny: Non-maximal Suppression

Where exactly is the edge? peak of gradient

Suppress lower gradient magnitude values: need to check **ACROSS** gradient

0	0	3	12	4	0
0	0	6	10	2	0
0	2	8	7	1	0
0	3	11	4	0	0

Estimate gradient magnitudes using gradient direction:



$$A = \frac{|G_r|}{|G_c|}$$

$$H_\alpha = AH_1 + (1 - A)H_2$$

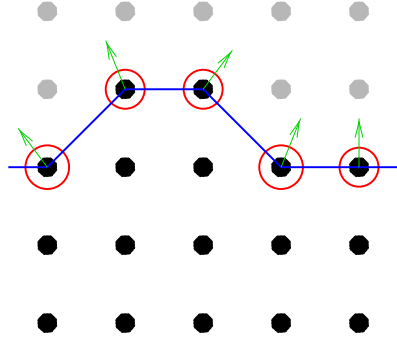
$$H_\beta = AH_4 + (1 - A)H_3$$

Suppress (set to 0) if $H < H_\alpha$ OR $H < H_\beta$

Canny: Hysteresis Tracking

Start edges at certainty: $H > \tau_{start}$

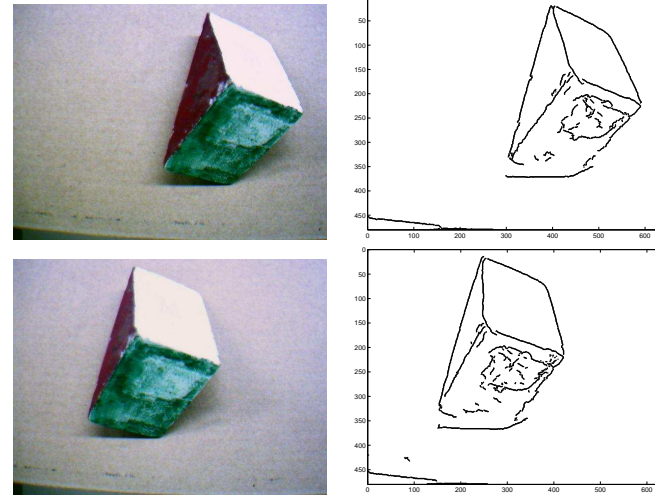
Reduce requirements at connected edges to get weaker edges: $H > \tau_{continue}$



AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 17

Stereo Canny Edges



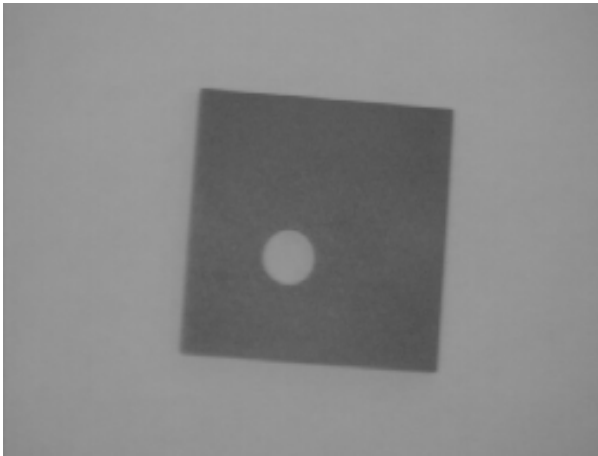
Matlab has Canny: `edge(left, 'canny', [0.08,0.2], 3);`

AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 18

Midlecture Problem

Where might the Canny edge detector find edges in this image?



AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 19

RANSAC: Random Sample and Consensus

Model-based feature detection: features based on some *a priori* model

Works even in much noise and clutter

Tunable failure rate

Assume

- Shape of feature determined by T true data points
- Hypothesized feature is valid if V data points nearby

AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 20

RANSAC Pseudocode

```

for i = 1 : Trials
  Select T data points randomly
  Estimate feature parameters
  if number of nearby data points > V
    return success
  end
end
return failure

```

RANSAC Termination Limit

p_{all-f} is probability of algorithm failing to detect a feature
 p_1 is probability of a data point belonging to a valid feature
 p_d is probability of a data point belonging to same feature
 Algorithm fails if $Trials$ consecutive failures

$$p_{all-f} = (p_{one-f})^{Trials}$$

Success if all needed T random data items are valid

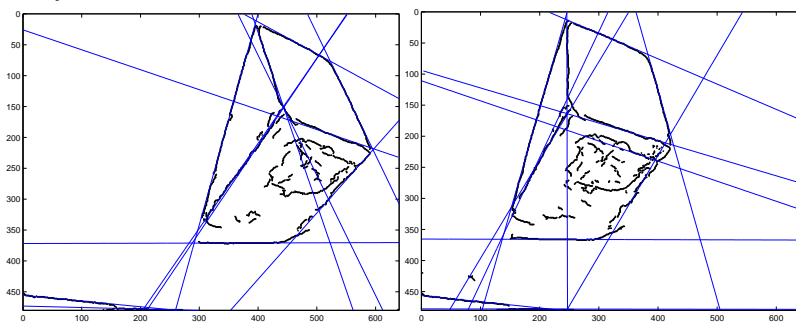
$$p_{one-f} = 1 - p_1(p_d)^{T-1}$$

Solving for expected number of trials:

$$Trials = \frac{\log(p_{all-f})}{\log(1 - p_1(p_d)^{T-1})}$$

RANSAC Line Detection

Line model: infinite line through 2 points ($T = 2$)
 $T = 2$ edge points randomly chosen
 Accept if $V = 80$ edge points within 3 pixels
 $p_{all-f} = 0.001$, $p_1 = 0.1$, $p_d = 0.01$, $Trials = 688$



Mostly accurate lines, but don't know endpoints

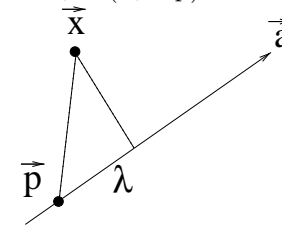
Finding line segments

Some random data crossings

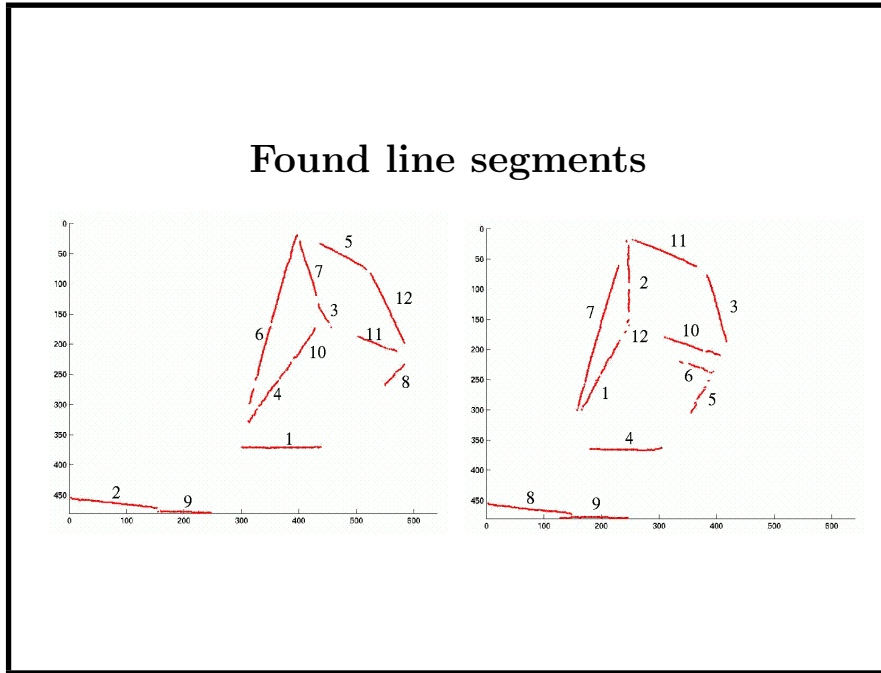
Want to find approximate observed start and end of true segment

1. Project points $\{\vec{x}_i\}$ onto ideal line thru point \vec{p} with direction

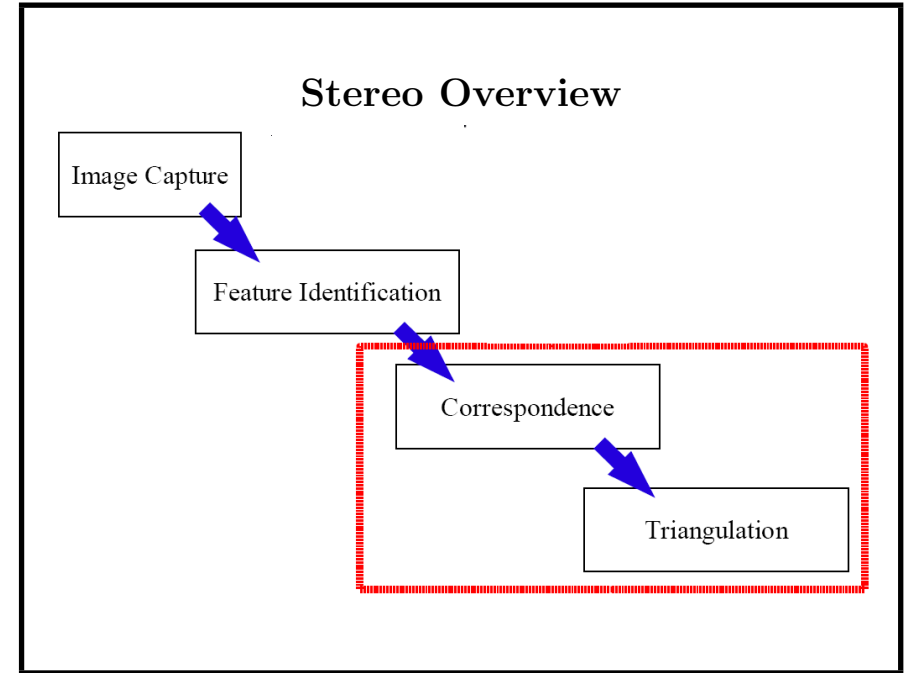
\vec{a} : $\lambda_i = (\vec{x}_i - \vec{p}) \cdot \vec{a}$. Projected point is $\vec{p} + \lambda_i \vec{a}$



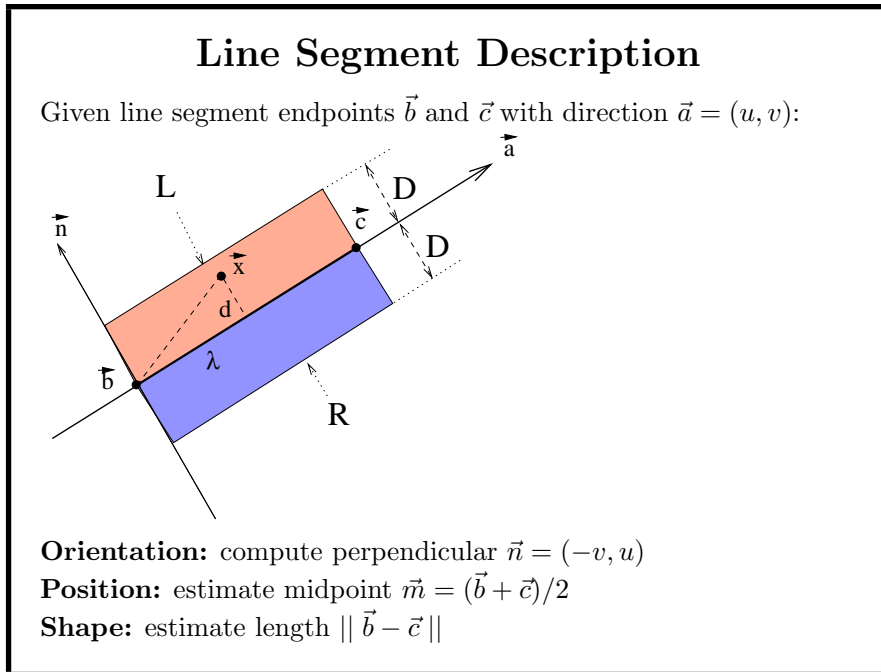
2. Remove points not having 43 neighbor points within 45 pixels distance
3. Endpoints are given by smallest and largest remaining λ_i .



AV: 3D recognition from binocular stereo



AV: 3D recognition from binocular stereo



AV: 3D recognition from binocular stereo

Estimate edge contrast of pixels within D of line in regions L & R:

For any pixel \vec{x}

If $0 \leq (\vec{x} - \vec{b}) \cdot \vec{a} \leq \|\vec{b} - \vec{c}\|$
 If $0 < (\vec{x} - \vec{b}) \cdot \vec{n} \leq D$, put \vec{x} in L
 If $-D \leq (\vec{x} - \vec{b}) \cdot \vec{n} < 0$, put \vec{x} in R

Estimate contrast $g =$ (average brightness of pixels in L)
 $-$ (average brightness of pixels in R)

Description is direction \vec{a} , midpoint \vec{m} , length $\|\vec{b} - \vec{c}\|$ and contrast g .

AV: 3D recognition from binocular stereo

Corresponding Line Properties

Edges	L dir	R dir	L len	R len
L1:R4	(0.00,-1.00)	(0.00,-1.00)	137	125
L5:R11	(0.45,0.89)	(0.38,0.93)	90	119
L6:R7	(0.96,-0.28)	(0.96,-0.29)	291	251
L7:R2	(0.95,0.32)	(1.00,0.00)	93	140
L8:R5	(0.70,-0.71)	(0.86,-0.50)	49	61
L10:R12	(0.81,-0.59)	(0.86,-0.52)	87	82
L11:R10	(-0.33,-0.94)	(-0.29,-0.96)	72	101
L12:R3	(0.89,0.45)	(0.96,0.28)	129	115

Edges	L mid	R mid	L con	R con
L1:R4	(370,369)	(364,242)	-91	-37
L5:R11	(55,477)	(40,309)	28	18
L6:R7	(158,355)	(181,194)	44	55
L7:R2	(74,416)	(90,246)	-117	-44
L8:R5	(250,567)	(278,371)	-66	-31
L10:R12	(208,402)	(205,221)	38	55
L11:R10	(199,536)	(195,358)	141	113
L12:R3	(140,554)	(132,400)	47	23

Binocular Stereo

Goal: build 3D scene description (eg. depth) given two 2D image descriptions

Useful for: obstacle avoidance, grasping, object location

Key principle: triangulation

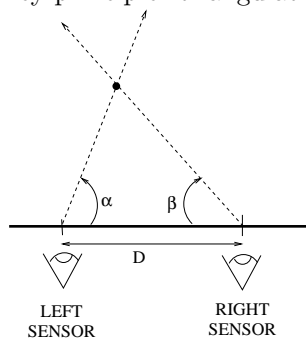
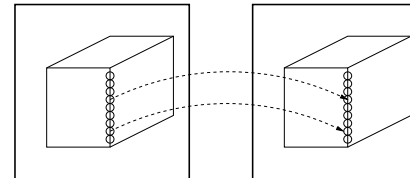


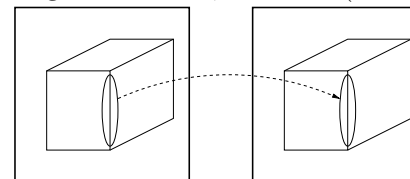
Image Features

What to match?

Edge fragments from edge detector. Usually many.



Larger structures, like lines (RANSAC, vertical for mobile robots)

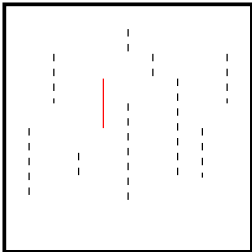


Larger easier to match but harder to get and less dependable
Human visual system thought to work at edge fragment level

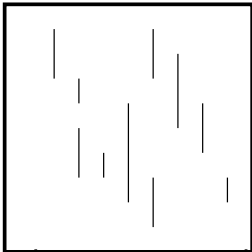
Stereo Correspondence Problem

Which feature in left image matches a given feature in the right?

LEFT



RIGHT



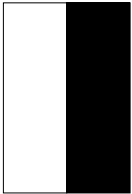
WHICH?

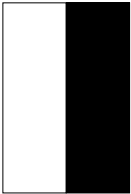
Different pairings give different depth results
Often considered the key problem of stereo


AV: 3D recognition from binocular stereo

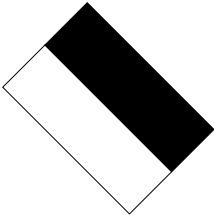
Fisher lecture 6 slide 33

Constraining Matches: Edge Direction


:


OK


BAD

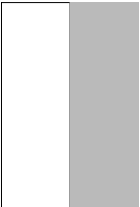

BAD

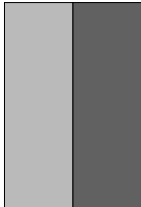
Match features with nearly same orientation

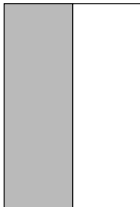
AV: 3D recognition from binocular stereo

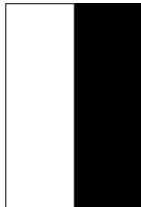
Fisher lecture 6 slide 34

Constraining Matches: Edge Contrast


:


OK


BAD



BAD


Match features with nearly same contrast across edge


AV: 3D recognition from binocular stereo

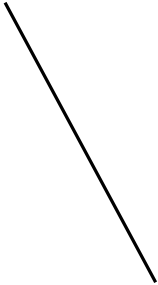
Fisher lecture 6 slide 35

Constraining Matches: Feature Shape


:


OK


BAD


BAD

Match features with nearly same length

AV: 3D recognition from binocular stereo

Fisher lecture 6 slide 36

Constraining Matches: Uniqueness and Smoothness

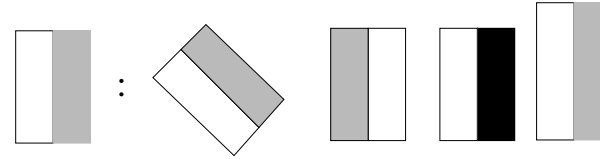
Smoothness: match features giving nearly same depth as neighbors

Uniqueness: a feature in one image can match from the other image:

- 0 - occlusion
- 1 - normal case
- 2+ - transparencies, wires, vines, etc from coincidental alignments

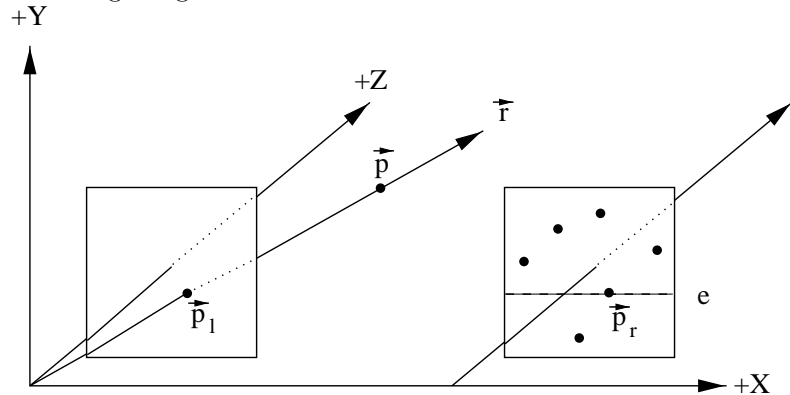
Midlecture Problem

Which stereo correspondence constraint would you use to reject these matches?

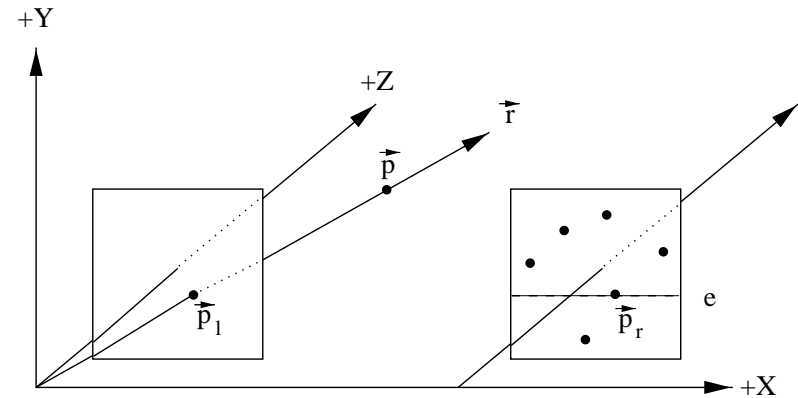


Constraining Matches: Epipolar Geometry

Feature \vec{p}_l in left image lies on a ray \vec{r} thru space.
 \vec{r} projects to an epipolar line e in the right image, along which the matching image feature must lie.



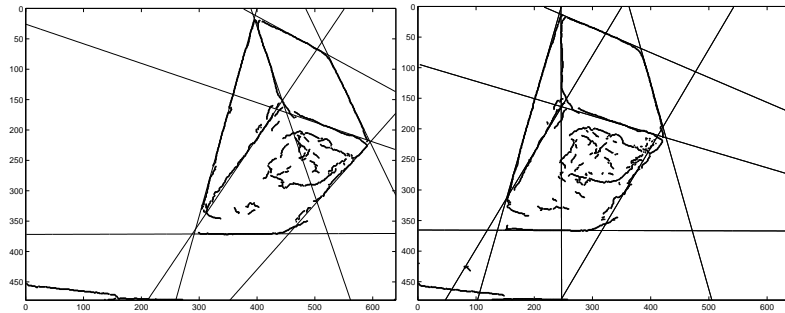
If images are 'rectified', then the epipolar line is an image row.
 Reduces 2D search to 1D search



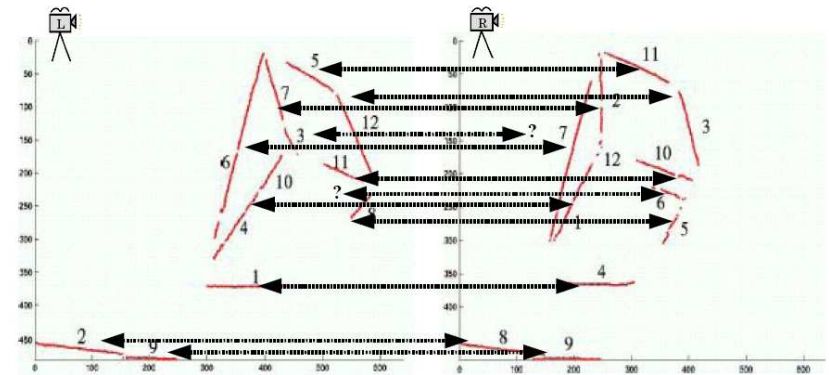
Images are linked by the **Fundamental matrix \mathbf{F}**
 Matched points satisfy $\vec{p}_l^T \mathbf{F} \vec{p}_r = 0$
 (Points are in homogeneous coordinates, \mathbf{F} is 3×3)

Stereo Matching Results

Matched segments:



Stereo Matching Results



Maximally consistent set of matches
Based on local and epipolar constraints

Stereo Matching Code

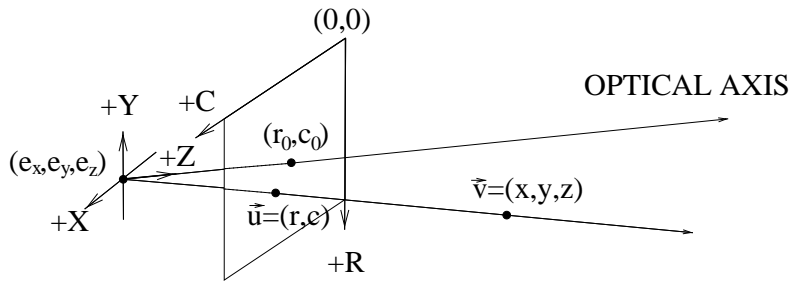
```
% find lines with matching properties
linepairs = zeros(100,2);
numpairs=0;
for l = 1 : llinecount
    for r = 1 : rlinecount
        if abs(llinem(1,1) - rlinem(r,1)) < Tep ... % epipolar
            & abs(llinem(1,2) - rlinem(r,2)) > Tdisl ... % min disp
            & abs(llinem(1,2) - rlinem(r,2)) < Tdish ... % max disp
            & abs(llinel(1) - rlinel(r)) < ...
                Tlen*(llinel(1) + rlinel(r)) ... % length
            & llinea(1,:)*rlinea(r,:) > Tort ... % direction
            & abs(llineg(1) - rlineg(r)) < Tcon % contrast
                numpairs = numpairs+1;
                linepairs(numpairs,1)=l;
                linepairs(numpairs,2)=r;
    end
end
```

```
pairs = 0; % Enforce uniqueness
pairlist = zeros(100,2);
changes = 1;
while changes
    changes = 0;
    for l = 1 : numpairs
        if linepairs(l,1) > 0
            testset = find(linepairs(:,1)==linepairs(l,1));
            if length(testset) == 1
                changes = 1;
                pairs = pairs + 1;
                pairlist(pairs,1) = linepairs(l,1);
                pairlist(pairs,2) = linepairs(l,2);
                linepairs(testset(1),:) = zeros(1,2); % clear taken
            end
        end
    end
end
```

Image Projection Geometry

Pinhole camera model: projects 3D point $\vec{v} = (x, y, z, 1)'$ onto image point $\vec{u}_i = (r_i, c_i, 1)'$: $\lambda \vec{u}_i = P_i \vec{v}$. $i = L, R$.

Notice use of homogeneous coordinates in 2D and 3D



P_i decomposes as

$$P_i = K_i R_i [I | -\vec{e}_i]$$

where R_i : orientation of camera (3 degrees of freedom)

$\vec{e}_i = (e_{xi}, e_{yi}, e_{zi})'$: camera center in world (3 DoF)

K_i : camera intrinsic calibration matrix =

$$\begin{bmatrix} f_i m_{ri} & s_i & r_{0i} \\ 0 & f_i m_{ci} & c_{0i} \\ 0 & 0 & 1 \end{bmatrix}$$

f_i : camera focal length in mm

m_{ri}, m_{ci} : row, col pixels/mm conversion on image plane

r_{0i}, c_{0i} : where optical axis intersects image plane

s_i : skew factor

=====

12 Degrees of Freedom per camera

Calibration

$R_L = R_R = I$ (by definition and controlled camera motions)

$\vec{e}_L = (0, 0, 0)'$ by definition and $\vec{e}_R = (55, 0, 0)'$ by calibrated motion

$s_L = s_R = 0$ (usual for CCD cameras)

$(r_{0L}, c_{0L})' = (r_{0R}, c_{0R})' = (HEIGHT/2, WIDTH/2)'$
(approximate)

$f_L m_{rL} = f_L m_{cL} = f_R m_{rR} = f_R m_{cR} = 832.5$ (same camera + calibration)

$K_L = K_R$ as same camera moved 55 mm to right

Actual Intrinsic Parameters Matrix

Slightly permuted because of different axis labeling and image axis directions

$$\begin{bmatrix} 0 & -f m_r & r_0 \\ f m_c & 0 & c_0 \\ 0 & 0 & 1 \end{bmatrix}$$

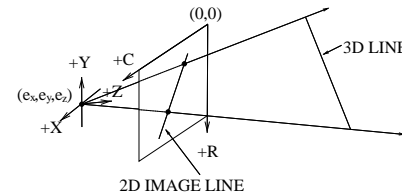
3D Line Calculation

Aim: recovery of 3D line positions

Assume: line successfully matched in L & R images

1. Compute 3D plane that goes through image line and camera origin
2. Compute intersection of 3D planes from 2 cameras (which gives a line)

3D plane passing thru 2D image line



2D image line is:

$$r \cos(\theta) + c \sin(\theta) + b = 0$$

From our projection relations $(\lambda r, \lambda c, \lambda)' = \mathbf{P}(x, y, z, 1)'$ and solving for r and c we get:

$$r = -f m_r \frac{y - e_y}{z - e_z} + r_0$$

$$c = f m_c \frac{x - e_x}{z - e_z} + c_0$$

Both geometric and algebraic analysis give the same result.

Substituting r and c into the line equation, we get the equation for the plane that contains the coplanar points $(x, y, z)'$:

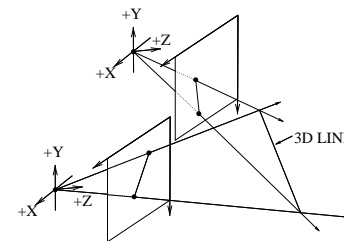
$$\vec{n}' \cdot (x, y, z)' + d = 0$$

$$\vec{n} = (f m_c \sin(\theta), -f m_r \cos(\theta), r_0 \cos(\theta) + c_0 \sin(\theta) + b)'$$

$$d = -\vec{n}'(e_x, e_y, e_z)'$$

Do this for both left and right cameras: gives two planes that intersect at the 3D edge.

3D Plane Intersection \rightarrow 3D Line



Planes: $\vec{x}' \vec{n}_L + d_L = 0$ and $\vec{x}' \vec{n}_R + d_R = 0$

Line given by vector direction \vec{v} through point \vec{a}

Direction: $\vec{v} = \vec{n}_L \times \vec{n}_R$

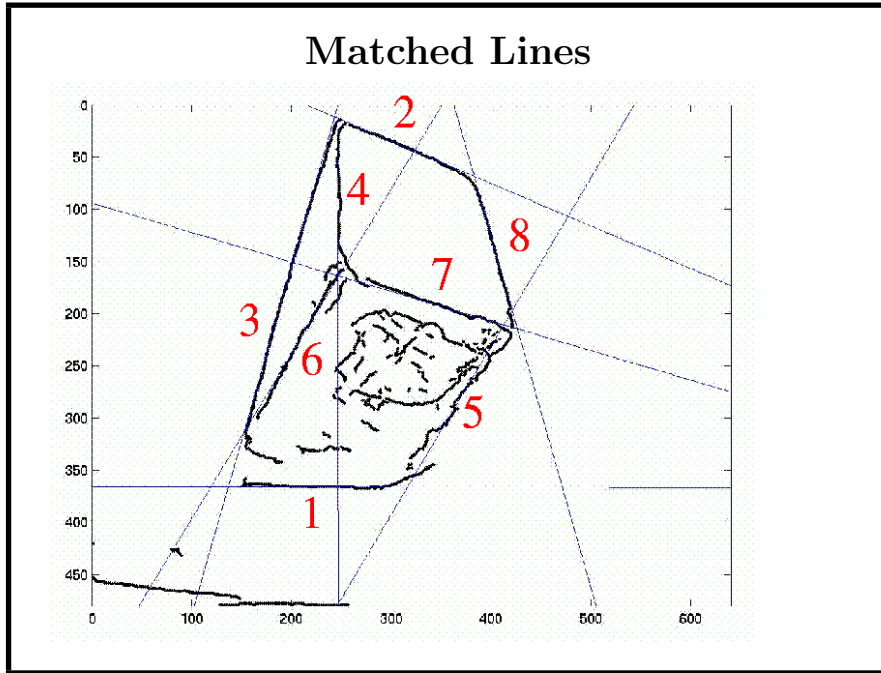
Choose point \vec{a} on line closest to origin, thus $\vec{a}' \vec{v} = 0$

Let: $M = [\vec{v}, \vec{n}_L, \vec{n}_R]'$

Let: $\vec{b} = [0, -d_L, -d_R]'$

Then:

$$\vec{a} = M^{-1} \vec{b}$$

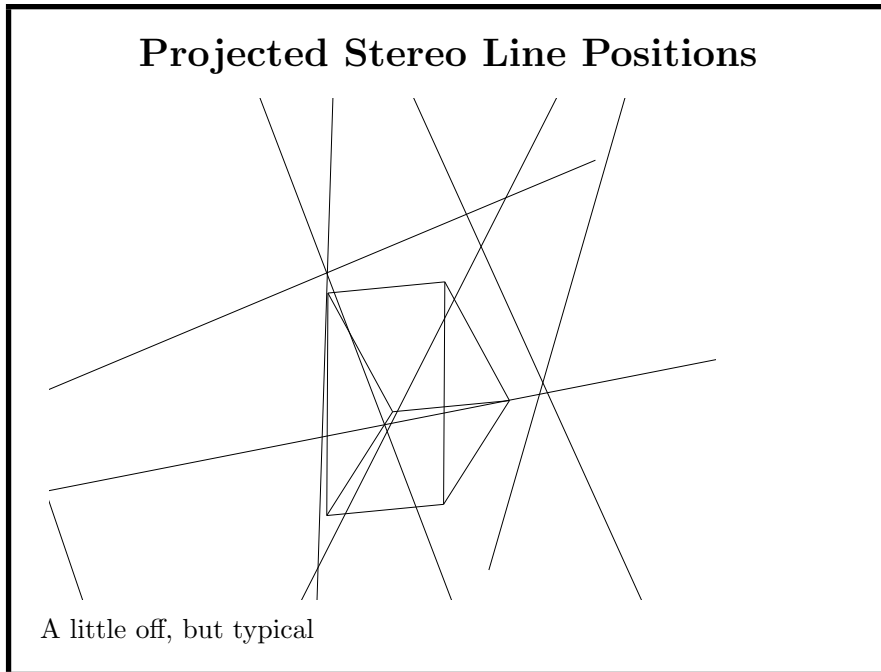


AV: 3D recognition from binocular stereo

3D Line Equations

Number	Pairs	direction	point
1	L1:R4	(-0.82, 0.08,-0.56)	(9.1,2.0,-13.0)
2	L5:R11	(0.61,-0.06, 0.78)	(-125.3,98.6,107.1)
3	L6:R7	(-0.28,-0.95,-0.03)	(0.9,-10.6,294.4)
4	L7:R2	(0.07,-0.62,-0.77)	(48.3,-97.0,82.9)
5	L8:R5	(-0.18,-0.45, 0.87)	(114.8,91.8,72.1)
6	L10:R12	(-0.50,-0.73, 0.44)	(71.5,77.0,208.8)
7	L11:R10	(0.79,-0.20, 0.57)	(-98.4,57.2,154.6)
8	L12:R3	(0.11,-0.69,-0.70)	(110.4,-123.6,140.1)

AV: 3D recognition from binocular stereo

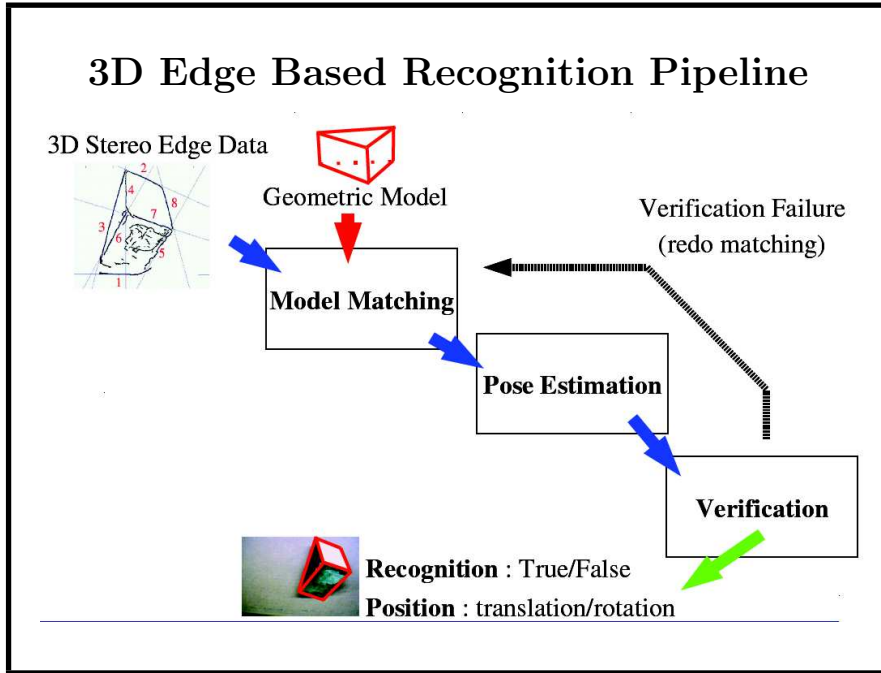


AV: 3D recognition from binocular stereo

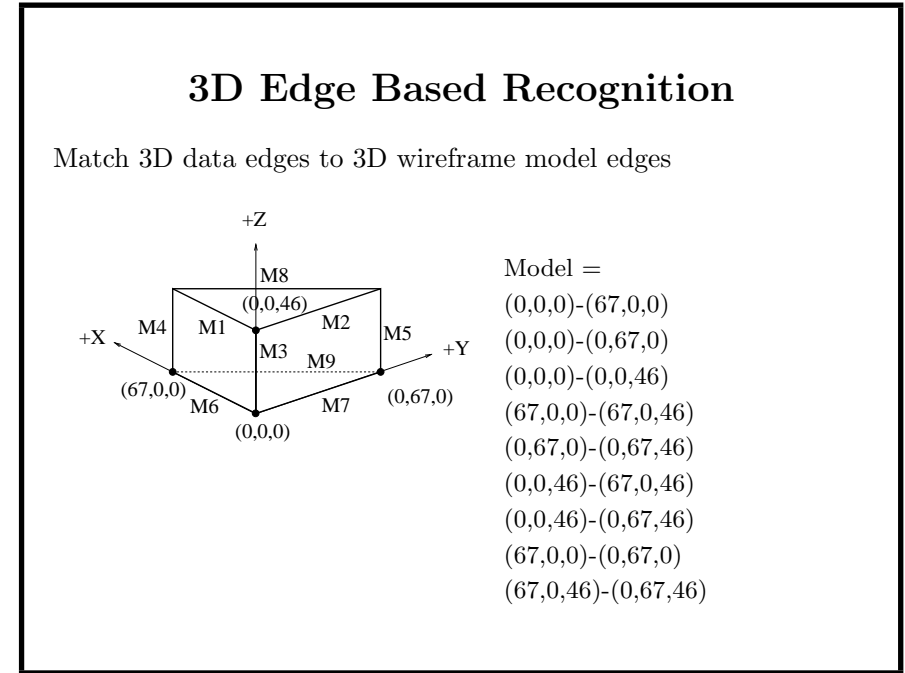
Angles Between Lines

3D line 1	3D line 2	Angle	True
2	3	1.4347	1.57
2	4	1.0221	1.57
2	5	0.9281	1.57
2	6	1.4863	1.57
2	7	0.3023	0.00
2	8	1.1186	1.57
3	4	0.9180	0.71
3	5	1.0964	0.71
3	6	0.5848	0.71
3	7	1.5221	1.57
3	8	0.8457	0.71
4	5	1.1527	1.57
4	6	1.4920	1.57
4	7	1.3026	1.57
4	8	0.1085	0.00
5	6	0.6152	0.00
5	7	1.1060	1.57
5	8	1.2453	1.57
6	7	1.5679	1.57
6	8	1.4276	1.57
7	8	1.3918	1.57

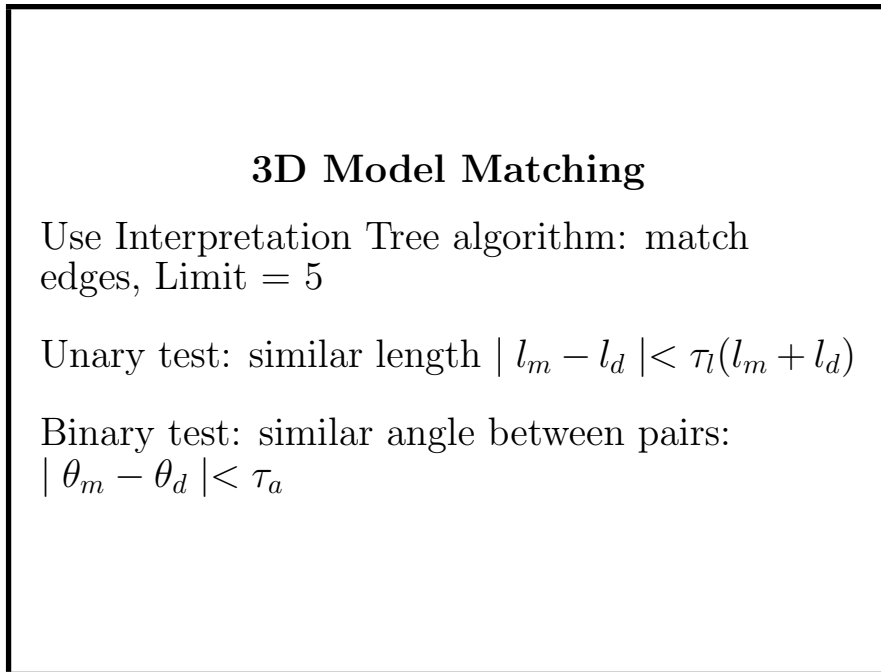
AV: 3D recognition from binocular stereo



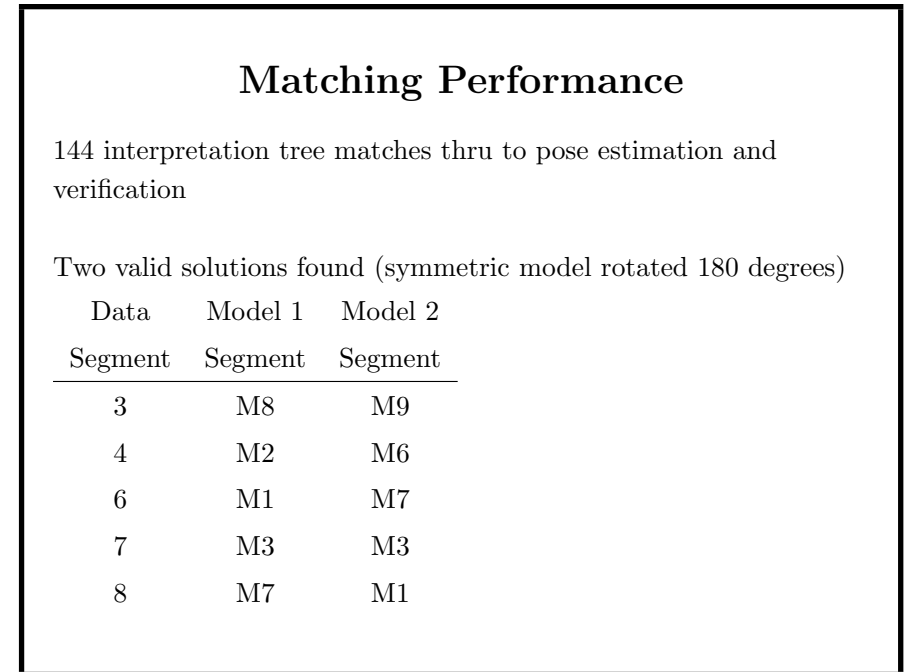
AV: 3D recognition from binocular stereo



AV: 3D recognition from binocular stereo



AV: 3D recognition from binocular stereo



AV: 3D recognition from binocular stereo

3D Pose Estimation

Given: matched line directions $\{(\vec{m}_i, \vec{d}_i)\}$ and points on corresponding lines (but not necessarily same point positions) $\{(\vec{a}_i, \vec{b}_i)\}$

Rotation (matrix R): estimate rotation from matched vectors (same as previous task) except:

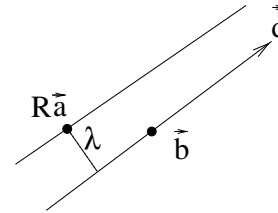
- 1) use line directions instead of surface normals
- 2) don't know which \pm direction for edge correspondence: try both for each matched segment
- 3) if $\det(R) = -1$ then need to flip symmetry
- 4) verify rotation by comparing rotated model and data line orientations

3D Translation Estimation

Given N paired model and data segments, with point \vec{a}_i on model segment i and \vec{b}_i on data segment i

Direction \vec{d}_i of data segment i

Previously estimated rotation R



$\vec{\lambda}_i = R\vec{a}_i + \vec{t} - \vec{b}_i - \vec{d}_i(\vec{d}_i'(R\vec{a}_i + \vec{t} - \vec{b}_i))$ is translation error to minimize

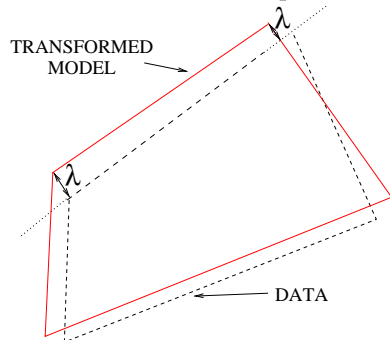
Goal: find \vec{t} that minimizes $\sum_i \vec{\lambda}_i' \vec{\lambda}_i$

$$\text{How: } \mathbf{L} = \sum_i (\mathbf{I} - \vec{d}_i \vec{d}_i') (\mathbf{I} - \vec{d}_i \vec{d}_i')$$

$$\vec{n} = \sum_i (\mathbf{I} - \vec{d}_i \vec{d}_i') (\mathbf{I} - \vec{d}_i \vec{d}_i') (R\vec{a}_i - \vec{b}_i)$$

$$\vec{t} = \mathbf{L}^{-1} \vec{n}$$

Verify translation by comparing perpendicular distance of transformed model endpoints to data line

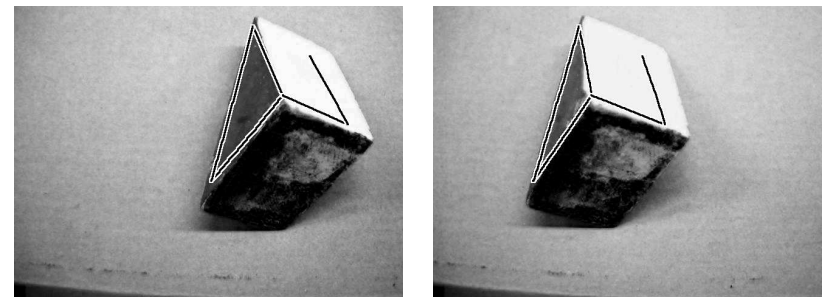


287 verify attempts (2 successes)

Matching Overlay

Two solutions for symmetric model

Left and right image with matched portion of model overlaid:



Calibration a bit off

Discussion

- Hard to find reliable edges/lines, but Canny finds most reasonable edges and RANSAC can put them together for lines
- Given enough stereo correspondence constraints, can get reasonably correct correspondences
- Large features help stereo matching but require more preprocessing
- Stereo geometry easy but needs accurate calibration: not always easy, but now possible to autocalibrate using 8 matched points
- Binocular feature matching stereo gives good 3D at corresponding features, but nothing in between: use scan line stereo?

Dense Depth Data

Problem: have depth only at triangulated feature locations

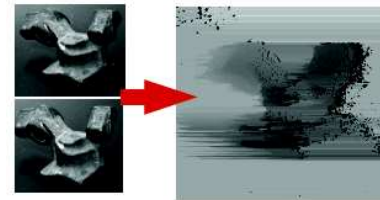
Solution 1: Linear interpolate known values at all other pixels

Solution 2: Correlation-based stereo

Use pixel neighborhoods as features

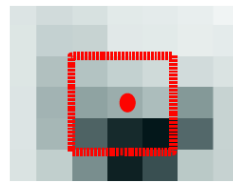
Triangulate depth at every pixel

But needs to find matching pixel - not easy



Correlation based stereo

- Use stereo image pair



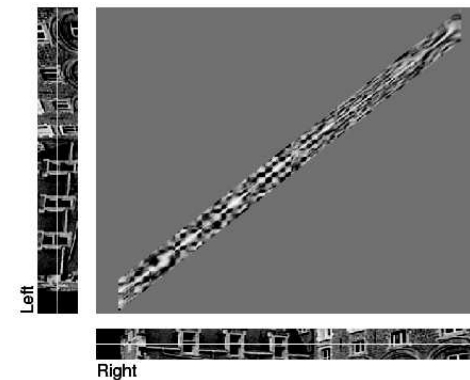
- Features are neighborhoods at each pixel
- Match using similarity metric: SSD - Sum of Squared Differences (of pixel values) of left image at (u, v) to right image at (r, s) :

$$SSD(u, v, r, s) = \sum_{i=-\frac{N}{2}}^{\frac{N}{2}} \sum_{j=-\frac{N}{2}}^{\frac{N}{2}} (L(u+i, v+j) - R(r+i, s+j))^2$$

Finding best match

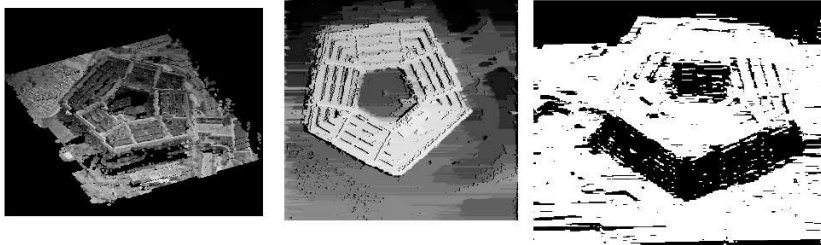
For each scanline on rectified image pair:

1. Build array of all possible matching scores



2. Dynamic programming finds lowest cost path (bright line thru middle of array above - optimisation problem)

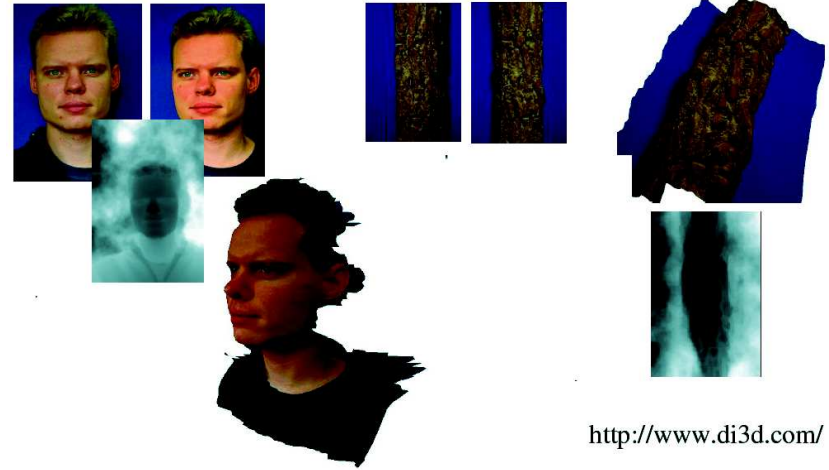
Dense Stereo Results



Technique = [Cox et al. 1996]

AV: 3D recognition from binocular stereo

Commercial Dense Stereo Results



<http://www.di3d.com/>

AV: 3D recognition from binocular stereo

What We Have Learned

Feature based stereo:

- Feature detection: Canny & RANSAC
- Stereo matching: local & epipolar constraints
- Triangulation & 3D: epipolar geometry
- Recognition: IT algorithm & verification

Correlation based stereo: similar but using pixel neighbourhoods

AV: 3D recognition from binocular stereo