

---

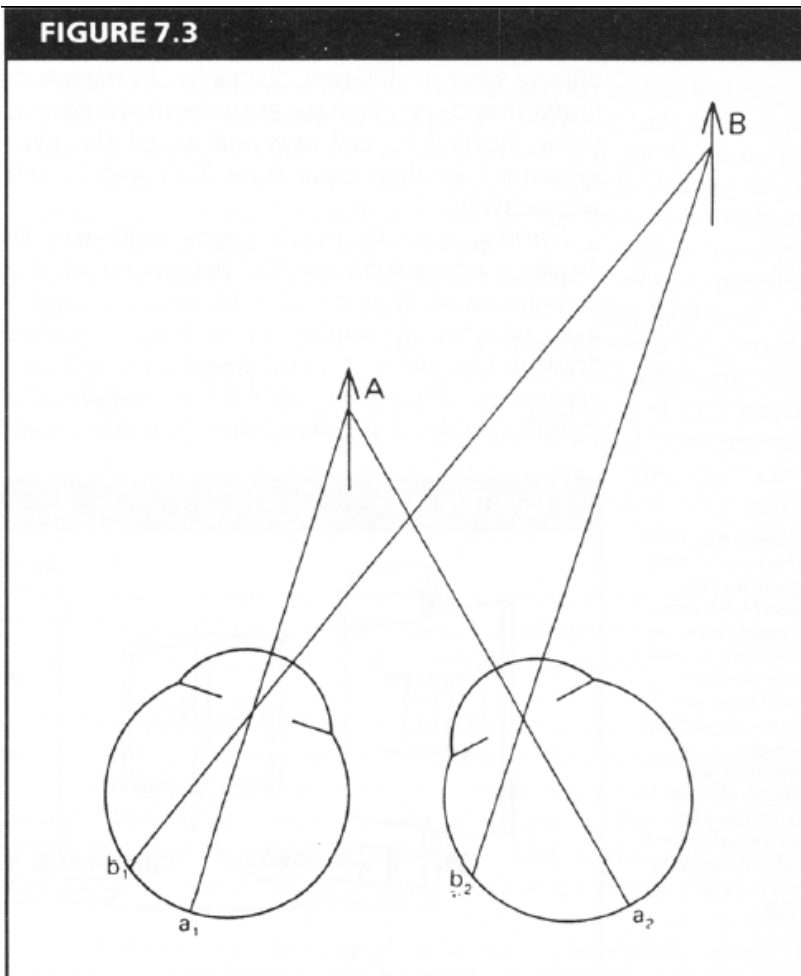
# Stereopsis

CS 554 – Computer Vision

Pinar Duygulu

Bilkent University

# Disparity

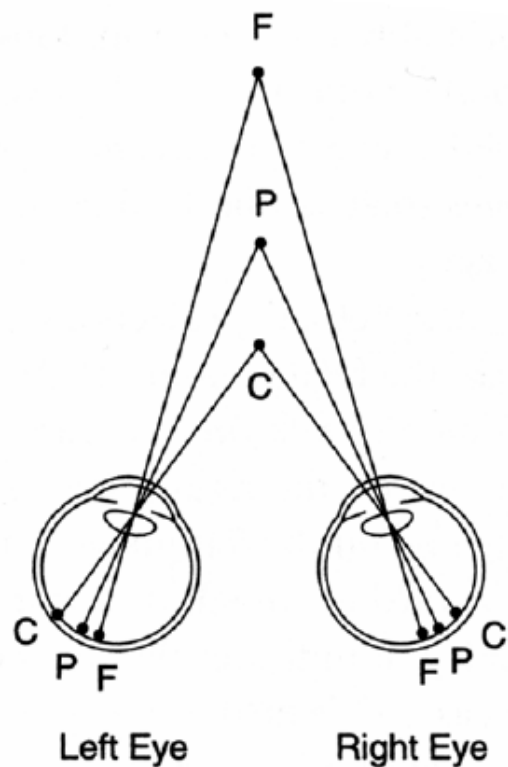


Disparity occurs when  
Eyes verge on one object;  
Others appear at different  
Visual angles

From Bruce and Green, Visual Perception,  
Physiology, Psychology and Ecology

Adapted from David Forsyth, UC Berkeley

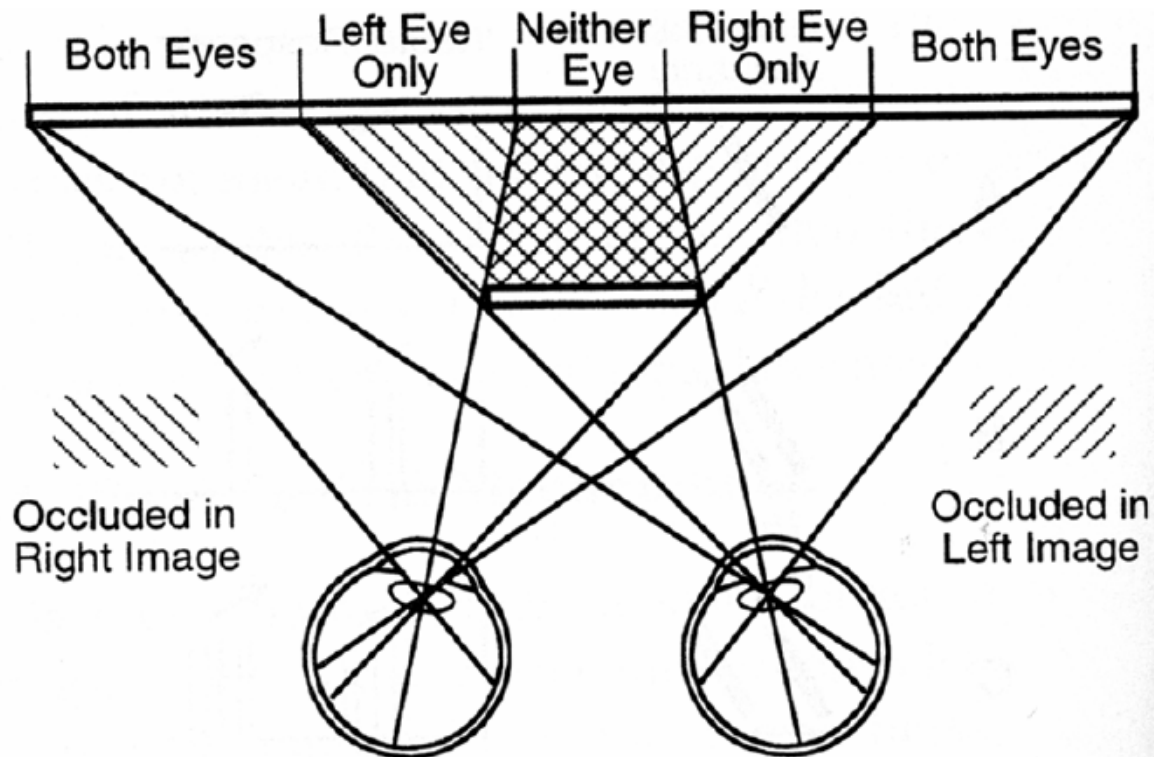
# Disparity



**Figure 5.3.2** Crossed versus uncrossed binocular disparity. When a point *P* is fixated, closer points (such as *C*) are displaced outwardly in crossed disparity, whereas farther points (such as *F*) are displaced inwardly in uncrossed disparity.

From Palmer, “Vision Science”, MIT Press

# Disparity



**Figure 5.3.23** Da Vinci stereopsis. Depth information also arises from the fact that certain parts of one retinal image have no corresponding parts in the other image. (See text for details.)

From Palmer, “Vision Science”, MIT Press

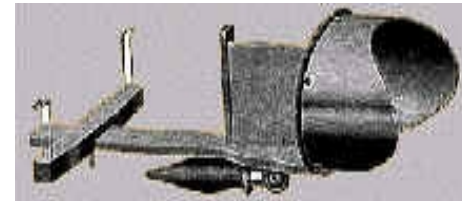
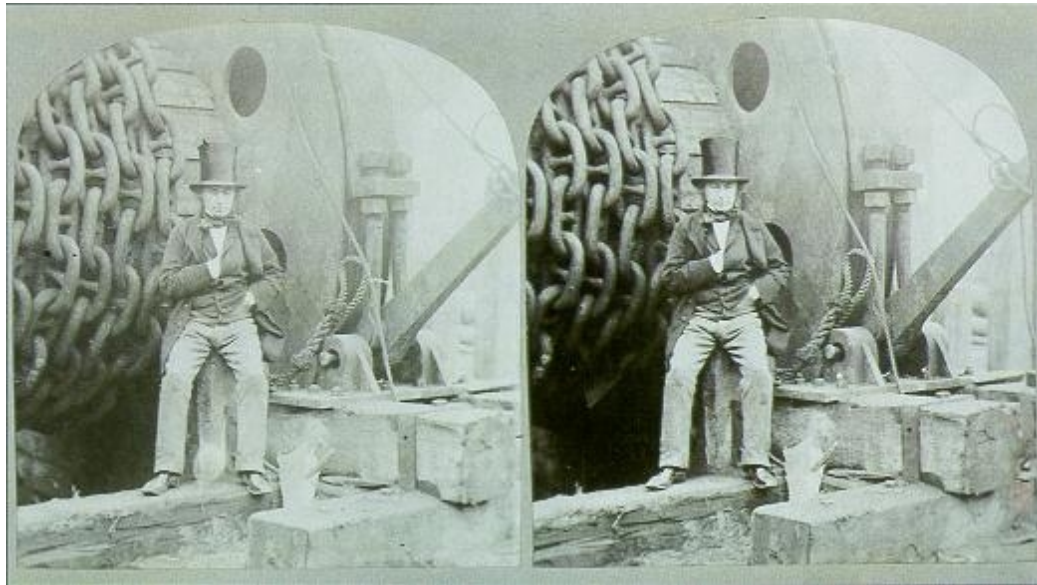
Adapted from David Forsyth, UC Berkeley

# Stereo Vision

---

- The whole process is called **stereo vision** and it is derived from the Greek word '*stereos*' which means form or solid i.e. having three dimensions.
- **Stereoscopy** is the science by which two photographs of the same object taken at slightly different angles are viewed together, giving an impression of depth and solidity as in ordinary human vision.
- **Stereo photography** is the art of taking two pictures of the same subject from two slightly different viewpoints and displaying them in such a way that each eye sees only one of the images.

# Stereo photography



- Capturing the image on film requires the photographer to take two pictures from slightly different viewpoints.
- In order to view the captured photographs, the images have to be displayed in such a way that each of the viewer's eyes sees only one image.

# Anaglyph



+



=



Left Eye Image  
(Red channel only)

Right Eye Image  
(Red channel removed)

Anaglyph  
(Left & Right  
images overlaid)

- Requires the viewer to wear glasses with red and green/cyan lenses.
- The left image has the blue and green colour channels removed to leave a purely red picture while the right image has the red channel removed.
- The two images are superimposed into one picture which produces a picture very like the original with a red and cyan fringes around objects where the stereo separation produces differences in the original images.
- The red and cyan lenses in the glasses let the eyes separate the two superimposed images into their individual components which the brain then combines to form a 3D-image.



# Freeview

---

- Free Viewing, the eyes should not converge but look parallel as if the image being looked at is in the distance.
- The brain is fooled into thinking that it has two separate images and creates a 3-D visualisation.
- Single Image Random Dots Stereogram (SIRDS)
- Single Image Stereogram (SIS)
- "Magic Eye" pictures are created by computer and rely on the fact that the brain depends on matching vertical edges to synchronise the left and right images.
- The picture is made up of columns of patterns, which vary slightly across the picture.
- The brain interprets the columns as left and right pairs and the slight differences between each column define the subject e.g. the fish.

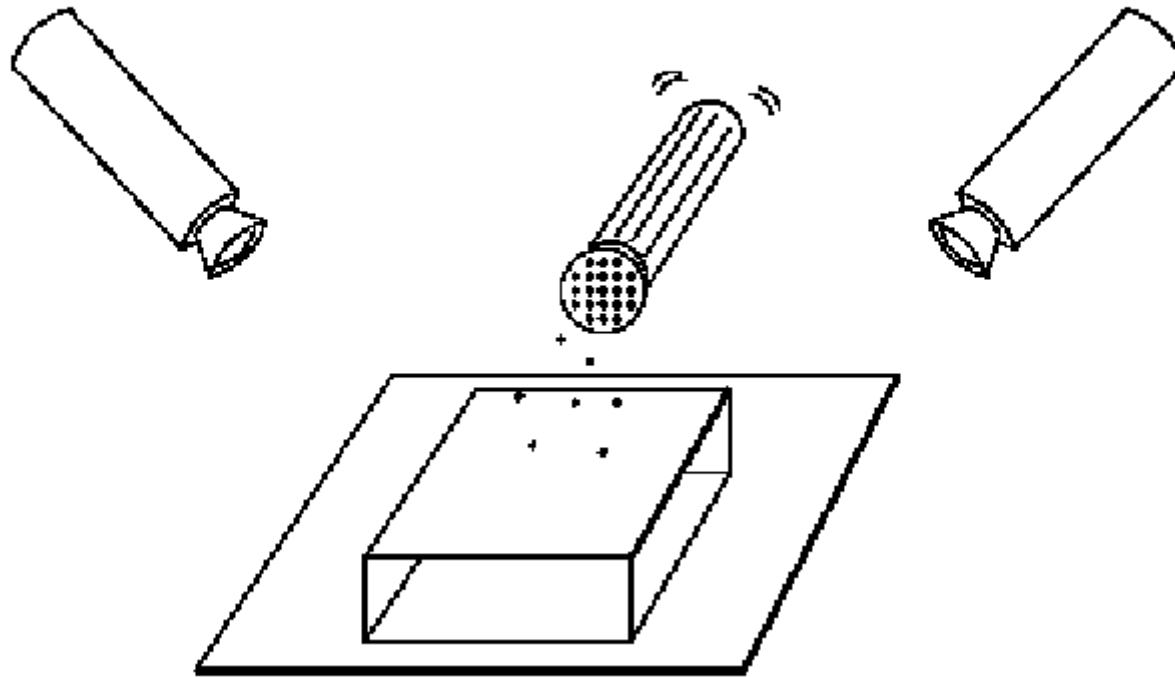






# Random dot stereograms

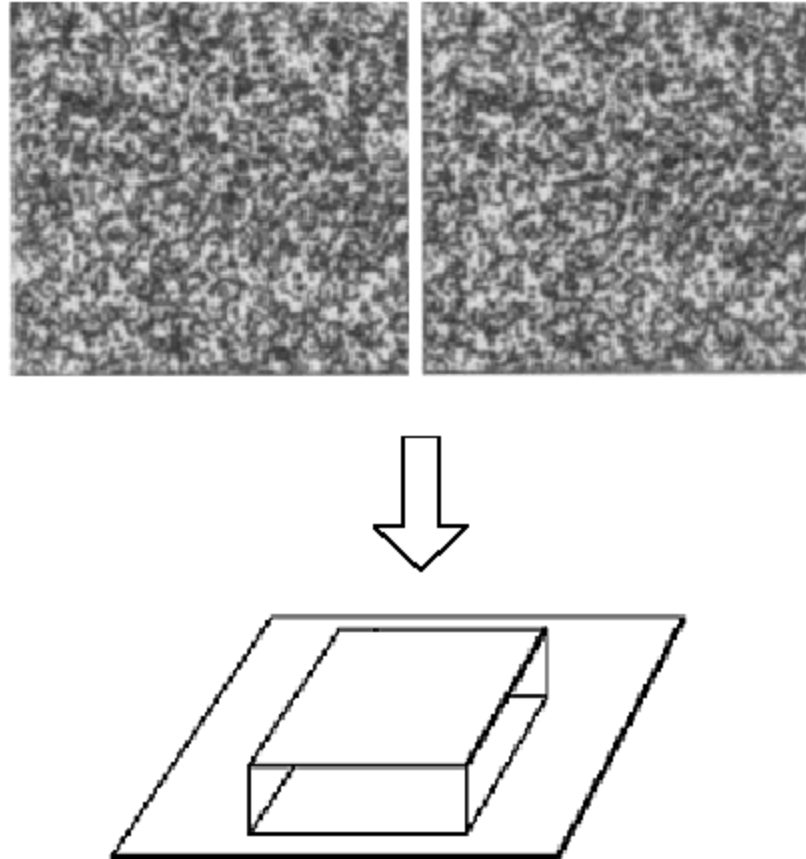
---



Adapted from David Forsyth, UC Berkeley

# Random dot stereograms

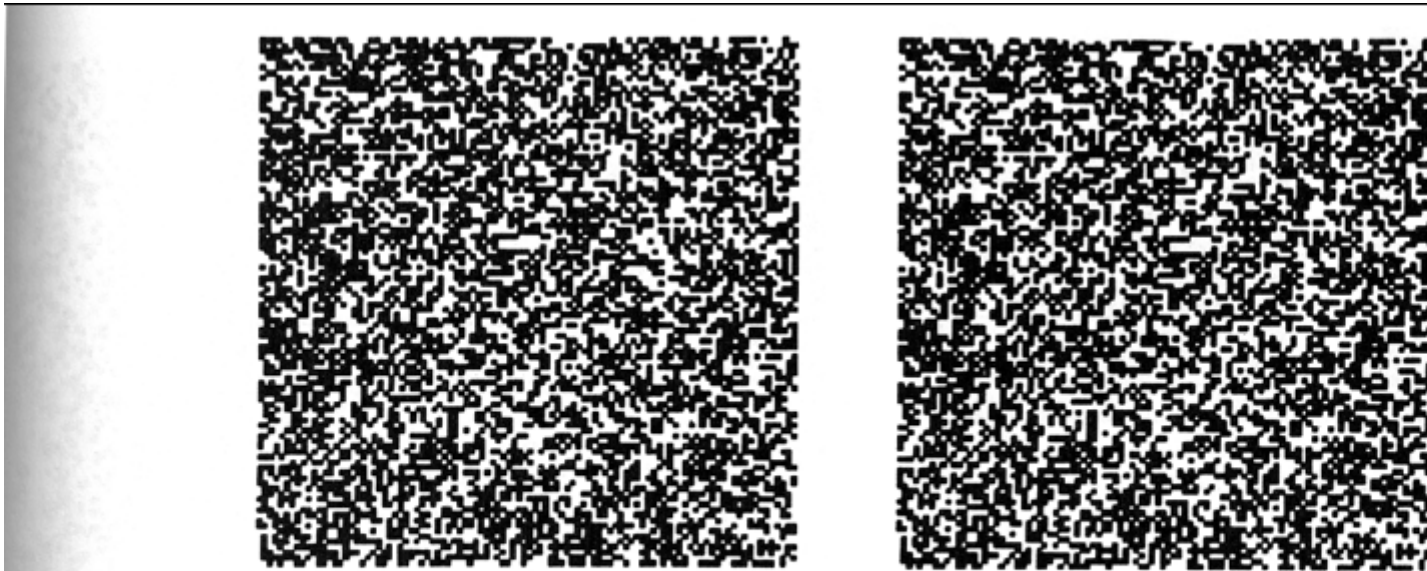
---



Adapted from Trevor Darrell, MIT



# Random dot stereograms



**Figure 5.3.8** A random dot stereogram. These two images are derived from a single array of randomly placed squares by laterally displacing a region of them as described in the text. When they are viewed with crossed disparity (by crossing the eyes) so

that the right eye's view of the left image is combined with the left eye's view of the right image, a square will be perceived to float above the page. (See pages 210–211 for instructions on fusing stereograms.)

square

# Random dot stereograms



**Figure 5.3.9** A random dot stereogram of a spiral surface. If these two images are fused with crossed convergence (see text on pages 210–211 for instructions), they can be perceived as a spiral

ramp coming out of the page toward your face. This perception arises from the small lateral displacements of thousands of tiny dots. (From Julesz, 1971.)

Spiral ramp

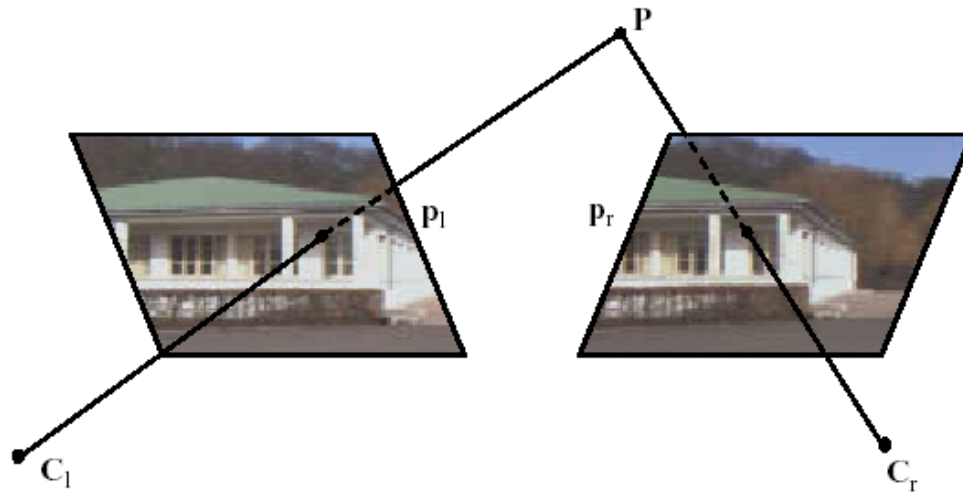
# Random dot stereograms

---

Human binocular fusion cannot be explained by peripheral processes directly associated with the physical retinas.

Instead, it must involve the central nervous system and an imaginary *cyclopean retina* that combines the left and right image stimuli as a single unit

# Stereo vision = correspondences + reconstruction



Stereovision involves two problems:

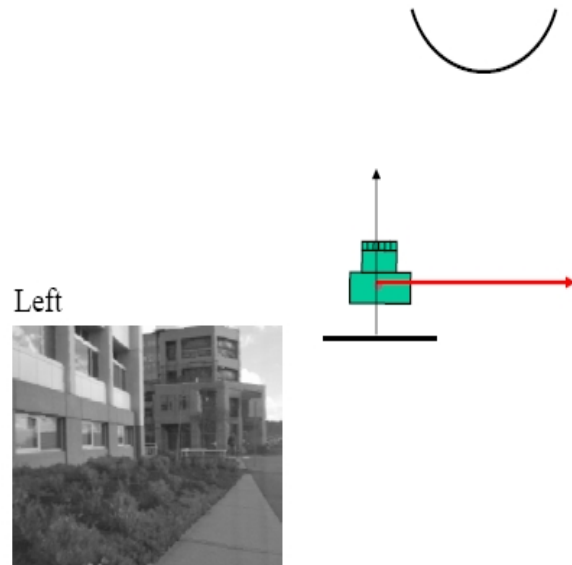
**Correspondence** : Given a point  $p_l$  in one image, find the corresponding point in the other image

**Reconstruction**: Given a correspondence  $(p_l, p_r)$  compute the 3D coordinates of the corresponding point in space,  $P$



# Binocular Stereo

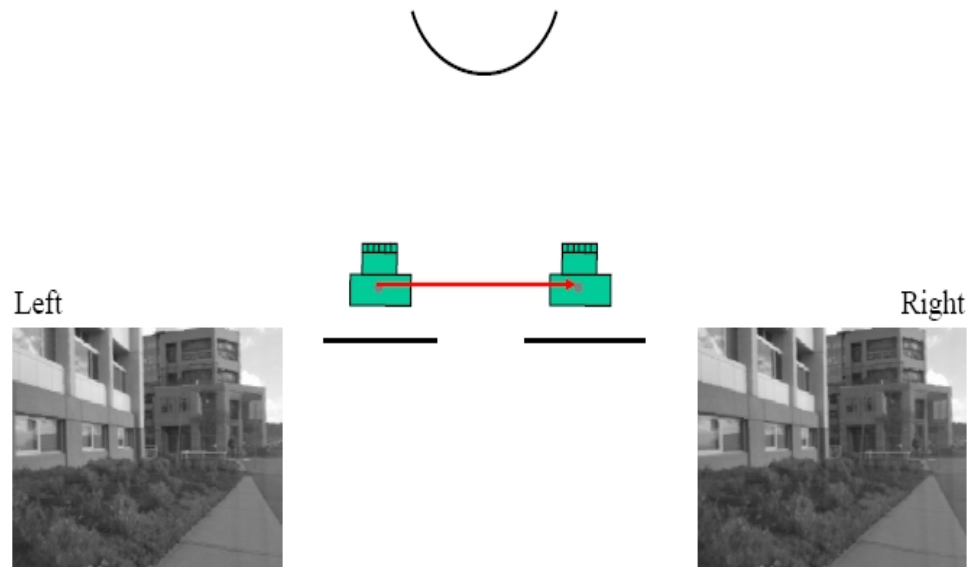
---



Adapted from Michael Black

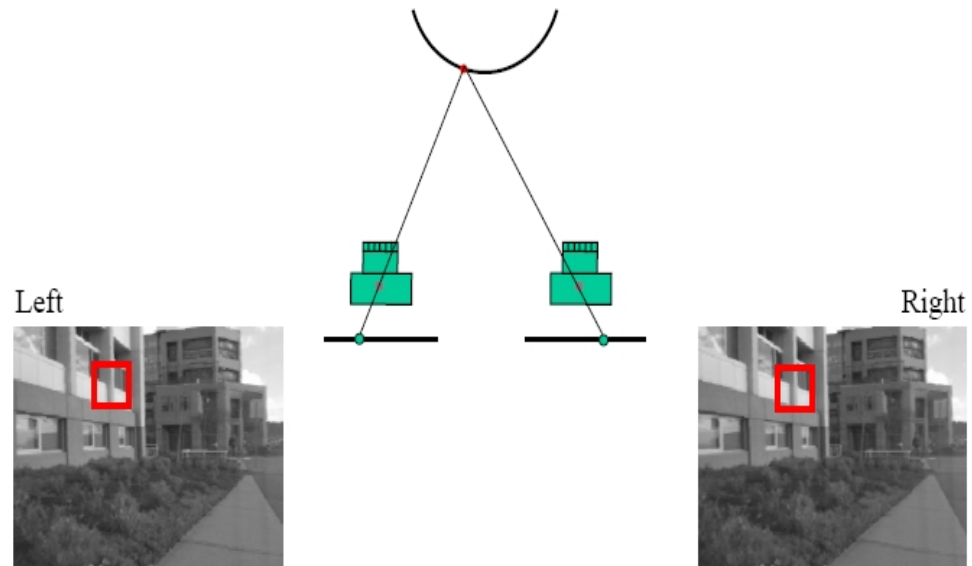
# Binocular Stereo

---



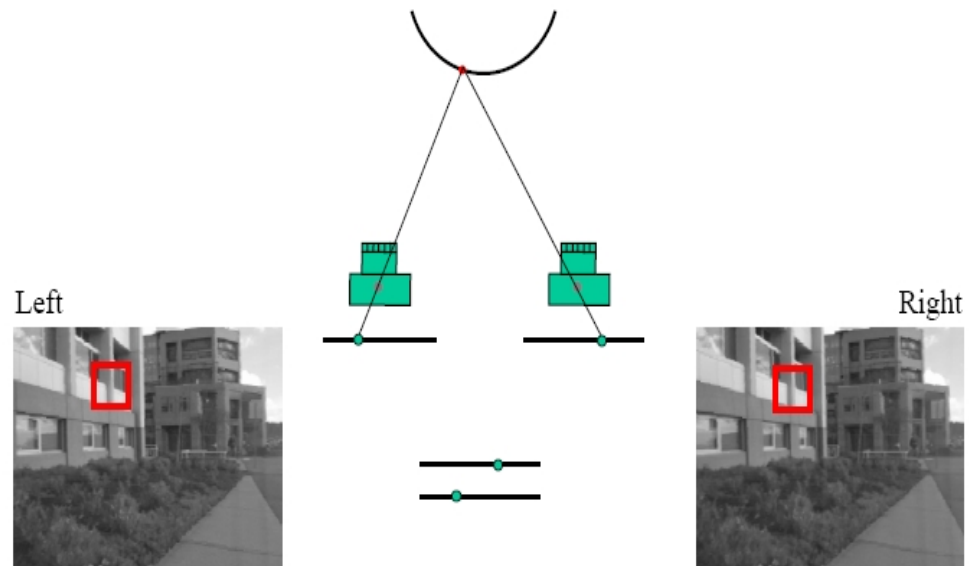
Adapted from Michael Black

# Binocular Stereo



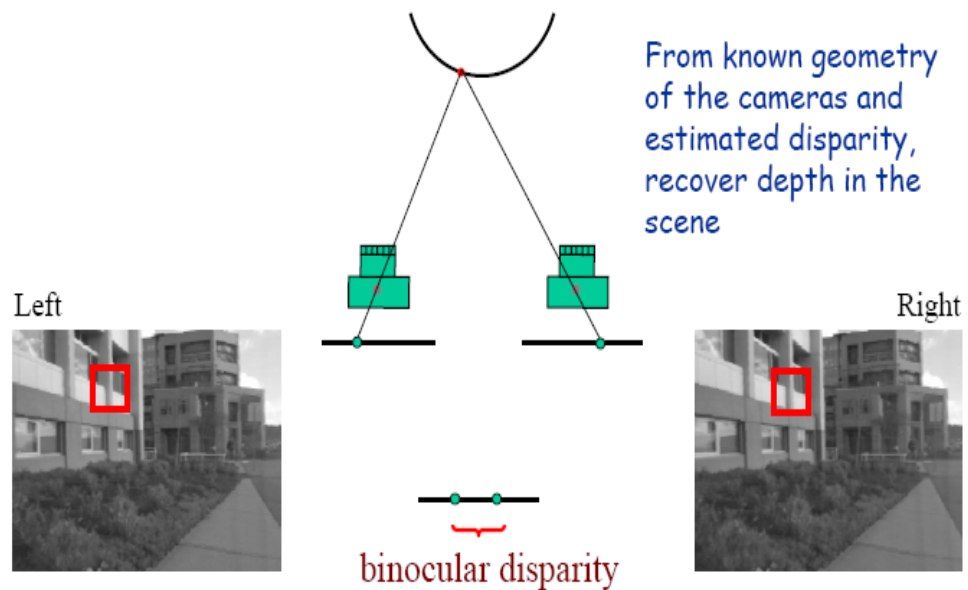
Adapted from Michael Black

# Binocular Stereo



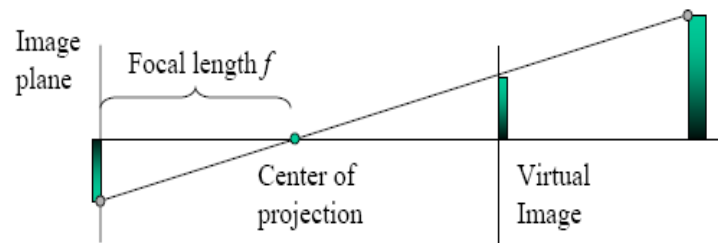
Adapted from Michael Black

# Binocular Stereo



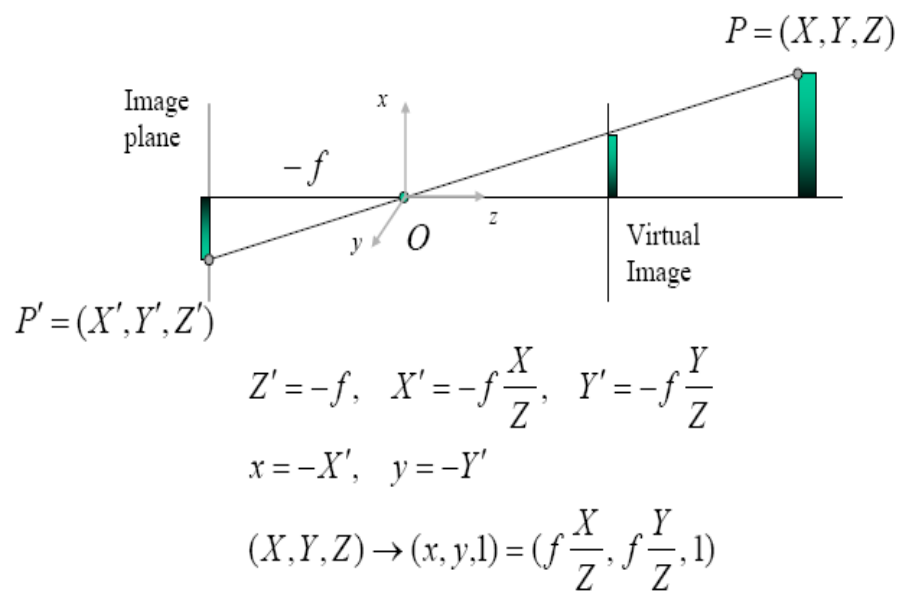
Adapted from Michael Black

# Depth Estimation



Adapted from Michael Black

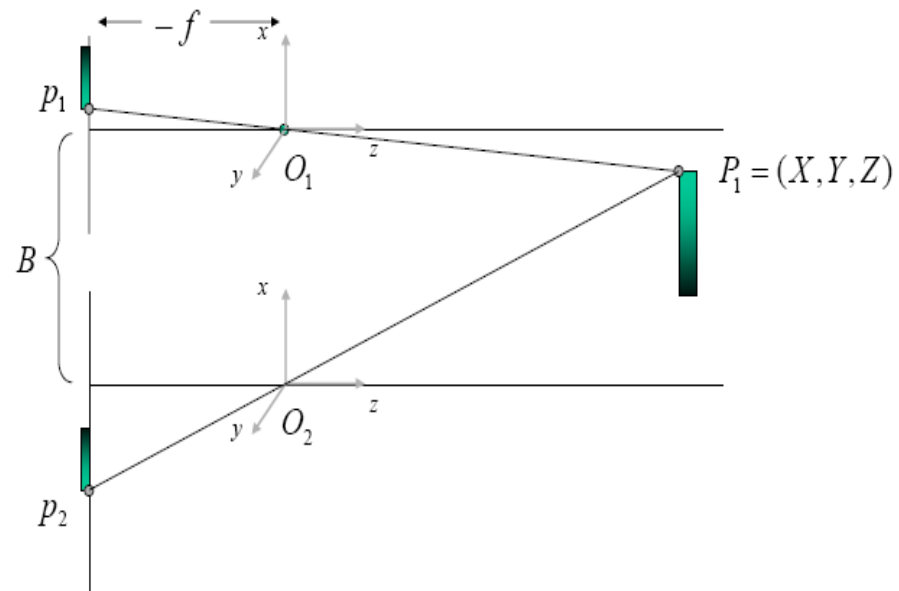
# Depth Estimation



Adapted from Michael Black

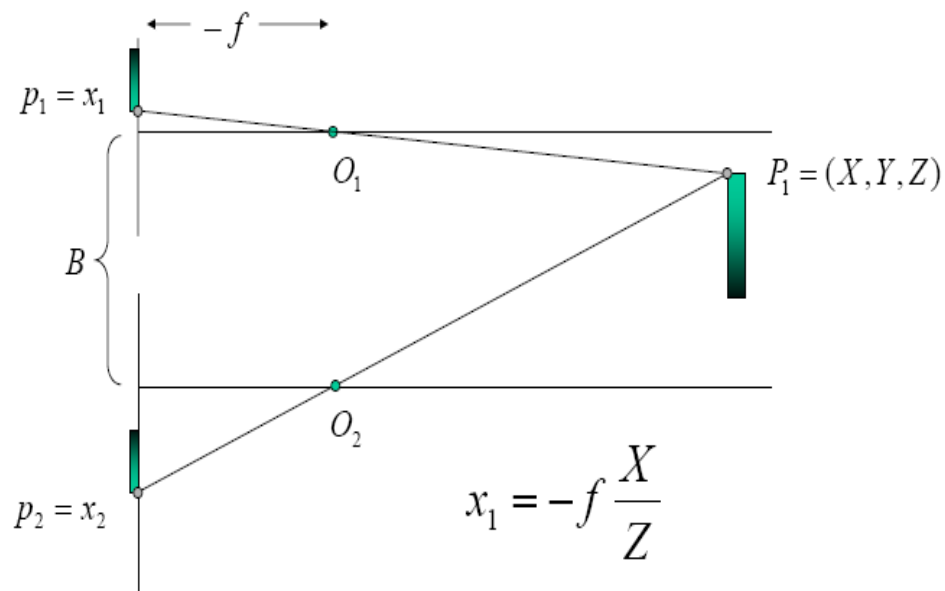


# Depth Estimation



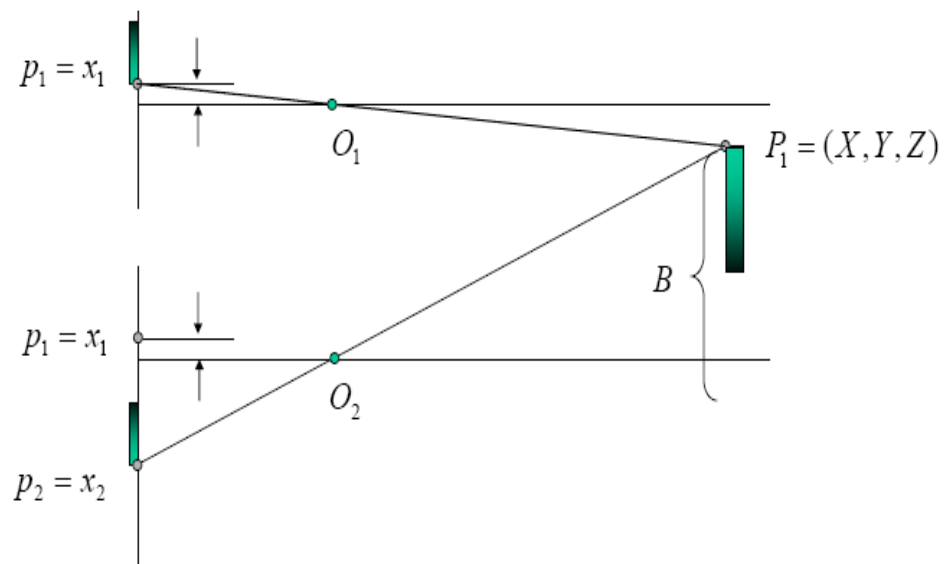
Adapted from Michael Black

# Depth Estimation



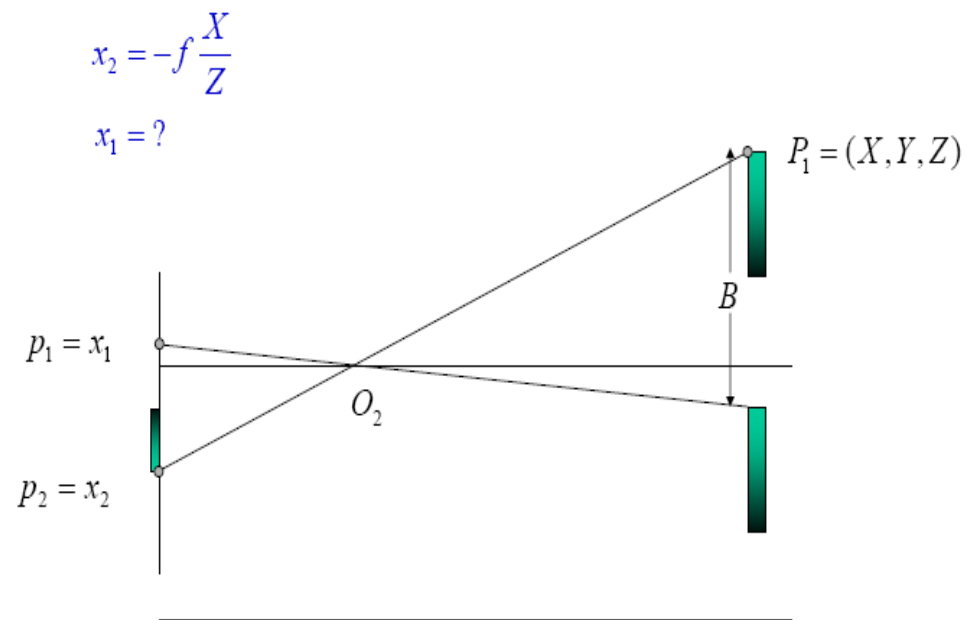
Adapted from Michael Black

# Depth Estimation

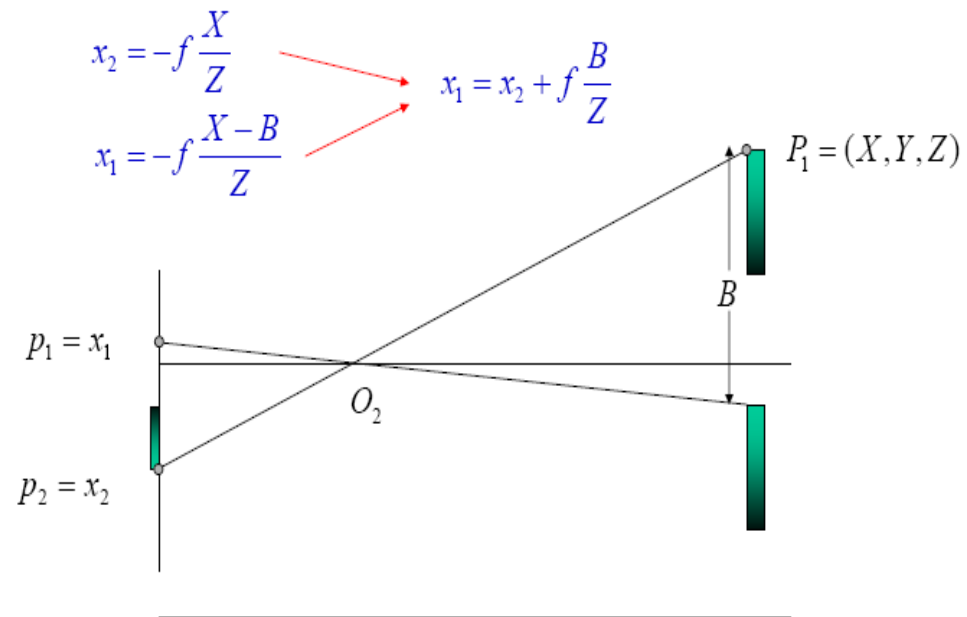


Adapted from Michael Black

# Depth Estimation

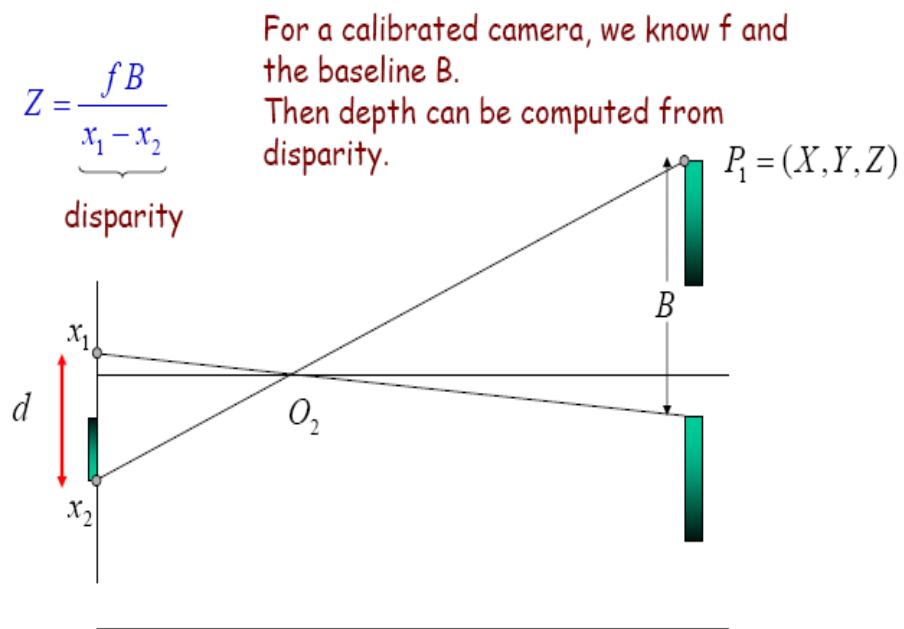


# Depth Estimation



Adapted from Michael Black

# Depth Estimation

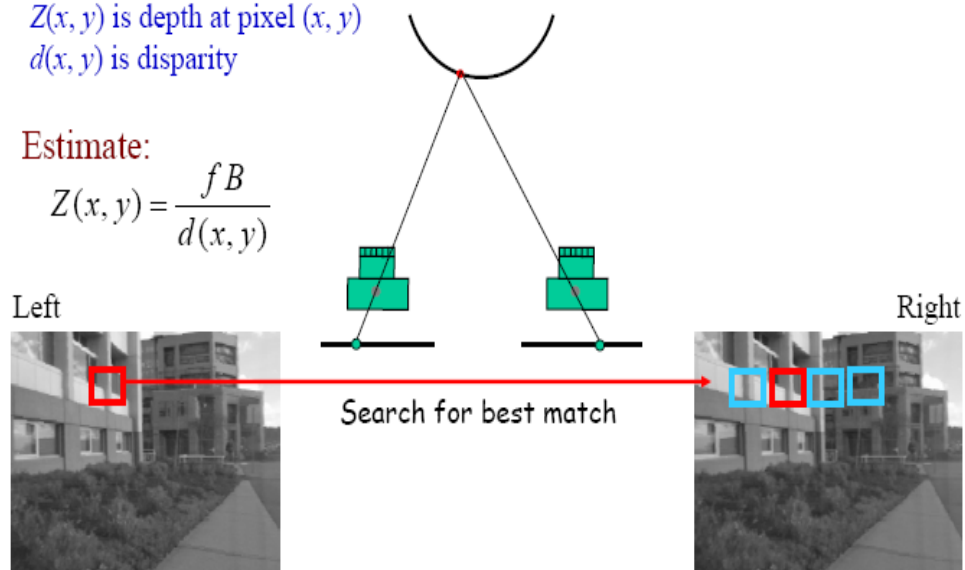


# Correspondence

$Z(x, y)$  is depth at pixel  $(x, y)$   
 $d(x, y)$  is disparity

Estimate:

$$Z(x, y) = \frac{fB}{d(x, y)}$$



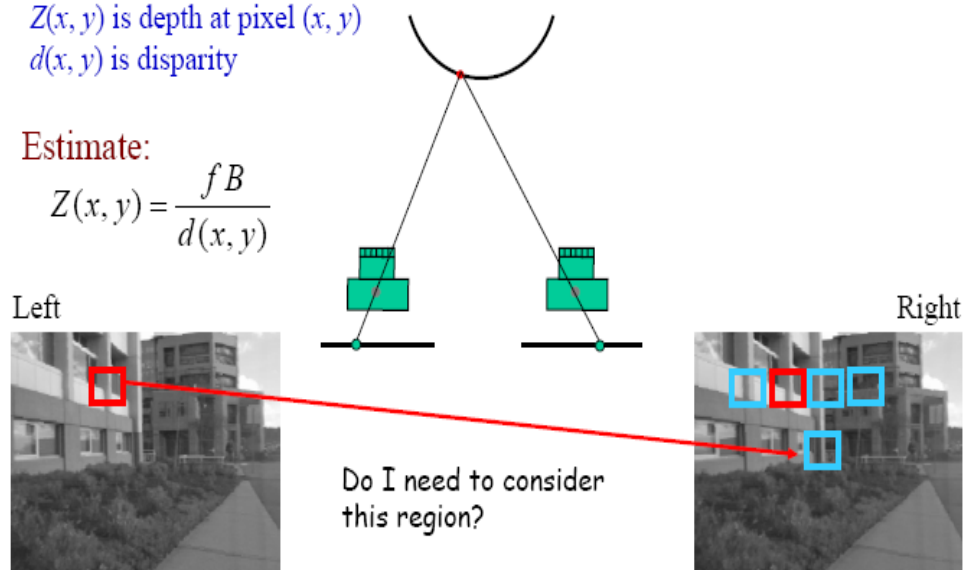


# Correspondence

$Z(x, y)$  is depth at pixel  $(x, y)$   
 $d(x, y)$  is disparity

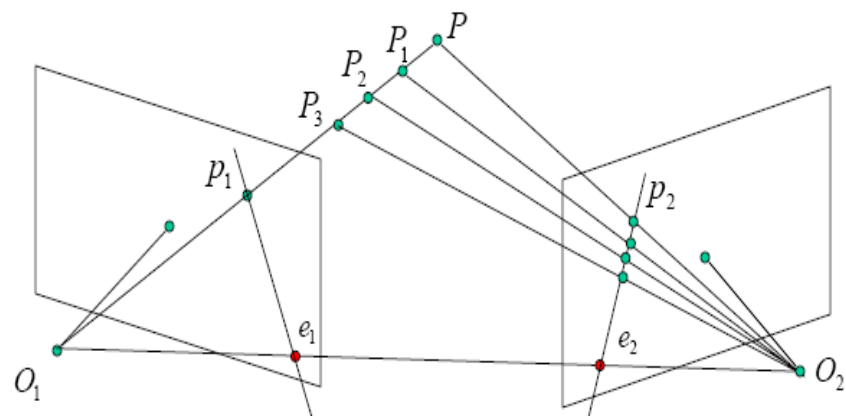
Estimate:

$$Z(x, y) = \frac{fB}{d(x, y)}$$

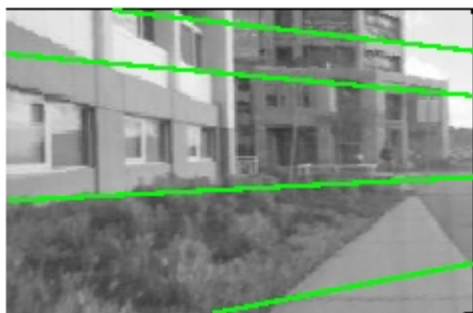


# Correspondence

Possible matches for  $p_1$  are constrained to lie along the epipolar line in the other image

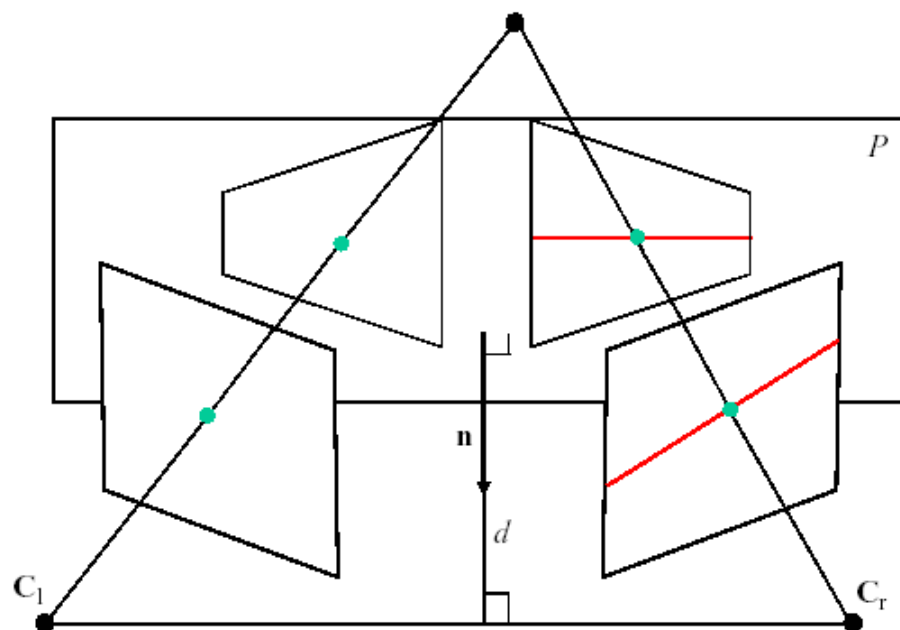


# Rectification



Searching along epipolar lines at arbitrary orientation is intuitively expensive. It would be nice to be able to always search along the rows of the right image. Fortunately, given the epipolar geometry of the stereo pair, there always exists a transformation that maps the images into a pair of images with the epipolar lines parallel to the rows of the image. This transformation is called *rectification*. Images are almost always rectified before searching for correspondences in order to simplify the search.

# Rectification



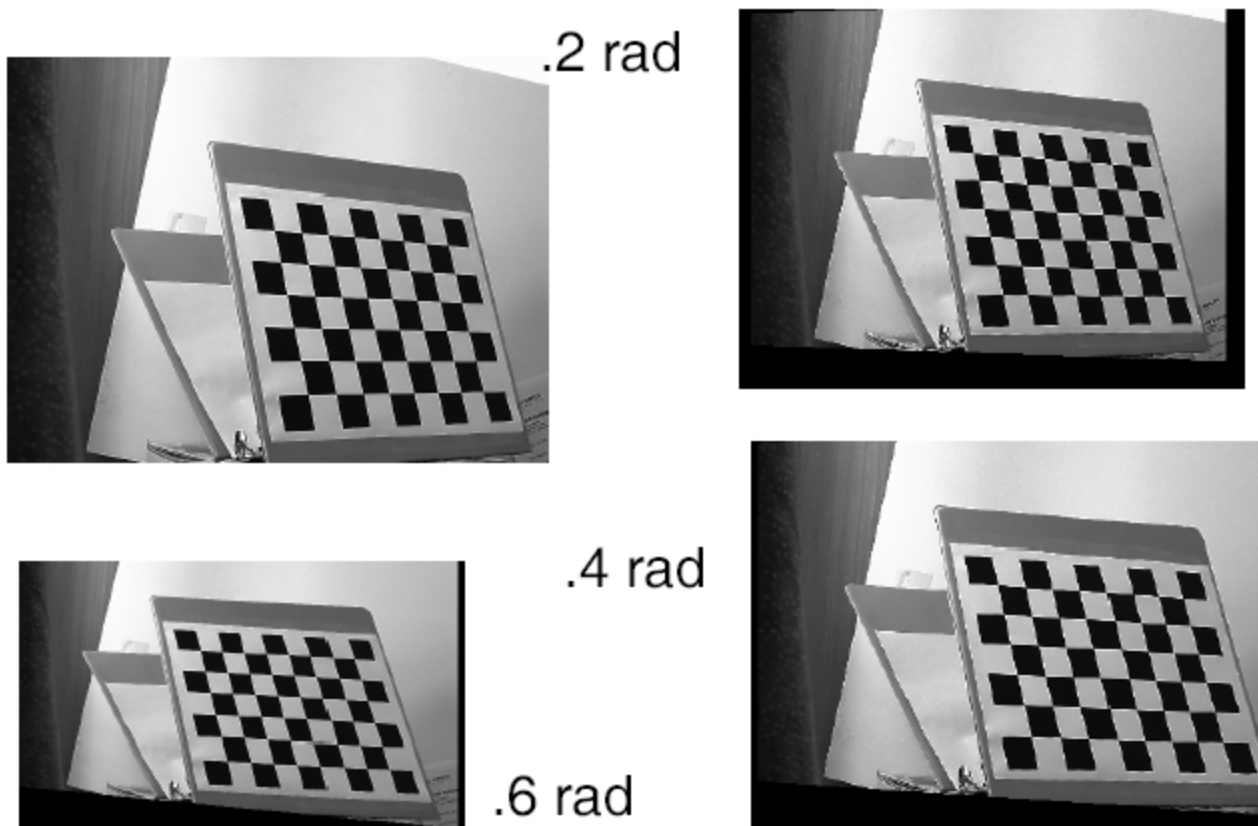
We know that, given a plane  $\mathbf{P}$  in space, there exists two homographies  $\mathbf{H}_l$  and  $\mathbf{H}_r$  that map each image plane onto  $\mathbf{P}$ . That is, if  $\mathbf{p}_l$  is a point in the left image, then the corresponding point in  $\mathbf{P}$  is  $\mathbf{H}\mathbf{p}$  (in homogeneous coordinates). If we map both images to a common plane  $\mathbf{P}$  such that  $\mathbf{P}$  is parallel to the line  $C_L C_R$ , then the pair of virtual (rectified) images is such that the epipolar lines are parallel. With proper choice of the coordinate system, the epipolar lines are parallel to the rows of the image.

The algorithm for rectification is then:

- Select a plane  $\mathbf{P}$  parallel to  $C_L C_R$
- Define the left and right image coordinate systems on  $\mathbf{P}$
- Construct the rectification matrices  $\mathbf{H}_l$  and  $\mathbf{H}_r$  from  $\mathbf{P}$  and the virtual image's coordinate systems.

# Rectification Results

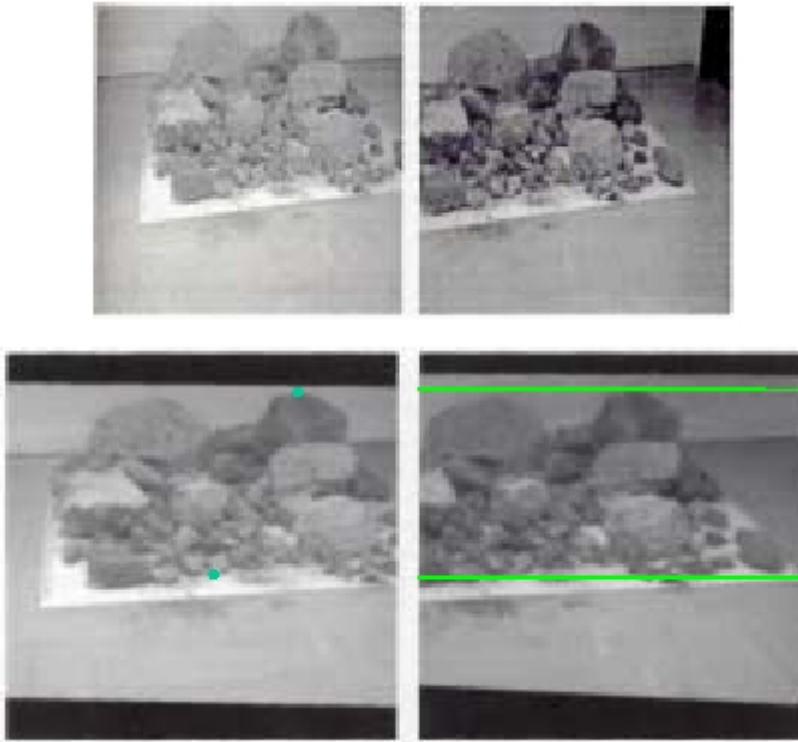
## Rectification Results



Adapted from G. Hager, JHU

# Rectification

---

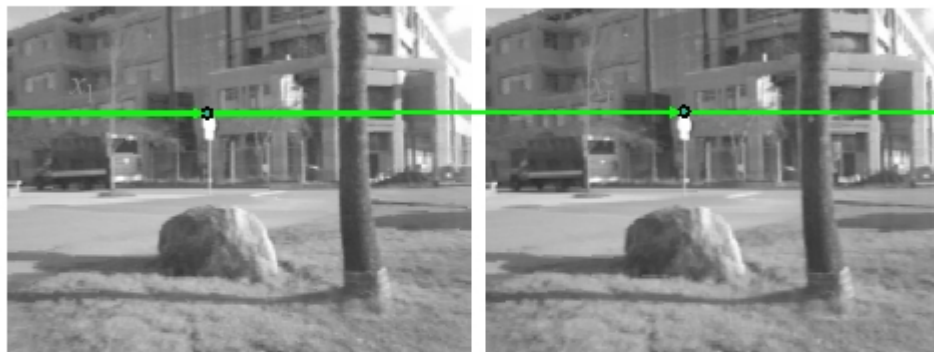


Adapted from Martial Hebert, CMU

# Disparity

Assuming that images are rectified to simplify things, given two corresponding points  $\mathbf{p}_l$  and  $\mathbf{p}_r$ , the difference of their coordinates along the epipolar line  $x_l - x_r$  is called the disparity  $d$ . The disparity is the quantity that is directly measured from the correspondence.

It turns out that the position of the corresponding 3-D point  $\mathbf{P}$  can be computed from  $\mathbf{p}_l$  and  $d$ , assuming that the camera parameters are known.

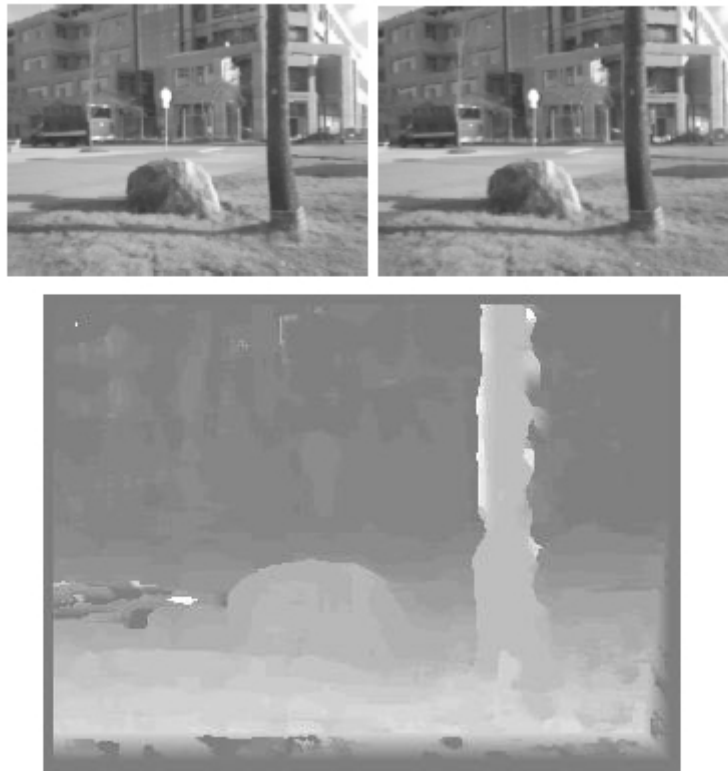


$$d = x_l - x_r$$



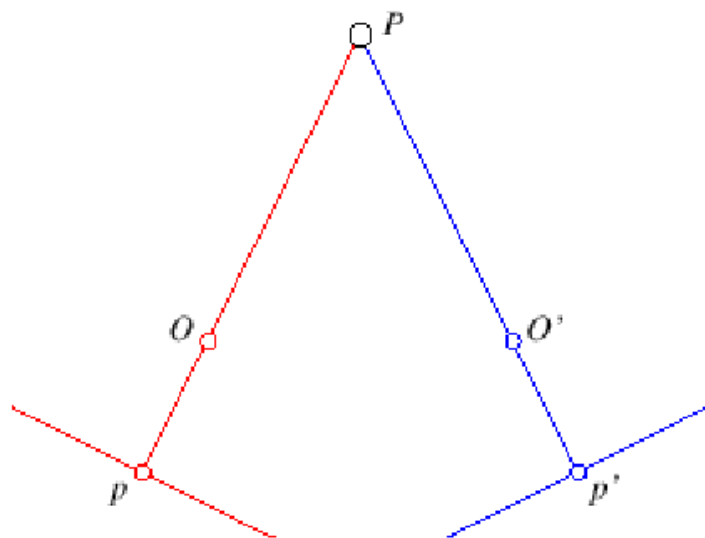
# Disparity

---

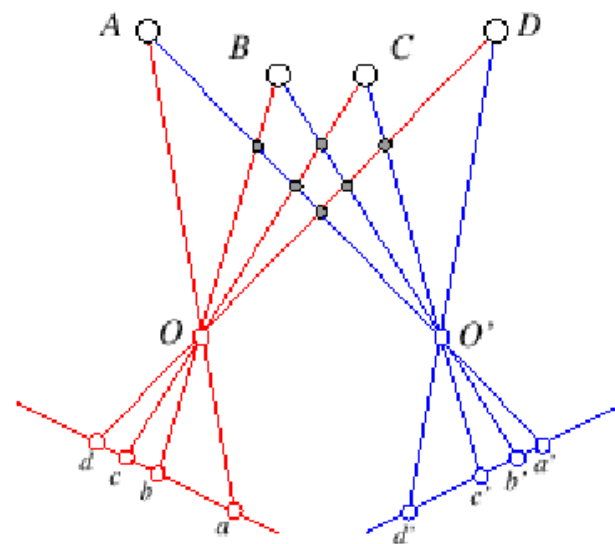


Larger disparity  $\rightarrow$  closer to camera

# Stereopsis

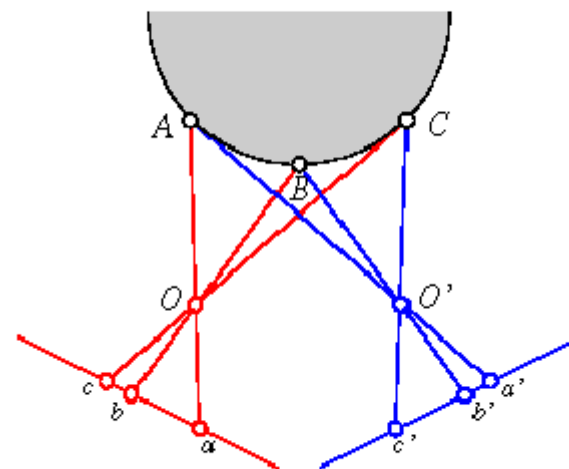


If a single image point is observed at any given time  
Stereo vision is easy



However, each picture consists of hundreds/thousands of  
pixels, therefore it is very hard to find the correct  
correspondences

# Ordering constraint

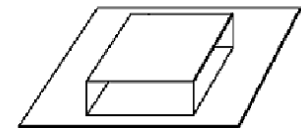
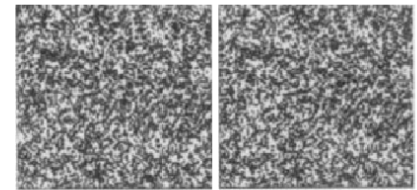


“It is reasonable to assume that the order of matching image features along a pair of epipolar lines is the inverse of the order of the corresponding surface attributes along the curve where the epipolar plane intersects the observed object’s boundary.”

This is the so-called *ordering constraint* introduced by [Baker and Binford, 1981; Ohta and Kanade, 1985].

# Correspondence is ambiguous (Marr & Poggio)

---



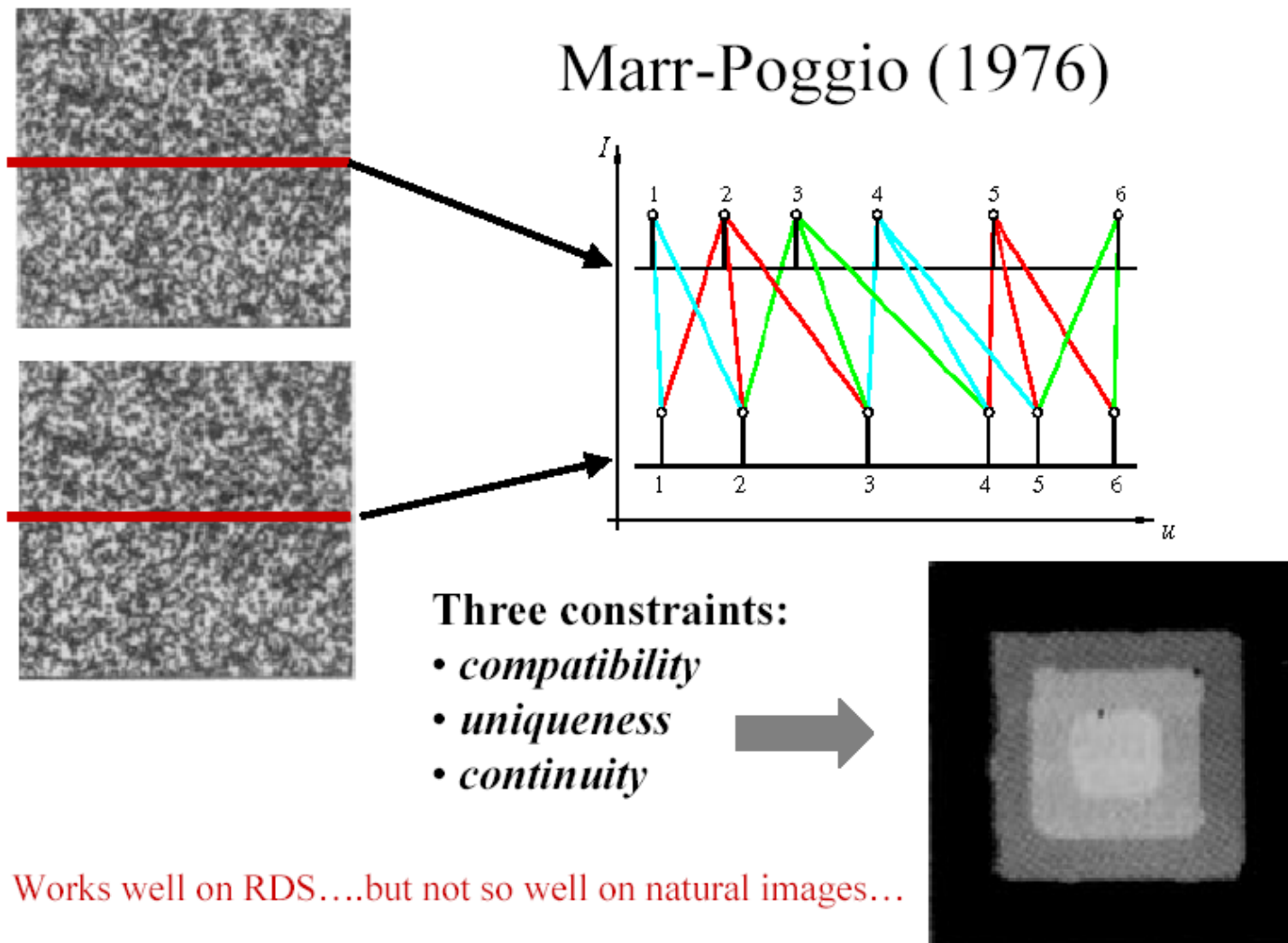
Three constraints :

Compatibility : black dots can only match black dots , or more generally, two image features can only match if they have possibly arisen from the same physical marking

Uniqueness : a black dot in one image matches at most one black dot in another image

Continuity : the disparity of matches varies smoothly almost everywhere in the image

# Correspondence is ambiguous

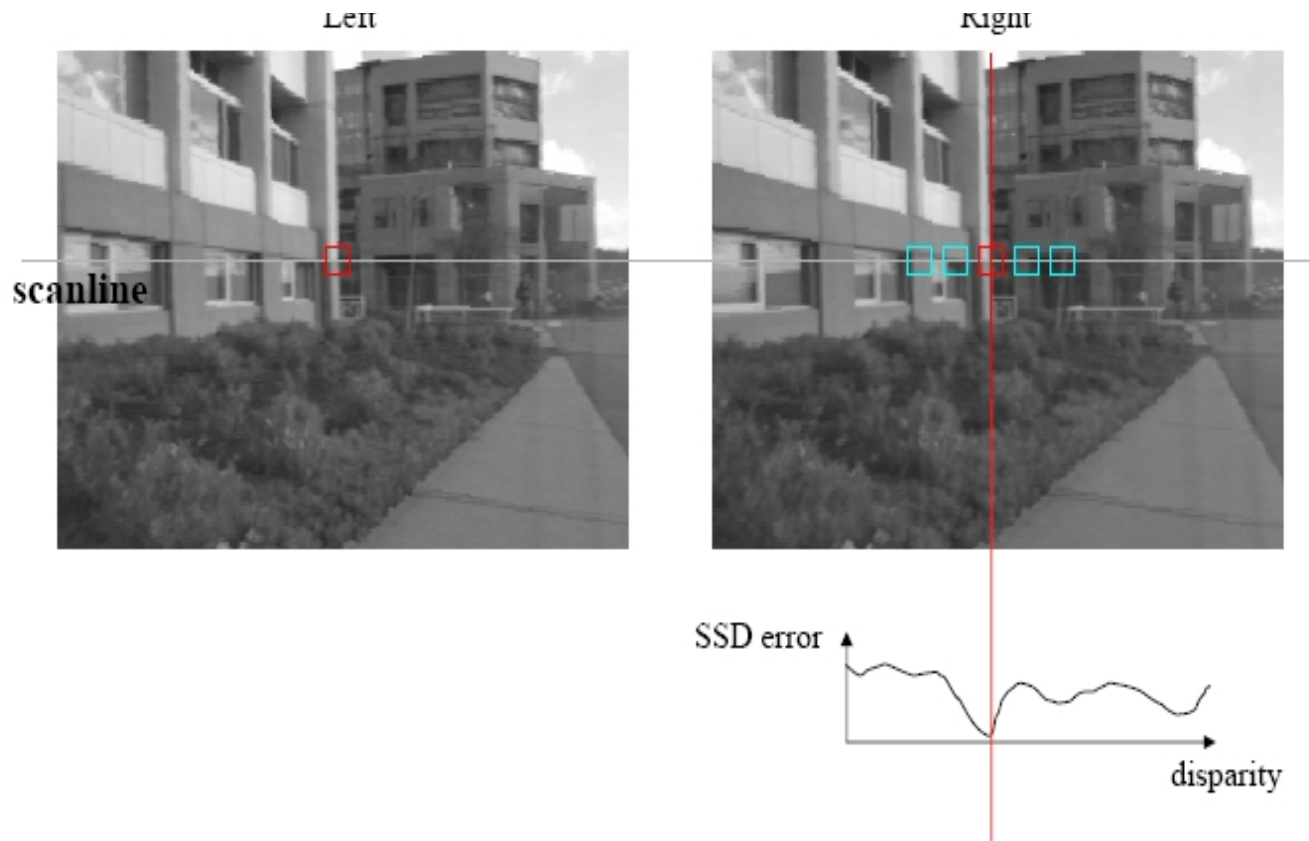


Adapted from Trevor Darrell, MIT

# Correspondence using window matching

Points are highly individually ambiguous...

More unique matches are possible with small regions of image.



Adapted from Michael Black

# Finding Correspondences

---

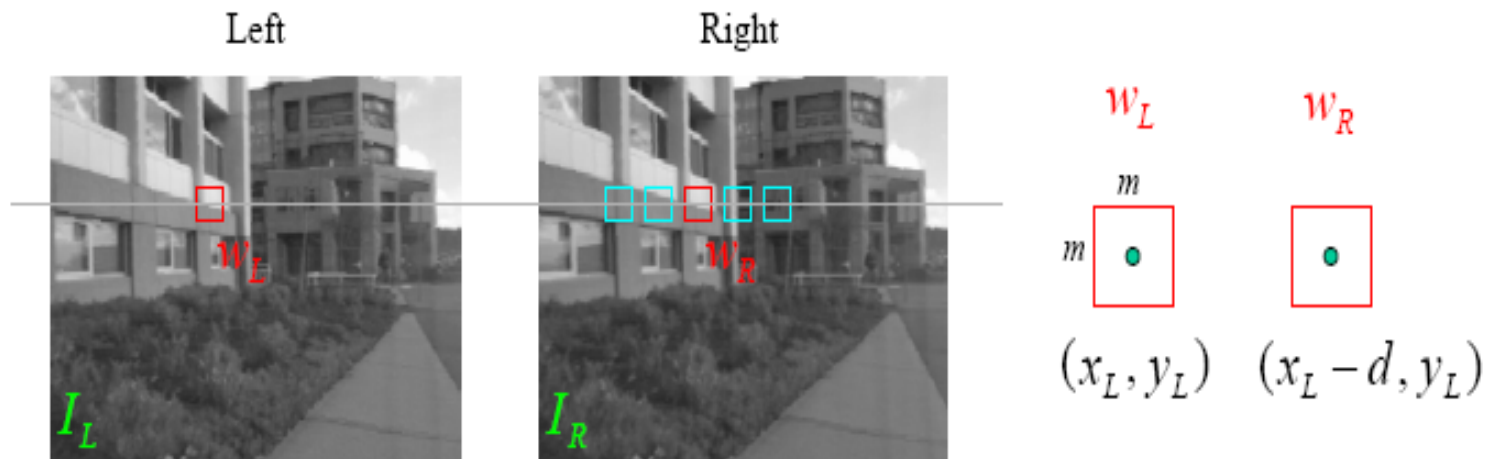


$W(\mathbf{p}_l)$



$W(\mathbf{p}_r)$

# Sum of squared distances



$w_L$  and  $w_R$  are corresponding  $m$  by  $m$  windows of pixels.

The SSD cost measures the intensity difference as a function of disparity :

$$SSD_r(x, y, d) = \sum_{(x', y') \in W_m(x, y)} (I_L(x', y') - I_R(x' - d, y'))^2$$



# Image Normalization

---

- Even when the cameras are identical models, there can be differences in gain and sensitivity.
- The cameras do not see exactly the same surfaces, so their overall light levels can differ.
- For these reasons and more, it is a good idea to normalize the pixels in each window:

$$\bar{I} = \frac{1}{|W_m(x,y)|} \sum_{(u,v) \in W_m(x,y)} I(u,v)$$

Average pixel

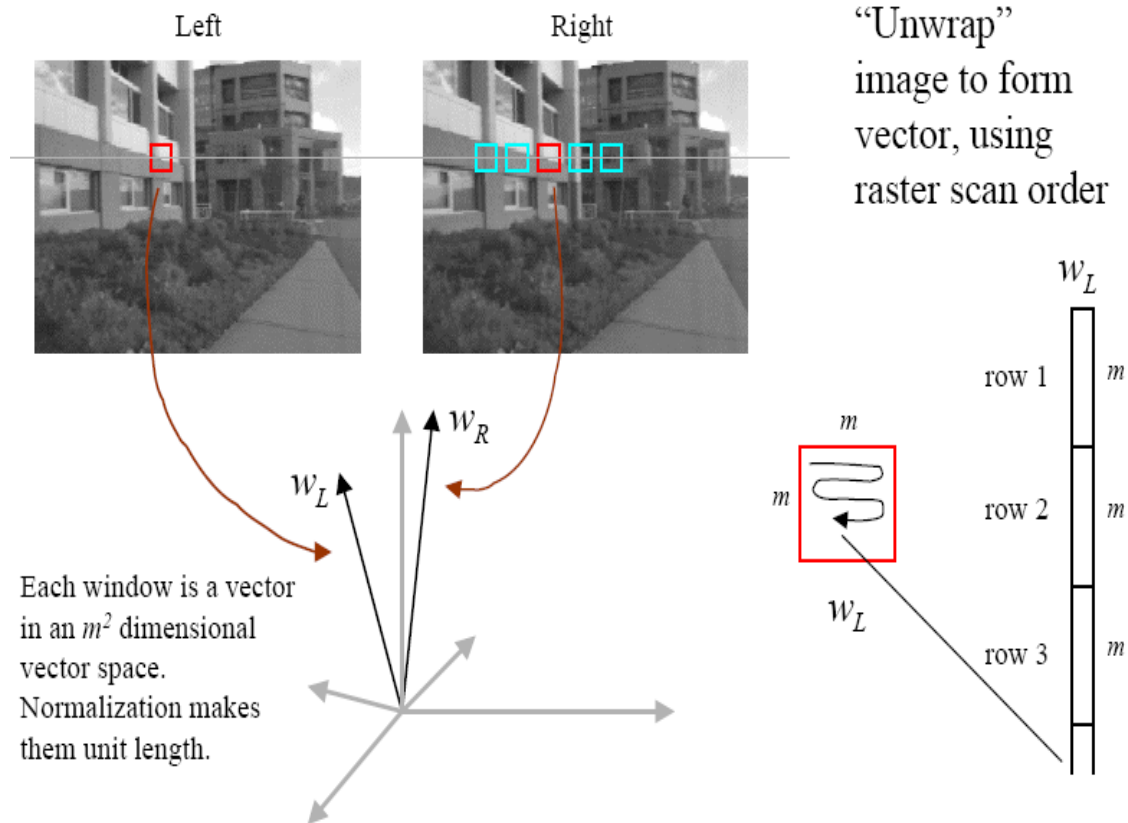
$$\|I\|_{W_m(x,y)} = \sqrt{\sum_{(u,v) \in W_m(x,y)} [I(u,v)]^2}$$

Window magnitude

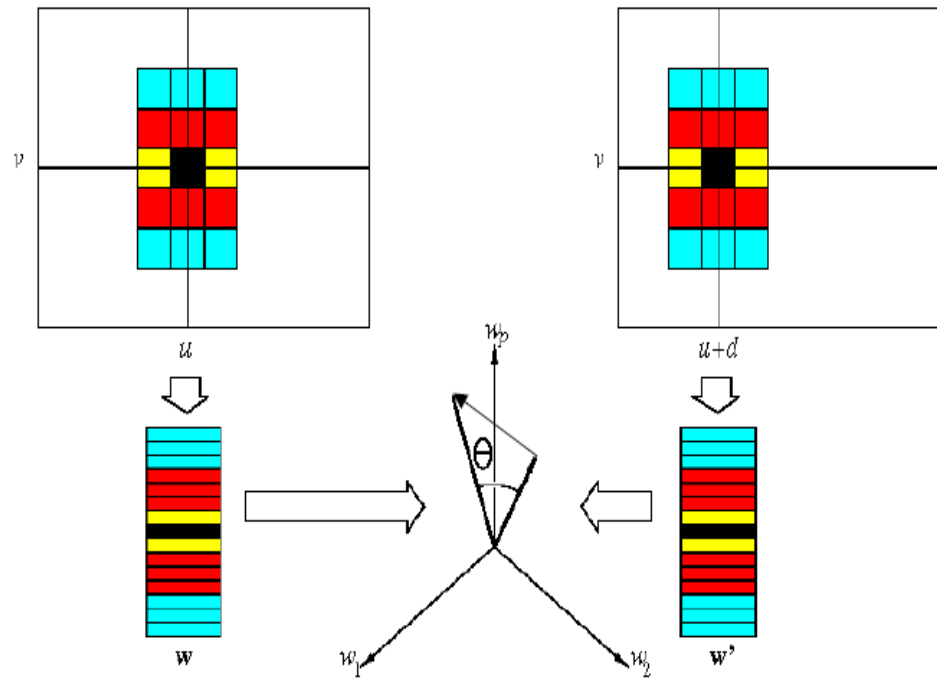
$$\hat{I}(x,y) = \frac{I(x,y) - \bar{I}}{\|I - \bar{I}\|_{W_m(x,y)}}$$

Normalized pixel

# Images as vectors

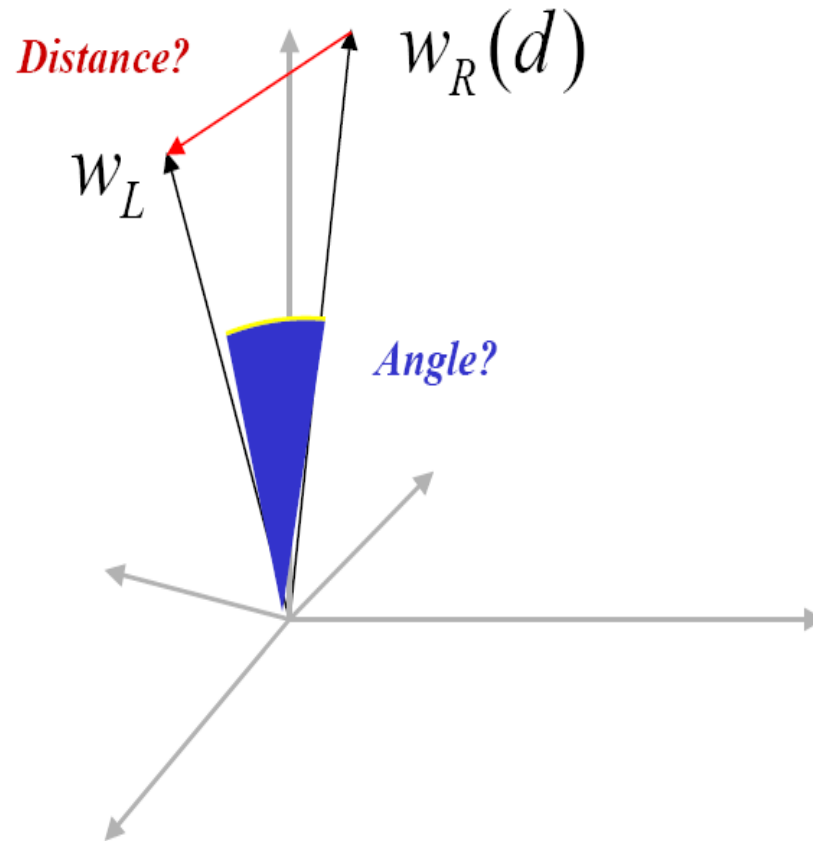


# Images as vectors

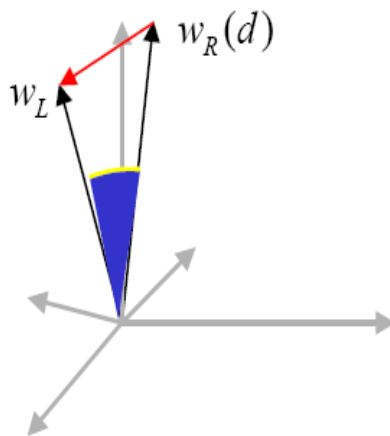


Adapted from Darrell

# Possible metric



# Possible metric



(Normalized) Sum of Squared Differences

$$C_{\text{SSD}}(d) = \sum_{(u,v) \in W_m(x,y)} [\hat{I}_L(u,v) - \hat{I}_R(u-d,v)]^2$$

$$= \|w_L - w_R(d)\|^2$$

Normalized Correlation

$$C_{\text{NC}}(d) = \sum_{(u,v) \in W_m(x,y)} \hat{I}_L(u,v) \hat{I}_R(u-d,v)$$

$$= w_L \cdot w_R(d) = \cos \theta$$

$$d^* = \arg \min_d \|w_L - w_R(d)\|^2 = \arg \max_d w_L \cdot w_R(d)$$

# Matching using correlation

---

Left



Disparity Map



Images courtesy of Point Grey Research

# Matching using correlation

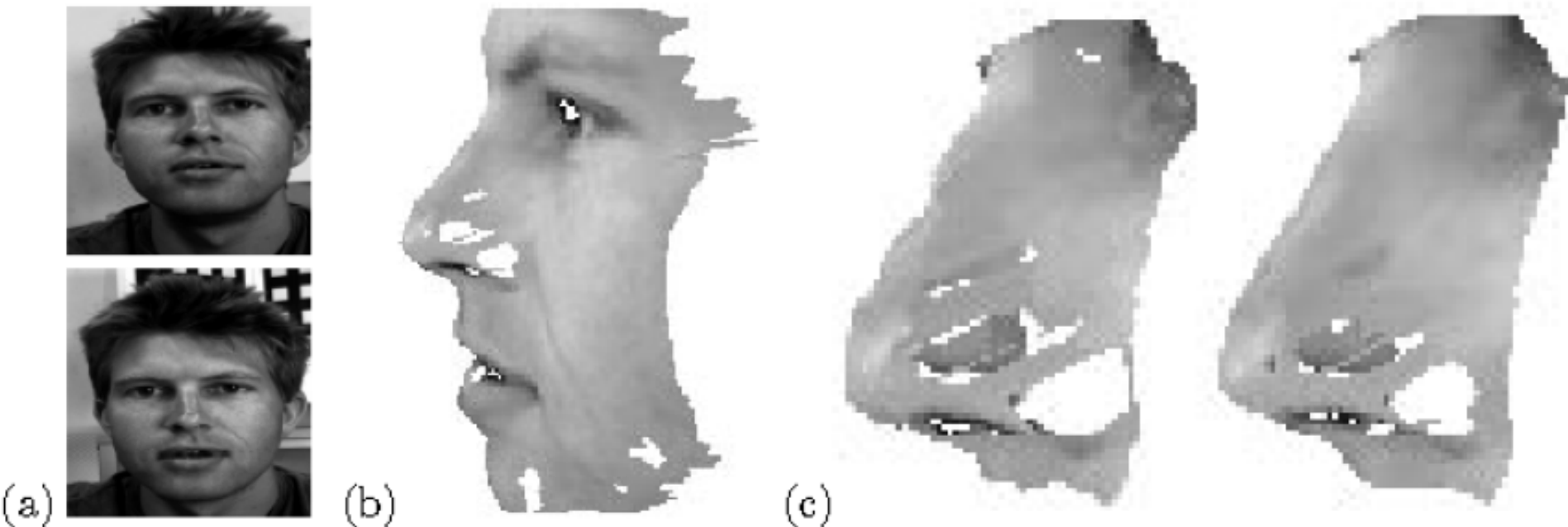


FIGURE 12.13: Correlation-based stereo matching: (a) a pair of stereo pictures; (b) a texture-mapped view of the reconstructed surface; (c) comparison of the regular (left) and refined (right) correlation methods in the nose region. Reprinted from [Devernay and Faugeras, 1994], Figures 5, 8 and 9.

# Problems with window methods

---

Patch too small?

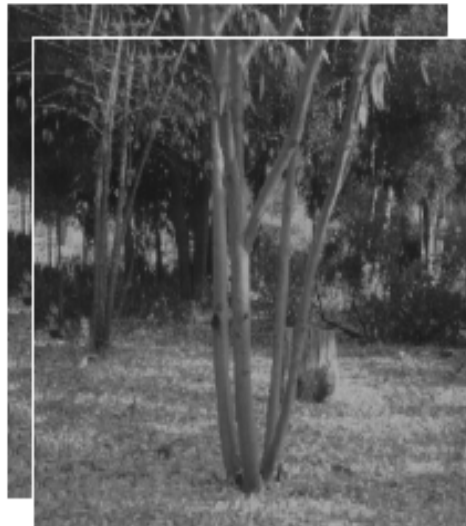
Patch too large?

*Can try variable patch size [Okutomi and Kanade],  
or arbitrary window shapes [Veksler and Zabih]*

Should match between physically meaningful  
quantities, and at multiple scales [Marr]...



# Window size



$W = 3$



$W = 20$

## Better results with *adaptive window*

- T. Kanade and M. Okutomi, [A Stereo Matching Algorithm with an Adaptive Window: Theory and Experiment](#), Proc. International Conference on Robotics and Automation, 1991.
- D. Scharstein and R. Szeliski. [Stereo matching with nonlinear diffusion](#). International Journal of Computer Vision, 28(2):155-174, July 1998

(Seitz)

# Stereo Results

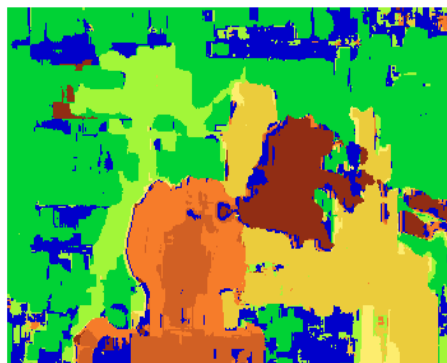
– Data from University of Tsukuba



Scene



Ground truth



Window-based matching  
(best window size)



Ground truth

(Seitz)

Adapted from Michael Black

# Multi – scale edge matching

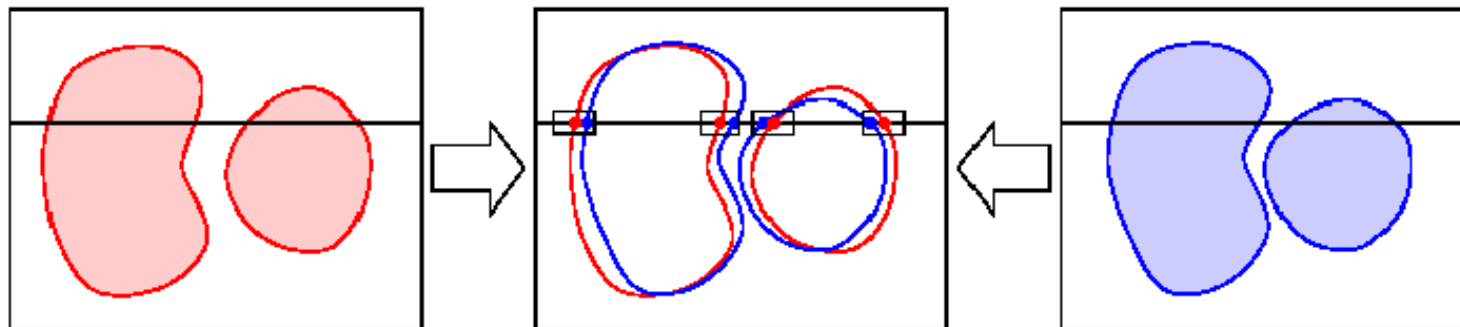
---

1. Convolve the two (rectified) images with  $\nabla^2 G_\sigma$  filters of increasing standard deviations  $\sigma_1 < \sigma_2 < \sigma_3 < \sigma_4$ .
2. Find zero crossings of the Laplacian along horizontal scanlines of the filtered images.
3. For each filter scale  $\sigma$ , match zero crossings with the same parity and roughly equal orientations in a  $[-w_\sigma, +w_\sigma]$  disparity range, with  $w_\sigma = 2\sqrt{2}\sigma$ .
4. Use the disparities found at larger scales to control eye vergence and cause unmatched regions at smaller scales to come into correspondence.

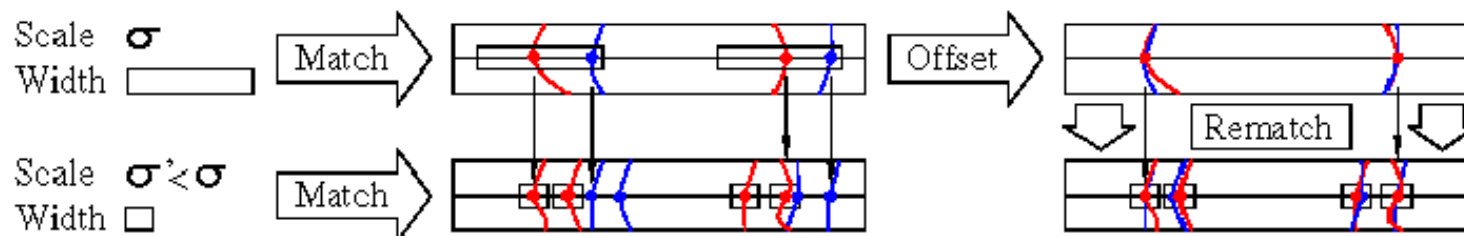
# Marr-Poggio Algorithm

Search for edges, a.k.a. “zero crossings”: (more during edge detection lectures...)

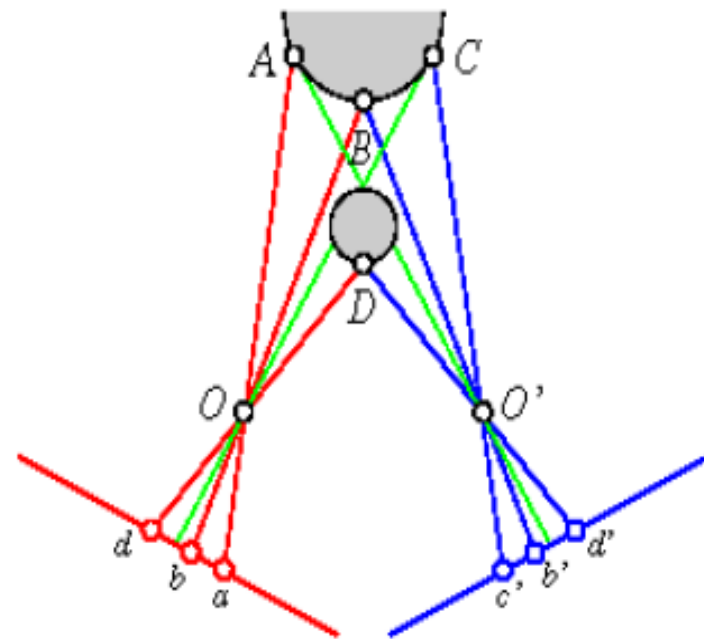
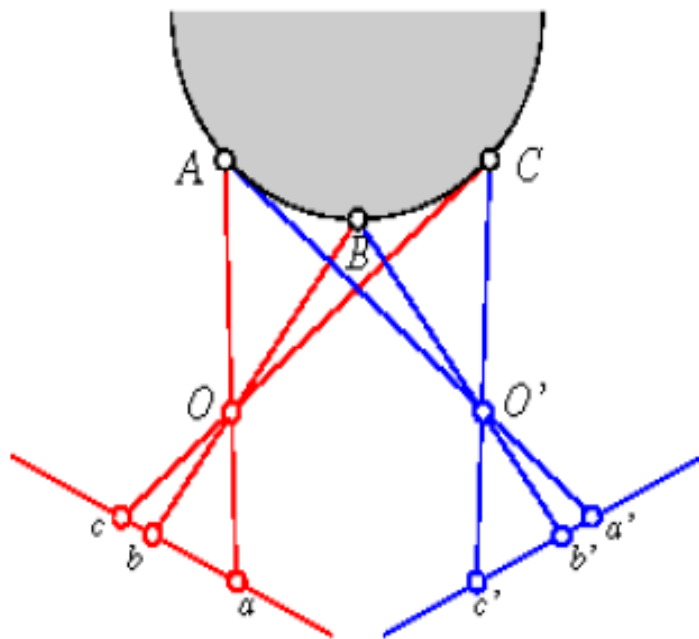
Matching zero-crossings at a single scale



Matching zero-crossings at multiple scales

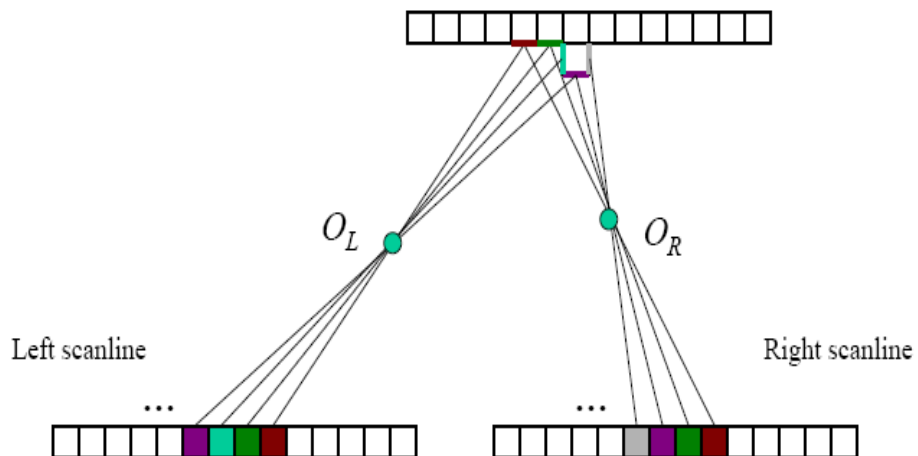


# Correspondence



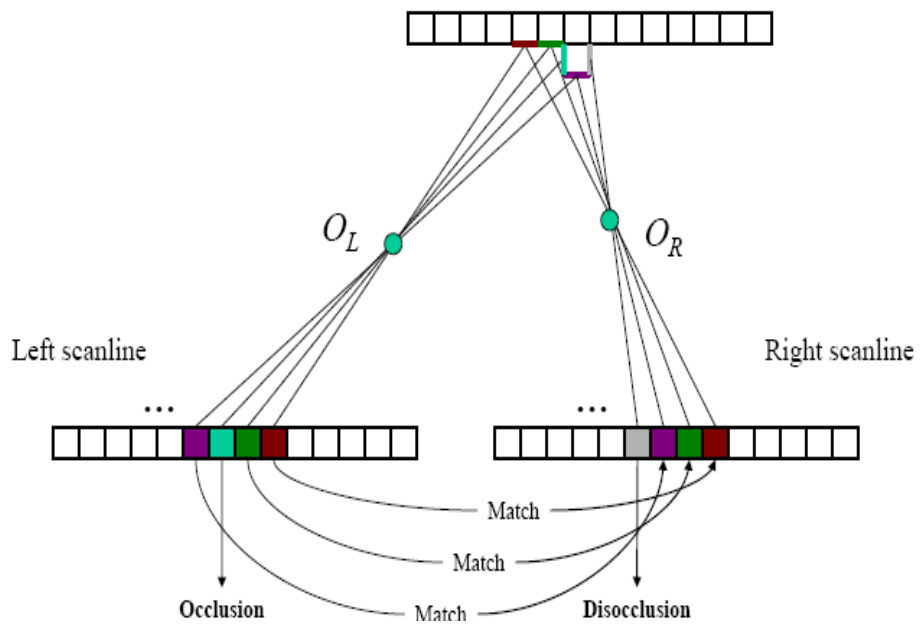
Oops!

# Correspondence



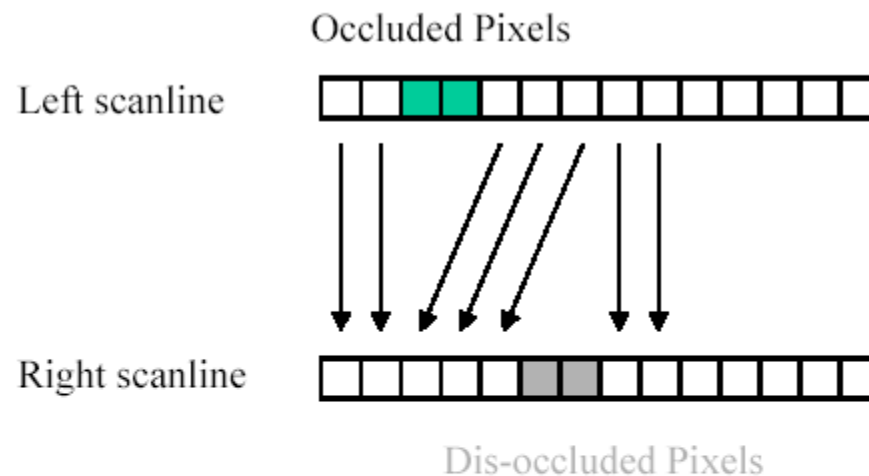
Adapted from Michael Black

# Correspondence



Adapted from Michael Black

# Search over correspondences

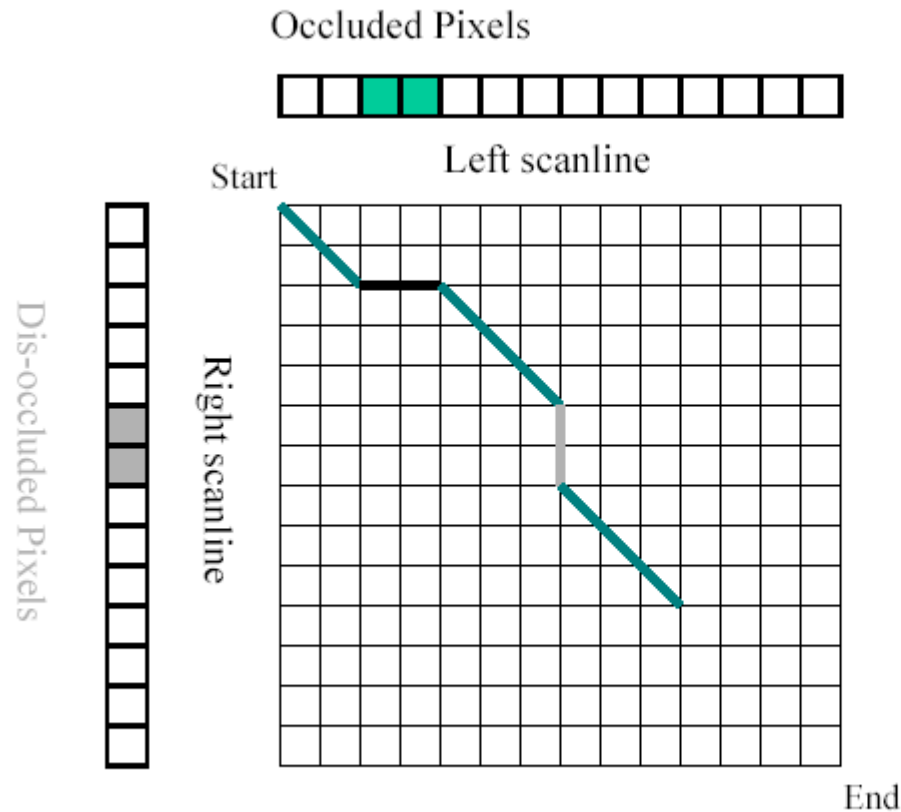


Three cases:

- Sequential – cost of match
- Occluded – cost of no match
- Disoccluded – cost of no match

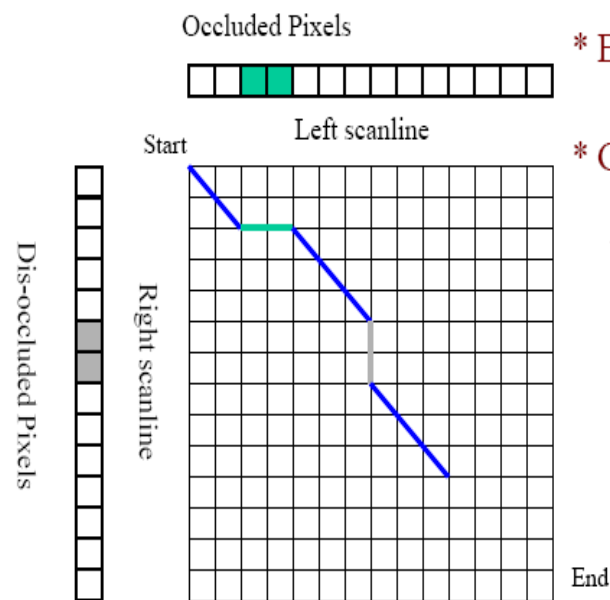


# Dynamic programming



Dynamic programming yields the optimal path through grid. This is the best set of matches that satisfy the ordering constraint

# Stereo Matching with Dynamic Programming

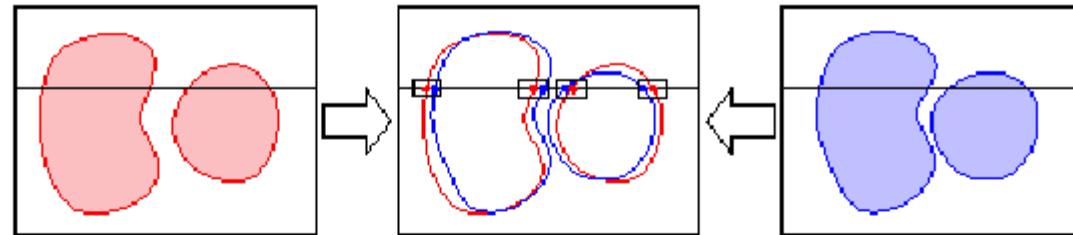


\* Enforces ordering constraint.

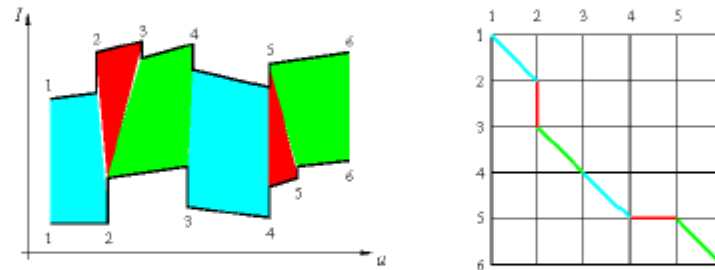
\* Given appropriate cost functions, solves for best path (matches, occlusions, disocclusions).

# DP vs. edges

Edges:



DP:

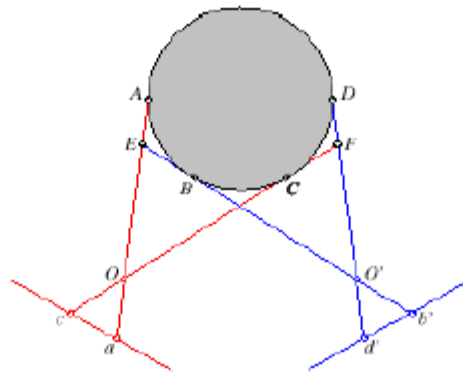


- Which method is better?
  - Edges are more “meaningful” [Marr]...but hard to find!
  - Edges tend to fail in dense texture (outdoors)
  - Correlation tends to fail in smooth featureless areas

# Computing correspondences

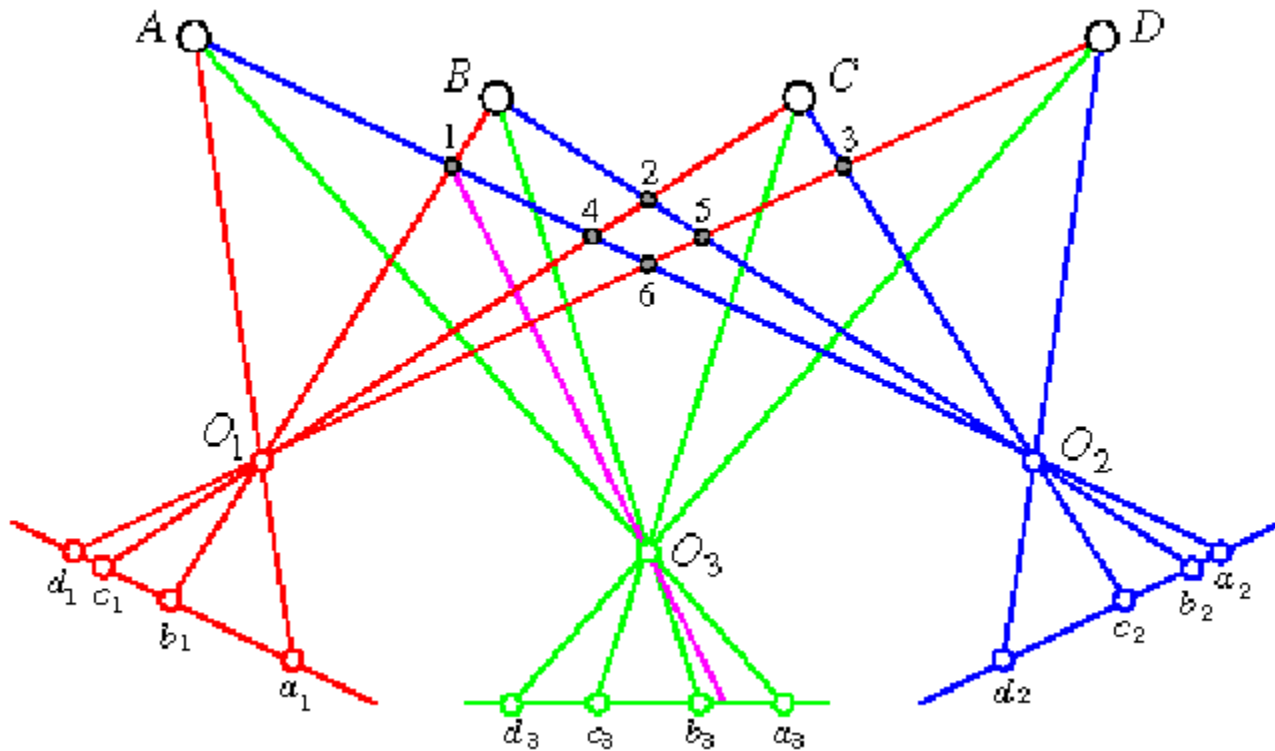
---

Both methods fail for smooth surfaces



There is currently no good solution to the correspondence problem

# Three (calibrated) views



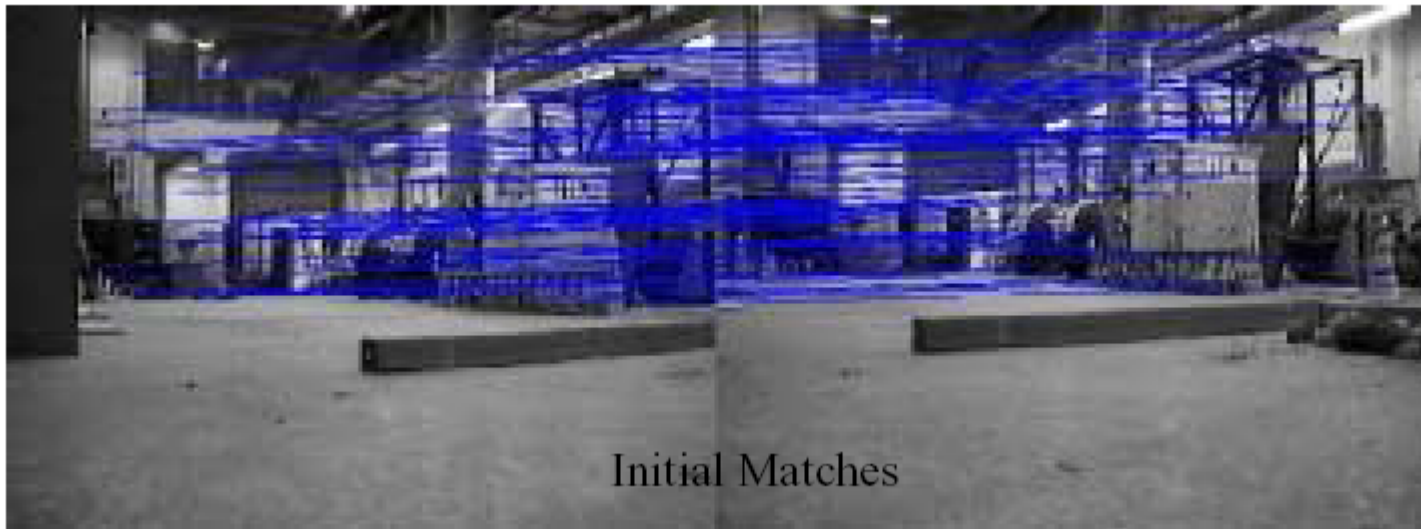
Adding a third camera eliminates the ambiguity inherent in two-view point matching

# RANSAC

---

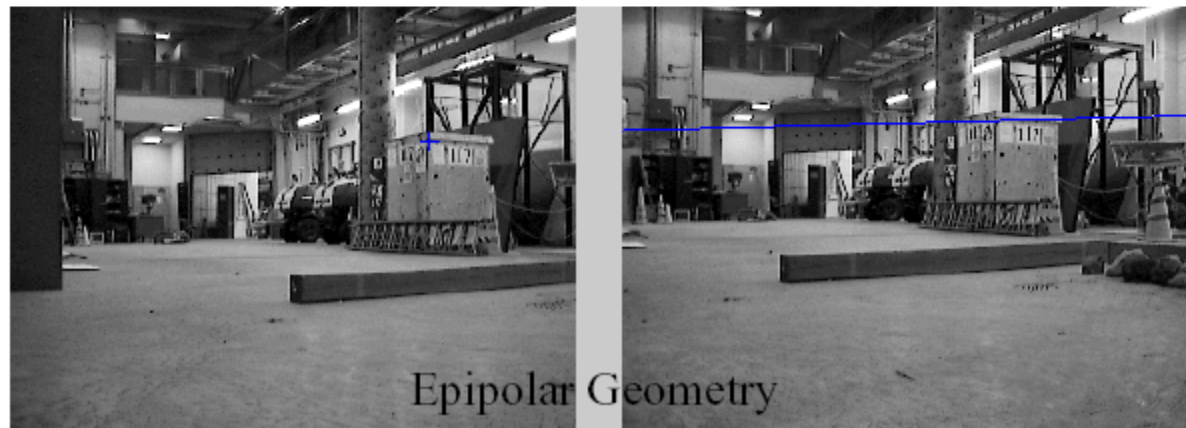
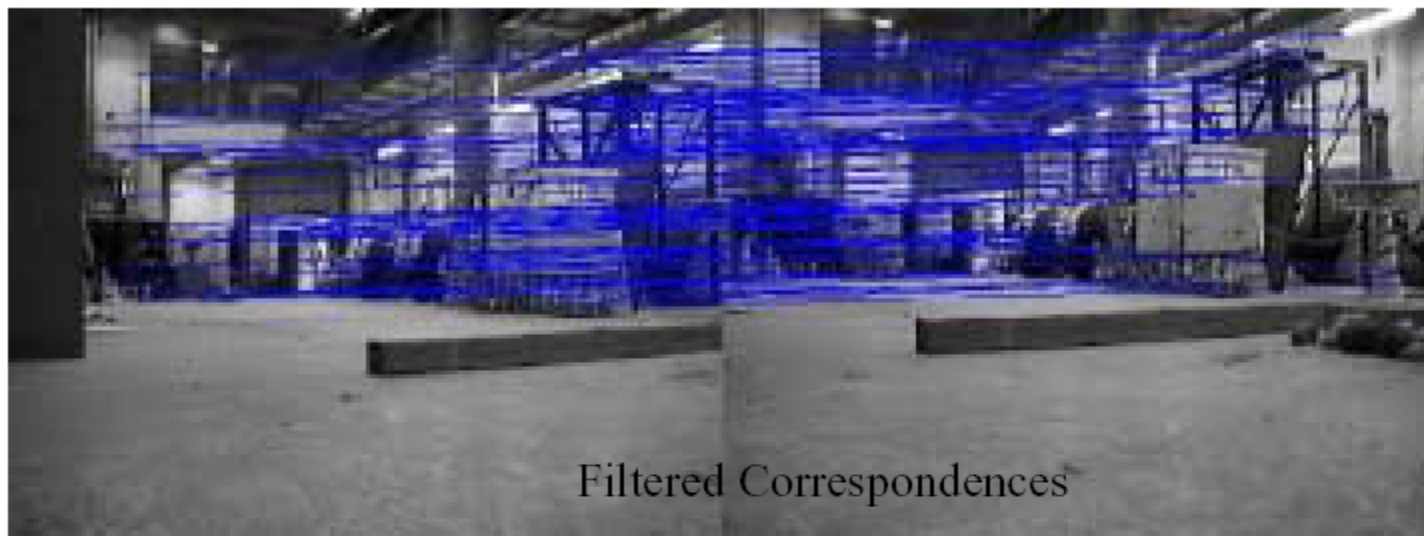
- Do  $k$  times:
  - Draw set of 8 correspondences
  - Fit  $F$  to the set
  - Count the number  $d$  of correspondences that are closer than  $t$  to the fitted epipolar lines
  - If  $d > d_{\min}$ , recompute fit error using all the correspondences
- Return best fit found

# RANSAC



Adapted from Martial Hebert, CMU

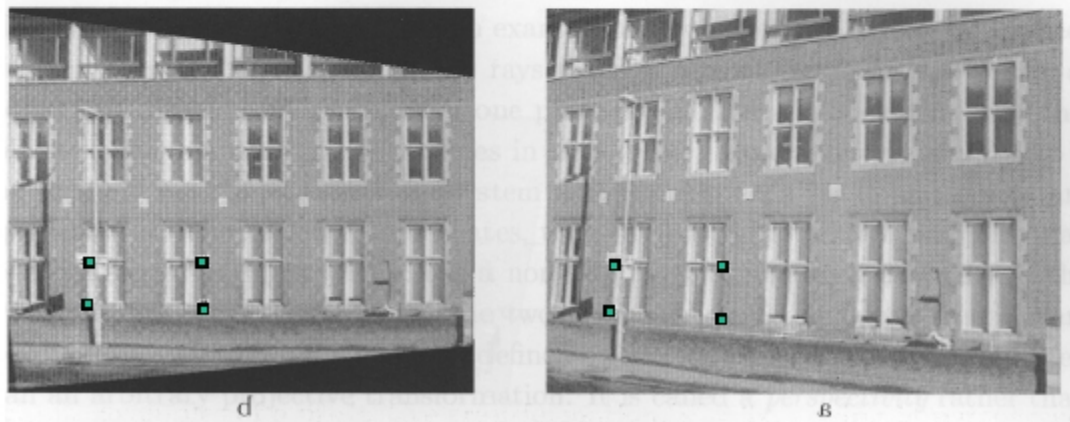
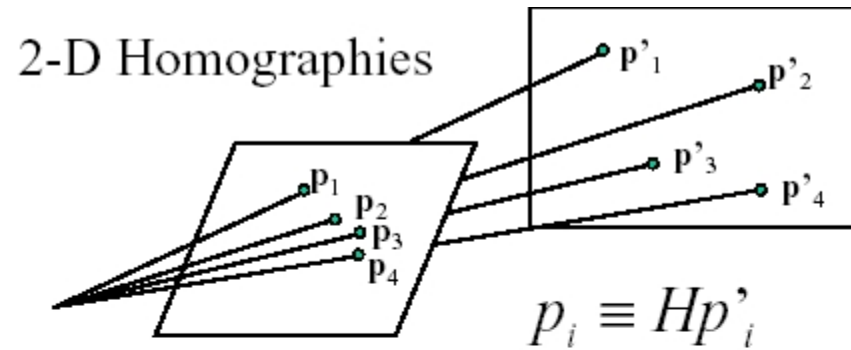
# RANSAC



Adapted from Martial Hebert, CMU

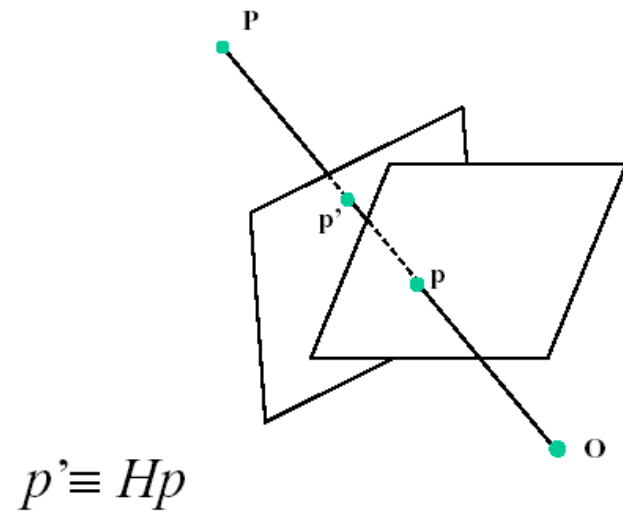
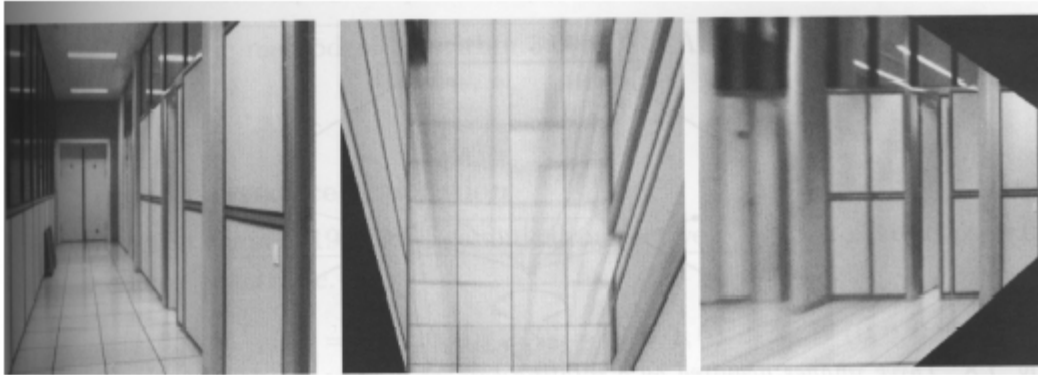


# 2D homographies



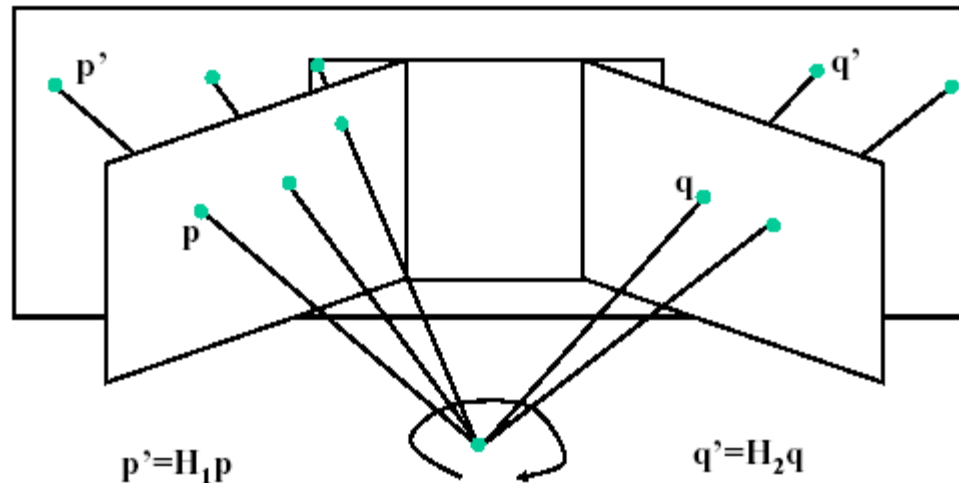
2D homographies transforms points from one plane to another

# 2D homographies



Adapted from Martial Hebert, CMU

# 2D homographies



If we choose the plane of one of the images from a set of images obtained by rotating the camera around the optical center, for each image, there exists an homography which maps the point  $p$  in the original image plane to the reference image plane. If we map all the points from all the images into the reference image plane, we obtain a single image, a mosaic, which contains the data from all the input images.

# 2D homographies

---



Adapted from Martial Hebert, CMU