



CS 554

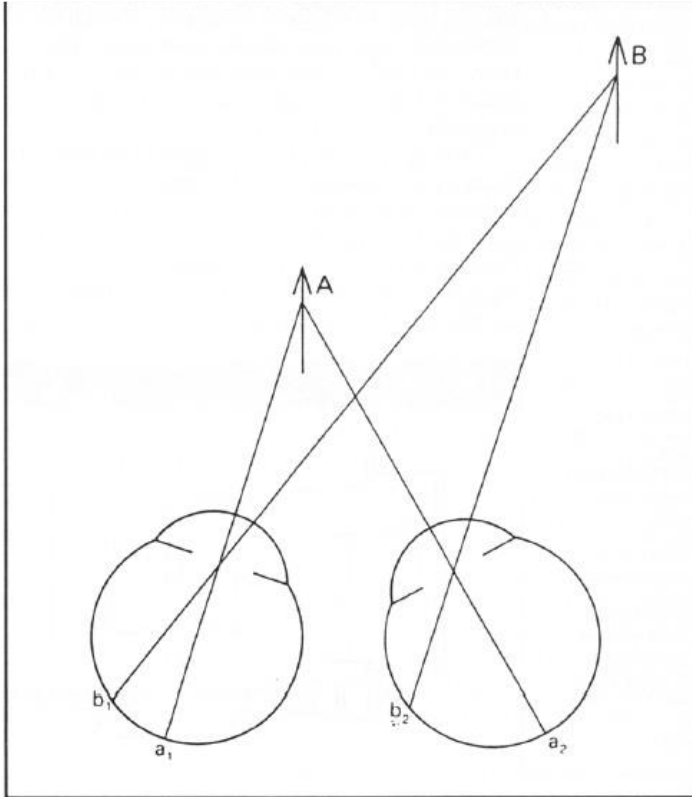
Computer Vision

Stereo Vision and Depth Estimation

Hamdi Dibeklioglu

Slide Credits: P. Duygulu Sahin, T. Darrell, M. Black,
D. Forsyth, and J. Ponce

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

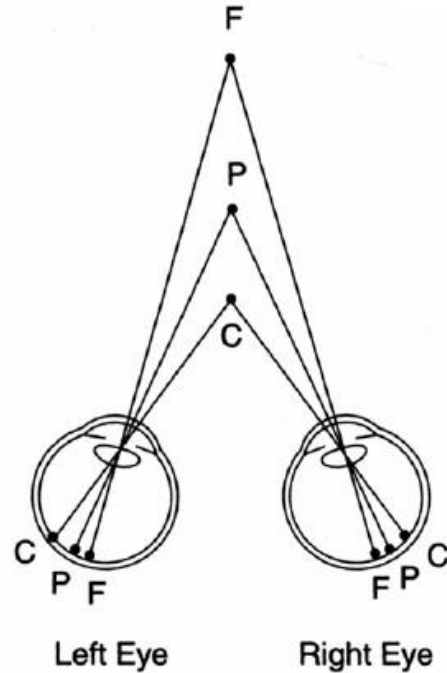


Figure 5.3.2 Crossed versus uncrossed binocular disparity. When a point *P* is fixated, closer points (such as *C*) are displaced outwardly in crossed disparity, whereas farther points (such as *F*) are displaced inwardly in uncrossed disparity.

From Palmer, “Vision Science”, MIT Press

Disparity

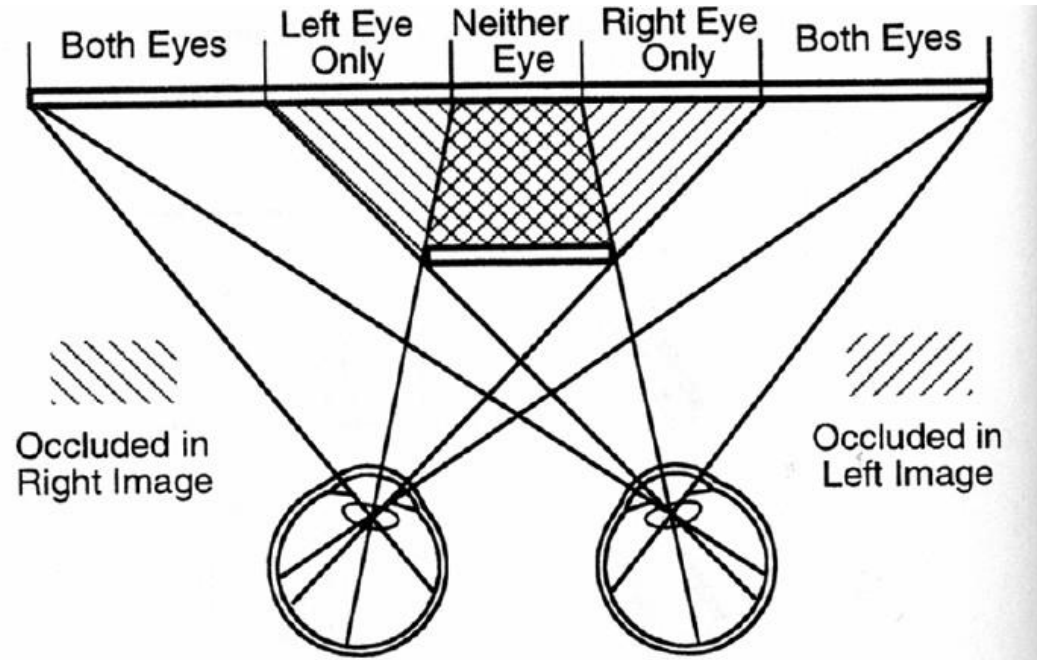


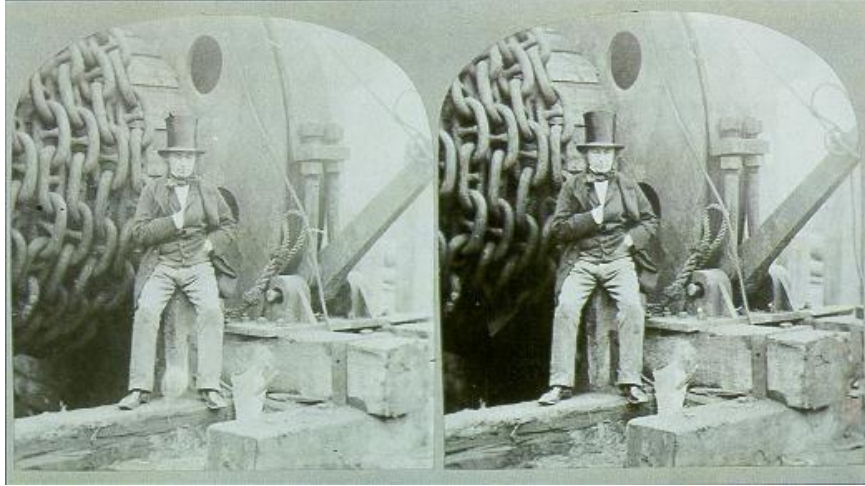
Figure 5.3.23 Da Vinci stereopsis. Depth information also arises from the fact that certain parts of one retinal image have no corresponding parts in the other image. (See text for details.)

From Palmer, "Vision Science", MIT Press

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



Stereoscope

- 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

- 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.



Left Eye Image
(Red channel only)

+



Right Eye Image
(Red channel removed)

=



Anaglyph
(Left & Right images overlaid)

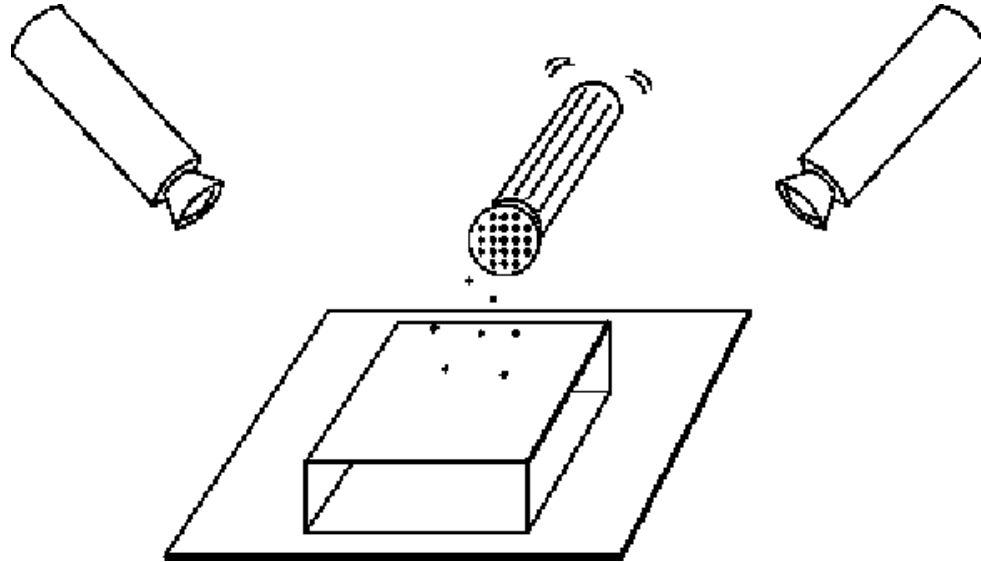
Freeview

- 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.

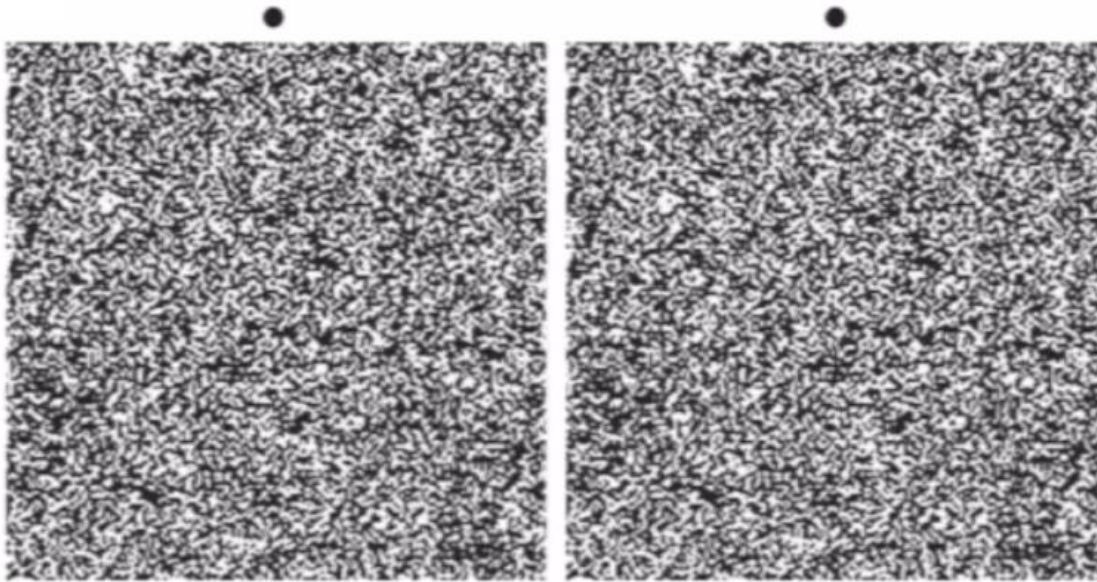




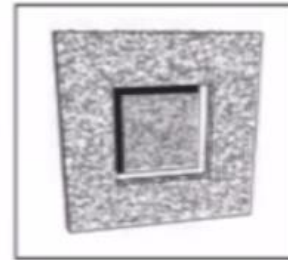
Random dot stereograms



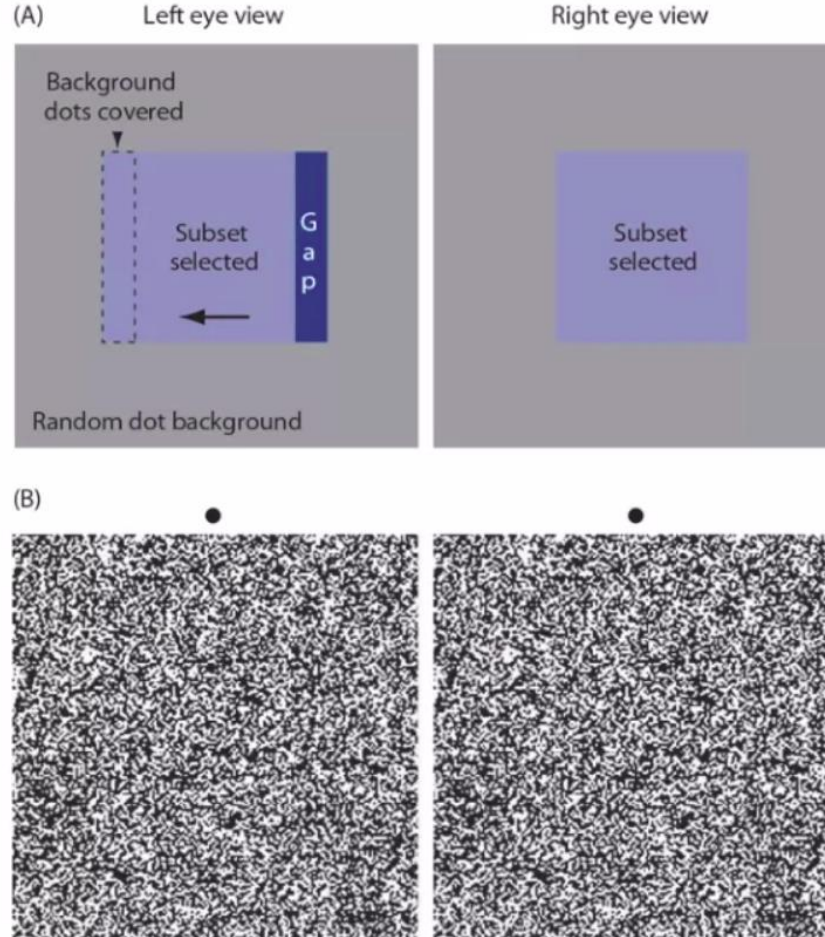
Random dot stereograms



Breaking
camouflage



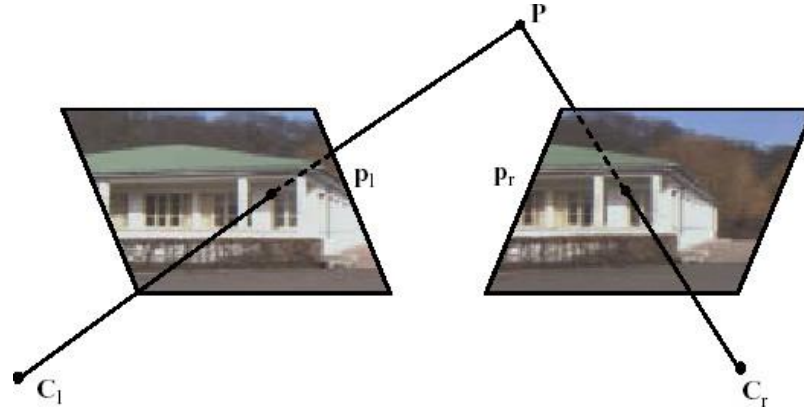
Random dot stereograms



Random dot stereograms

- When viewed monocularly, they appear random; when viewed stereoscopically, see 3D structure.
- Human binocular fusion not directly associated with the physical retinas; must involve the central nervous system (V2, for instance)
- Imaginary "cyclopean retina" that combines the left and right image stimuli as a single unit
- High level scene understanding not required for stereo

Stereo vision = correspondences + reconstruction



Stereovision involves two problems:

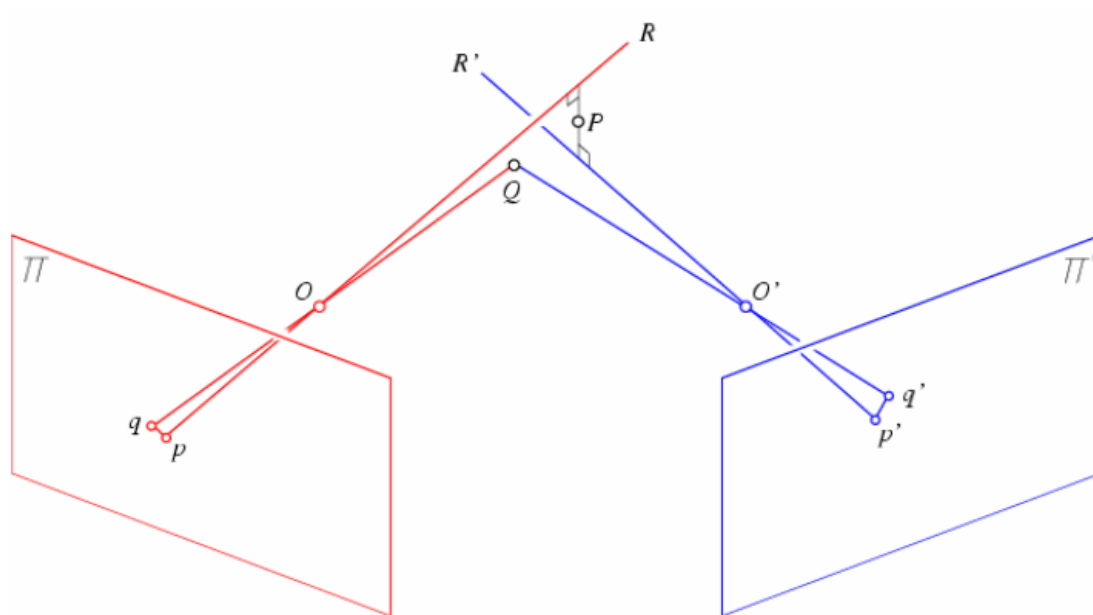
Correspondence:

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

Reconstruction

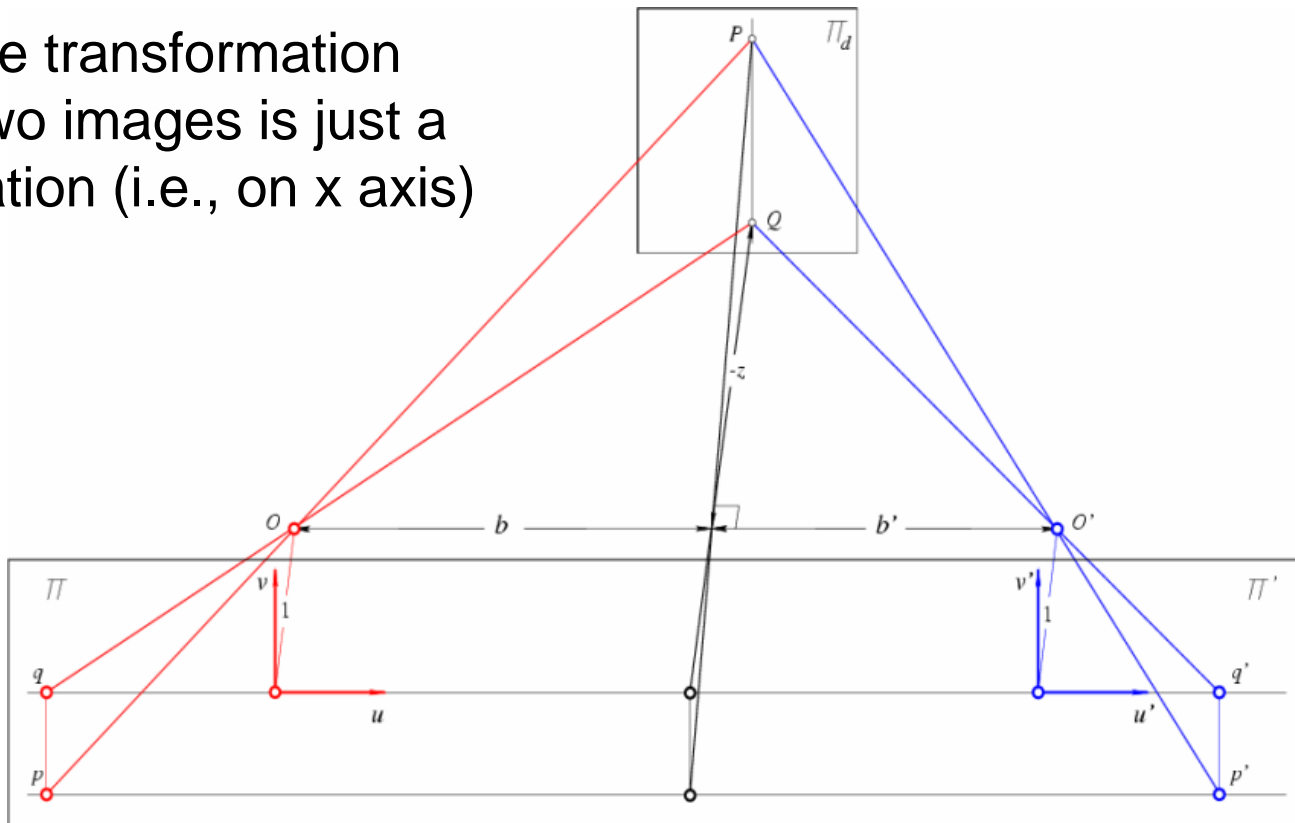


Linear Method: find P such that
$$\begin{cases} \mathbf{p} \times \mathcal{M}\mathbf{P} = 0 \\ \mathbf{p}' \times \mathcal{M}'\mathbf{P} = 0 \end{cases} \iff \begin{pmatrix} [\mathbf{p}_\times] \mathcal{M} \\ [\mathbf{p}'_\times] \mathcal{M}' \end{pmatrix} \mathbf{P} = 0$$

Non-Linear Method: find Q minimizing $d^2(p, q) + d^2(p', q')$

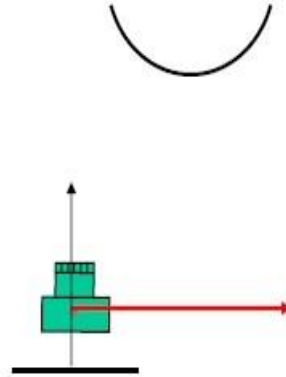
Let's simplify the task:

Assume the transformation between two images is just a 1-D translation (i.e., on x axis)

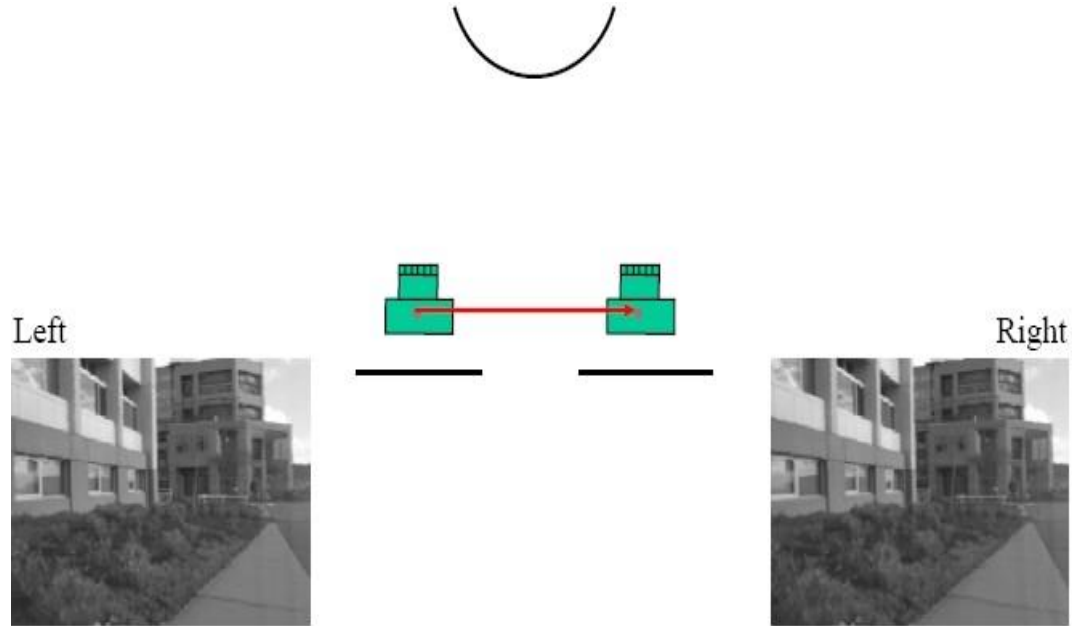


Binocular Stereo

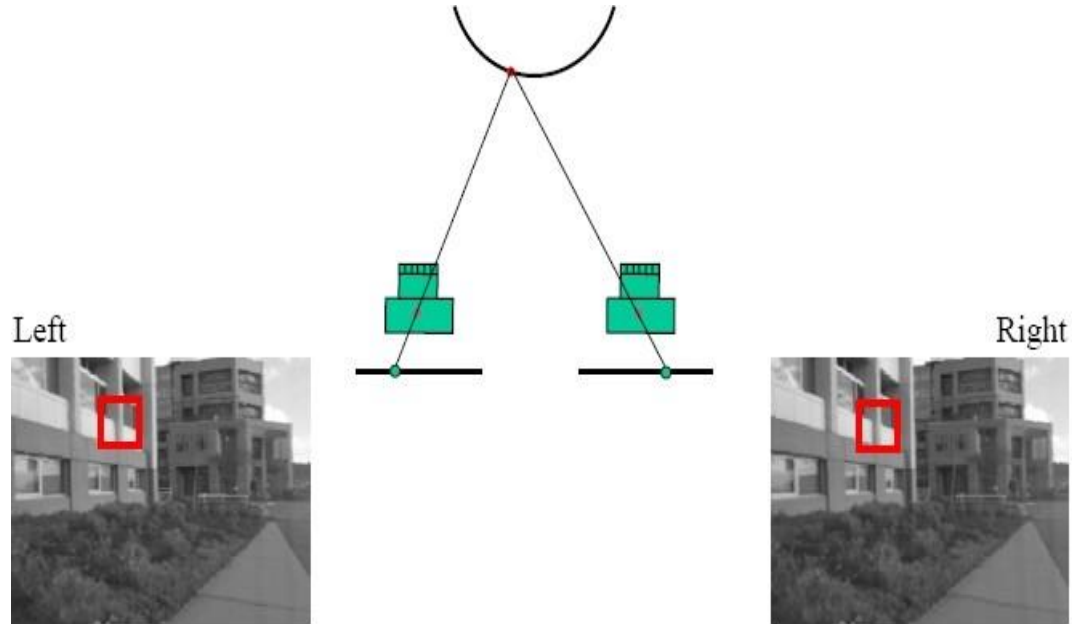
Left



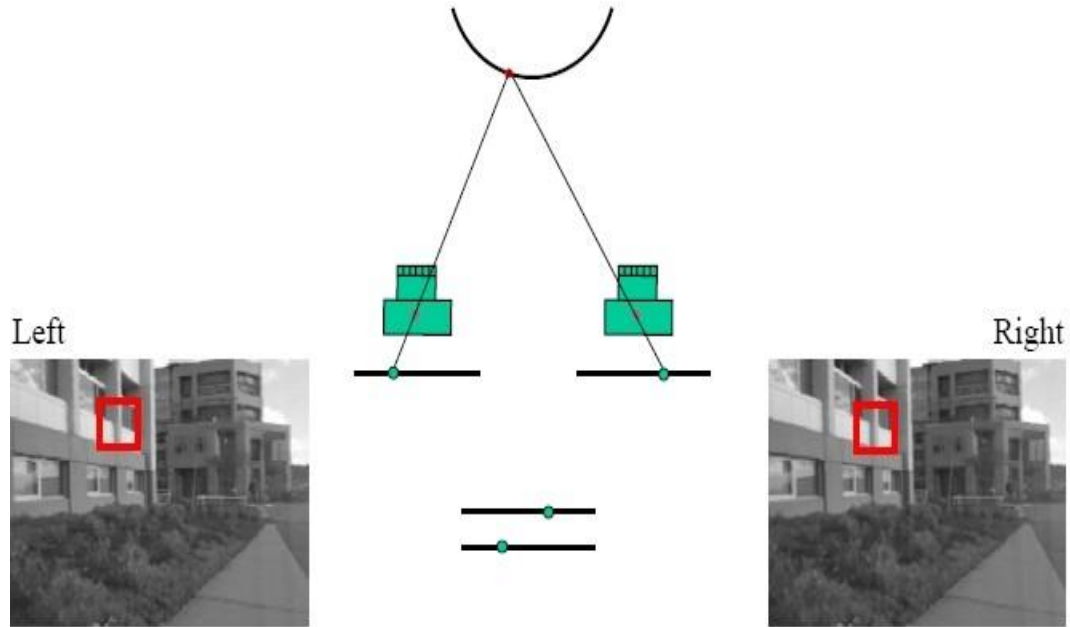
Binocular Stereo



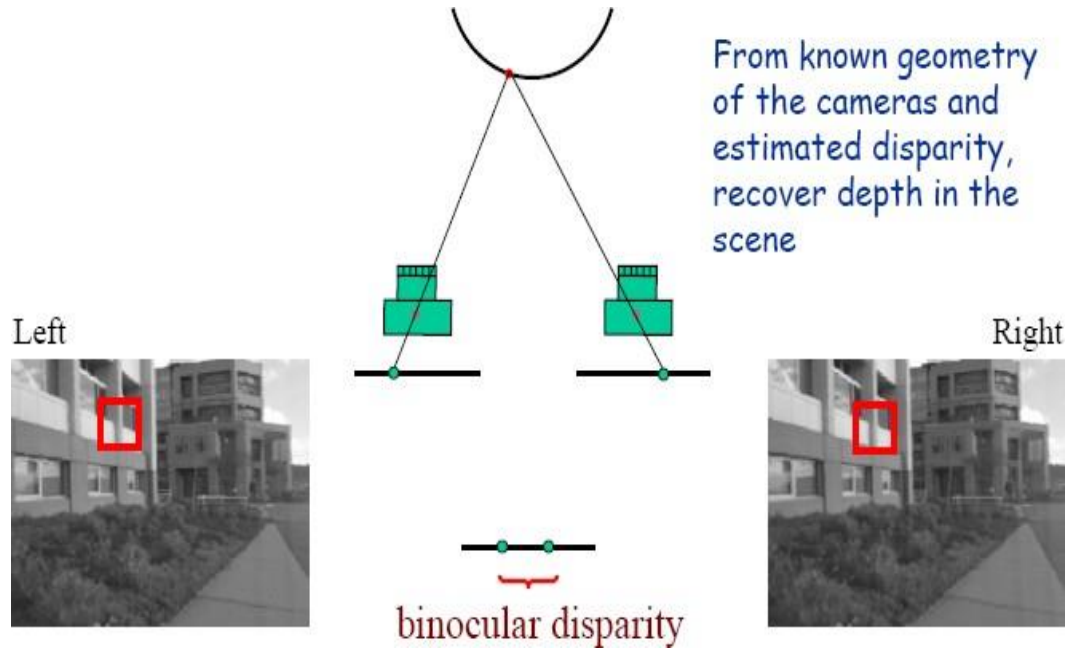
Binocular Stereo



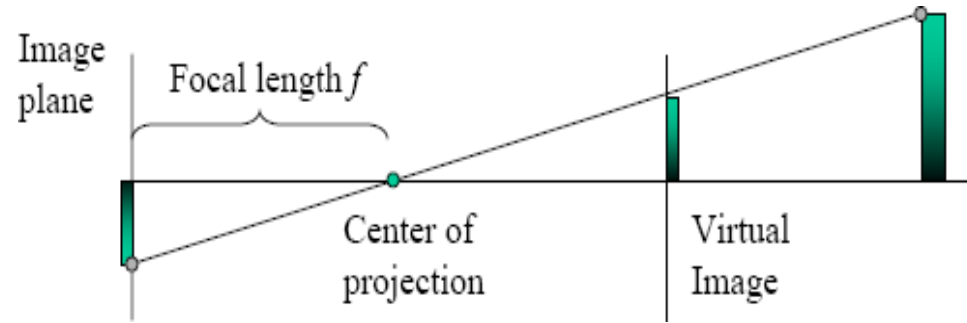
Binocular Stereo



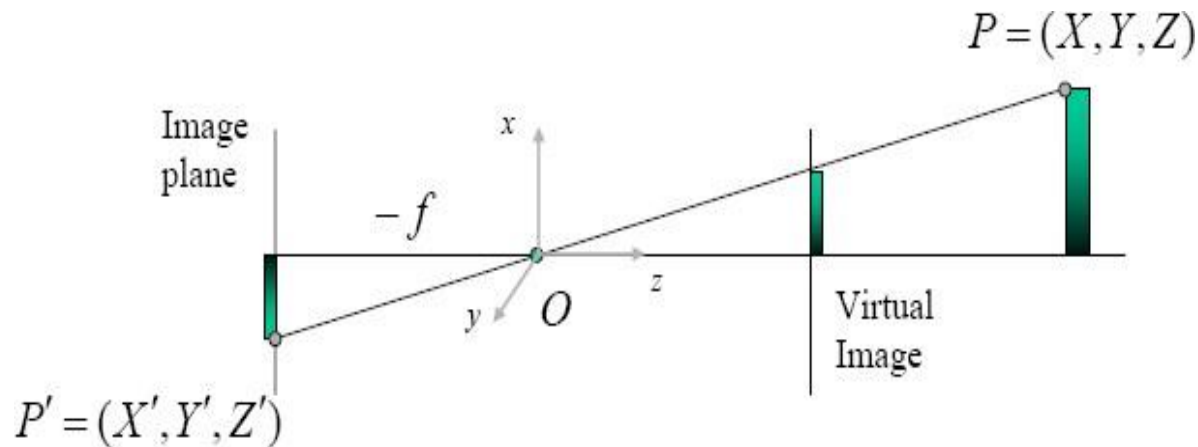
Binocular Stereo



Depth Estimation



Depth Estimation

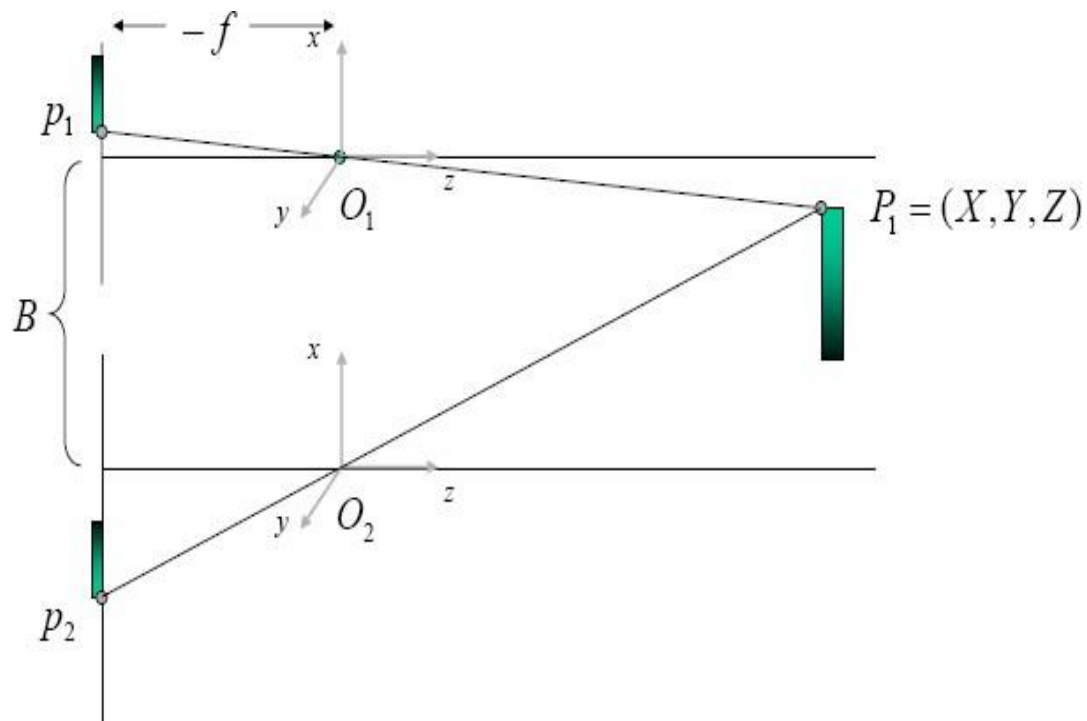


$$Z' = -f, \quad X' = -f \frac{X}{Z}, \quad Y' = -f \frac{Y}{Z}$$

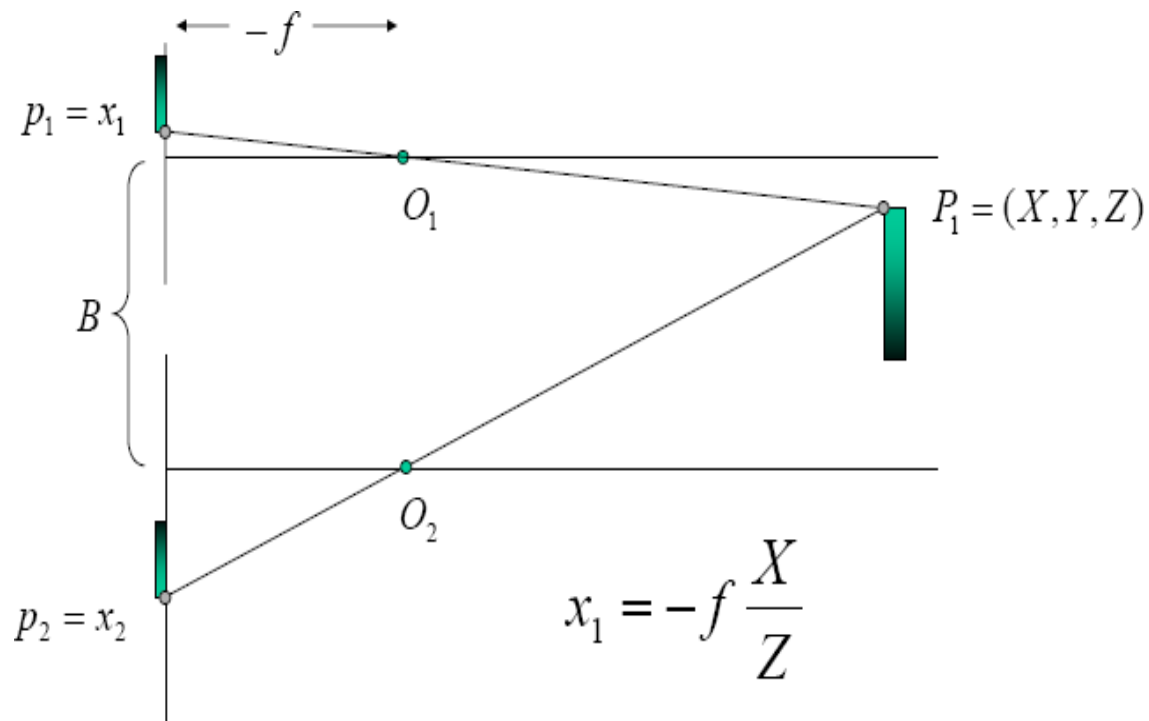
$$x = -X', \quad y = -Y'$$

$$(X, Y, Z) \rightarrow (x, y, 1) = \left(f \frac{X}{Z}, f \frac{Y}{Z}, 1\right)$$

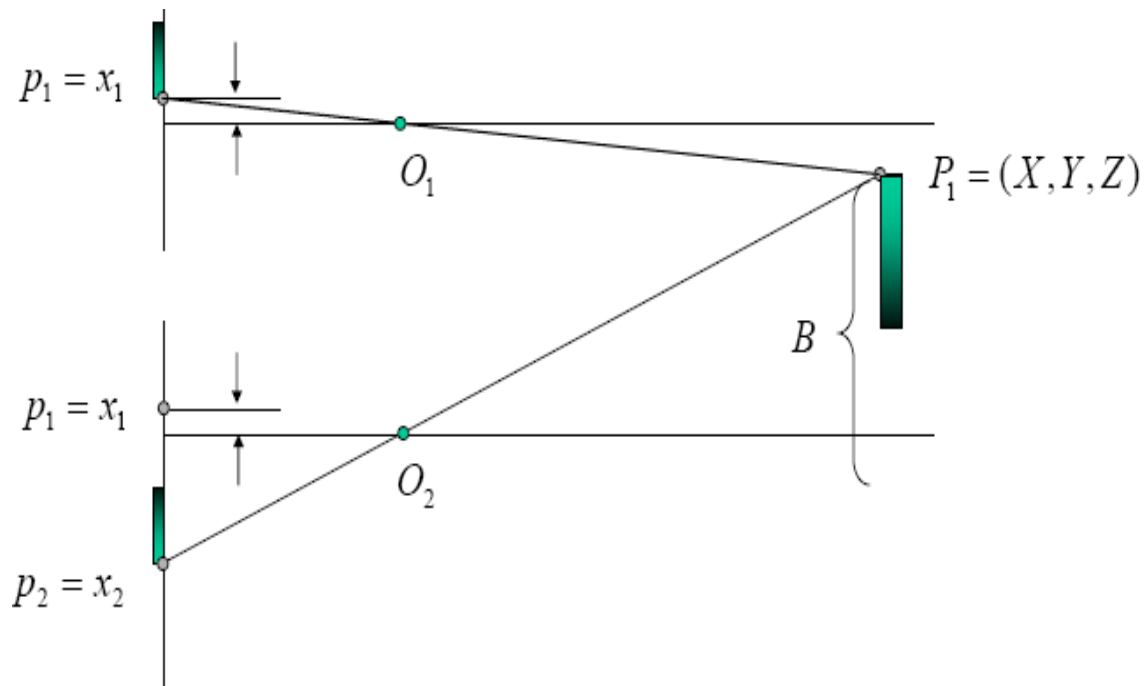
Depth Estimation



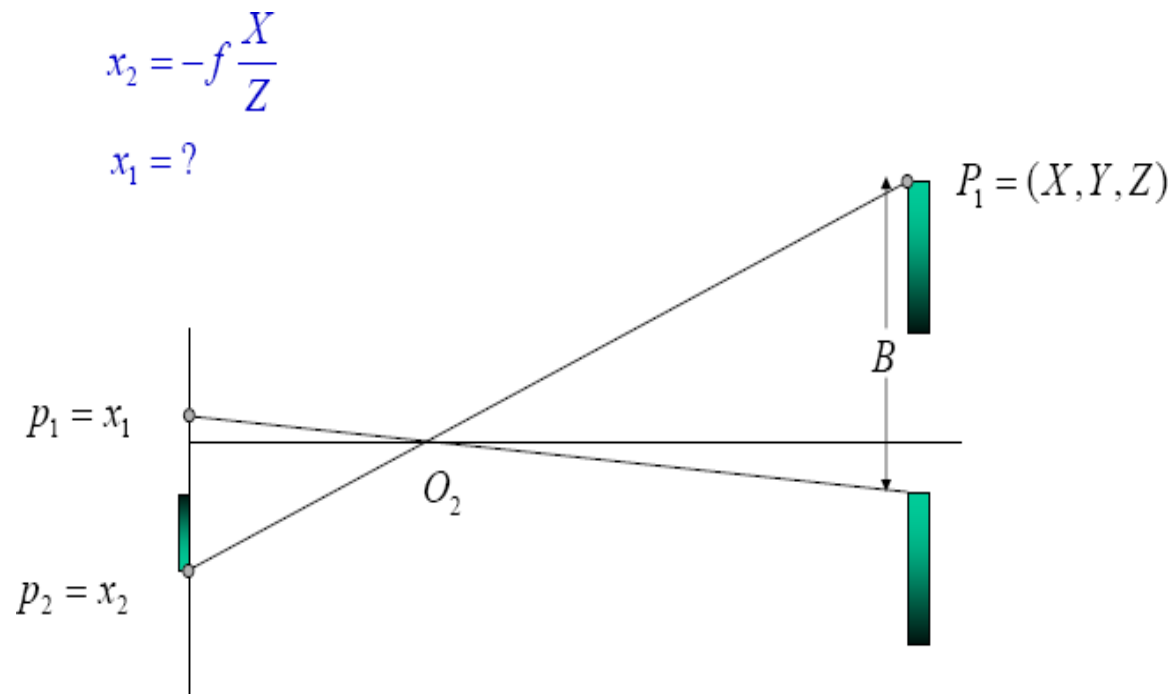
Depth Estimation



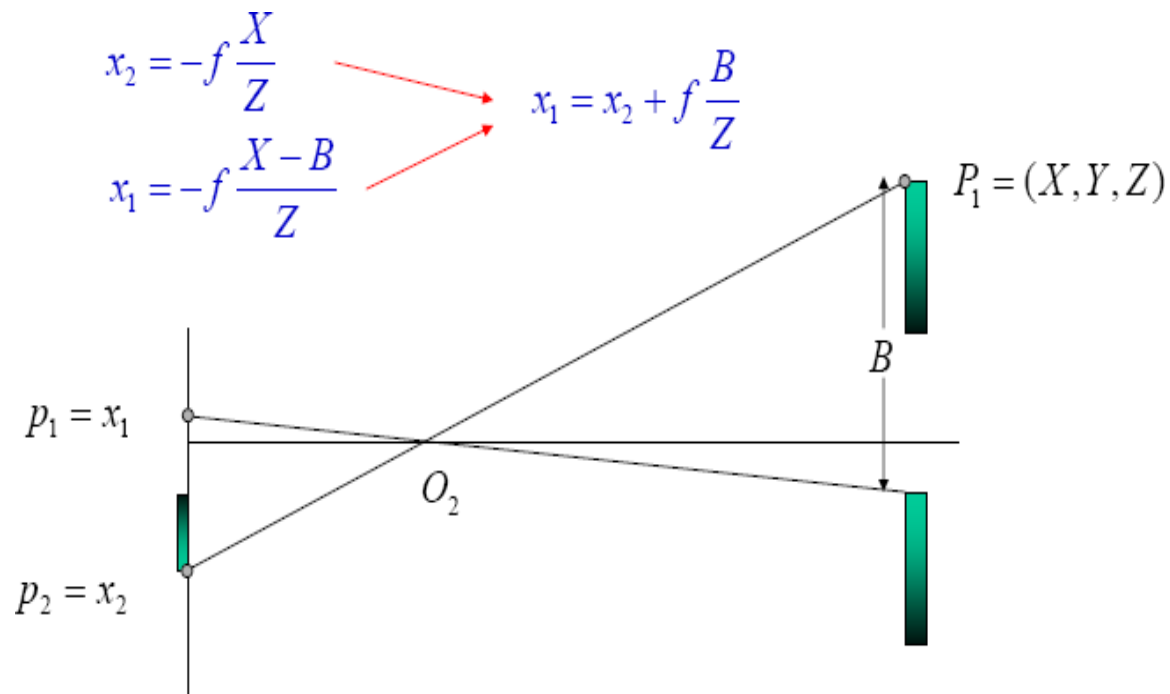
Depth Estimation



Depth Estimation



Depth Estimation

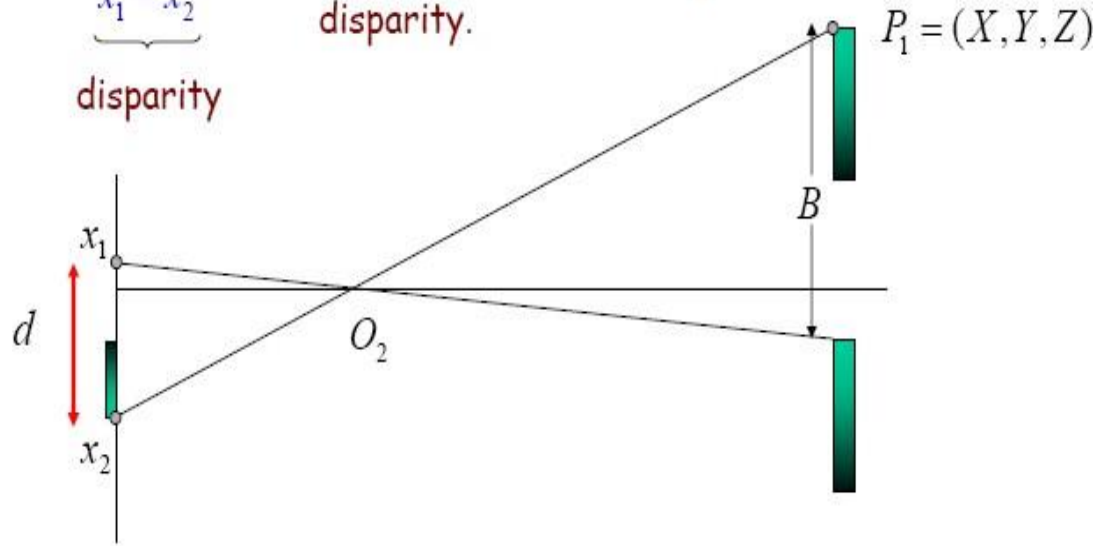


Depth Estimation

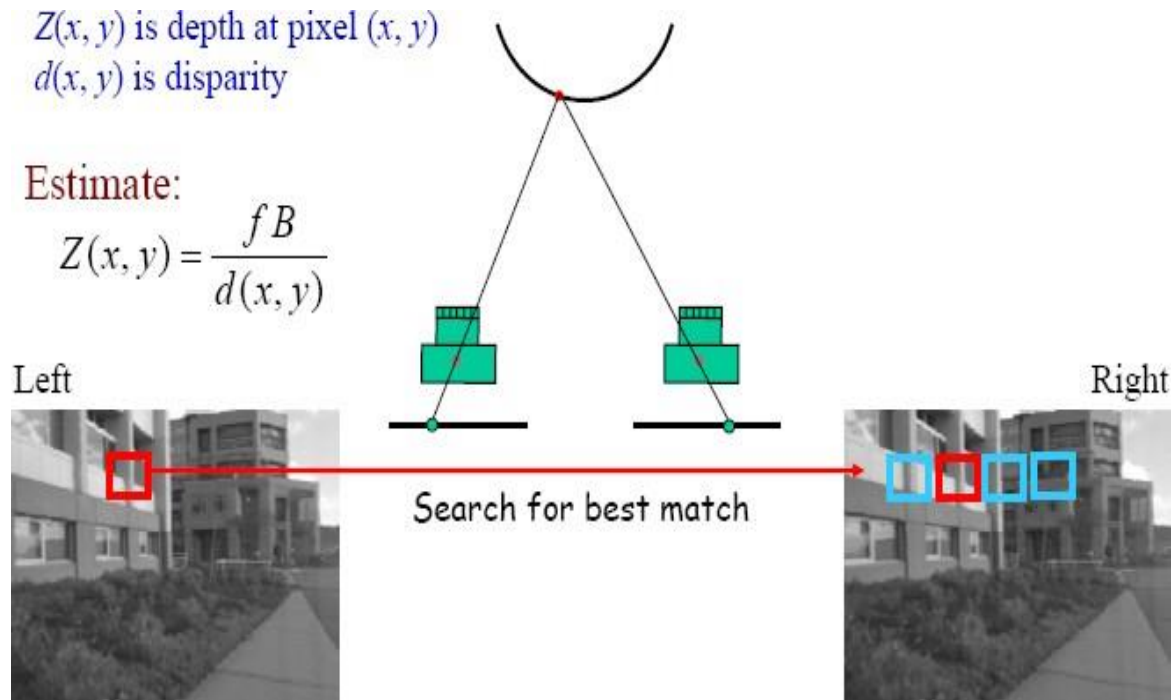
$$Z = \frac{fB}{\underbrace{x_1 - x_2}_{\text{disparity}}}$$

disparity

For a calibrated camera, we know f and the baseline B . Then depth can be computed from disparity.



Correspondence



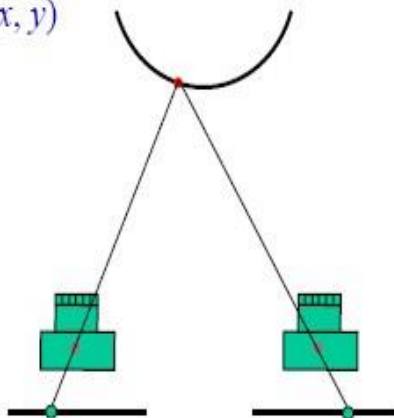
Correspondence

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

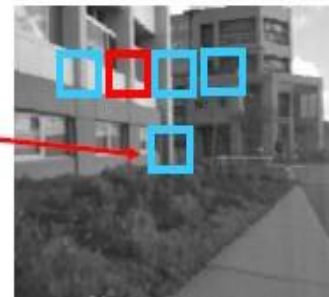
Estimate:

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

Left



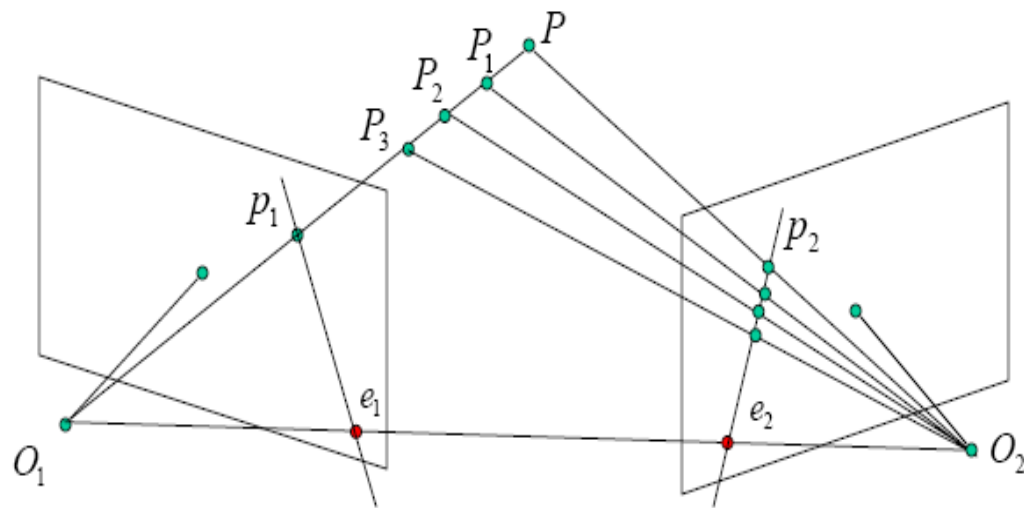
Right



Do I need to consider
this region?

Correspondence

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



Recap: Camera Geometry & Homography

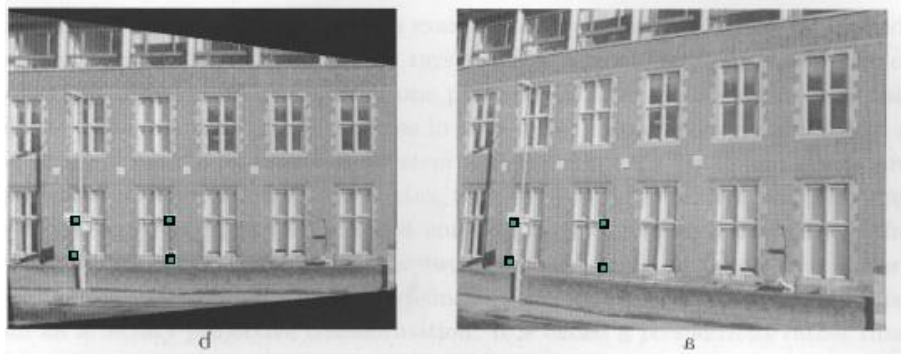
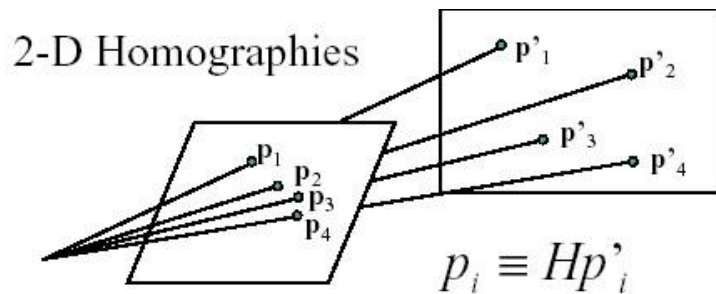
$$p = \mathbf{H}p' \quad \rightarrow \text{Planar homography } \mathbf{H}$$

$$p^T \boldsymbol{\varepsilon} p' = 0 \quad \rightarrow \text{Essential matrix } \boldsymbol{\varepsilon}$$

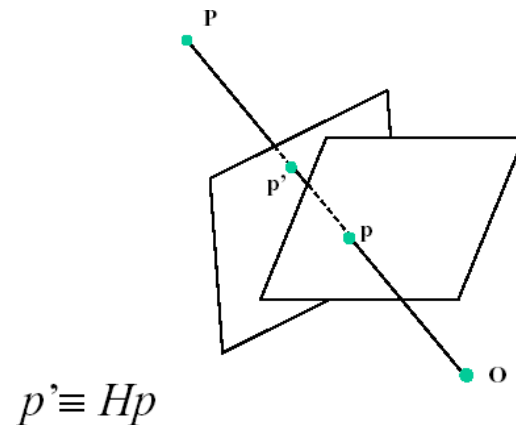
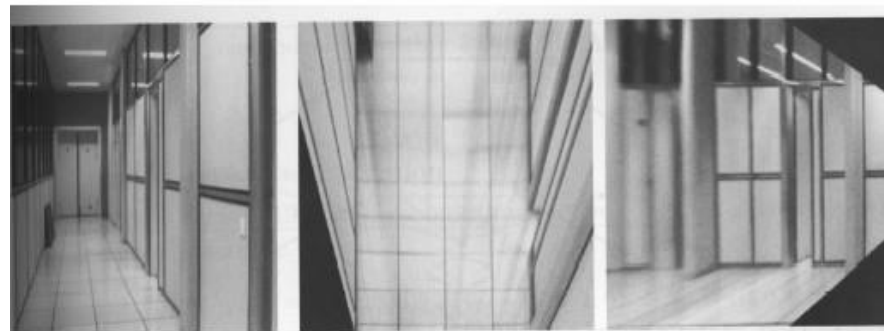
$$p^T \boldsymbol{\mathcal{F}} p' = 0 \quad \rightarrow \text{Fundamental matrix } \boldsymbol{\mathcal{F}}$$

where, $\boldsymbol{\mathcal{F}} = \mathcal{K}^{-T} \boldsymbol{\varepsilon} \mathcal{K}'^{-1}$,
and $\mathcal{K}, \mathcal{K}'$ are calibration matrices

Recap: Homography



2D homographies transforms points from one plane to another



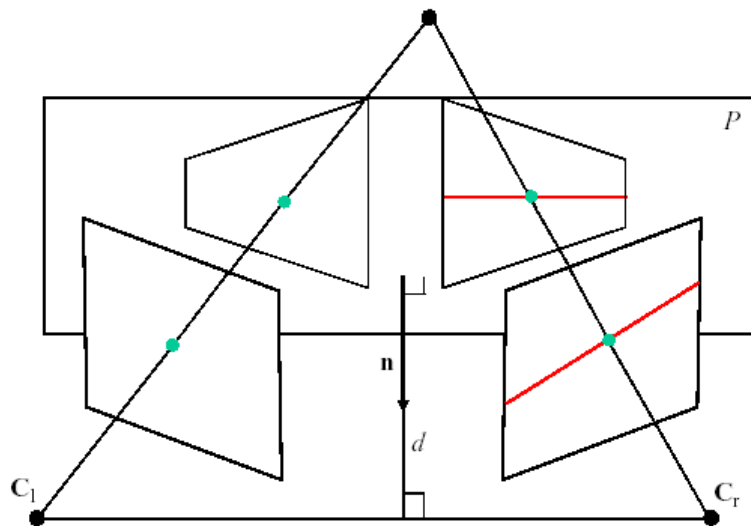
Rectification



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

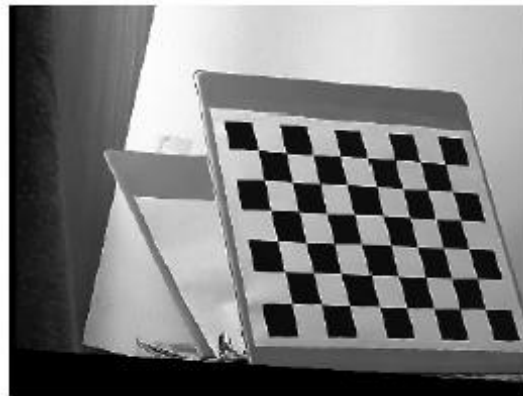
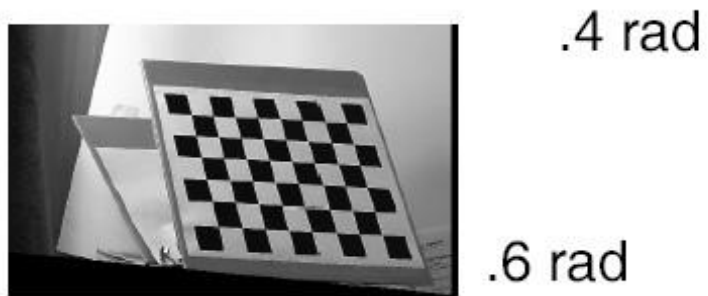
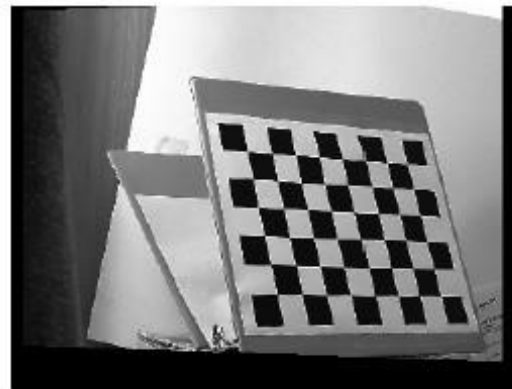
We know that, given a plane \mathbf{P} in space, there exists two homographies \mathbf{H}_l and \mathbf{H}_r that map each image plane onto \mathbf{P} . That is, if \mathbf{p}_l is a point in the left image, then the corresponding point in \mathbf{P} is $\mathbf{H}_l \mathbf{p}_l$ (in homogeneous coordinates). If we map both images to a common plane \mathbf{P} such that \mathbf{P} is parallel to the line $\mathbf{C}_l \mathbf{C}_r$ 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.



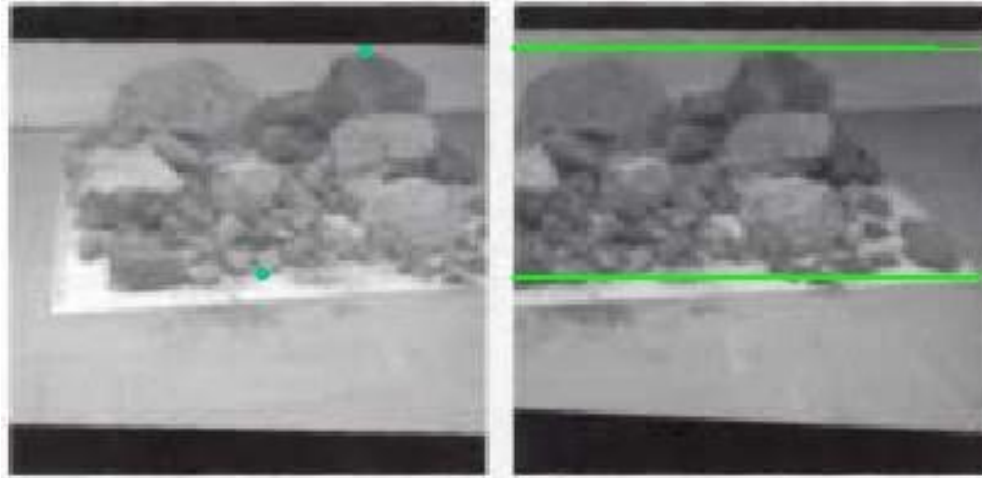
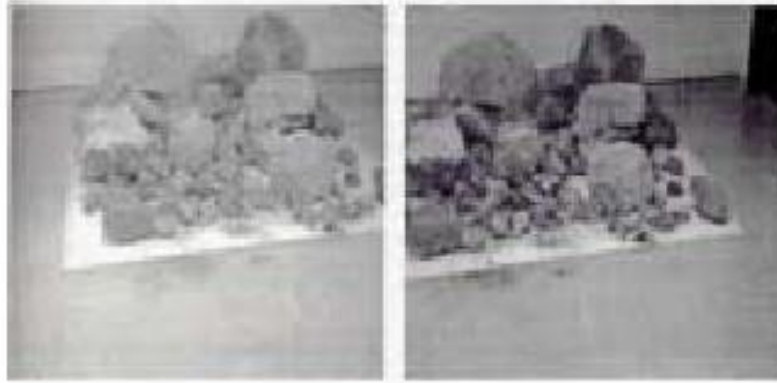
The algorithm for rectification is then:

- Select a plane \mathbf{P} parallel to $\mathbf{C}_r \mathbf{C}_l$
- 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 Martiel Hebert

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.



$$d = x_l - x_r$$

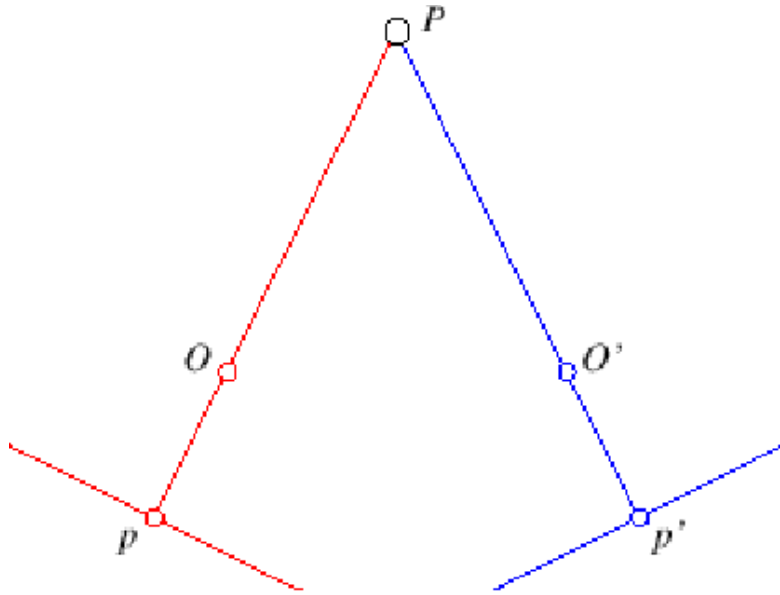
Disparity



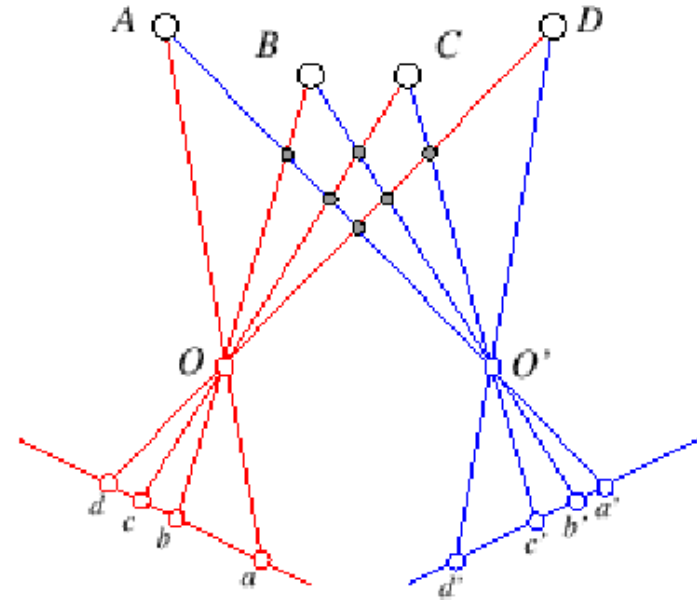
Larger disparity \rightarrow closer to camera

Adapted from Martiel Hebert

Stereopsis



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

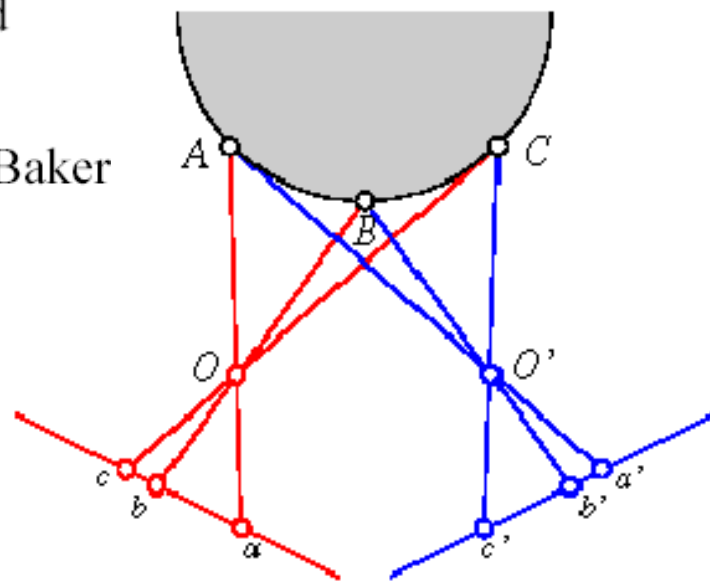


However, each picture consists of hundreds/thousands of pixels,
therefore it is very hard to find the correct correspondences

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].



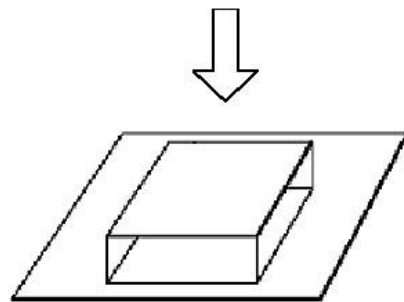
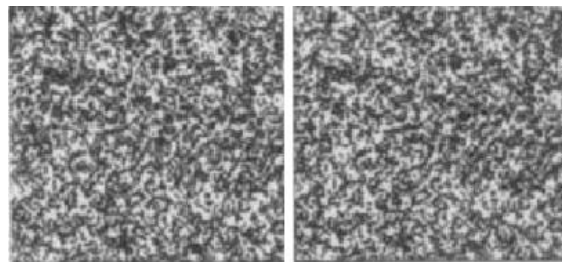
Correspondence is ambiguous (Marr & Poggio)

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

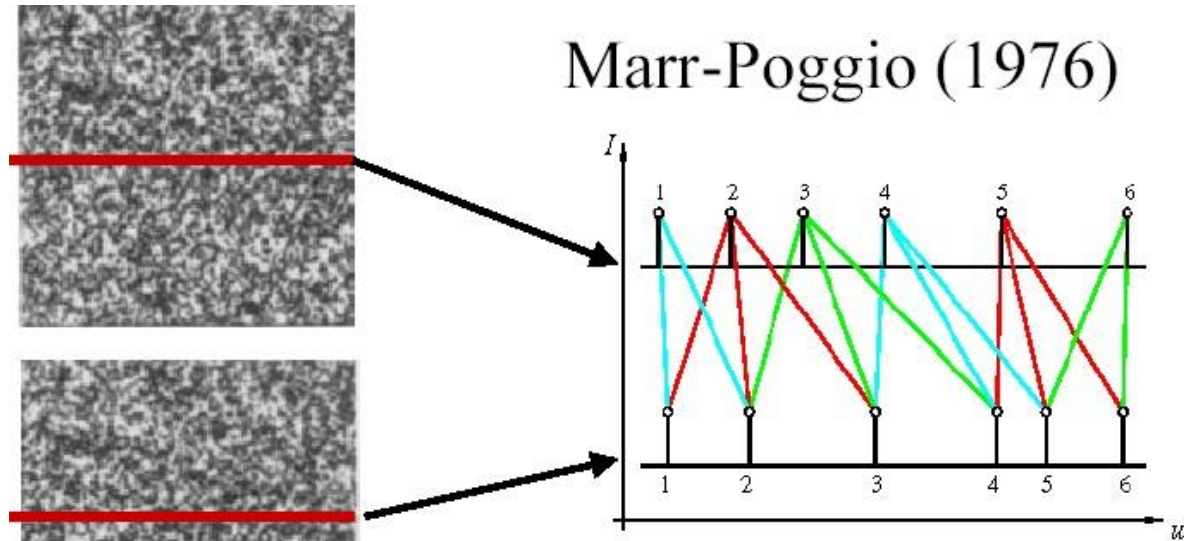
Uniqueness: a black dot in one image matches at most one black dot in another image

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



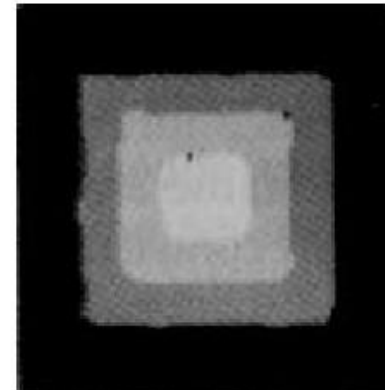
Correspondence is ambiguous

Adapted from Trevor Darrell, MIT



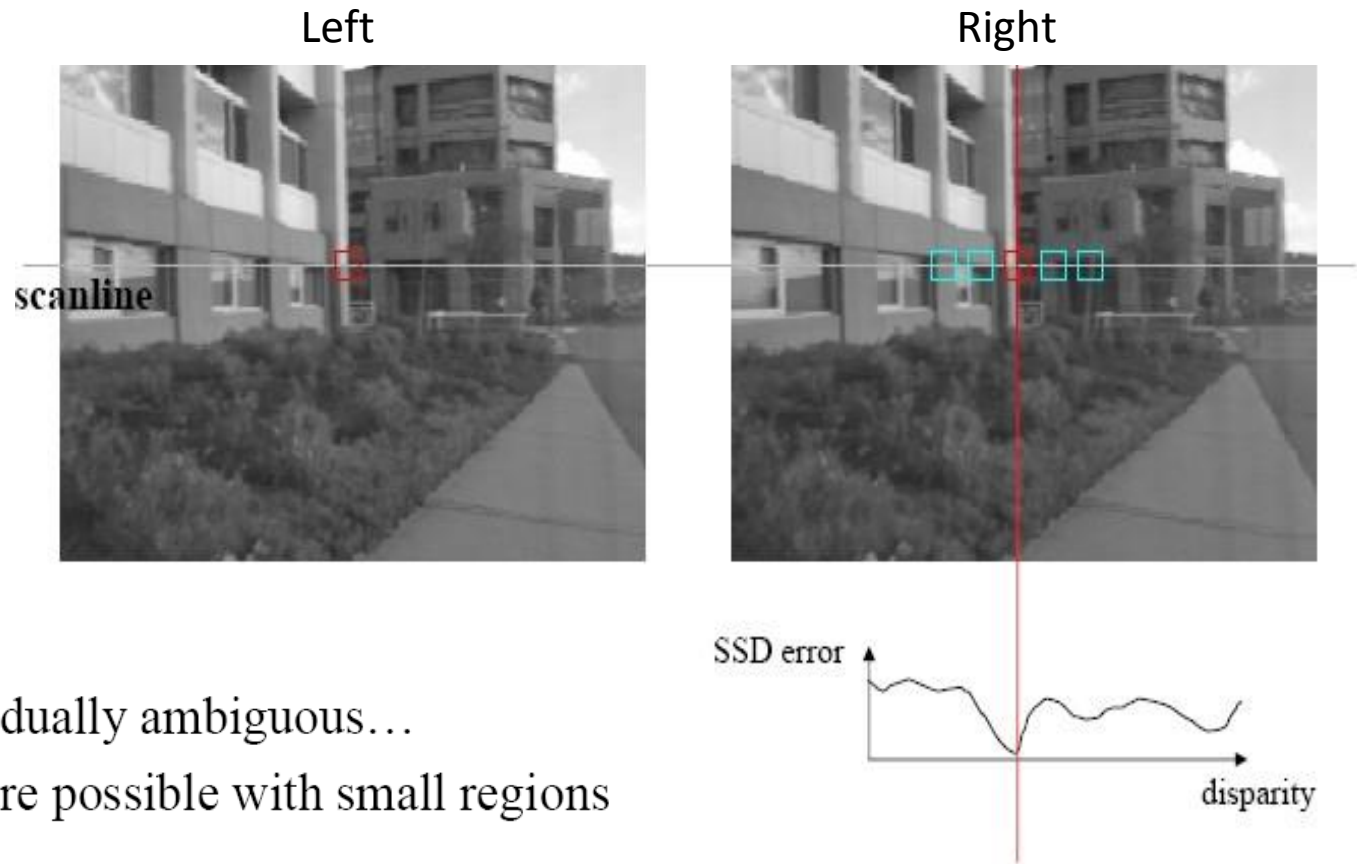
Three constraints:

- *compatibility*
- *uniqueness*
- *continuity*



Works well on RDS....but not so well on natural images...

Correspondence using window matching



Points are highly individually ambiguous...

More unique matches are possible with small regions of image.

Finding Correspondences



$W(\mathbf{p}_l)$



$W(\mathbf{p}_r)$

Adapted from Martiel Hebert

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

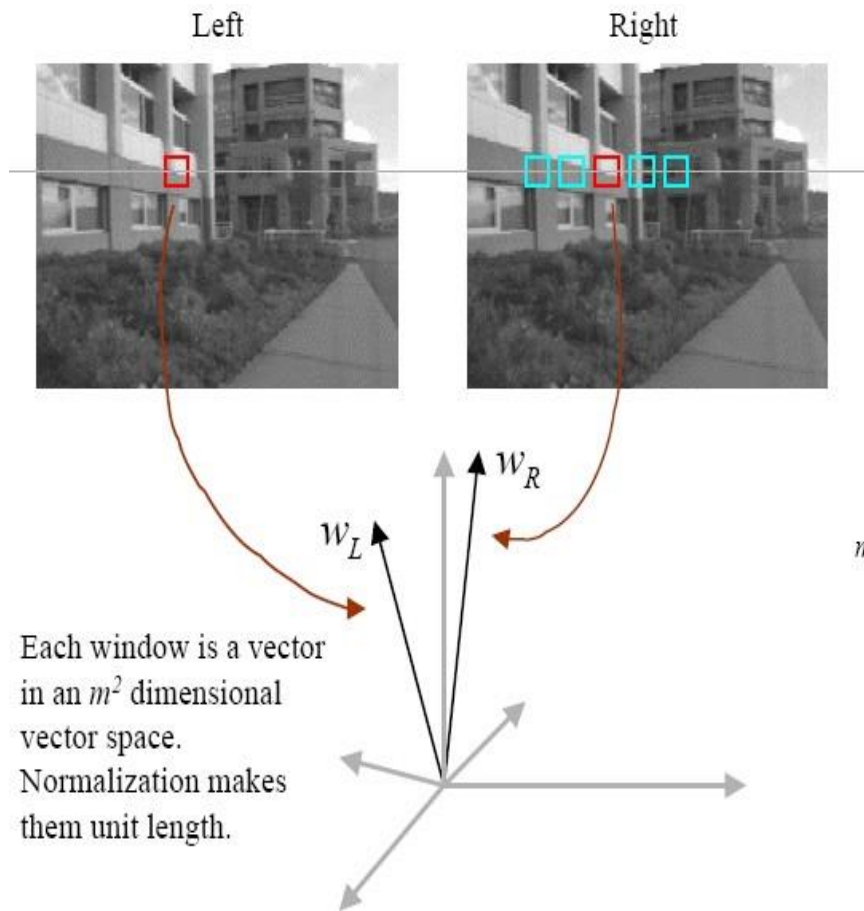
$$\|I\|_{W_m(x,y)} = \sqrt{\sum_{(u,v) \in W_m(x,y)} [I(u,v)]^2}$$

Window magnitude

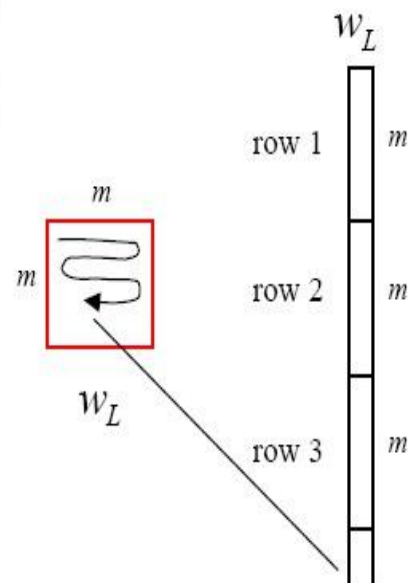
$$\hat{I}(x,y) = \frac{I(x,y) - \bar{I}}{\|I - \bar{I}\|_{W_m(x,y)}}$$

Normalized pixel

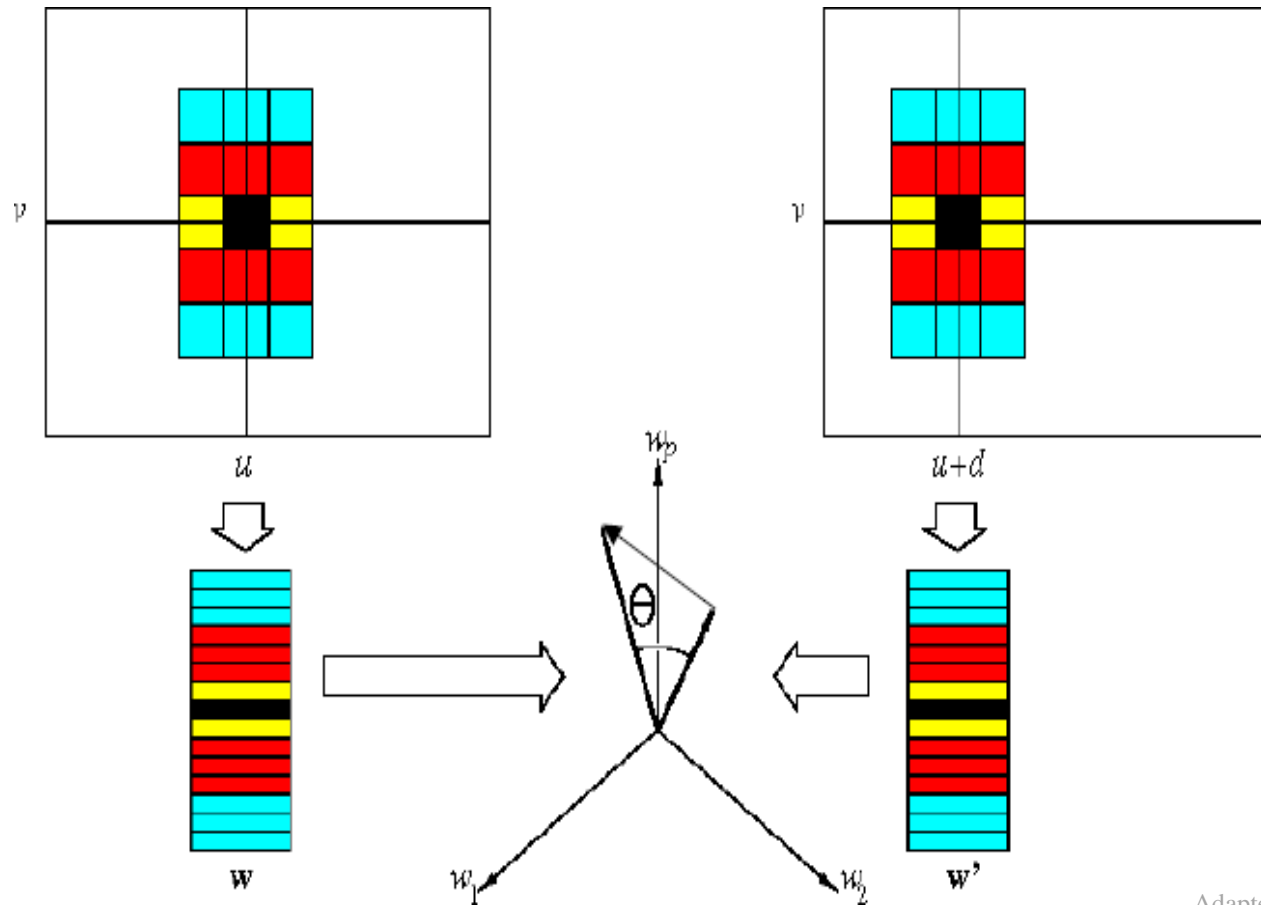
Images as vectors



“Unwrap”
image to form
vector, using
raster scan order

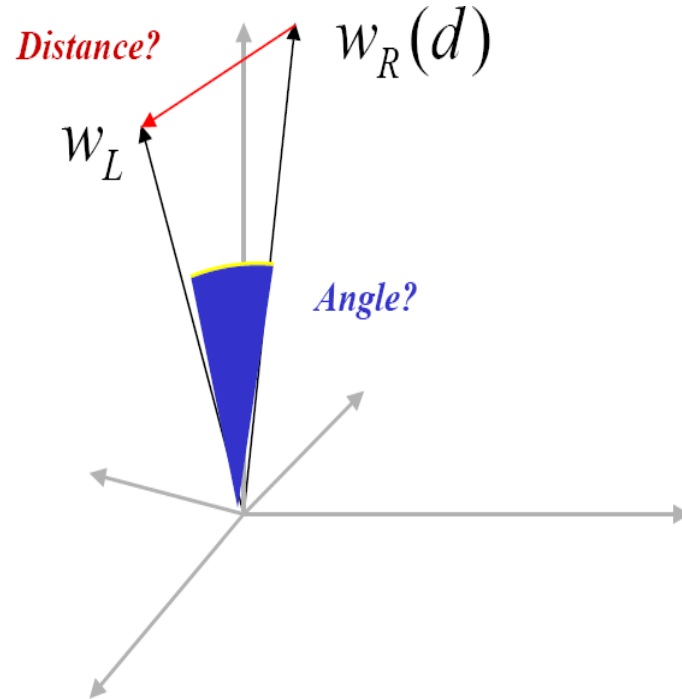


Images as vectors

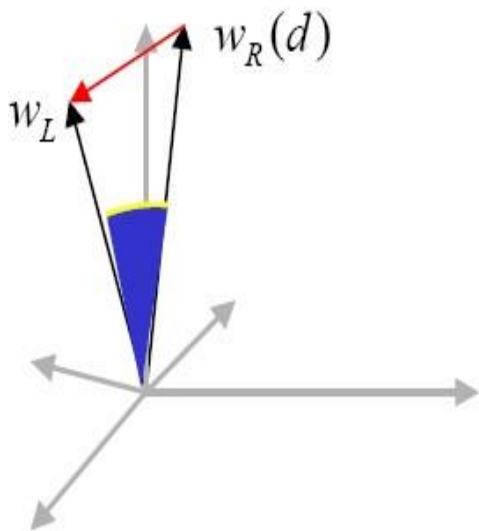


Adapted from Trevor Darrell, MIT

Possible metric



Possible metric



(Normalized) Sum of Squared Differences

$$\begin{aligned} C_{\text{SSD}}(d) &= \sum_{(u,v) \in \mathcal{W}_m(x,y)} [\hat{I}_L(u,v) - \hat{I}_R(u-d,v)]^2 \\ &= \|w_L - w_R(d)\|^2 \end{aligned}$$

Normalized Correlation

$$\begin{aligned} C_{\text{NC}}(d) &= \sum_{(u,v) \in \mathcal{W}_m(x,y)} \hat{I}_L(u,v) \hat{I}_R(u-d,v) \\ &= w_L \cdot w_R(d) = \cos \theta \end{aligned}$$

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

Matching using correlation

Left



Disparity Map



Images courtesy of Point Grey Research

Matching using correlation

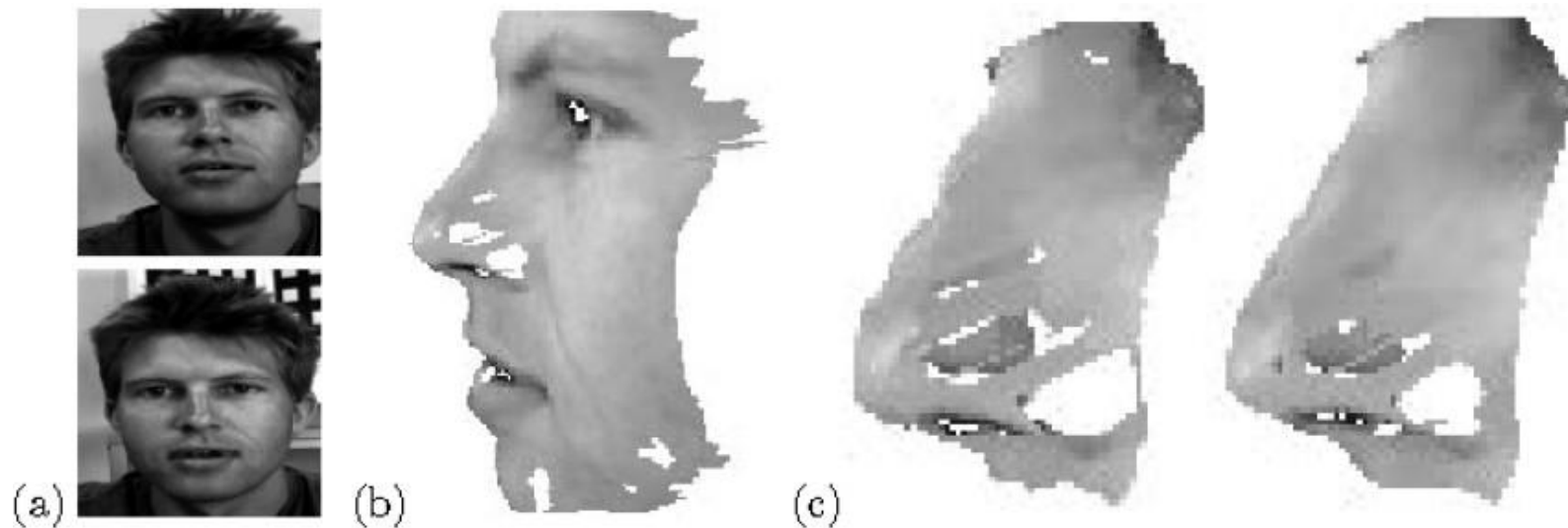


FIGURE 12.13: Correlation-based stereo matching: (a) a pair of stereo pictures; (b) a texture-mapped view of the reconstructed surface; (c) comparison of the regular (left) and refined (right) correlation methods in the nose region. Reprinted from [Devernay and Faugeras, 1994], Figures 5, 8 and 9.

Problems with window methods

Patch too small?

Patch too large?

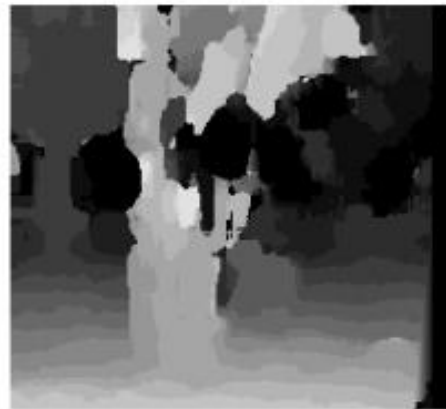
*Can try variable patch size [Okutomi and Kanade],
or arbitrary window shapes [Veksler and Zabih]*

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

Window size



$W = 3$



$W = 20$

- Effect of window size
 - Smaller window
 - + More details
 - More noise
 - Larger window
 - + Less noise
 - Less details

Better results with *adaptive window*

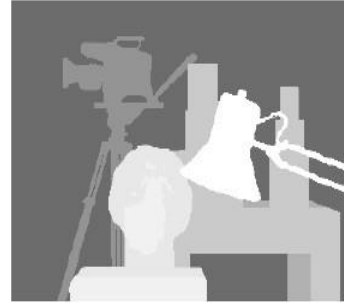
- T. Kanade and M. Okutomi,
A Stereo Matching Algorithm with an Adaptive Window: Theory and Experiment., Proc.
International Conference on Robotics and
Automation, 1991.
- D. Scharstein and R. Szeliski.
Stereo matching with nonlinear diffusion.
International Journal of Computer Vision, 28(2):
155-174,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

Stereo Results

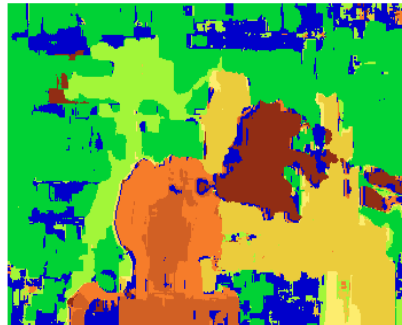
– Data from University of Tsukuba



Scene



Ground truth



Window-based matching
(best window size)



Ground truth

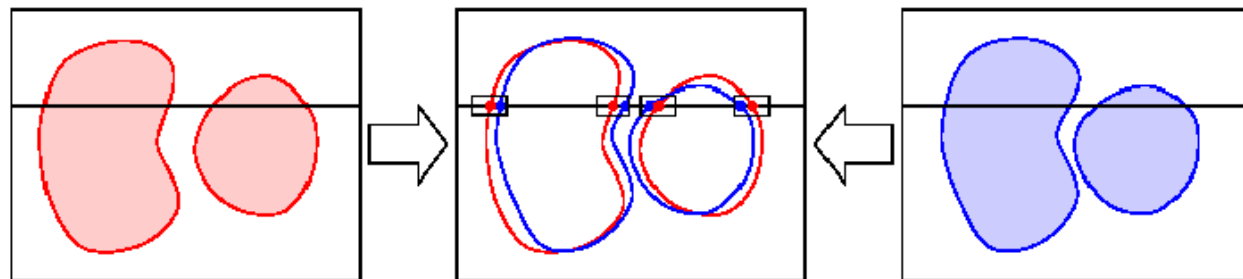
(Seitz)

Adapted from Michael Black

Multi-Scale Edge Matching (Marr, Poggio and Grimson, 1979-81)

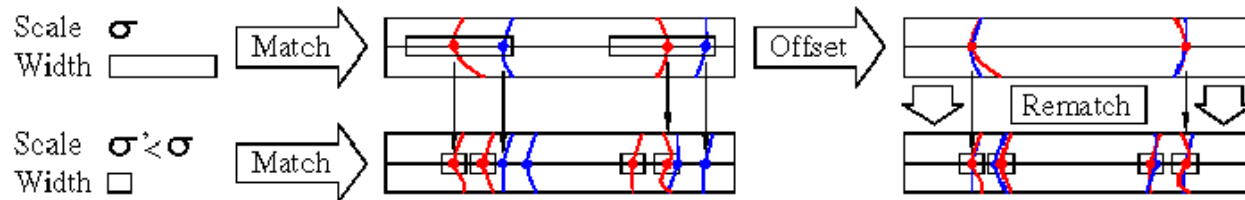
Search for edges, a.k.a “zero crossings”

Matching zero-crossings at a single scale



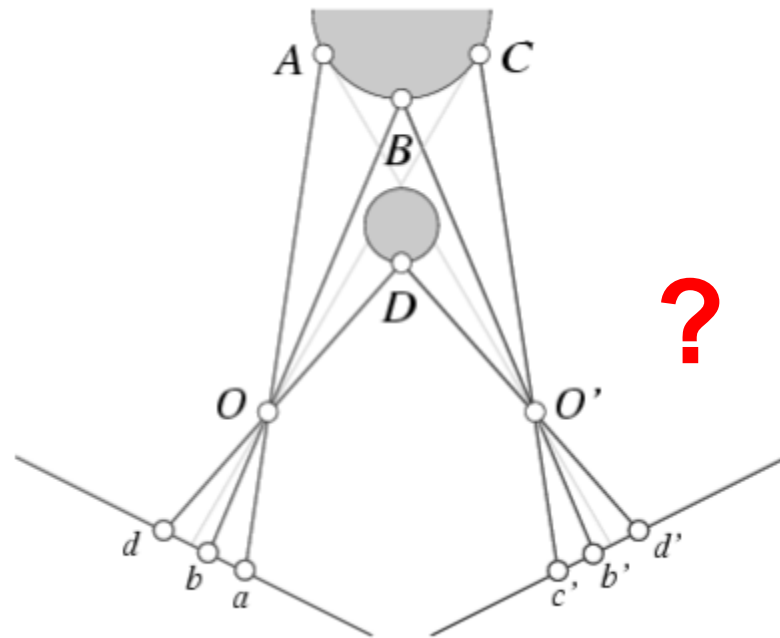
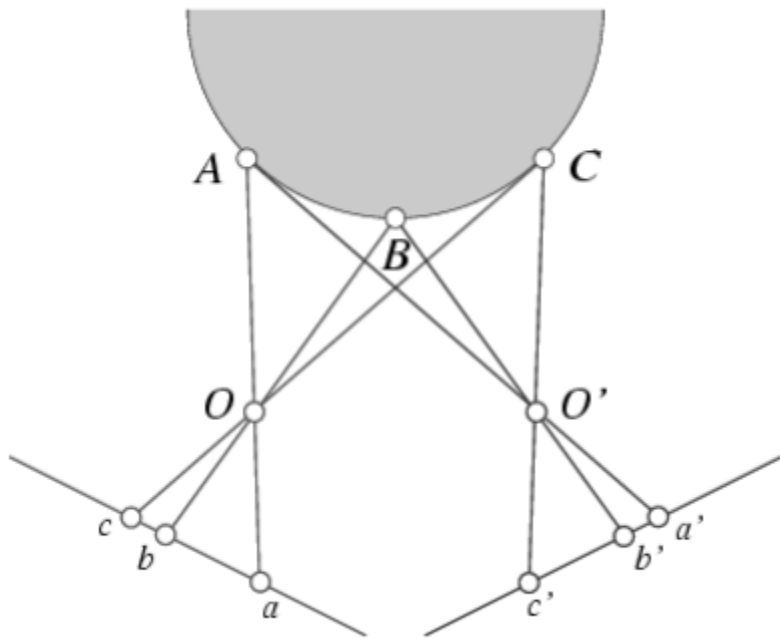
- Edges are found by repeatedly smoothing the image and detecting the zero crossings of the second derivative (Laplacian).

Matching zero-crossings at multiple scales



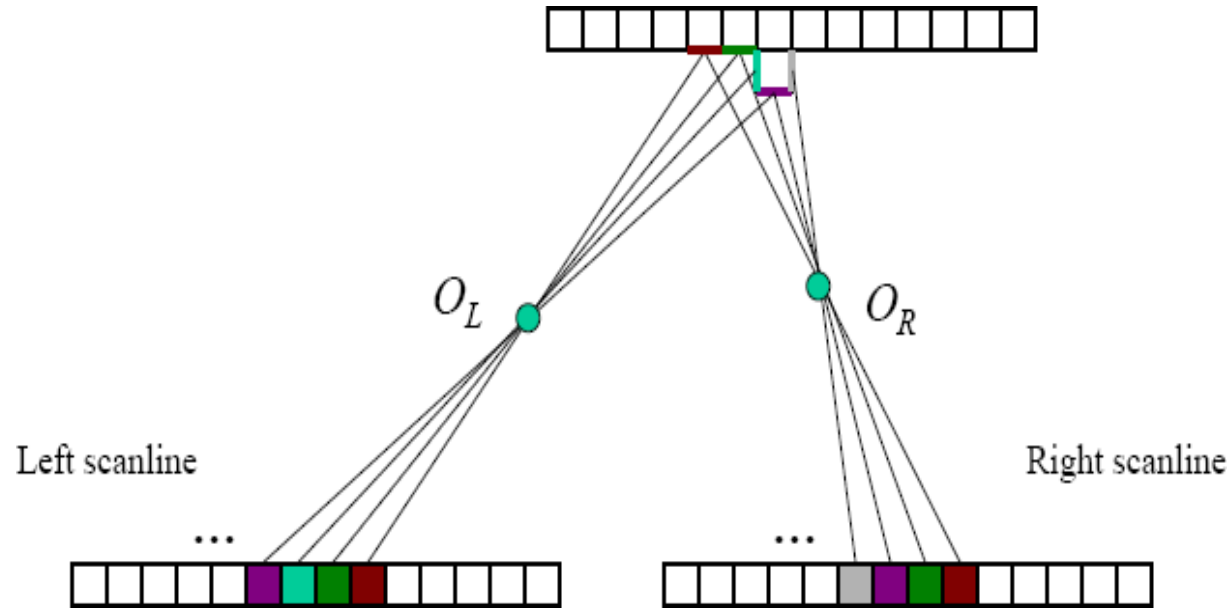
- Matches at coarse scales are used to offset the search for matches at fine scales (equivalent to eye movements).

Correspondence

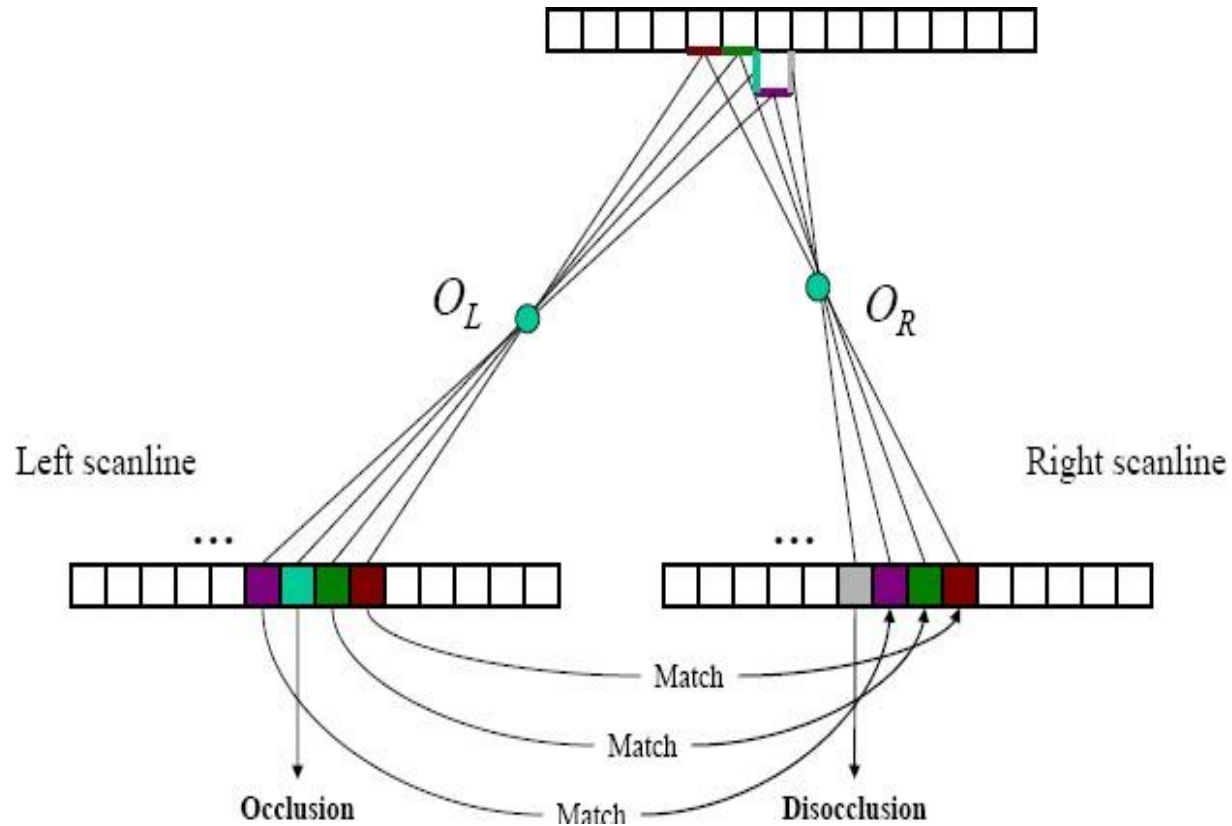


- In general the points are in the same order on both epipolar lines.
- But it is not always the case..

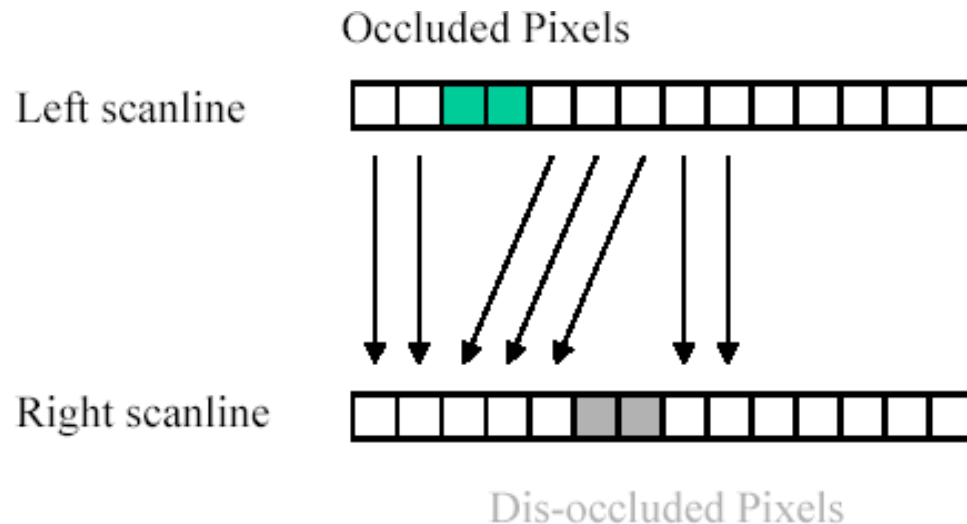
Correspondence



Correspondence



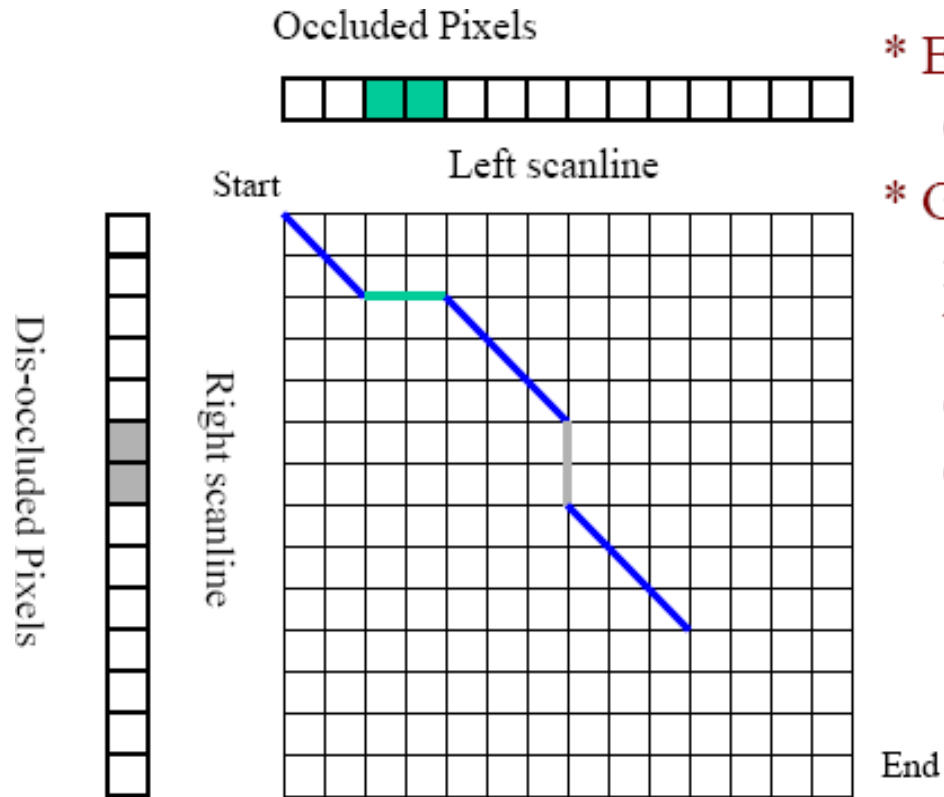
Search over correspondences



Three cases:

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

Stereo Matching with Dynamic Programming



- * Enforces ordering constraint.
- * Given appropriate cost functions, solves for best path (matches, occlusions, disocclusions).

Approaches to Find Correspondences

- **Intensity Correlation-based approaches**

- (+) dense disparity (disparity at each pixel)
- (-) tends to fail in smooth featureless areas
- (-) foreshortening

Solution: warp windows?

- **Edge / feature matching approaches**

- (+) solve the foreshortening problem
- (-) sparse disparity

Solution: interpolate intermediate disparities.

- (-) requires feature detection

- **Dynamic programming**

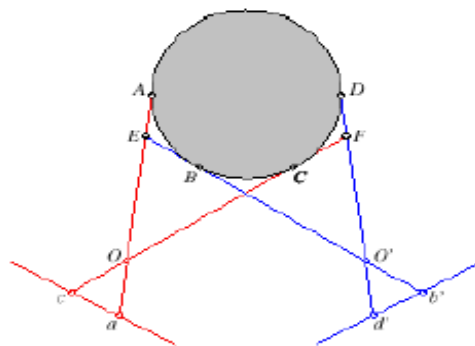
- (+) use both features and intensities
- (-) emphasize ordering constraint

- **Energy minimization / Graph cuts**

- **Probabilistic approaches**

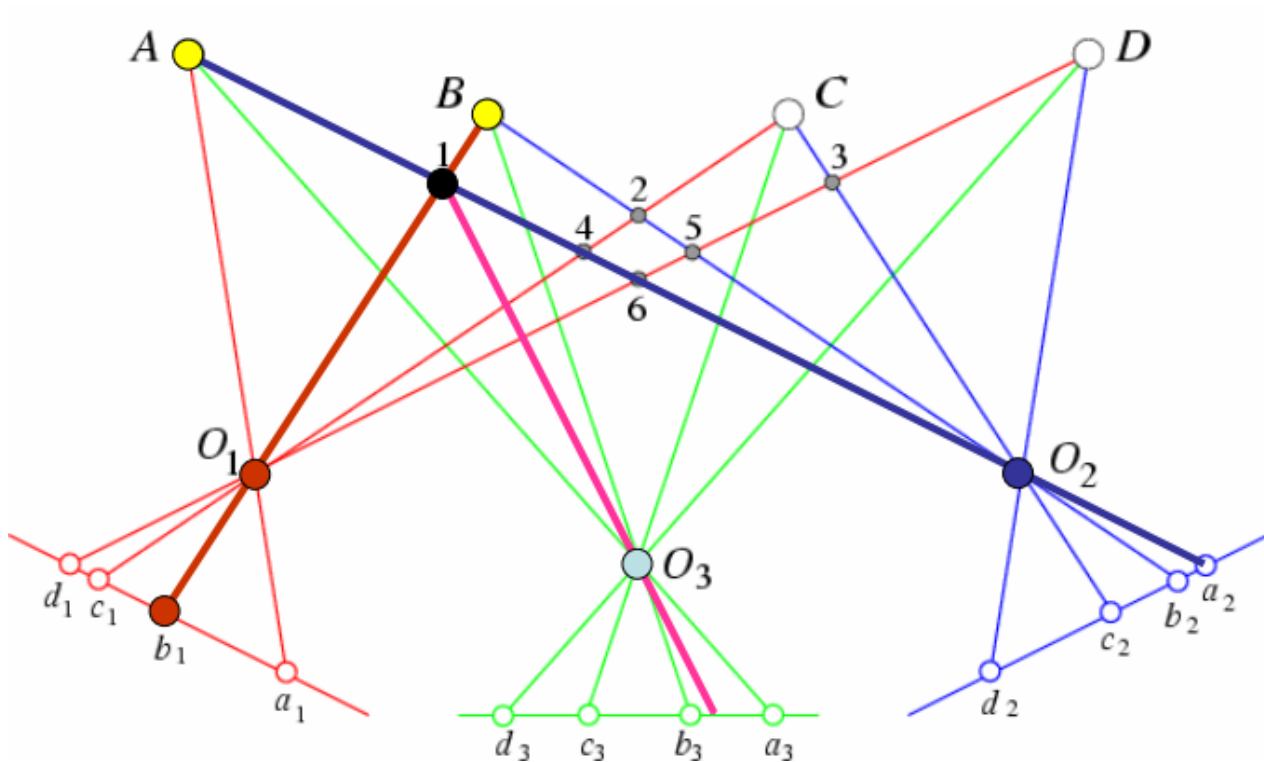
Computing correspondences

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

Reading material: Section 11 of Szeliski