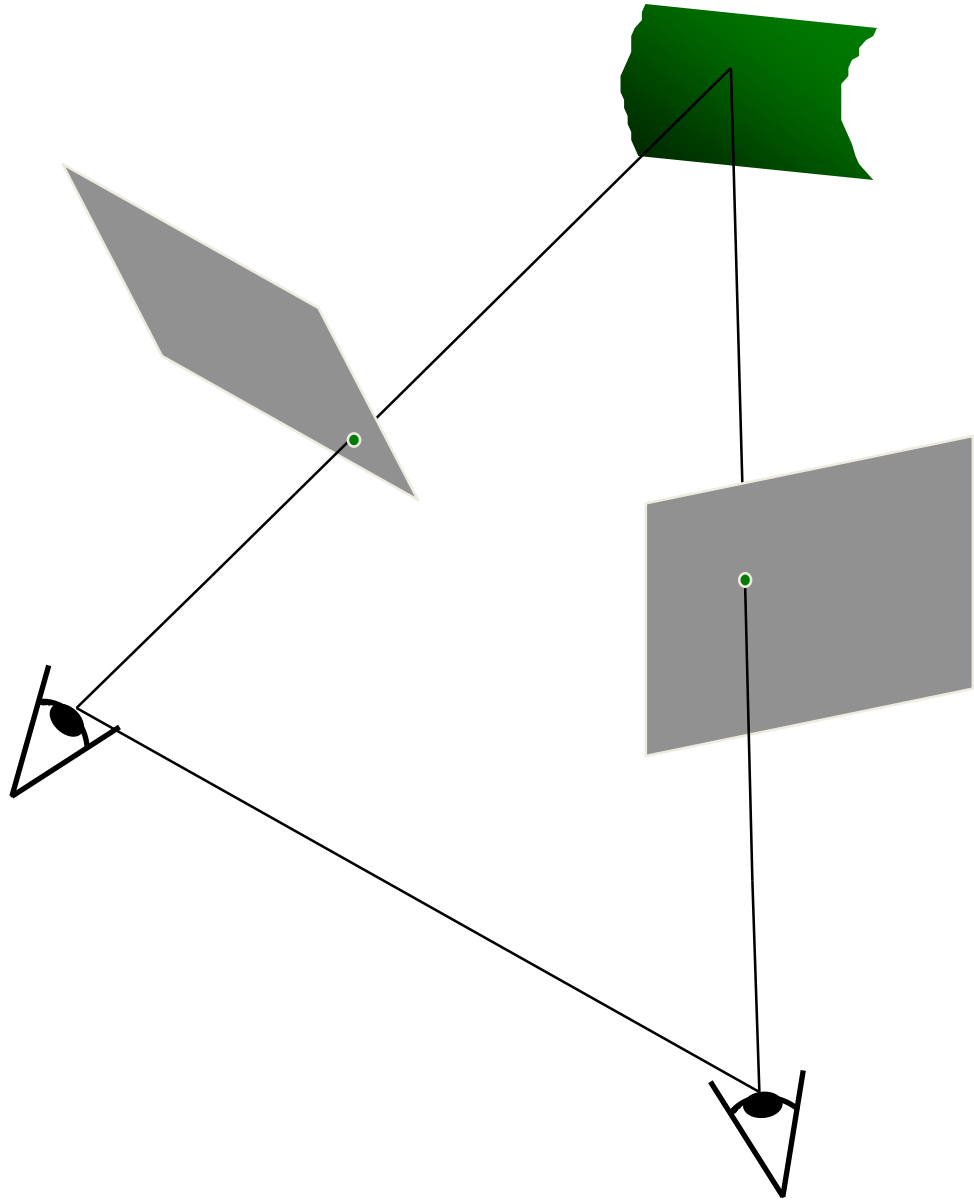# Stereo  Matching

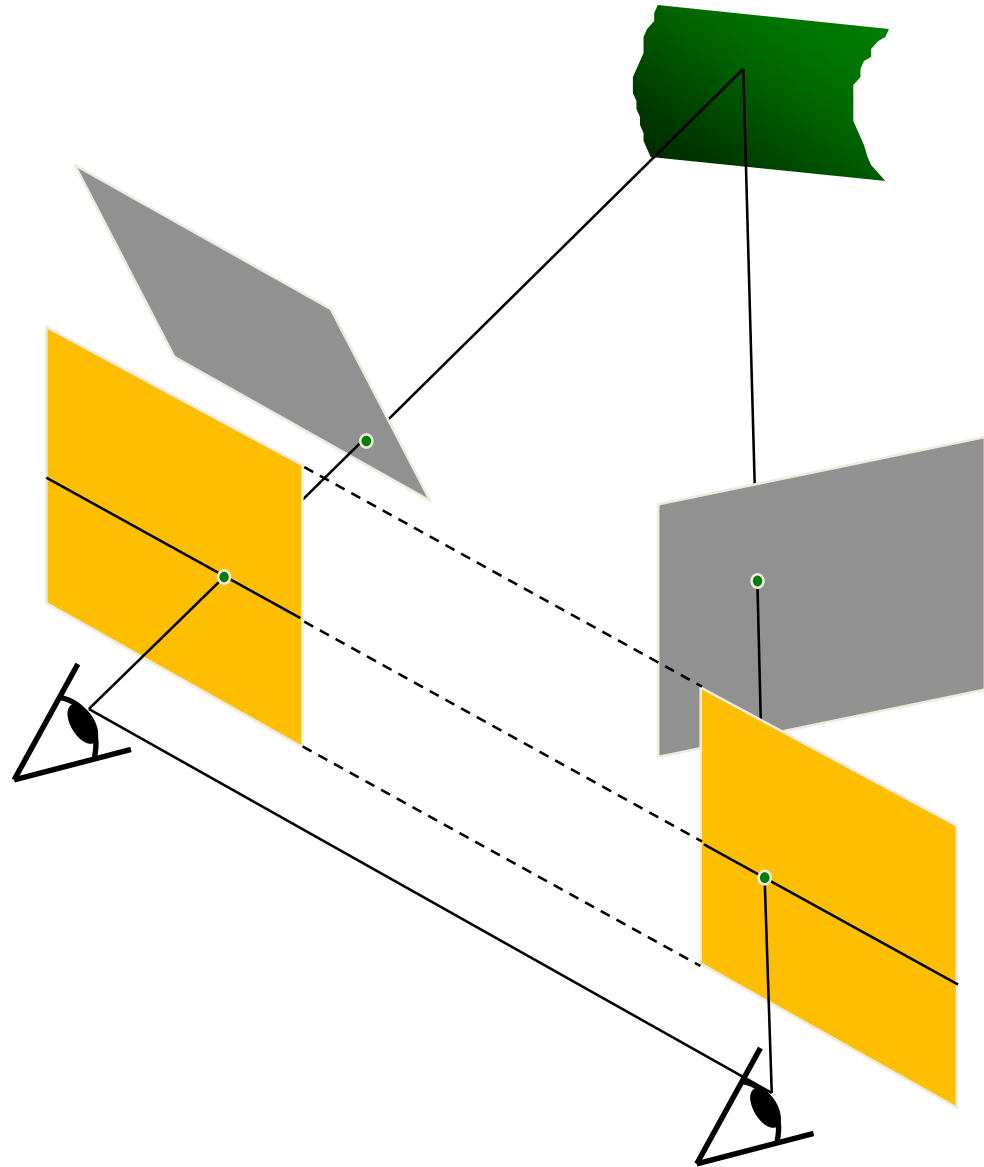**Jianping Fan**
**Department of Computer Science**
**UNC-Charlotte**

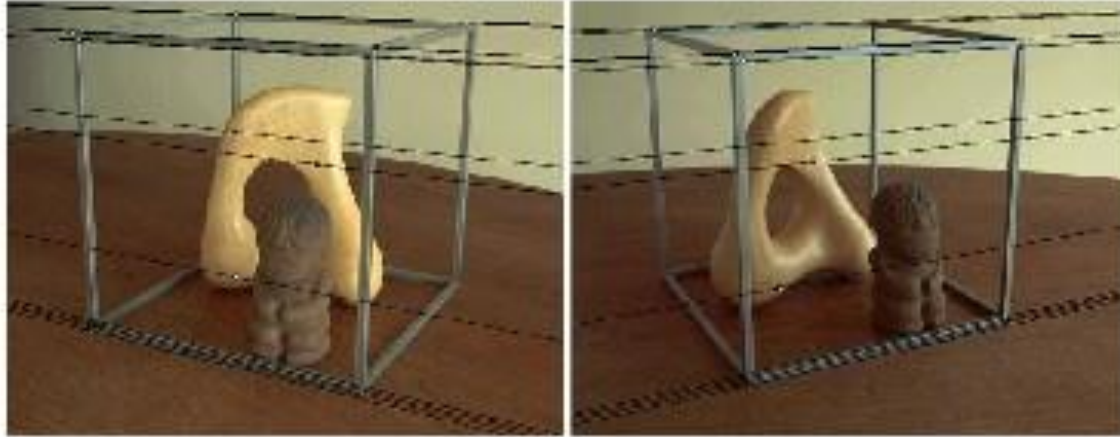# Stereo image rectification

# Stereo image rectification

- Reproject image planes onto a common plane parallel to the line between camera centers

- Pixel motion is horizontal after this transformation

- Two homographies (3x3 transform), one for each input image reprojection

➢ C. Loop and Z. Zhang. Computing Rectifying Homographies for Stereo Vision. IEEE Conf. Computer Vision and Pattern Recognition, 1999.

# Rectification example

# The correspondence problem

- Epipolar geometry constrains our search, but we still have a difficult correspondence problem.

# Fundamental Matrix + Sparse correspondence

# Fundamental Matrix + Dense correspondence

The Visual Turing Test for Scene Reconstruction
Supplementary Video

Qi Shan[+]    Riley Adams[+]    Brian Curless[+]
Yasutaka Furukawa[*]         Steve Seitz[+*]

[+]University of Washington    [*]Google

3DV 2013

# SIFT + Fundamental Matrix + RANSAC

Despite their scale invariance and robustness to appearance changes, SIFT features are *local* and do not contain any global information about the image or about the location of other features in the image. Thus feature matching based on SIFT features is still prone to errors. However, since we assume that we are dealing with rigid scenes, there are strong geometric constraints on the locations of the matching features and these constraints can be used to clean up the matches. In particular, when a rigid scene is imaged by two pinhole cameras, there exists a $3 \times 3$ matrix $F$, the *Fundamental matrix*, such that corresponding points $x_{ij}$ and $x_{ik}$ (represented in homogeneous coordinates) in two images $j$ and $k$ satisfy[10]:

$$x_{ij}^{\top} F x_{ij} = 0. \tag{3}$$

A common way to impose this constraint is to use a greedy randomized algorithm to generate suitably chosen random estimates of $F$ and choose the one that has the largest support among the matches, i.e., the one for which the most matches satisfy (3). This algorithm is called Random Sample Consensus (RANSAC)[6] and is used in many computer vision problems.
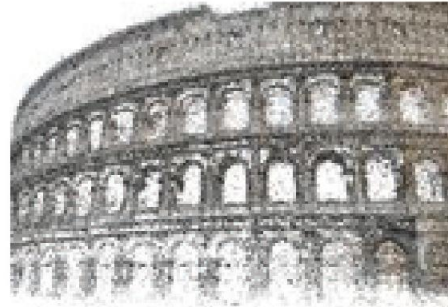
# Sparse to Dense Correspondence



|  | Input images | SfM points | MVS points |
|---|---|---|---|
| Colosseum | | | |
| St. Peter's | | | |

# Structure from motion (or SLAM)

- Given a set of corresponding points in two or more images, compute the camera parameters and the 3D point coordinates



Camera 1
$R_1, t_1$

Camera 2
$R_2, t_2$

Camera 3
$R_3, t_3$

# Structure from motion ambiguity

- If we scale the entire scene by some factor $k$ and, at the same time, scale the camera matrices by the factor of $1/k$, the projections of the scene points in the image remain exactly the same:

$$\mathbf{x} = \mathbf{P}\mathbf{X} = \left( \frac{1}{k}\mathbf{P} \right)(k\,\mathbf{X})$$

It is impossible to recover the absolute scale of the scene!
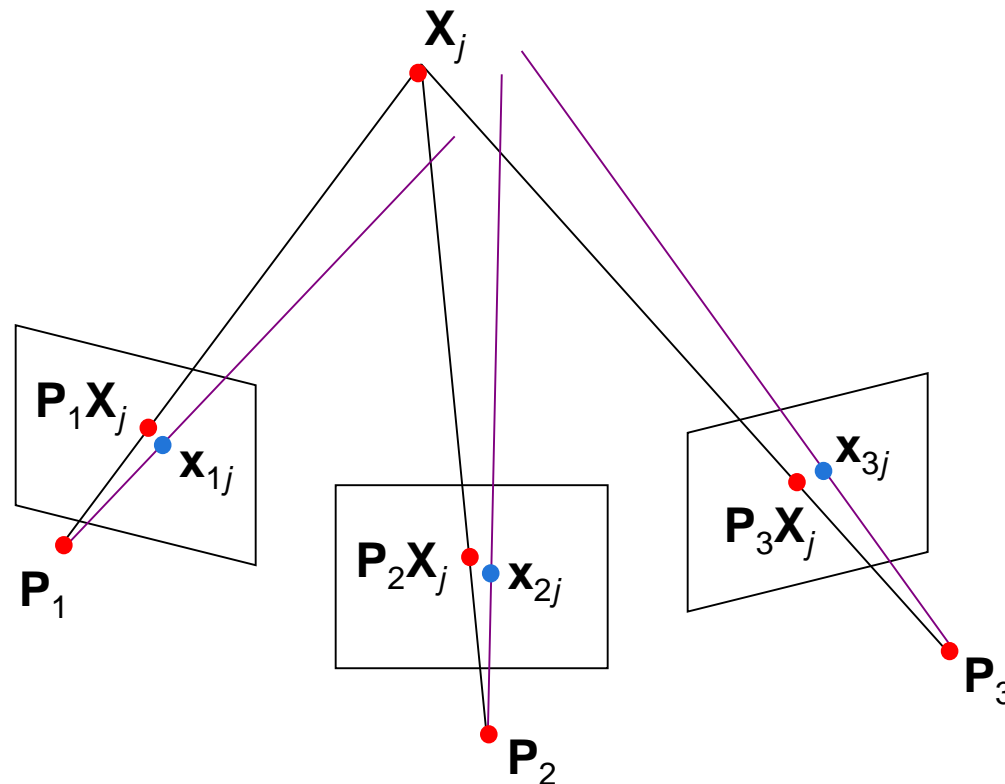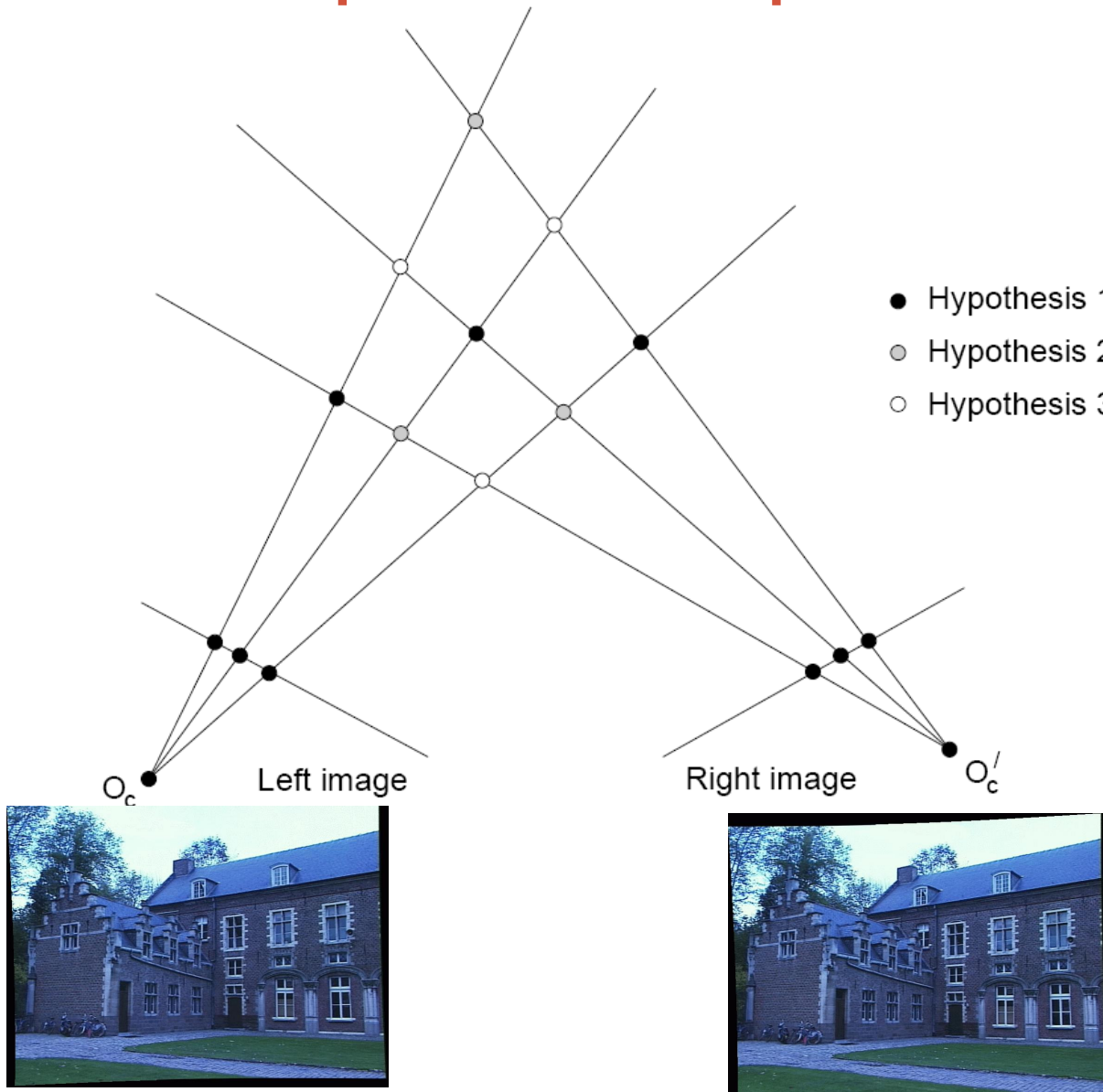
# How do we know the scale of image content?

# Bundle adjustment

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

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

# Correspondence problem



- ● Hypothesis 1
- ◉ Hypothesis 2
- ○ Hypothesis 3

Left image

Right image

$O_c$    $O_c'$

Multiple match hypotheses satisfy epipolar constraint, but which is correct?
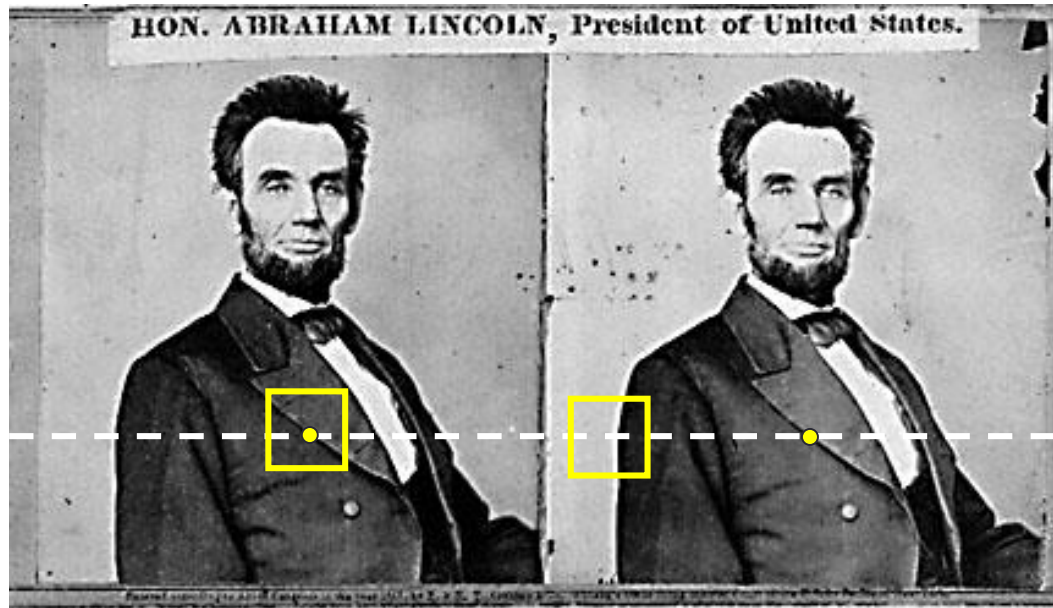
Figure from Gee & Cipolla 1999

# Correspondence problem

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

- To find matches in the image pair, we will assume
  - Most scene points visible from both views
  - Image regions for the matches are similar in appearance
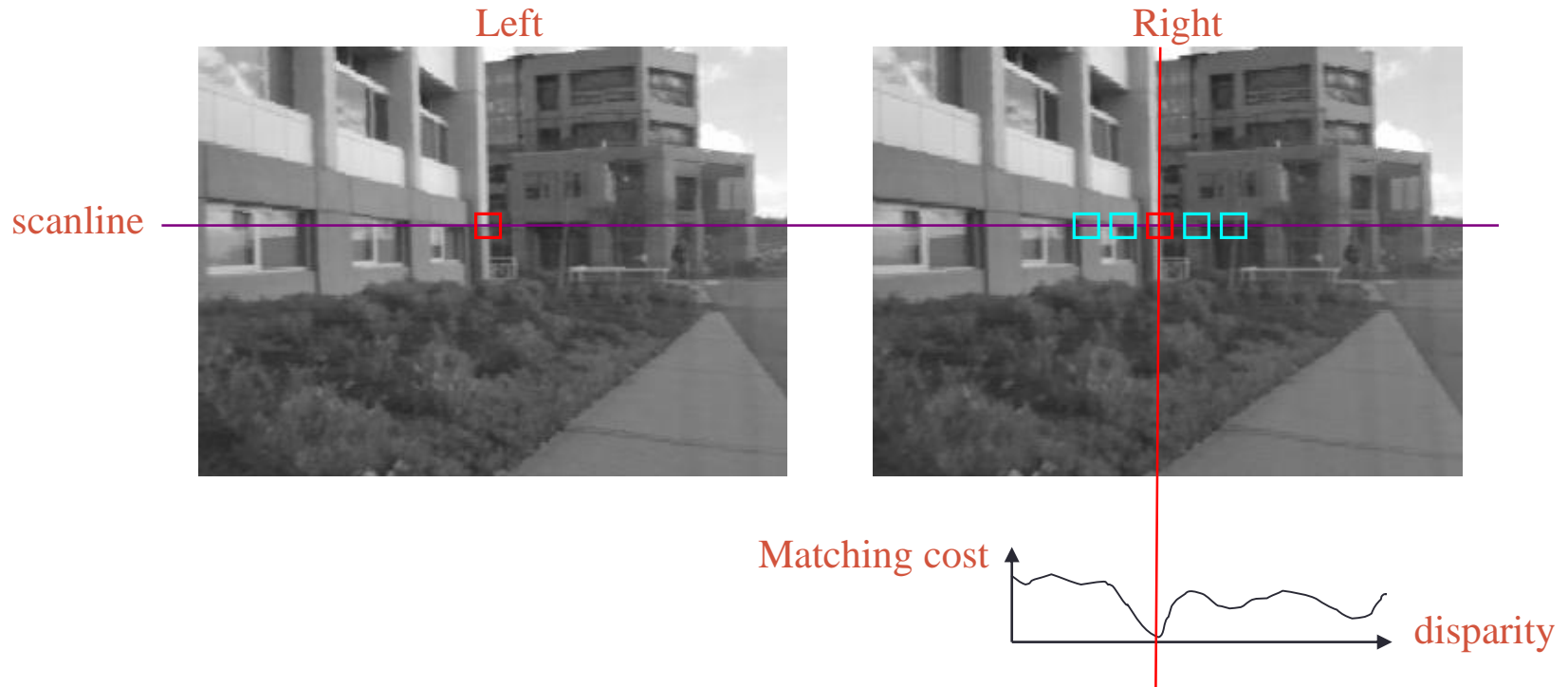
# Dense correspondence search



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, normalized correlation)

# Correspondence search with similarity constraint



- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

# Correspondence search with similarity constraint



Left

Right

scanline

SSD

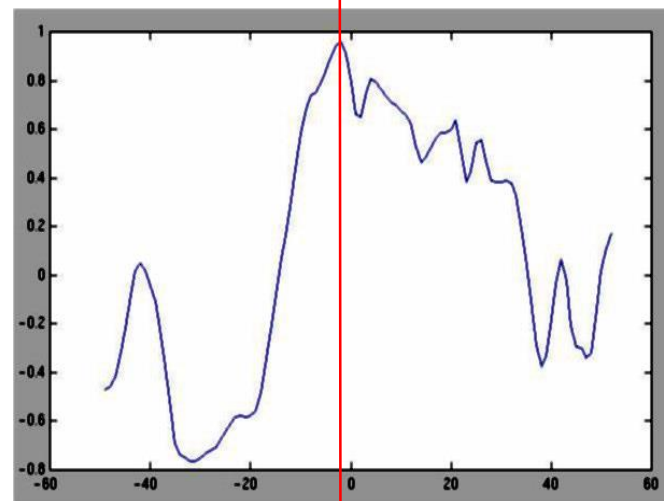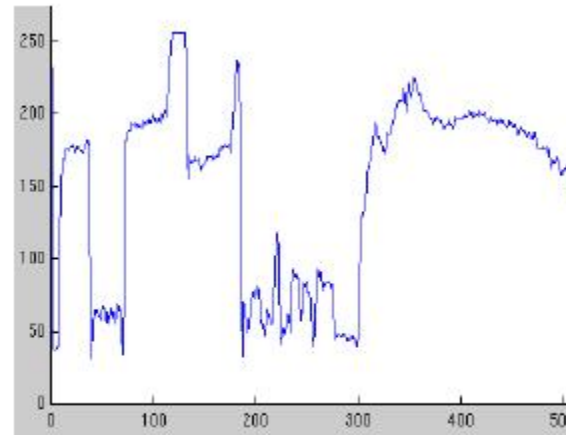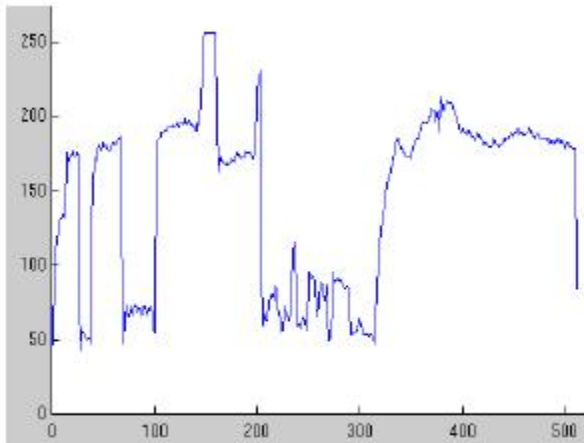# Correspondence search with similarity constraint

Left

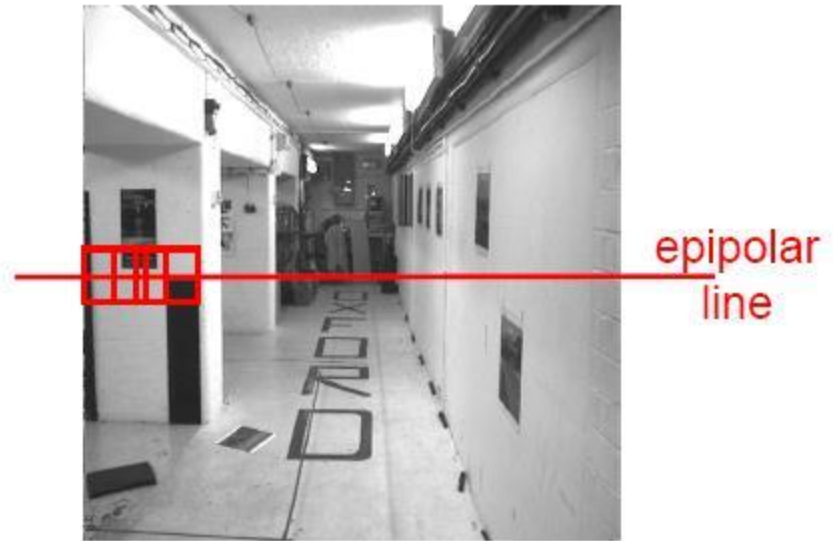Right

scanline



Norm. corr

# Correspondence problem
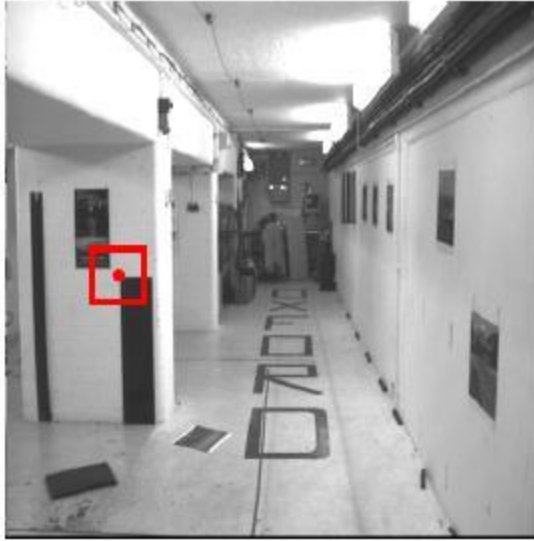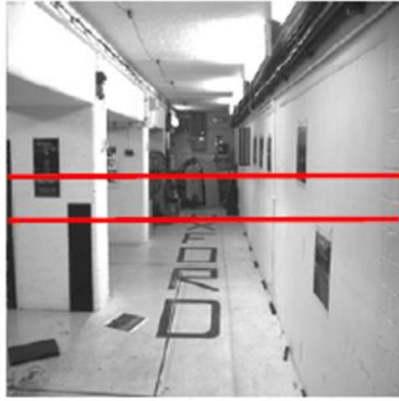


Intensity profiles

- Clear correspondence between intensities, but also noise and ambiguity

# Correspondence problem



epipolar line

Neighborhoods of corresponding points are similar in intensity patterns.

# Correlation-based window matching



left image band (x)
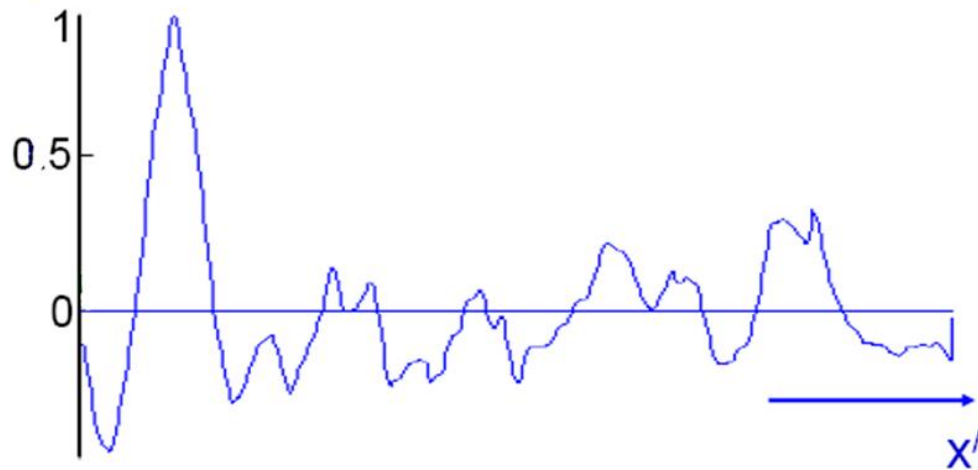
# Correlation-based window matching



left image band (x)

right image band (x')

# Correlation-based window matching



left image band (x)

right image band (x′)

cross correlation

disparity = x′ - x

# Correlation-based window matching



target region

left image band (x)

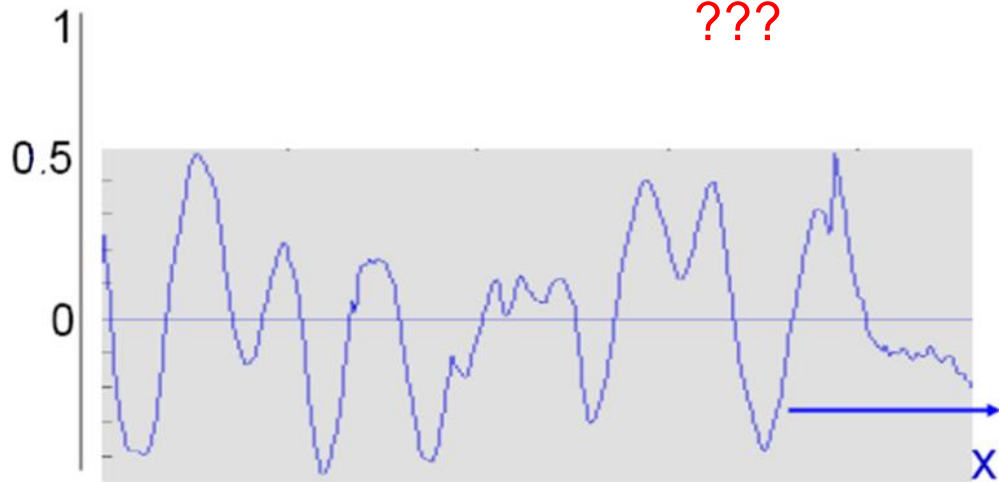right image band (x')

# Correlation-based window matching

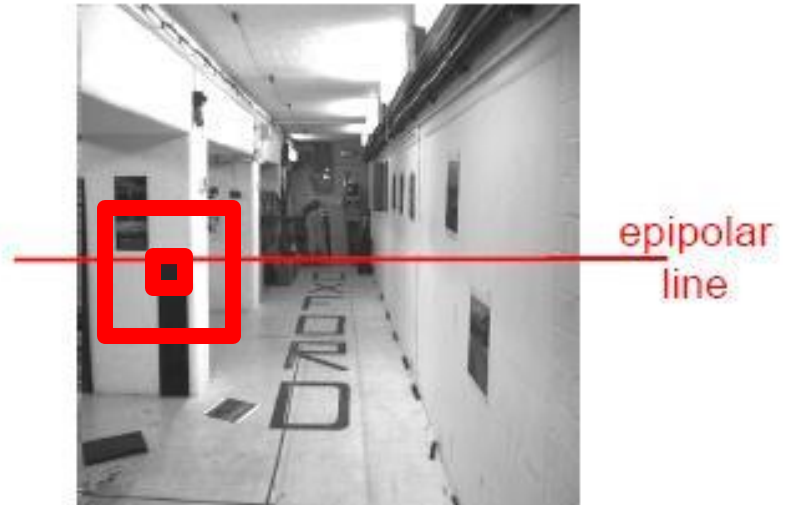

target region

left image band (x)

right image band (x')

???

cross correlation

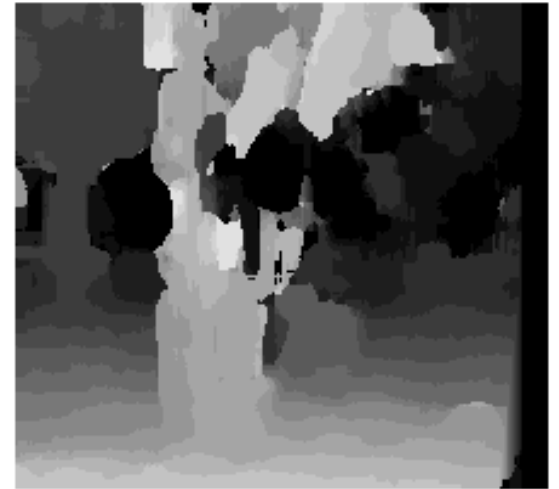Textureless regions are non-distinct; high ambiguity for matches.

# Effect of window size



epipolar
line

Source: Andrew Zisserman

# Effect of window size



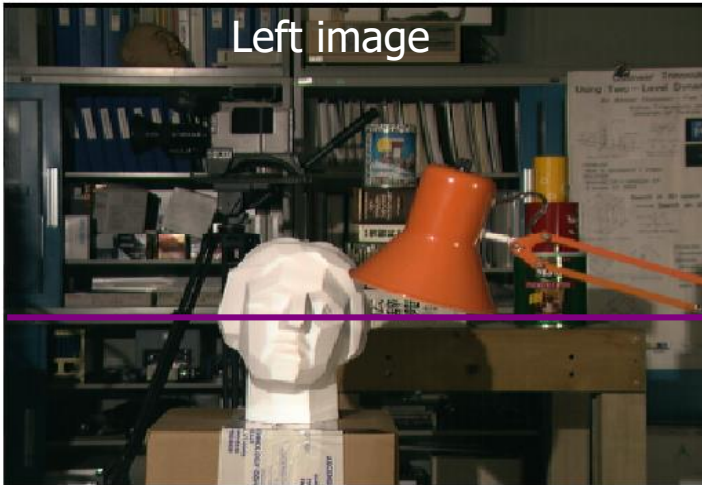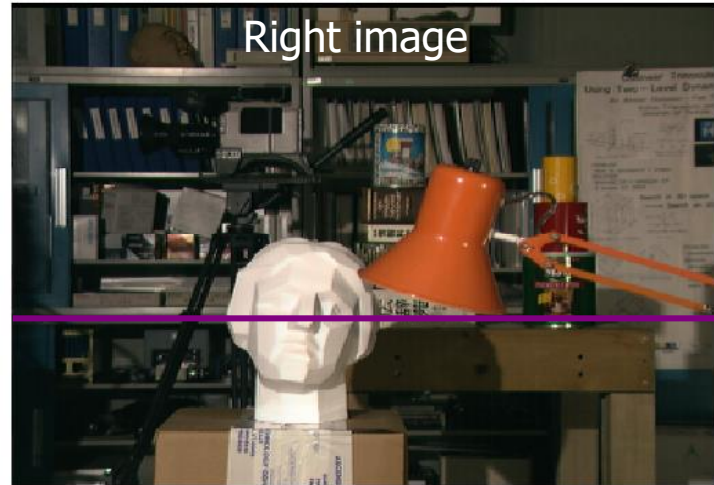W = 3               W = 20

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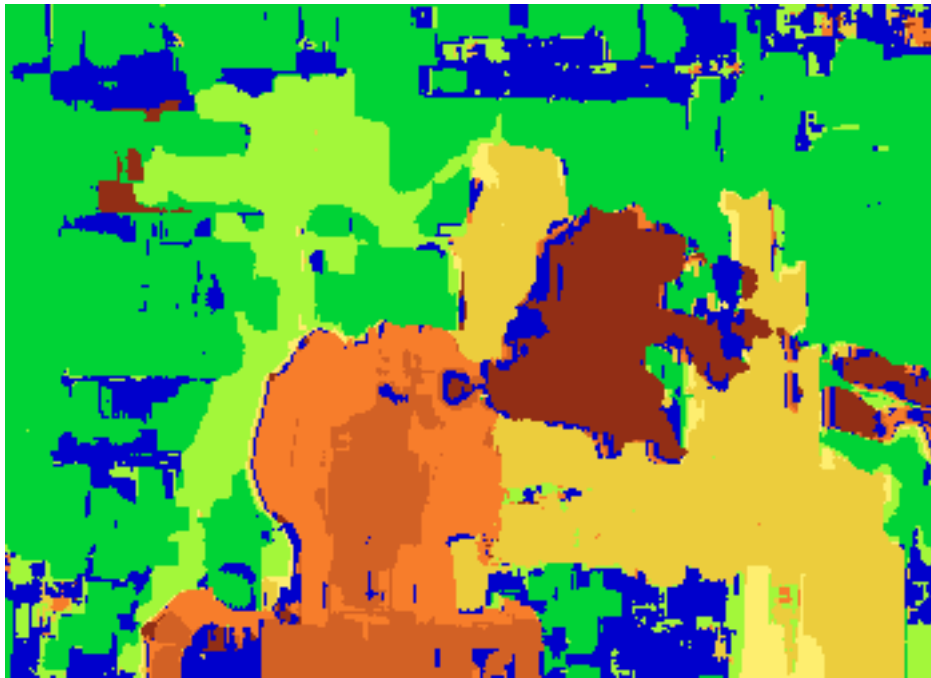
Left image
Right image

# Results with window search
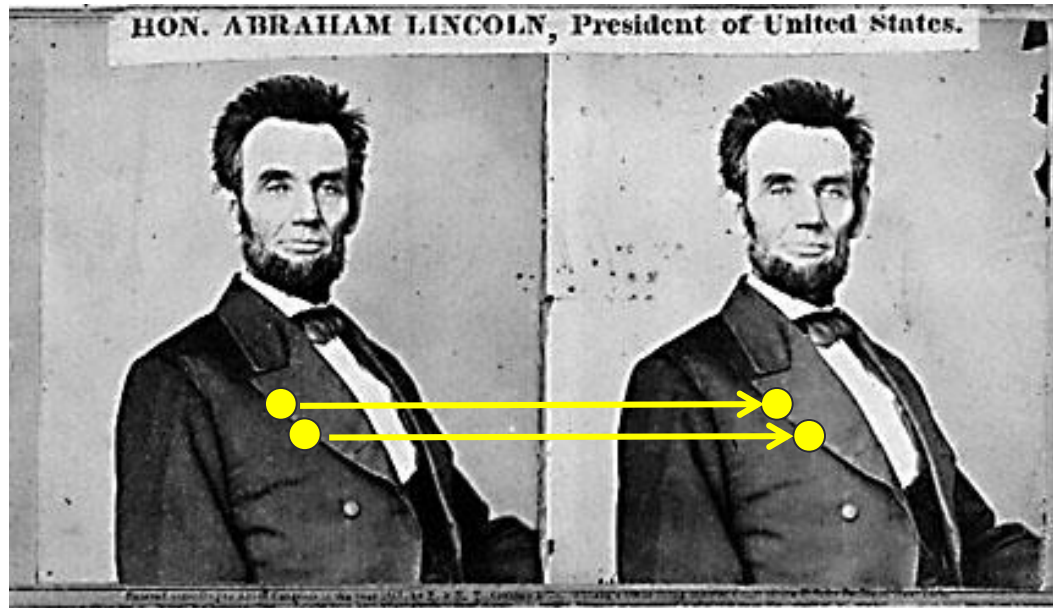


Window-based matching
(best window size)

Ground truth

# Better solutions

- Beyond individual correspondences to estimate disparities:

- Optimize correspondence assignments jointly
  - Scanline at a time (DP)
  - Full 2D grid (graph cuts)

# Stereo as energy minimization



- What defines a good stereo correspondence?
    1. Match quality
        - Want each pixel to find a good match in the other image
    2. Smoothness
        - If two pixels are adjacent, they should (usually) move about the same amount

# Stereo matching as energy minimization



$I_1$  $W_1(i)$

$I_2$  $W_2(i+D(i))$

$D$  $D(i)$

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

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

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

- 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

# Better results…



Graph cut method                    Ground truth

Boykov et al., Fast Approximate Energy Minimization via Graph Cuts,
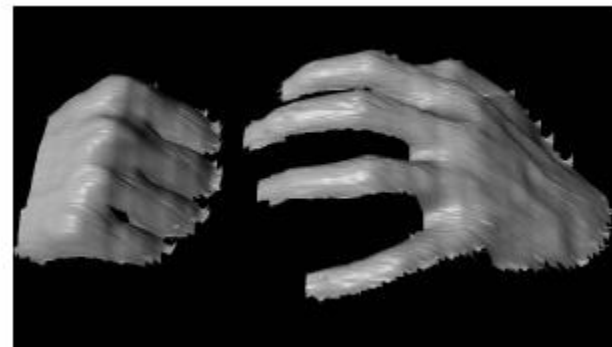    International Conference on Computer Vision, September 1999.

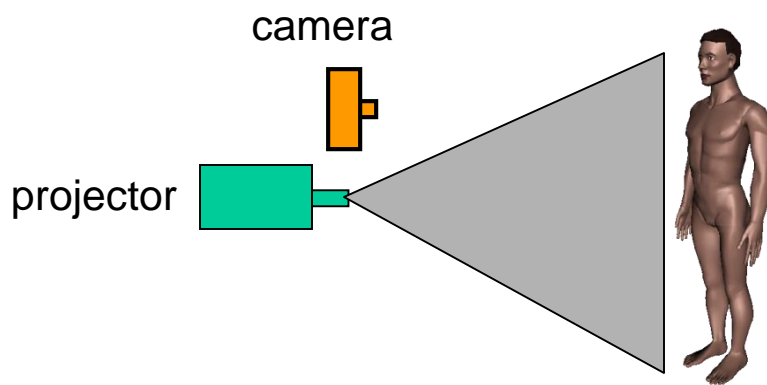For the latest and greatest:  http://www.middlebury.edu/stereo/

# Challenges

- Low-contrast ; textureless image regions
- Occlusions
- Violations of brightness constancy (e.g., specular reflections)
- Really large baselines (foreshortening and appearance change)
- Camera calibration errors

# Active stereo with structured light



- Project "structured" light patterns onto the object
  - Simplifies the correspondence problem
  - Allows us to use only one camera



camera

projector

L. Zhang, B. Curless, and S. M. Seitz. Rapid Shape Acquisition Using Color Structured Light and Multi-pass Dynamic Programming. *3DPVT* 2002

# Kinect: Structured infrared light

# iPhone X