



SAIR

Spatial AI & Robotics Lab

# CSE 473/573-A

## L16: STEREO MATCHING

Chen Wang

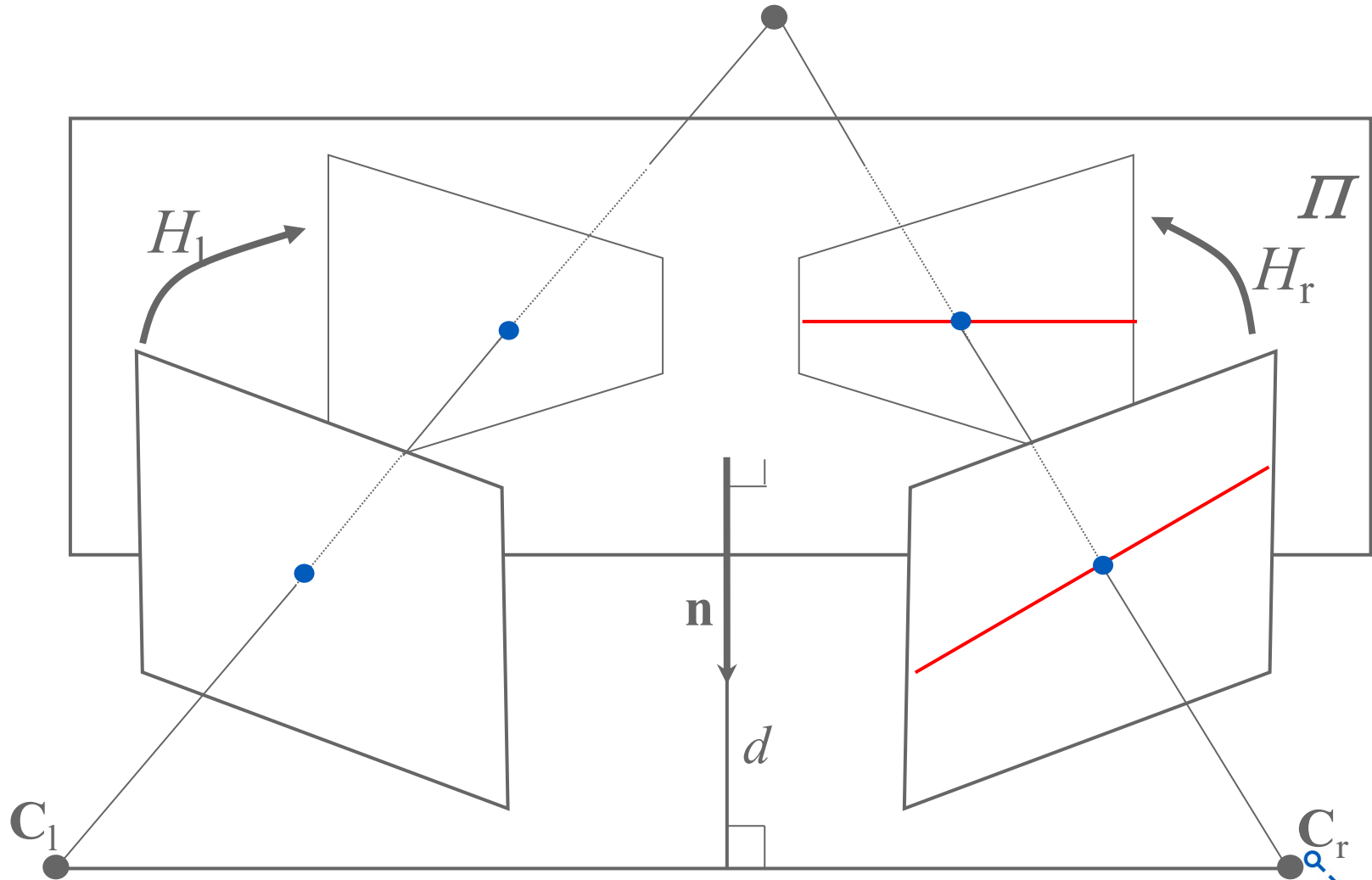
Spatial AI & Robotics Lab

Department of Computer Science and Engineering

 **University at Buffalo** The State University of New York

Many Slides from Lana Lazebnik

# Assume Rectified Stereo Images





SAIR

Spatial AI & Robotics Lab

# STEREO VISION

## Essential / Fundamental Matrix

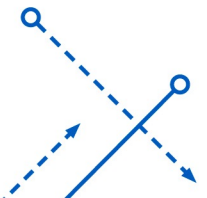
# Coordinates in 2-D (Recap)

- Cartesian / homogeneous coordinates of point  $p$

$$p = \begin{bmatrix} x \\ y \end{bmatrix} \quad \begin{bmatrix} u \\ v \\ w \end{bmatrix} \quad \begin{aligned} x &= u / w \\ y &= v / w \end{aligned}$$

- Homogeneous coordinate vector are equivalent if they are proportional to each other

$$\begin{bmatrix} u \\ v \\ w \end{bmatrix} \equiv \begin{bmatrix} u' \\ v' \\ w' \end{bmatrix} \Leftrightarrow \begin{bmatrix} u \\ v \\ w \end{bmatrix} \equiv \lambda \begin{bmatrix} u' \\ v' \\ w' \end{bmatrix} \quad \lambda \neq 0$$



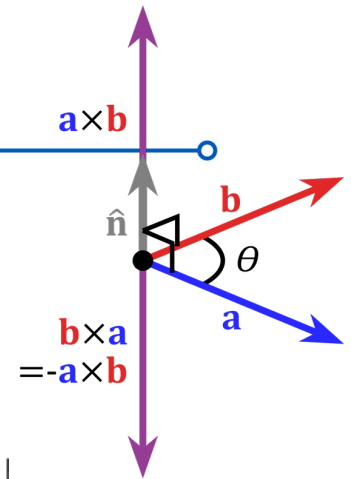
# Cross product (Recap)

- Cross product of two 3D vector

$$\mathbf{a} \times \mathbf{b} = \hat{n} |\mathbf{a}| |\mathbf{b}| \sin \theta$$

$$\begin{aligned} \mathbf{a} \times \mathbf{b} &= (a_1 \mathbf{i} + a_2 \mathbf{j} + a_3 \mathbf{k}) \times (b_1 \mathbf{i} + b_2 \mathbf{j} + b_3 \mathbf{k}) \\ &= a_1 b_1 (\mathbf{i} \times \mathbf{i}) + a_1 b_2 (\mathbf{i} \times \mathbf{j}) + a_1 b_3 (\mathbf{i} \times \mathbf{k}) + \\ &\quad a_2 b_1 (\mathbf{j} \times \mathbf{i}) + a_2 b_2 (\mathbf{j} \times \mathbf{j}) + a_2 b_3 (\mathbf{j} \times \mathbf{k}) + \\ &\quad a_3 b_1 (\mathbf{k} \times \mathbf{i}) + a_3 b_2 (\mathbf{k} \times \mathbf{j}) + a_3 b_3 (\mathbf{k} \times \mathbf{k}) \end{aligned}$$

$$\mathbf{a} \times \mathbf{b} = \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{vmatrix}$$

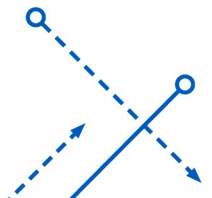


- Skew-symmetric Matrix

$$\mathbf{a} \times \mathbf{b} = [\mathbf{a}]_{\times} \mathbf{b} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad [\mathbf{a}]_{\times} = \begin{bmatrix} 0 & -a_3 & a_2 \\ a_3 & 0 & -a_1 \\ -a_2 & a_1 & 0 \end{bmatrix}$$

- Properties:

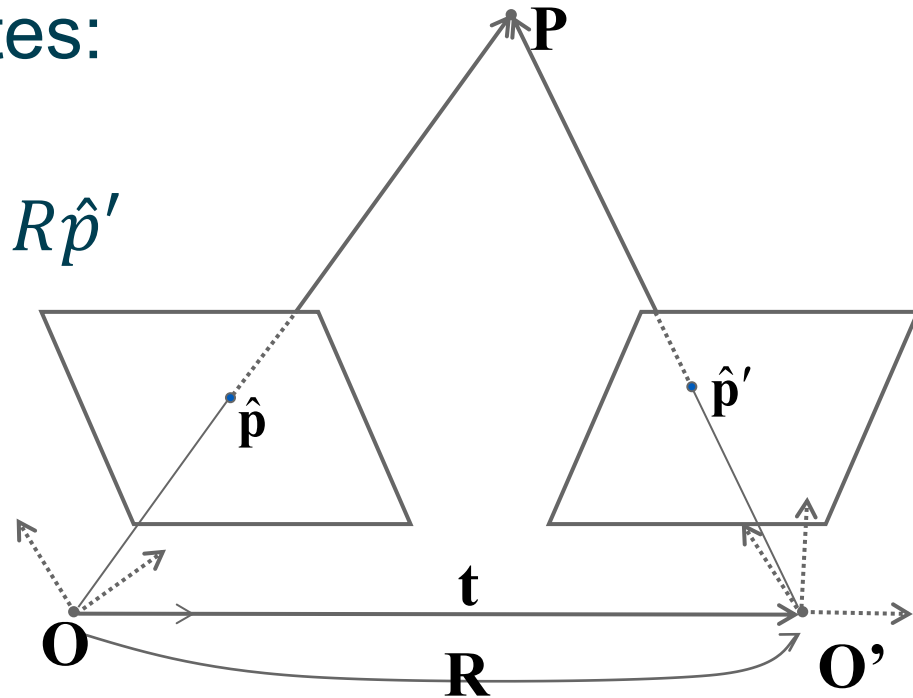
$$\mathbf{a}^T (\mathbf{a} \times \mathbf{b}) = \mathbf{b}^T (\mathbf{a} \times \mathbf{b}) = 0$$



# Essential Matrix

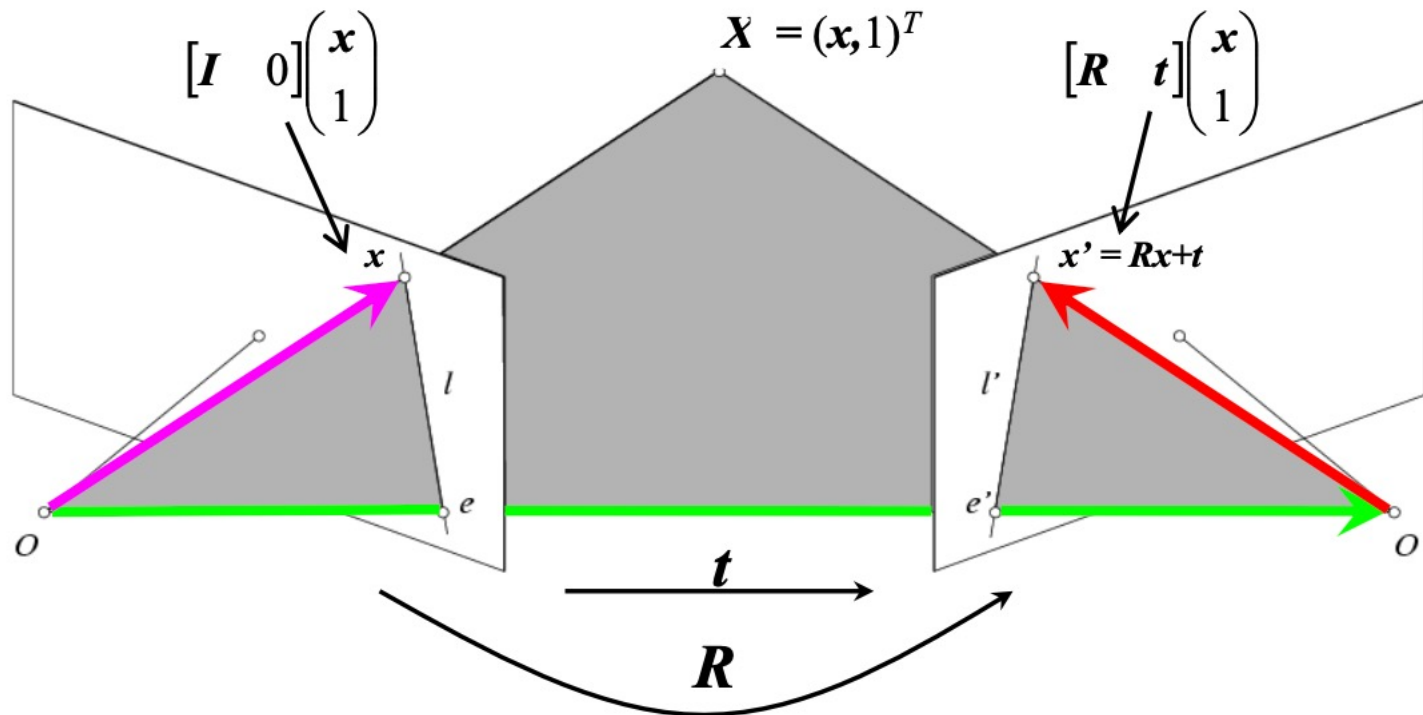
- Due to cross product properties
  - $\overrightarrow{OP} \cdot (\overrightarrow{OO'} \times \overrightarrow{O'P}) = 0$
- In homogeneous coordinates:
  - transform  $O$  to align  $O'$
- Then direction of  $\hat{p}'$  in  $O$  is  $R\hat{p}'$ 
  - $\hat{p}^T (t \times R\hat{p}') = 0$
  - $\hat{p}^T ([t]_{\times} R) \hat{p}' = 0$

$$\hat{p}^T \mathbf{E} \hat{p}' = 0$$
$$\mathbf{E} = [t]_{\times} \mathbf{R}$$



Essential Matrix

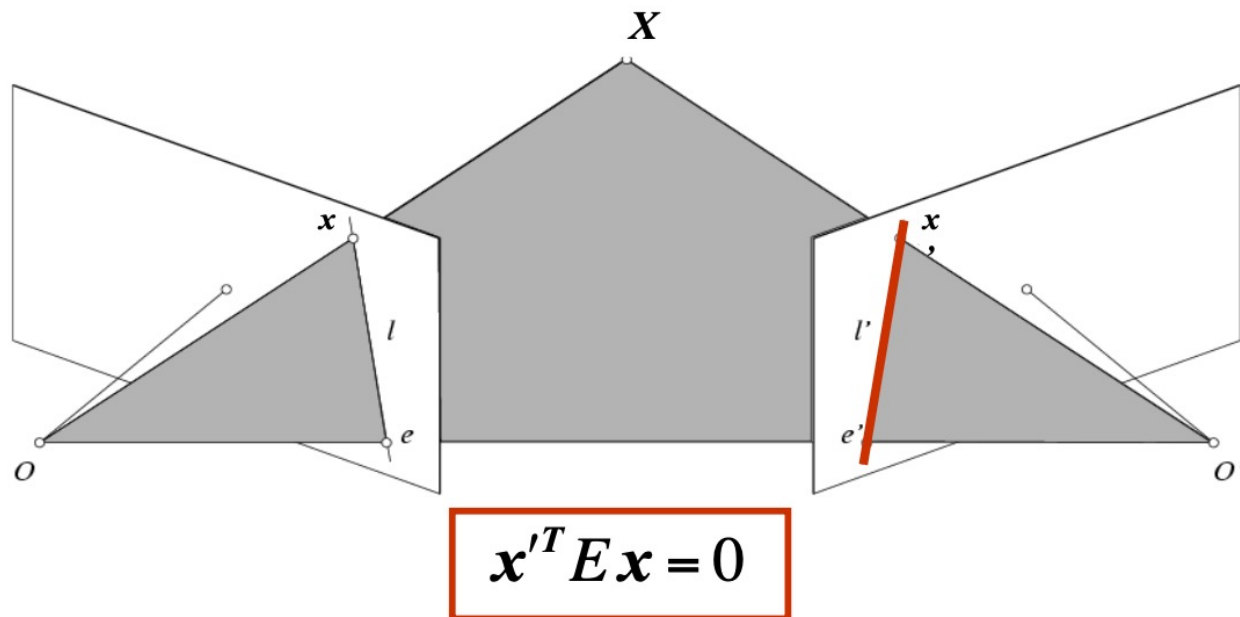
# Epipolar constraint: Calibrated case



$$x' \cdot [t \times (Rx)] = 0 \quad \Rightarrow \quad x'^T [t_{\times}] Rx = 0 \quad \Rightarrow \quad x'^T E x = 0$$

**Essential Matrix**  
(Longuet-Higgins, 1981)

# Epipolar constraint: Calibrated case



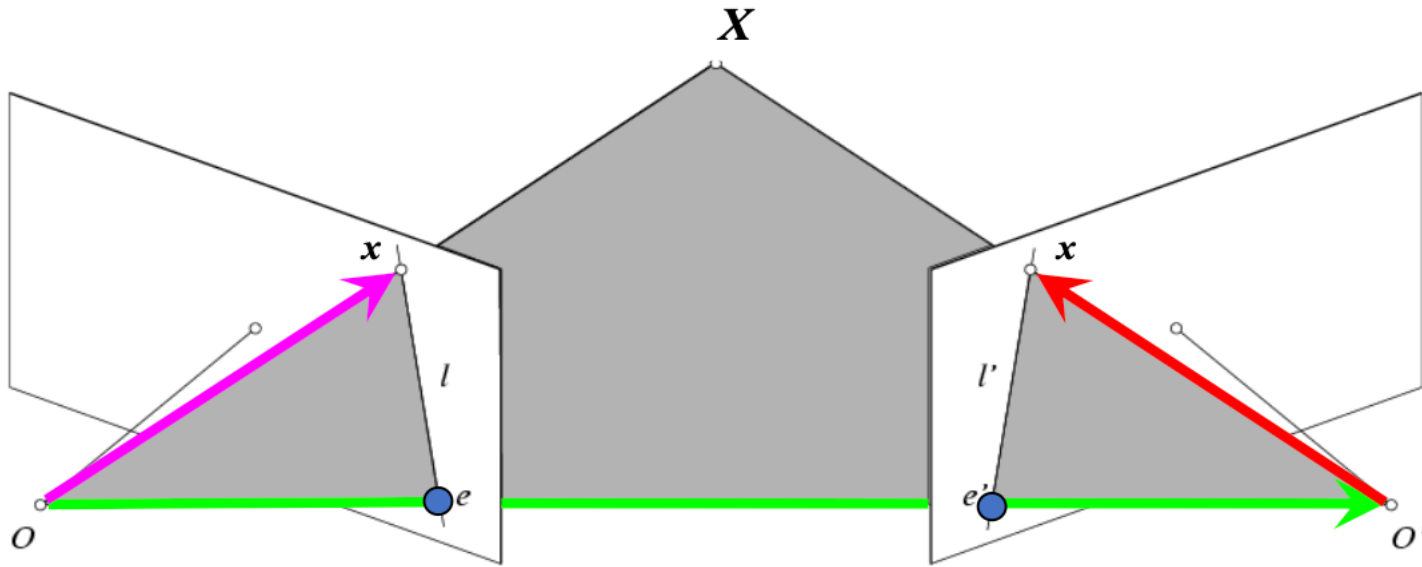
- $E\mathbf{x}$  is the epipolar line associated with  $\mathbf{x}'$  ( $l' = E\mathbf{x}$ )
- $E^T \mathbf{x}'$  is the epipolar line associated with  $\mathbf{x}$  ( $l = E^T \mathbf{x}'$ )
- $E\mathbf{e} = 0$  and  $E^T \mathbf{e}' = 0$
- $E$  is **singular** (rank two)
- $E = [t_{\times}]R$  has five degrees of freedom

• Recall: a line is given by  $ax + by + c = 0$  or

$$\mathbf{l}^T \mathbf{x} = 0 \quad \text{where} \quad \mathbf{l} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



# Epipolar constraint: Uncalibrated case

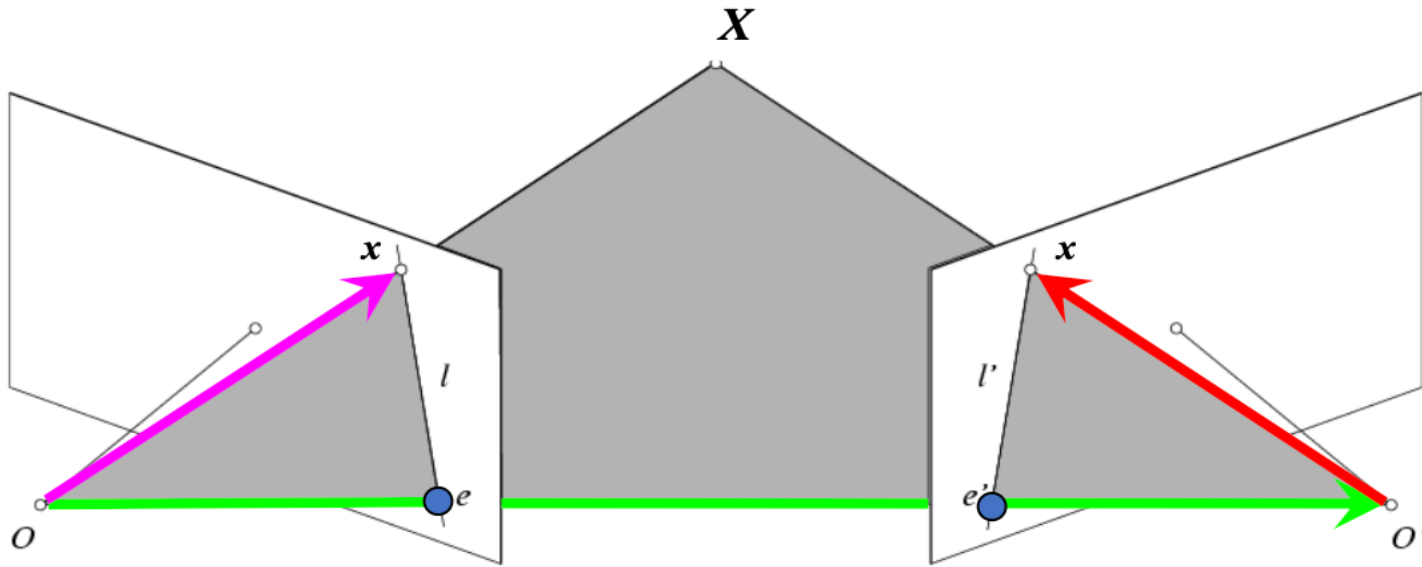


- The calibration matrices  $K$  and  $K'$  are unknown.
- We can write the epipolar constraint in terms of unknown normalized coordinates:

$$\hat{x}'^T E \hat{x} = 0$$

$$\hat{x} = K^{-1} x, \quad \hat{x}' = K'^{-1} x'$$

# Epipolar constraint: Uncalibrated case



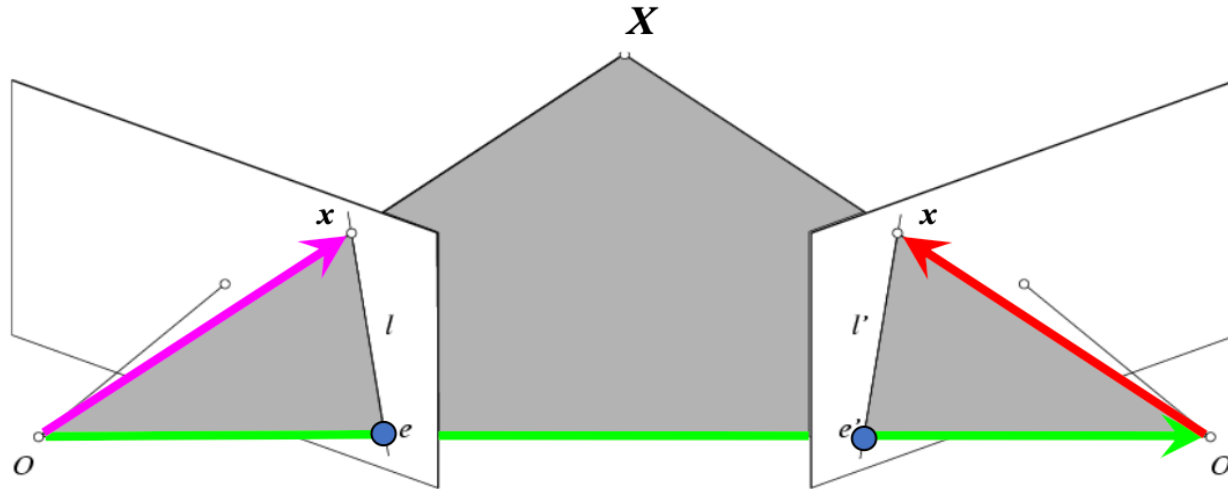
$$\hat{x}'^T E \hat{x} = 0 \quad \Rightarrow \quad x'^T F x = 0 \quad \text{with} \quad F = K'^{-T} E K^{-1}$$

$$\hat{x} = K^{-1} x$$

$$\hat{x}' = K'^{-1} x'$$

**Fundamental Matrix**  
(Faugeras and Luong, 1992)

# Epipolar constraint: Uncalibrated case



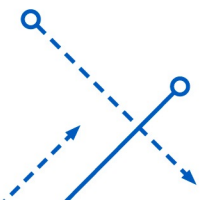
$$\hat{x}'^T E \hat{x} = 0 \quad \Rightarrow \quad x'^T F x = 0 \quad \text{with} \quad F = K'^{-T} E K^{-1}$$

- $Fx$  is the epipolar line associated with  $x'$  ( $l' = Fx$ )
- $F^T x'$  is the epipolar line associated with  $x$  ( $l = F^T x'$ )
- $Fe = 0$  and  $F^T e' = 0$
- $F$  is singular (rank two)
- $F$  has seven degrees of freedom  $(5E + 2K)$

# Estimating the Fundamental Matrix

---

- 8-point algorithm
  - Least squares solution using SVD on equations from 8 pairs of correspondences.
- 7-point algorithm
  - least squares to solve for null space (two vectors) using SVD and 7 pairs of correspondences.
  - Solve for linear combination of null space vectors that satisfies  $\det(F) = 0$
- Minimize reprojection error
  - Non-linear least squares
- Note: estimation of  $F$  (or  $E$ ) is degenerate for a planar scene.



# 8-point algorithm

- Solve a system of homogeneous linear equations
  - a. Write down the system of equations

$$\mathbf{x}^T F \mathbf{x}' = 0$$

$$uu'f_{11} + uv'f_{12} + uf_{13} + vu'f_{21} + vv'f_{22} + vf_{23} + u'f_{31} + v'f_{32} + f_{33} = 0$$

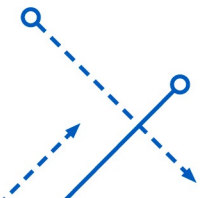
$$A\mathbf{f} = \begin{bmatrix} u_1u_1' & u_1v_1' & u_1 & v_1u_1' & v_1v_1' & v_1 & u_1' & v_1' & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u_nu_n' & u_nv_n' & u_n & v_nu_n' & v_nv_n' & v_n & u_n' & v_n' & 1 \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{bmatrix} = \mathbf{0}$$

- b. Solve  $\mathbf{f}$  from  $A\mathbf{f} = \mathbf{0}$  using SVD.

$$[U, S, V] = \text{svd}(A);$$

$$\mathbf{f} = V(:, \text{end});$$

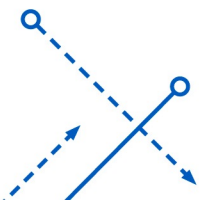
$$F = \text{reshape}(\mathbf{f}, [3 \ 3])';$$



# Q & A

---

- Why do we need 4 points for homography but 7/8 points for fundamental matrix calculation?
  - In the case of fundamental matrix, each point relates to only one constraint, while in homograph, each point is related to two constraints.
- Why can 7 points solve fundamental matrix?
  - In fact, the fundamental matrix only has 7 degrees of freedom. In this case, the rank-2 constraint must be enforced during the computations.

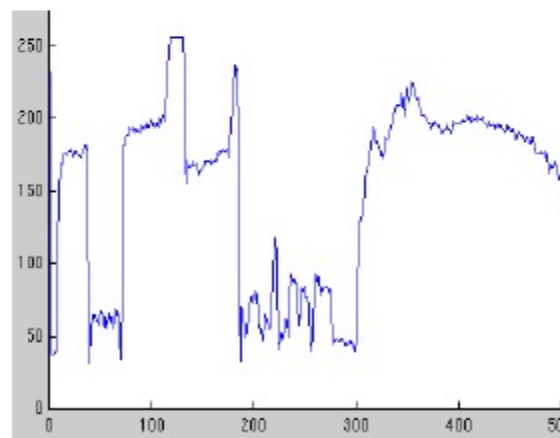
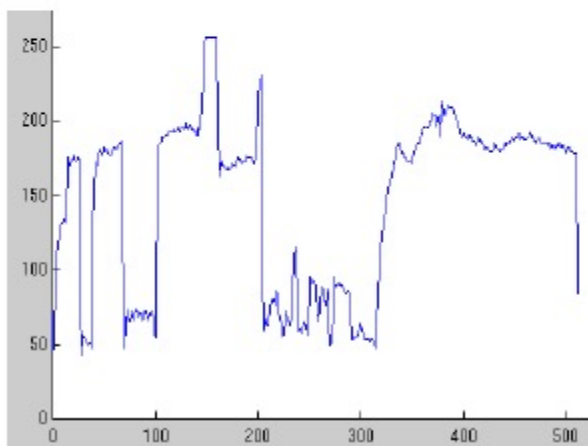


# Correspondence Search

---

- Other “soft” constraints (To cover)
  - 1. Similarity
  - 2. Uniqueness
  - 3. Disparity gradient
  - 4. Ordering
- To find matches in the image pair, we will assume
  - Most scene points visible from both views
  - Matched regions are similar in appearance

# Intensity profiles



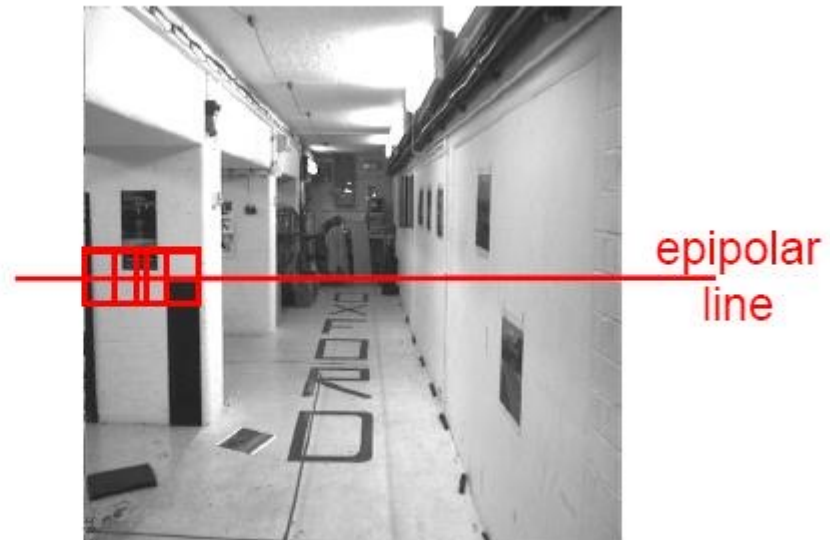
Intensity  
profiles

- Clear correspondence between intensities, but also noise and ambiguity



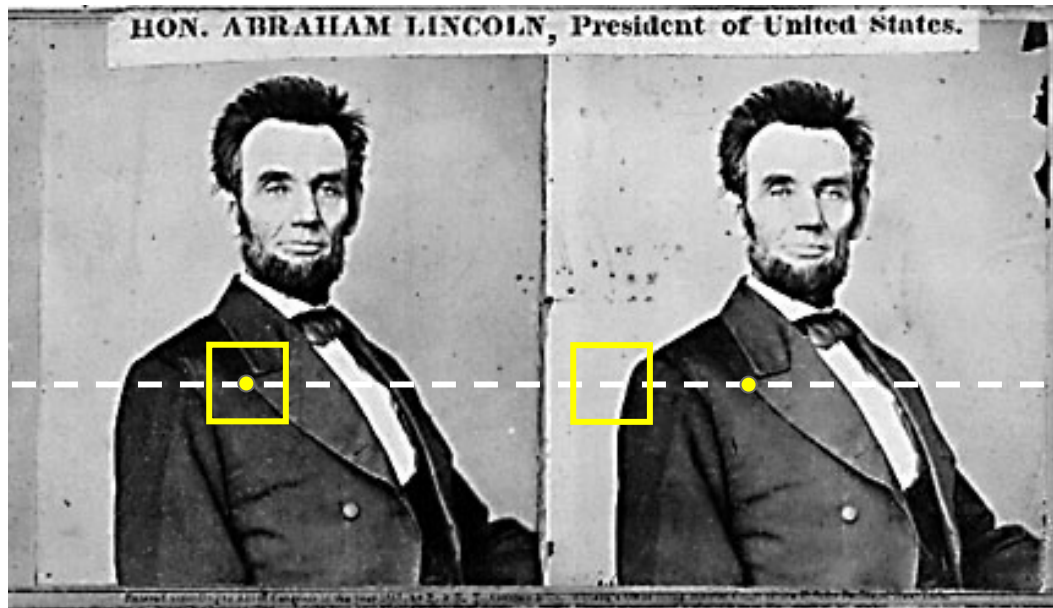
# Dense correspondence search

- Neighborhoods of corresponding points are similar in intensity patterns.



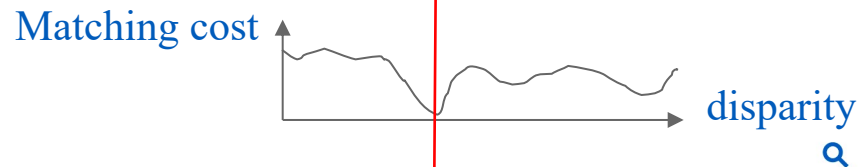
# Dense correspondence search

- For each epipolar line
  - For each pixel / window in the left image
    - Compare with every pixel / window on same epipolar line
    - Pick position with minimum match cost
      - SSD, normalized correlation



# Similarity constraints

- 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

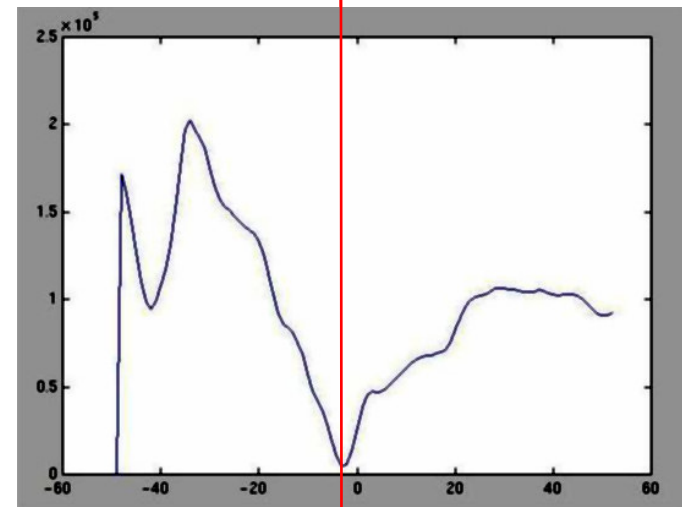


# Similarity constraints: SSD

Left

Right

scanline



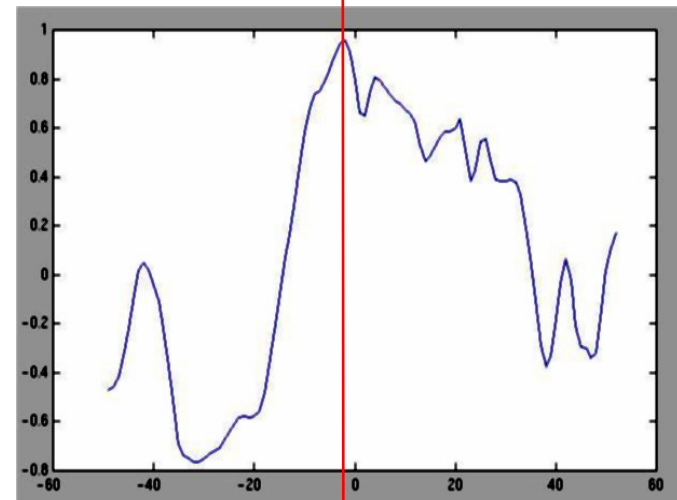
SSD

# Similarity constraints: Norm. Corr.

Left

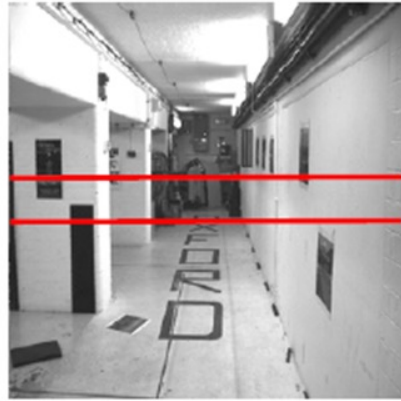
Right

scanline



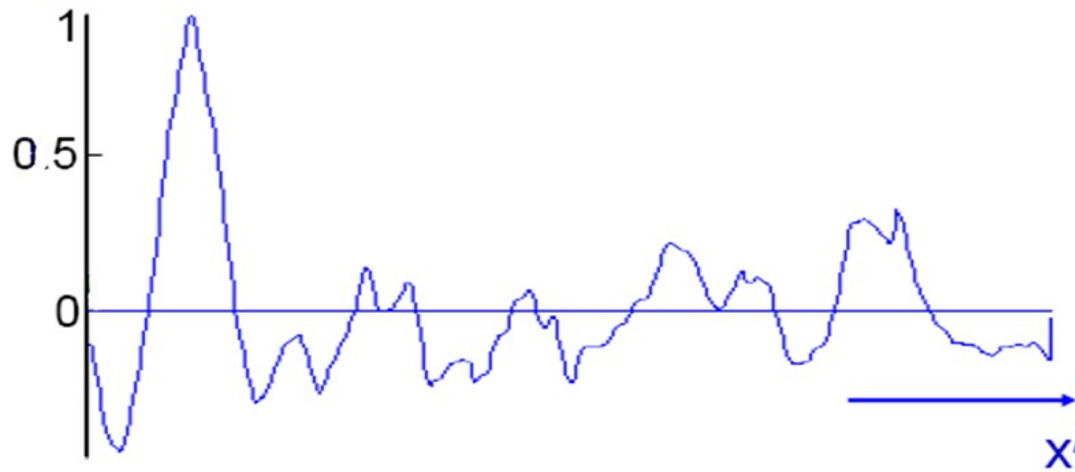
Norm. corr

# Correlation-based window matching



left image band ( $x$ )

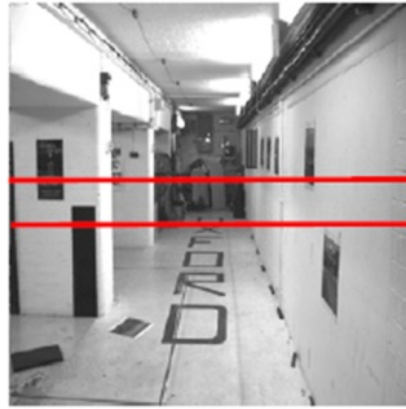
right image band ( $x'$ )



cross  
correlation

disparity =  $x' - x$

# Correlation-based window matching

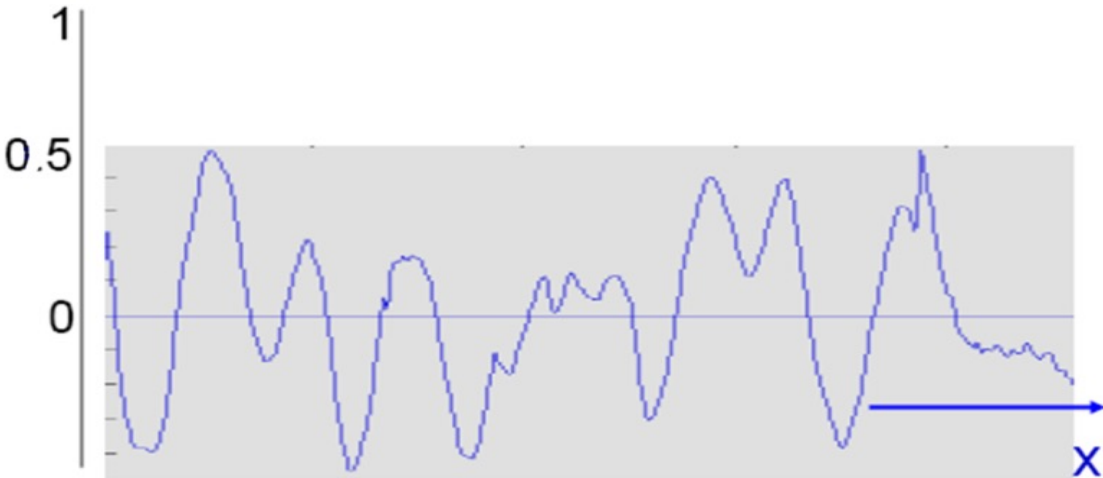


target region



left image band ( $x$ )

right image band ( $x'$ )



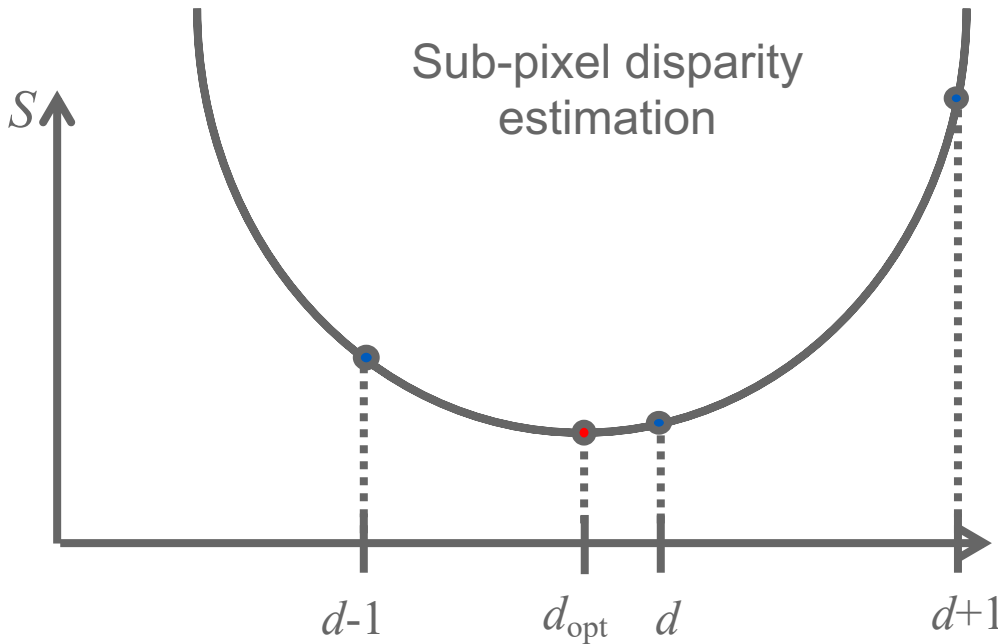
cross  
correlation

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



# Sub-pixel disparity estimation

- Let  $S$  be the SSD

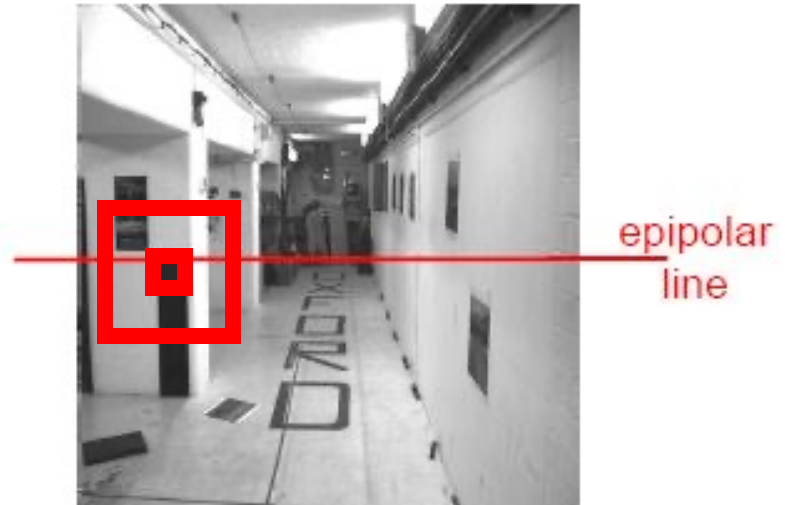


- $S(d) = ad^2 + bd + c$ 
  - $S(0) = c$
  - $S(1) = a + b + c$
  - $S(-1) = a - b + c$
- Solving this, we obtain:
  - $a = (S(1) + S(-1) - 2S(0))/2$
  - $b = (S(1) - S(-1))/2$
  - $c = S(0)$
- $S'(d) = 2ad + b = 0$

$$d_{opt} = \frac{(S(-1) - S(1))}{2(S(1) + S(-1) - 2S(0))}$$



# Effect of window size

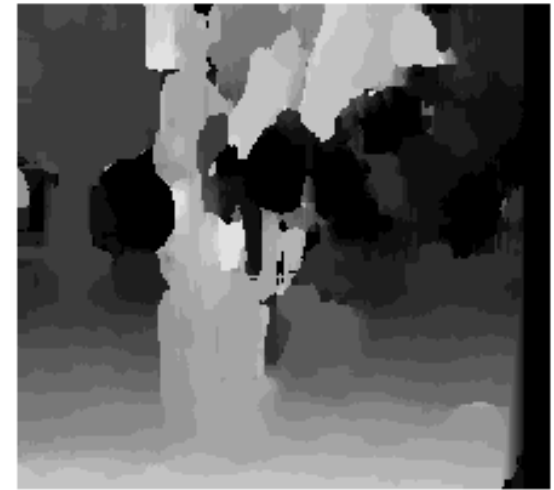


# Effect of window size

- large enough to have sufficient intensity variation
- small enough to contain only pixels with about the same disparity.



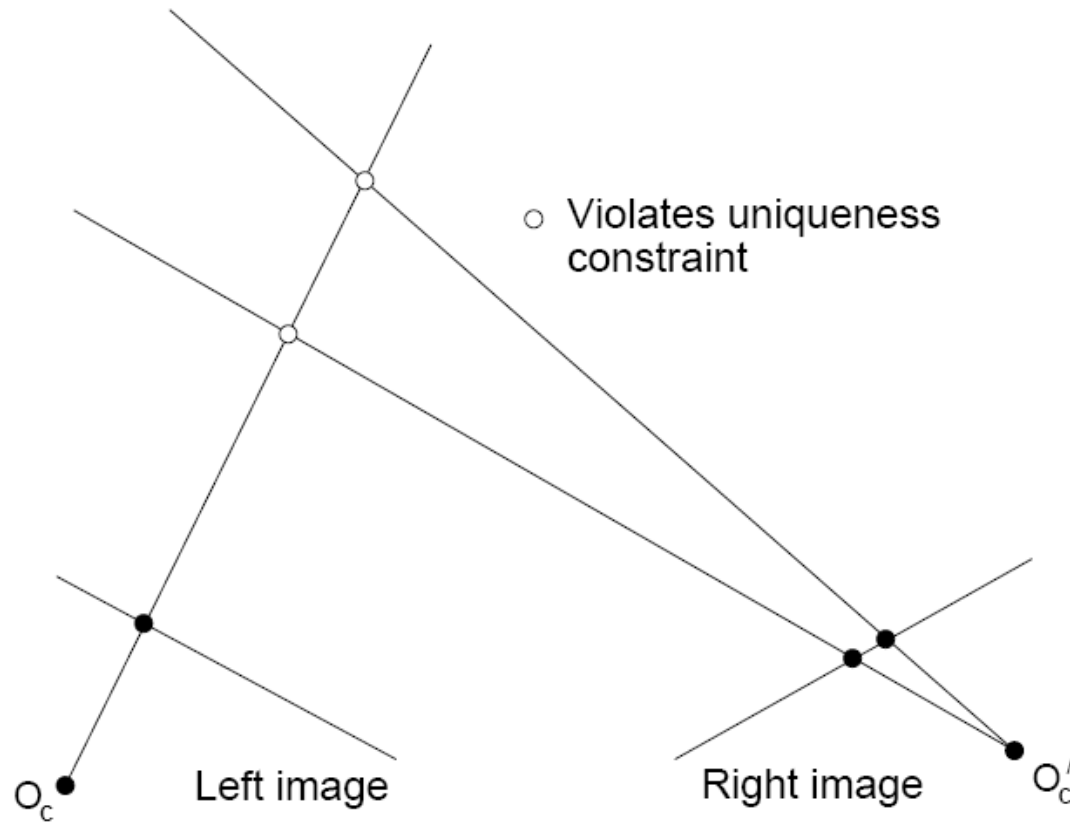
$W = 3$



$W = 20$

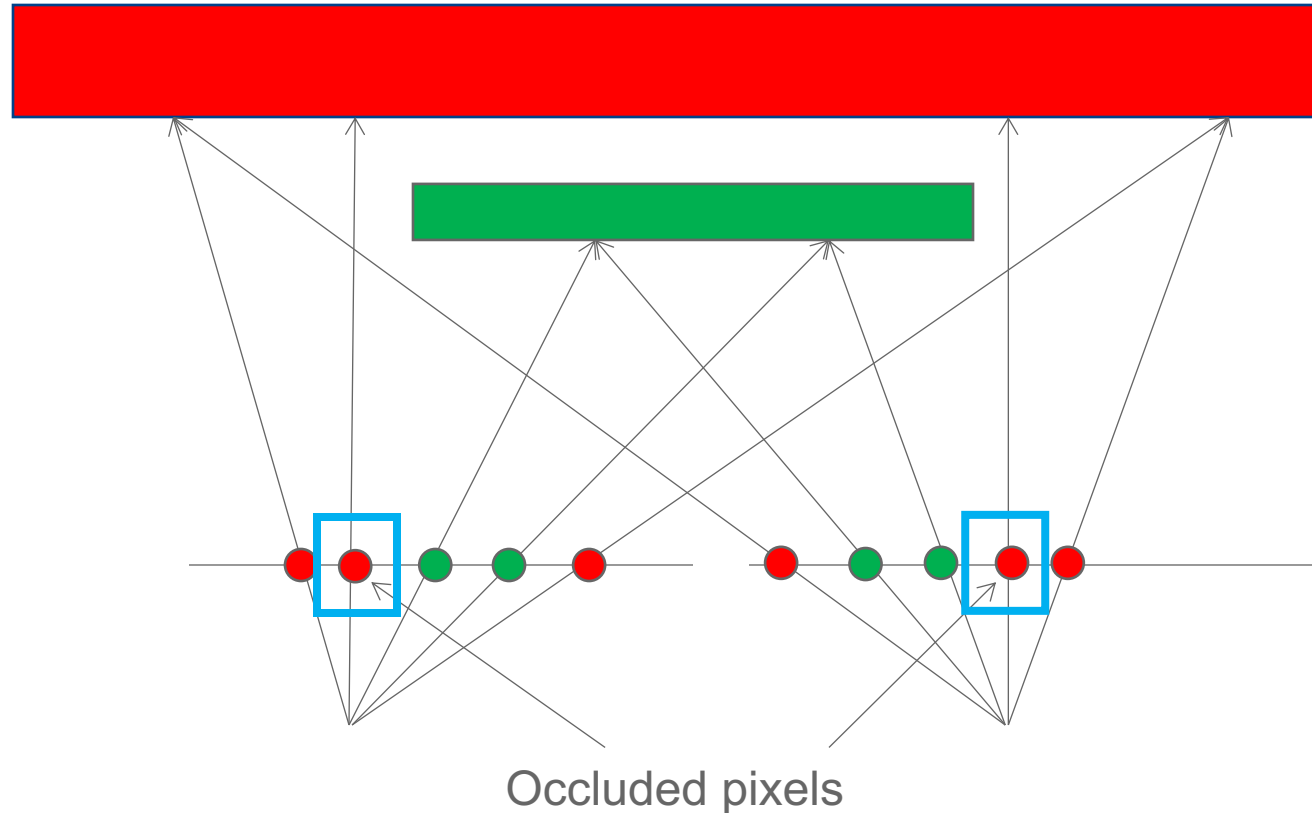
# Uniqueness constraint

- Up to one match in right image for every point in left image



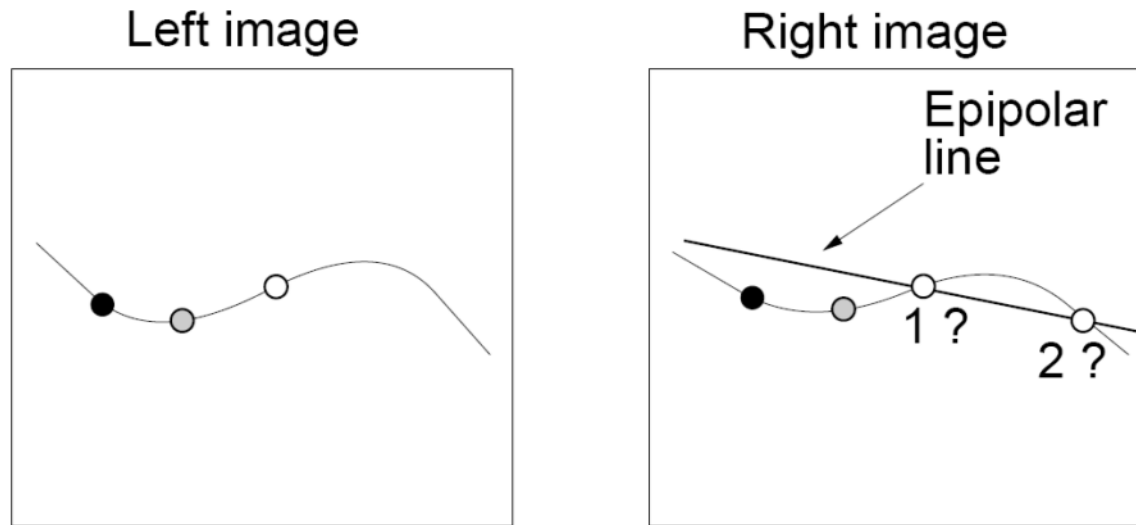
# Problem: Occlusion

- Uniqueness says “up to one match” per pixel
- What if there is no match?



# Disparity gradient constraint

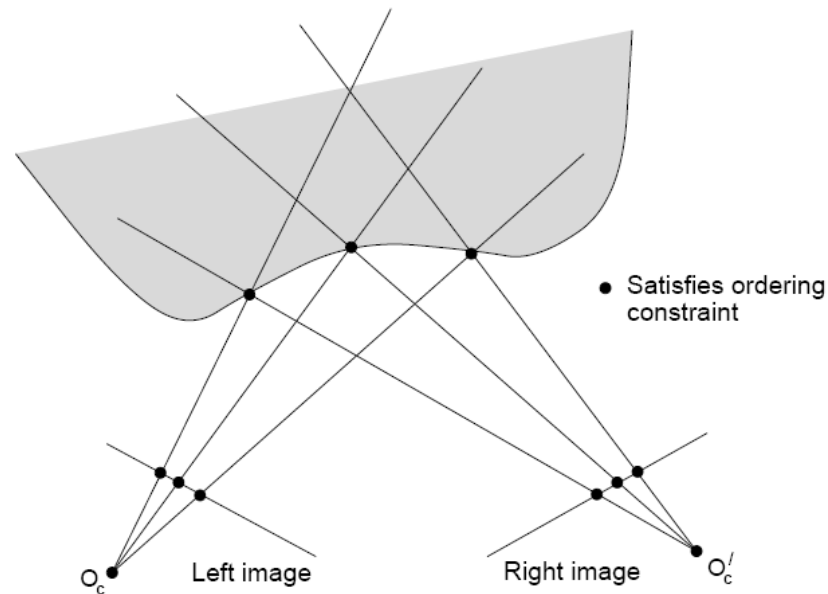
- Assume piecewise continuous surface, so we want disparity estimates to be locally smooth



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

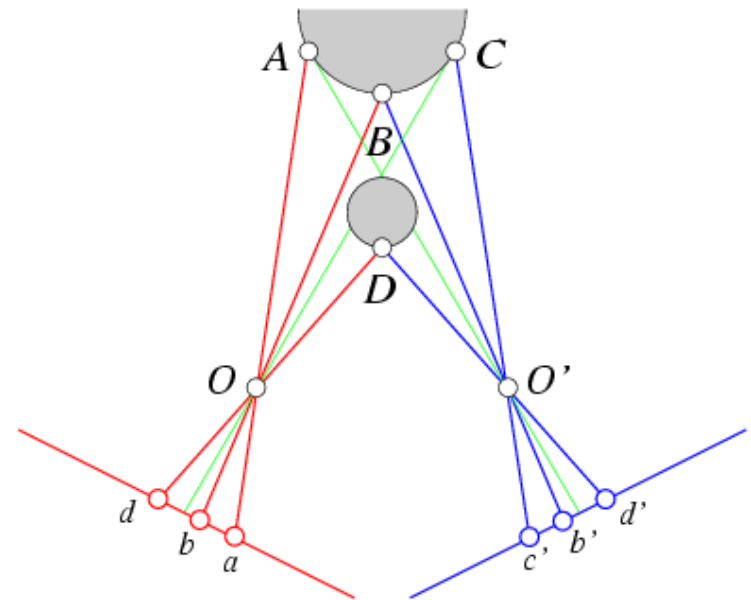
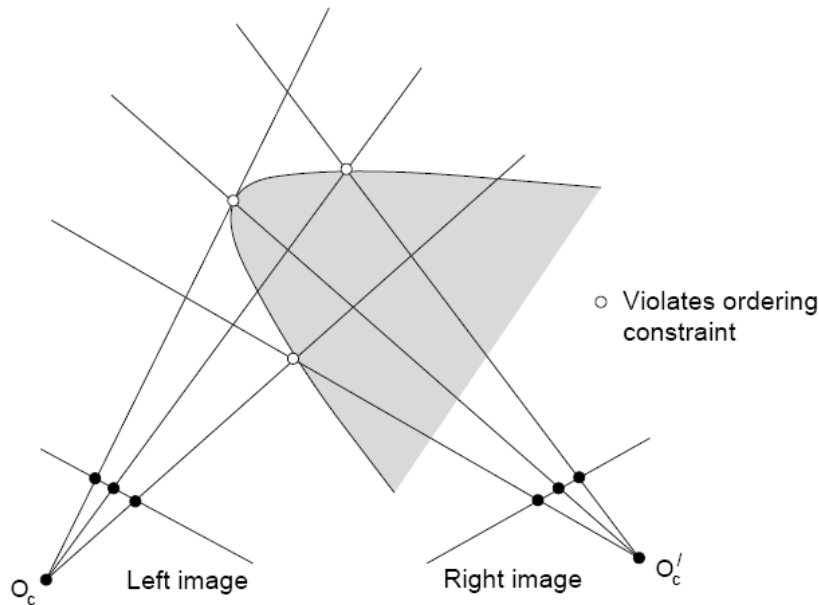
# Ordering constraint

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

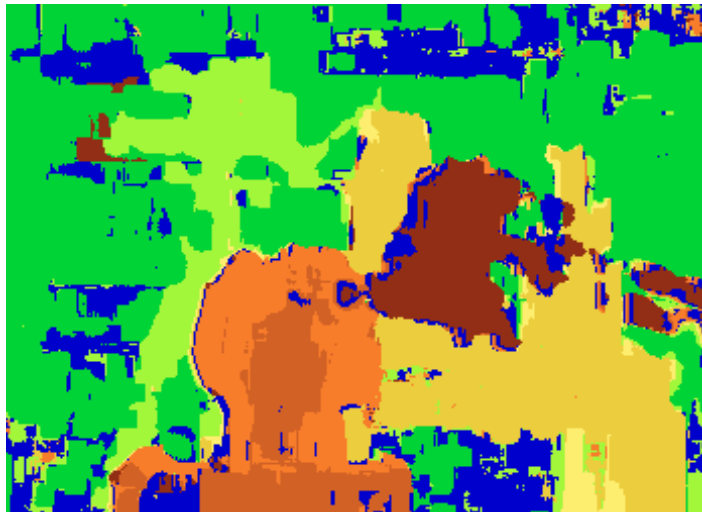
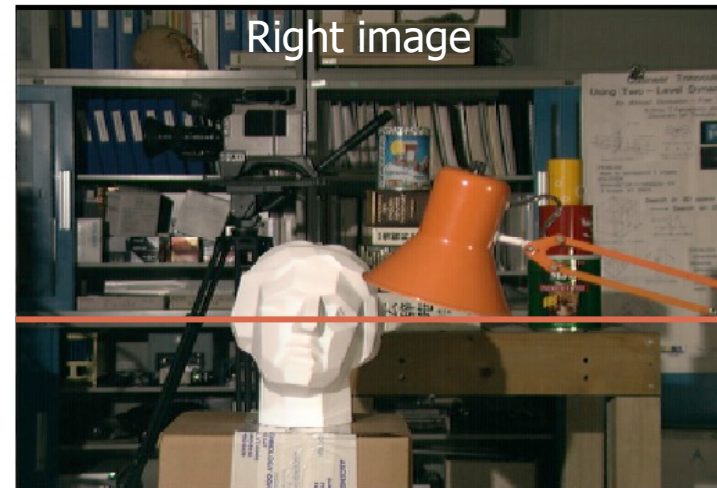
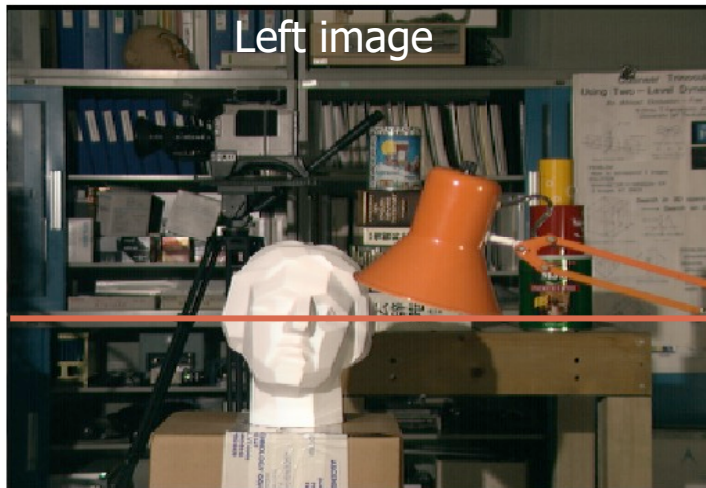


# Ordering constraint

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



# Results with window search



Window-based matching  
(best window size)



Ground truth



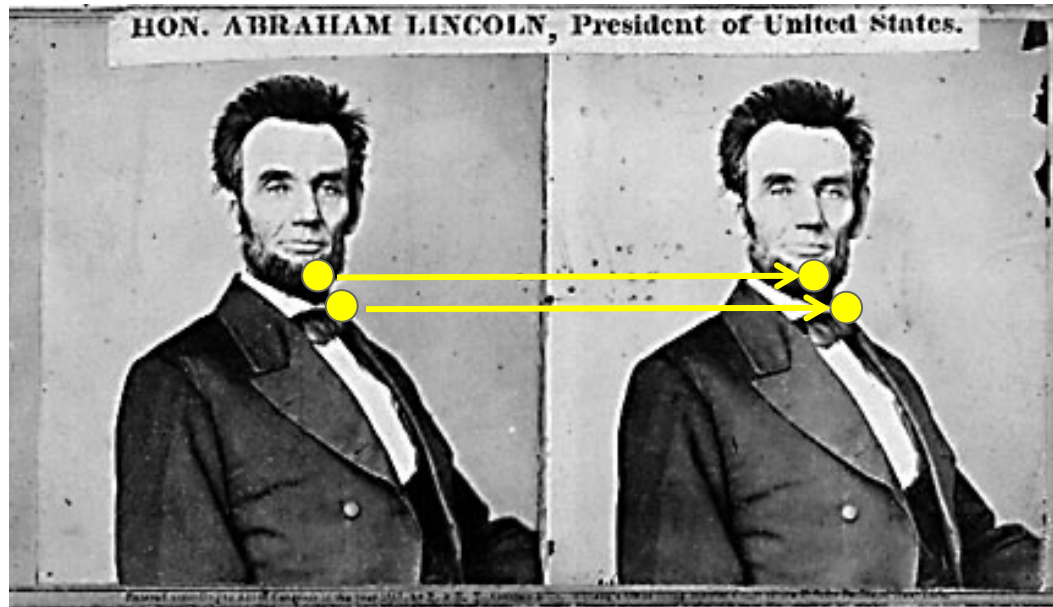
# Better solutions

---

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

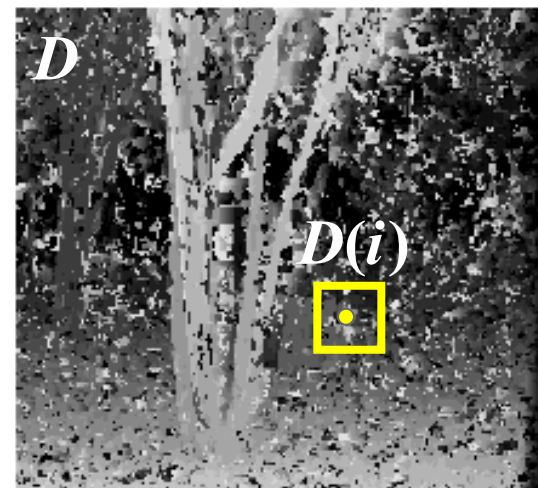
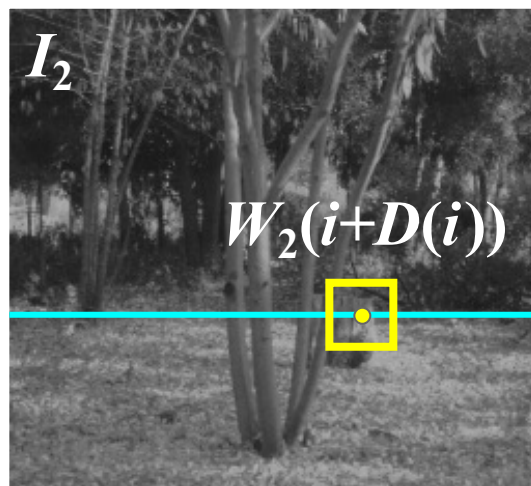
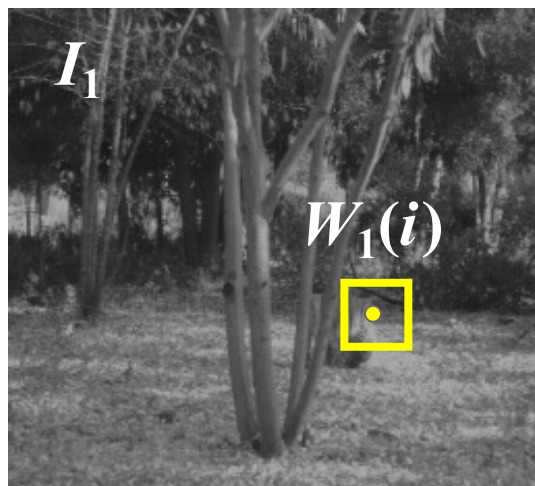
# Stereo as energy minimization

- What defines a good stereo correspondence?
  - Match quality
    - Want each pixel to find a good match in the other image
  - Smoothness
    - Adjacent pixels often move about the same amount.



# Stereo matching as energy minimization

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



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

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

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

# Better results...



Graph cut method



Ground truth

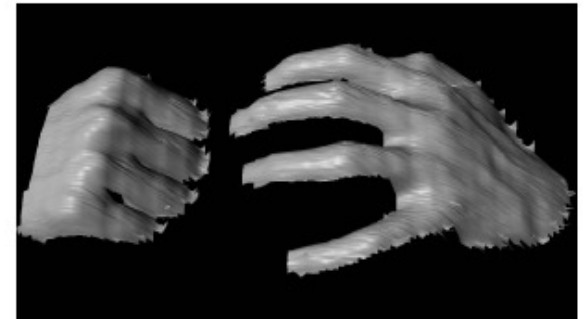
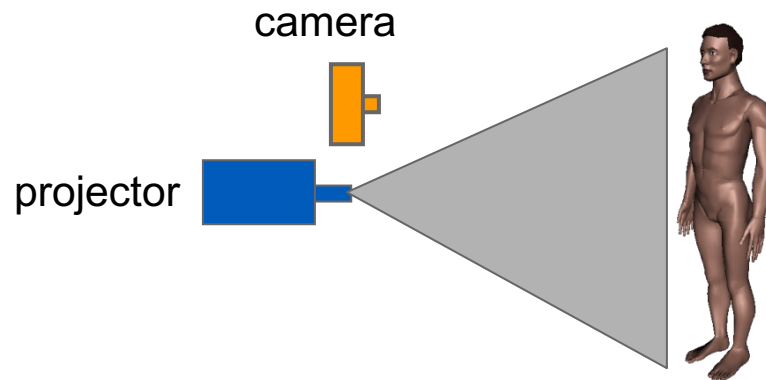
# Challenges

---

- 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



# Kinect: Structured infrared light





# iPhone X

- IR Emitter
- 30,000 points
- 2D IR snapshot

