

Lecture 6

Stereo Systems

Multi-view geometry



1891

Professor Silvio Savarese

Computational Vision and Geometry Lab

Lecture 6

Stereo Systems

Multi-view geometry



- Stereo systems
 - Rectification
 - Correspondence problem
- Multi-view geometry
 - The SFM problem
 - Affine SFM

Reading:

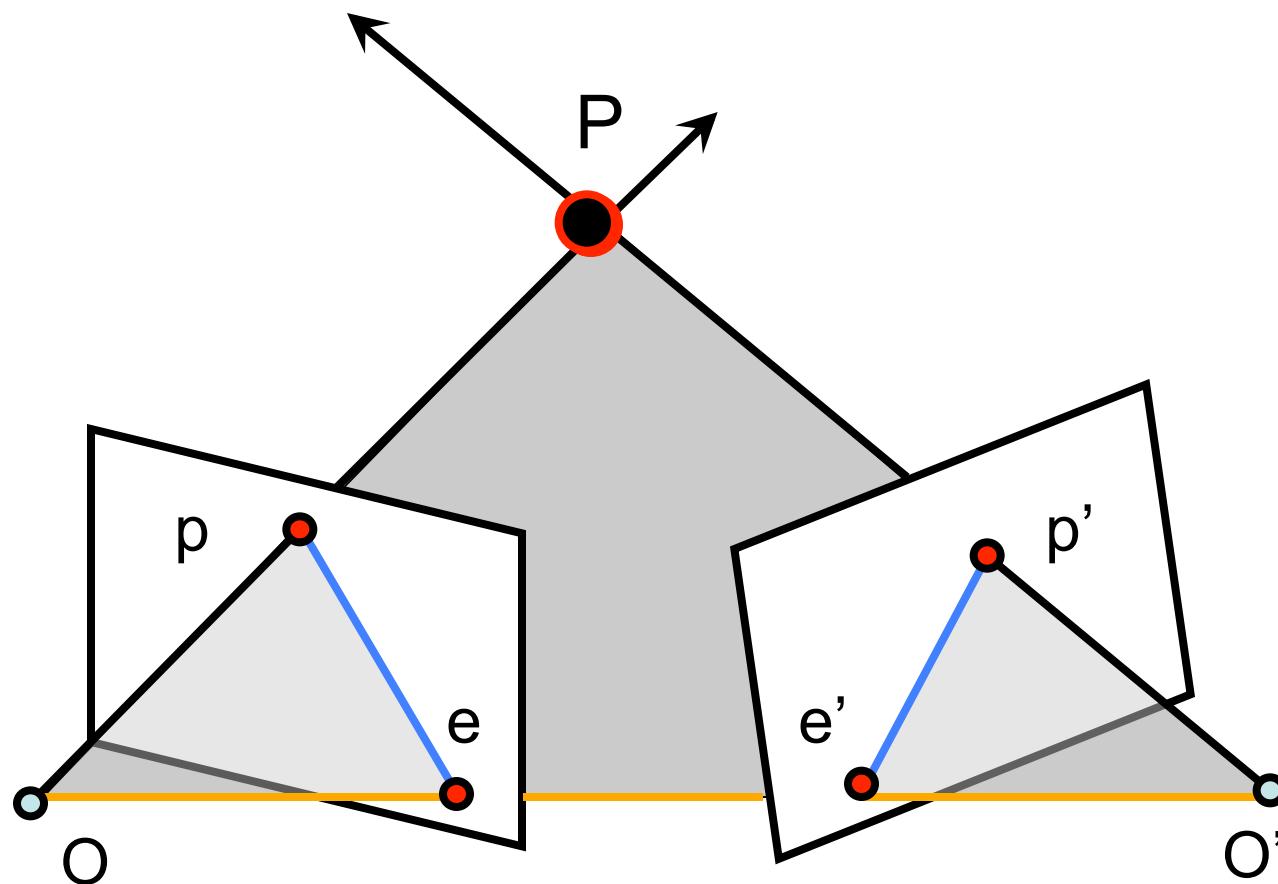
[AZ] Chapter: 9 “Epip. Geom. and the Fundam. Matrix Transf.”

[AZ] Chapter: 18 “N view computational methods”

[FP] Chapters: 7 “Stereopsis”

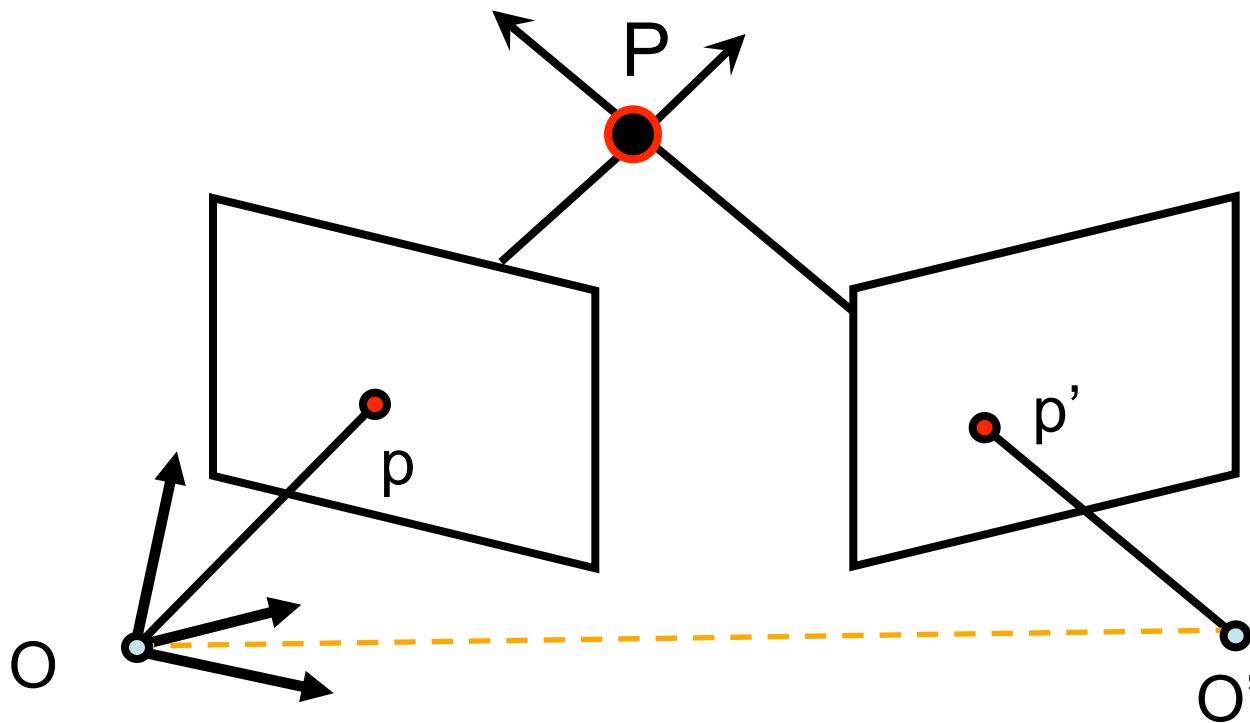
[FP] Chapters: 8 “Structure from Motion”

Epipolar geometry



- Epipolar Plane
- Baseline
- Epipolar Lines
- Epipoles e, e'
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Constraint

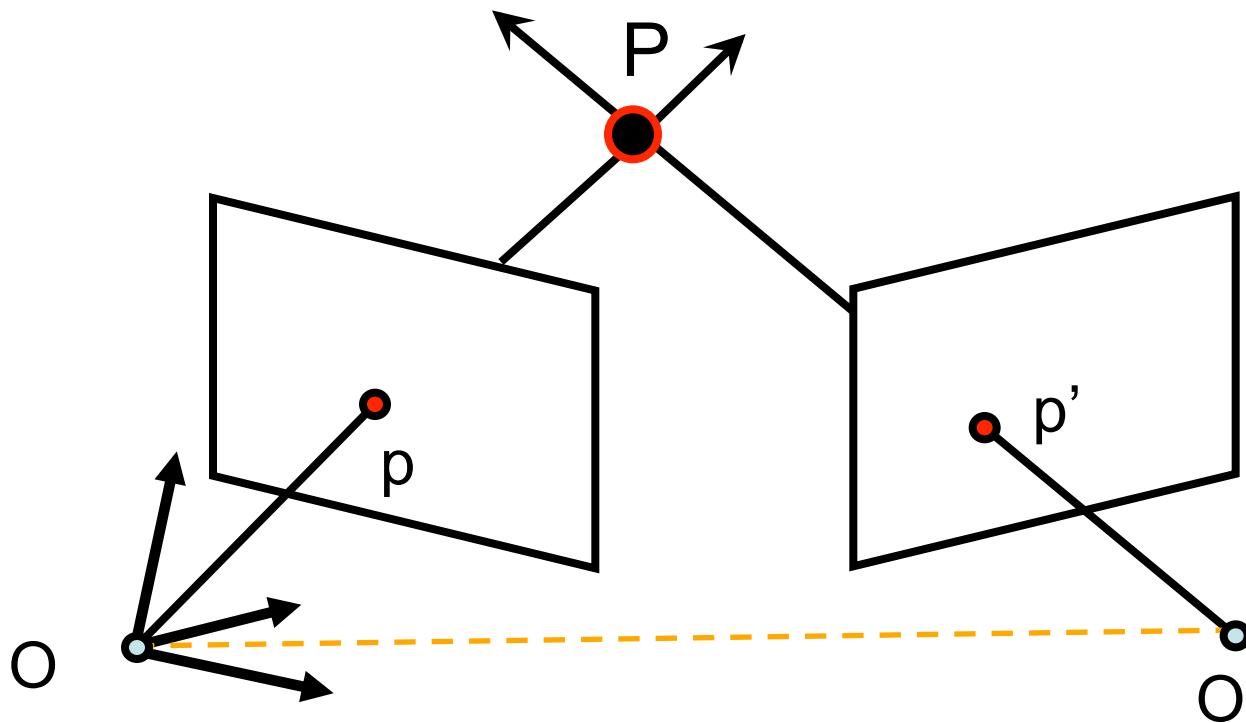


$$p^T E p' = 0$$

$$E = [T_x] \cdot R$$

E = Essential Matrix
(Longuet-Higgins, 1981)

Epipolar Constraint

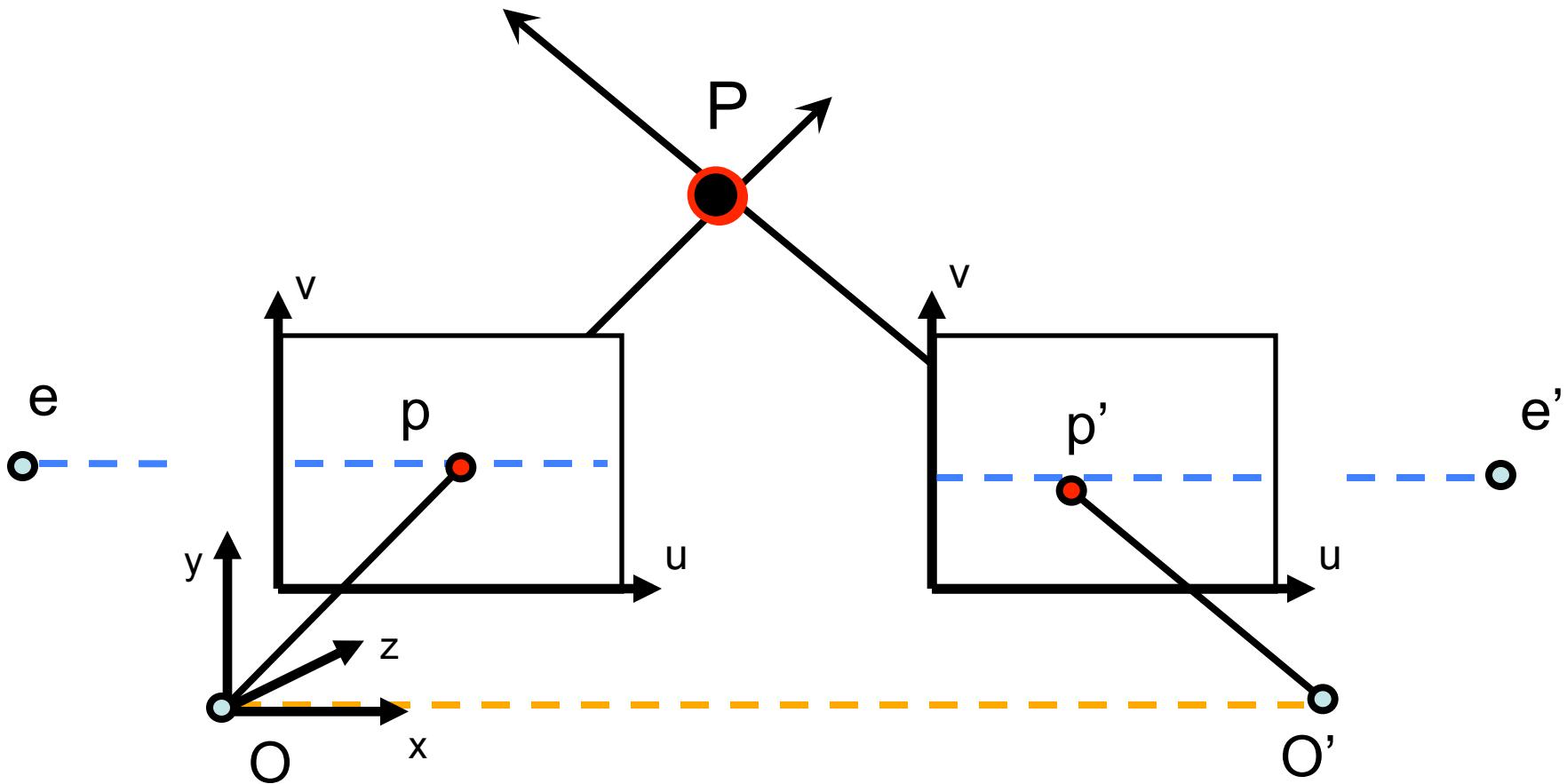


$$p^T F p' = 0$$

$$F = K^{-T} \cdot [T_x] \cdot R \cdot K'^{-1}$$

F = Fundamental Matrix
(Faugeras and Luong, 1992)

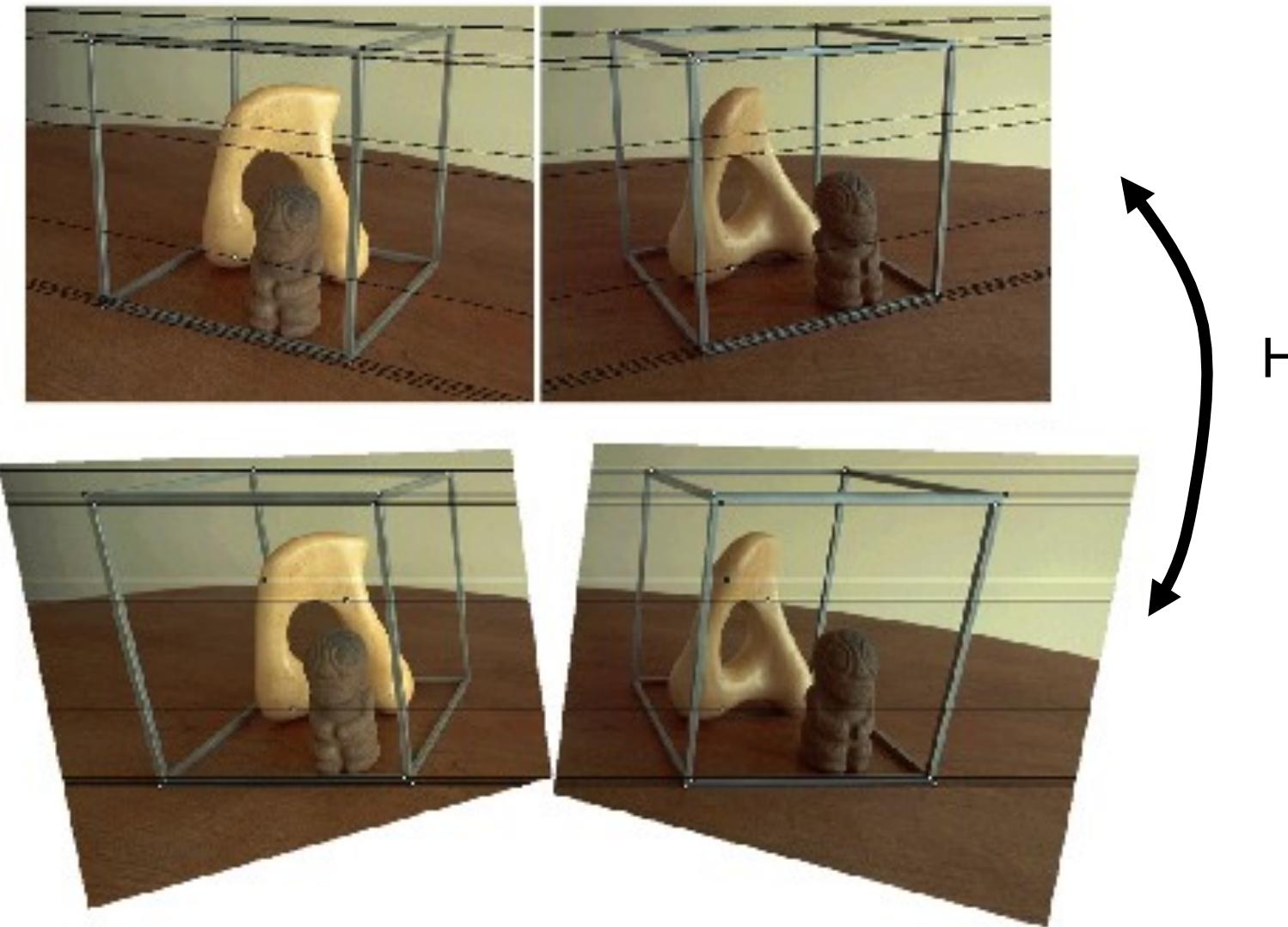
Example: Parallel image planes



- Epipolar lines are horizontal
- Epipoles go to infinity
- v -coordinates are equal

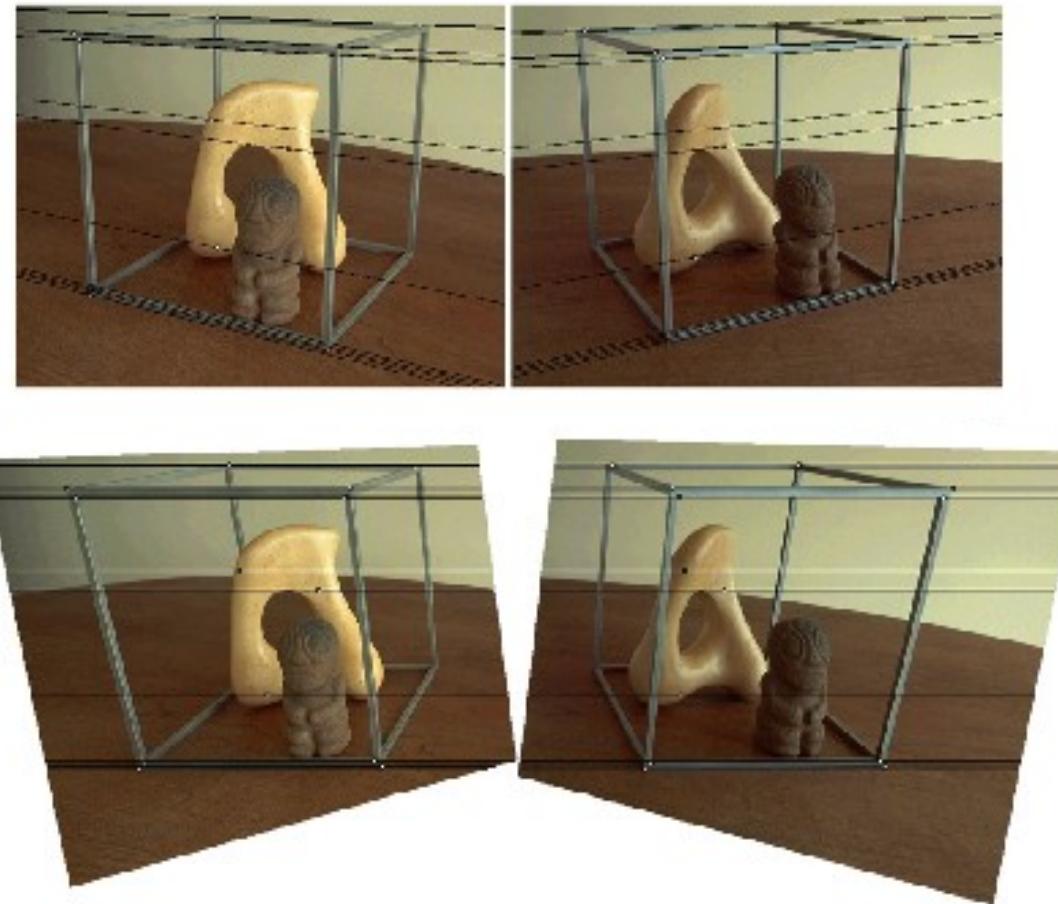
$$p = \begin{bmatrix} p_u \\ p_v \\ 1 \end{bmatrix} \quad p' = \begin{bmatrix} p'_u \\ p'_v \\ 1 \end{bmatrix}$$

Rectification: making two images “parallel”



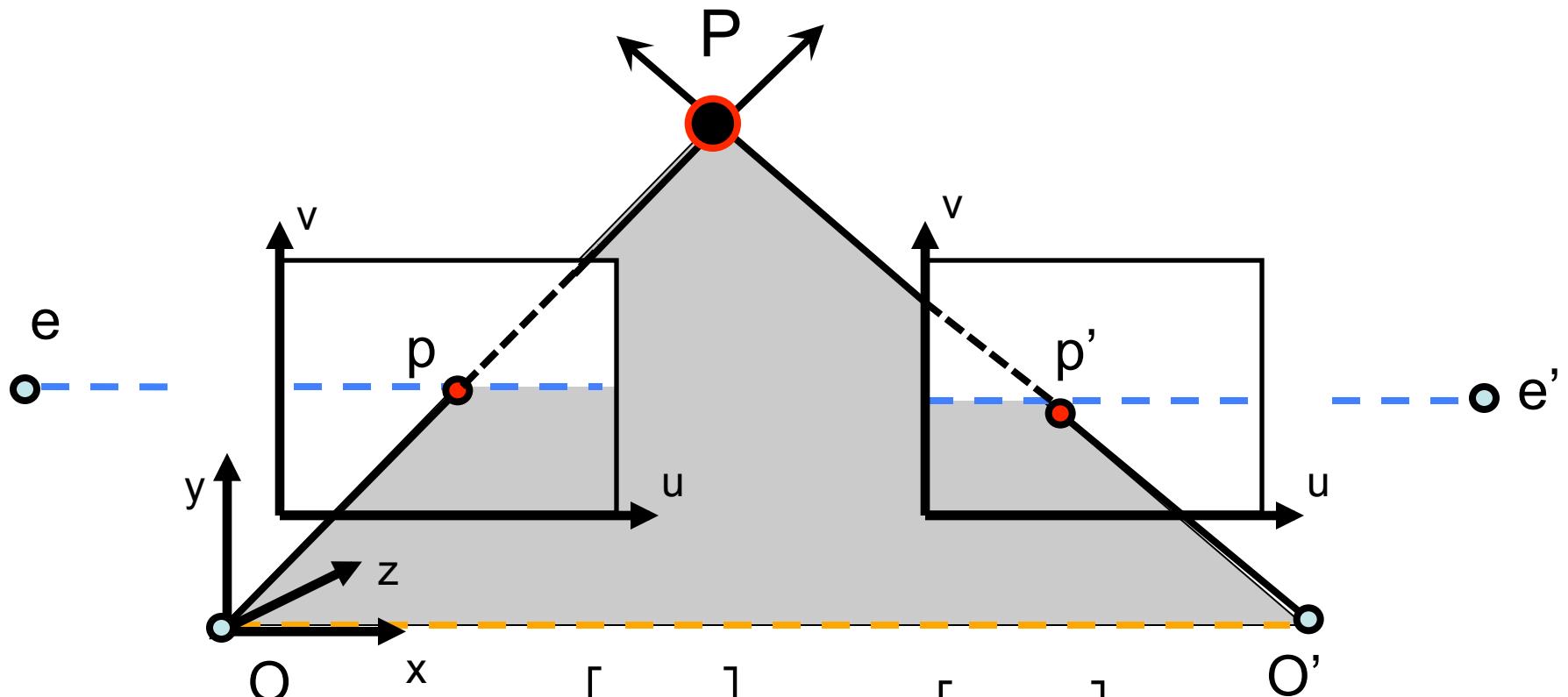
Courtesy figure S. Lazebnik

Why are parallel images useful?



- Makes triangulation easy
- Makes the correspondence problem easier

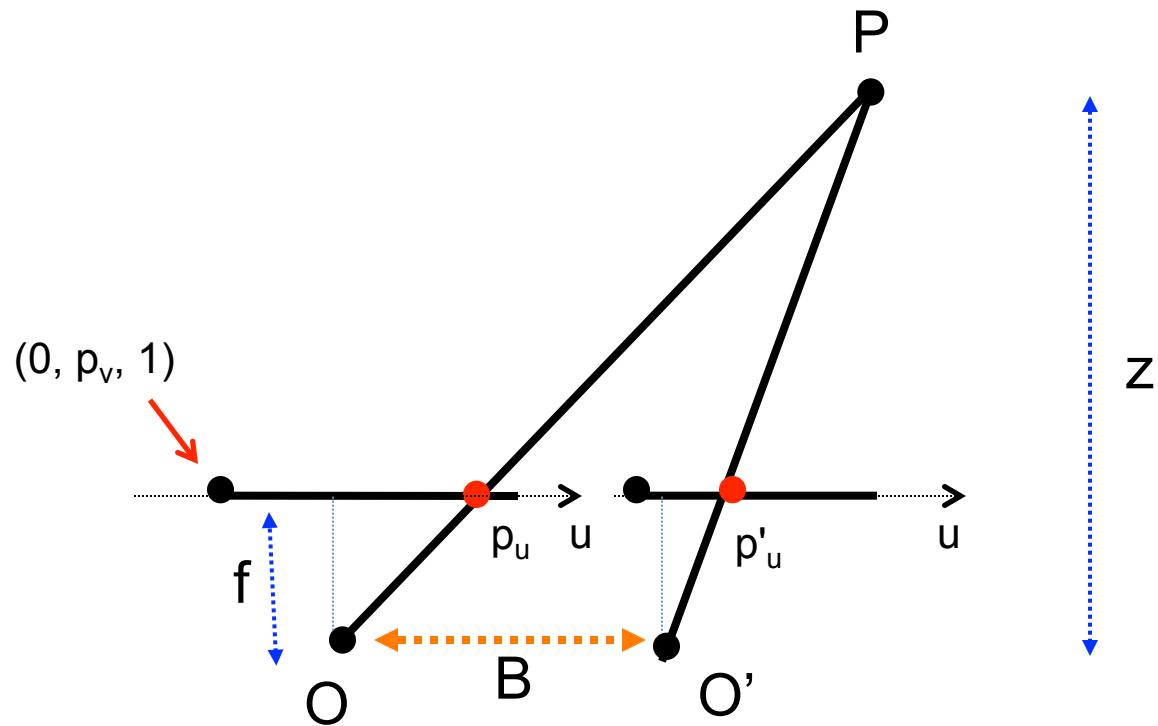
Point triangulation



$$p = \begin{bmatrix} p_u \\ p_v \\ 1 \end{bmatrix}$$

$$p' = \begin{bmatrix} p'_u \\ p'_v \\ 1 \end{bmatrix}$$

Computing depth

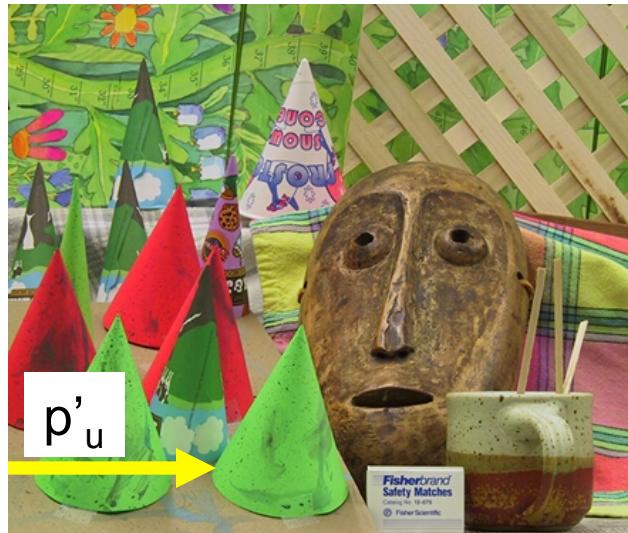
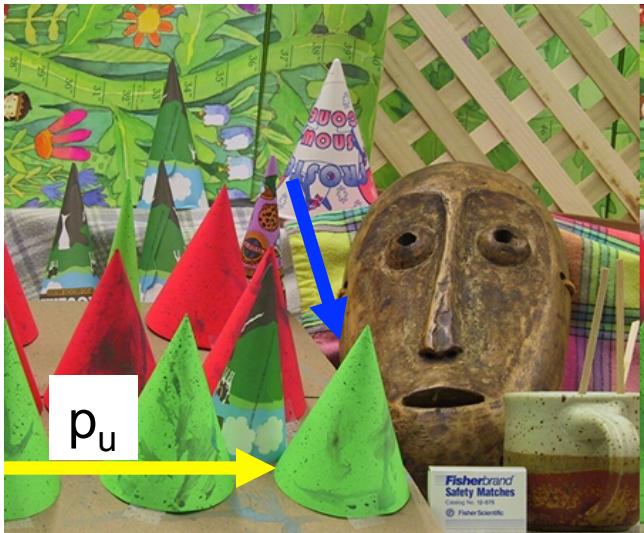


$$\text{disparity} = p_u - p'_u \propto \frac{B \cdot f}{z} \quad [\text{Eq. 1}]$$

Note: Disparity is inversely proportional to depth

Disparity maps

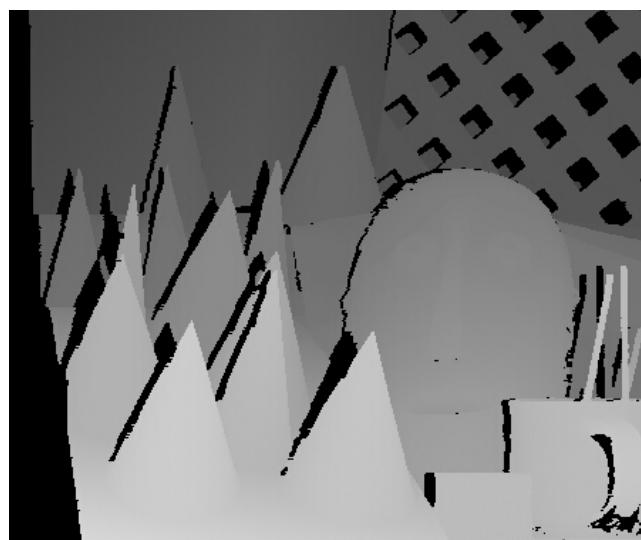
<http://vision.middlebury.edu/stereo/>



$$p_u - p'_u \propto \frac{B \cdot f}{z}$$

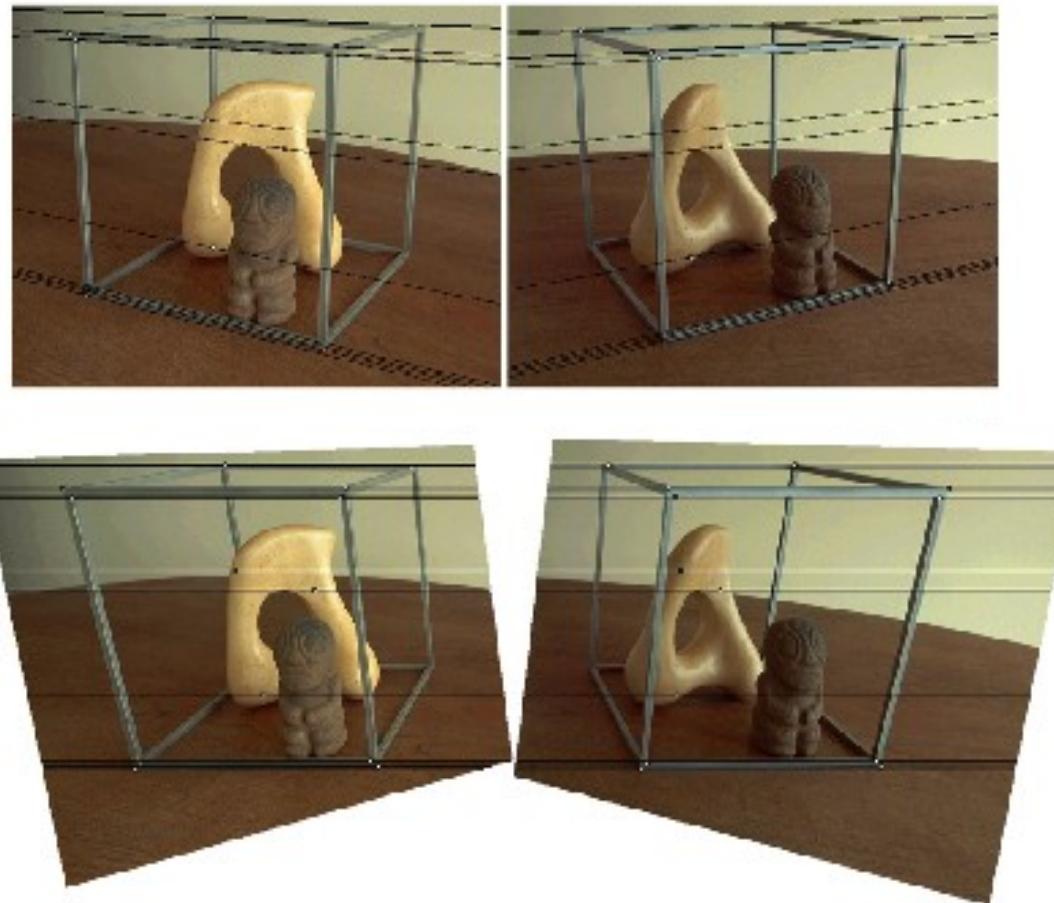
[Eq. 1]

Stereo pair



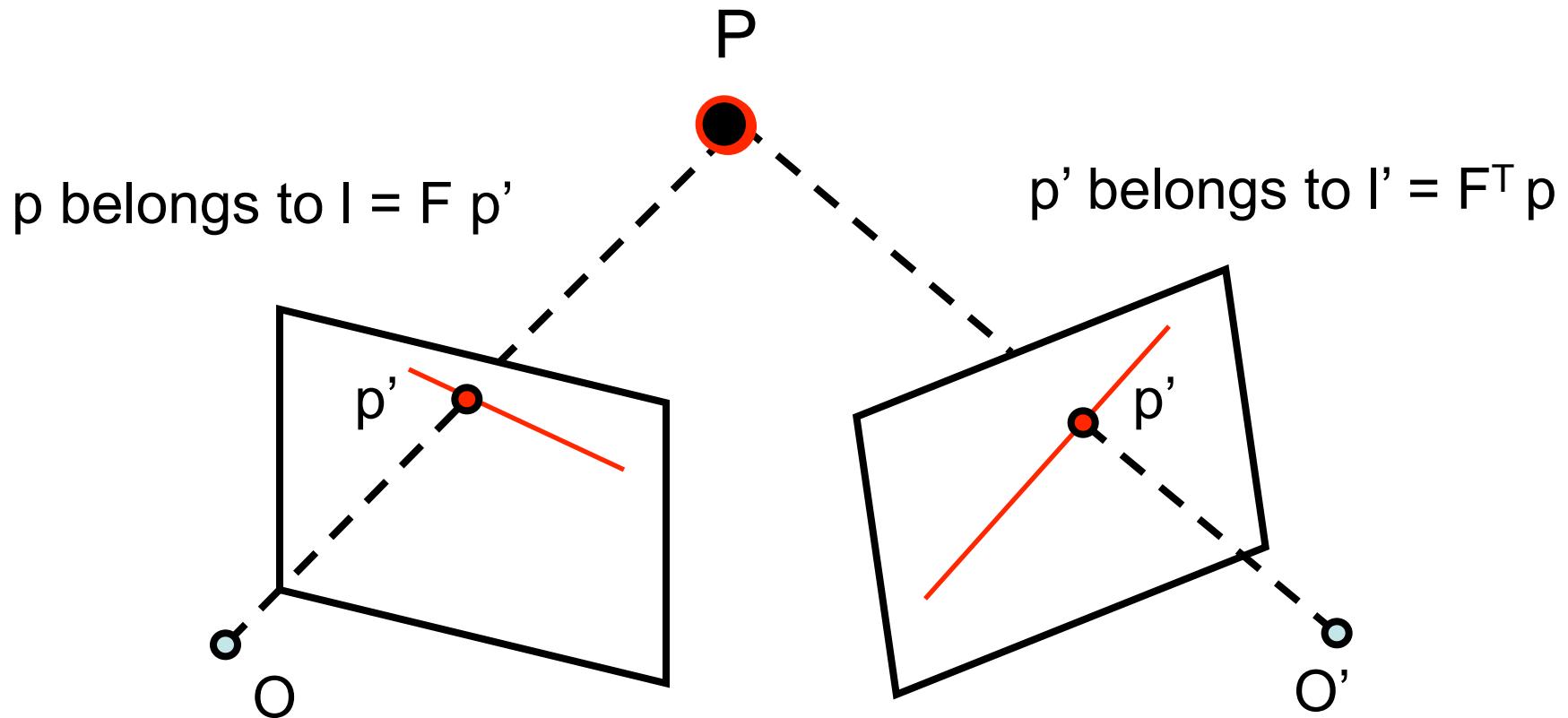
Disparity map / depth map

Why are parallel images useful?



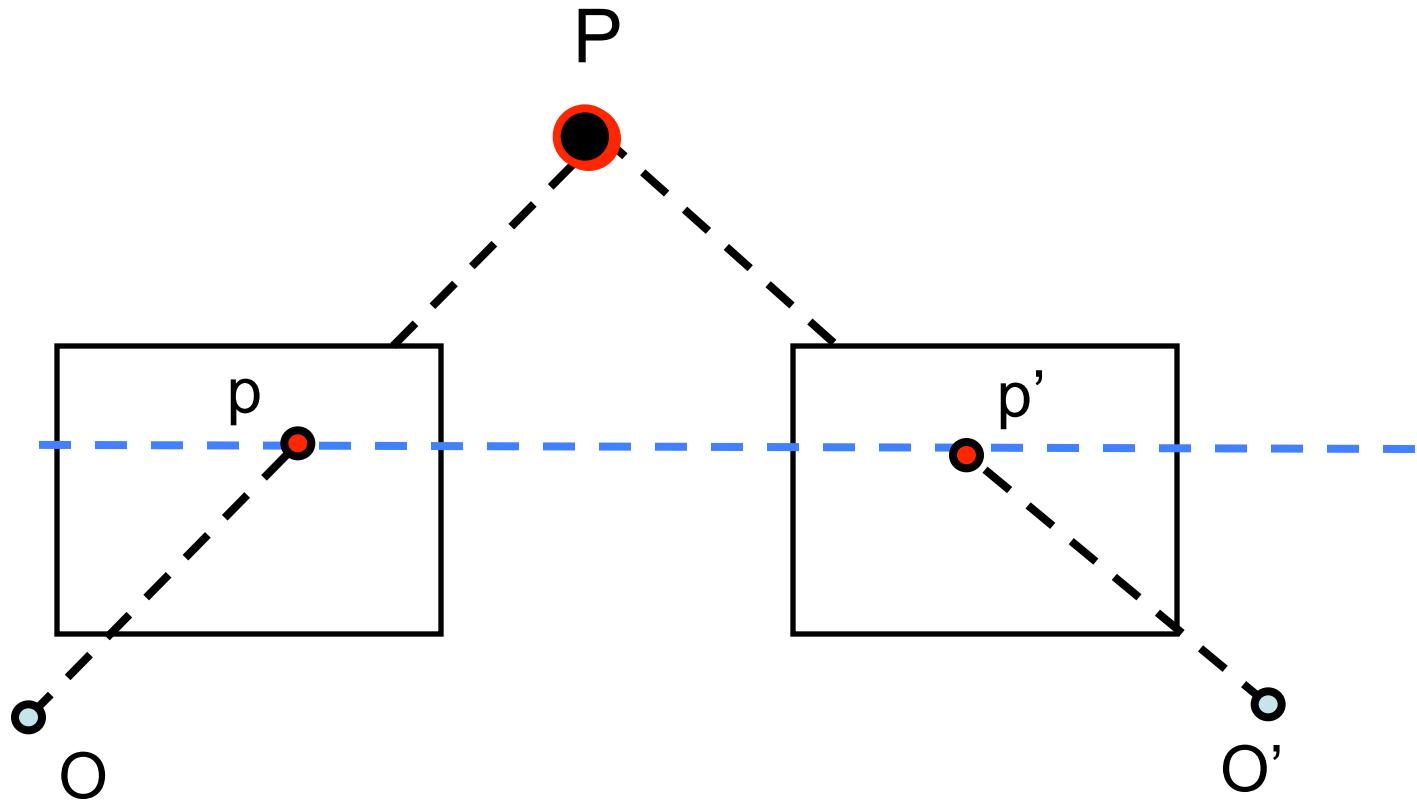
- Makes triangulation easy
- Makes the correspondence problem easier

Correspondence problem



Given a point in 3D, discover corresponding observations in left and right images [also called binocular fusion problem]

Correspondence problem



When images are rectified, this problem is much easier!

Correspondence problem

- A Cooperative Model (Marr and Poggio, 1976)
- Correlation Methods (1970--)
- Multi-Scale Edge Matching (Marr, Poggio and Grimson, 1979-81)

[FP] Chapters: 7

Correlation Methods (1970--)

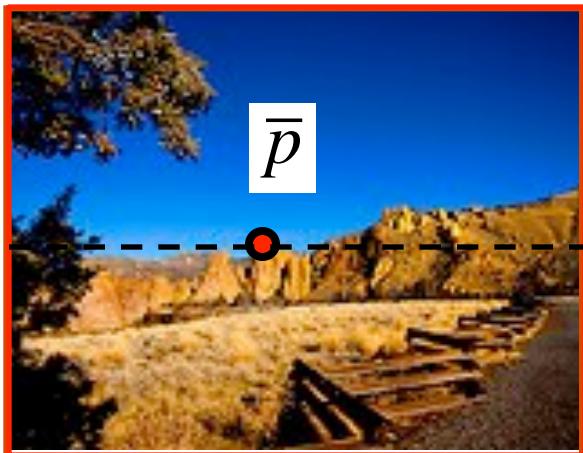


image 1

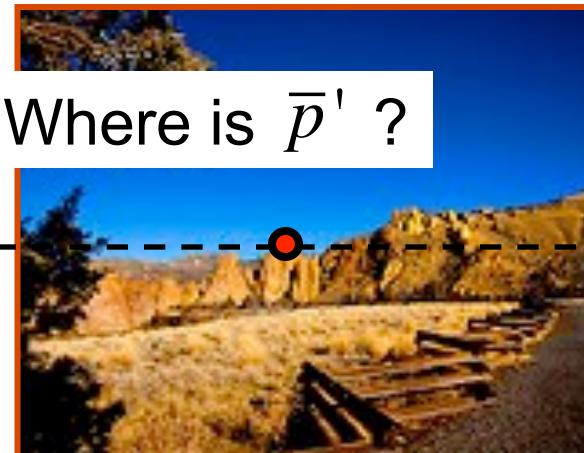
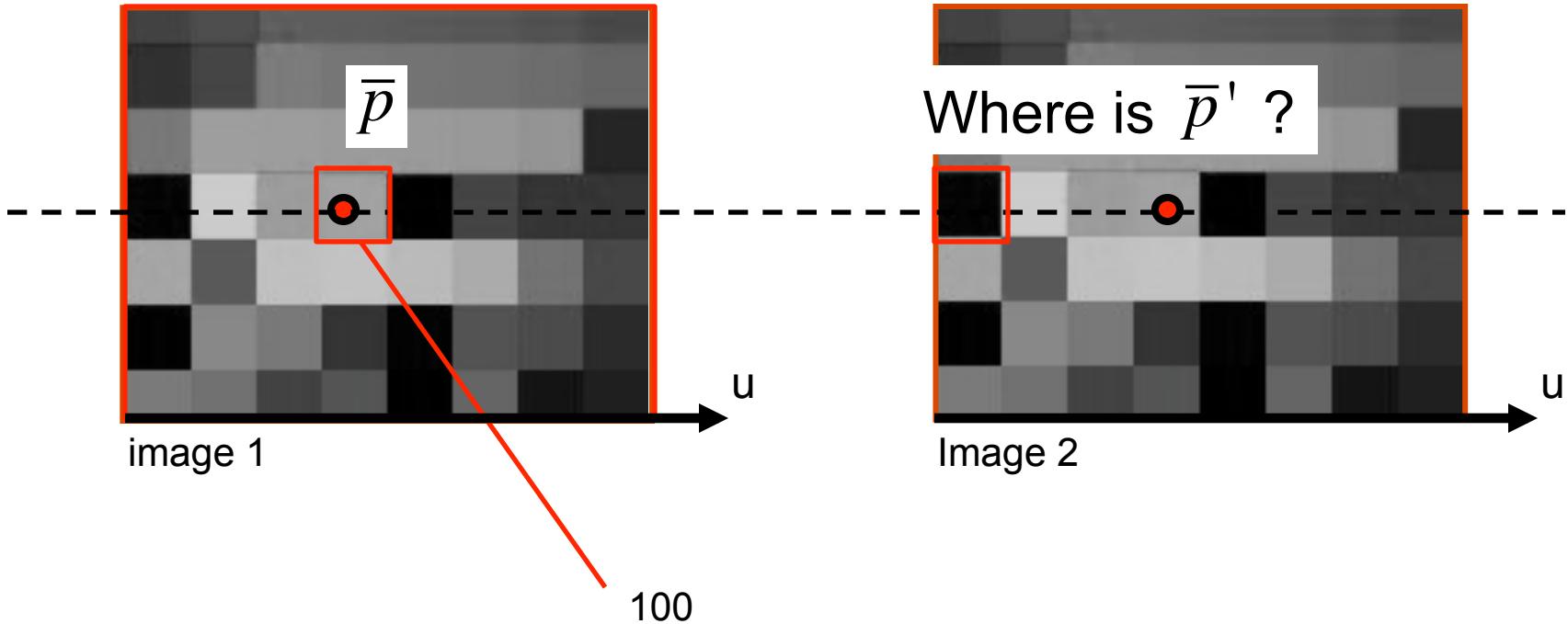


Image 2

$$\bar{p} = \begin{bmatrix} \bar{u} \\ \bar{v} \\ 1 \end{bmatrix}$$

$$\bar{p}' = \begin{bmatrix} \bar{u}' \\ \bar{v} \\ 1 \end{bmatrix}$$

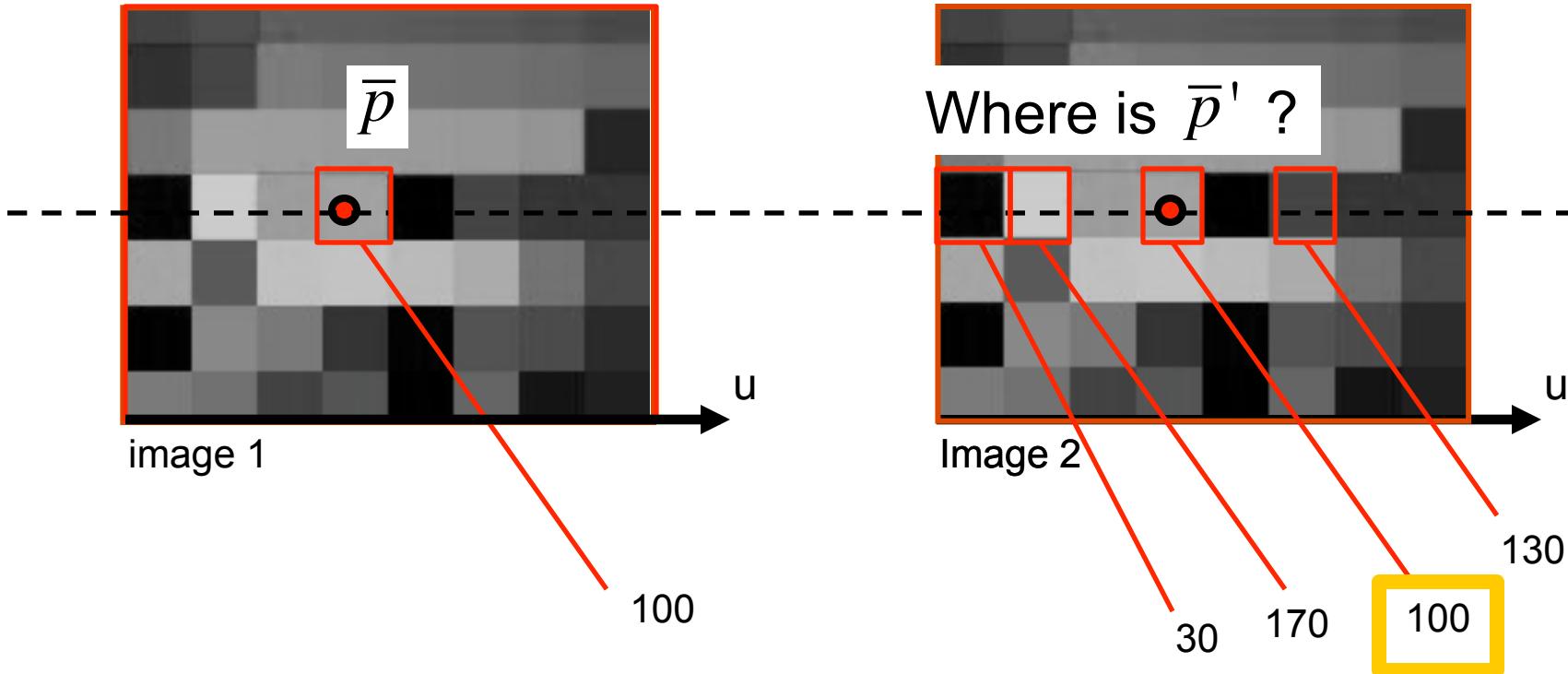
Correlation Methods (1970--)



$$\bar{p} = \begin{bmatrix} \bar{u} \\ \bar{v} \\ 1 \end{bmatrix}$$

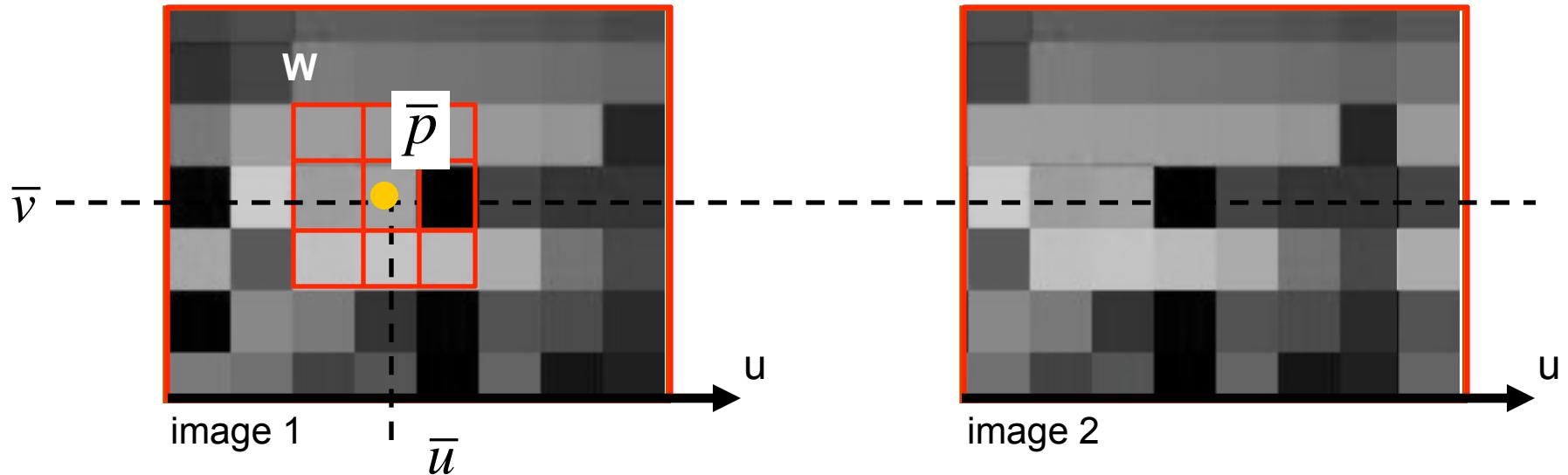
$$\bar{p}' = \begin{bmatrix} \bar{u}' \\ \bar{v} \\ 1 \end{bmatrix}$$

Correlation Methods (1970--)



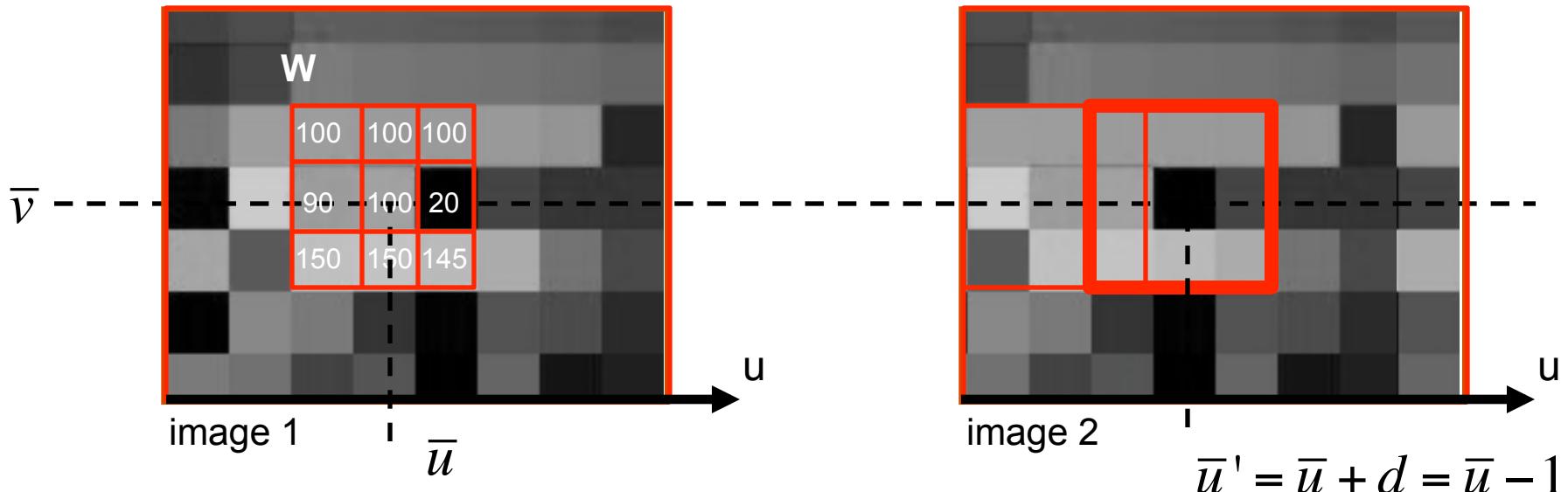
What's the problem with this?

Window-based correlation



- Pick up a window \mathbf{W} around $\bar{p} = (\bar{u}, \bar{v})$
- Build vector \mathbf{w}

Window-based correlation



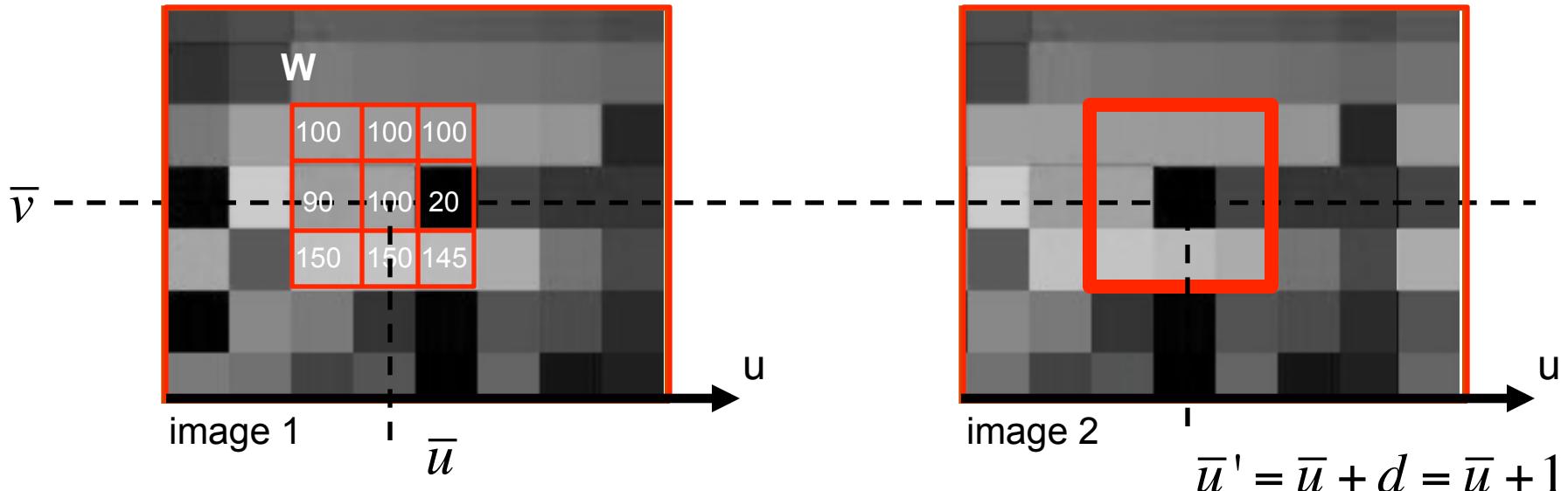
Example: \mathbf{W} is a 3×3 window in red

\mathbf{w} is a 9×1 vector

$$\mathbf{w} = [100, 100, 100, 90, 100, 20, 150, 150, 145]^T$$

- Pick up a window \mathbf{W} around $\bar{p} = (\bar{u}, \bar{v})$
- Build vector \mathbf{w}
- Slide the window \mathbf{W} along $v = \bar{v}$ in image 2 and compute $\mathbf{w}'(u)$ for each u
- Compute the dot product $\mathbf{w}^T \mathbf{w}'(u)$ for each u and retain the max value

Window-based correlation



Example: **W** is a 3x3 window in red

w is a 9x1 vector

$$\mathbf{w} = [100, 100, 100, 90, 100, 20, 150, 150, 145]^T$$

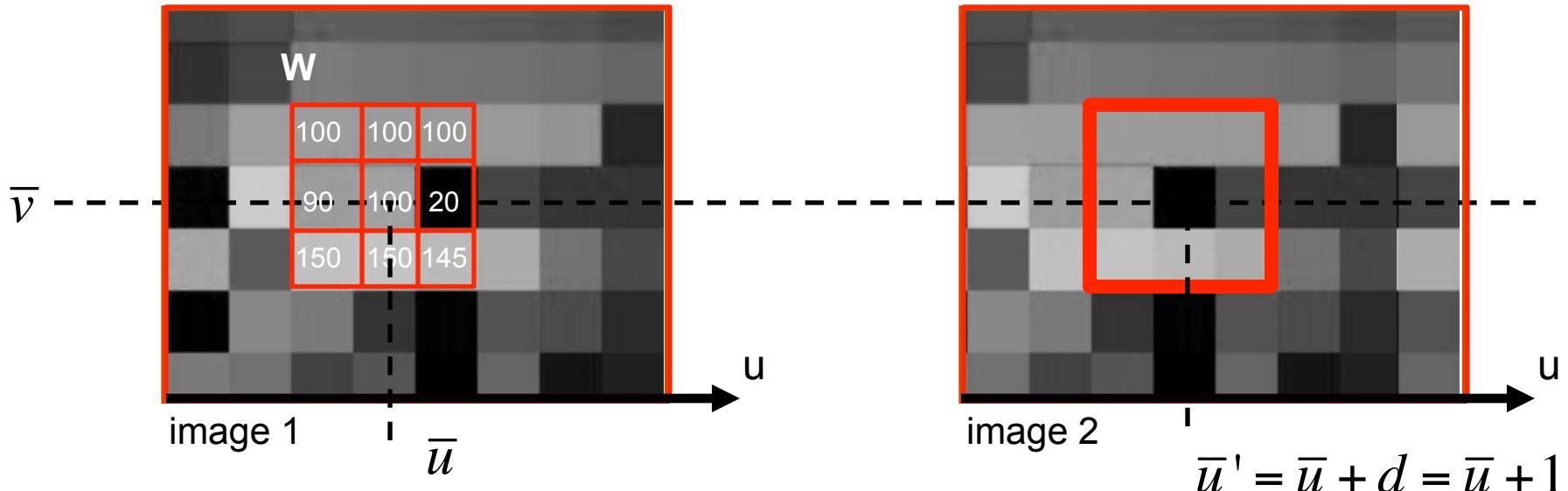
What's the problem with this?

Changes of brightness/exposure



Changes in the mean and the variance of intensity values in corresponding windows!

Normalized cross-correlation



Find u that maximizes:

$$\frac{(w - \bar{w})^T (w'(u) - \bar{w}')} {\| (w - \bar{w}) \| \| (w'(u) - \bar{w}') \|} \quad [\text{Eq. 2}]$$

\bar{w} = mean value within \mathbf{W}
located at $u^{\bar{u}}$ in image 1

$\bar{w}'(u)$ = mean value within \mathbf{W}
located at u in image 2

Example

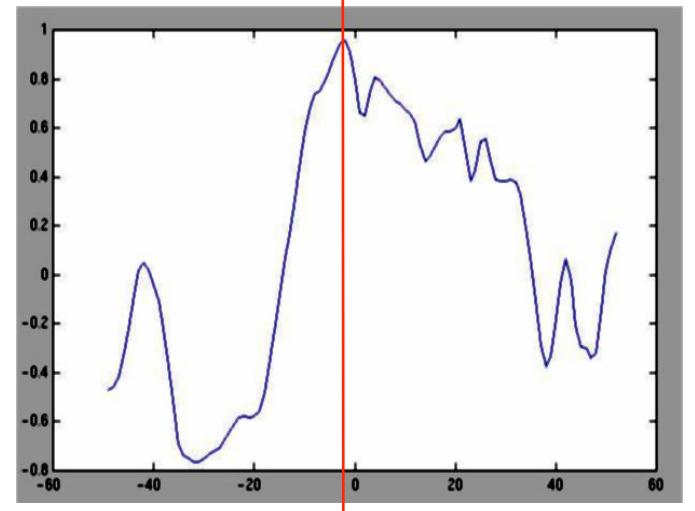
Image 1



Image 2



NCC

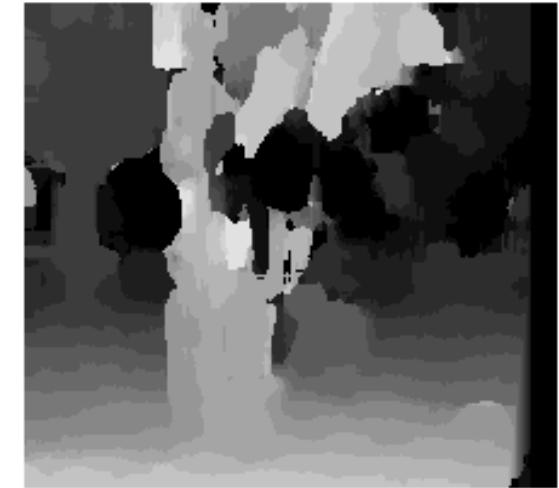


u

Effect of the window's size



Window size = 3

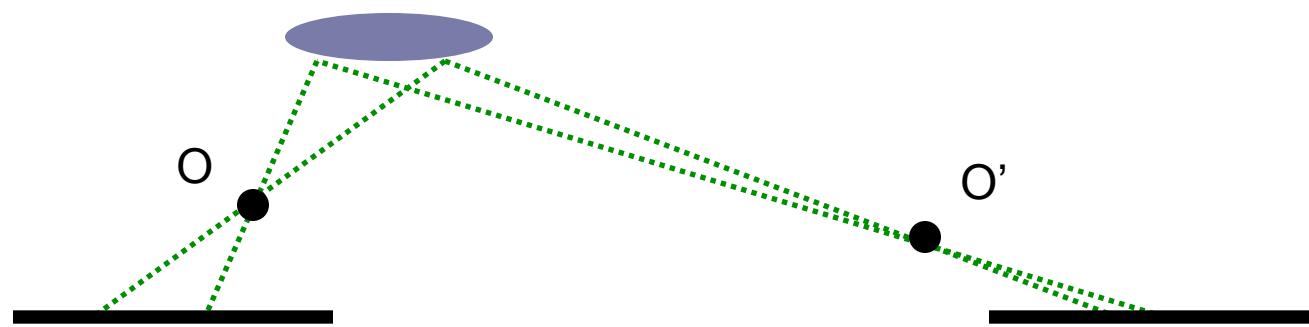


Window size = 20

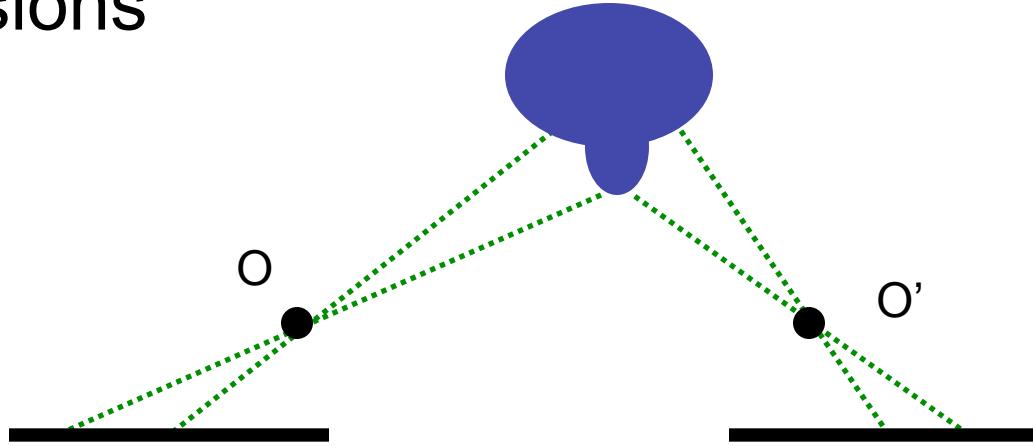
- Smaller window
 - More detail
 - More noise
- Larger window
 - Smoother disparity maps
 - Less prone to noise

Issues

- Fore shortening effect

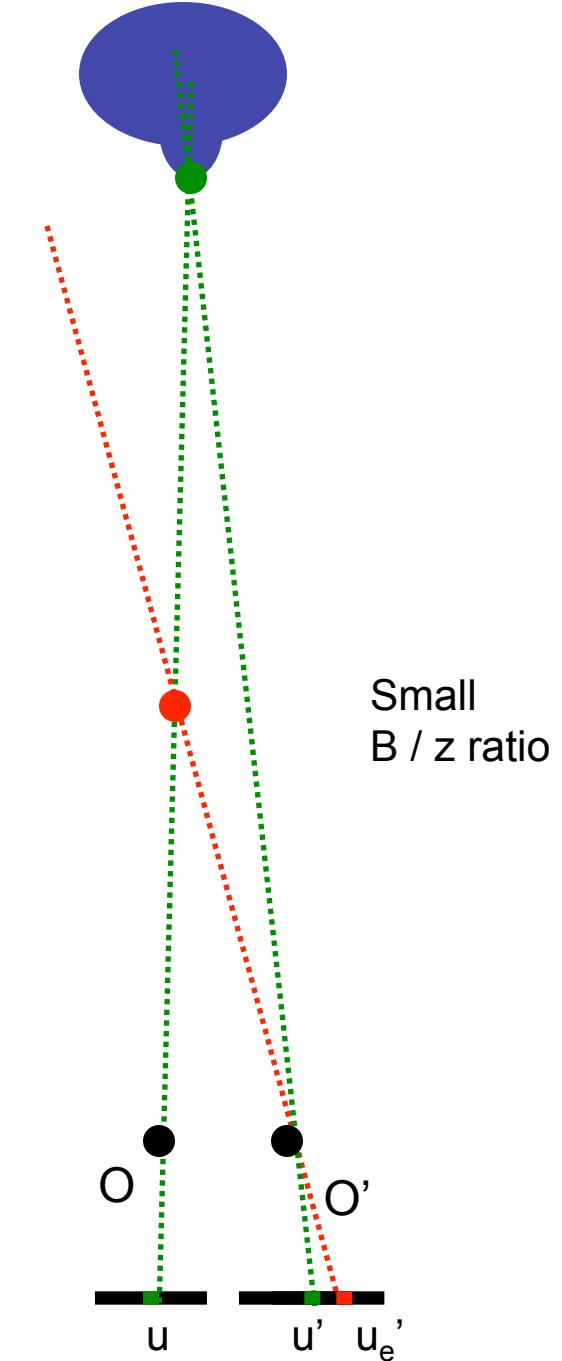
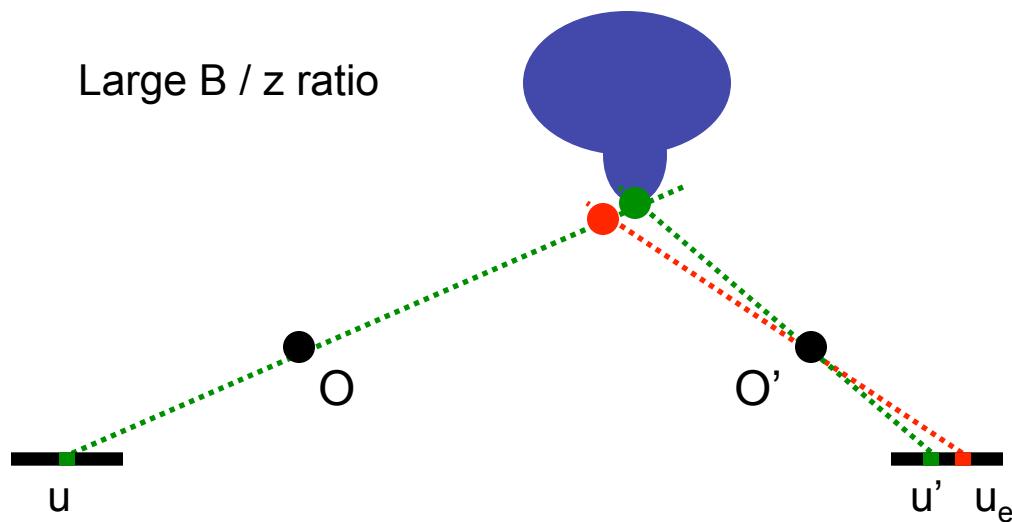


- Occlusions



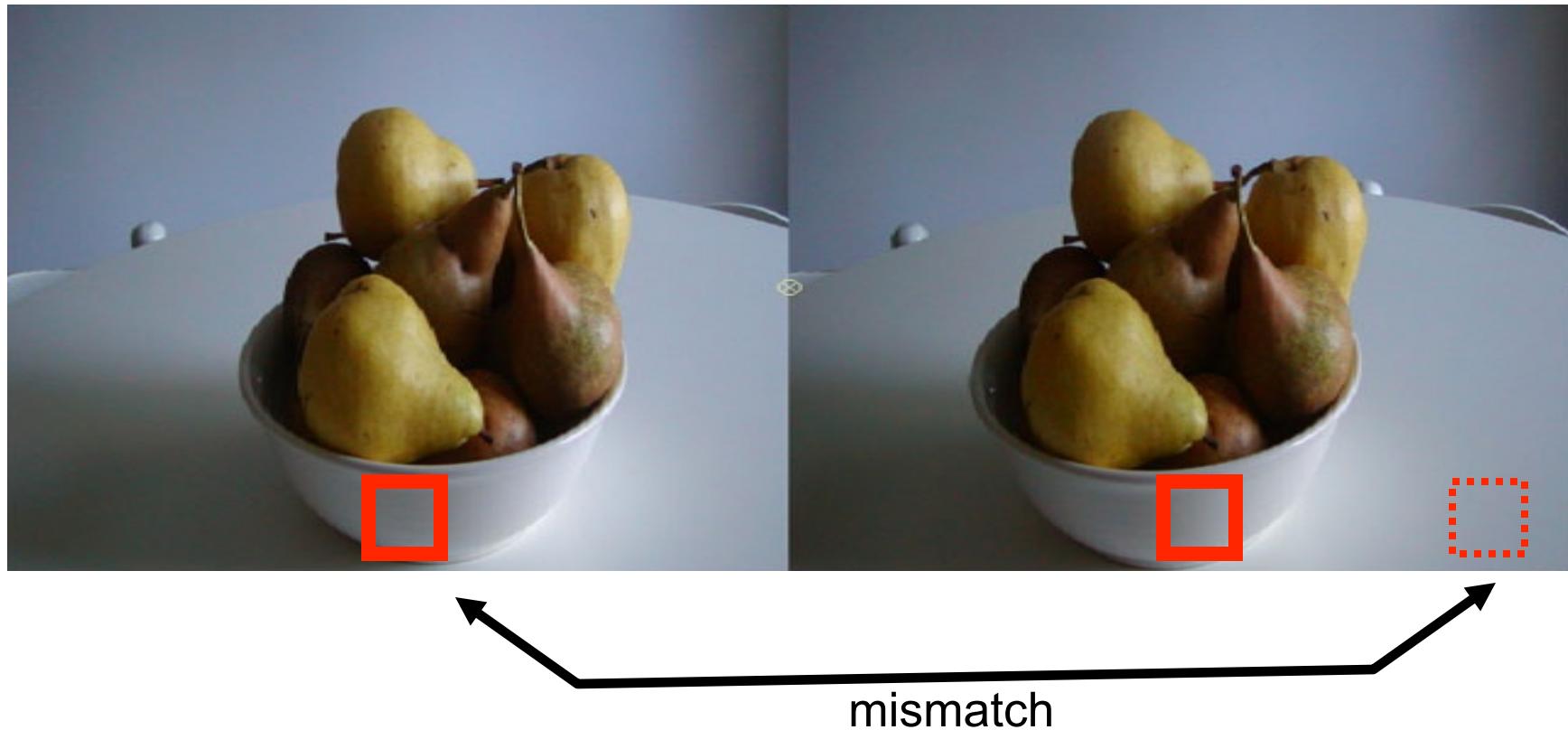
Issues

- To reduce the effect of foreshortening and occlusions, it is desirable to have small B / z ratio!
- However, when B/z is small, small errors in measurements imply large error in estimating depth



Issues

- Homogeneous regions



Issues

- Repetitive patterns



Correspondence problem is difficult!

- Occlusions
- Fore shortening
- Baseline trade-off
- Homogeneous regions
- Repetitive patterns

Apply non-local constraints to help
enforce the correspondences

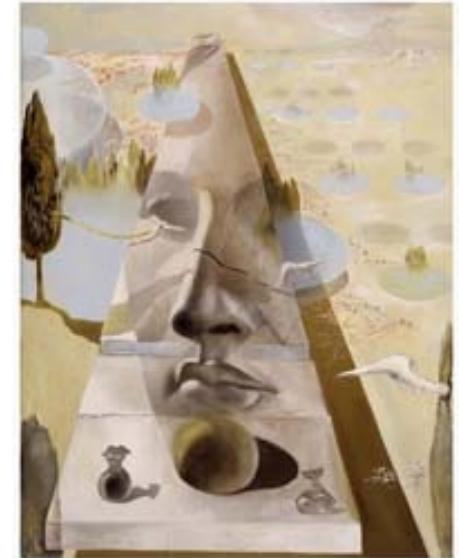
Non-local constraints

- Uniqueness
 - For any point in one image, there should be at most one matching point in the other image
- Ordering
 - Corresponding points should be in the same order in both views
- Smoothness
 - Disparity is typically a smooth function of x (except in occluding boundaries)

Lecture 6

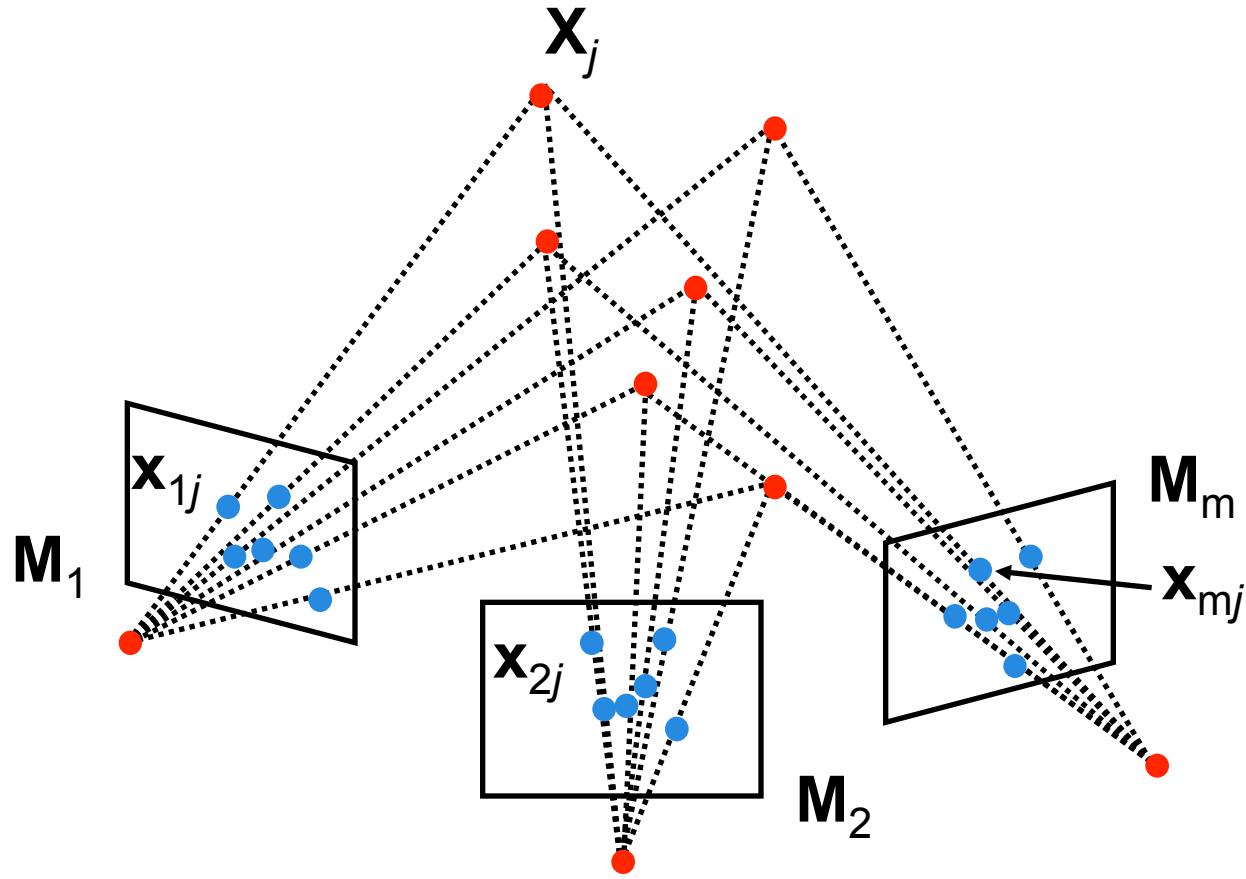
Stereo Systems

Multi-view geometry



- Stereo systems
 - Rectification
 - Correspondence problem
- Multi-view geometry
 - The SFM problem
 - Affine SFM

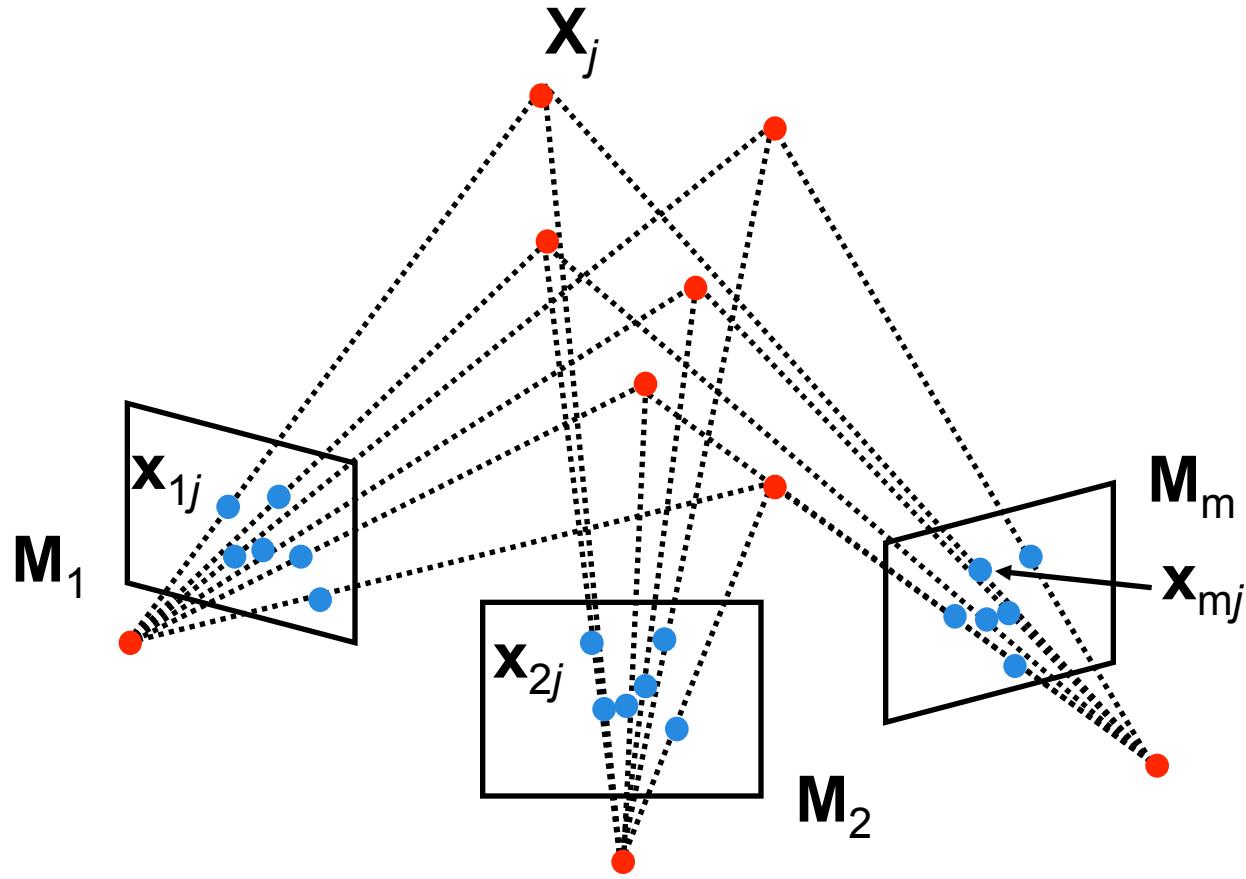
Structure from motion problem



Given m images of n fixed 3D points

$$\bullet \mathbf{x}_{ij} = \mathbf{M}_i \mathbf{X}_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$

Structure from motion problem

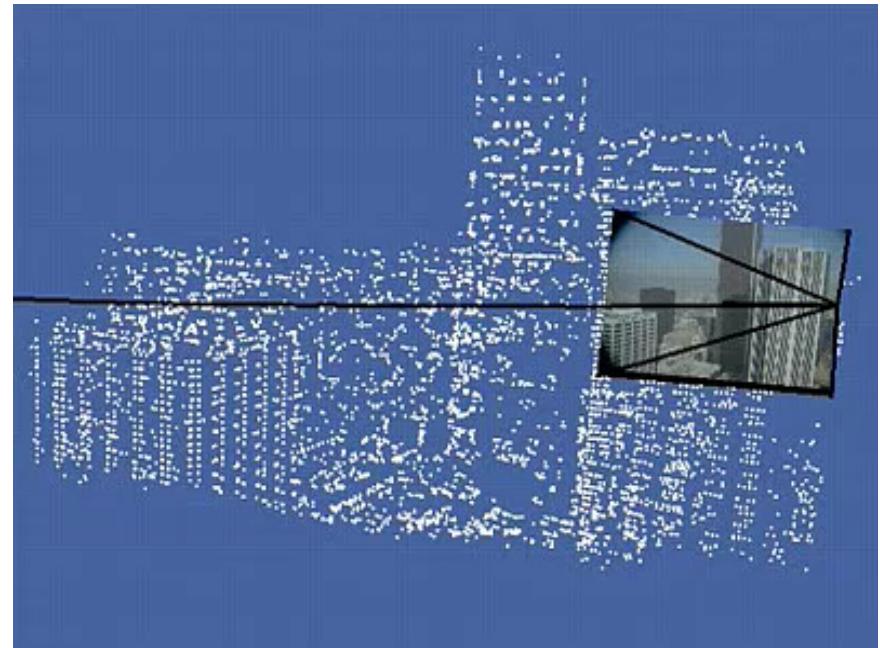


From the $m \times n$ observations \mathbf{x}_{ij} , estimate:

- m projection matrices \mathbf{M}_i
- n 3D points \mathbf{X}_j

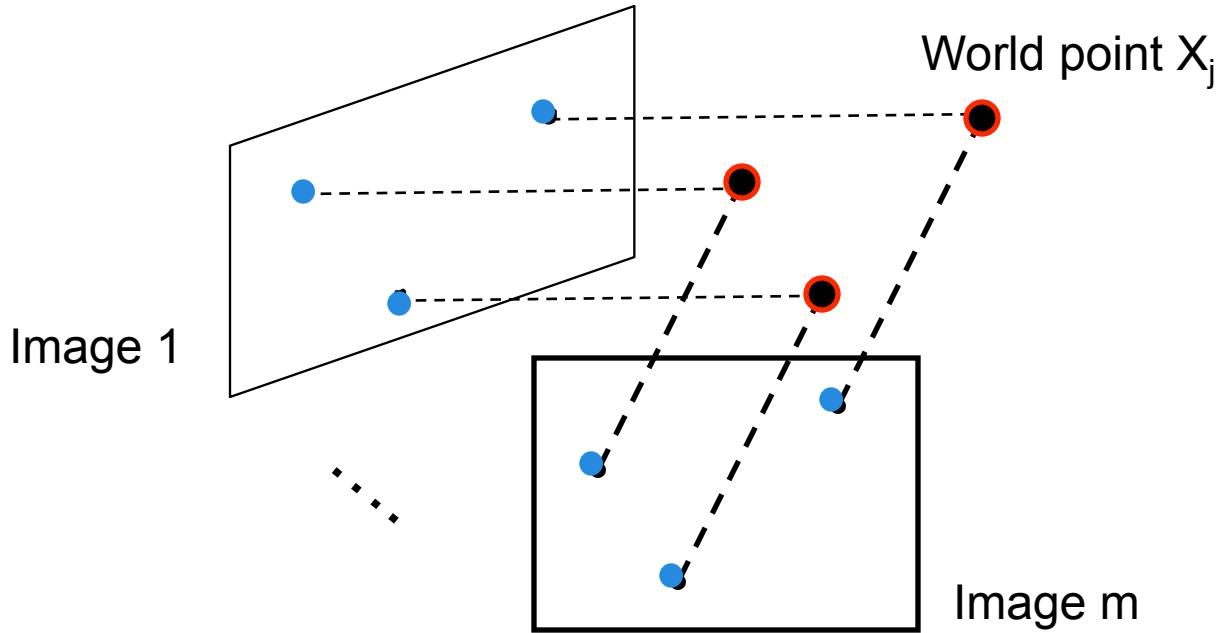
motion
structure

Structure from motion problem



Courtesy of Oxford **Visual Geometry Group**

Affine structure from motion (simpler problem)



From the $m \times n$ observations \mathbf{x}_{ij} , estimate:

- m projection matrices \mathbf{M}_i (affine cameras)
- n 3D points \mathbf{X}_j

Perspective

$$\mathbf{X} = M \mathbf{X} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix} \mathbf{X} = \begin{bmatrix} \mathbf{m}_1 \mathbf{X} \\ \mathbf{m}_2 \mathbf{X} \\ \mathbf{m}_3 \mathbf{X} \end{bmatrix}$$

$$M = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{v} & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix}$$

$$\mathbf{x}^E = \left(\frac{\mathbf{m}_1 \mathbf{X}}{\mathbf{m}_3 \mathbf{X}}, \frac{\mathbf{m}_2 \mathbf{X}}{\mathbf{m}_3 \mathbf{X}} \right)$$

Affine

$$\mathbf{X} = M \mathbf{X} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix} \mathbf{X} = \begin{bmatrix} \mathbf{m}_1 \mathbf{X} \\ \mathbf{m}_2 \mathbf{X} \\ 1 \end{bmatrix}$$

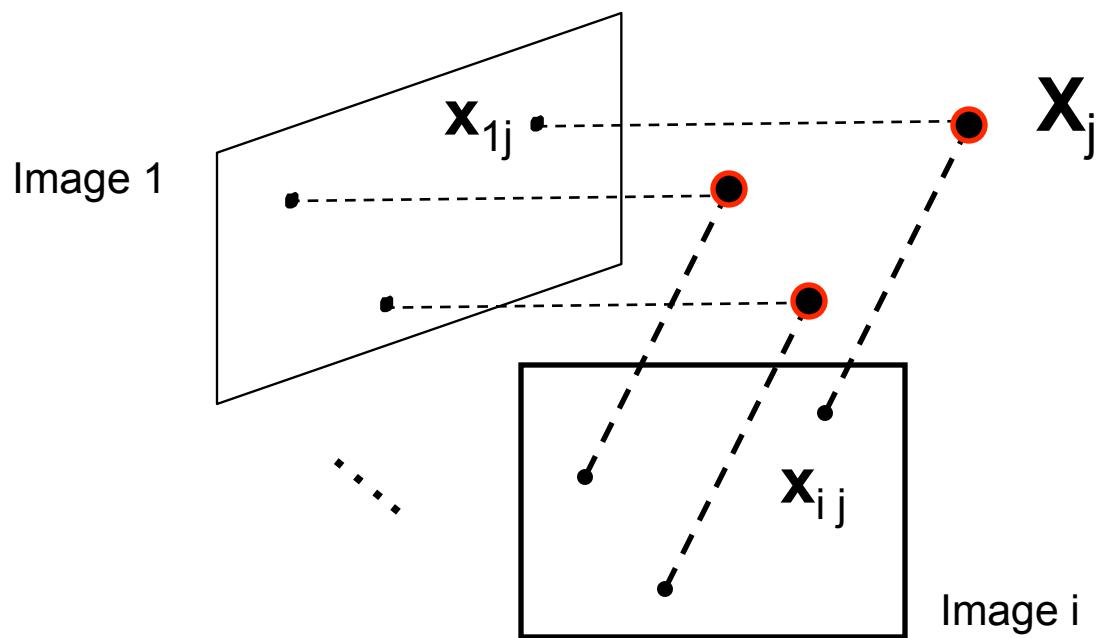
$$M = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{2 \times 3} & \mathbf{b}_{2 \times 1} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\mathbf{x}^E = (\mathbf{m}_1 \mathbf{X}, \mathbf{m}_2 \mathbf{X}) = \begin{bmatrix} \mathbf{A} & \mathbf{b} \end{bmatrix} \mathbf{X} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \mathbf{AX}^E + \mathbf{b}$$

↑ ↑
magnification [Eq. 3]

$$\mathbf{X}^E = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

Affine cameras



Camera matrix \mathbf{M}^a for the affine case (in Euclidean space)

$$\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i = \mathbf{M}_i^a \begin{bmatrix} \mathbf{X}_j \\ 1 \end{bmatrix}; \quad \mathbf{M}^a = \begin{bmatrix} \mathbf{A} & \mathbf{b} \end{bmatrix}$$

[Eq. 4]

The Affine Structure-from-Motion Problem

Given m images of n fixed points \mathbf{X}_j we can write

$$\mathbf{x}_{ij} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i \quad \text{for } i = 1, \dots, m \quad \text{and } j = 1, \dots, n$$

N. of cameras N. of points

Problem: estimate m matrices \mathbf{A}_i , m matrices \mathbf{b}_i ,
and the n positions \mathbf{X}_j from the $m \times n$ observations \mathbf{x}_{ij} .

How many equations and how many unknowns?

$2m \times n$ equations in $8m+3n$ unknowns

With $m=2$ cameras, I need at least $n=16$ 3D points
With $m=3$ cameras, I need at least $n=8$ 3D points

Two approaches:

- Algebraic approach (affine epipolar geometry; estimate F ; cameras; points)
- Factorization method

A factorization method – Tomasi & Kanade algorithm

C. Tomasi and T. Kanade

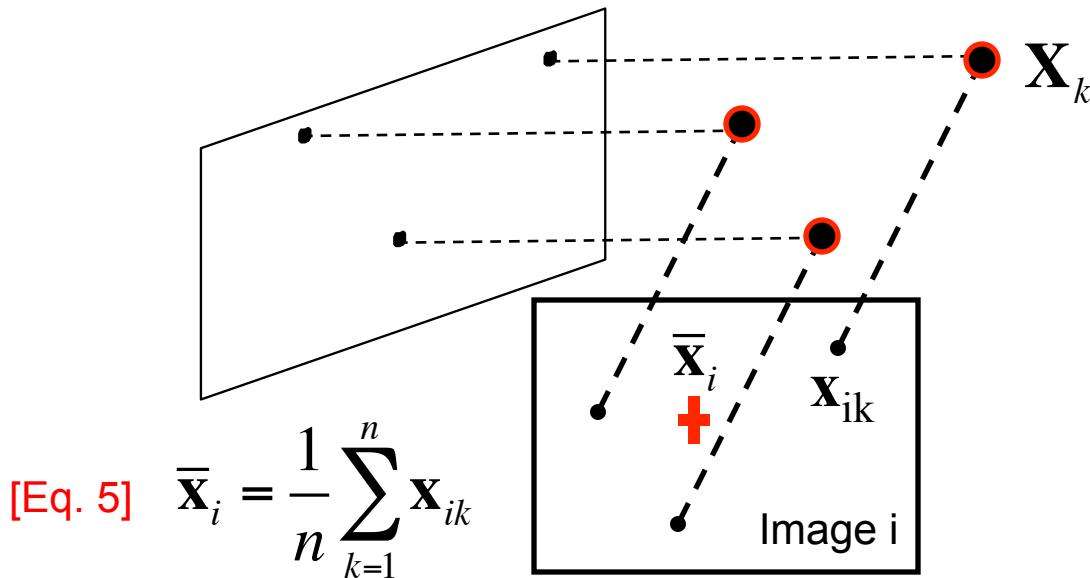
[Shape and motion from image streams under orthography: A factorization method.](#) *IJCV*, 9(2):137-154,
November 1992.

- Data centering
- Factorization

A factorization method - Centering the data

Centering: subtract the centroid of the image points

$$[\text{Eq. 6}] \quad \hat{\mathbf{x}}_{ij} = \mathbf{x}_{ij} - \boxed{\frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik}} \quad \bar{\mathbf{x}}_i$$



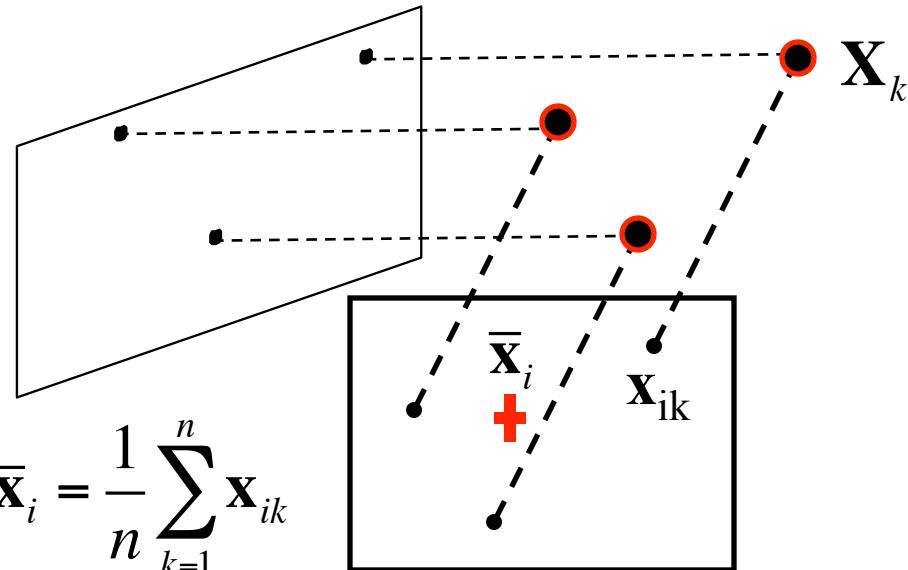
A factorization method - Centering the data

Centering: subtract the centroid of the image points

$$[\text{Eq. 6}] \quad \hat{\mathbf{x}}_{ij} = \mathbf{x}_{ij} - \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i - \frac{1}{n} \sum_{k=1}^n \mathbf{A}_i \mathbf{X}_k - \frac{1}{n} \sum_{k=1}^n \mathbf{b}_i$$

$$\mathbf{x}_{ik} = \mathbf{A}_i \mathbf{X}_k + \mathbf{b}_i$$

[Eq. 4]

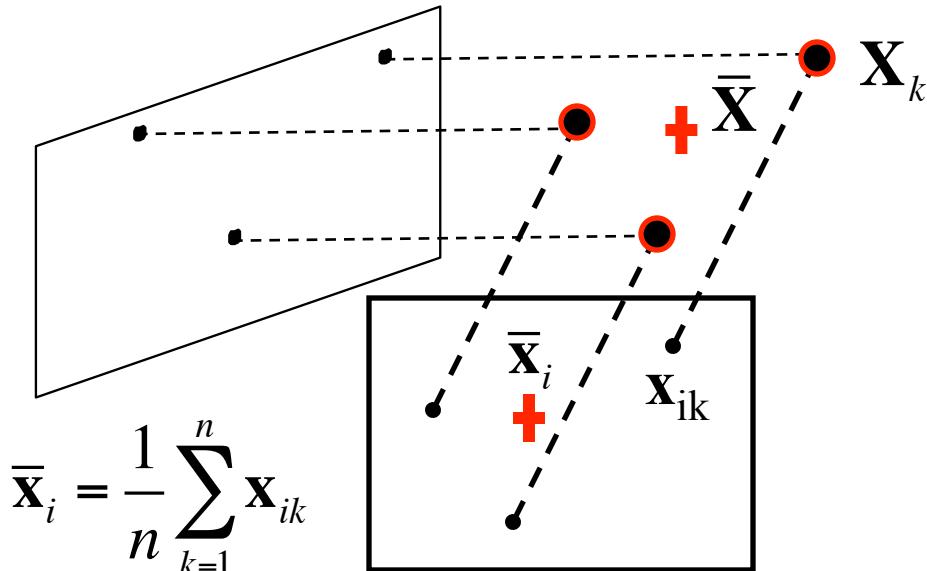


$$[\text{Eq. 5}] \quad \bar{\mathbf{x}}_i = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik}$$

A factorization method - Centering the data

Centering: subtract the centroid of the image points

$$\begin{aligned}\hat{\mathbf{x}}_{ij} &= \mathbf{x}_{ij} - \frac{1}{n} \sum_{k=1}^n \mathbf{x}_{ik} = \mathbf{A}_i \mathbf{X}_j + \mathbf{b}_i - \frac{1}{n} \sum_{k=1}^n \mathbf{A}_i \mathbf{X}_k - \frac{1}{n} \sum_{k=1}^n \mathbf{b}_i \\ \mathbf{x}_{ik} &= \mathbf{A}_i \mathbf{X}_k + \mathbf{b}_i \\ &= \mathbf{A}_i \left(\mathbf{X}_j - \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \right) = \mathbf{A}_i (\mathbf{X}_j - \bar{\mathbf{X}}) \\ &= \mathbf{A}_i \hat{\mathbf{X}}_j\end{aligned}\quad [\text{Eq. 8}]$$



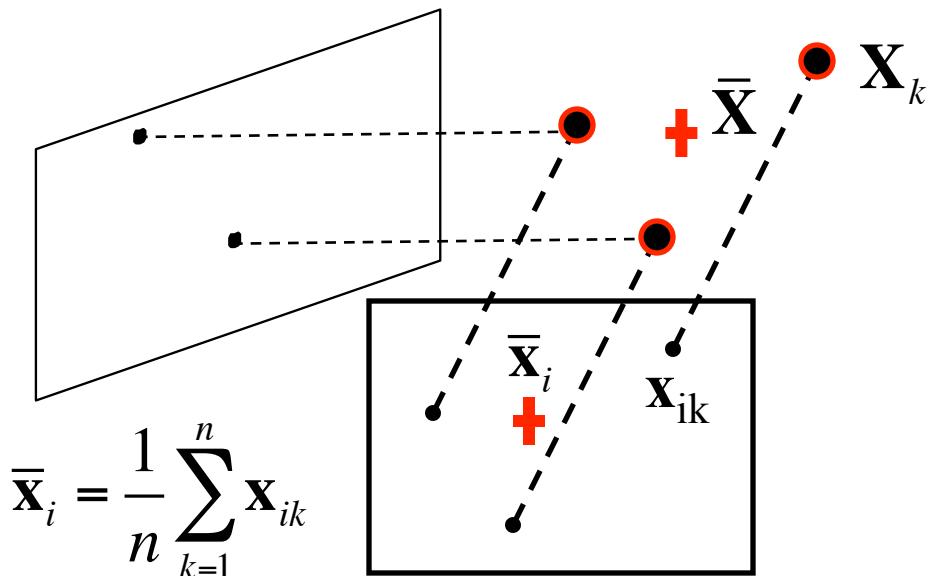
$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \quad [\text{Eq. 7}]$$

Centroid of 3D points

A factorization method - Normalizing the data

Thus, after centering, each **normalized** observed point is related to the 3D point by

$$\hat{\mathbf{X}}_{ij} = \mathbf{A}_i \hat{\mathbf{X}}_j \quad [\text{Eq. 8}]$$



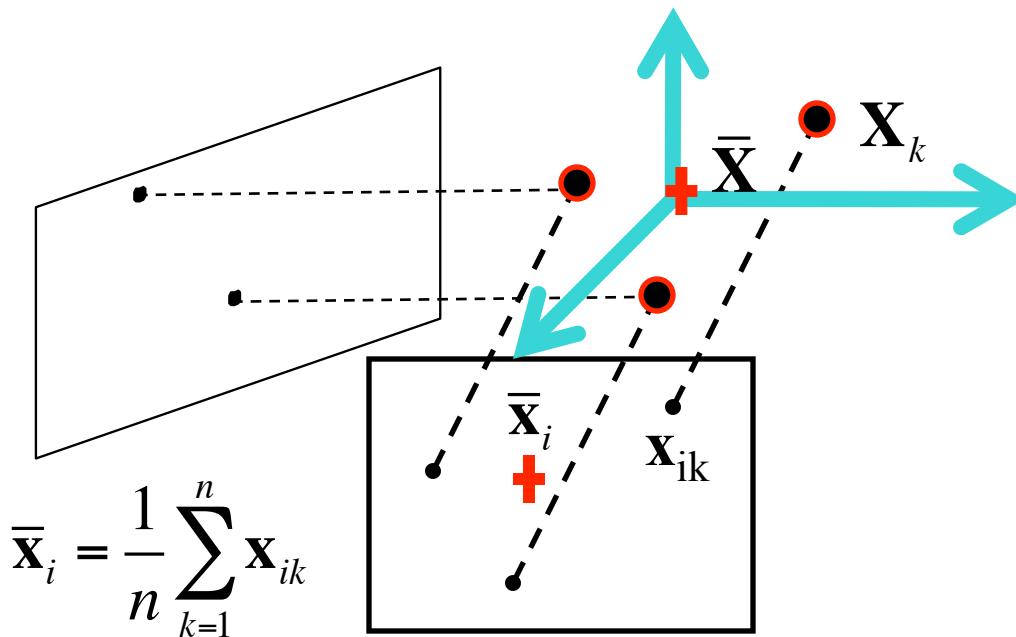
$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \quad [\text{Eq. 7}]$$

Centroid of 3D points

A factorization method - Normalizing the data

If the centroid of points in 3D = center of the world reference system

$$\hat{\mathbf{X}}_{ij} = \mathbf{A}_i \hat{\mathbf{X}}_j = \mathbf{A}_i \mathbf{X}_j \quad [\text{Eq. 9}]$$



$$\bar{\mathbf{X}}_i = \frac{1}{n} \sum_{k=1}^n \mathbf{X}_{ik}$$

$$\bar{\mathbf{X}} = \frac{1}{n} \sum_{k=1}^n \mathbf{X}_k \quad [\text{Eq. 7}]$$

Centroid of 3D points

A factorization method - factorization

Let's create a $2m \times n$ data (measurement) matrix:

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ & & \ddots & \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix}$$

cameras
($2m$)

points (n)

Each $\hat{\mathbf{x}}_{ij}$ entry is a 2×1 vector!

A factorization method - factorization

Let's create a $2m \times n$ data (measurement) matrix:

$$\mathbf{D} = \begin{bmatrix} \hat{\mathbf{x}}_{11} & \hat{\mathbf{x}}_{12} & \cdots & \hat{\mathbf{x}}_{1n} \\ \hat{\mathbf{x}}_{21} & \hat{\mathbf{x}}_{22} & \cdots & \hat{\mathbf{x}}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\mathbf{x}}_{m1} & \hat{\mathbf{x}}_{m2} & \cdots & \hat{\mathbf{x}}_{mn} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_n \end{bmatrix}$$

points ($3 \times n$)

cameras
($2m \times 3$)

M S

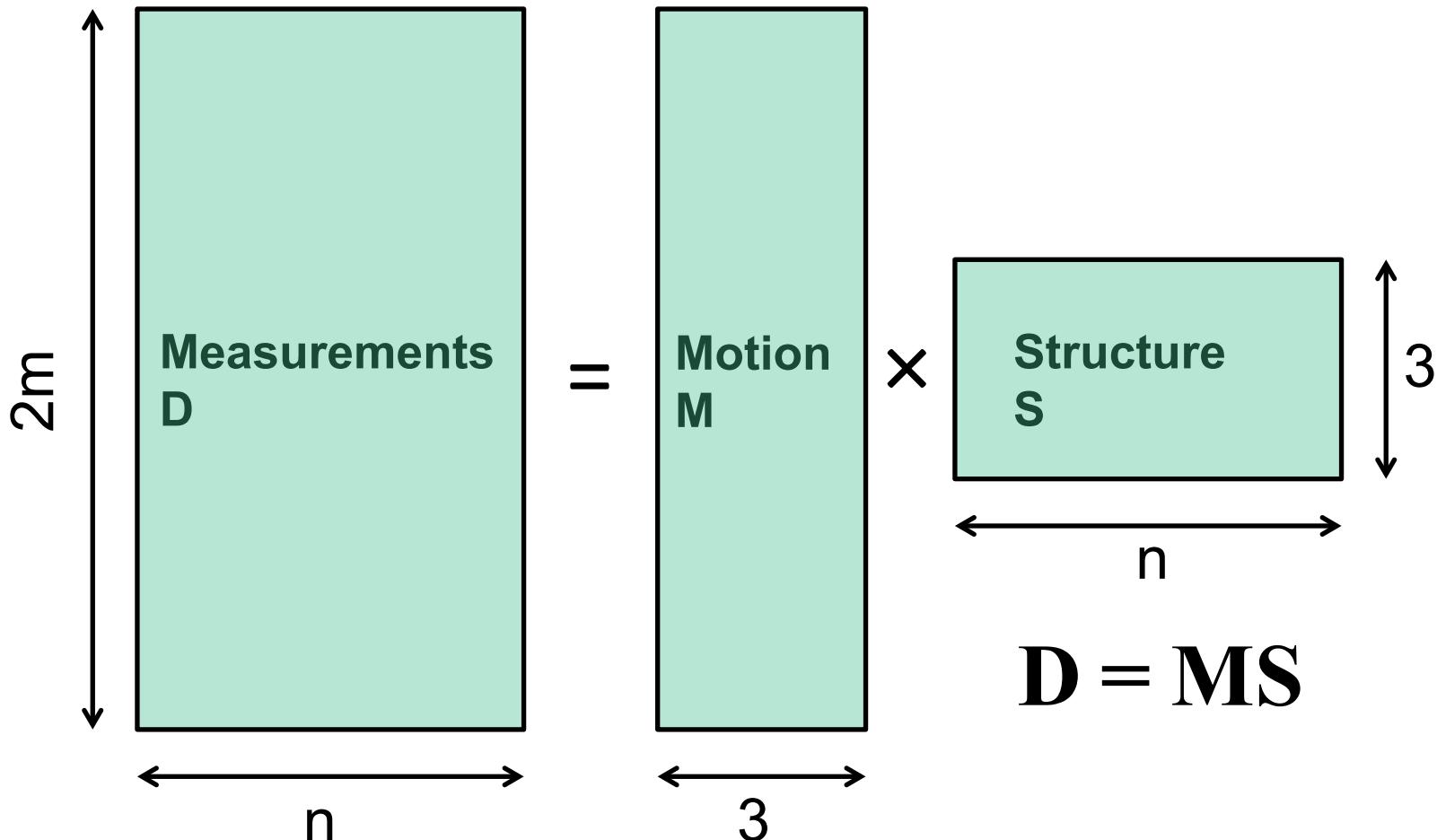
$(2m \times n)$

[Eq. 10]

Each $\hat{\mathbf{x}}_{ij}$ entry is a 2×1 vector!
 \mathbf{A}_i is 2×3 and \mathbf{X}_j is 3×1

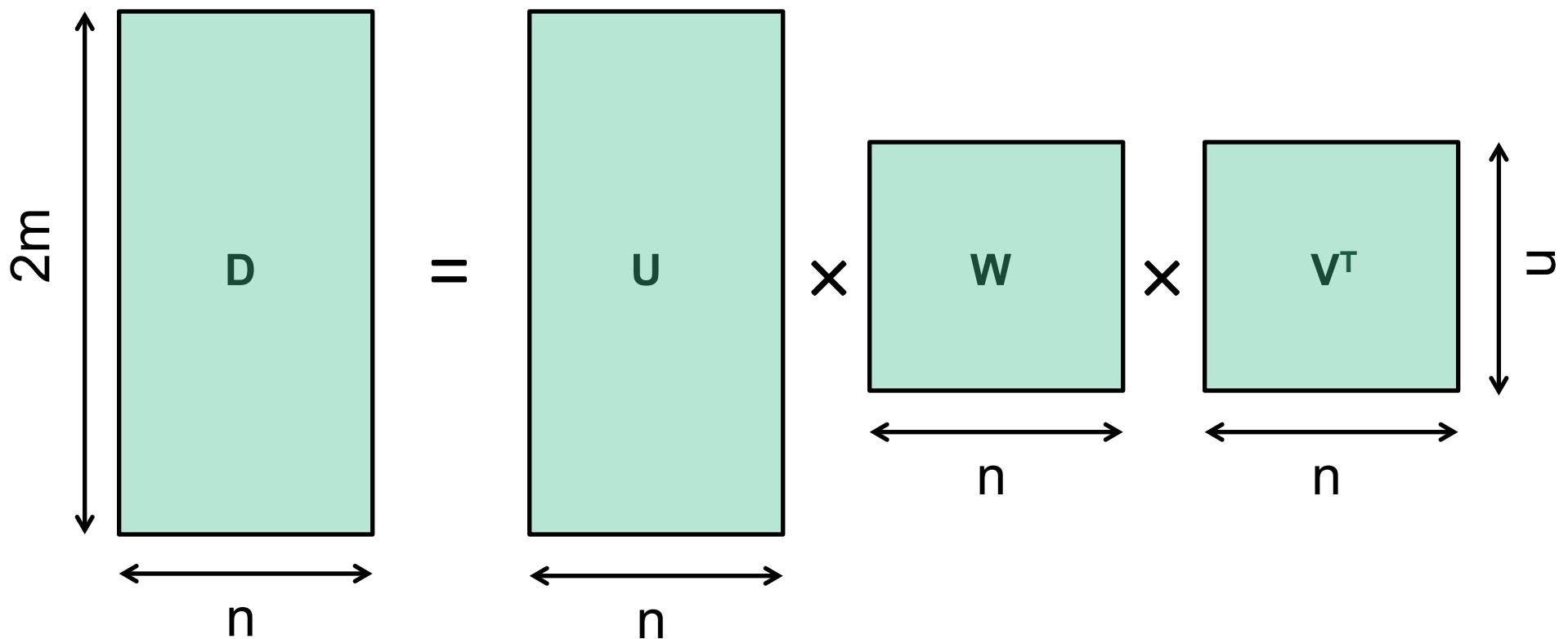
The measurement matrix $\mathbf{D} = \mathbf{M} \mathbf{S}$ has rank 3
(it's a product of a $2m \times 3$ matrix and $3 \times n$ matrix)

Factorizing the Measurement Matrix



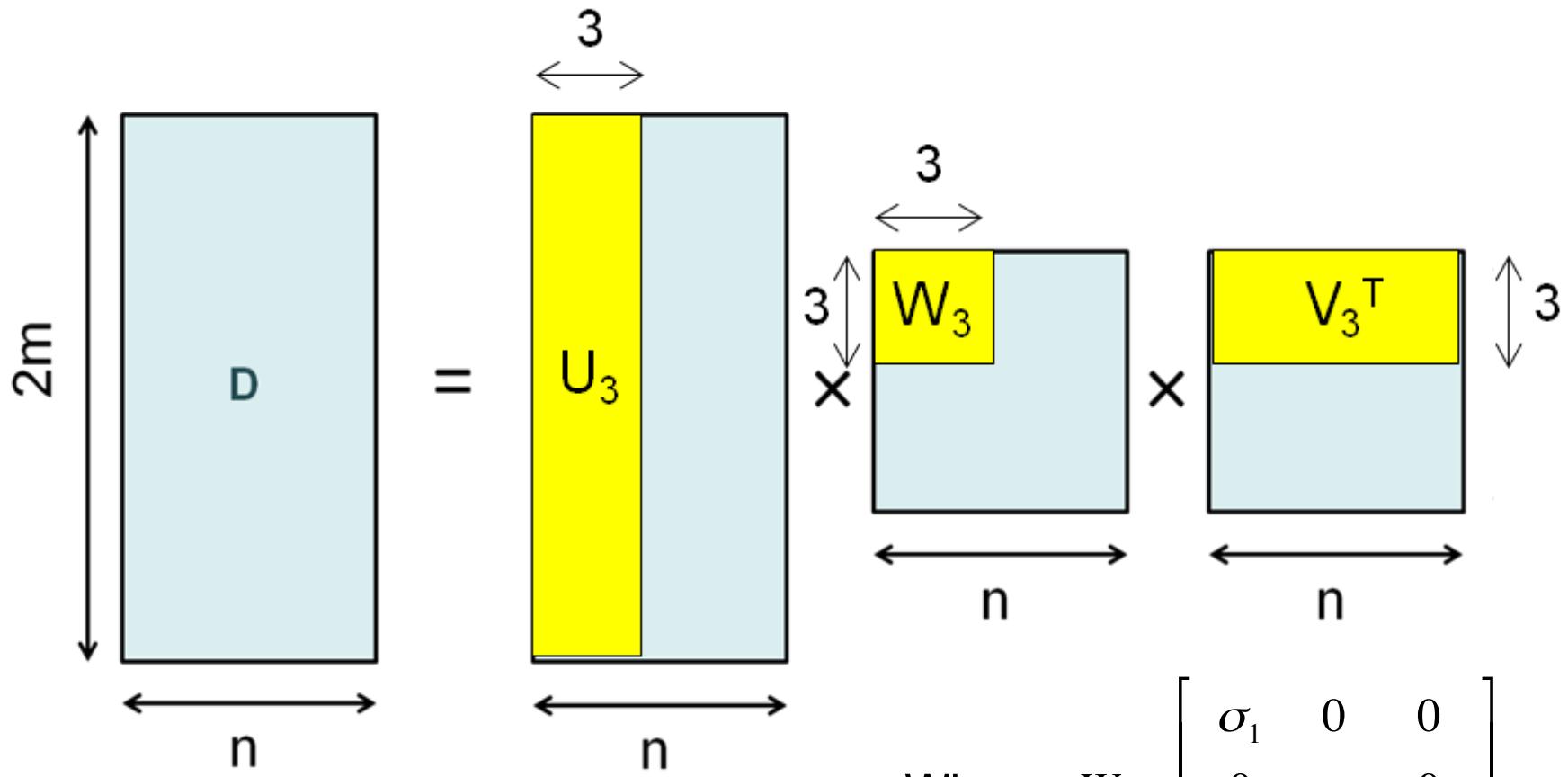
Factorizing the Measurement Matrix

- How to factorize D? By computing the Singular value decomposition of D!



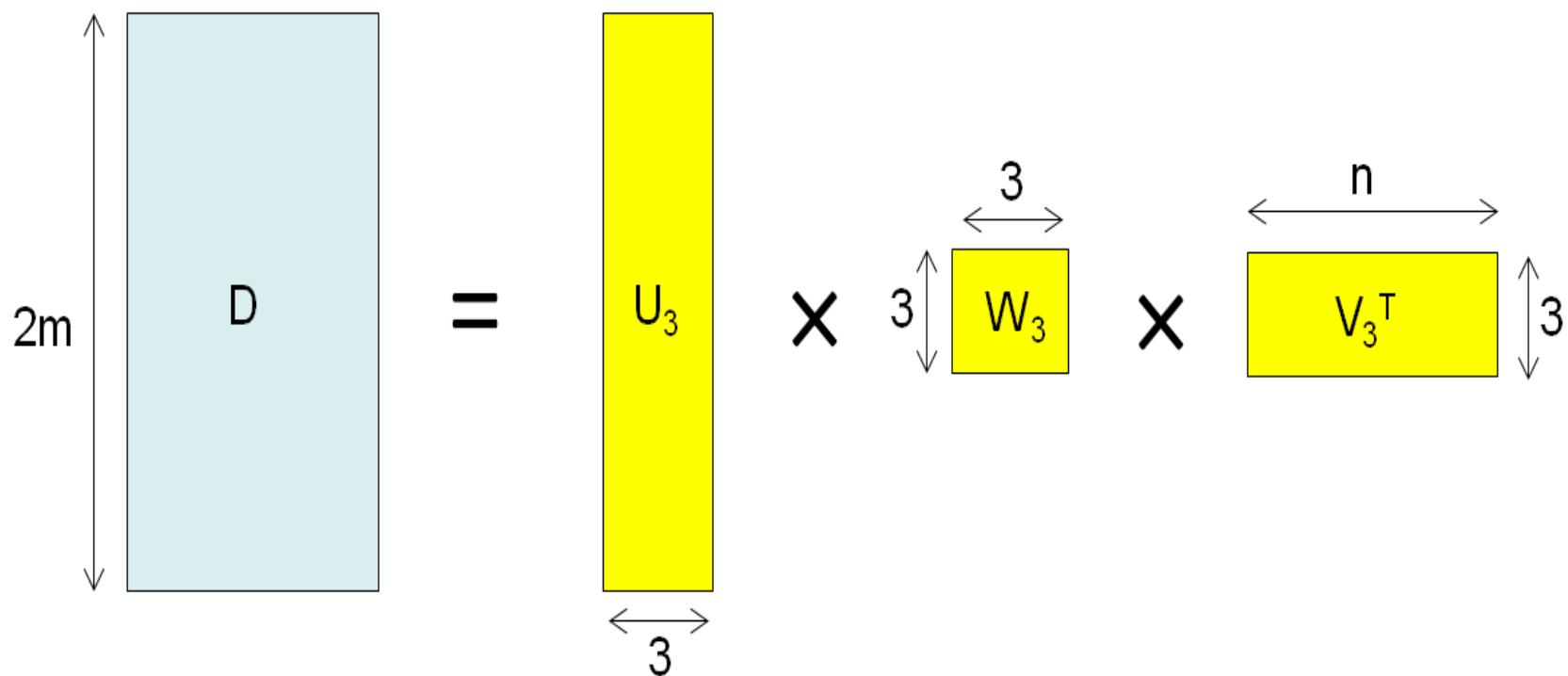
Factorizing the Measurement Matrix

Since rank (D)=3, there are only 3 non-zero singular values σ_1 , σ_2 and σ_3

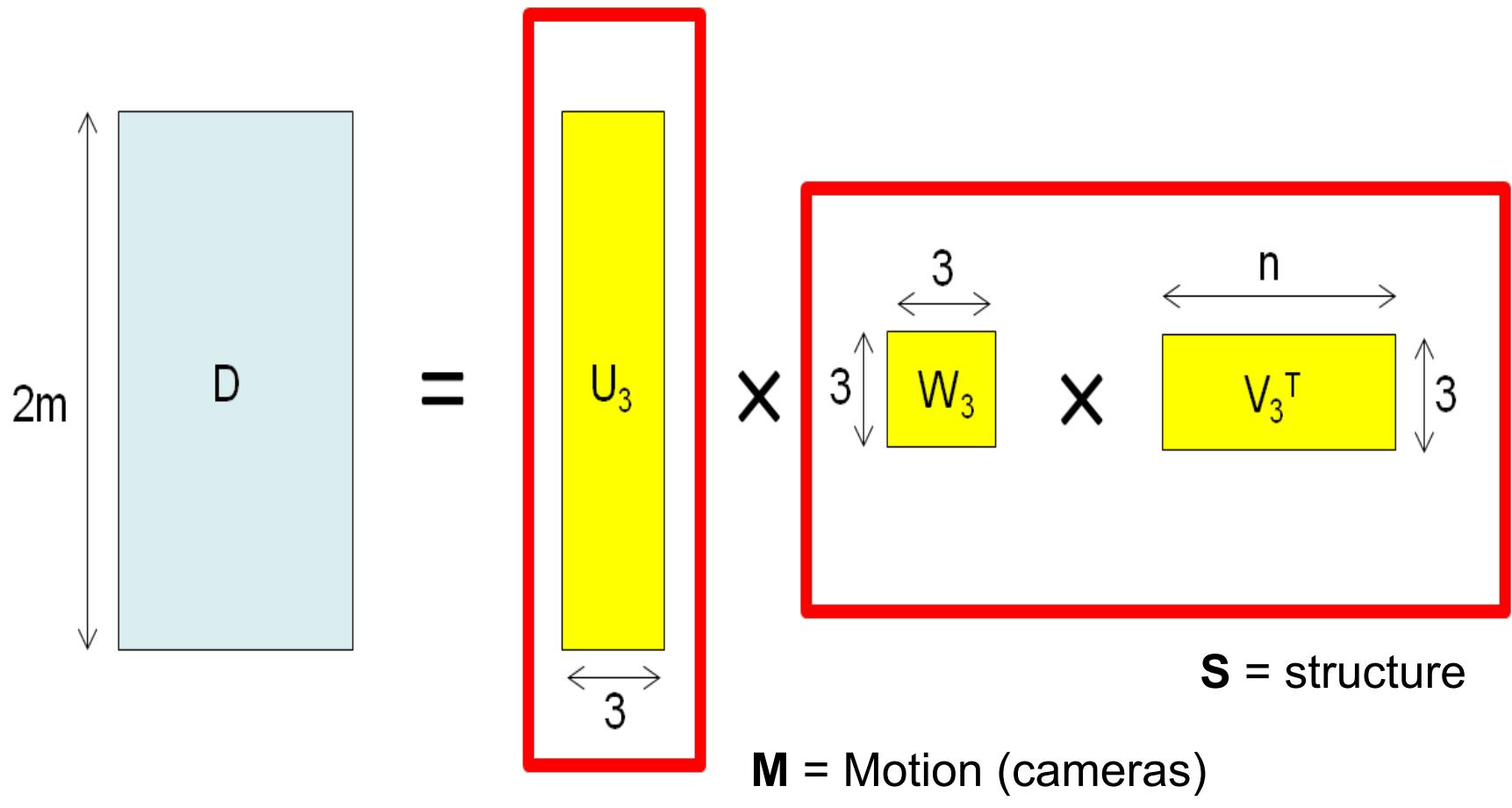


$$\text{Where } W_3 = \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix} \quad [\text{Eq. 11}]$$

Factorizing the Measurement Matrix



Factorizing the Measurement Matrix



$$D = U_3 W_3 V_3^T = U_3 (W_3 V_3^T) = M S \quad [\text{Eq. 12}]$$

Factorizing the Measurement Matrix

$$\mathbf{D} = \mathbf{U}_3 \mathbf{W}_3 \mathbf{V}_3^T = \mathbf{U}_3 (\mathbf{W}_3 \mathbf{V}_3^T) = \mathbf{M} \mathbf{S} \quad [\text{Eq. 12}]$$

What is the issue here? \mathbf{D} has rank>3 because of:

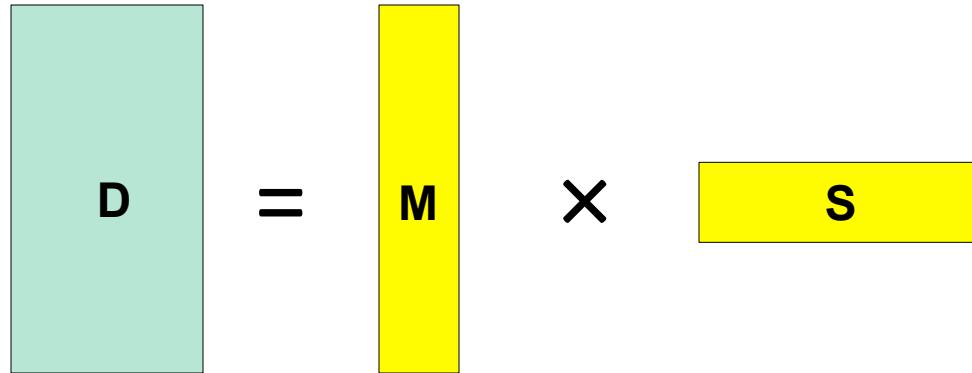
- measurement noise
- affine approximation

Theorem: When \mathbf{D} has a rank greater than 3, $\mathbf{U}_3 \mathbf{W}_3 \mathbf{V}_3^T$ is the best possible rank-3 approximation of \mathbf{D} in the sense of the Frobenius norm.

$$\mathbf{D} = \mathbf{U}_3 \mathbf{W}_3 \mathbf{V}_3^T \quad \left\{ \begin{array}{l} \mathbf{M} \approx \mathbf{U}_3 \\ \mathbf{S} \approx \mathbf{W}_3 \mathbf{V}_3^T \end{array} \right.$$

$$\|A\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} = \sqrt{\sum_{i=1}^{\min\{m, n\}} \sigma_i^2}$$

Affine Ambiguity

$$\begin{matrix} D \\ = \\ M \\ \times \\ S \end{matrix}$$


- The decomposition is not unique. We get the same D by using any 3×3 matrix C and applying the transformations $M \rightarrow MC$, $S \rightarrow C^{-1}S$
- We can enforce some Euclidean constraints to resolve this ambiguity (more on next lecture!)

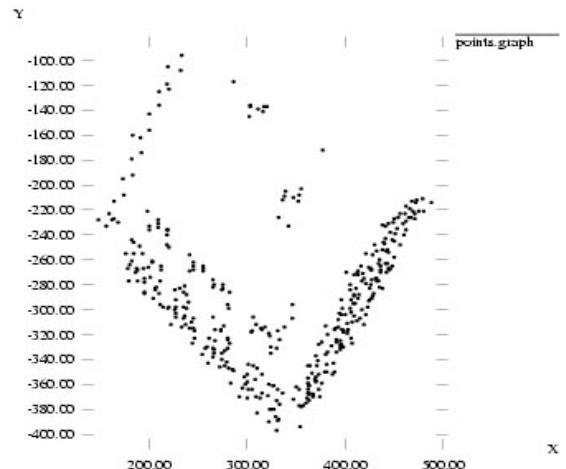
Reconstruction results



1



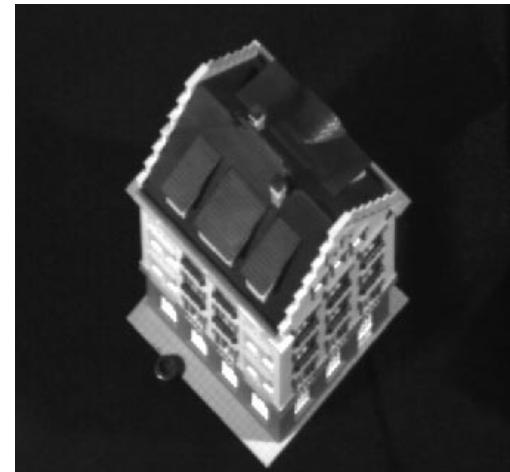
60



120



150



Next lecture

Multiple view geometry

Perspective structure from Motion