

# Structure from Motion

## Lecture-15

# Shape From X

- Recovery of 3D (shape) from one or two (2D images).

# Shape From X

- Stereo
- Motion
- Shading
- Photometric Stereo
- Texture
- Contours
- Silhouettes
- Defocus

# Applications

- Object Recognition
- Robotics
- Computer Graphics
- Image Retrieval
- Geo-localization
- Archeology
- Sports

# Microsoft Kinect sensor

- Data Captured using Microsoft Kinect sensor

RGB Camera

IR Camera 1



IR Camera 2

- Approximately 50,000 gesture samples

# Gesture Lexicons



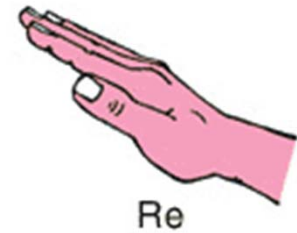
**Diving Signals**



**Referee Signals**



**Nurse Gesture**



**Music Notes**



**Gestures from Depth camera** ▲





▼ **Gestures from RGB camera**





## Test case 1: Torso motion adds noise (devel 01– 10 gestures)



 Instances of Successful Recognition  
 Instances of Failed Recognition

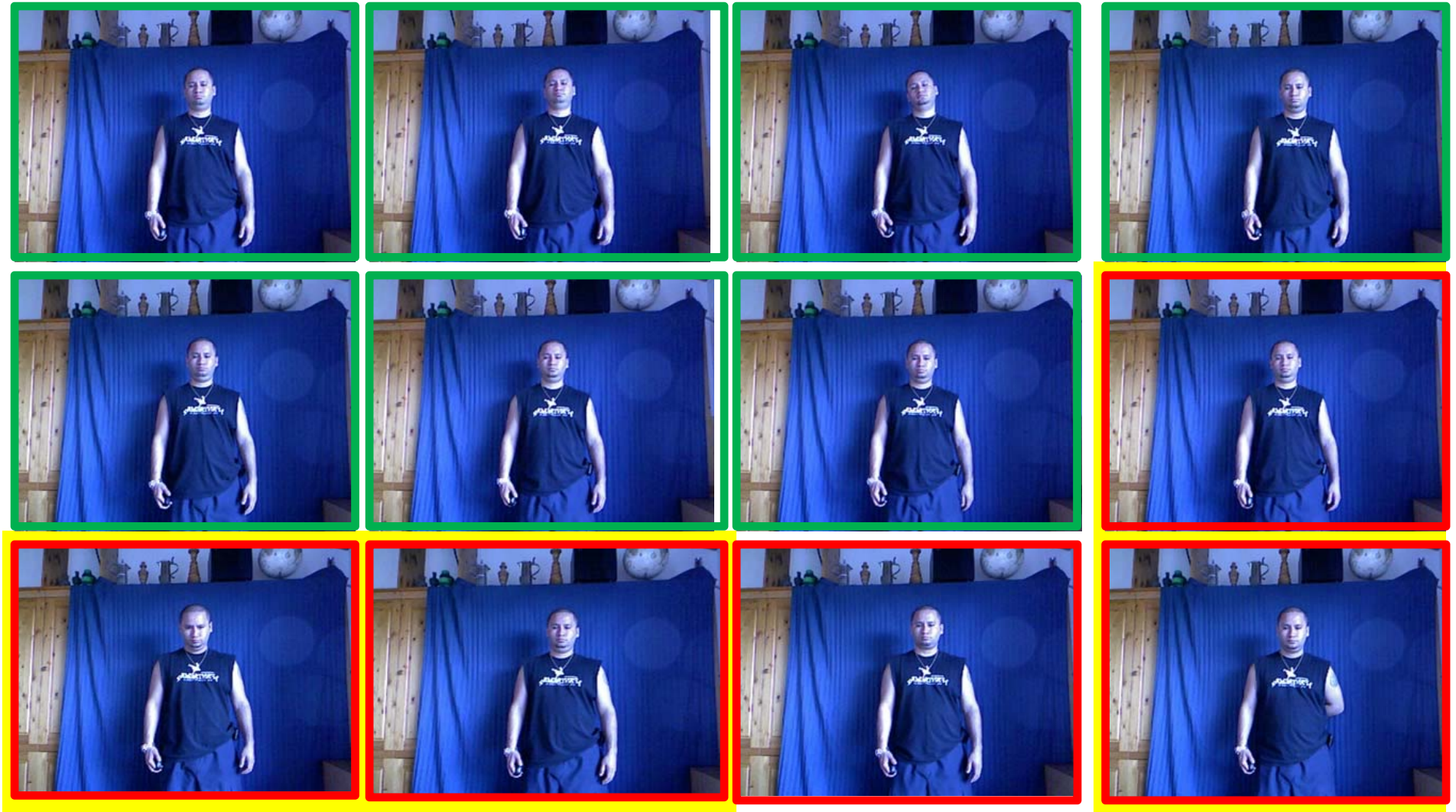
## Test Case 2: Improvisations (devel 06 – 9 gestures)





 Instances of Successful Recognition  
 Instances of Failed Recognition



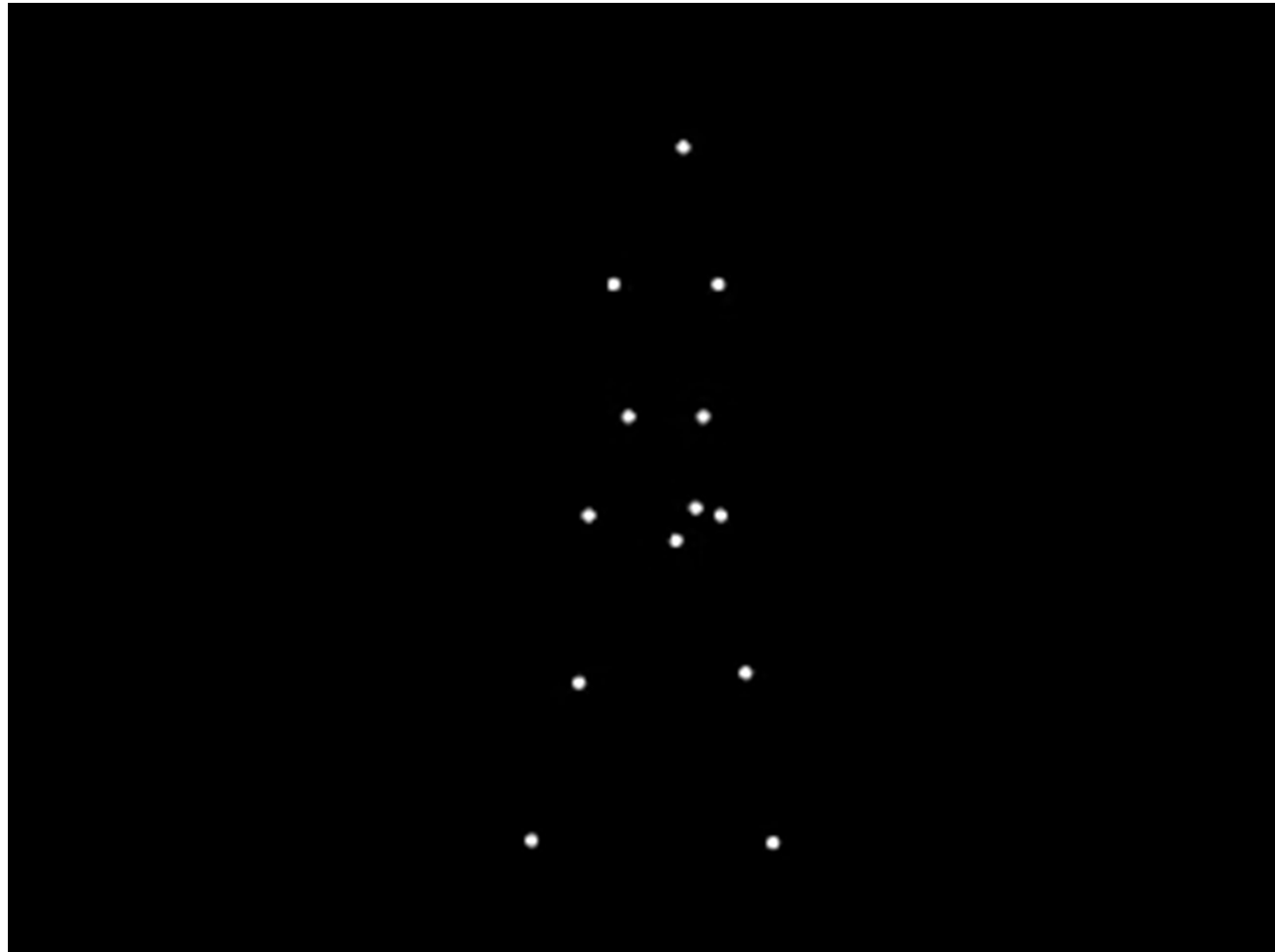
## Test Case 3: Subtle differences (devel 09 – 10 gestures)



-  Instances of Successful Recognition
-  Instances of Failed Recognition

# Moving Light Display

Humans are able  
to recover 3D  
from motion



# Shape from Motion



(a)



(b)



(c)



(d)

# Problem

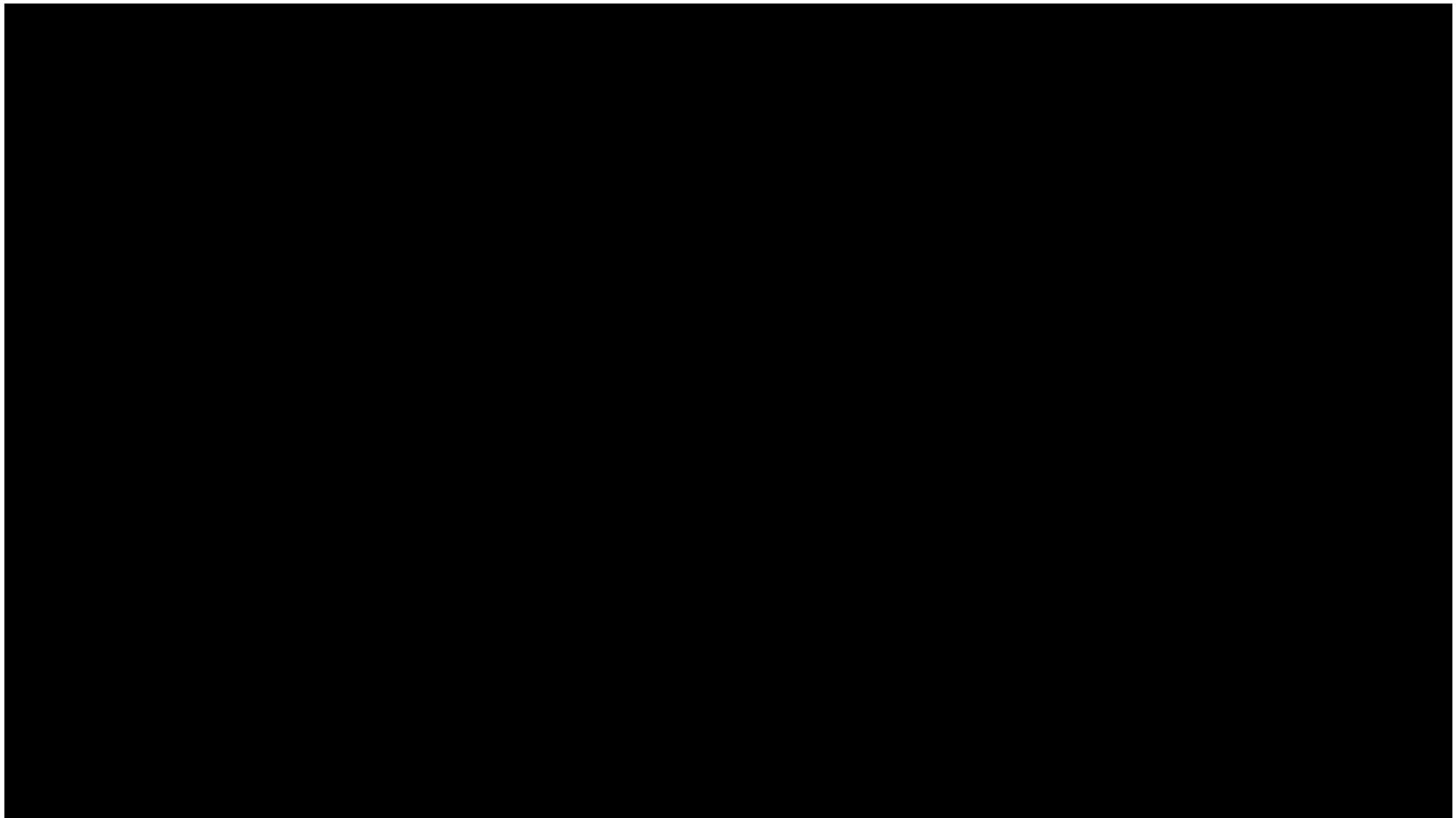
- Given optical flow or point correspondences, compute 3-D motion (translation and rotation) and shape (depth).

# Structure from Motion

- S. Ullman
- Hanson & Riseman
- Webb & Aggarwal
- T. Huang
- Heeger and Jepson
- Chellappa
- Faugeras
- Zisserman
- Kanade
- Pentland
- Van Gool
- Pollefeys
- Seitz & Szeliski
- Shahsua
- Irani
- Vidal & Yi Ma
- Medioni
- Fleet
- Tian & Shah
- -



# Photosynth



# Tomasi and Kanade Factorization Orthographic Projection

# Assumptions

- The camera model is orthographic.
- The positions of “P” points in “F” frames ( $F \geq 3$ ), which are not all coplanar, and have been tracked.
- The entire sequence has been acquired before starting (batch mode).
- Camera calibration not needed, if we accept 3D points up to a scale factor.

# Input



1



100



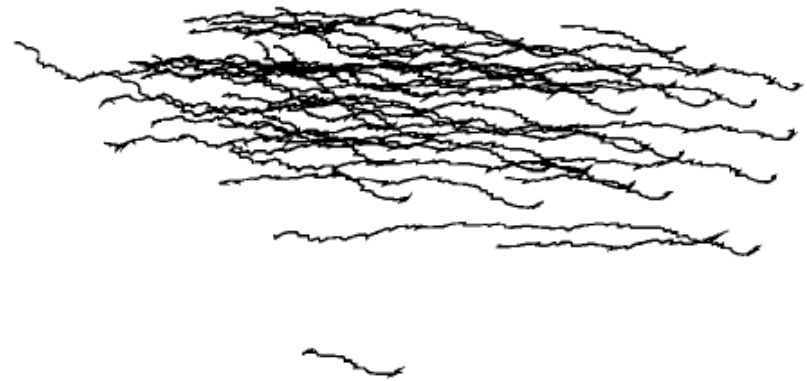
150



210

(a)

Images



(b)

KLT Tracks

# Feature Points

Image points  
(This is not optical flow)  $\{(u_{fp}, v_{fp}) \mid f = 1, \dots, F, p = 1, \dots, P\}$

$$W = \begin{bmatrix} u_{11} \dots u_{1P} \\ \vdots \\ u_{F1} \dots u_{FP} \\ v_{11} \dots v_{1P} \\ \vdots \\ v_{F1} \dots v_{FP} \end{bmatrix} \qquad W = \begin{bmatrix} U \\ - \\ V \end{bmatrix}$$



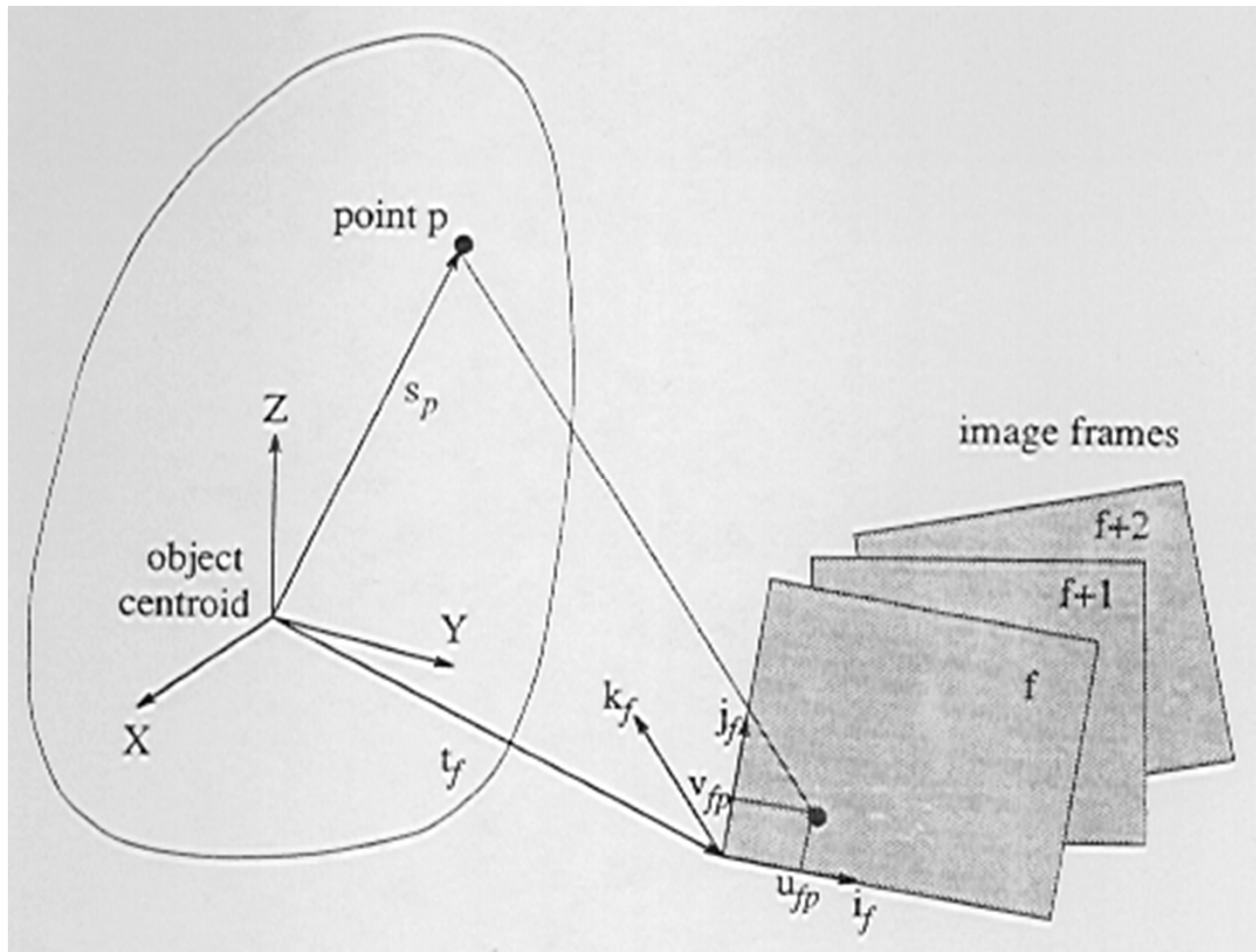
# Mean Normalize Feature Points

$$a_f = \frac{1}{P} \sum_{p=1}^P u_p \qquad b_f = \frac{1}{P} \sum_{p=1}^P v_p$$

$$\tilde{u}_{fP} = u_{fP} - a_{fP} \quad (\text{A})$$

$$\tilde{v}_{fP} = v_{fP} - b_{fP}$$

# Orthographic Projection



Copyright Mubarak Shah 2003

# Orthographic Projection

$$s_p = (X_p, Y_p, Z_p)$$

3D world point

$$u_{fP} = i_f^T (s_P - t_f) \quad (\text{C})$$

Orthographic projection

$$v_{fP} = j_f^T (s_P - t_f)$$

$$k_f = i_f \times j_f \quad i, j, k \text{ are unit vectors along } X, Y, Z$$

$$\tilde{u}_{fp} = u_{fp} - a_f \quad a_f = \frac{1}{P} \sum_{p=1}^P u_p$$

$$= i_f^T (s_p - t_f) -$$

$$= i_f^T \left[ s_p - \frac{1}{P} \sum_{q=1}^P s_q \right]$$

$$= i_f^T s_p$$

Copyright Mubarak Shah 2003

If Origin of world is at the  
centroid of object points,  
Second term is zero

$$\tilde{\mathbf{u}}_{fP} = \dot{\mathbf{i}}_f^T \mathbf{S}_P$$

$$\tilde{\mathbf{v}}_{fP} = \dot{\mathbf{j}}_f^T \mathbf{S}_P$$

$$\tilde{\mathbf{W}} = \begin{bmatrix} \tilde{\mathbf{U}} \\ - \\ \tilde{\mathbf{V}} \end{bmatrix}$$



$$\tilde{\mathbf{u}}_{fP} = \dot{\mathbf{i}}_f^T \mathbf{S}_P \quad (\mathbf{B})$$

$$\tilde{\mathbf{v}}_{fP} = \dot{\mathbf{j}}_f^T \mathbf{S}_P$$

$$\tilde{\mathbf{W}} = \begin{bmatrix} \tilde{\mathbf{U}} \\ - \\ \tilde{\mathbf{V}} \end{bmatrix}$$

$$\tilde{\mathbf{W}} = \begin{bmatrix} \tilde{\mathbf{u}}_{11} \dots \tilde{\mathbf{u}}_{1p} \\ \vdots \\ \tilde{\mathbf{u}}_{F1} \dots \tilde{\mathbf{u}}_{FP} \\ \tilde{\mathbf{v}}_{11} \dots \tilde{\mathbf{v}}_{1p} \\ \vdots \\ \tilde{\mathbf{v}}_{F1} \dots \tilde{\mathbf{v}}_{FP} \end{bmatrix}$$

$$\tilde{\mathbf{W}} = \begin{bmatrix} \dot{\mathbf{i}}_1^T \\ \vdots \\ \dot{\mathbf{i}}_F^T \\ \dot{\mathbf{j}}_1^T \\ \vdots \\ \dot{\mathbf{j}}_F^T \end{bmatrix} \begin{bmatrix} s_1 & \dots & s_P \end{bmatrix} = \mathbf{R}\mathbf{S}$$

**3XP**

**2FX3**

Rank of  $\mathbf{S}$  is 3, because points in 3D space are not  
Co-planar

# Rank Theorem

Without noise, the registered measurement matrix  $\tilde{W}$  is at most of rank three.

$$\tilde{W} = \begin{bmatrix} i_1^T \\ \vdots \\ i_F^T \\ j_1^T \\ \vdots \\ j_F^T \end{bmatrix} \begin{bmatrix} s_1 & \dots & s_p \end{bmatrix} = RS$$

Because  $W$  is a product of two matrices.  
The maximum rank of  $S$  is 3.

# Linearly Independence

A finite subset of  $n$  vectors,  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ , from the vector space  $V$ , is *linearly dependent* if and only if there exists a set of  $n$  scalars,  $a_1, a_2, \dots, a_n$ , not all zero, such that

$$a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \cdots + a_n\mathbf{v}_n = \mathbf{0}.$$

# Rank of a Matrix

- The **column rank** of a matrix  $A$  is the maximum number of linearly independent column vectors of  $A$ .
- The **row rank** of a matrix  $A$  is the maximum number of linearly independent row vectors of  $A$ .
- The column rank of  $A$  is the dimension of the column space of  $A$
- The row rank of  $A$  is the dimension of the row space of  $A$ .

## Example (Row Echelon)

$$A = \begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 1 \\ -2 & -3 & 1 \\ 3 & 5 & 0 \end{bmatrix} R_2 \rightarrow 2r_1 + r_2$$

Rank is 2



## How to find Translation?

$$\tilde{u}_{fp} = u_{fP} - a_f \quad \text{From (A)}$$

$$u_{fp} = \tilde{u}_{fP} + a_f \quad \tilde{u}_{fp} = i_f^T s_P$$

$$u_{fp} = i_f s_p + a_f \quad \text{(E)} \quad u_{fp} = i_f^T (s_p - t_f) \quad \text{From (B)}$$

Comparing above two eqs

$$a_f = -t_f i_f^T \quad \text{(D)}$$

$a_f$  is projection of camera translation along x-axis

## How to find Translation

$$u_{fp} = i_f s_p + a_f \quad v_{fp} = j_f s_p + b_f$$

$$\mathbf{W} = \mathbf{R}\mathbf{S} + \mathbf{t}\mathbf{e}_p^T \quad a_f = -t_f i_f^T$$

$\begin{matrix} 2\text{FXP} & 2\text{FX3} & 3\text{XP} & 2\text{FX1} & 1\text{XP} \end{matrix}$

From (D)

$$\mathbf{t} = (a_1, \dots, a_f, b_1, \dots, b_f)^T$$

$$\mathbf{e}_p^T = (1, \dots, 1)$$

# How to find Translation

Projected camera translation can be computed:

$$-i_f^T t_f = a_f = \frac{1}{P} \sum_{p=1}^P u_p \quad \text{From (D)}$$

$$-j_f^T t_f = b_f = \frac{1}{P} \sum_{p=1}^P v_p$$

# Noisy Measurements

- Without noise, the matrix  $\tilde{W}$  must be at most of rank 3. When noise corrupts the images, however,  $\tilde{W}$  will not be of rank 3. Rank theorem can be extended to the case of noisy measurements.

# Singular Valued Decomposition

SVD

$$\tilde{W} = O_1 \Sigma O_2$$

2FXP                      2FXP      PXP      PXP

# Singular Value Decomposition (SVD)

Theorem: Any  $m$  by  $n$  matrix  $A$ , for which  $m \geq n$ , can be written as

$$A = O_1 \Sigma O_2$$

$m \times n \quad \quad m \times n \quad n \times n \quad n \times n$

$\Sigma$  is diagonal

$O_1, O_2$  are orthogonal

$$O_1^T O_1 = O_2^T O_2 = I$$

# Approximate Rank

$$\tilde{W} = O_1 \Sigma O_2$$

$$O_1 = \begin{matrix} & \begin{matrix} 3 & \mathbf{P-3} \end{matrix} \\ \begin{matrix} O'_1 & O''_1 \end{matrix} \end{matrix} \quad \begin{matrix} 2F \\ \end{matrix}$$

$$\Sigma = \begin{matrix} & \begin{matrix} 3 & \mathbf{P-3} \end{matrix} \\ \begin{bmatrix} \Sigma' & 0 \\ 0 & \Sigma'' \end{bmatrix} \end{matrix} \quad \begin{matrix} 3 \\ \mathbf{P-3} \end{matrix}$$

$$O_1 \Sigma O_2 = O'_1 \Sigma' O'_2 + O''_1 \Sigma'' O''_2$$

$$O_2 = \begin{matrix} & \begin{matrix} 3 \\ \mathbf{P-3} \end{matrix} \\ \begin{bmatrix} O'_2 \\ O''_2 \end{bmatrix} \end{matrix}$$

# Approximate Rank

$$\tilde{W} = O_1 \Sigma O_2 = O_1' \Sigma' O_2' + O_1'' \Sigma'' O_2''$$

$$\hat{W} = O_1' \Sigma' O_2'$$

The best rank 3 approximation to the ideal registered measurement matrix.



# Rank Theorem for noisy measurement

The best possible shape and rotation estimate is obtained by considering only 3 greatest singular values of  $\tilde{W}$  together with the corresponding left, right eigenvectors.

# Approximate Rank

$$\hat{R} = O_1' [\Sigma']^{1/2}$$

Approximate Rotation matrix

$$\hat{S} = [\Sigma']^{1/2} O_2'$$

Approximate Shape matrix

$$\hat{W} = \hat{R} \hat{S}$$

This decomposition is not unique

$$\hat{W} = (\hat{R} Q)(Q^{-1} \hat{S})$$

$Q$  is *any* 3X3 invertable matrix

# Results

[..\..\CAP6411\Fall2002\tomasiTr92Figures.pdf](#)

# Hotel Sequence



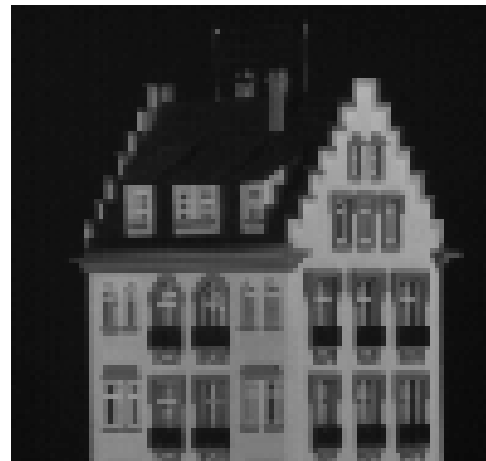
1



60

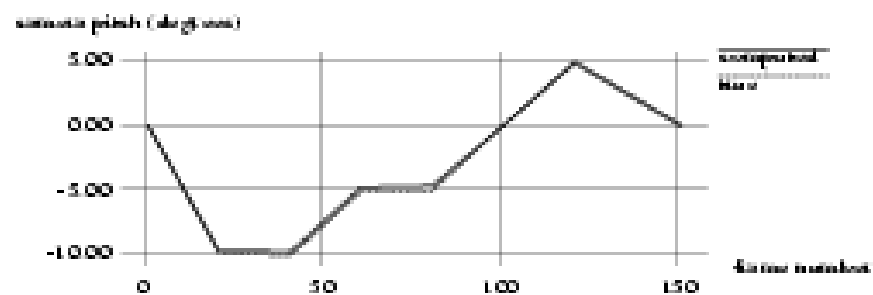
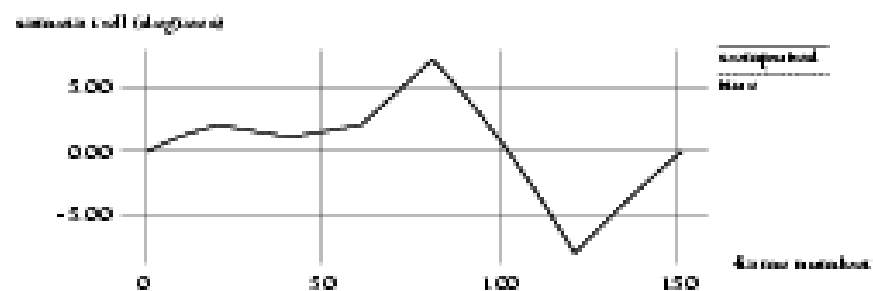
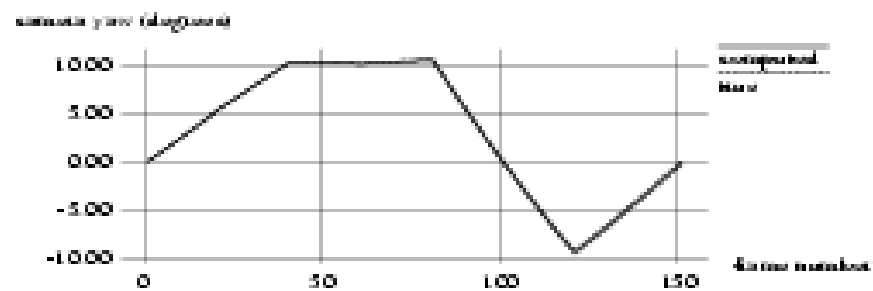


120

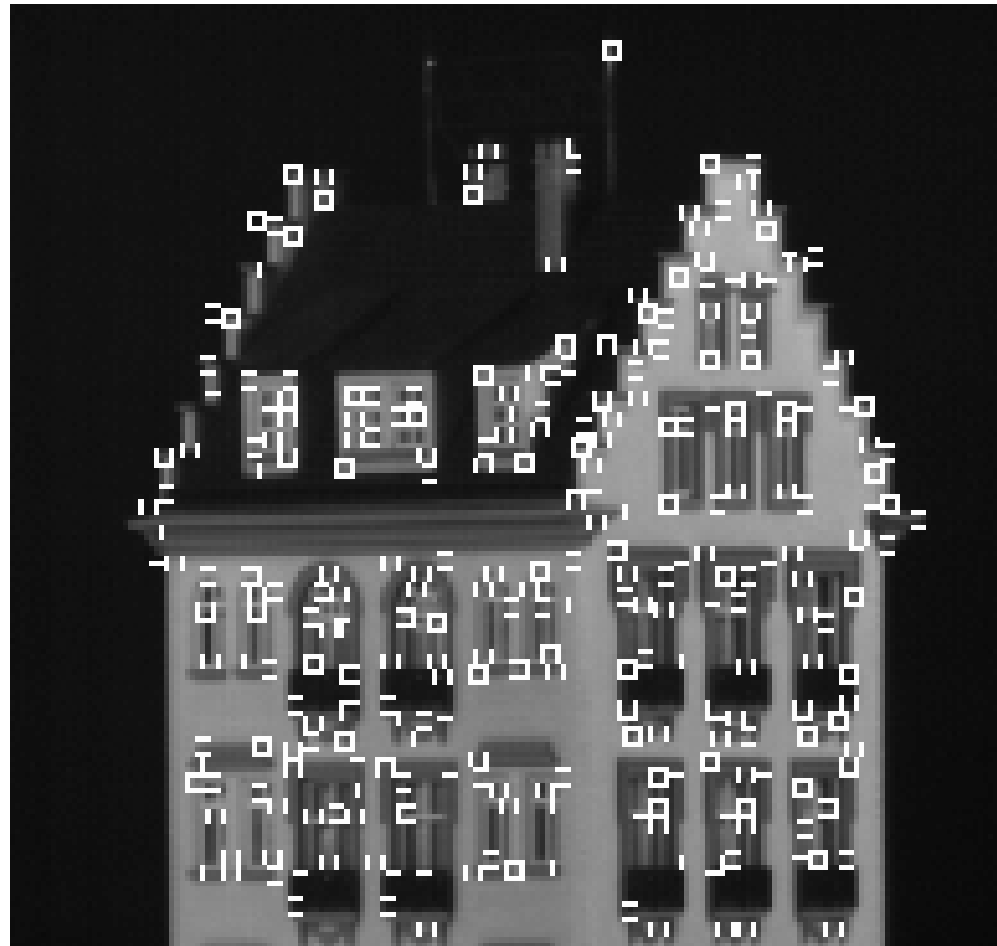


150

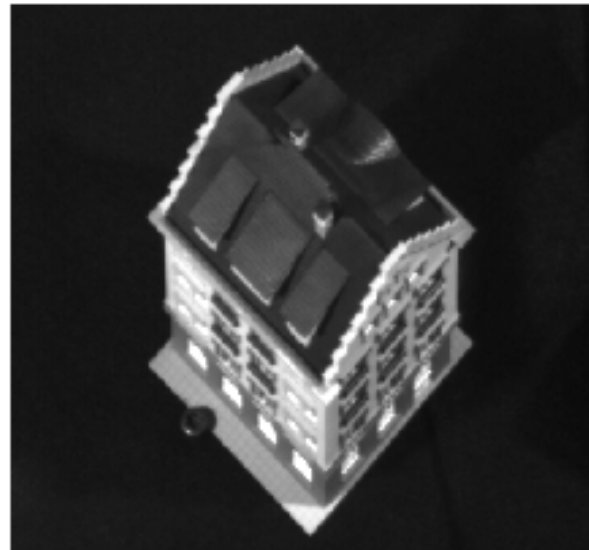
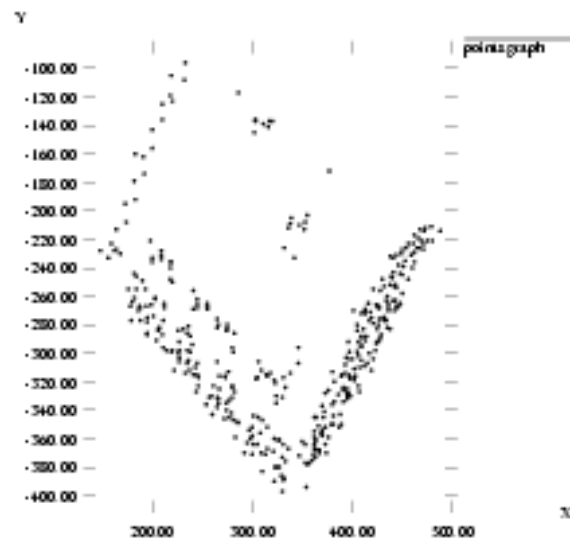
# Results (rotations)



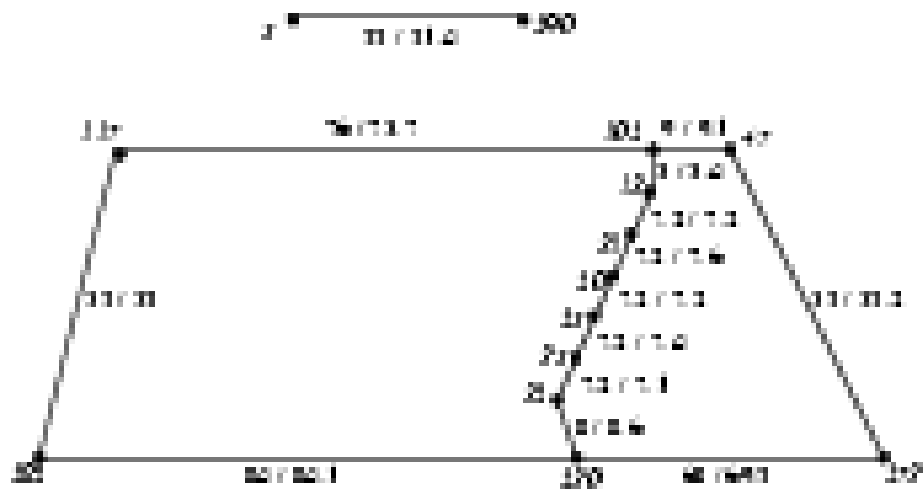
# Selected Features



# Reconstructed Shape



# Comparison





# House Sequence



1



60

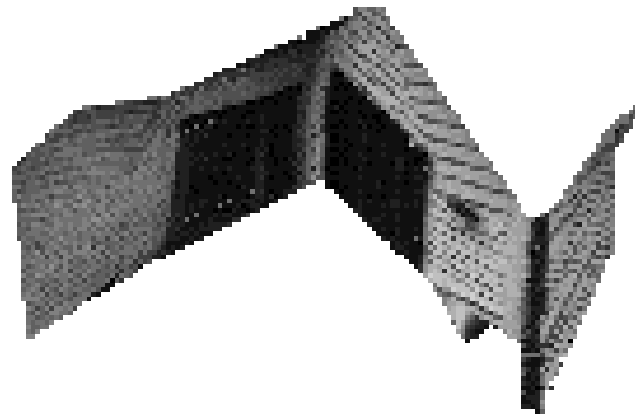
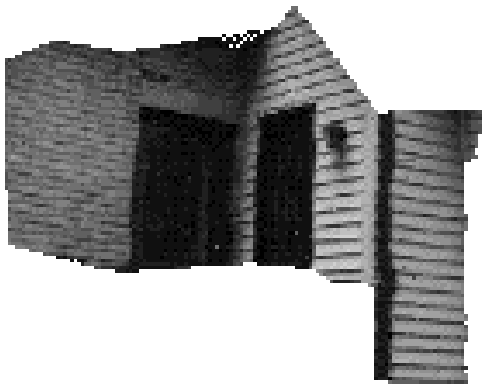


120



180

# Reconstructed Walls



## Further Reading

- C. Tomasi and T. Kanade. Shape and motion from image streams under orthography---a factorization method. *International Journal on Computer Vision*, 9(2):137-154, November 1992.
- Computer Vision: Algorithms and Applications, Richard Szeliski, Section 7.3