

VIDEO SURVEILLANCE & MONITORING

Mubarak Shah
Computer Vision Lab
University of Central Florida

Main Steps



plane



walking



car

person



running



ship



Waving

Object detection

Object tracking

Object categorization
and classification

Events or Activities
Recognition



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

BAYESIAN MODELING OF DYNAMIC SCENES FOR OBJECT DETECTION

Omar Javed, Khurram Shafique, and
Mubarak Shah

IEEE Workshop on Motion & Video Computing
2002



University of
Central Florida

VISION

Interest Region Detection Using Background Subtraction

- Aim:
 - Mark pixels in the image corresponding to interesting objects
 - Completely unsupervised learning



- General Approach:
 - Build a per-pixel model of background
 - Find deviations from the model

Background Subtraction

- Important Problems in Realistic situations:
 - Quick illumination changes



- Relocation of background objects.



Background Subtraction

- Initialization with moving objects



- Shadows



Background Modeling

THE HIERARCHICAL APPROACH

- Which features to use?
 - Color
 - Gradient

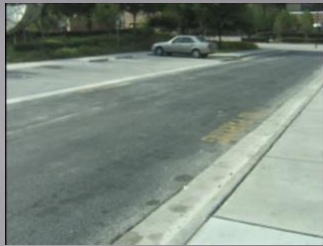


Image-1

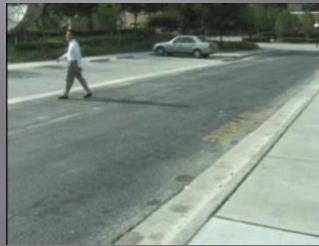
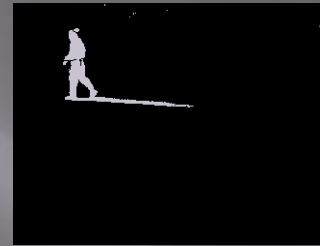
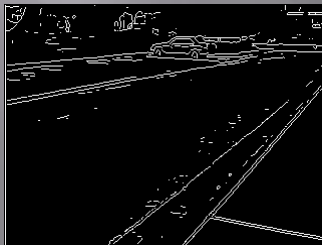


Image-153



Color based (Image 153)



Gradient
Image-1



Gradient
Image-153



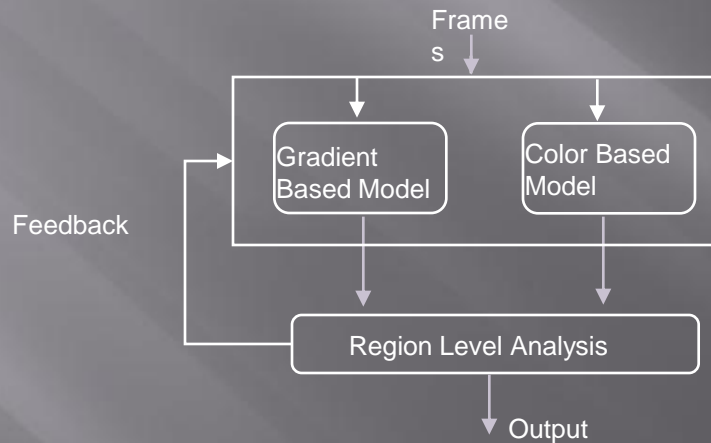
Gradient based
(image 153)



Background Modeling

THE HIERARCHICAL APPROACH

- Fusion of features
 - Validation of pixel level results at the regional level.



- Updating of pixel level models based on feedback from regional level process.

Pixel Level Subtraction

- Color based subtraction
 - Per-pixel mixture of Gaussians.
- Gradient based subtraction
 - Gradient feature vector, $\Delta = [\Delta_m, \Delta_d]$

where

$$\Delta_m = \sqrt{f_x^2 + f_y^2}$$

$$\Delta_d = \tan^{-1}\left(\frac{f_y}{f_x}\right)$$

- Distribution of the gradient feature vector?



Subtraction in the Gradient Domain

- Let $x_{i,j}^t$ be the latest value that matched k^{th} distribution belonging to background at pixel (i,j)

- The gray level value will be given as

$$g_{i,j}^t = \alpha R + \beta G + \gamma B$$

- Assuming independence of color channels,

$$g_{i,j}^t \sim N(\mu_{i,j}^t, (\sigma_{i,j}^t)^2)$$



Subtraction in the Gradient Domain

- Let us define

$$f_x = g_{i+1,j}^t - g_{i,j}^t$$

$$f_y = g_{i,j+1}^t - g_{i,j}^t$$

- Assuming independence between neighboring pixels.

$$f_x \sim N(\mu_{f_x}, (\sigma_{f_x})^2)$$

$$f_y \sim N(\mu_{f_y}, (\sigma_{f_y})^2)$$

- The covariance is given by

$$\text{Cov}(f_x, f_y) = \text{Cov}(g_{i+1,j}^t - g_{i,j}^t, g_{i,j+1}^t - g_{i,j}^t) = (\sigma_{i,j}^t)^2$$



Subtraction in the Gradient Domain

- Distribution of feature vector $[\Delta_m, \Delta_d]$

$$F(\Delta_m, \Delta_d) = \frac{\Delta_m}{2\pi\sigma_{fx}\sigma_{fy}\sqrt{1-\rho^2}} \exp\left(-\frac{z}{2(1-\rho^2)}\right)$$

– where

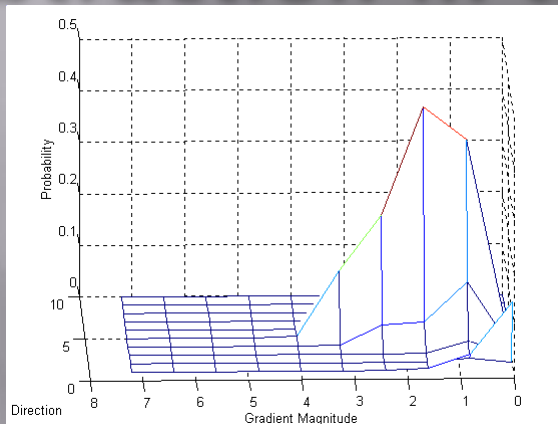
$$z = \left(\frac{\Delta_m \cos \Delta_d - m_{fx}}{\sigma_{fx}}\right)^2 - 2\rho\left(\frac{\Delta_m \cos \Delta_d - m_{fx}}{\sigma_{fx}}\right)\left(\frac{\Delta_m \sin \Delta_d - m_{fy}}{\sigma_{fy}}\right) + \left(\frac{\Delta_m \sin \Delta_d - m_{fy}}{\sigma_{fy}}\right)^2$$

$$\rho = \frac{\sigma_{i,j}^2}{\sigma_{fx}\sigma_{fy}}$$

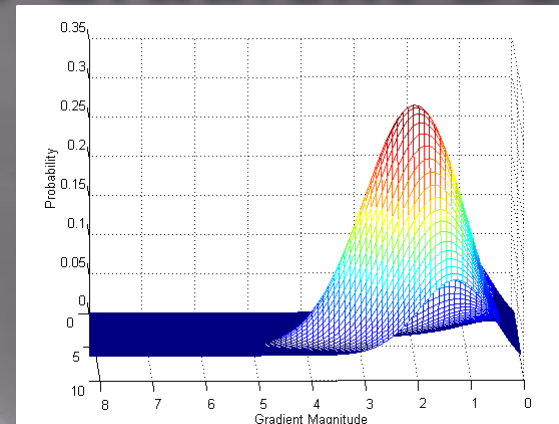
- If $F(\Delta_m, \Delta_d) < T_g$ then pixel is marked as foreground.



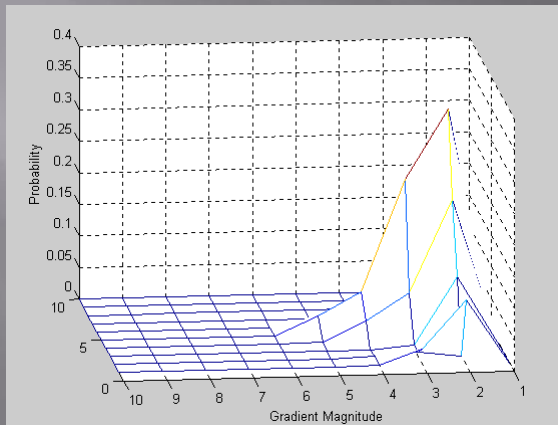
Subtraction in the Gradient Domain



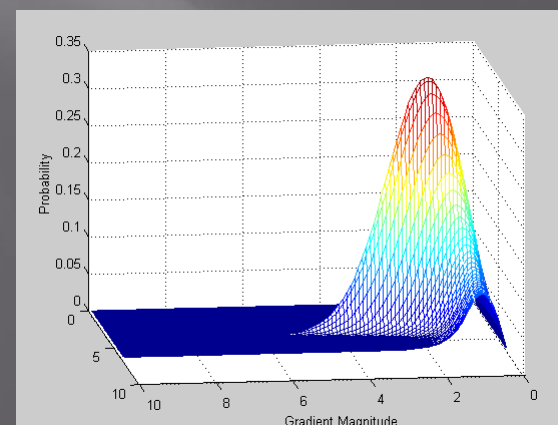
Sample Histogram of $\Delta = [\Delta_m, \Delta_d]$



Parametric Distribution of $\Delta = [\Delta_m, \Delta_d]$



Sample Histogram of $\Delta = [\Delta_m, \Delta_d]$



Parametric Distribution of $\Delta = [\Delta_m, \Delta_d]$



Pixel Level Subtraction

- Examples of pixel wise subtraction in the color and gradient domain.



Image-1



Image-85



Color based
(Image 85)



Gradient based
(Image 85)

Region Level Processing

- A region with edges on its boundary that are different from the background is a valid region.
- A region R is accepted as a valid region if

$$\frac{\sum_{(i,j) \in \partial R} (\nabla I(i,j) G(i,j))}{|\partial R|} \geq \rho_B$$

Where

- ∂R is the set of boundary pixels of a connected region R in color based results
- $G(i,j)$ is the gradient based subtraction output at pixel (i,j)
- $\nabla I(i,j)$ is the edge map

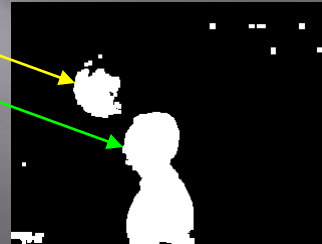


Region Level Processing

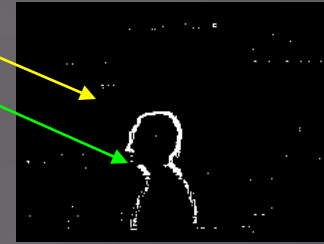
- For each color based region, presence of “edge difference” pixels at the boundaries is checked.



Image

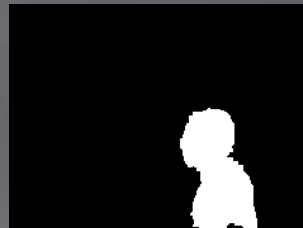


Color based



Gradient

- Regions with small number of edge difference pixel are removed, color model is updated.

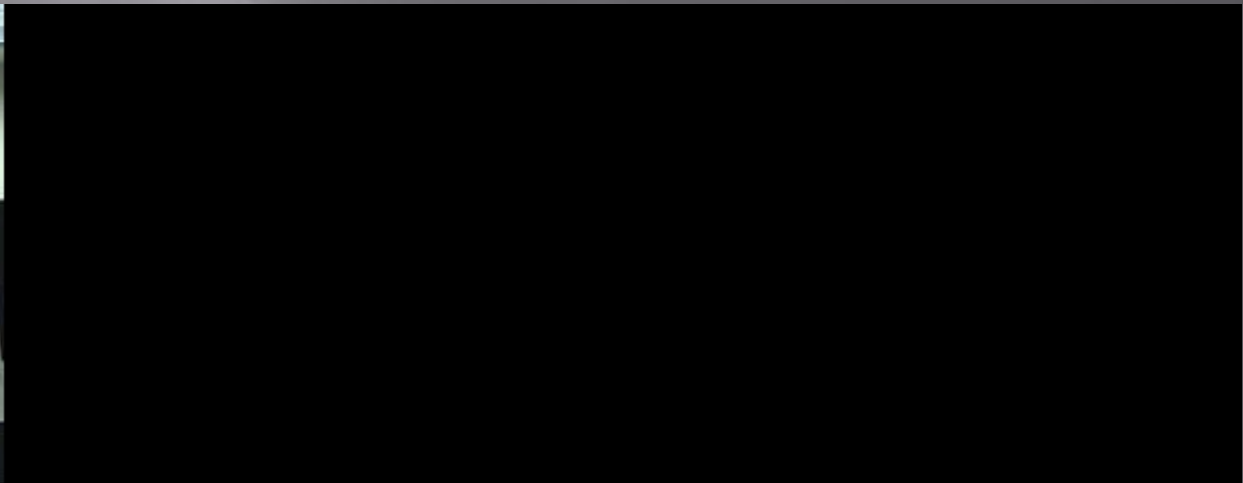


Final

Background Modeling

THE HIERARCHICAL APPROACH

- Local Illumination Change.



Mixture of
Gaussians
(S&G) Method

Hierarchical
Subtraction



University of
Central Florida

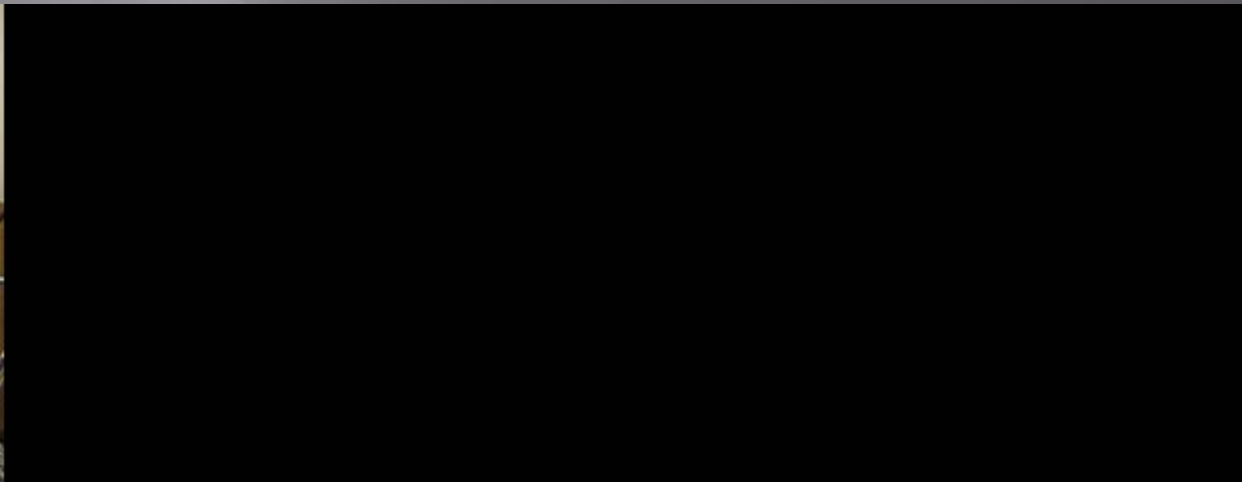
VISION

Copyright Mubarak Shah, UCF

Background Modeling

THE HIERARCHICAL APPROACH

- Relocation of background object.



Mixture of
Gaussians
(S&G) Method

Hierarchical
Subtraction



University of
Central Florida

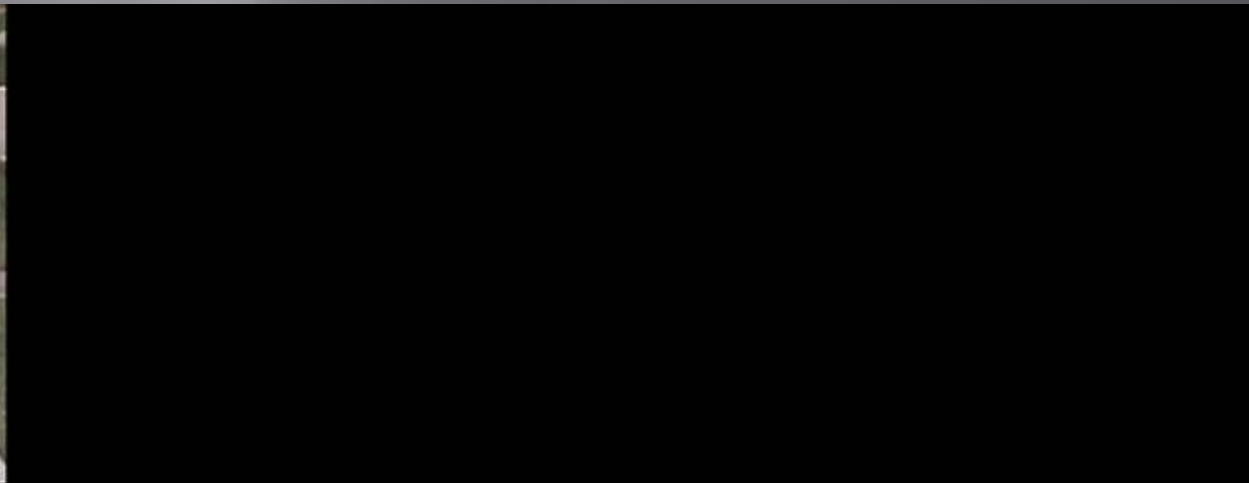
VISION

Copyright Mubarak Shah, UCF

Background Modeling

THE HIERARCHICAL APPROACH

- Moving with initialization & illumination change.



Mixture of
Gaussians
(S&G) Method

Hierarchical
Subtraction



University of
Central Florida

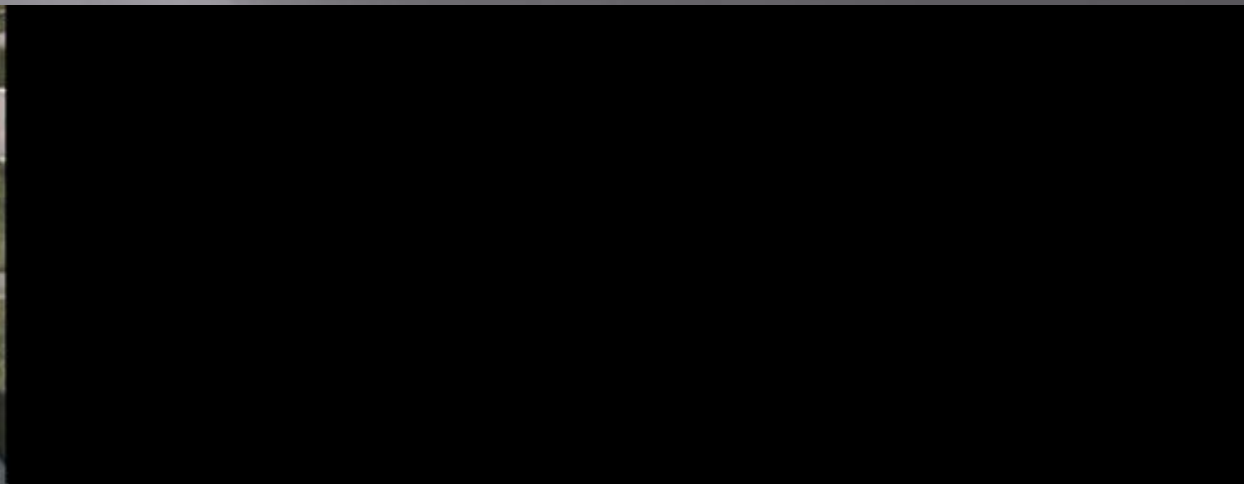
VISION

Copyright Mubarak Shah, UCF

Background Modeling

THE HIERARCHICAL APPROACH

- Quick Illumination change.



Mixture of
Gaussians
(S&G) Method

Hierarchical
Subtraction



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

BAYESIAN MODELING OF DYNAMIC SCENES FOR OBJECT DETECTION

Yaser Sheikh and Mubarak Shah

*IEEE Conference on
Computer Vision and Pattern Recognition
2005*



University of
Central Florida

VISION

What is a Dynamic Scene?

(and why should I be interested?)



Periodic



Temporal texture



Nominal camera motion



Temporal texture



Overview

- ▣ Modeling the Background
 - Non-parametric density estimation
 - Joint Domain-Range Feature Space
- ▣ Modeling the Foreground
 - Competitive detection (back ground vs foreground)
- ▣ MAP-MRF estimation framework
 - Efficient minimization using graph cuts



The Background Distribution

- ▣ **Object Detection:** Given an image, what is the probability of observing a pixel color at a certain location?
 $P(\text{pixel} \mid \text{background video})$ instead of
 $P(\text{pixel} \mid \text{background pixel})$
- ▣ Analysis on $\mathbf{x}_i \in \mathbb{R}^5$, the feature space:
 $[R, G, B, X, Y]$
- ▣ This is our background *model*
- ▣ Use Kernel Density Estimation on this 5 dimensional space
- ▣ Probabilistic Low Level Descriptor



Kernel Density Estimation

(a.k.a. Parzen Windows)

▣ Definition of KDE

$$f(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h_i} K\left(\frac{x - x_i}{h_i}\right)$$

where K is a d -variate kernel function usually satisfying .

- The Gaussian kernel is a common choice.



Kernel Density Estimation

- ▣ Characteristics
 - Nonparametric technique
 - Effective multi-modal data representation
- ▣ Background model is represented in the 5-space with the set $\psi_b = \{y_1, y_2 \dots y_n\}$, where n pixels have been observed thus far.

$$K(x, y) = \frac{1}{n} \sum_{i=1}^n K(x, y_i)$$



Temporal Persistence

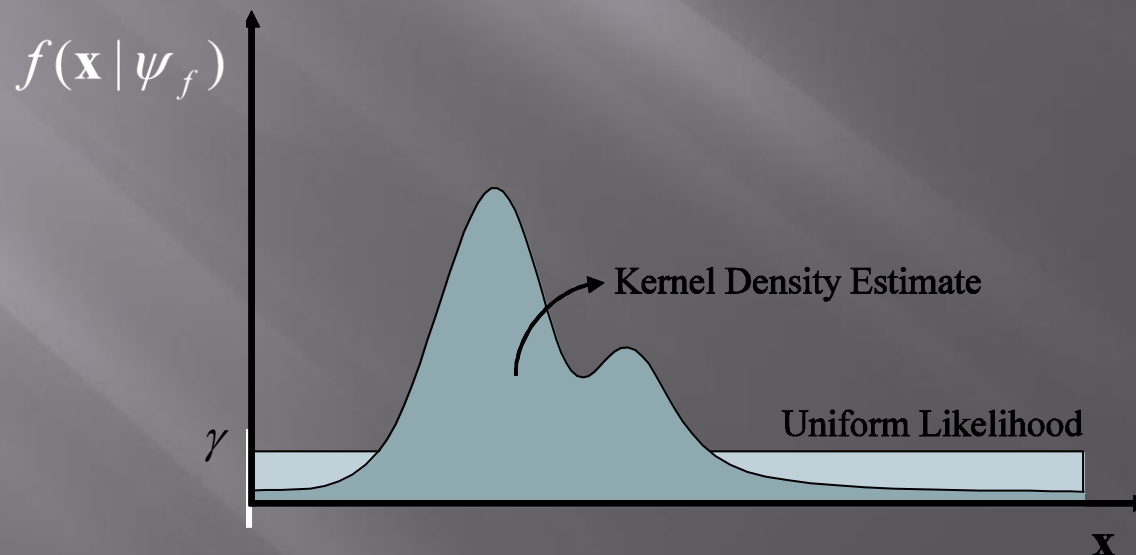
- ▣ **Intuition:** Objects tend to maintain constant colors, and tend to persist in the same proximity (smooth motion)
- ▣ Frame at $t-1$ contains substantial evidence for detection at t
- ▣ Given the detected foreground in previous frames, what is the probability of the object belonging to that foreground?



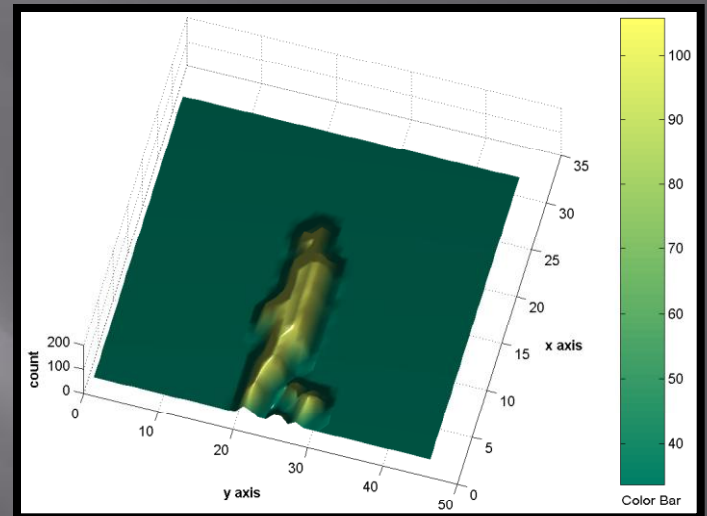
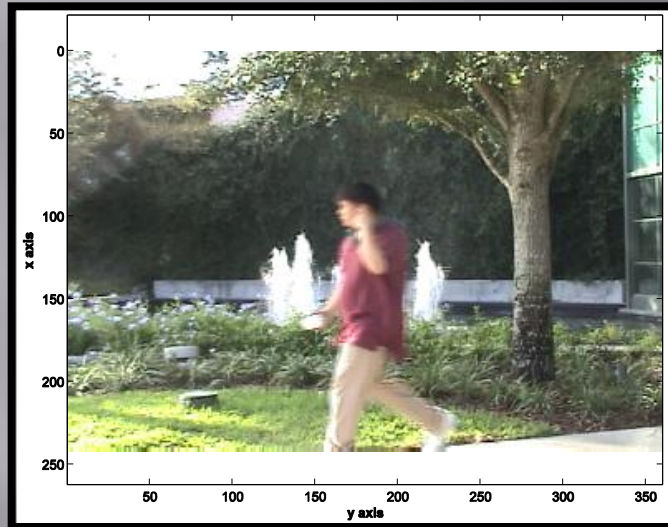
Foreground Model

▣ Foreground Density Estimator

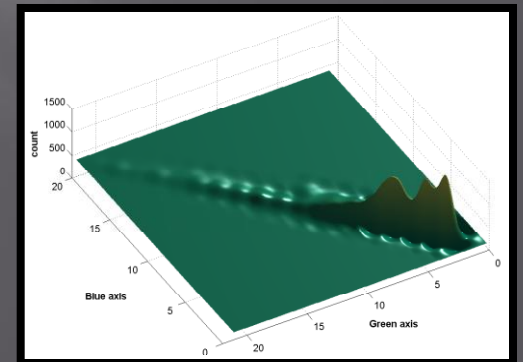
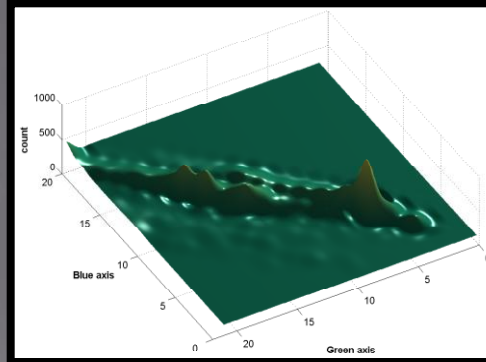
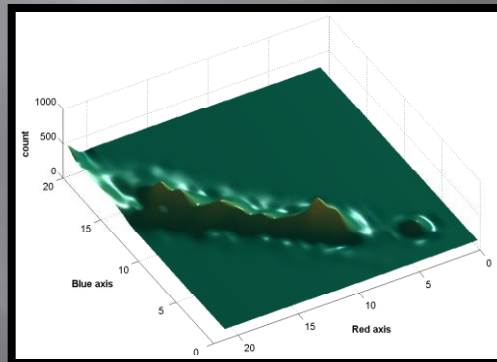
$$\hat{f}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m K(\mathbf{x} - \mathbf{x}_i)$$



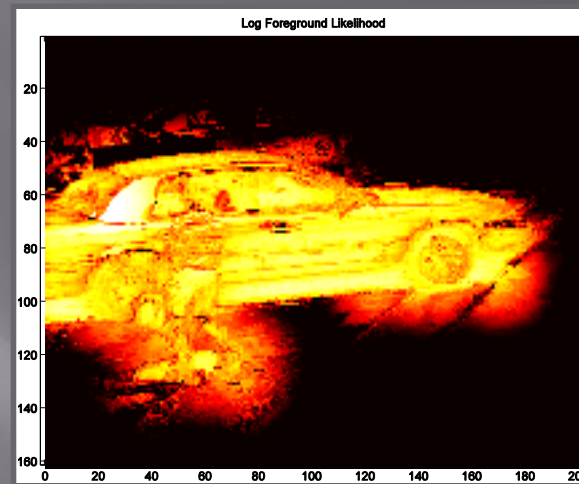
Marginals of the Foreground Distribution



BR
BG
GR



Classification vs Detection



foreground



ratio



background

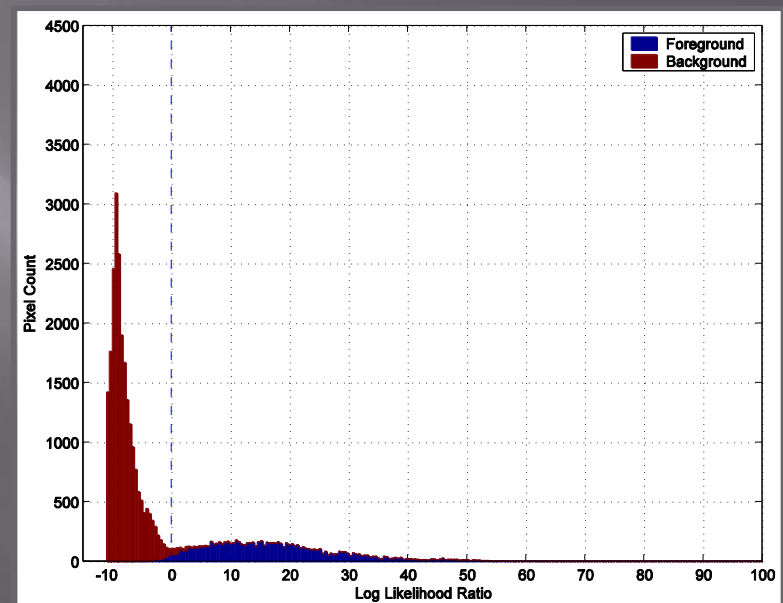
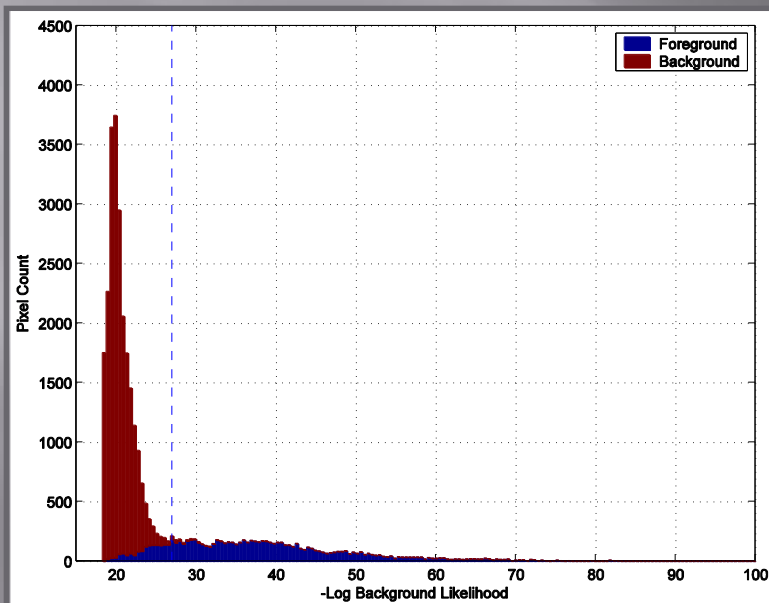


University of
Central Florida

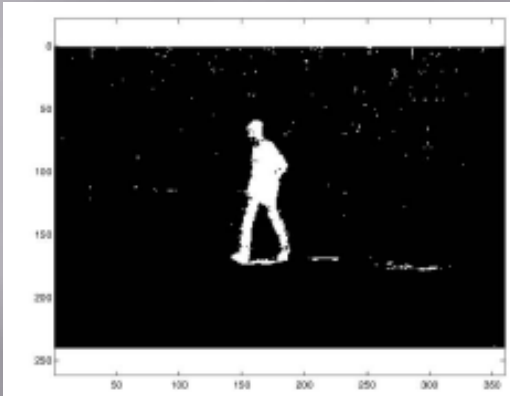
VISION

Copyright Mubarak Shah, UCF

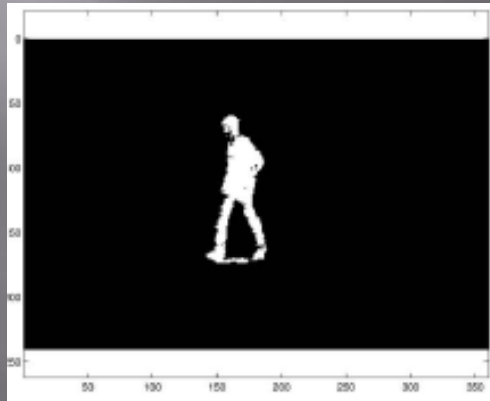
Improved Discrimination



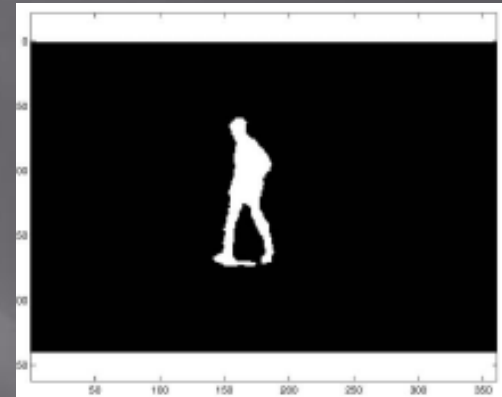
Detection



Thresholding on
background model



Thresholding on
Likelihood model

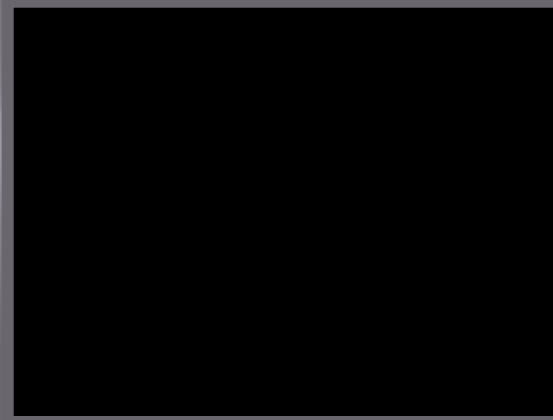


MAP MRF

The Watery Sequence



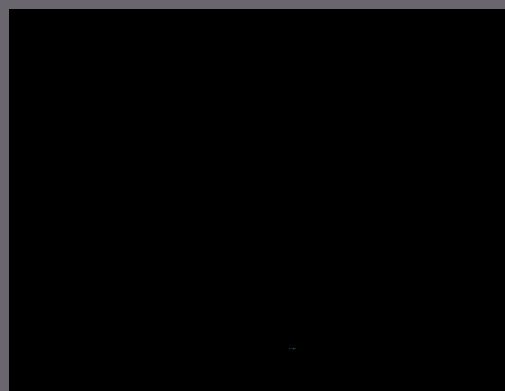
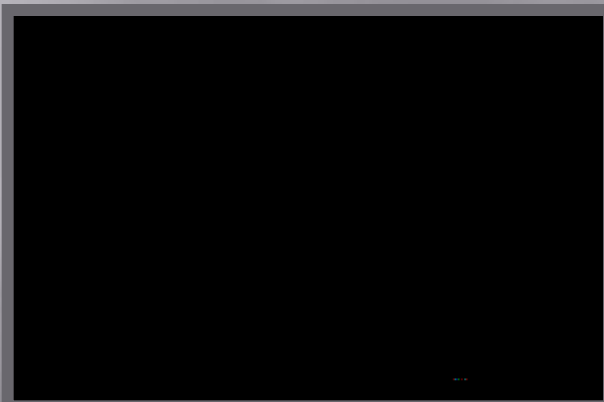
**Mixture Of
Gaussians**



Proposed



The Fountain Sequence

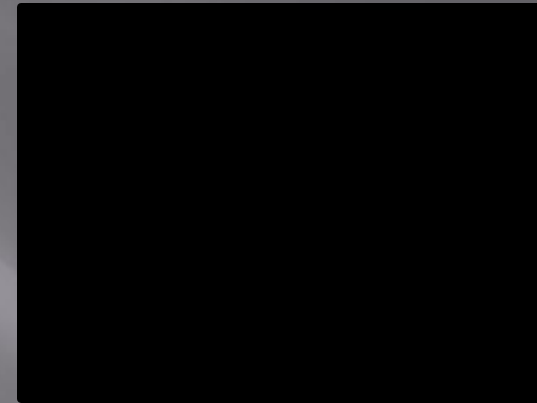


Proposed

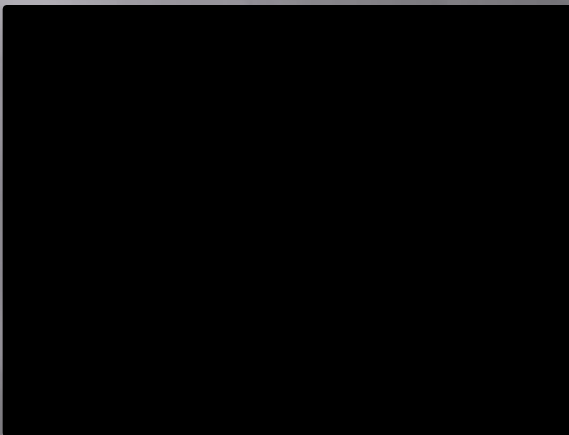
Mixture Of Gaussians



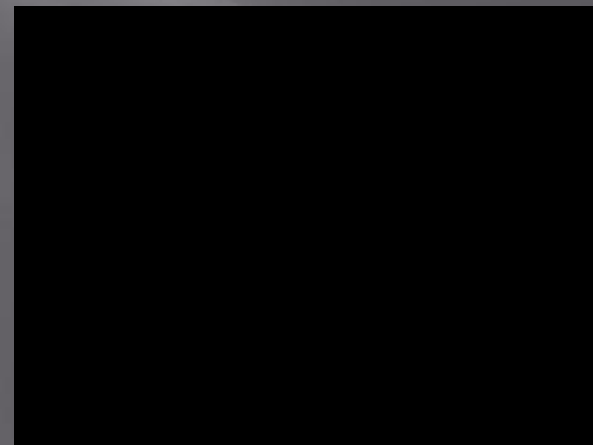
Nominal Motion Sequence



Proposed



**Ground
Truth**



**Mixture Of
Gaussians**

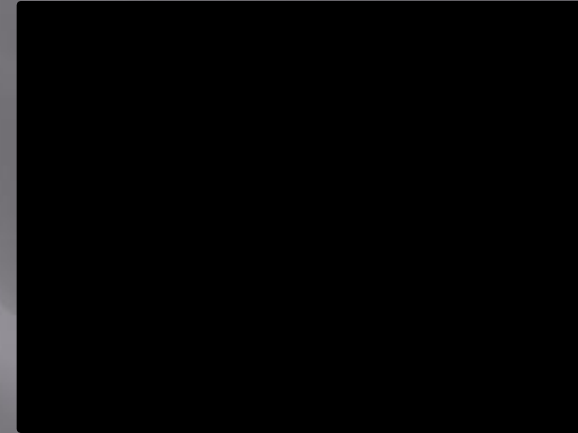


Quantitative Analysis (Pixel Level)

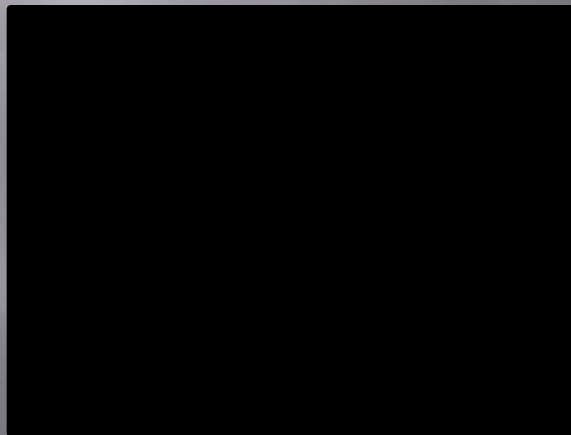
- ▣ Manually segmented 500 frames of the nominal motion sequence
- ▣ Pixel-wise comparison between manual detection and proposed detection scheme



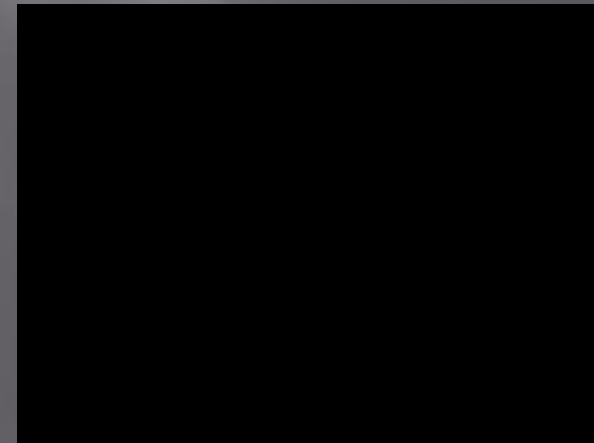
Nominal Motion Sequence



Proposed

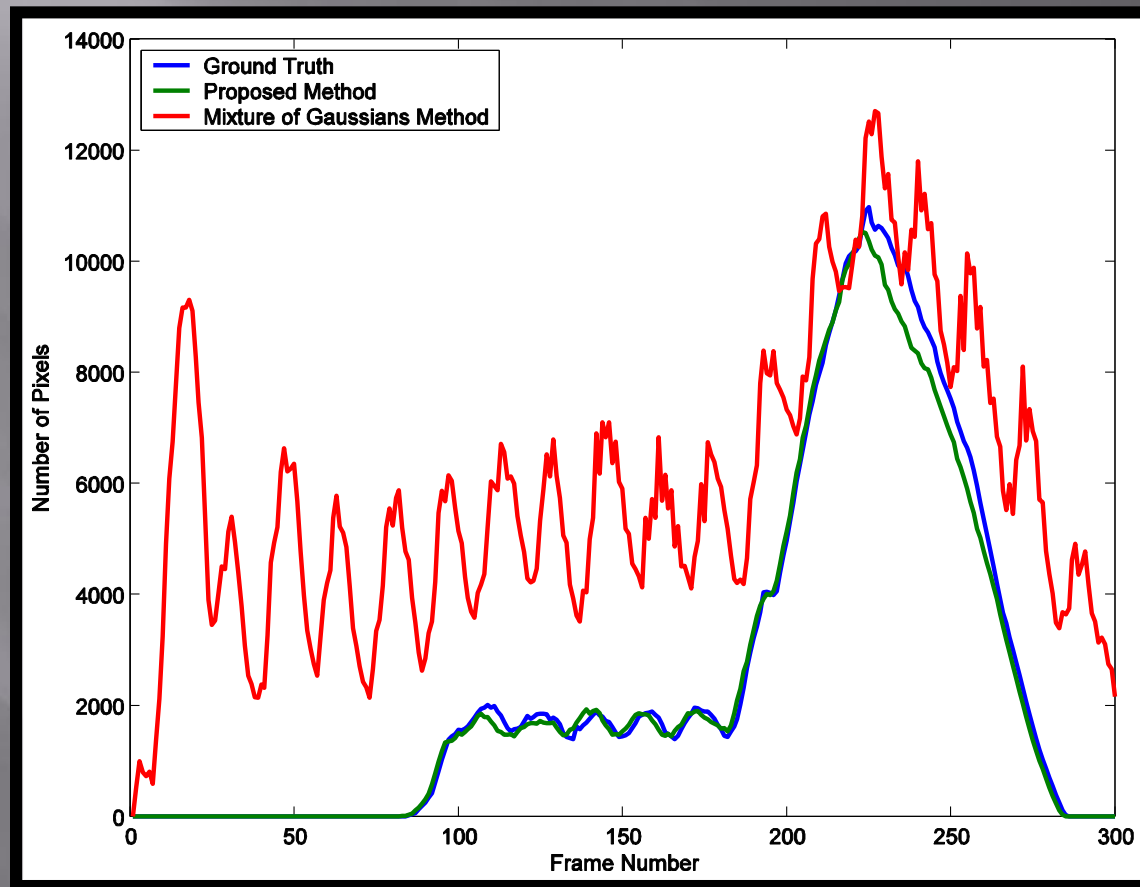


**Ground
Truth**



**Mixture Of
Gaussians**

Pixel-wise Analysis



Quantitative Analysis

Object Level

	Objects	Detected	Mis-detected	Detection Rate	Mis-Detection Rate
Sequence 1	84	84	0	100.00%	0.00%
Sequence 2	115	114	1	99.13%	0.87%
Sequence 3	161	161	0	100.00%	0.00%
Sequence 4	94	94	0	100.00%	0.00%
Sequence 5	170	169	2	99.41%	1.18%

VISUAL TRACKING



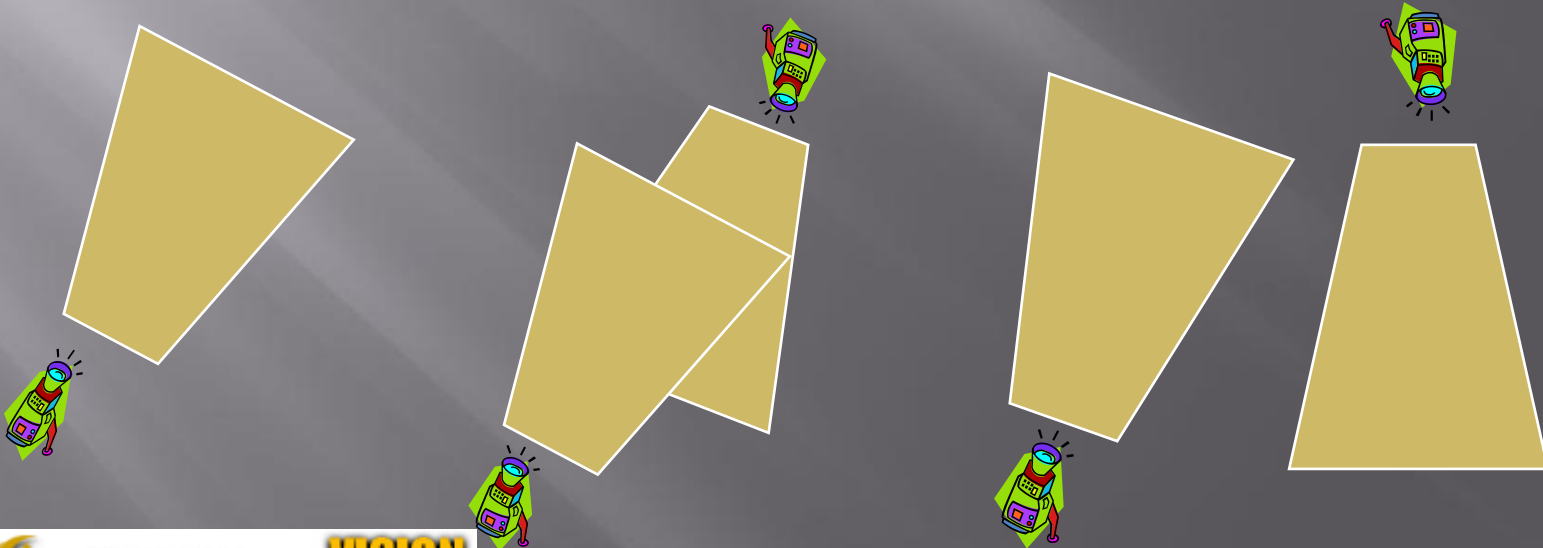
University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

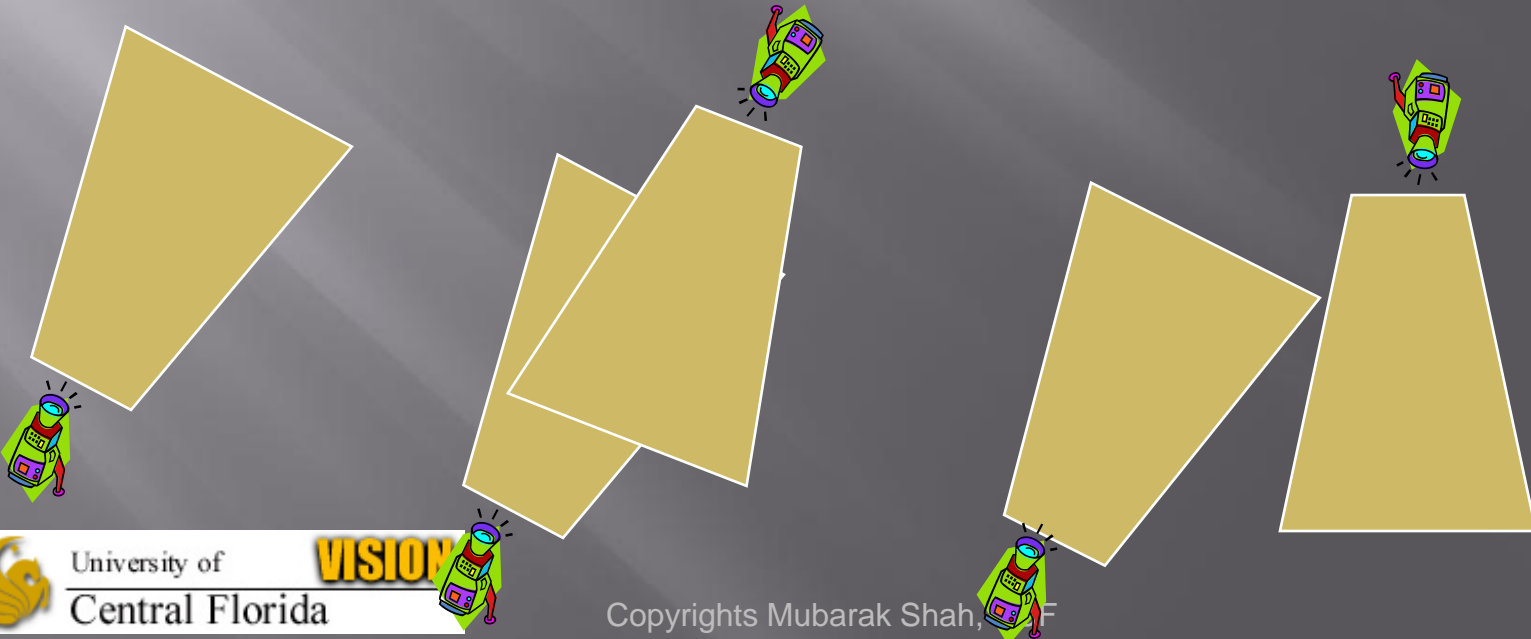
Camera Configurations

- ▣ Stationary Camera
 - Single camera (KNIGHT)
 - ▣ Javed et al. [ECCV 2002](#)
 - Multiple cameras with overlapping field of view
 - ▣ Khan et al. [PAMI 2003](#), Khan et al ECCV 2006
 - Multiple cameras with non-overlapping field of view
 - ▣ Javed et al. ICCV 2003, [CVPR 2005](#), [AAAI-2007](#)



Camera Configurations

- ▣ Moving Cameras
 - Single camera
 - ▣ Yilmaz et al. [PAMI 2004](#), Khan et al [CVPR 2007](#)
 - Multiple cameras with overlapping field of view
 - ▣ Sheikh et al. [ICCV 2005](#), Yilmaz et al ICCV2005
 - Multiple cameras with non-overlapping field of view
 - ▣ Sheikh et al. [CVPR 2007](#)



Tracking IN Single Fixed Camera



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

KNIGHT

Crime Scene Detection System for The Orlando Police Department



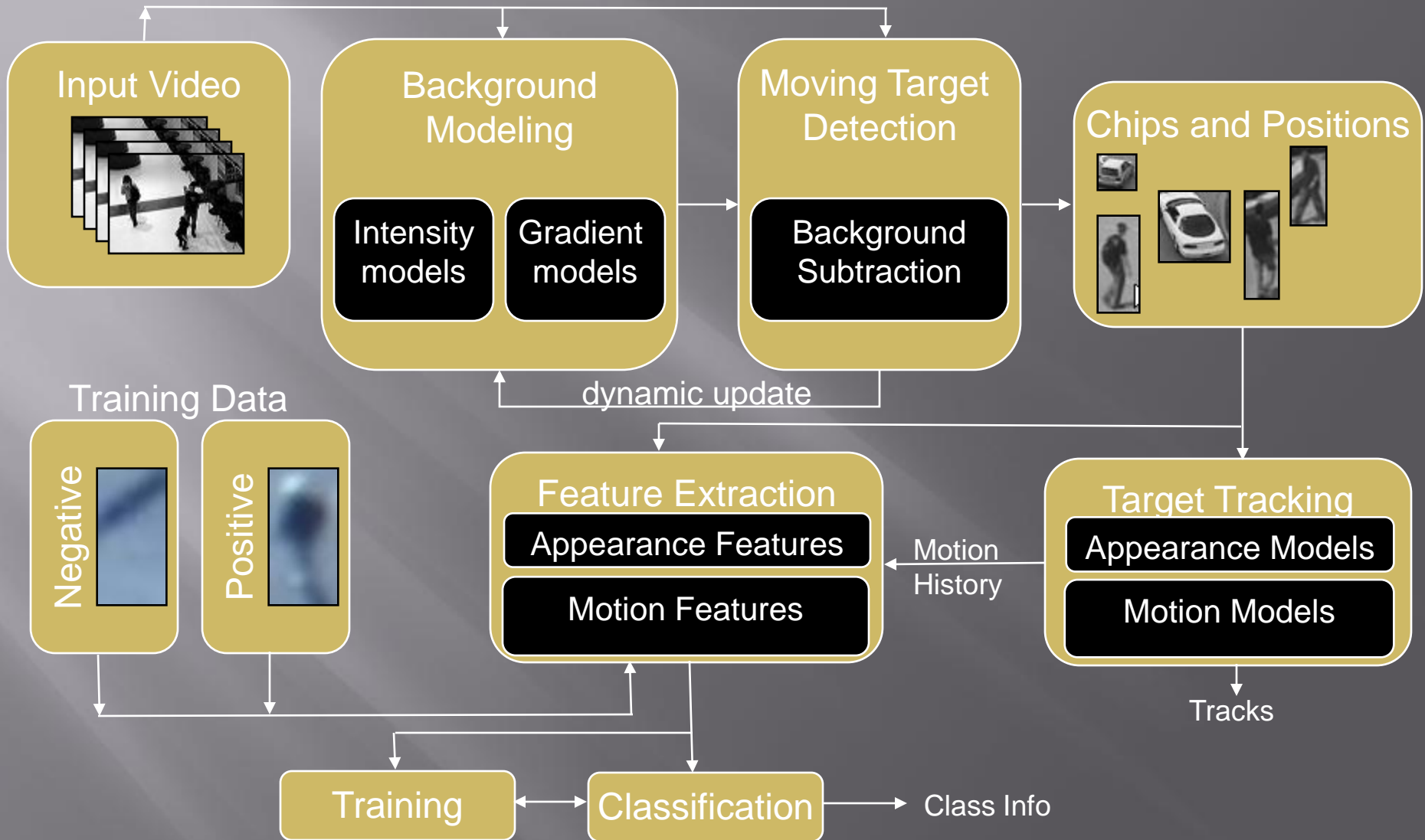
University of
Central Florida

VISION

Views

Copyrights Mubarak Shah, UCF

KNIGHT: Video Surveillance System



KNIGHT: Single Fixed Camera Tracking



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

KNIGHT: Single Fixed Camera Tracking (Occlusion)



Visual Monitoring of Railroad Crossings

▣ Detected Violators



03-01-2005, 11:32:33



03-01-2005, 11:32:19



03-08-2005, 13:18:23



02-22-2005, 11:01:42



03-02-2005, 12:42:57



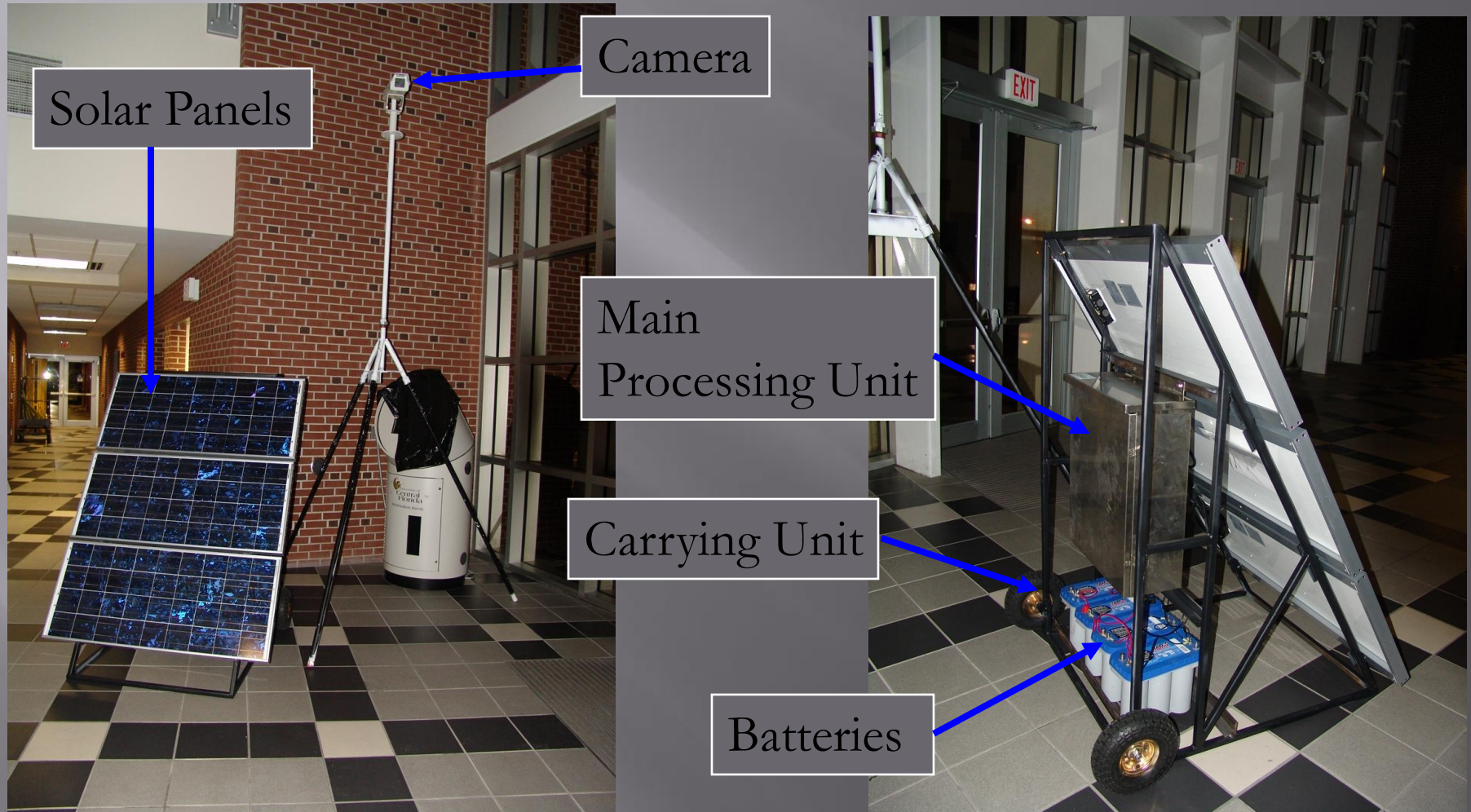
03-11-2005, 15:17:53



University of
Central Florida

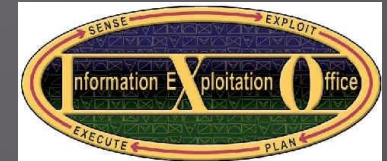
Copyright Mubarak Shah, UCF

The Portable System



Nighttime Video Surveillance

DARPA Phase II STTR



Space and Naval Warfare Systems Center San Diego



University of
Central Florida

VISION

Single Moving Camera Tracking



MULTIPLE CAMERA TRACKING



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

TRACKING ACROSS MULTIPLE FIXED OVERLAPPING CAMERA

Sohaib Khan and Mubarak Shah

PAMI 2003

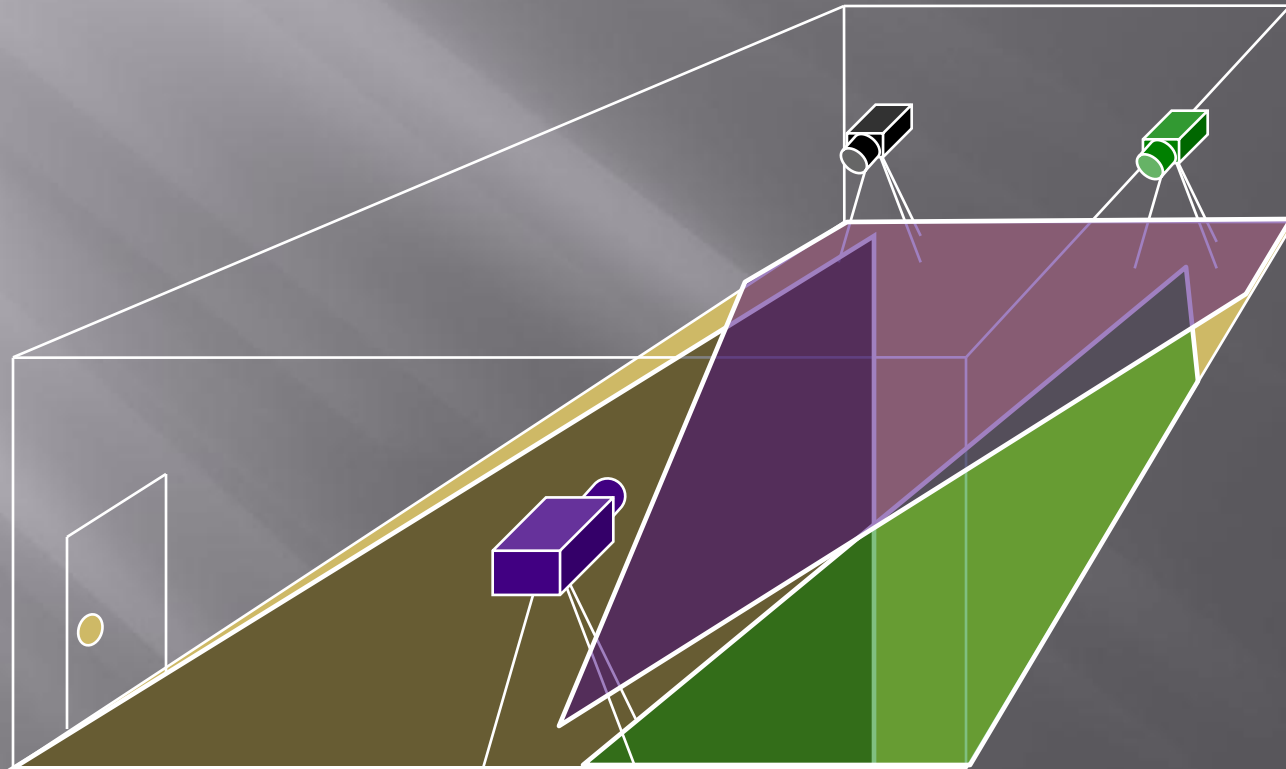


University of
Central Florida

VISION

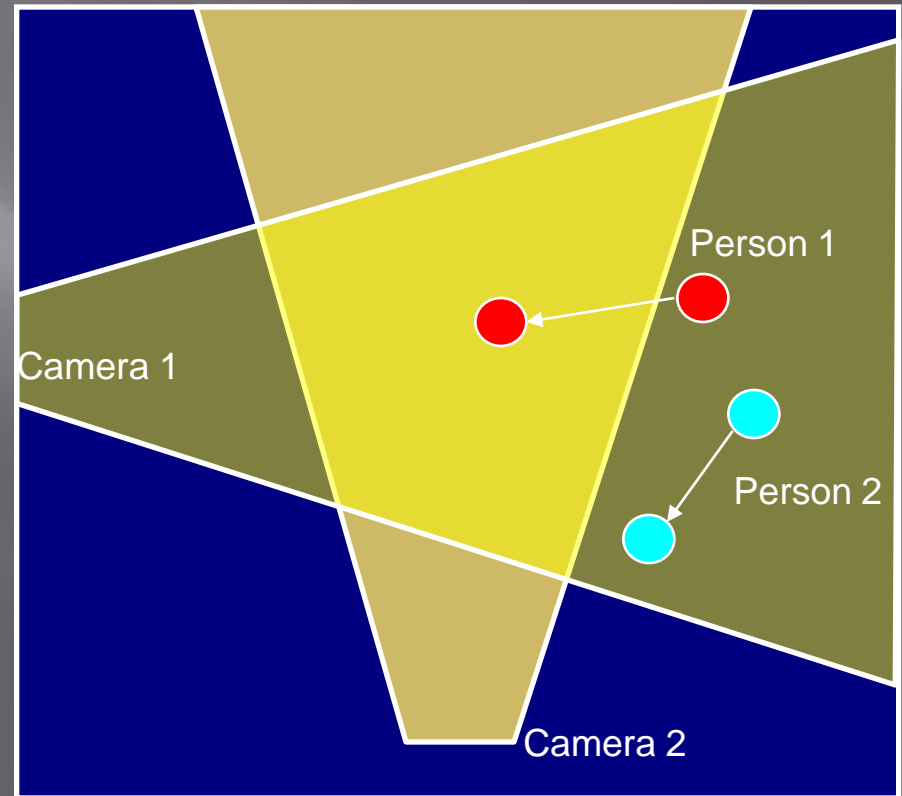
Copyrights Mubarak Shah, UCF

Completely Covering an Environment



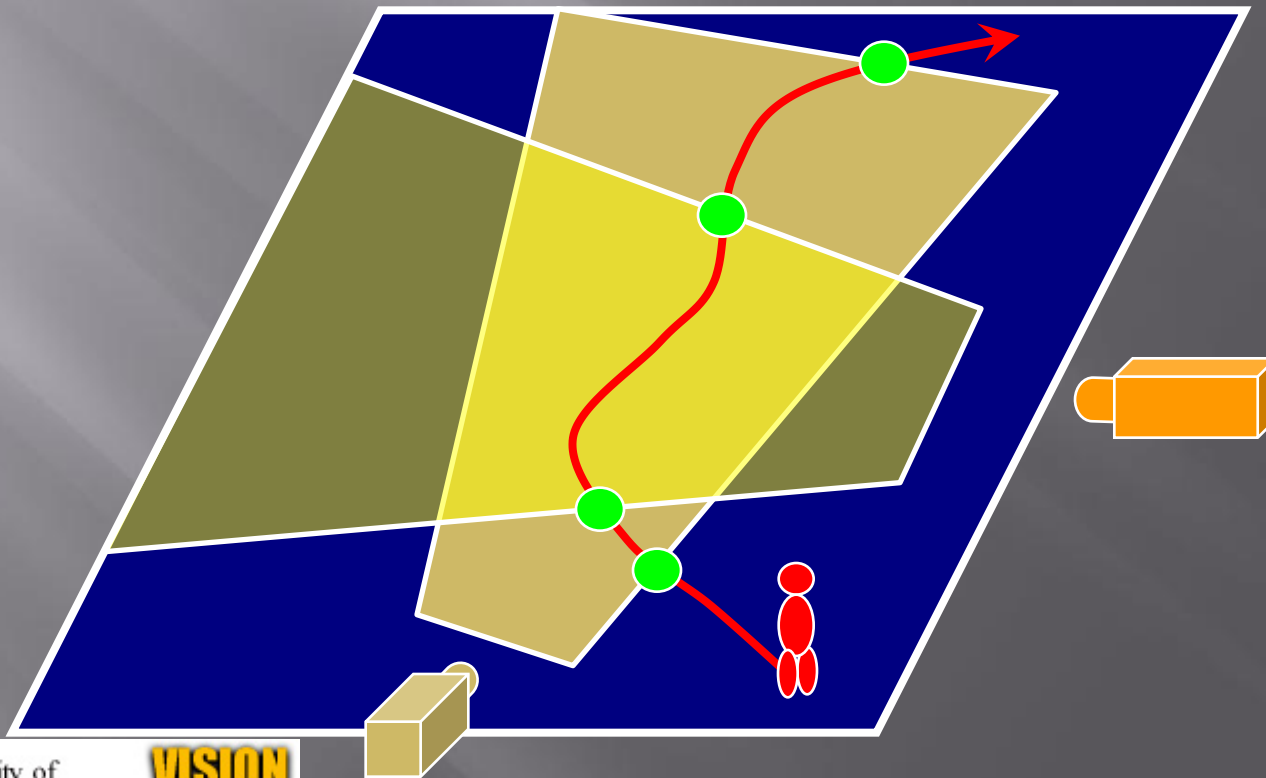
Ambiguity in Consistent Labeling

- What is the label of new person entering FOV of Camera 2?



View Events

- ▣ A view-event is an instant in time when an object enters or leaves the FOV of a camera



FOV lines

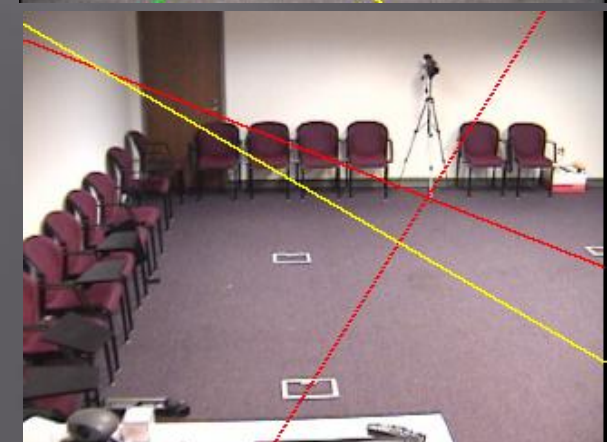
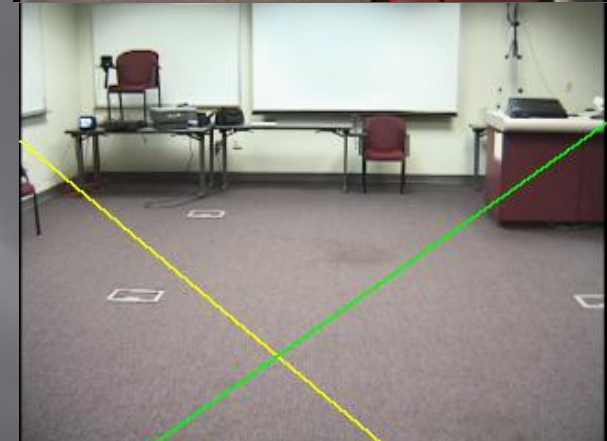
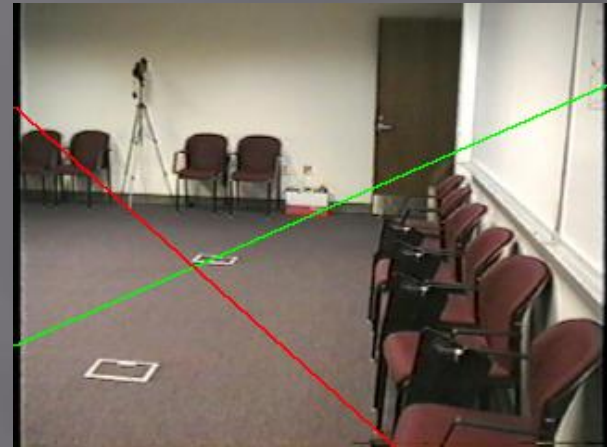
- ▣ Track the bottom of the bounding box
- ▣ Two non-ambiguous correspondences can mark edge of FOV line
- ▣ This line can be used for future ambiguities



Experiments - 1

- ▣ Three Cameras, indoor environment
- ▣ Training:
 - One person walked in the room for about 20 sec





University of
Central Florida

VISION

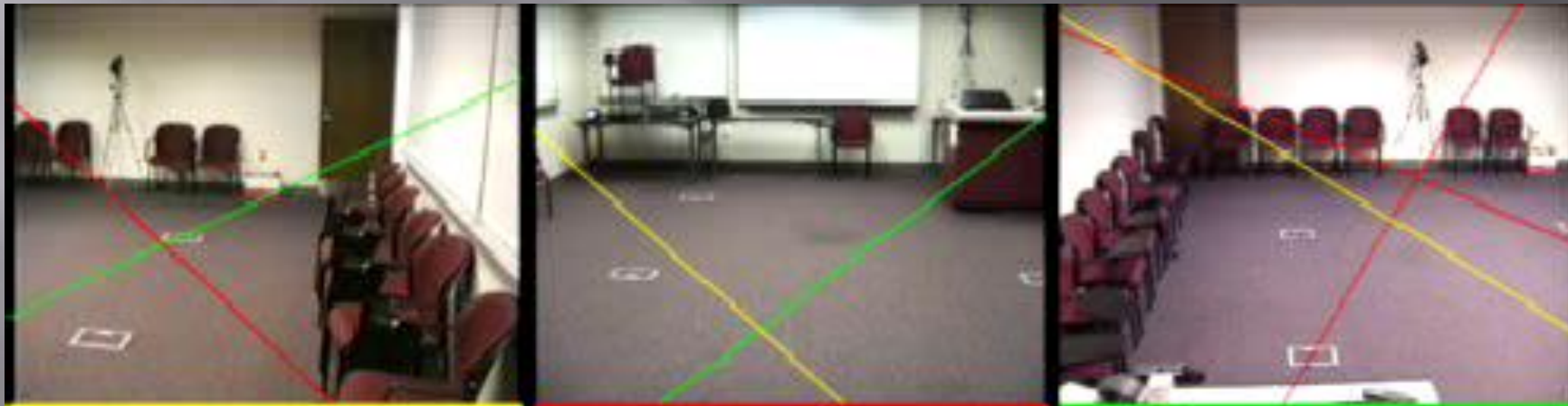
Copyrights Mubarak Shah, UCF

Results

Camera1

Camera2

Camera3

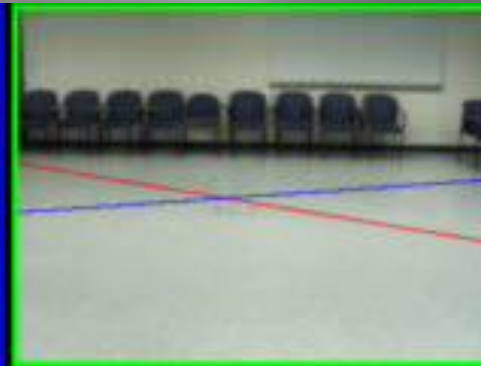


Multiple Fixed & Overlapping Cameras Tracking

Camera1

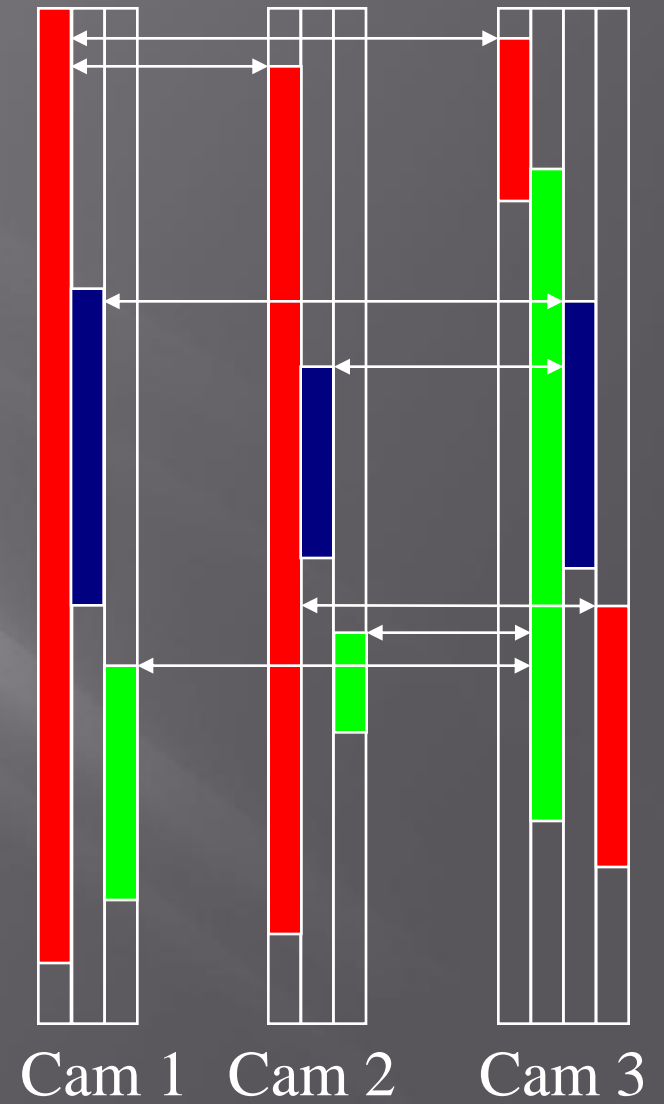
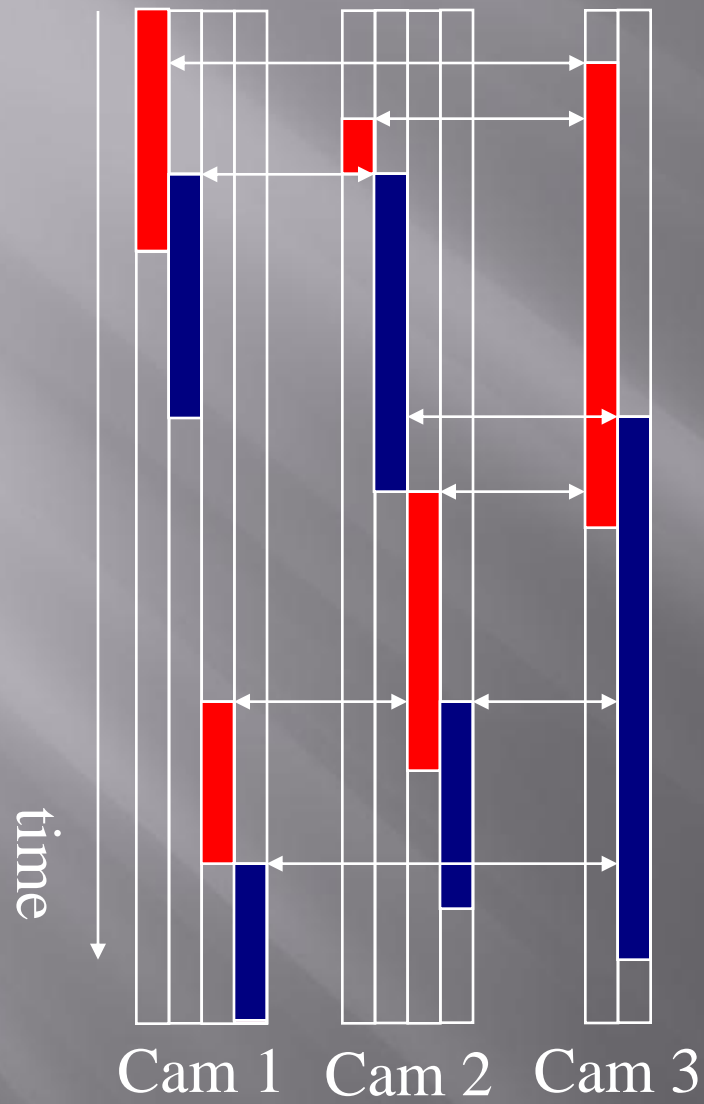


Camera2



Camera3





TRACKING ACROSS MULTIPLE FIXED NON- OVERLAPPING CAMERA

Omar Javed, Khurram Shafique and
Mubarak Shah

ICCV2003, CVPR2005



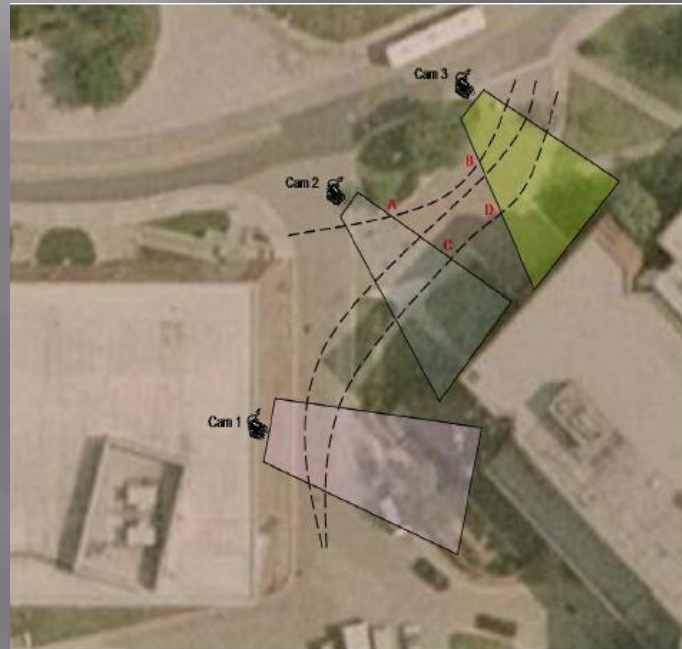
University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Tracking Across Multiple Cameras

- Task Definition:
 - To maintain the identity of objects as they move across multiple cameras



Motivation

- Wide area surveillance requires tracking over disjoint views.
 - Camera Resolution
 - Occlusion due to scene structures
- Luxury of calibrated cameras is not available in most cases.



Introduction

- Problems
 - Successive observations of object might be widely separated in space and time.
 - Appearance of objects change across cameras
 - Difference in illumination
 - Difference in camera parameters (focal length, gain, response function)
 - Difference in pose of object



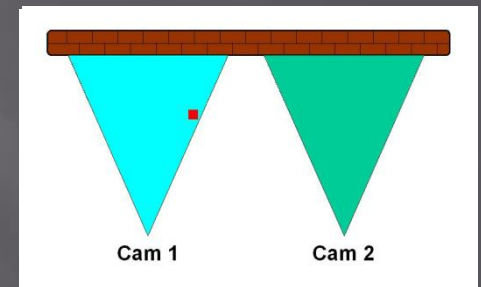
Key Observations

- ▣ Due to physical and practical constraints, some paths are more likely to be taken than others, e.g., roads, walkways, corridors.
- ▣ The observed motion pattern provides clues about inter-camera relationship.
- ▣ The transformation between the color distributions of two views of an object lies in a low dimensional space.



Features for Multi-Camera Tracking

- Location of exit and entry from one camera to another.
- Time Interval between observations
 - Magnitude of motion of object.
 - Direction of motion
 - Location of exit and entry
- Appearance.



Probabilistic Framework

- We have a system of l cameras C_1, C_2, \dots, C_l .
- There are n objects p_1, p_2, \dots, p_n that generate a sequence of tracks in the system of cameras in successive time instances.
- $O_j = \{O_{j,1}, O_{j,1}, \dots, O_{j,n}\}$ be the set of observations that were observed by C_j .
- Each observation $O_{j,a}$ is composed of two independent feature sets, which are appearance $O_{j,a}(A)$ and space-time (location, time, velocity) features $O_{j,a}(ST)$.
- Single camera tracks are available

Probabilistic Framework

MAP Solution: $K' = \arg \max_{K \in \Sigma} P(K | O)$

For $K = \{k_{a,b}^{c,d}\}$ in solution space Σ

$$P(K | O) = P(K | O_1, O_2, \dots, O_r) = \prod_{k_{i,a}^{j,b} \in K} P(k_{i,a}^{j,b} | O_{i,a}, O_{j,b})$$

Assuming independence

$$P(K | O) = \prod_{k_{i,a}^{j,b} \in K} \frac{P(O_{i,a}(A), O_{j,b}(A) | k_{i,a}^{j,b}) P(O_{i,a}(ST), O_{j,b}(ST) | k_{i,a}^{j,b}) P(k_{i,a}^{j,b})}{P(O_{i,a}, O_{j,b})}$$

Using Bayes Law and Assuming Independence
between appearance and spatio-temporal observations



Probabilistic Framework

- Maximizing the following term will give us the solution

$$K' = \arg \max_{K \in \Sigma} \sum_{k_{i,a}^{j,b} \in K} P(O_{i,a}(A), O_{j,b}(A) | k_{i,a}^{j,b}) P(O_{i,a}(ST), O_{j,b}(ST) | k_{i,a}^{j,b}) P(C_i, C_j)$$



Learning Probability Density Functions

- Learning phase
 - Assumption of known correspondences.
 - One solution is to use only appearance matching for correspondences. Note that only weak or ambiguous matches can be discarded during the training phase.
 - Estimate spatio-temporal and appearance pdfs from the observed data.

Estimating Spatio-Temporal Pdf

- Features:
 - Entry and Exit Locations and Cameras
 - Velocity
 - Inter-Camera Travel Time
- Parzen Windows for density estimation
- For a sample S , consisting of ' n ' data points x_1, x_2, \dots, x_n , the Parzen estimate is given as

$$\hat{\rho}(x) = \frac{1}{n} |H|^{-\frac{1}{2}} \sum_{i=1}^n K(H^{-\frac{1}{2}}(x - x_i))$$



Estimating Appearance Pdf

- Goal: Learn the change in appearance (color) of an object as it moves from one camera to other.



Same person in two different cameras



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Effectiveness of Subspace Learning

$O_{i,a}$

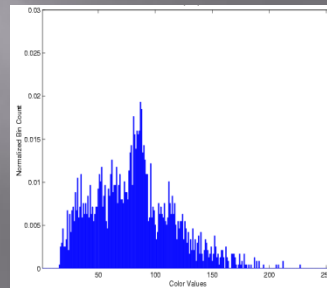


$O_{j,b}$

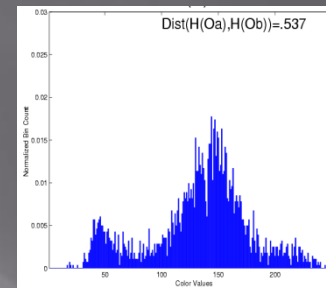


Same person in
cameras C_i & C_j

$H(O_{i,a})$



$H(O_{j,b})$



Red Channel
Histogram Similarity

Bhattacharray Distance = $d = 0.537$

$O_{i,a}$

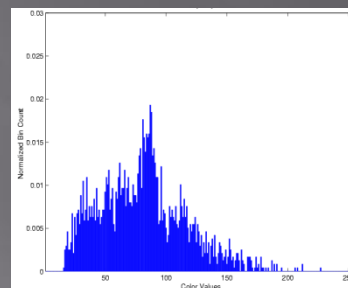


$O_{j,c}$

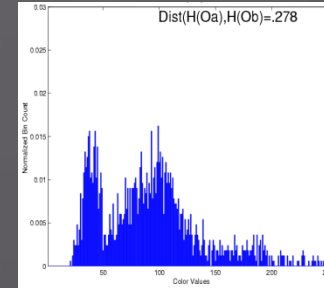


Different persons in
cameras C_i & C_j

$H(O_{i,a})$



$H(O_{j,c})$



Red Channel
Histogram Similarity

$d = 0.278$



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Effectiveness of Subspace Learning

$O_{i,a}$

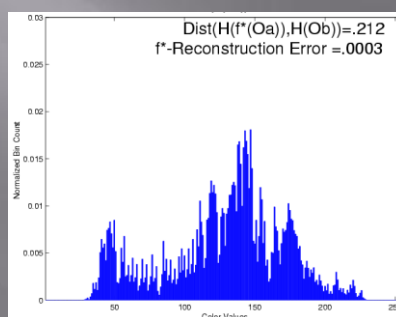


$O_{j,b}$

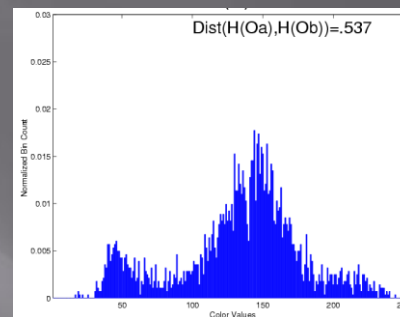


Same person in
cameras C_i & C_j

$H(f_1(O_{i,a}))$



$H(O_{i,b})$



Histogram Similarity
after Transformation

$$d = 0.212$$

$O_{i,a}$

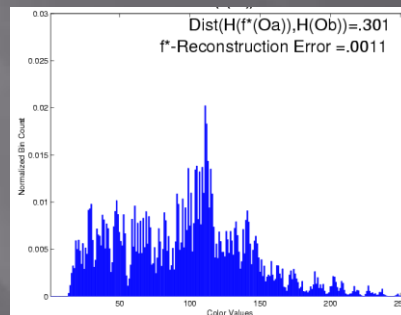


$O_{j,c}$

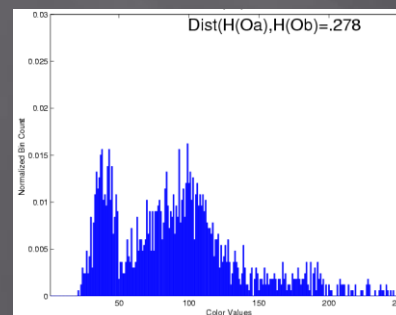


Different persons in
cameras C_i & C_j

$H(f_2(O_{i,a}))$



$H(O_{i,c})$



Histogram Similarity
after Transformation

$$d = 0.301$$

Copyright Mubarak Shah, UCF



University of
Central Florida

VISION

Effectiveness of Subspace Learning

$O_{i,a}$

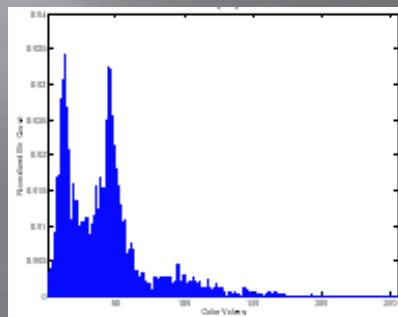


$O_{j,b}$

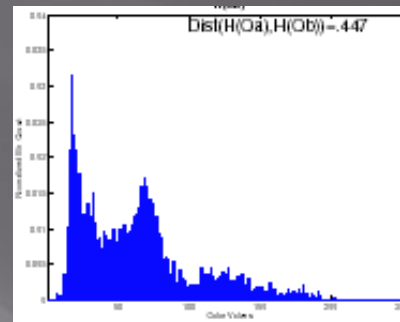


Same person in cameras C_i & C_j

$H(O_{i,a})$



$H(O_{j,b})$



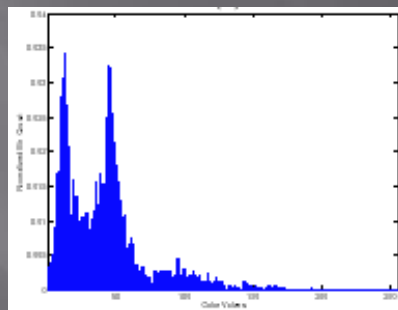
Blue Channel Histogram Similarity

$$d = 0.447$$

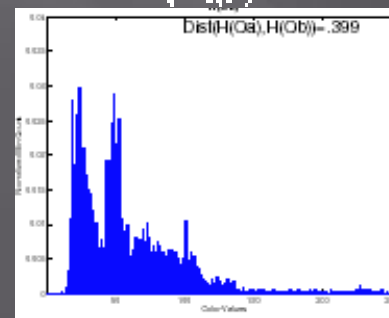


Different persons in cameras C_i & C_j

$H(O_{i,a})$



$H(O_{j,c})$



Blue Channel Histogram Similarity

$$d = 0.399$$



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Effectiveness of Subspace Learning

$O_{i,a}$

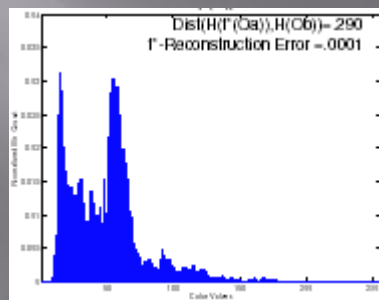


$O_{j,b}$

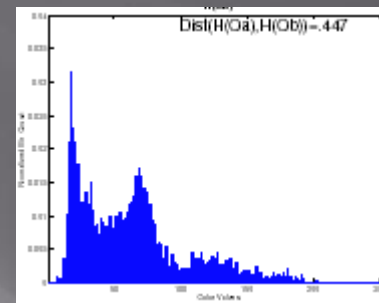


Same person in cameras C_i & C_j

$H(f_1(O_{i,a}))$



$H(O_{i,b})$



Histogram Similarity after Transformation

$$d = 0.290$$

$O_{i,a}$

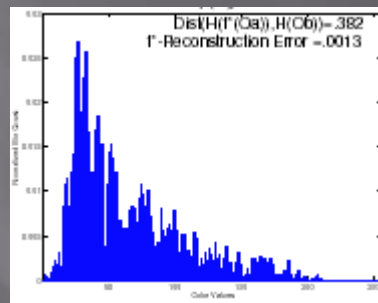


$O_{j,c}$

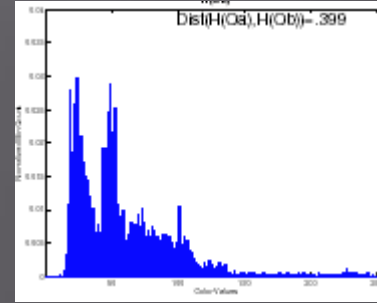


Different persons in cameras C_i & C_j

$H(f_2(O_{i,a}))$



$H(O_{i,c})$



Histogram Similarity after Transformation

$$d = 0.382$$

Copyright Mubarak Shah, UCF

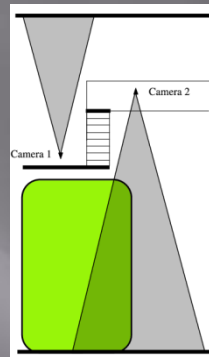


University of
Central Florida

VISION

Results

Camera Setup
for Experiment # 1



A Clip from the test sequence



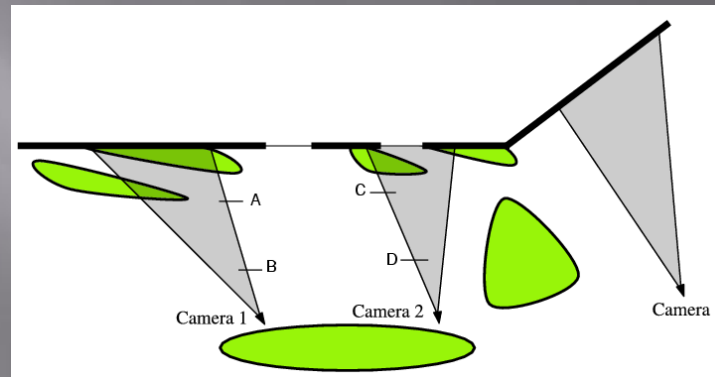
University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Results

Camera Setup for Experiment # 2



A Clip from the test sequence



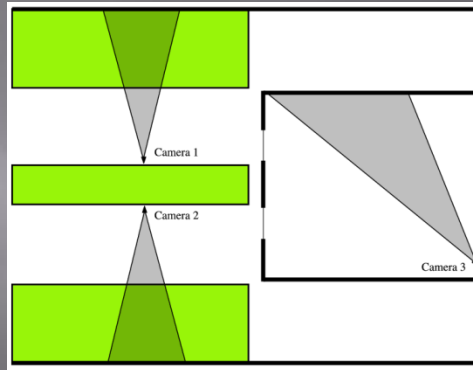
University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Results

Camera Setup
for Experiment # 2



A Clip from the test sequence

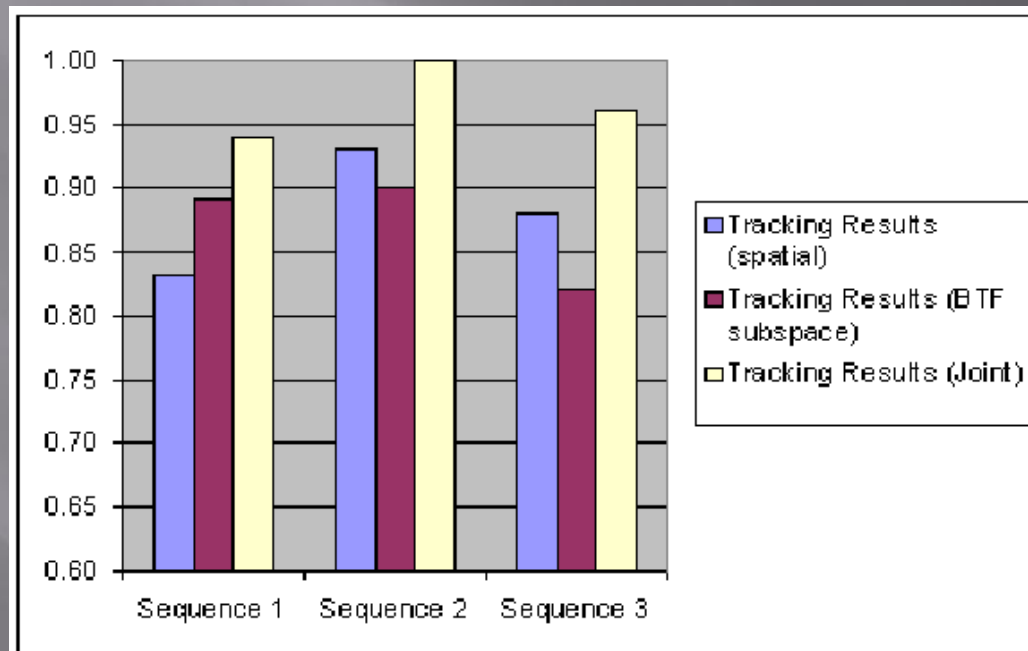


University of
Central Florida

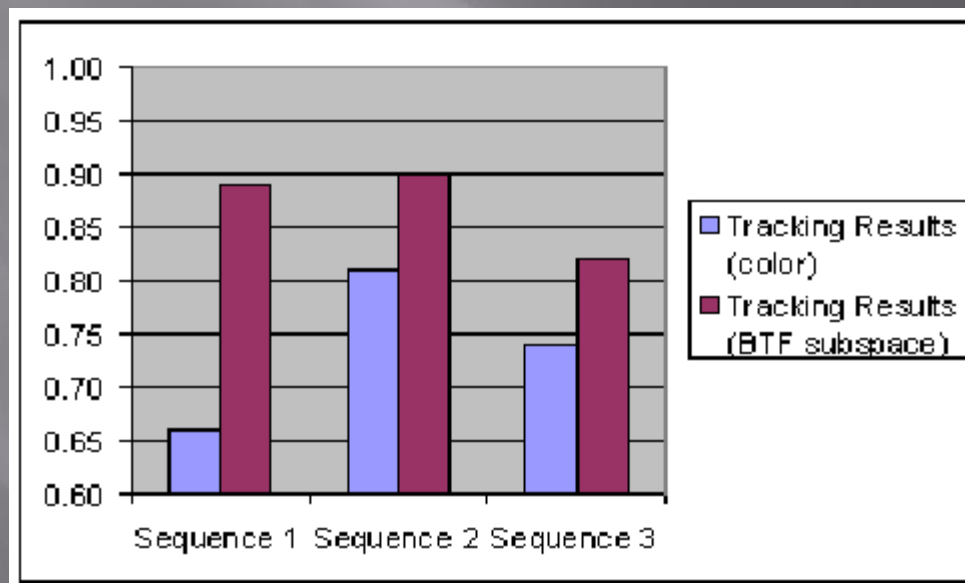
VISION

Copyright Mubarak Shah, UCF

Results (Tracking Accuracy)



Results (Comparison with direct Color Matching)



TRACKING ACROSS MULTIPLE MOVING AIRBORNE CAMERAS

Yaser Sheikh & Mubarak Shah
ICCV 2005

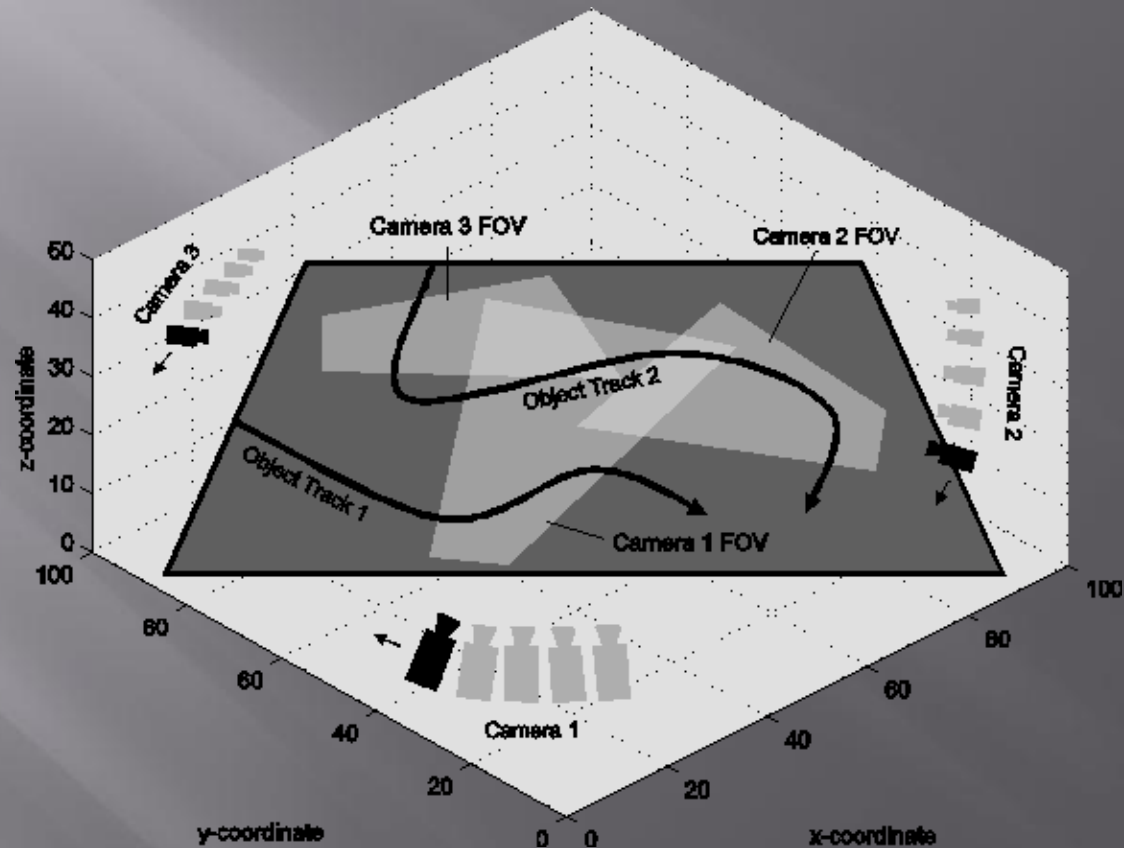


University of
Central Florida

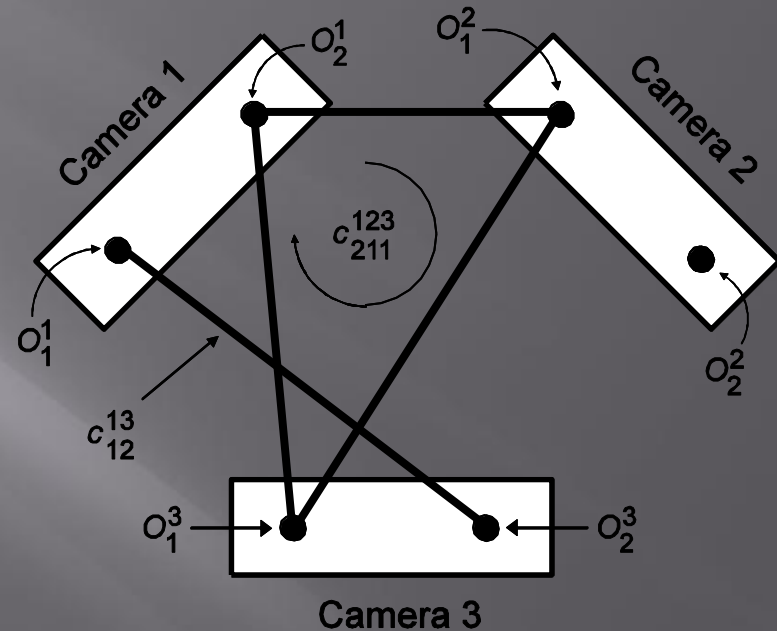
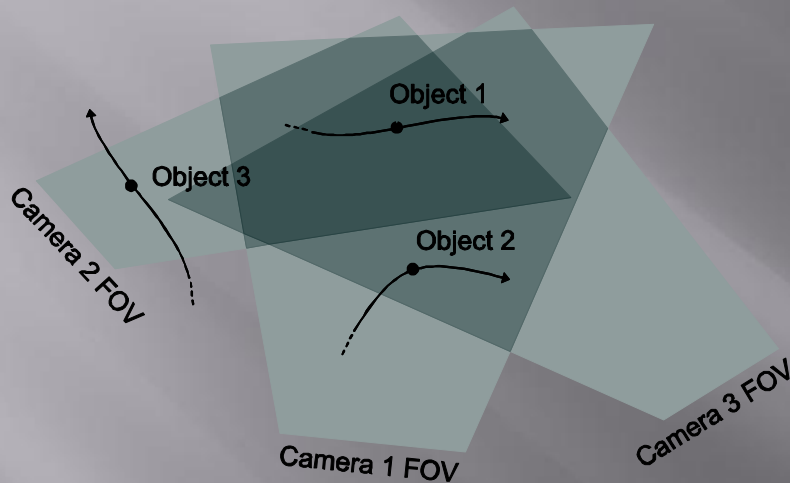
VISION

Copyrights Mubarak Shah, UCF

Motion in the Forest

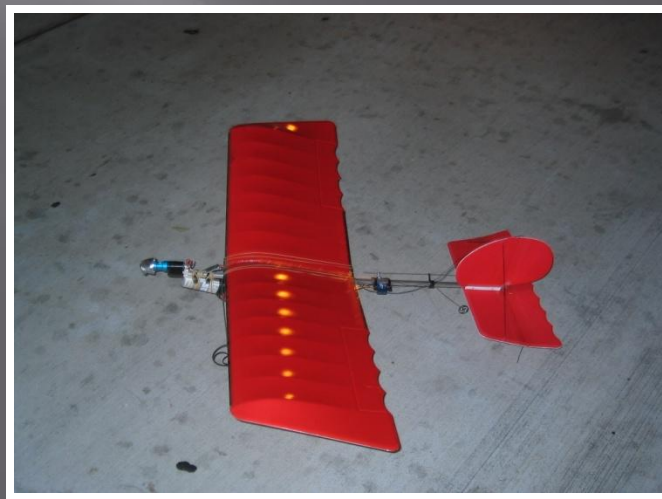
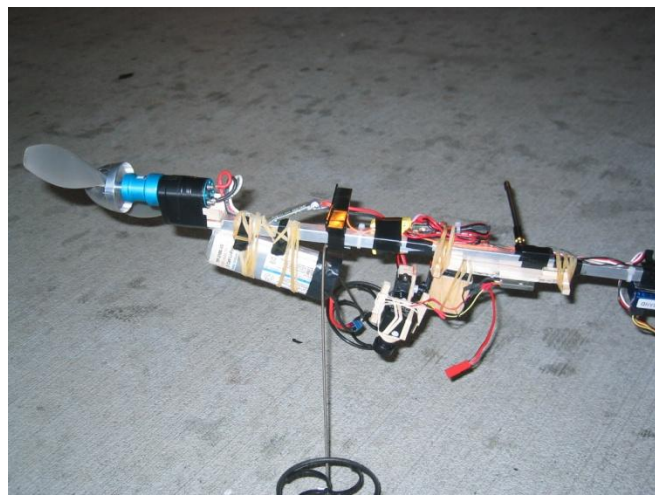
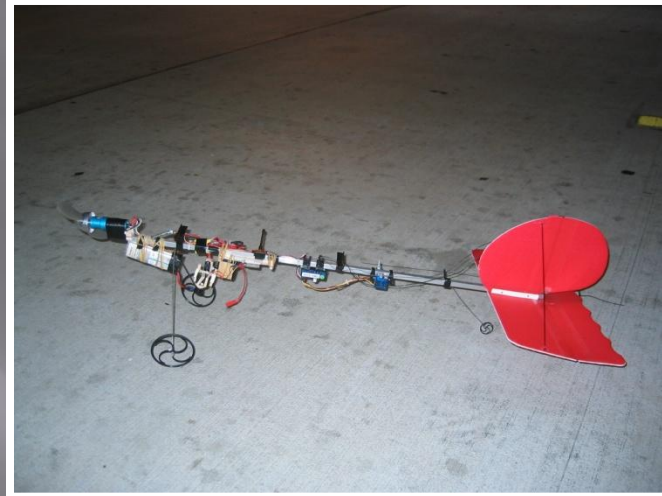


Graph Theoretic Formulation



- Three Objects, three cameras
- Object 1 is visible in all three cameras
- Object 2 is visible in Camera 1 and Camera 3
- Object 3 is visible only in Camera 2

UCF's UAVs

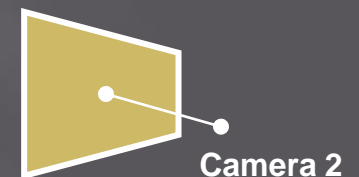
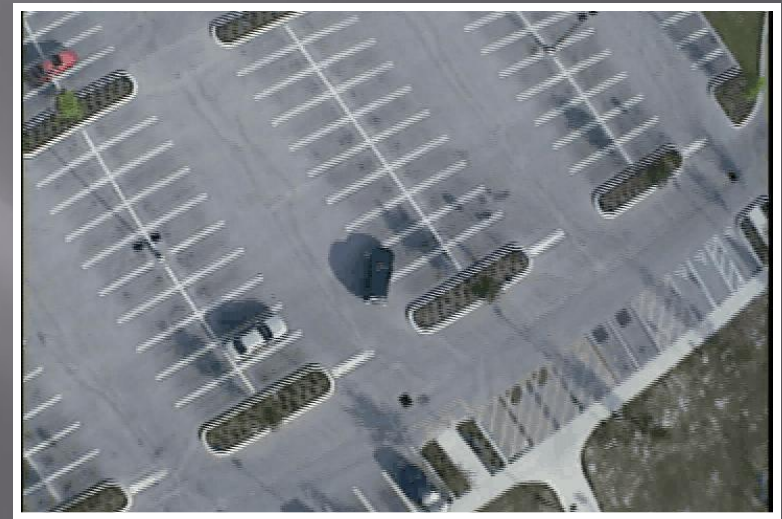
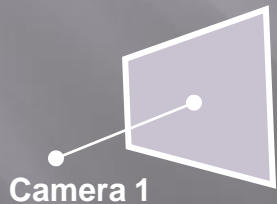


University of
Central Florida

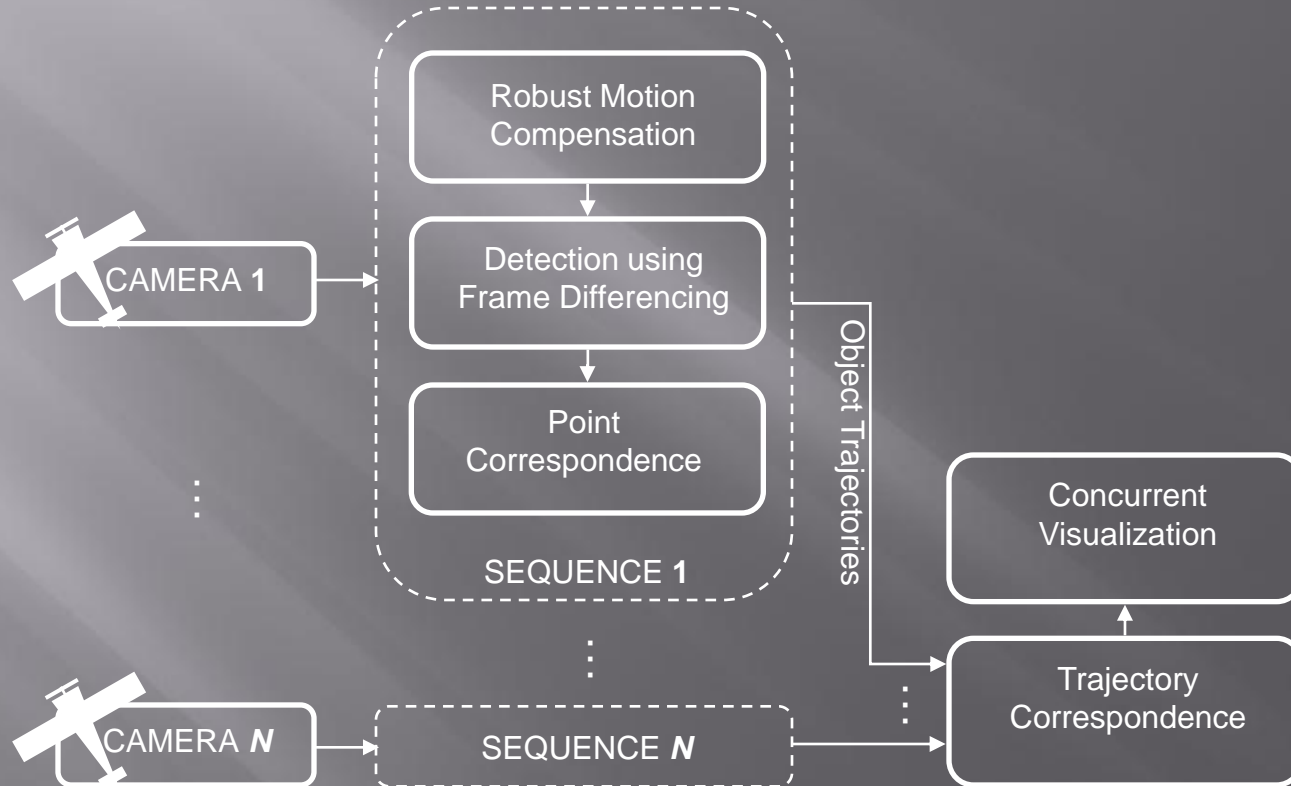
VISION

Copyright Mubarak Shah, UCF

Example Input Data



System Overview



Assumptions

- ▣ Assumption of Scene Planarity
 - Validated by altitude of aerial vehicle
 - Reasonable deviations do not affect solution
- ▣ Spatiotemporal Overlap of Fields of View
 - Objects are simultaneously visible in two cameras at a time (for all pairs of cameras).



Compensating Camera Motion



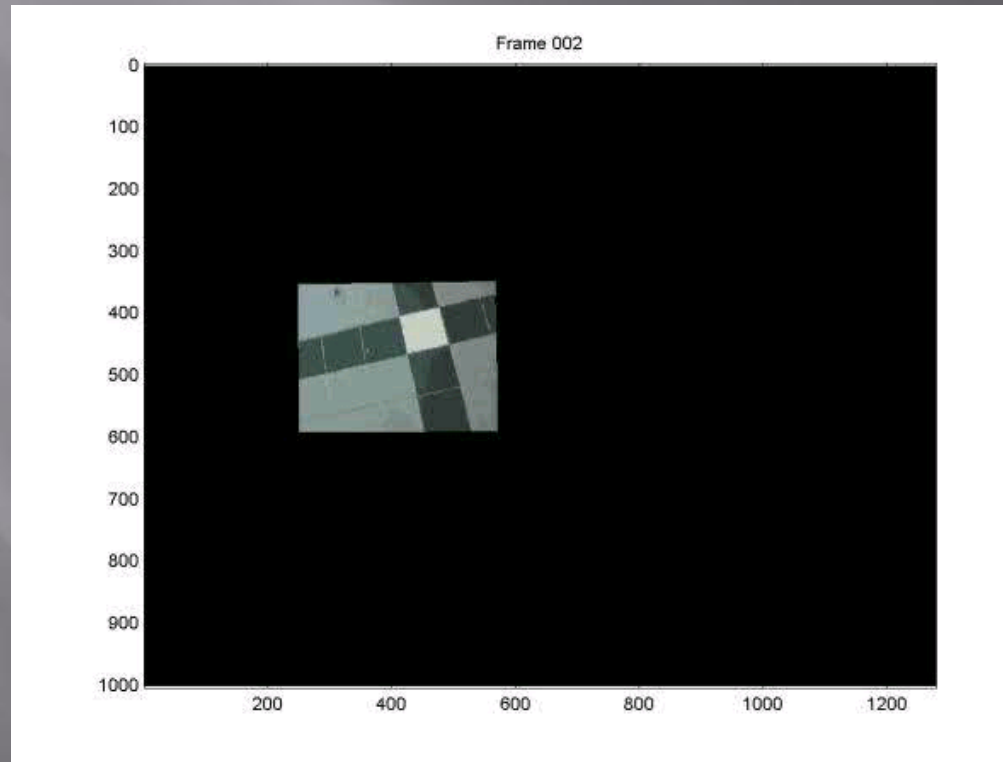
Original Sequence



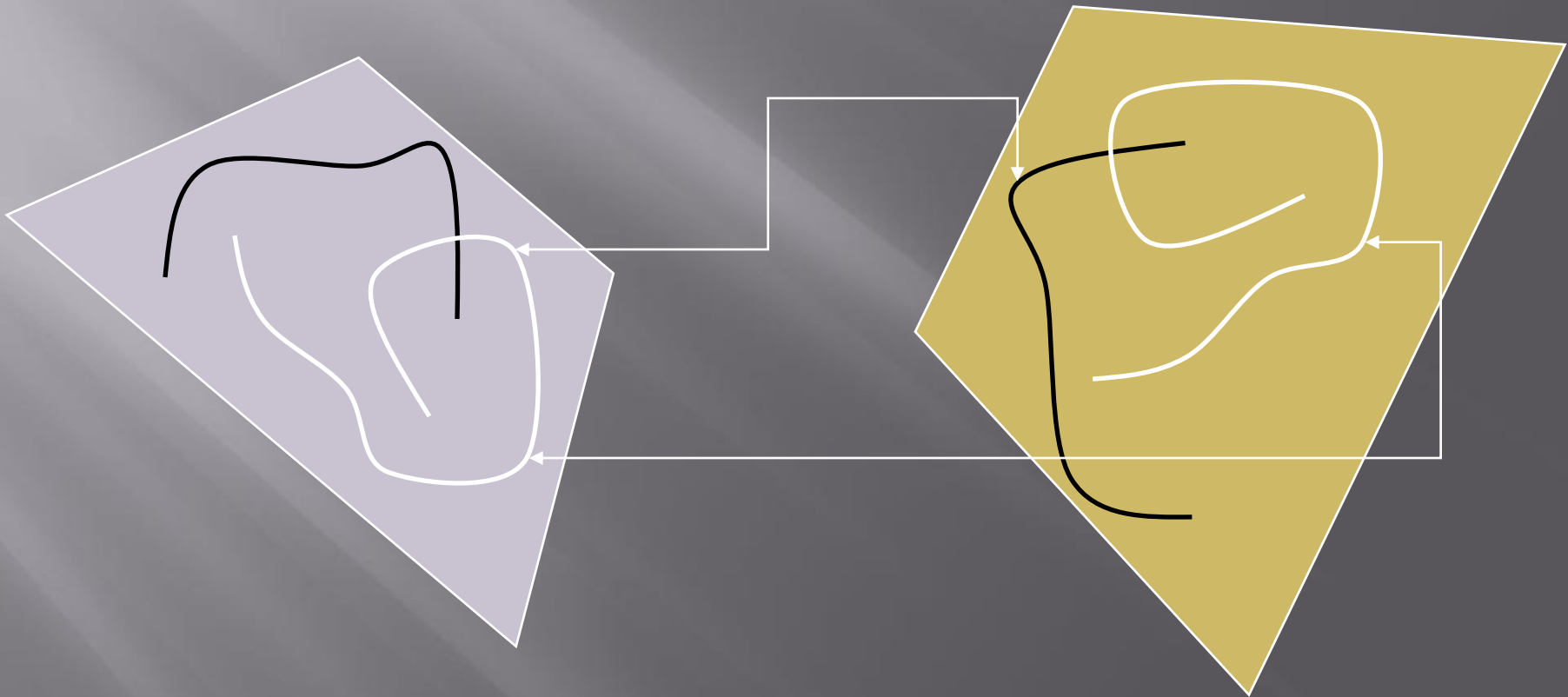
Motion Compensated

Detection and Tracking

- ▣ We are not interested in solving the problem of tracking *within* cameras in this work



Similarity Metric



Data Modeling

- Data Two trajectories, \mathbf{x}_1 and \mathbf{x}_2 . Each point modeled as a random variables with independent Gaussian noise,

$$\mathbf{x}_i = \mathbf{H}_i \mathbf{x} + \mathbf{e}_i$$

Homography \mathbf{H}_{12} exists between the two trajectories since we model the scene as a plane.

- The pair of tracks \mathbf{x}_1 and \mathbf{x}_2 related *exactly* by \mathbf{H}_{12}
- Maximum Likelihood Estimate of \mathbf{H}

$$\arg \max_{\mathbf{H}} \prod_i \frac{1}{2\pi\sigma^2} e^{-(\mathbf{x}_i - \mathbf{H} \mathbf{x}_j)^T \mathbf{C}^{-1} (\mathbf{x}_i - \mathbf{H} \mathbf{x}_j)}$$

$$\arg \max_{\mathbf{H}} \prod_i \left(\frac{1}{2\pi\sigma^2} \right) e^{-(\mathbf{x}_i - \mathbf{H} \mathbf{x}_j)^T \mathbf{C}^{-1} (\mathbf{x}_i - \mathbf{H} \mathbf{x}_j)}$$



Correspondence

- We wish to compute,

$$\prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{d_i^2}{2\sigma^2}\right)$$

- Taking the log,

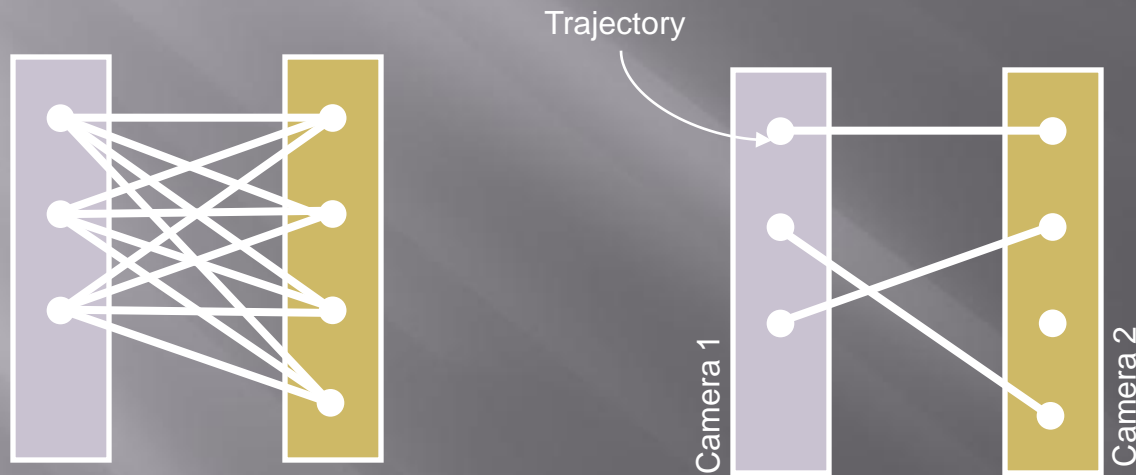
$$-\sum_{i=1}^N \left[\ln(\sqrt{2\pi}\sigma) + \frac{d_i^2}{2\sigma^2} \right]$$

- Intuitively this is using the estimate of the statistical *mean* of the reprojection error at each point.
- If outliers exists, robust estimators, such as the *median* of the reprojection error can be used



Correspondence Across 2 Cameras

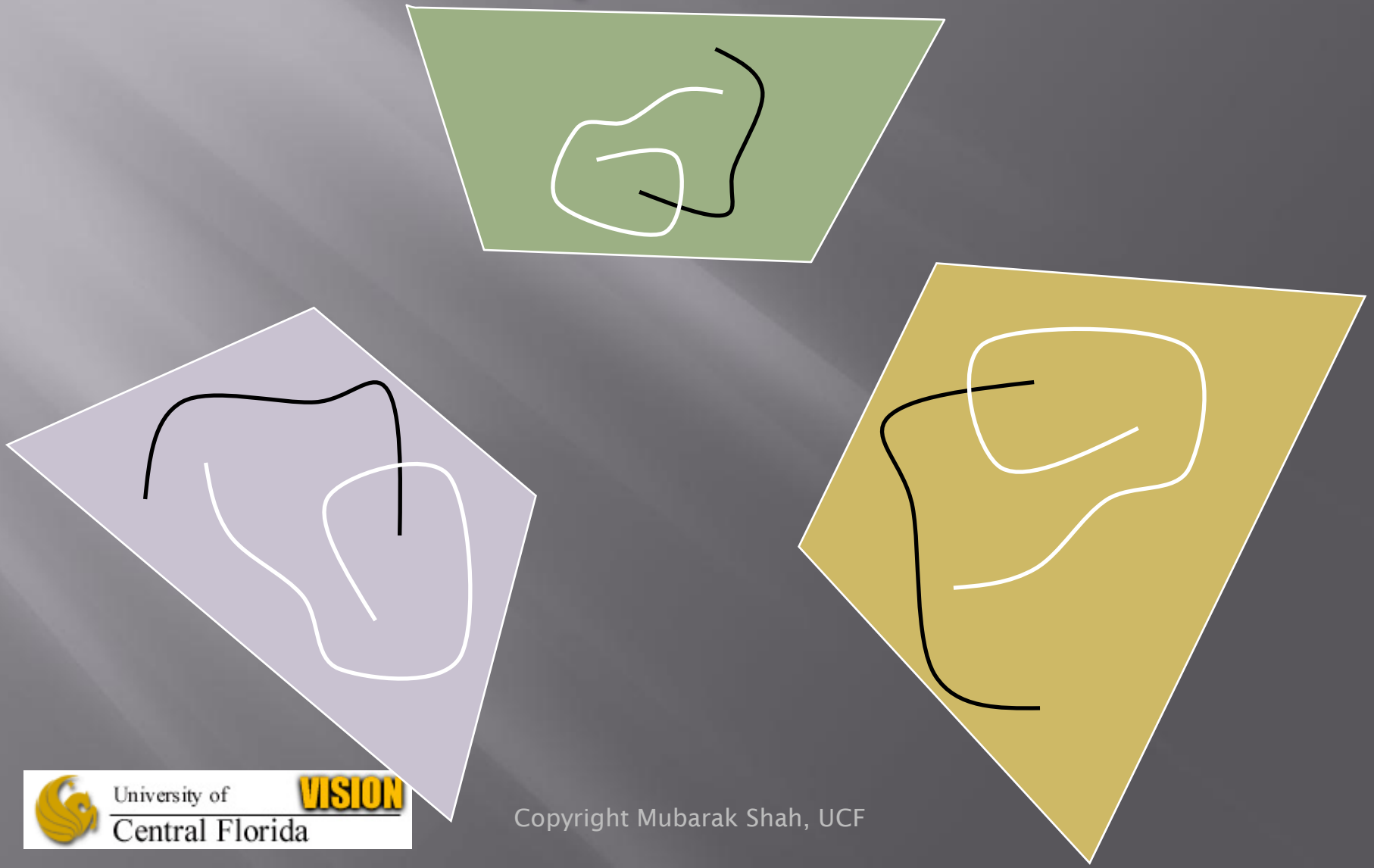
- Edge weight is $P(C_{ij}^{pq} | \mathbf{x}_t^p, \mathbf{x}_j^q)$



- Optimal correspondence in a ML sense can be found using bipartite matching



Trajectory Correspondence Across Multiple Cameras



Global Correspondence

- For multiple cameras, we need to find the correspondence C such that,

$$\sum_{p,q \in \mathcal{C}} \sum_{i,j \in \mathcal{C}} c_{ij}^{pq} \mathbf{x}_i^p \cdot \mathbf{x}_j^q = \max$$

- Where,

\mathbf{x}_i^p

Trajectory i in camera p

c_{ij}^{pq}

Correspondence hypotheses b/w \mathbf{x}_i^p and \mathbf{x}_j^q

C

Global correspondence hypotheses

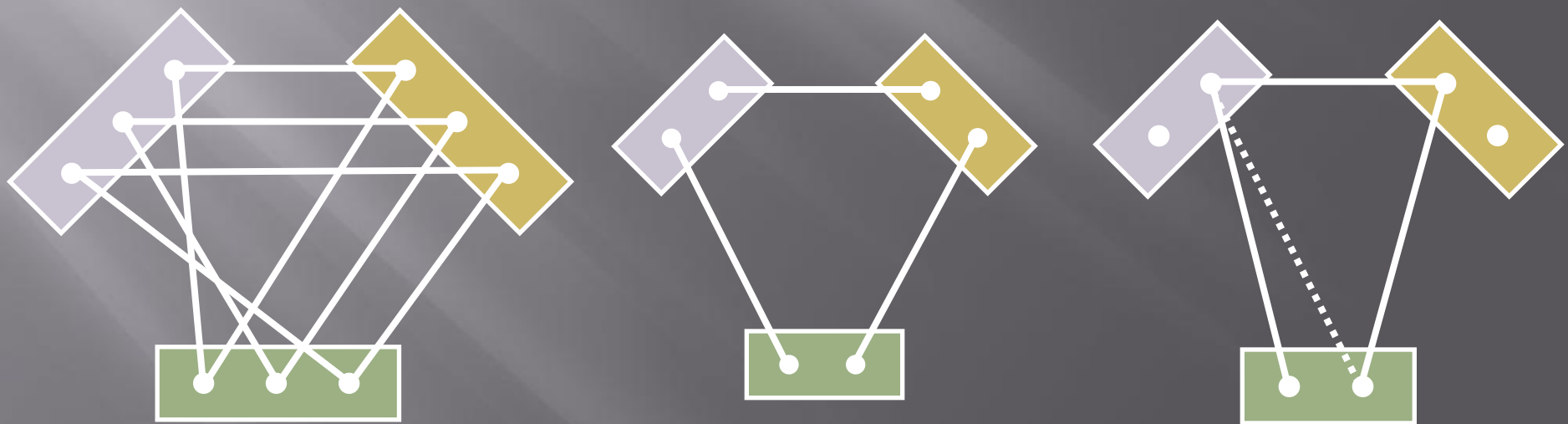
\mathcal{C}

Space of global correspondence solutions



Correspondence Across Cameras

- ▣ For greater than two cameras complexity increases
- ▣ In addition another constraint needs to be satisfied: transitive closure



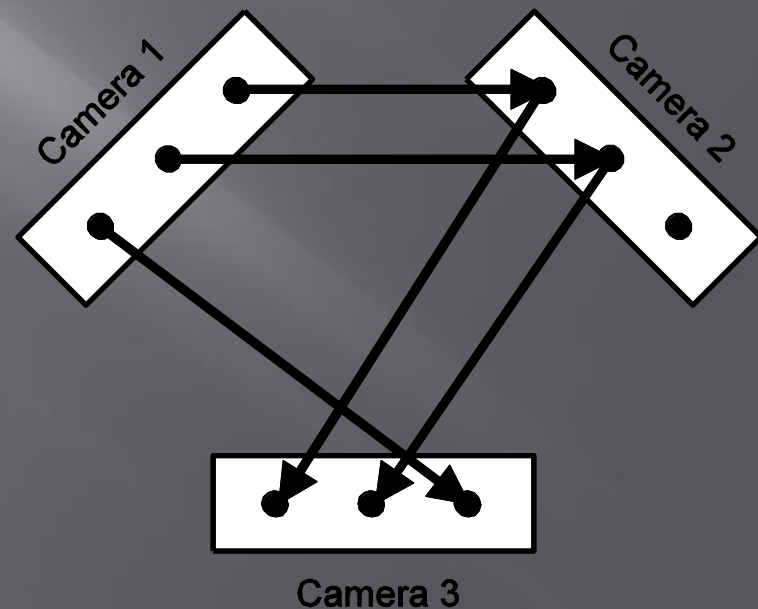
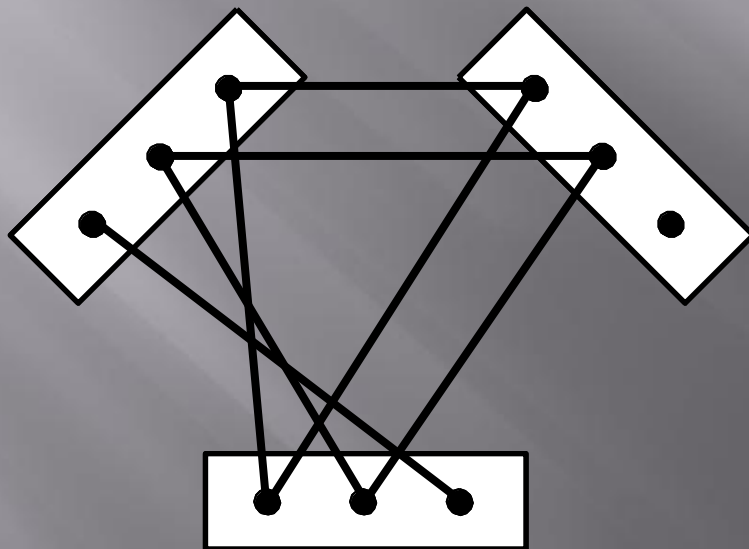
Correspondence Across Cameras

- ▣ For two sequences, the optimal solution can be obtained by maximum bipartite matching
- ▣ For $k \geq 3$, the problem of finding correspondences is known to be NP Hard.
- ▣ Therefore, we consider a reformulation of the problem as a directed graph.
- ▣ The direction comes from arbitrarily enumerating the cameras.

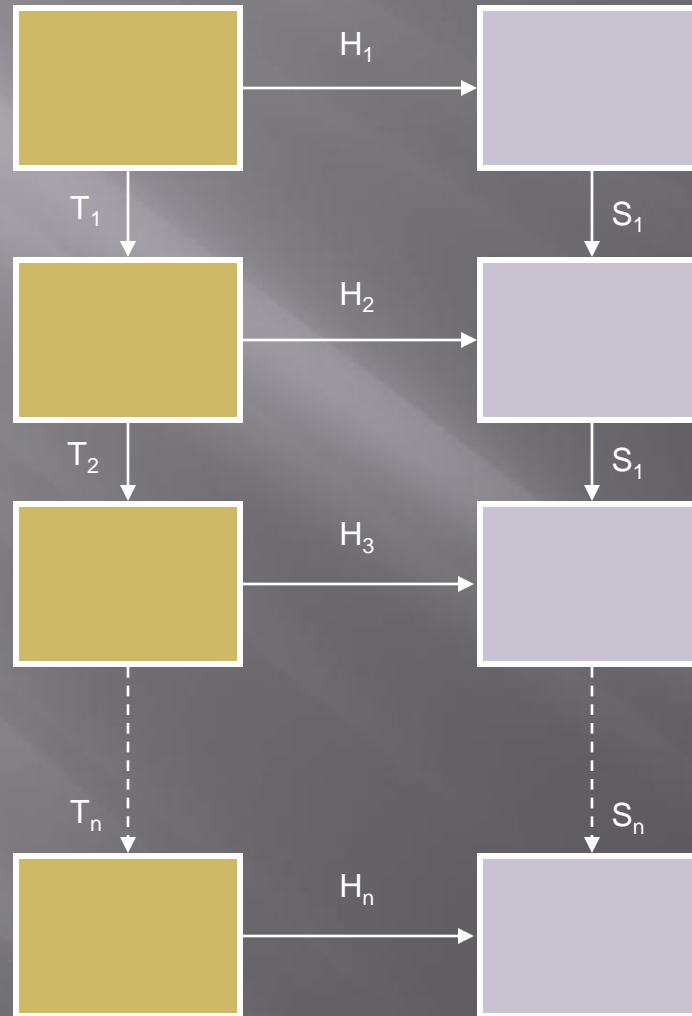


Correspondence Across Cameras

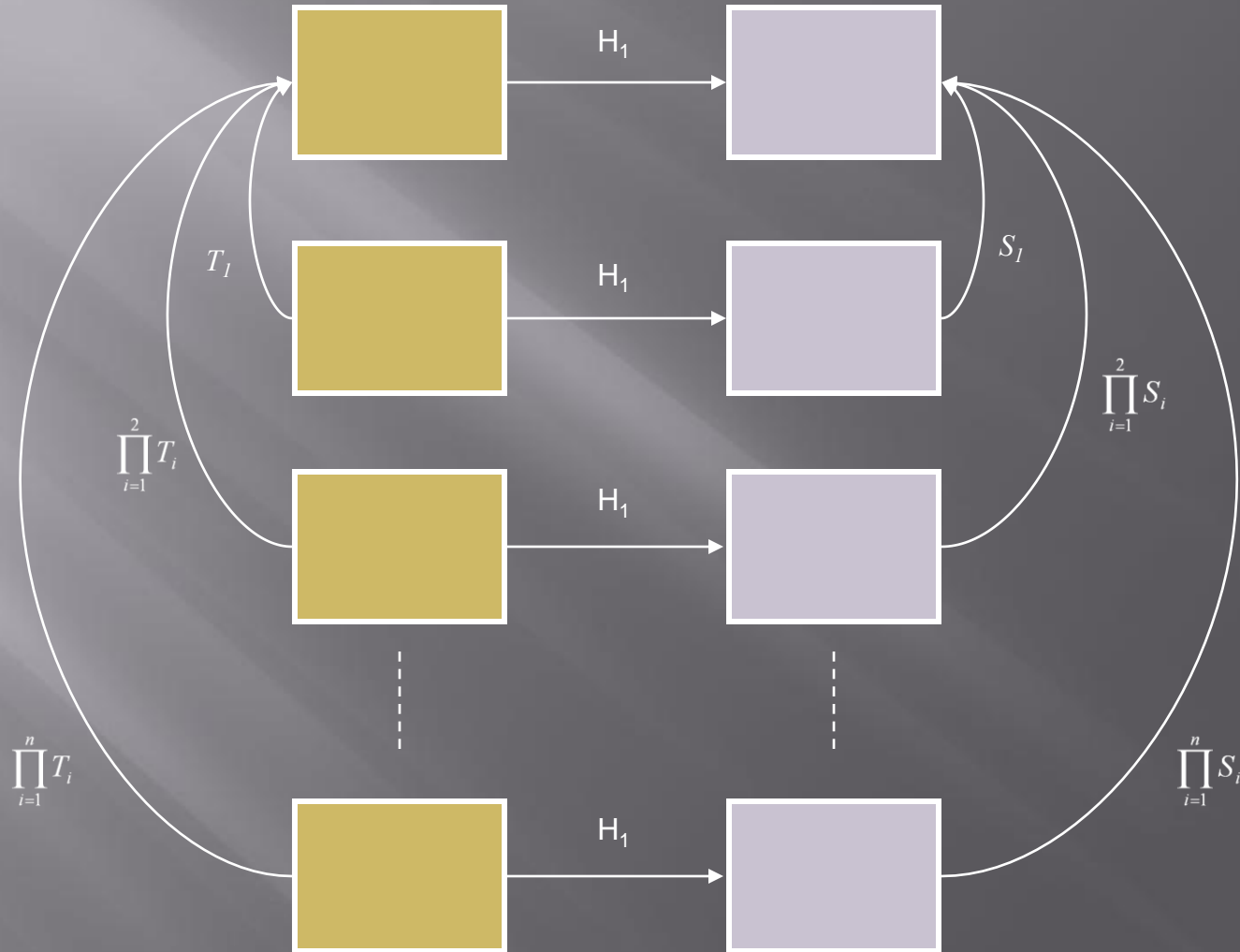
- A polynomial time solution exists to approximate the maximum matching (Shafique and Shah, TPAMI 2005)



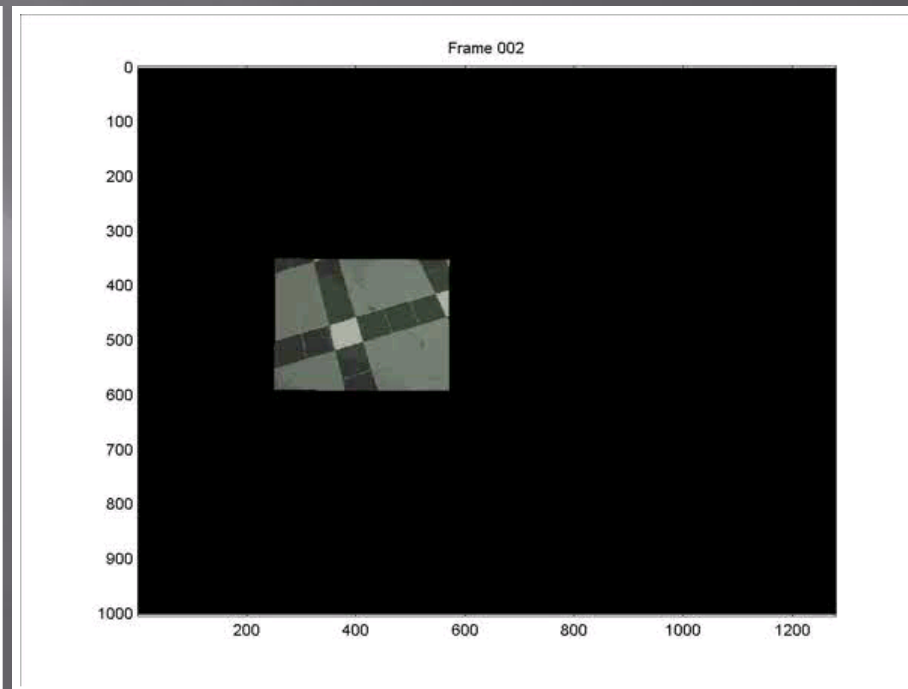
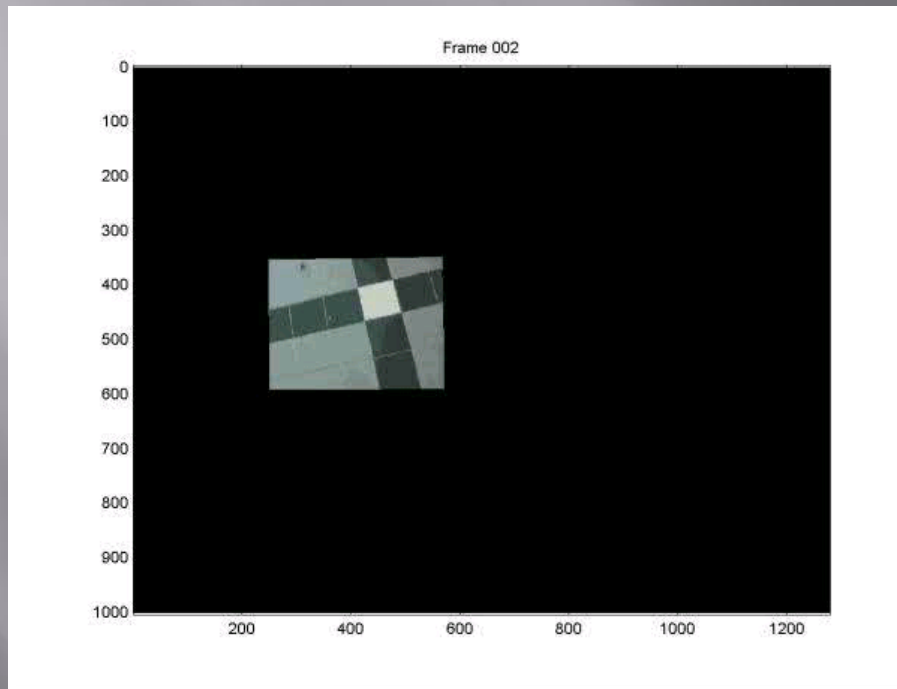
Concurrent Visualization



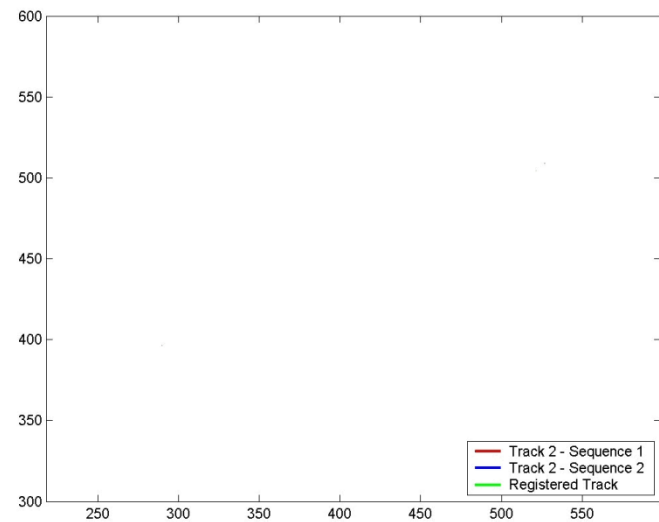
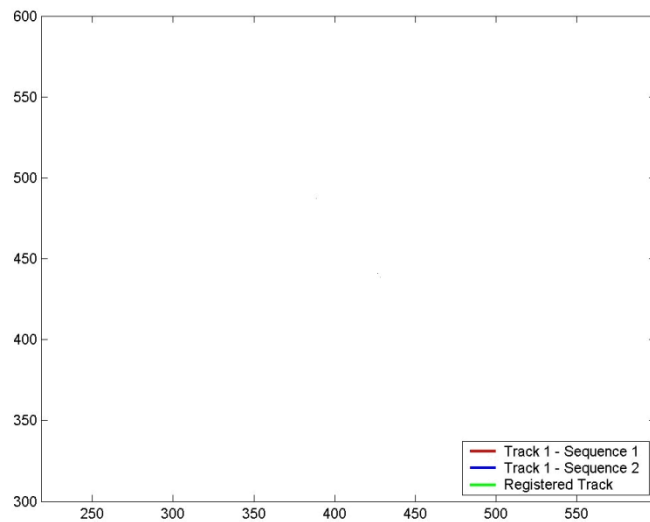
Concurrent Visualization



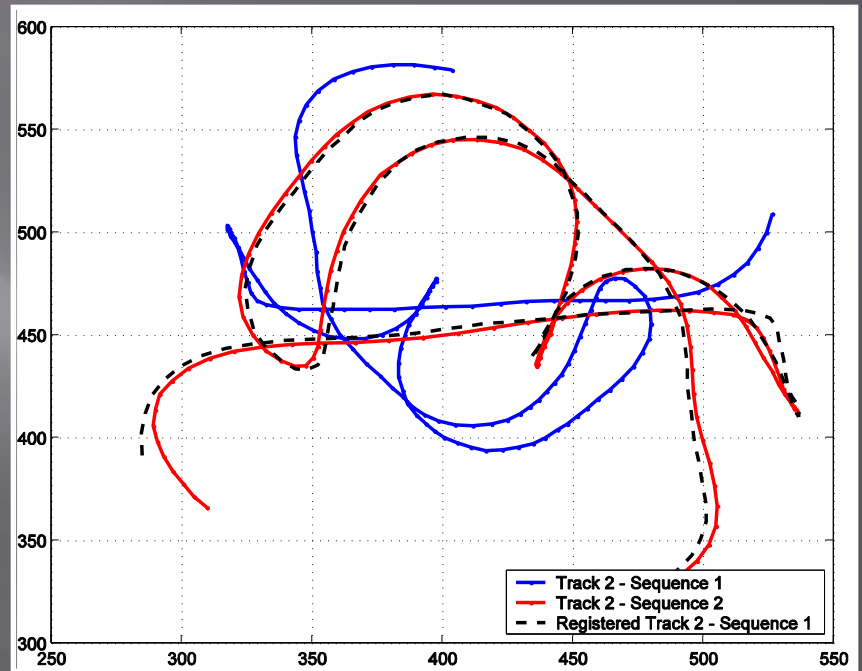
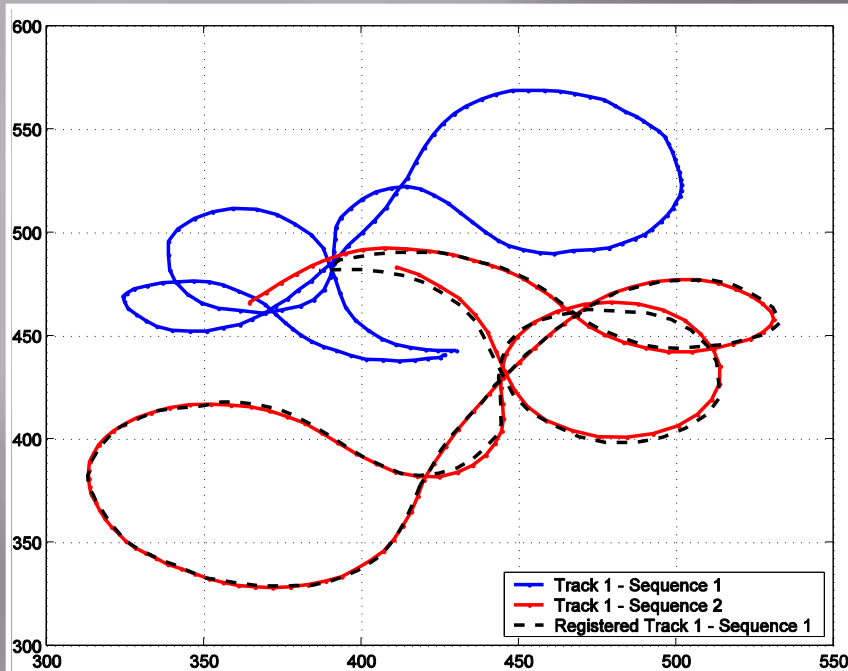
Controlled Sequence 1



Tracks Corresponded

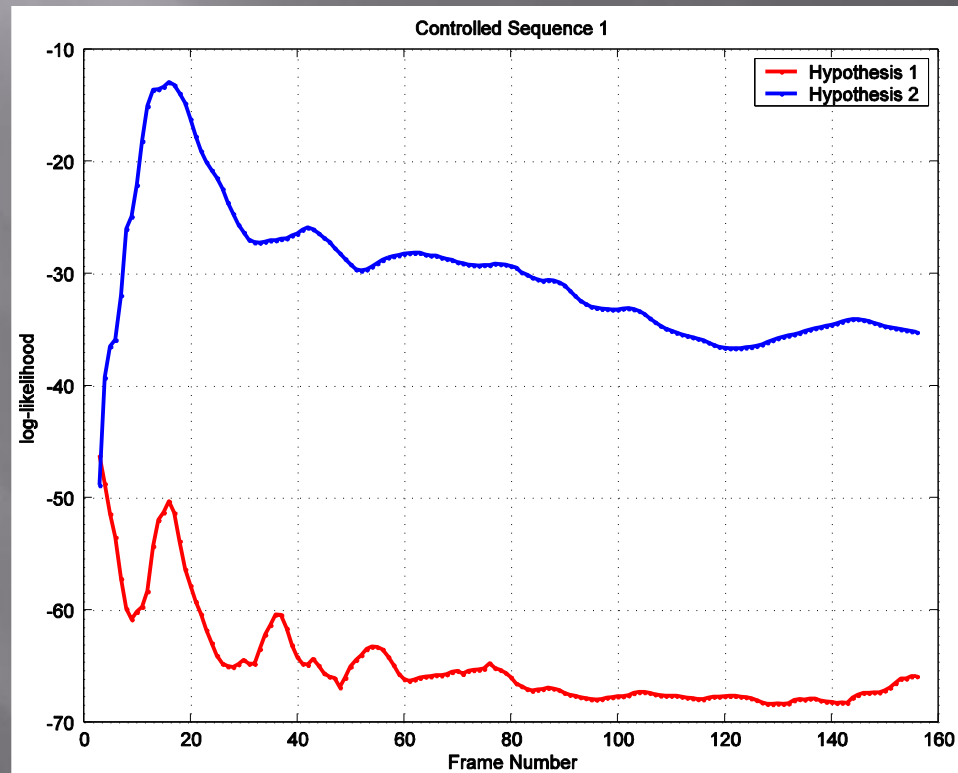


Still Picture...

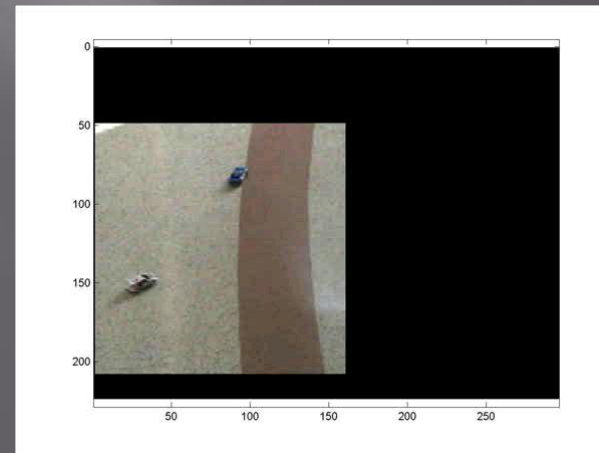
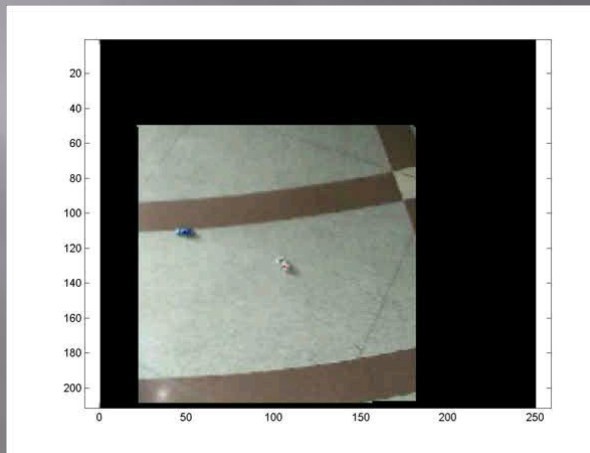
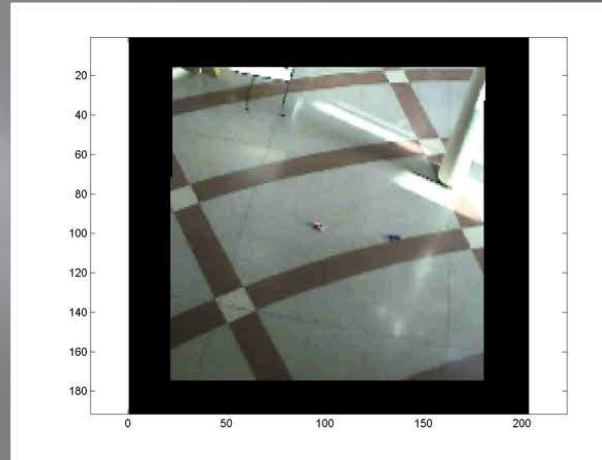


Hypotheses

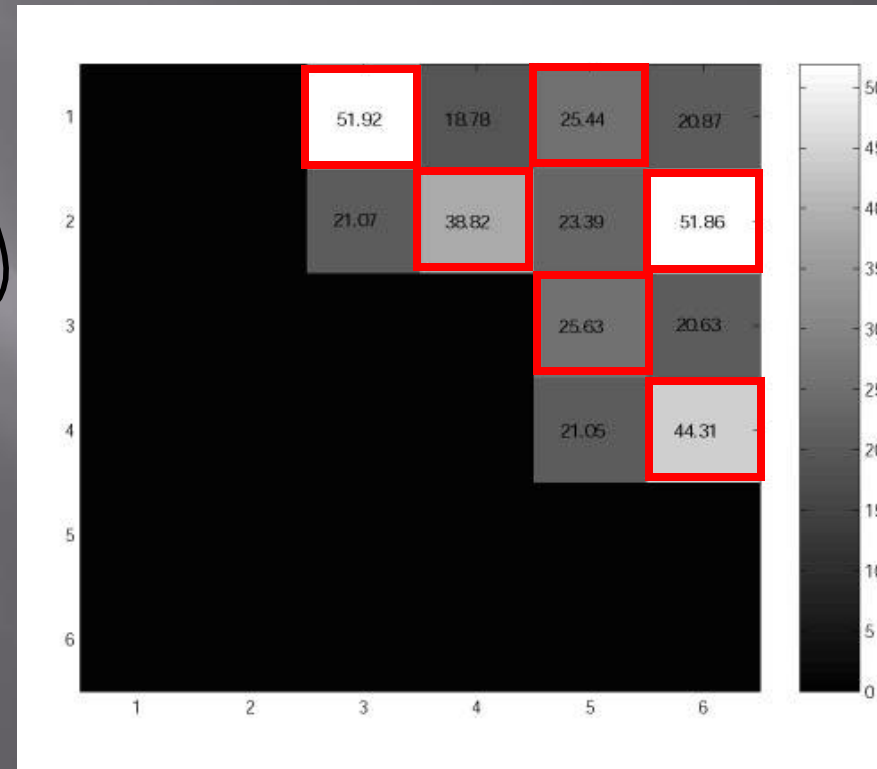
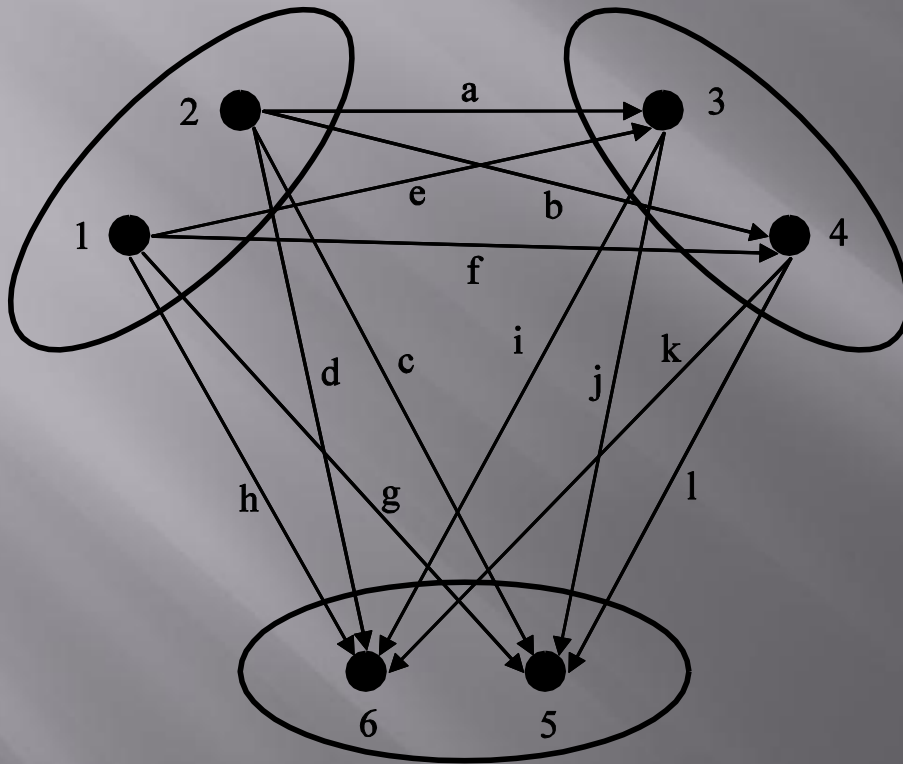
- Two hypotheses are tested



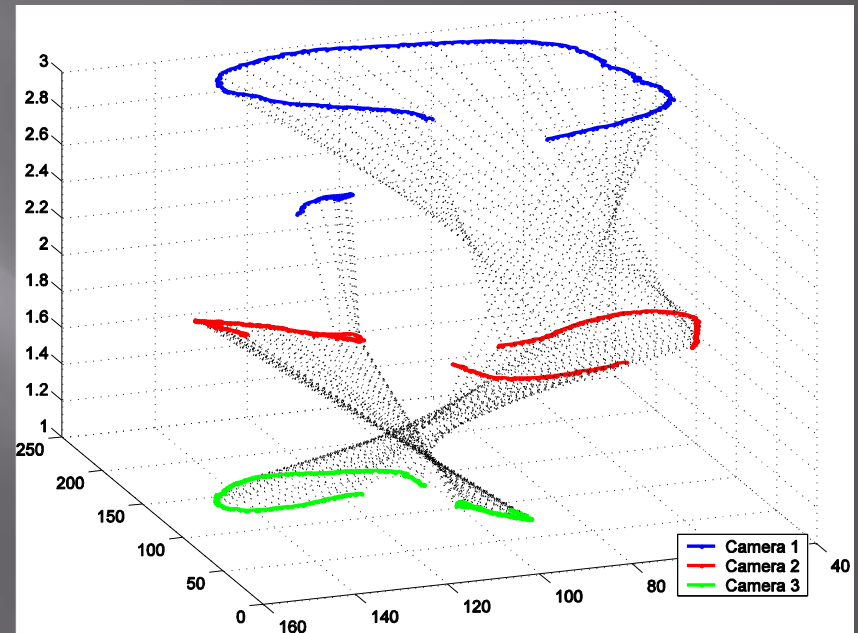
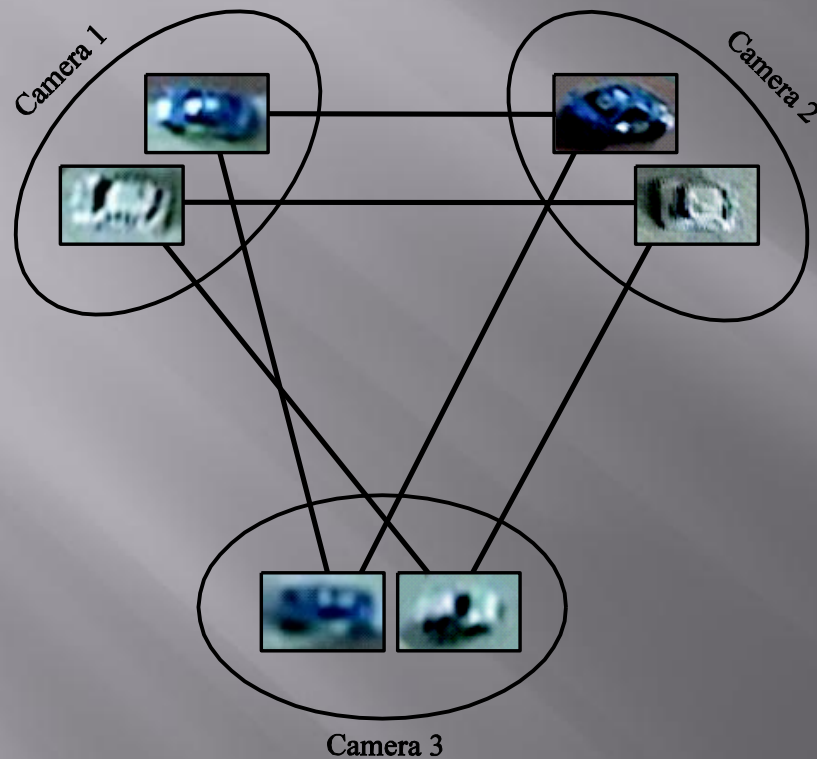
Controlled Sequence 2



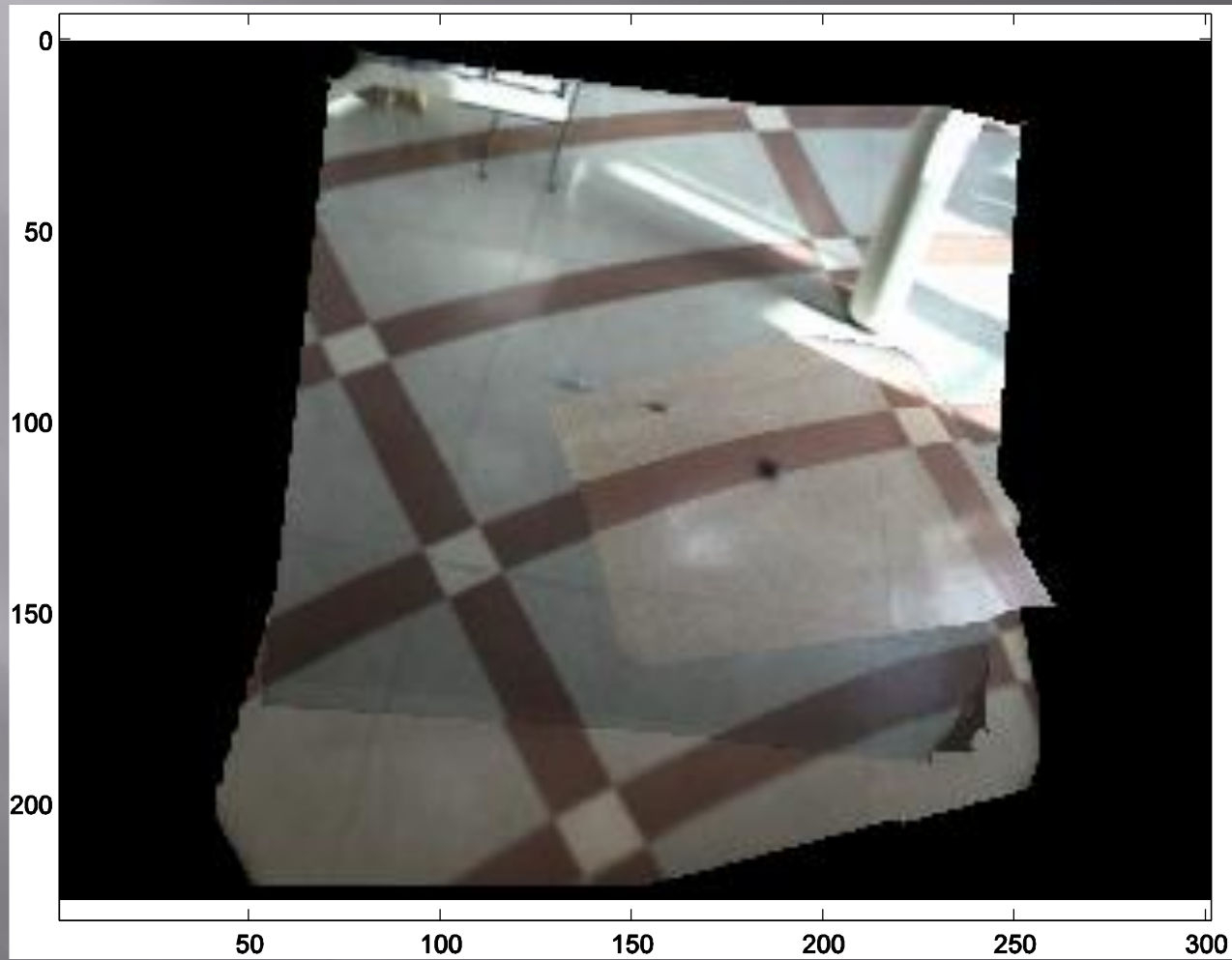
Controlled Sequence 2



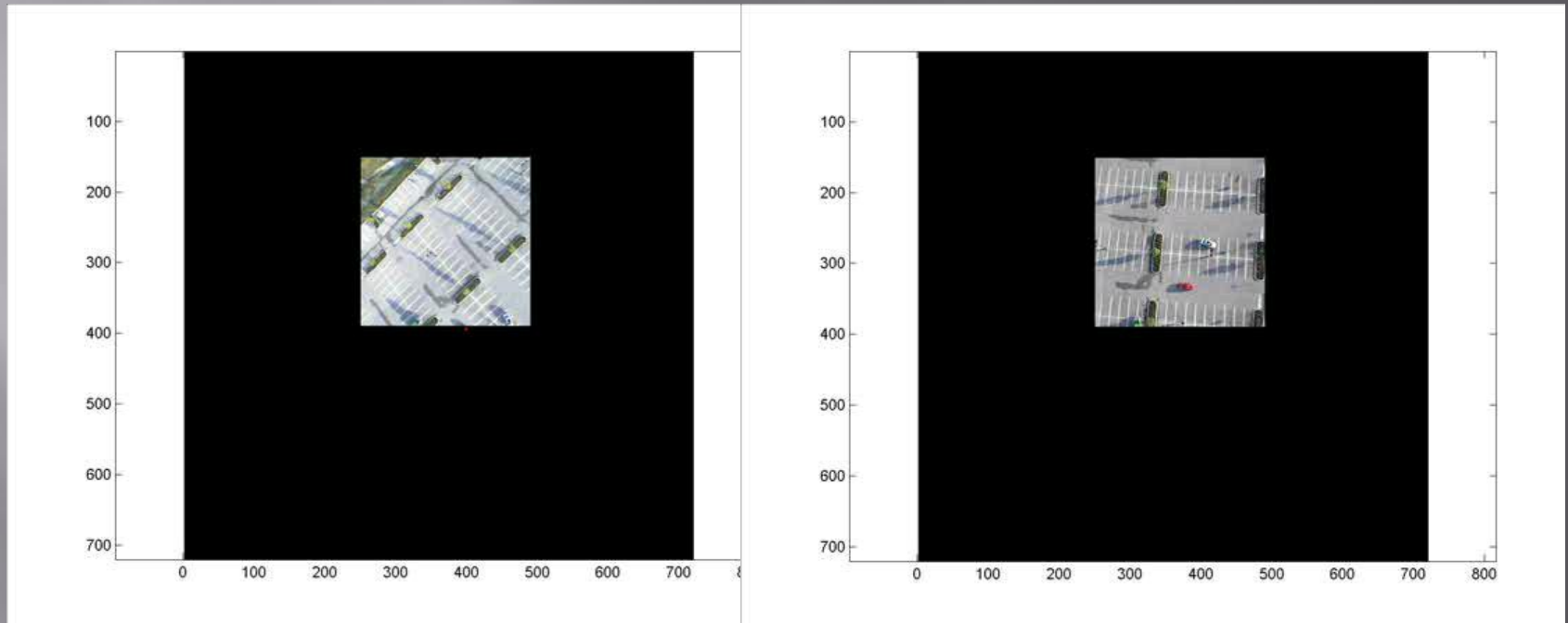
Controlled Sequence 2



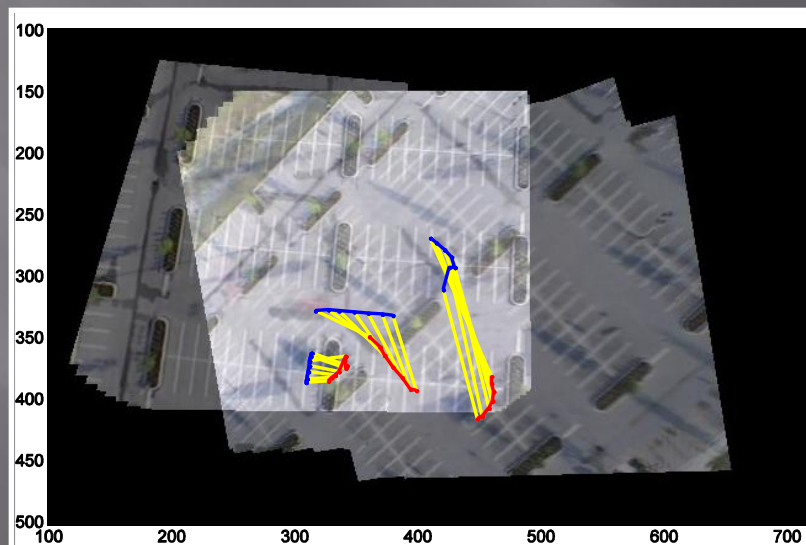
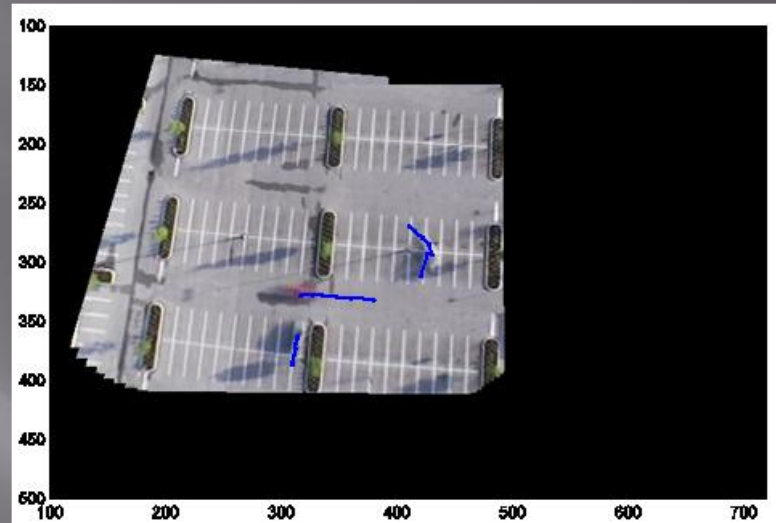
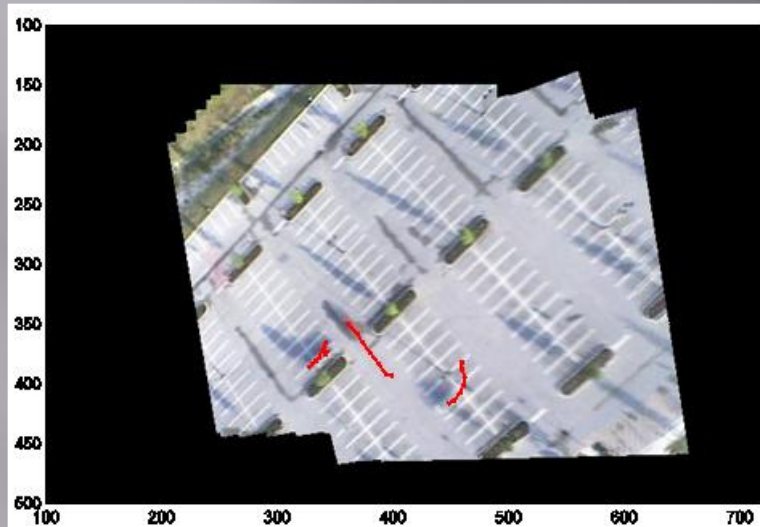
Controlled Sequence 2



UAV Sequence 1



Correspondence

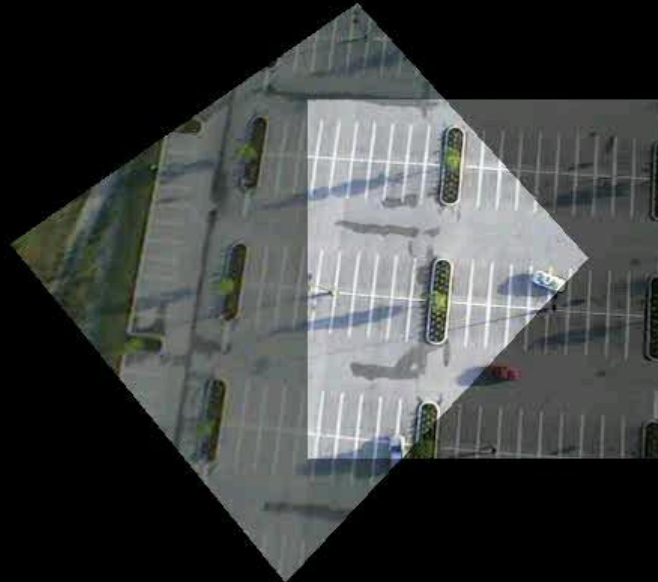


University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Results - Overlaid



Concurrent Visualization

- ▣ Mosaics are a compact representation of a single video (planar)
- ▣ Concurrent Mosaics can be used for the visualization of multiple videos
- ▣ Assuming a Lambertian scene, we can define a color transference function between the mosaics



Concurrent Visualization

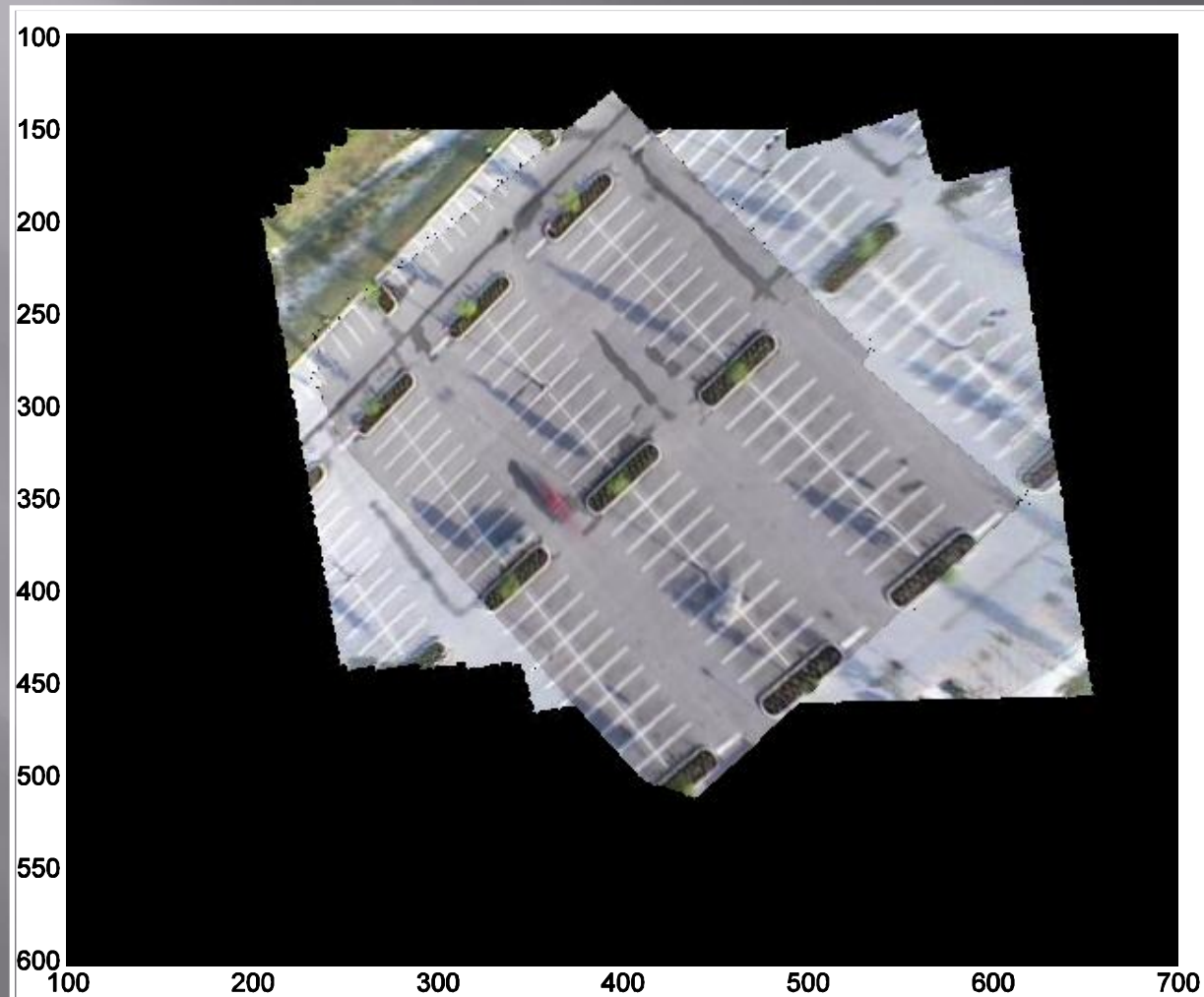
- We approximate the color transference function by a cubic trivariate polynomial

$$\begin{aligned} r(C_d) &= \sum_{ijk} a_{ijk} r_i r_j r_k \\ g(C_d) &= \sum_{ijk} a_{ijk} g_i g_j g_k \\ b(C_d) &= \sum_{ijk} a_{ijk} b_i b_j b_k \end{aligned}$$

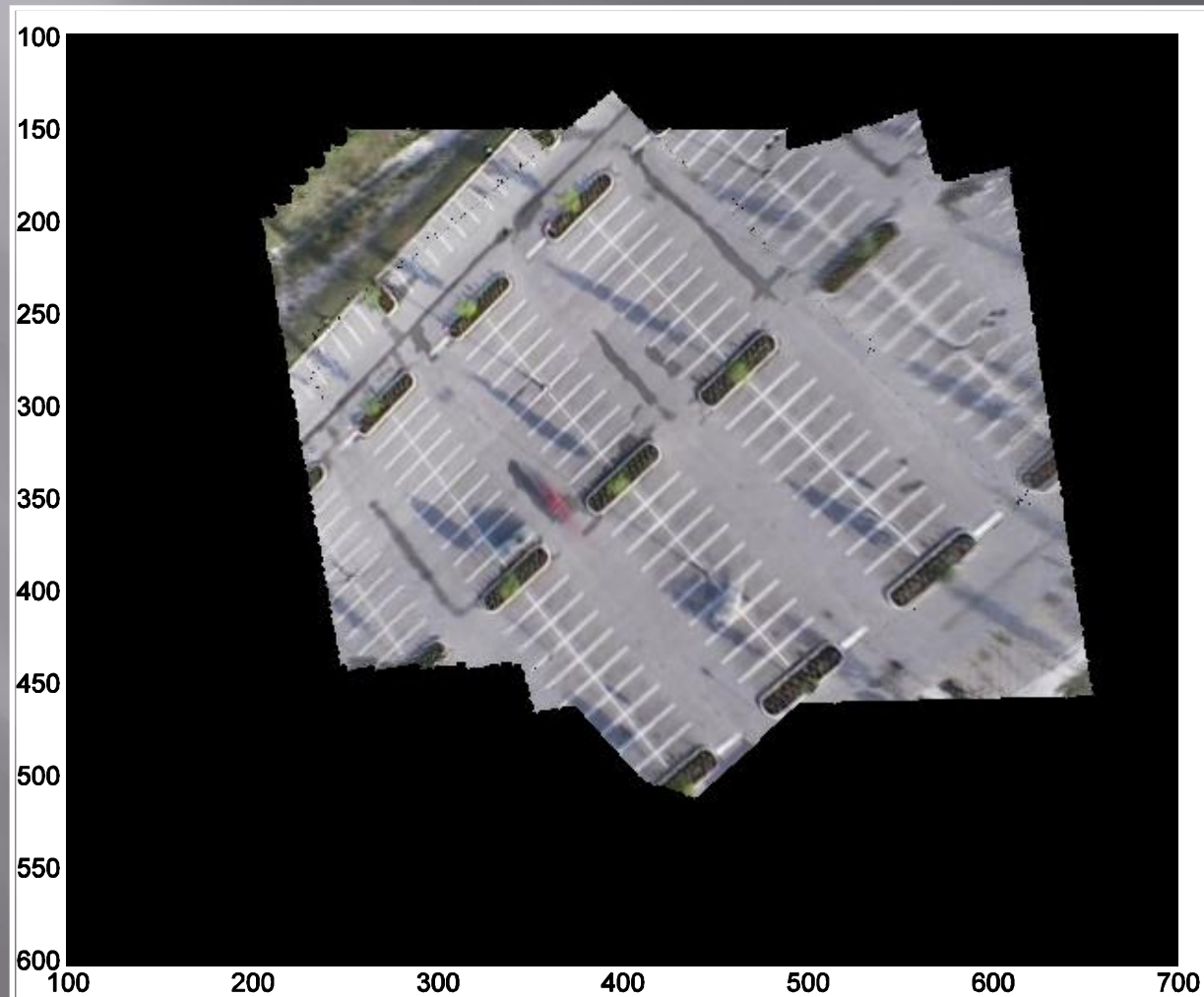
- The transfer functions values are estimated by multiple regression



Concurrent Mosaic



Blended Concurrent Mosaic

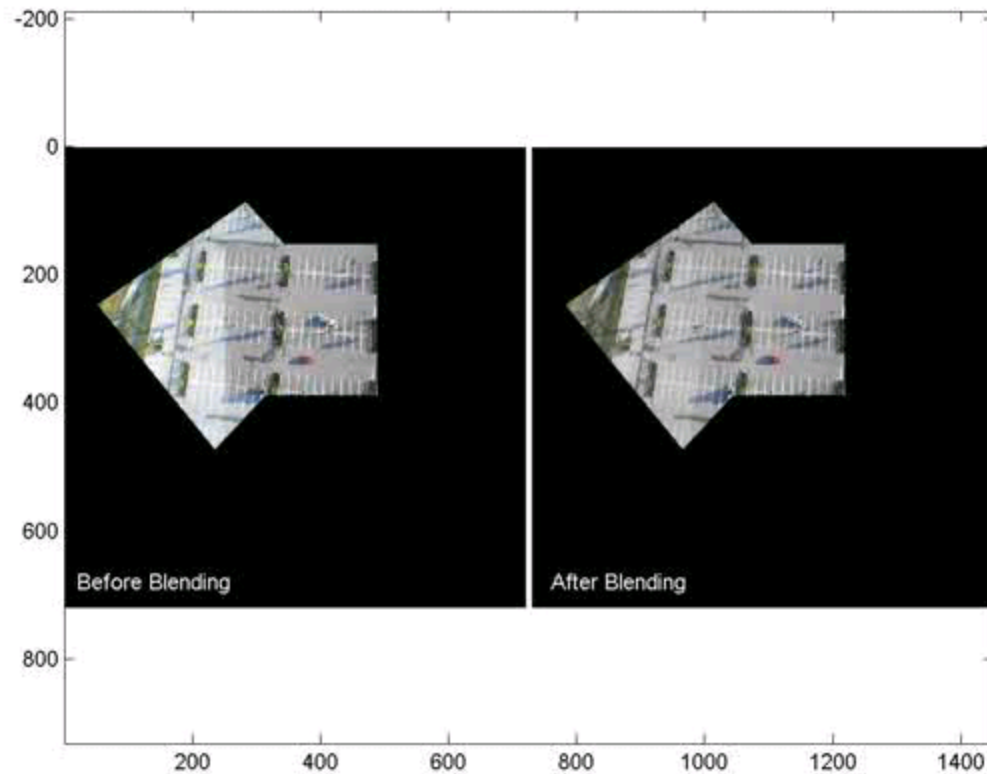


University of
Central Florida

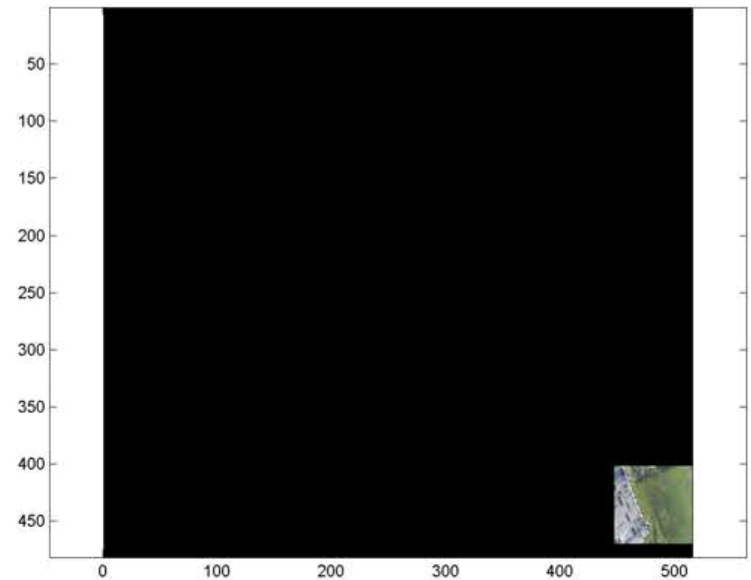
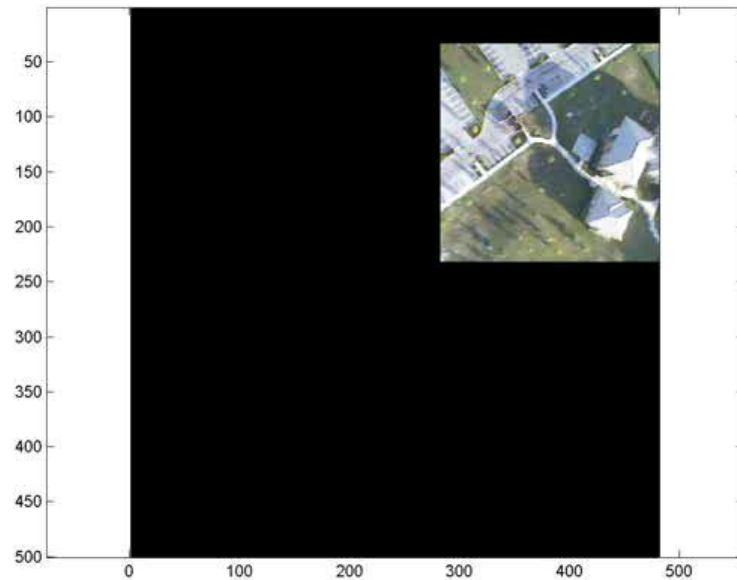
VISION

Copyright Mubarak Shah, UCF

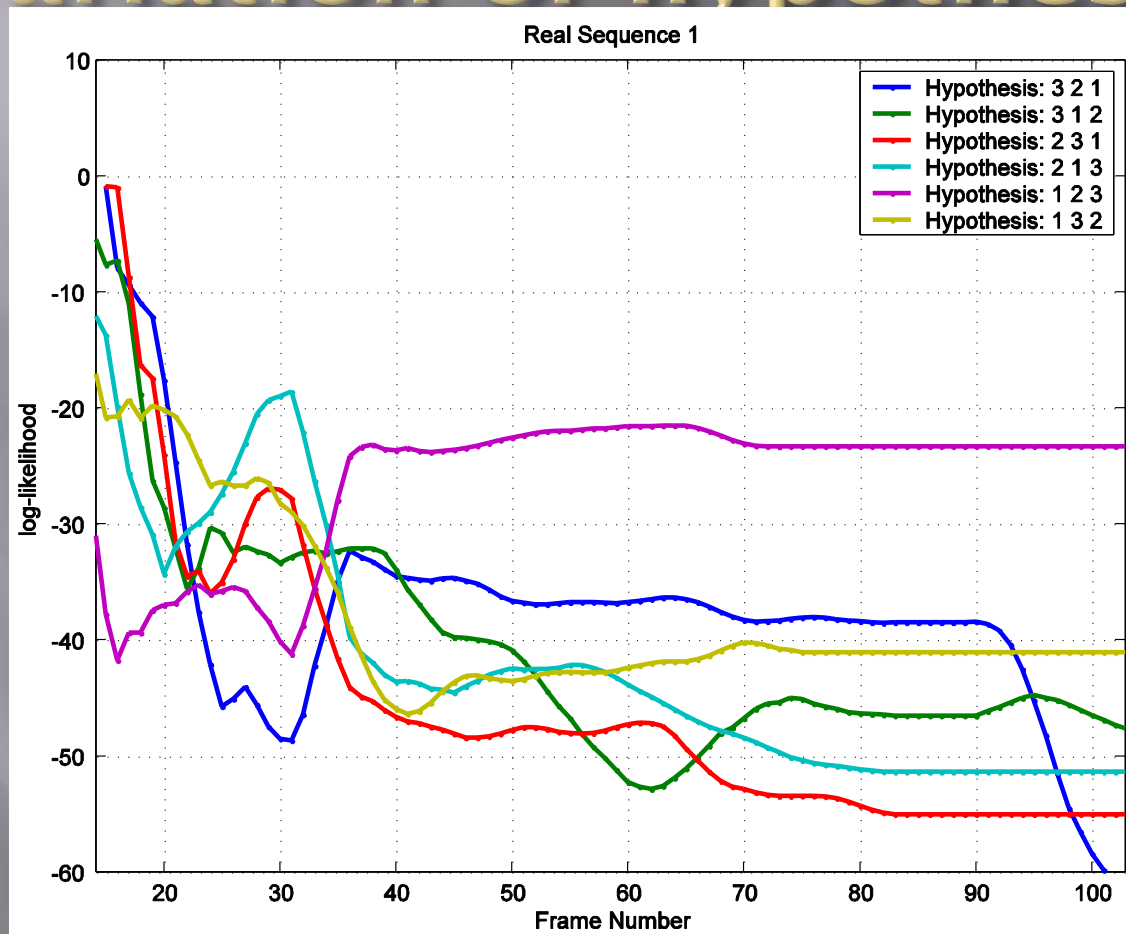
Blending



UAV Sequence 2



Variation of hypotheses



Collinear Motion

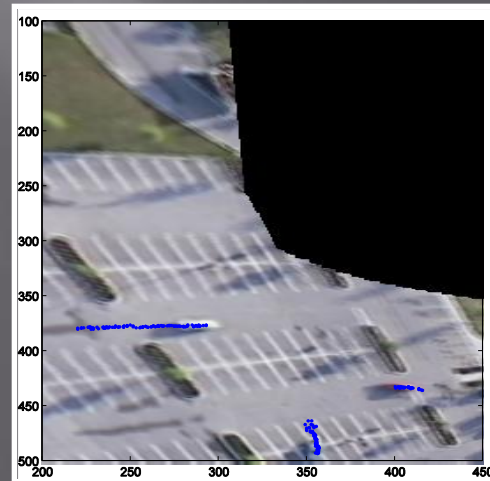
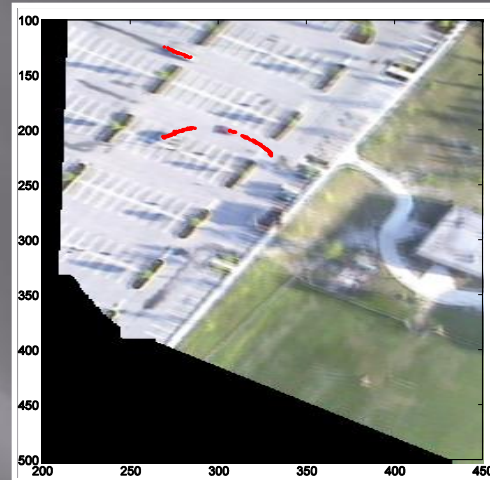
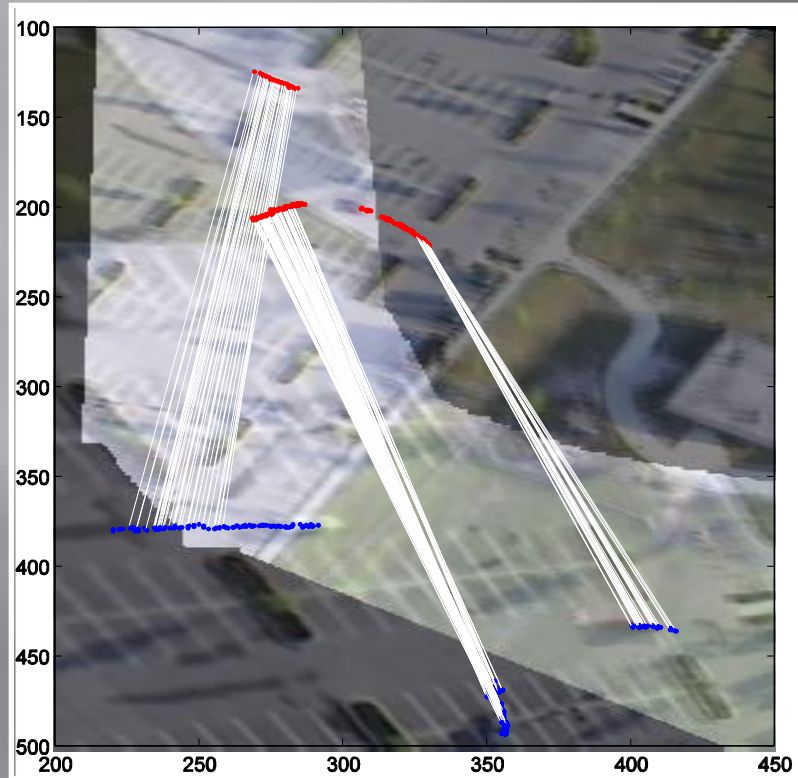


University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Correspondence



Results - Overlaid



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Results - Blended

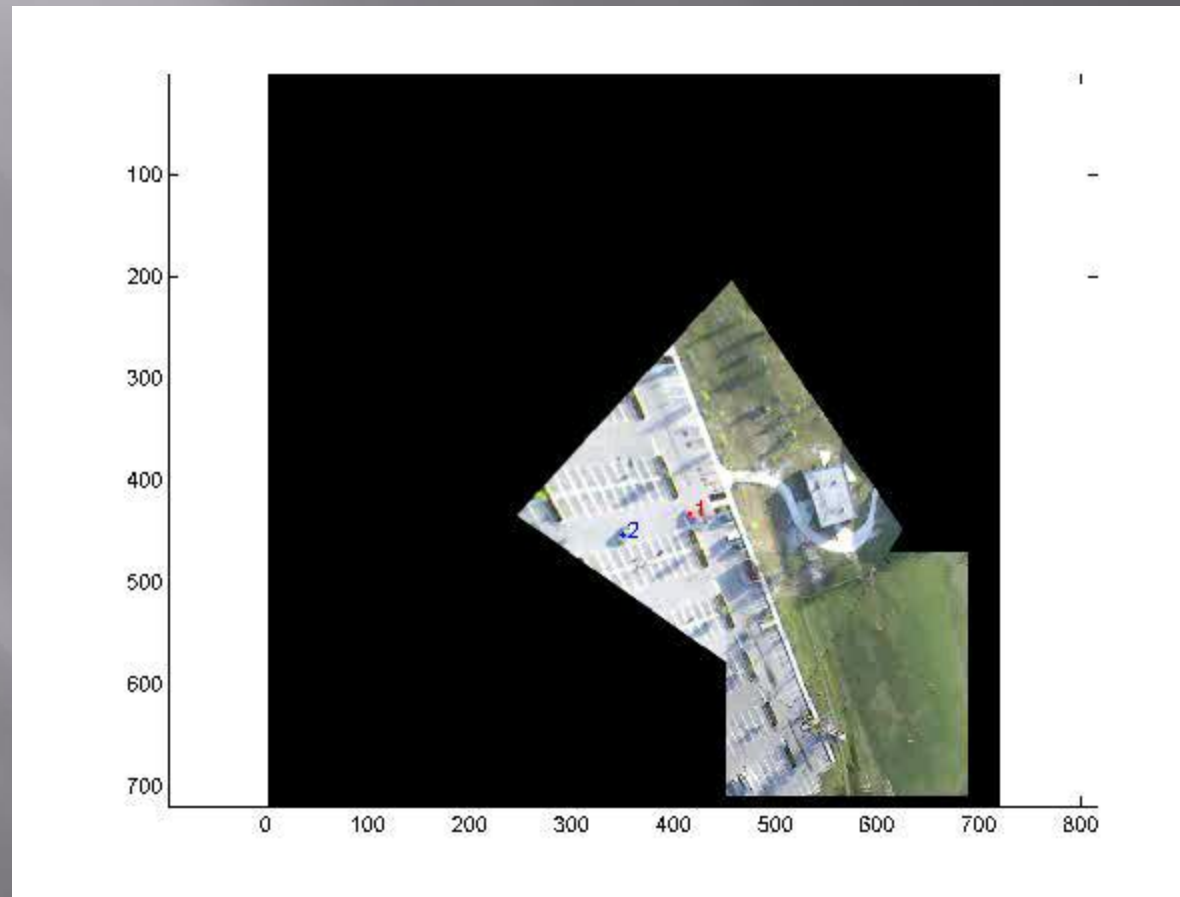


University of
Central Florida

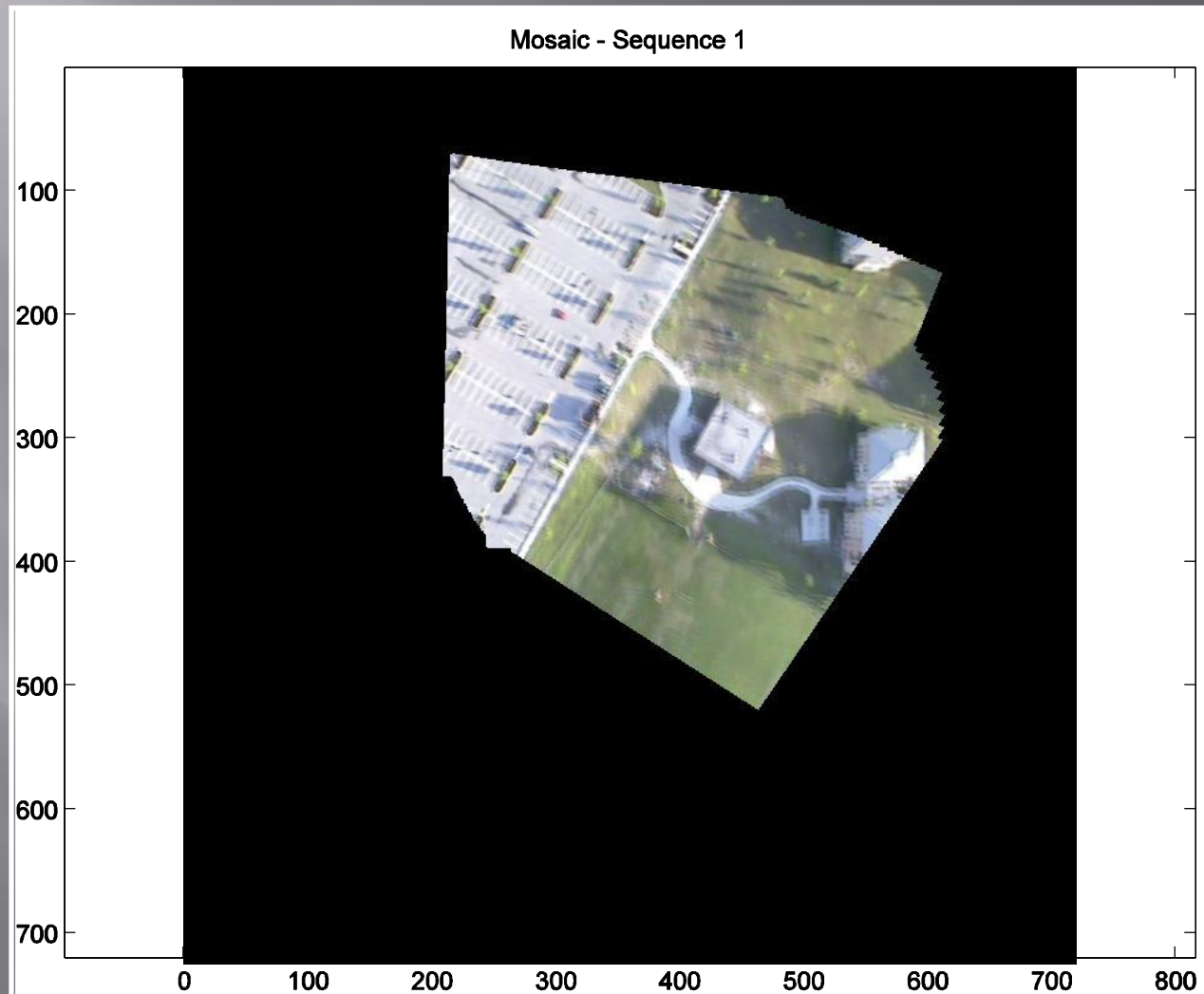
VISION

Copyright Mubarak Shah, UCF

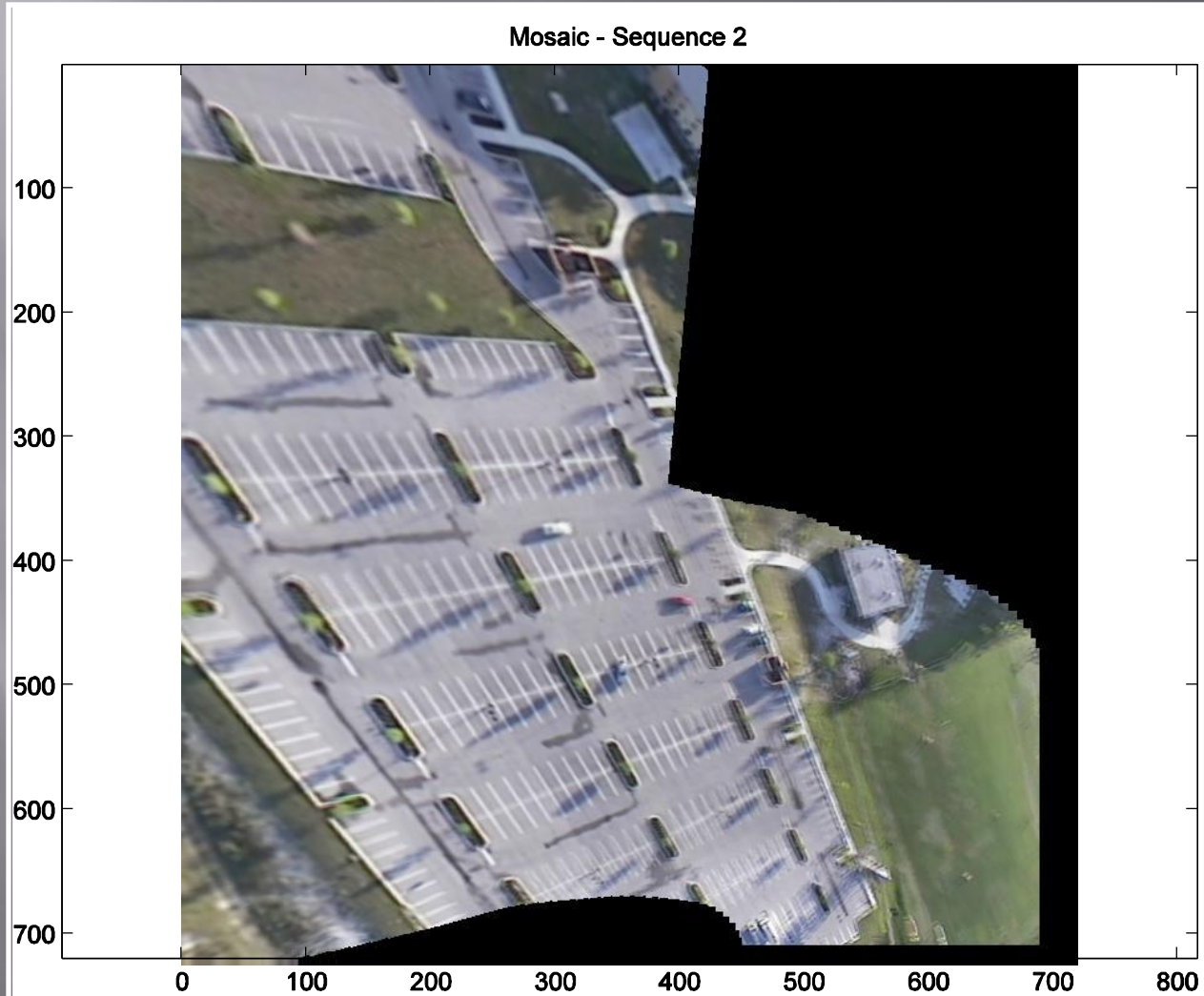
Corresponded Tracks



Mosaic 1



Mosaic 2

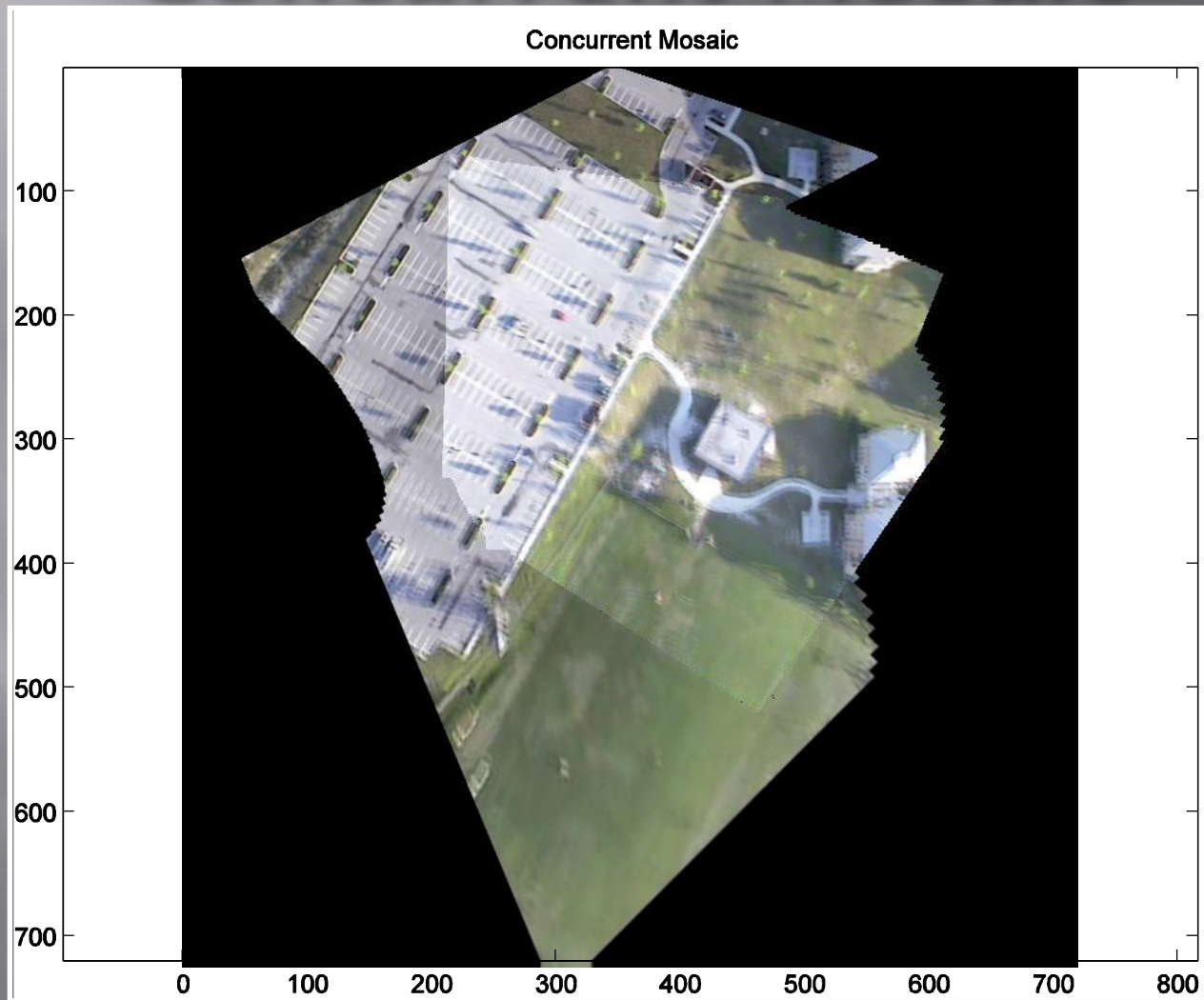


University of
Central Florida

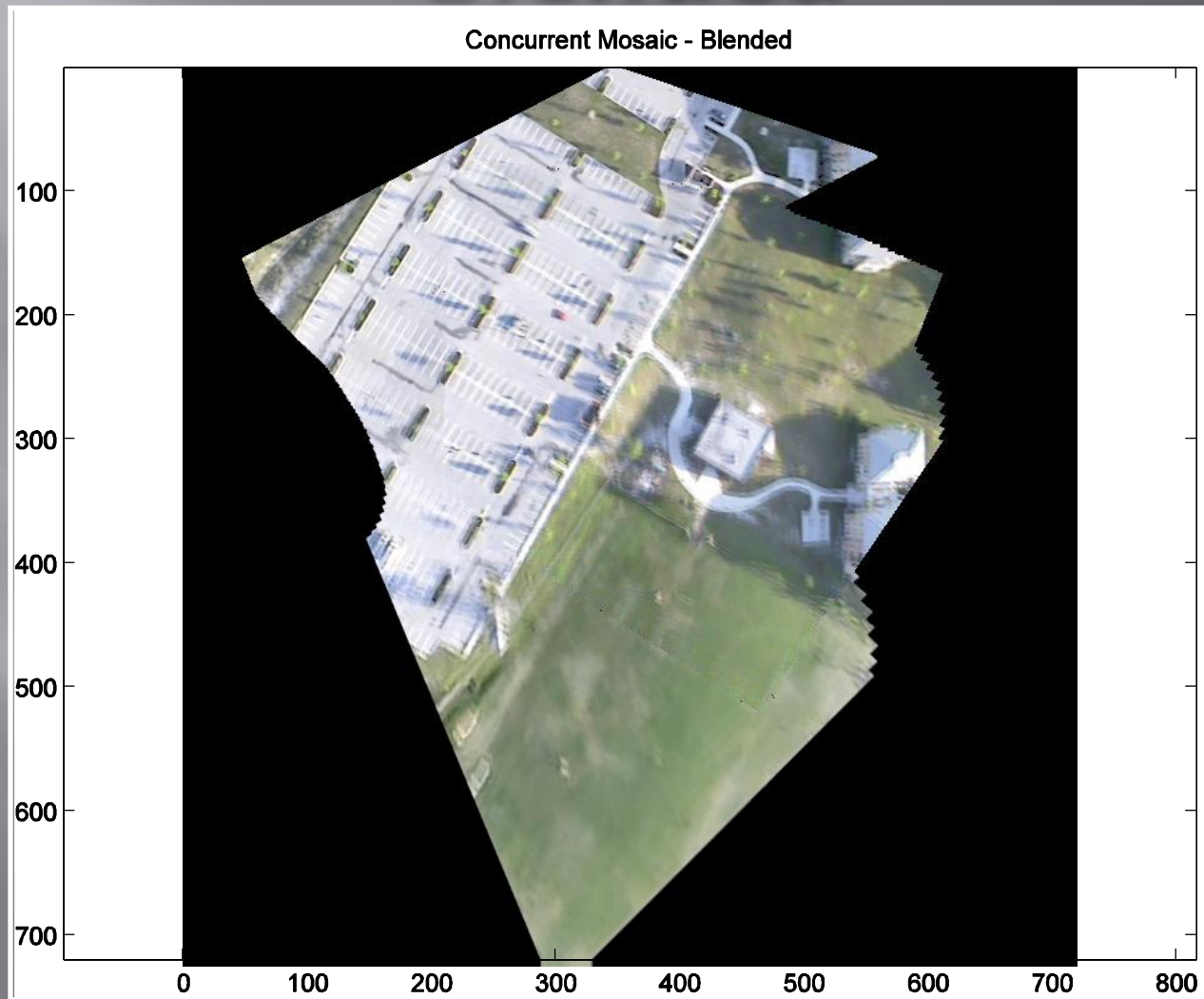
VISION

Copyright Mubarak Shah, UCF

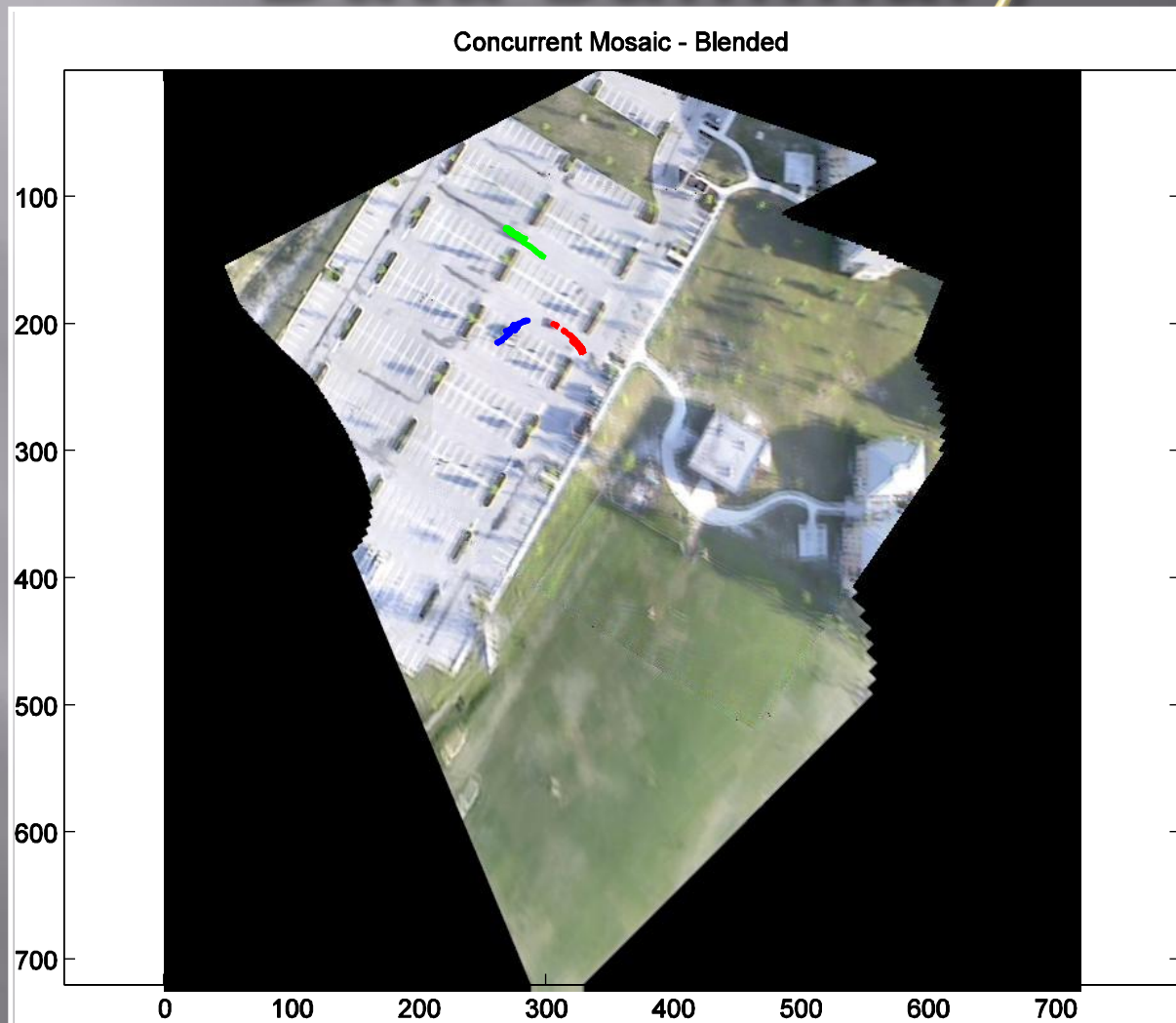
Concurrent Mosaic



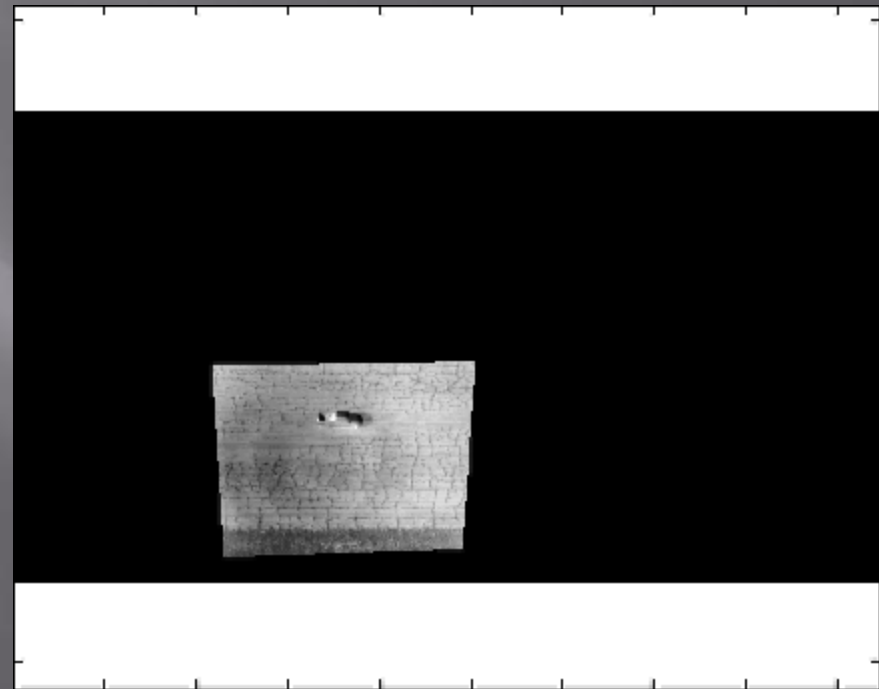
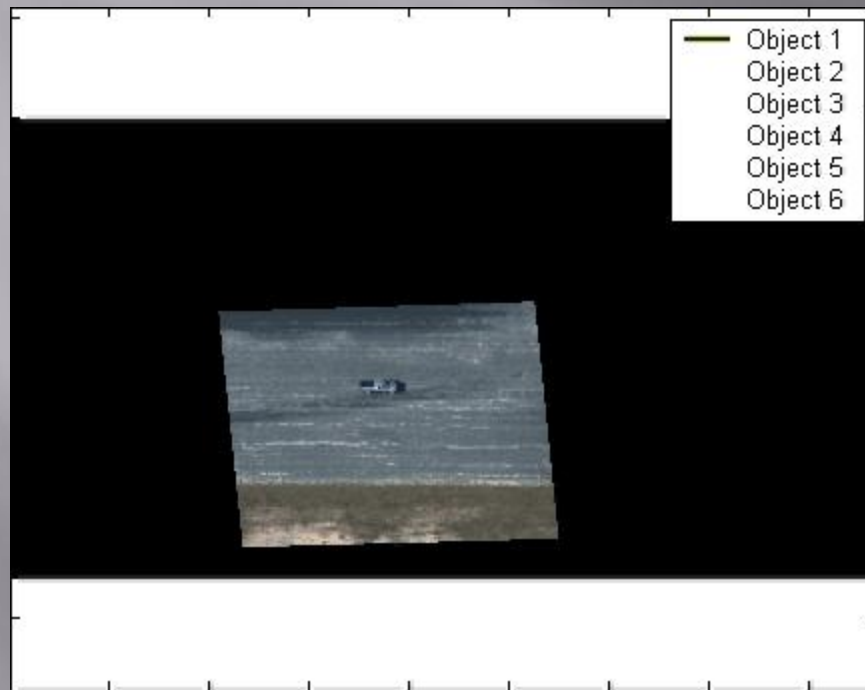
Blended



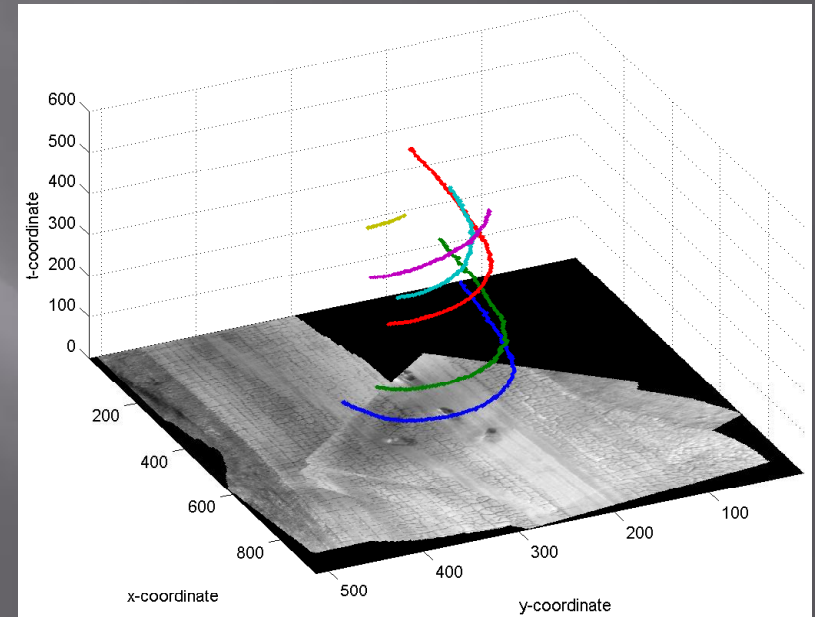
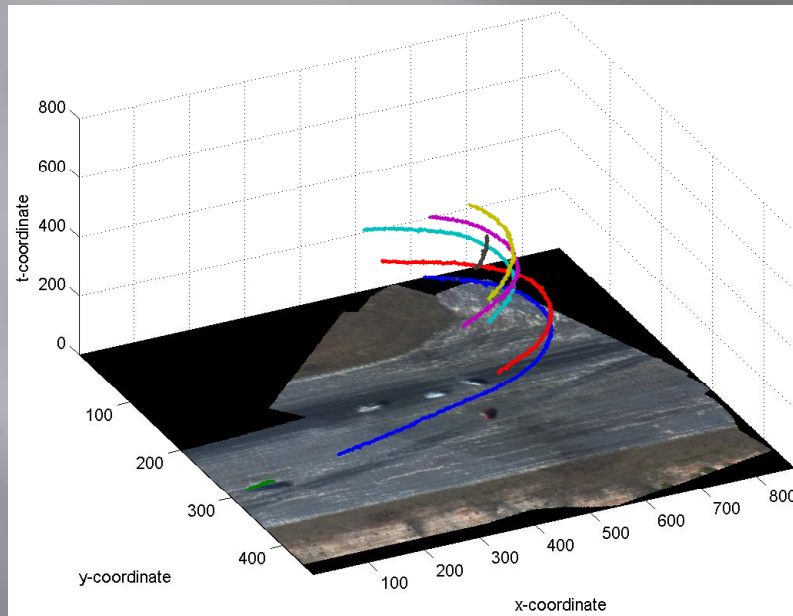
Data Summary



UAV Sequence 3



UAV Sequence 3



Trajectory Association Across Non-overlapping Cameras in Planar Scene

Yaser Sheikh, Xin Li and Mubarak Shah

CVPR 2007



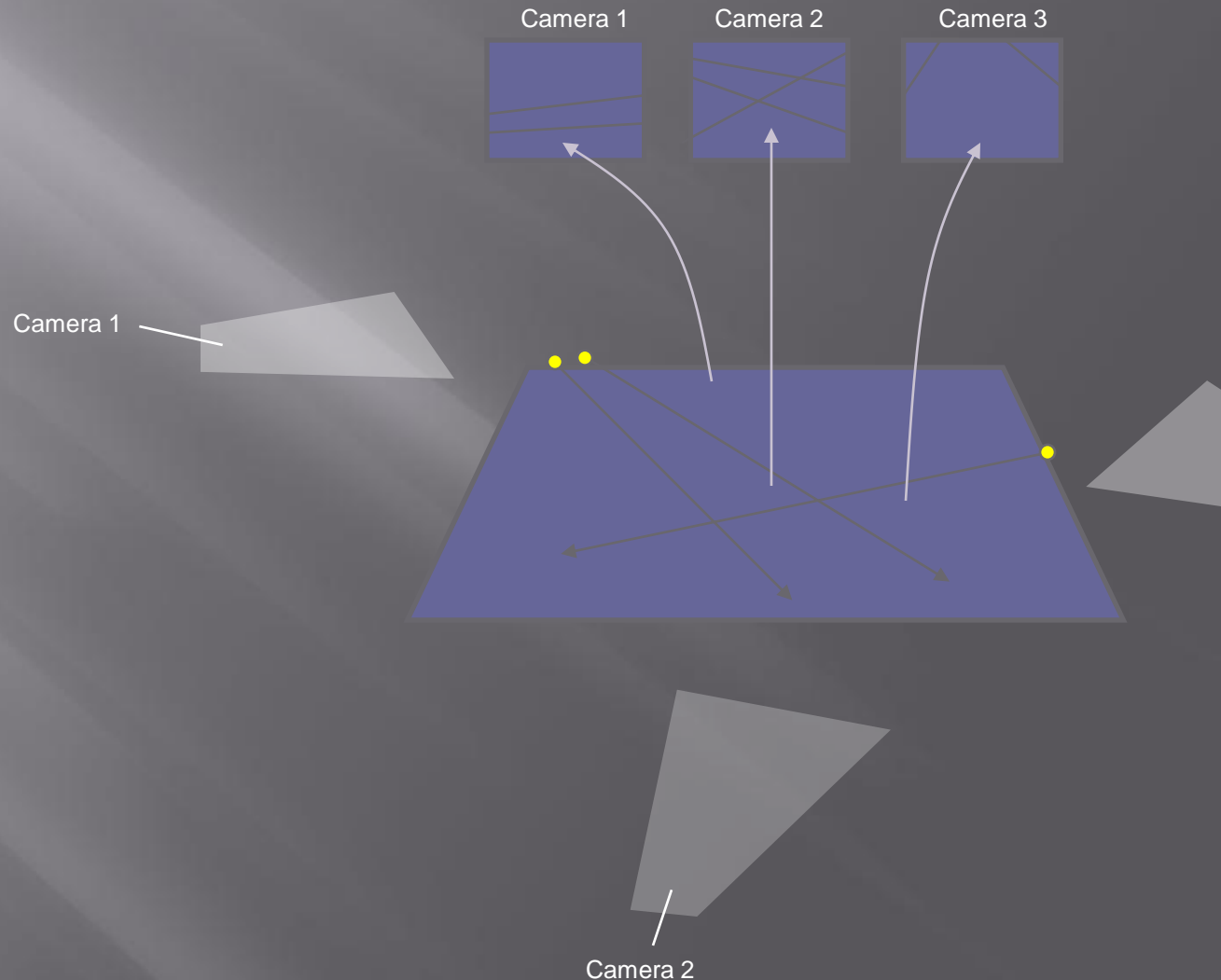
University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Data Model

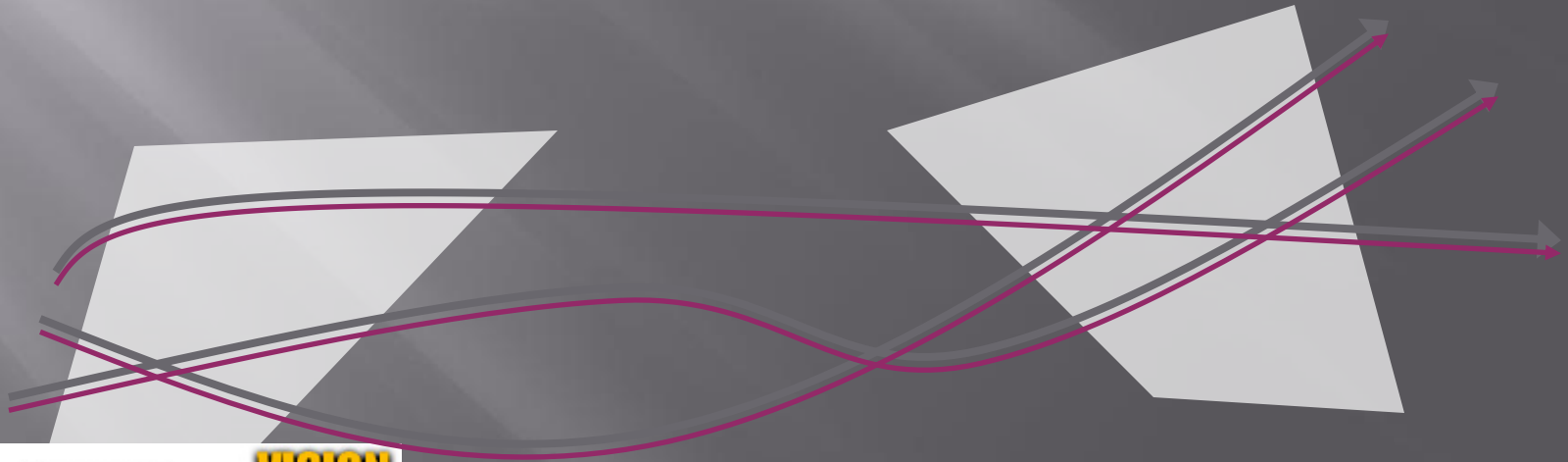
- ▣ K : # of objects
- ▣ N : # of cameras
- ▣ H : World Plane to image plane homography
- ▣ c_j^i : Association of object j in camera i
- ▣ **Goals:**
 - ▣ Association
 - ▣ Homographies
 - ▣ True Trajectories



Global Association

General Case: Zero spatiotemporal overlap

- ▣ The kinematics of the object are modeled
- ▣ Allows us to constrain spatial relation of non-overlapping FoVs



Kinematic Polynomial Models

- Model: The trajectories of each object define a polynomial in t .
- In general,



- Linear

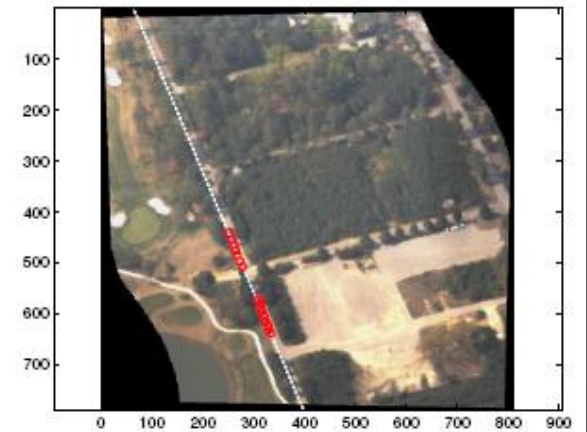
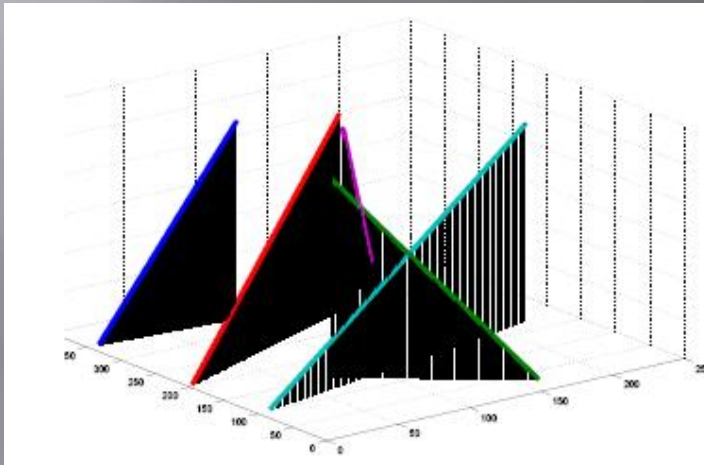


- Quadratic

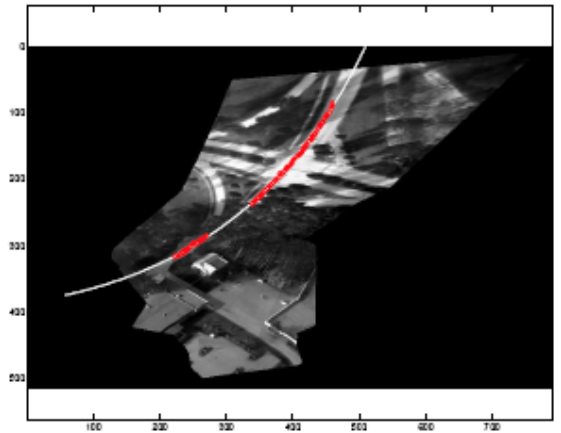
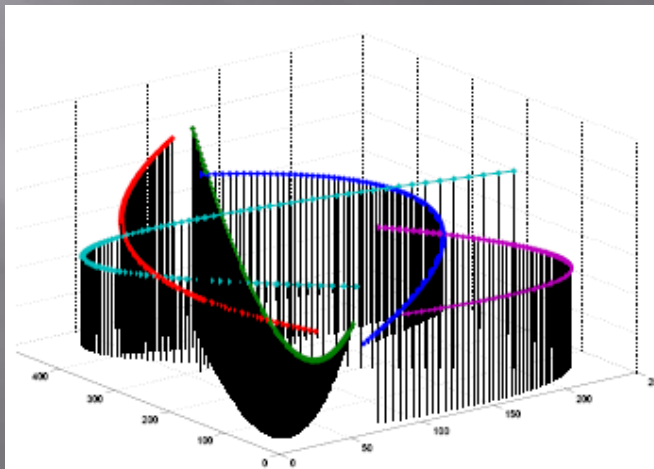


Kinematic Polynomial Models

Linear

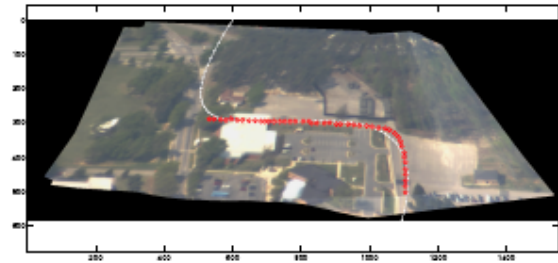
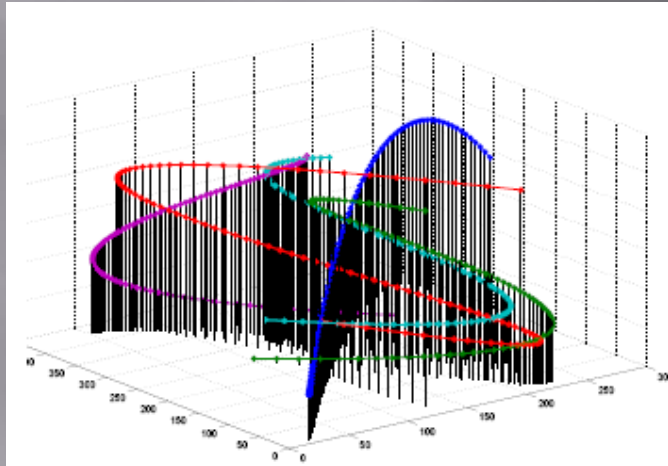


Quadratic



Kinematic Polynomial Models

Cubic



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Number of Unknowns

- ▣ Linear (constant velocity)
 - $4*k+9*N$
- ▣ Quadratic (constant acceleration)
 - $6*k+9*N$
- ▣ Cubic
 - $8*k+9*N$



Homographies

- ▣ To get the imaged point at time t we have,

$$\begin{bmatrix} \lambda x_j^i(t) \\ \lambda y_j^i(t) \\ \lambda \end{bmatrix} = \mathbf{H}^i \begin{bmatrix} x_j(t) \\ y_j(t) \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} \lambda x_j^i(t) \\ \lambda y_j^i(t) \\ \lambda \end{bmatrix} = \mathbf{H}^i \begin{bmatrix} p_{1j}t + p_{0j} \\ q_{1j}t + q_{0j} \\ 1 \end{bmatrix} = \mathbf{H}^i \mathbf{P}_j \begin{bmatrix} t \\ 1 \end{bmatrix}$$



Maximum Likelihood Parameter Estimation

▣ Question:

Given the data, under this model, what is the optimal estimate of association, parameters and homographies?

▣ Formally:

Find the Maximum Likelihood estimates of $\Theta = (\{P_k\}_K, \{H^n\}_N)$ and $C = \{c\}_K^N$ for data $\underline{X} = \{\underline{x}\}_K^N$.

Problem 1: Define the likelihood function

Problem 2: Provide a maximization algorithm

Results

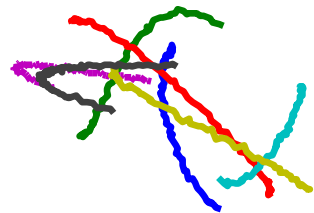
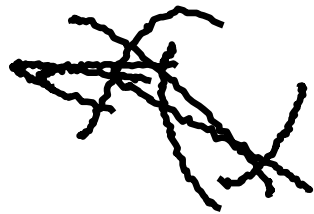
- ▣ Simulations
 - Linear, Quadratic and Cubic
- ▣ Reacquisition of objects in a single camera
- ▣ Association across cameras



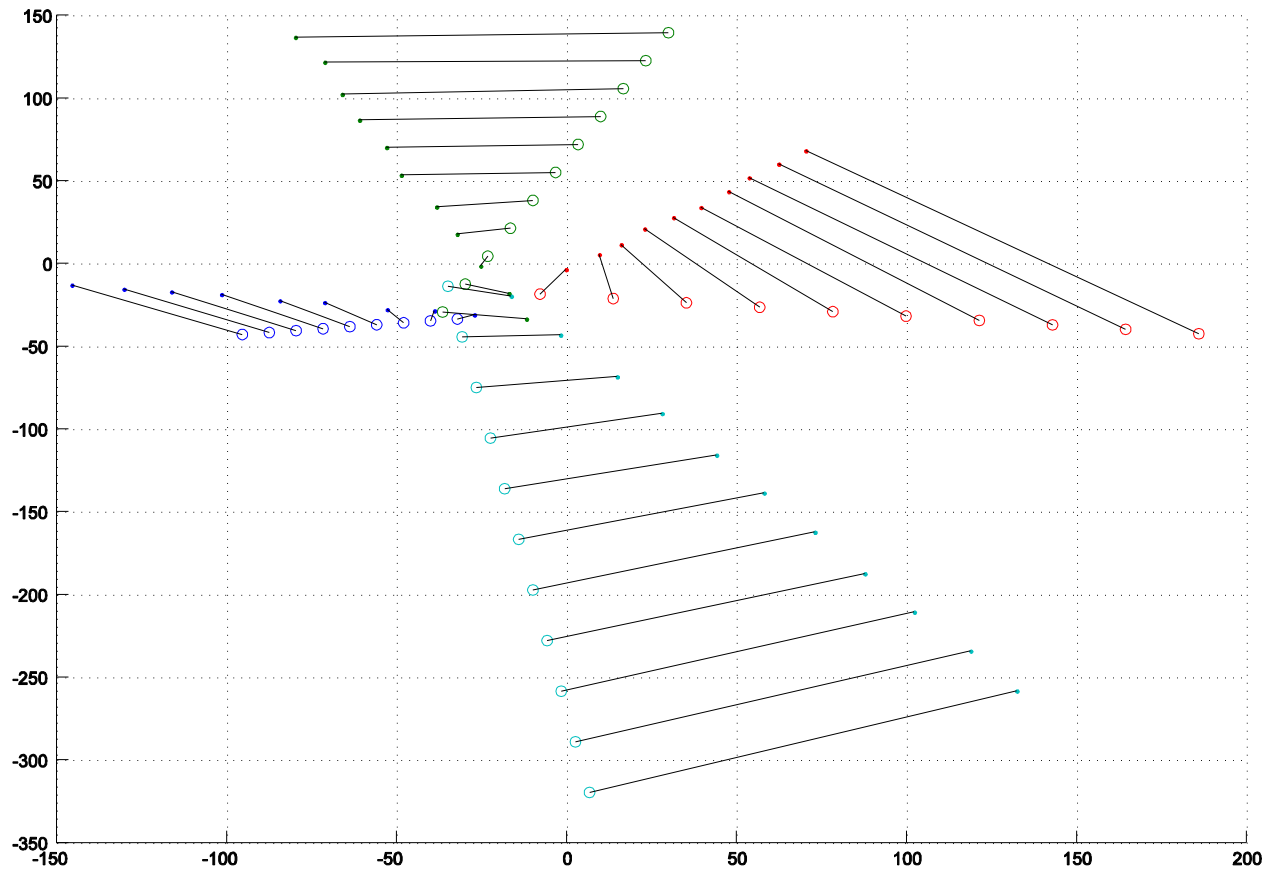
Simulations

- ▣ Measurement noise
- ▣ Number of Objects
- ▣ Number of Cameras
- ▣ Total Time duration
- ▣ Observation Start and Duration (per camera)
- ▣ Parameter noise

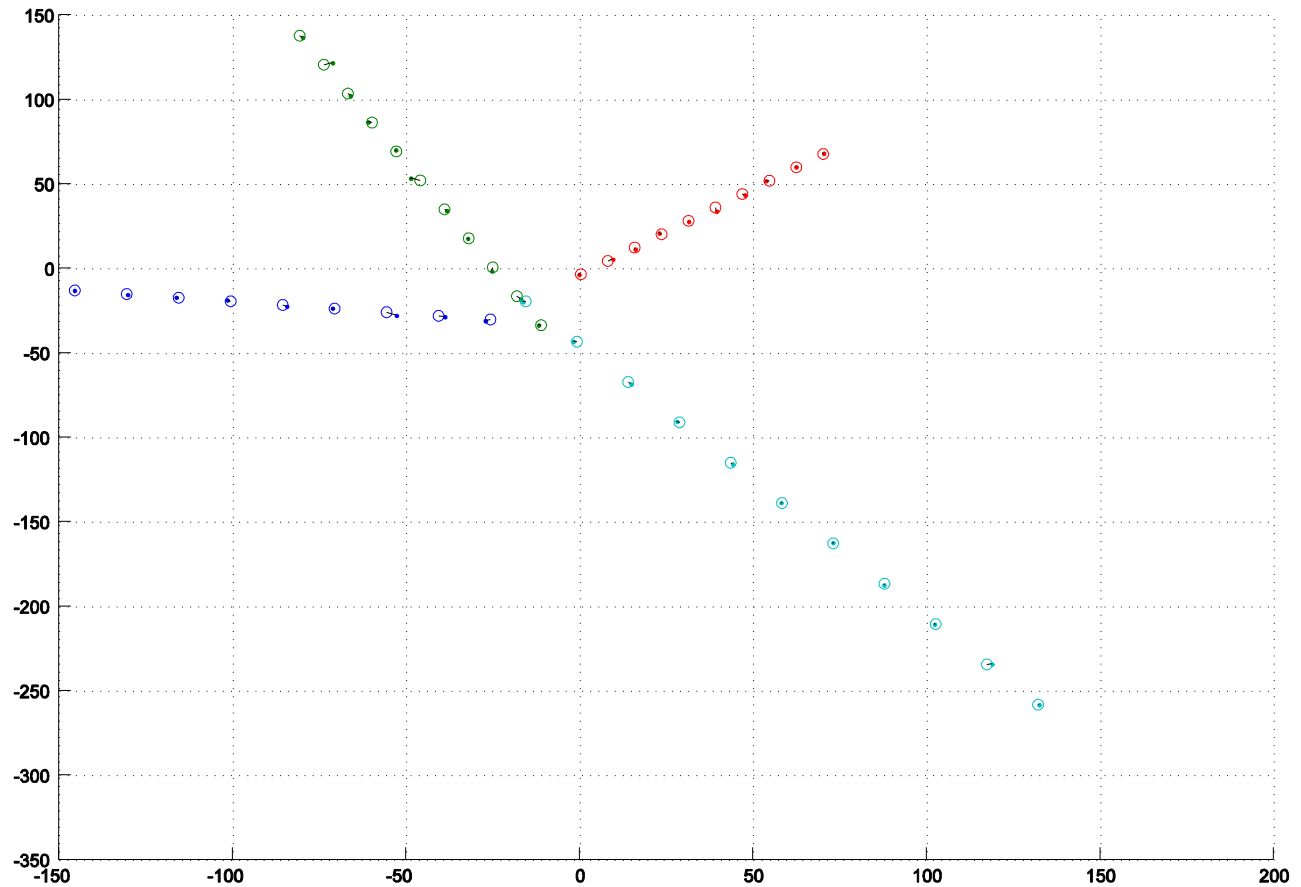
Simulations



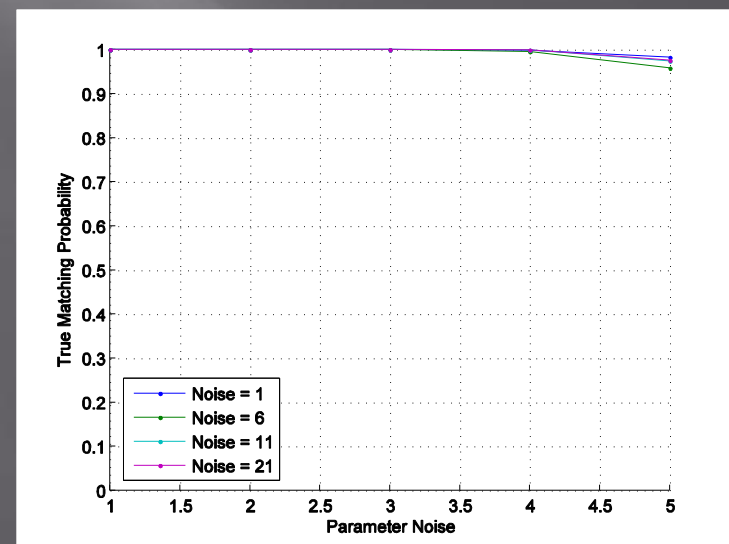
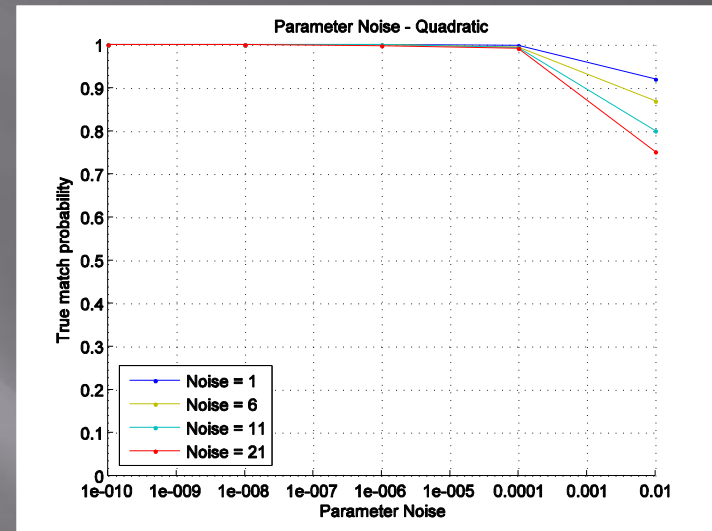
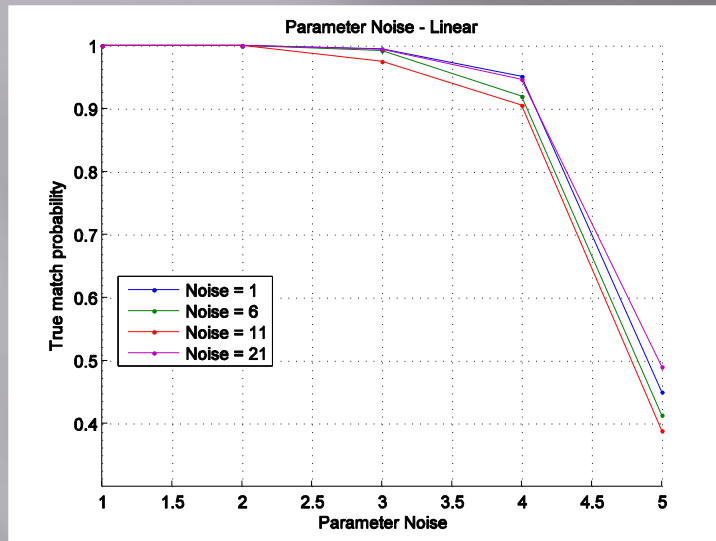
Simulations E.g.



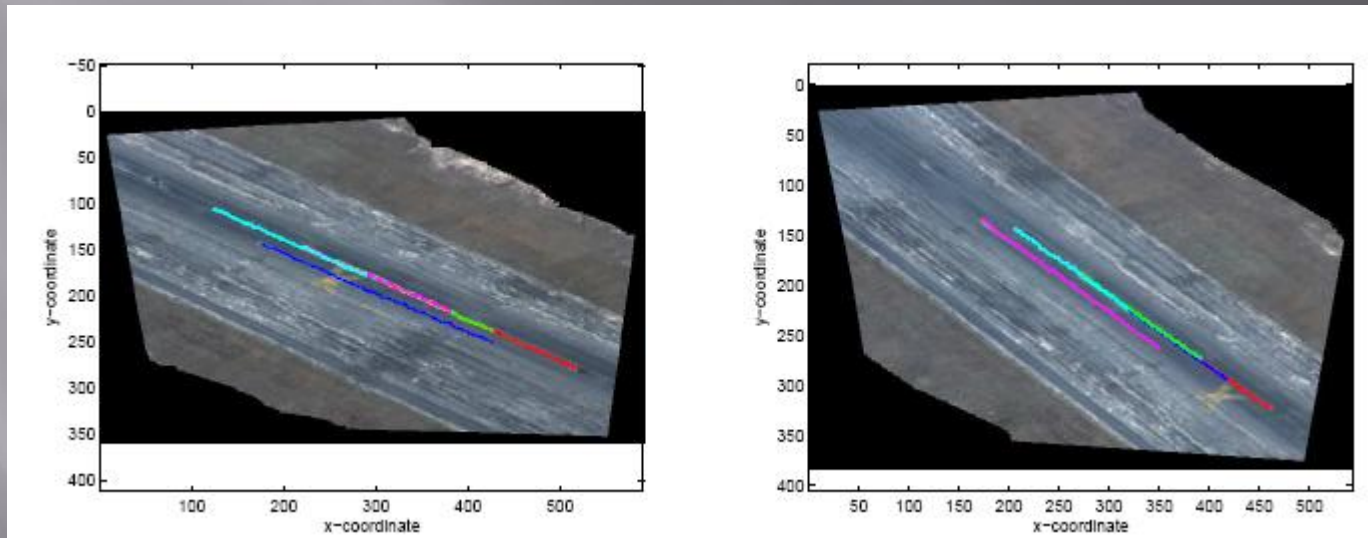
Simulations E.g.



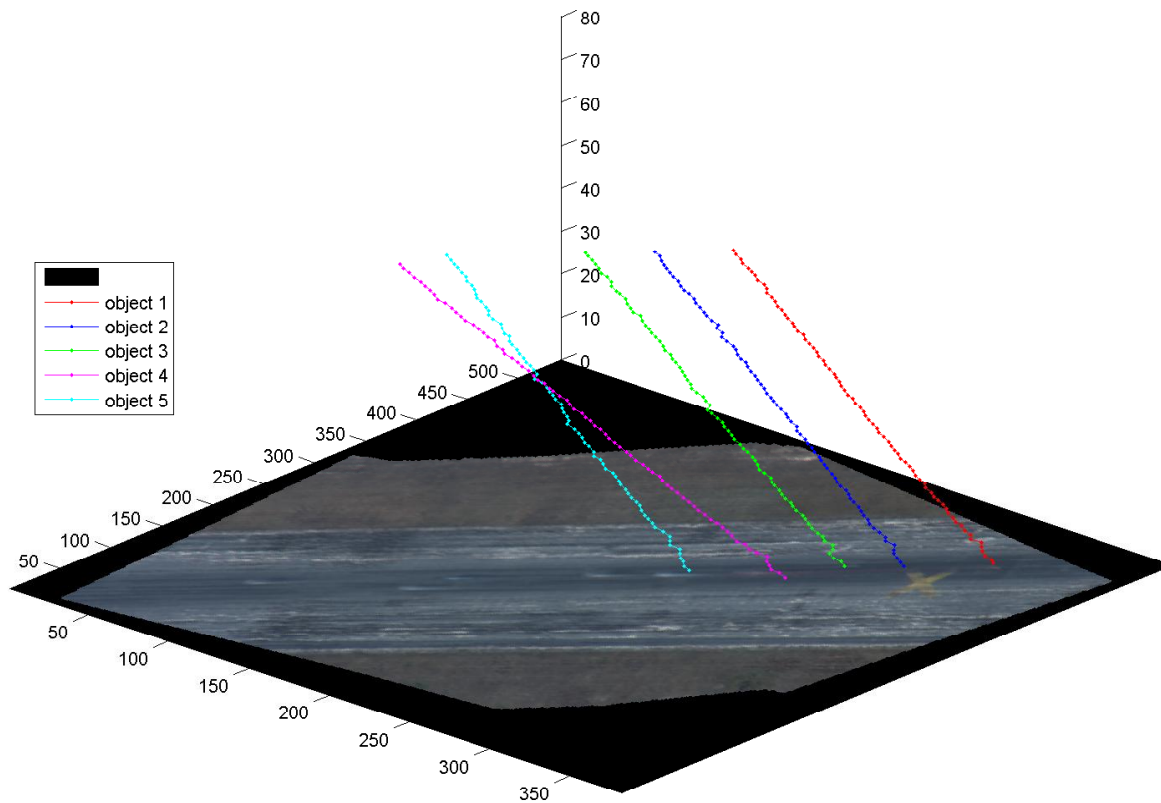
Parameter Noise



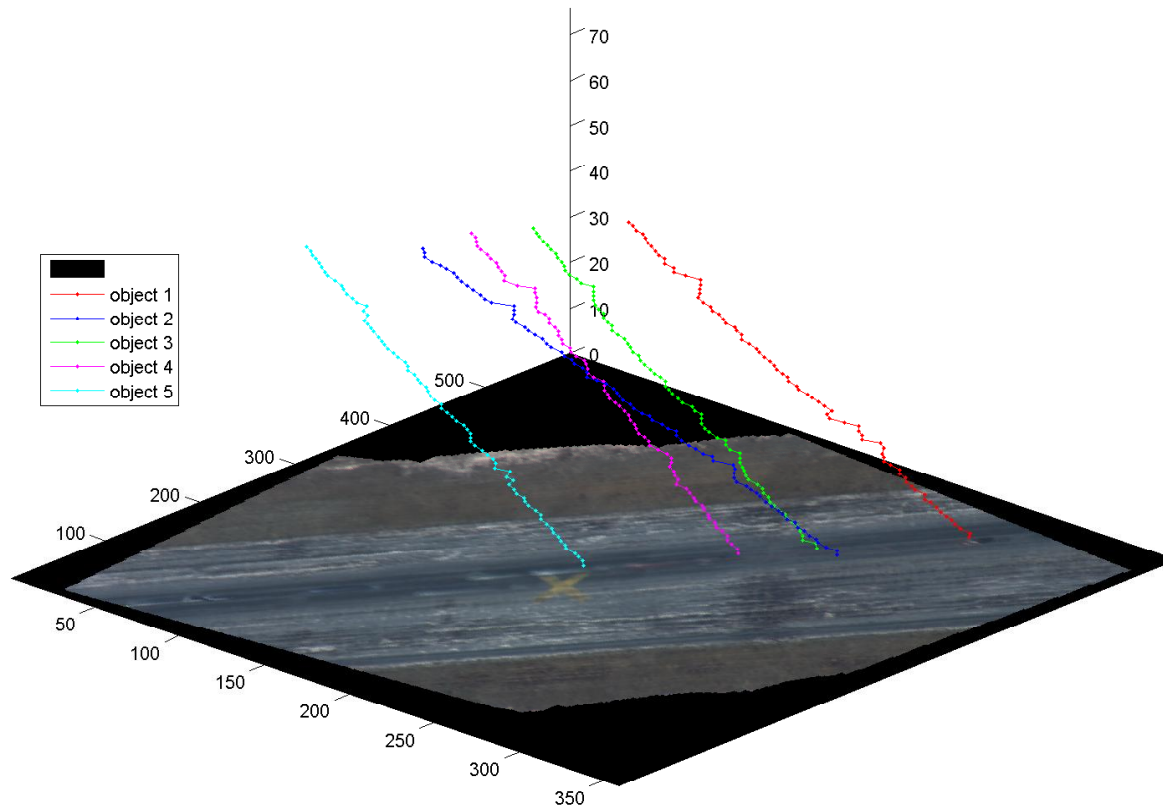
Re-association Experiment I



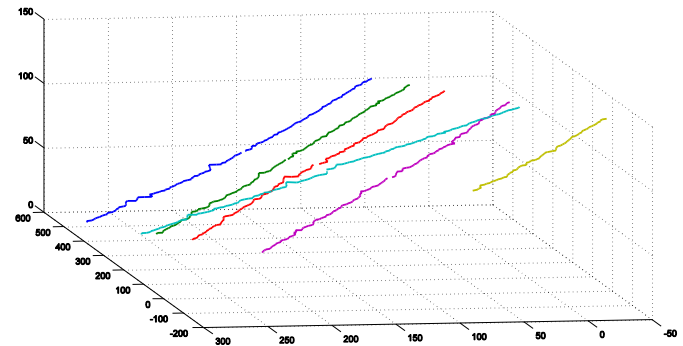
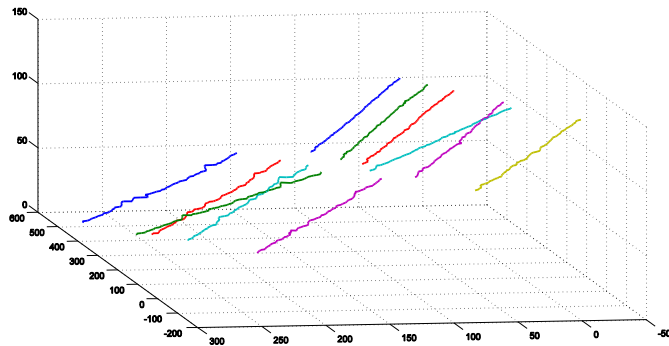
Re-association Experiment 1



Re-association Experiment I

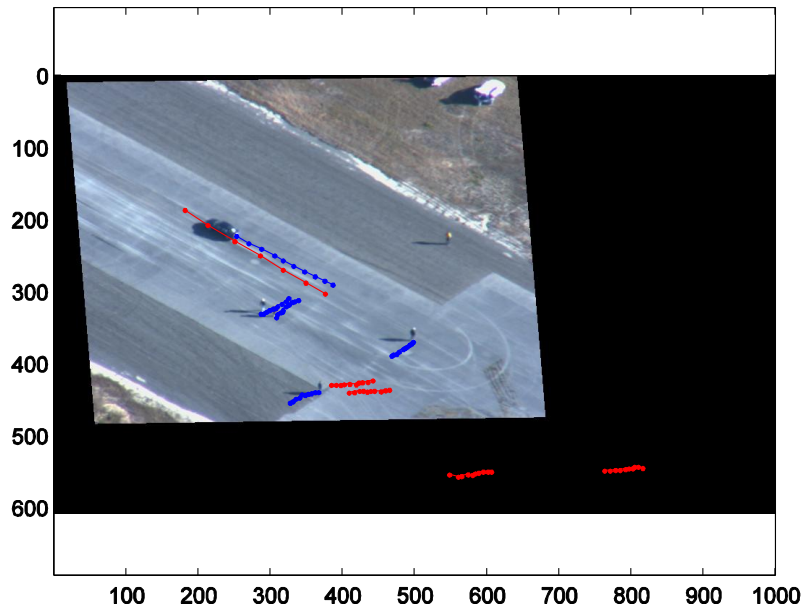


Before and After

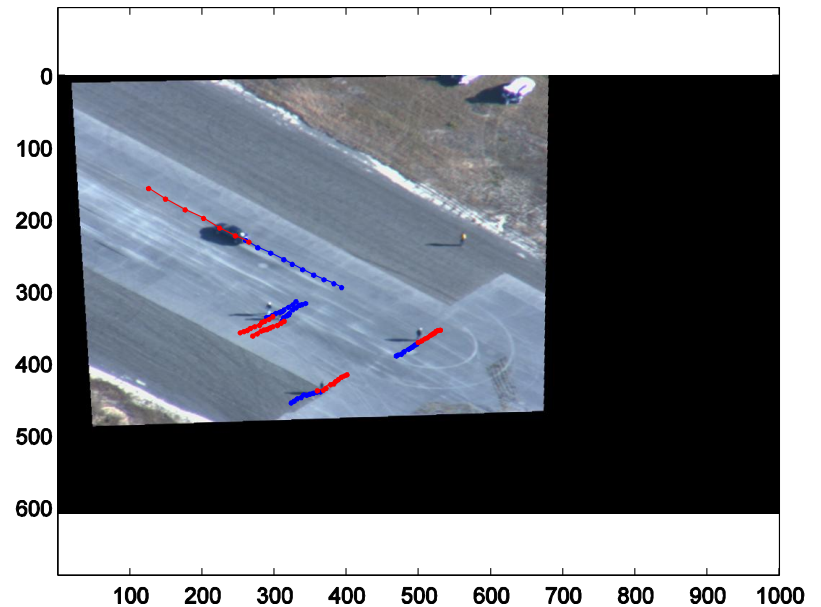


Experiment-2

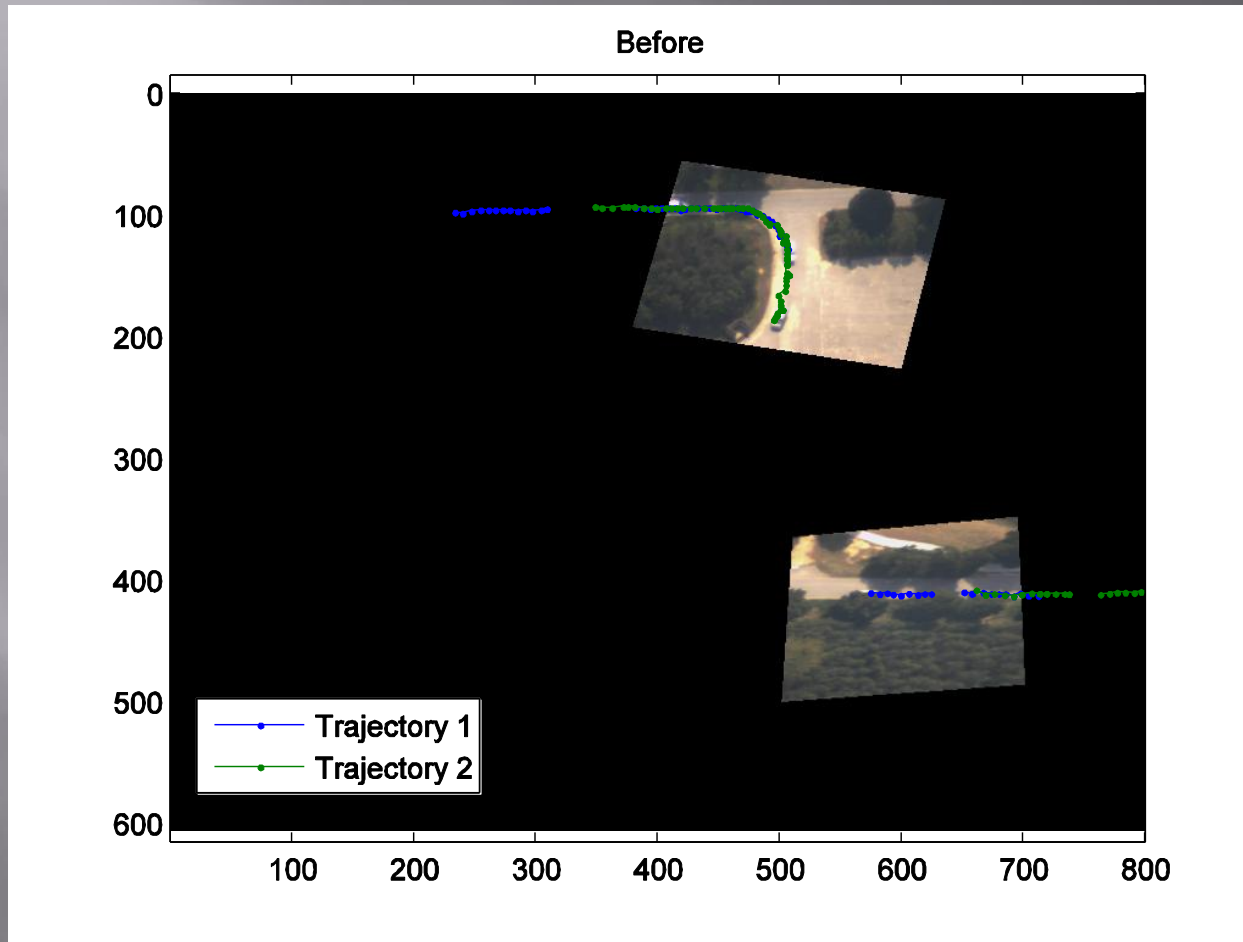
Before - Camera 1



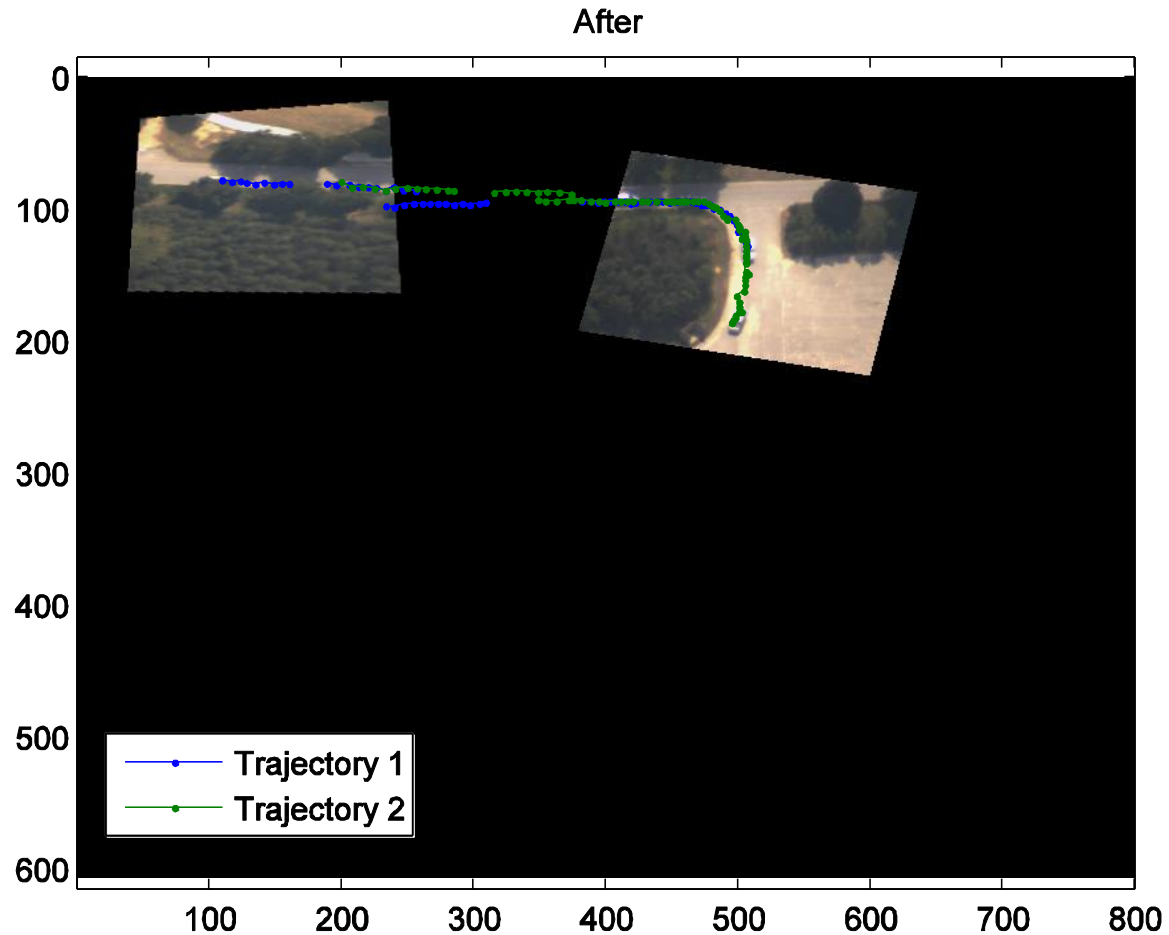
After - Camera 1



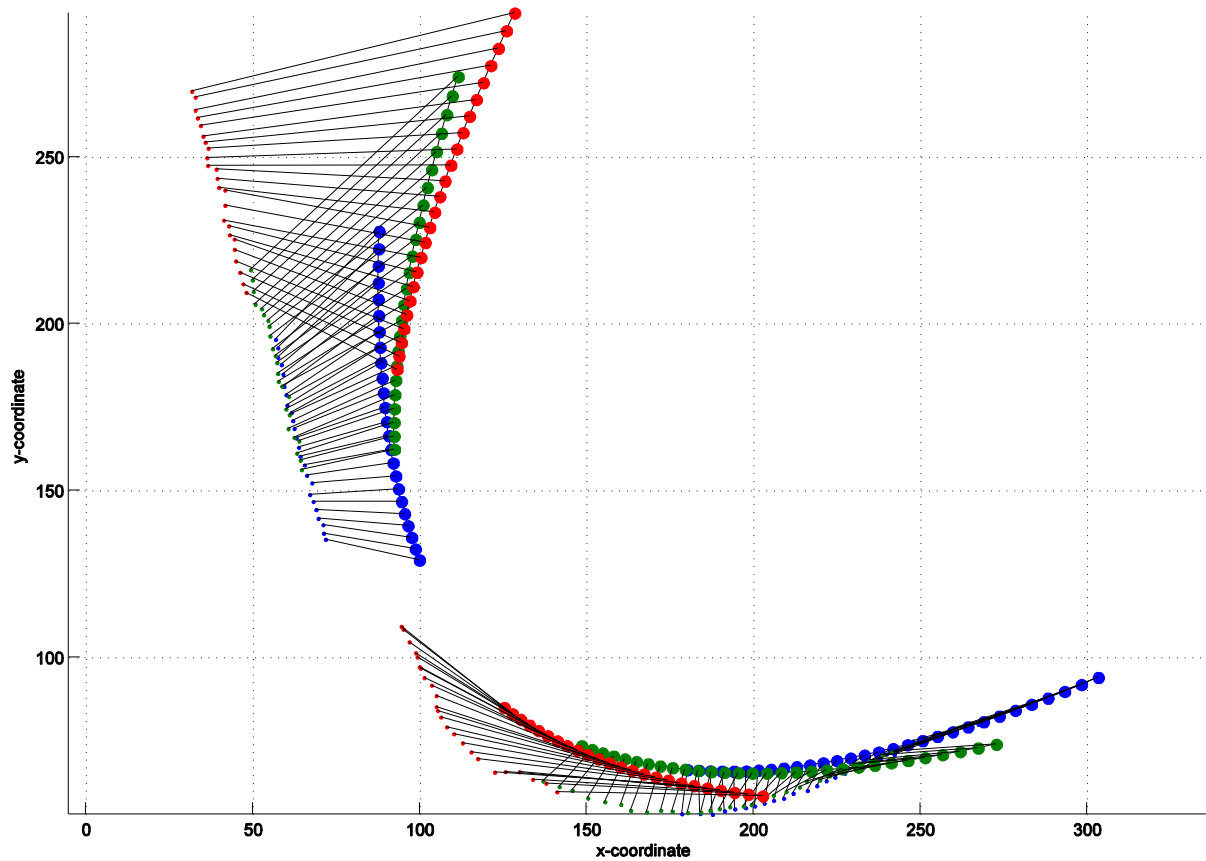
Experiment 3



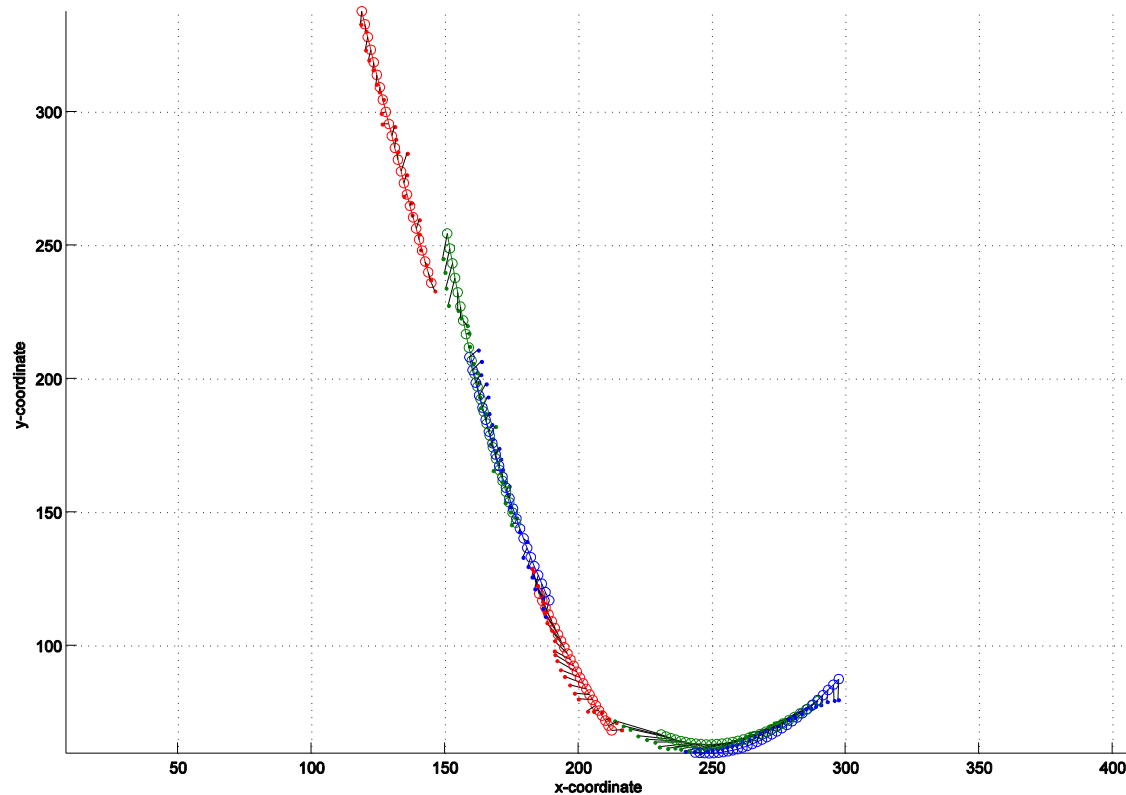
Experiment 3



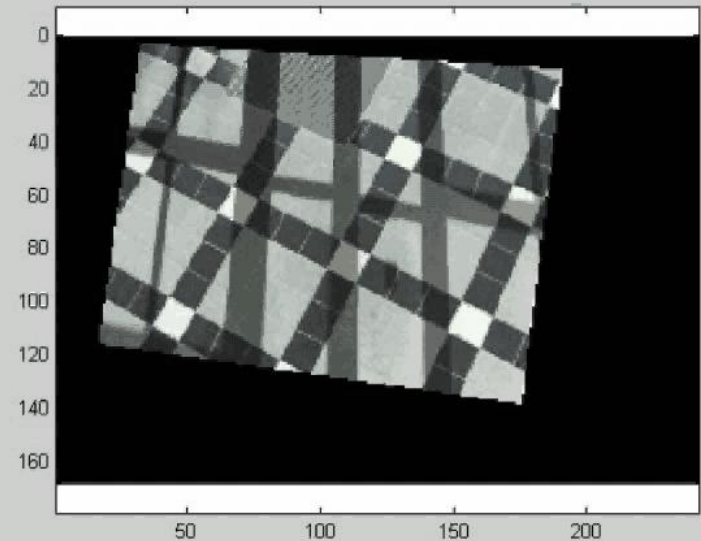
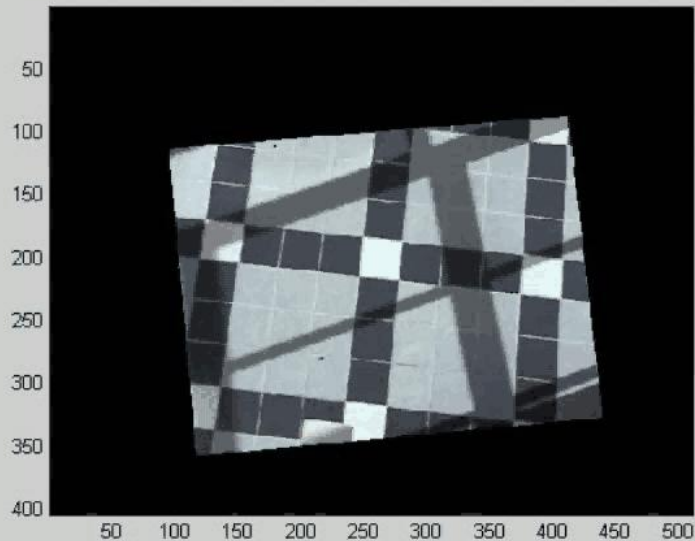
Before and After



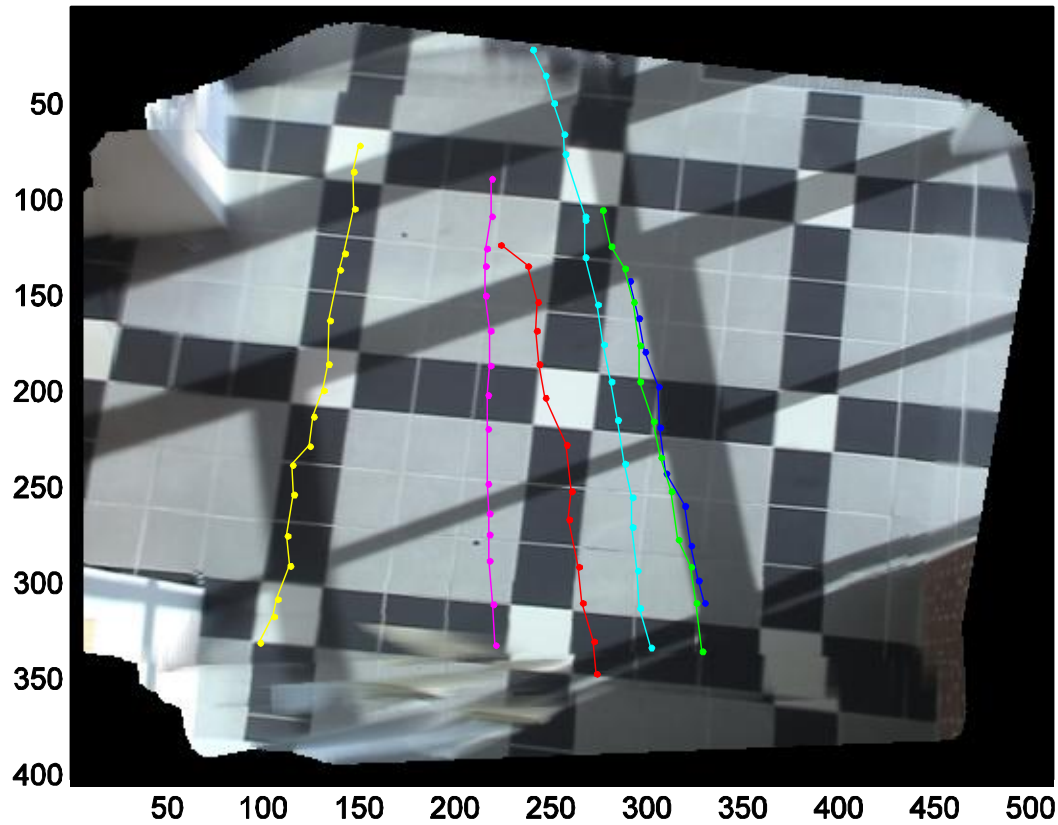
Before and After



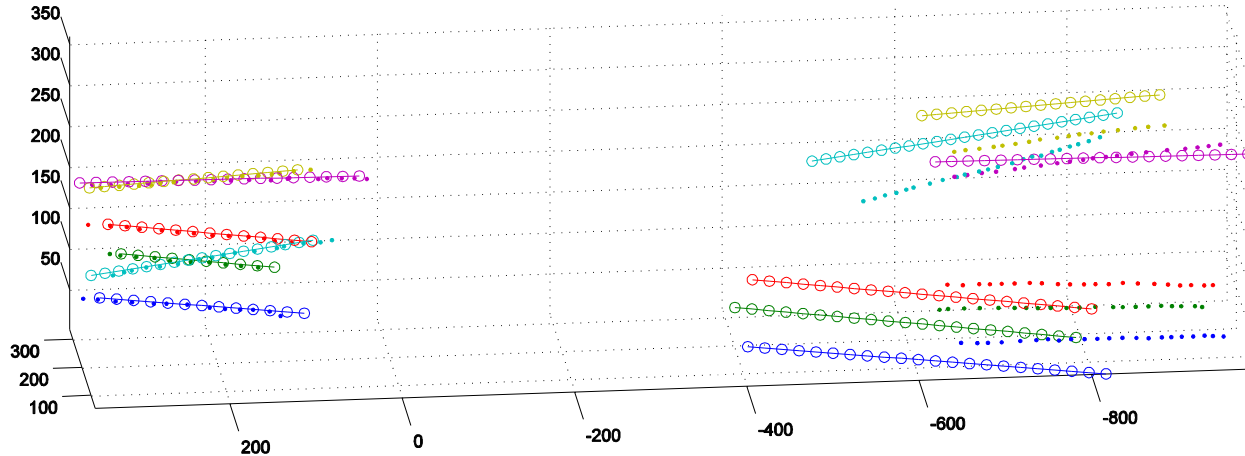
CSB Sequence



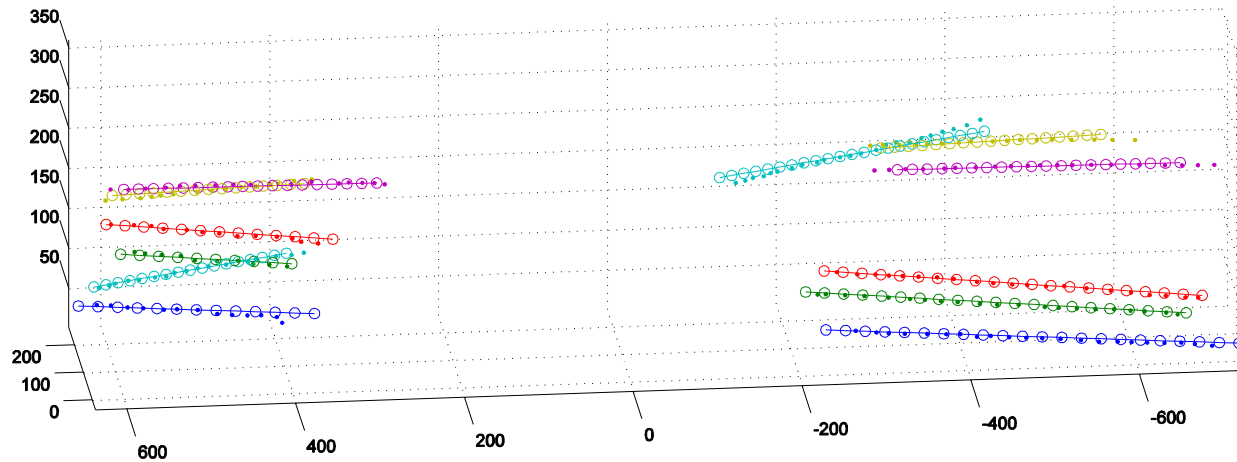
Experiment 5



Experiment 5



Experiment 5



HUMAN ACTION RECOGNITION

Part III



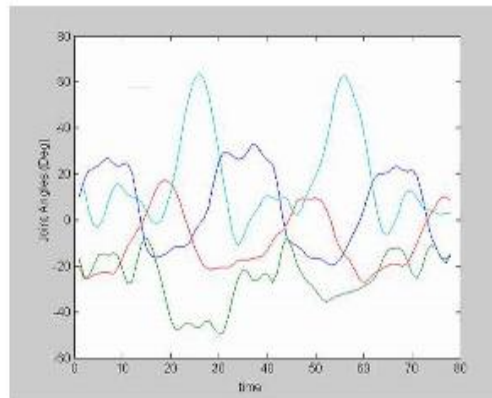
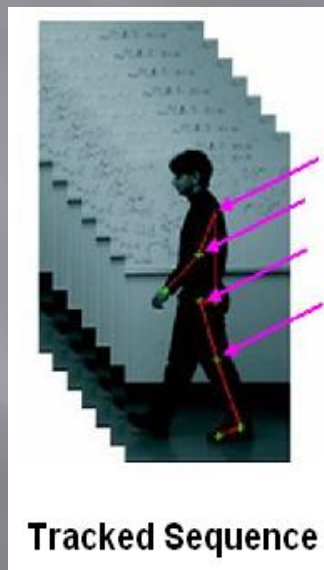
University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Actions in the Computer Vision Literature

- Categorization on the basis of *representation*



Contents

- ▣ View Invariant Action Recognition
- ▣ Actions As Objects
- ▣ Anthropometric Representation for Invariant Action Recognition
- ▣ Action Recognition In Two Moving Cameras
- ▣ Chaotic Invariants for Human Action Recognition



VIEW INVARIANT REPRESENTATION & RECOGNITION OF ACTIONS

Cen Rao, Alper Yilmaz, Alexi Gritai

IJCV 2002

ICCV 2003



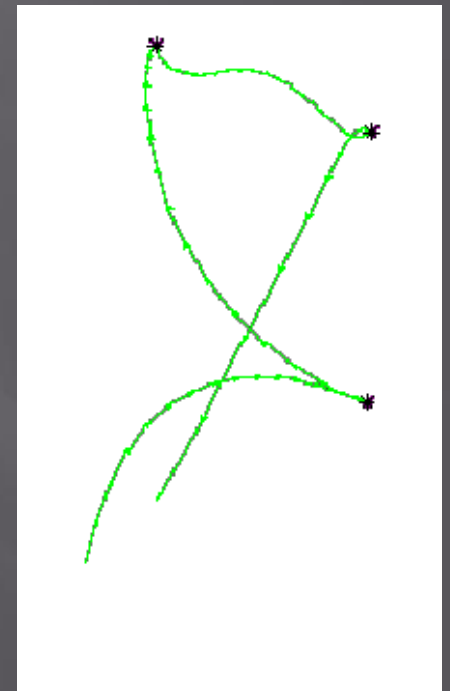
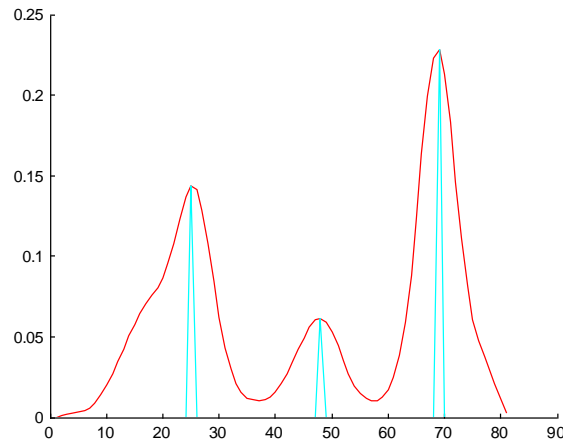
University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Representation of Actions

- Dynamic Instants:
 - Maximum in spatiotemporal curvature represents an important change of motion characteristic.
- Intervals

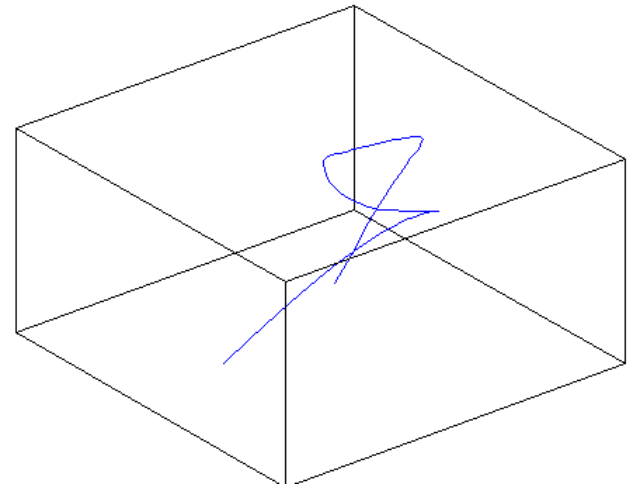
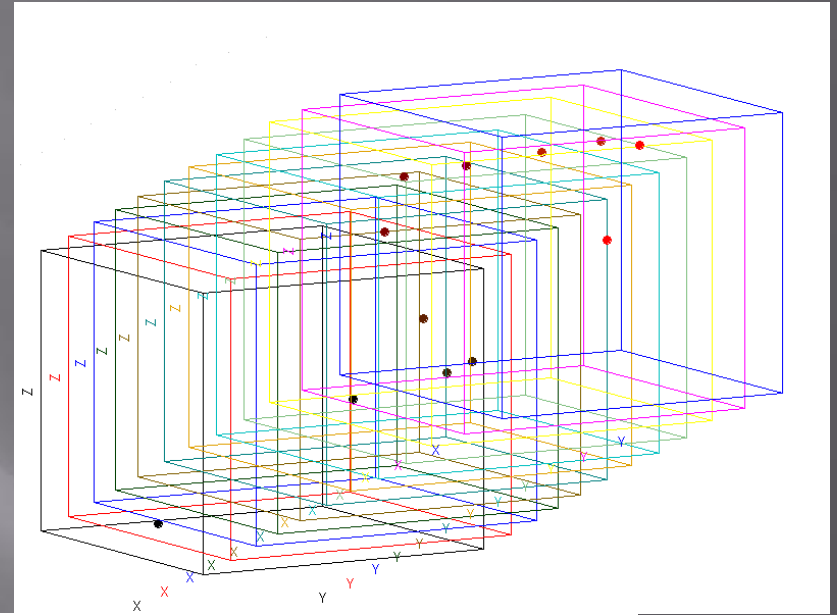


Action Trajectory in 4D



Sampling in time

Ignore the
time index

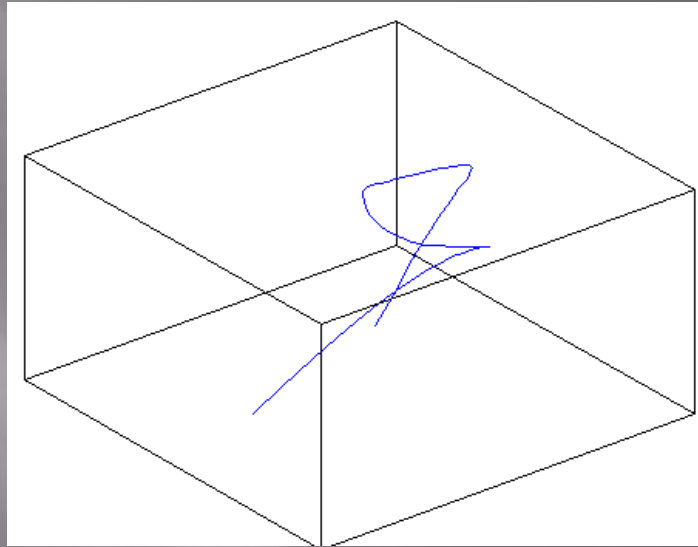


University of
Central Florida

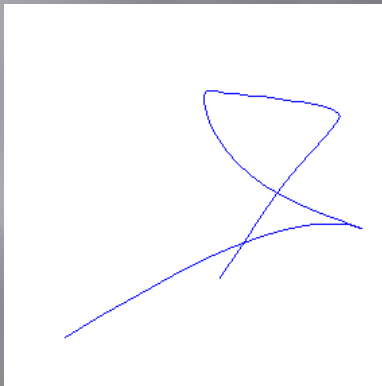
VISION

Copyrights Mubarak Shah, UCF

Viewing
directionn1

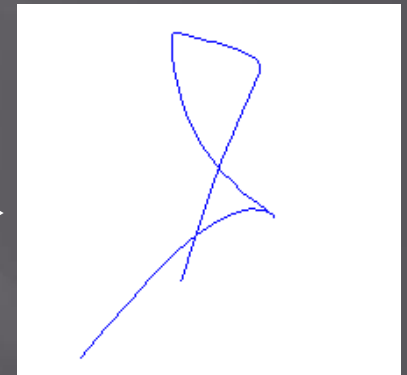


Viewing
directionn2



2D trajectory

Matching?



2D trajectory



University of
Central Florida

VISION

Affine View Invariant Matching

Rank Theorem (Tomasi & Kanade)

- S is a set of 3-D points and Π s are projection matrices for different viewpoints, then we can arrange image coordinates of points in an observation matrix, M , as follows:

$$M = \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ y_1^1 & y_2^1 & \dots & y_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ y_1^2 & y_2^2 & \dots & y_n^2 \end{bmatrix} = P S \quad \begin{matrix} \Pi_1 \\ \Pi_2 \\ \vdots \\ \Pi_n \end{matrix} \begin{bmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ y_1^1 & y_2^1 & \dots & y_n^1 \\ z_1^1 & z_2^1 & \dots & z_n^1 \end{bmatrix}$$

$$\Pi_i = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \end{bmatrix}$$

M is 4 by n , P is 4×3 and S is $3 \times n$,
then the rank of M is at most 3.



Generalized Affine Rank Theorem

- ▣ A set of image points match *if and only if* M is of rank at most 3. (Shapiro & Zisserman, Seitz & Dyer)
- ▣ A set of “instants” match *if and only if* M of rank at most 3. Therefore, the similarity measure is:

$$M = \begin{bmatrix} \mu_1^i & \mu_2^i & \dots & \mu_n^i \\ v_1^i & v_2^i & \dots & v_n^i \\ \mu_1^j & \mu_2^j & \dots & \mu_n^j \\ v_1^j & v_2^j & \dots & v_n^j \end{bmatrix} \quad dist = 0$$



Perspective View-Invariant Matching

- ▣ Fundamental matrix captures the relationship between the corresponding points in two views.

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix}^T F \begin{bmatrix} u'_i \\ v'_i \\ 1 \end{bmatrix} = 0, \quad F = \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix}$$



Perspective View-invariant Measure

- Consider the fundamental matrix constraint and rearrange the constraint as following:

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix}^T F \begin{bmatrix} u'_i \\ v'_i \\ 1 \end{bmatrix} = 0, \quad M = \begin{bmatrix} u_1 v'_1 & u_1 v'_2 & u_1 & v'_1 & v'_2 & 1 & u_2 v'_1 & u_2 v'_2 & u_2 & v'_1 & v'_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

M is 9 by n matrix

To solve the equation, the rank(M) must be 8.

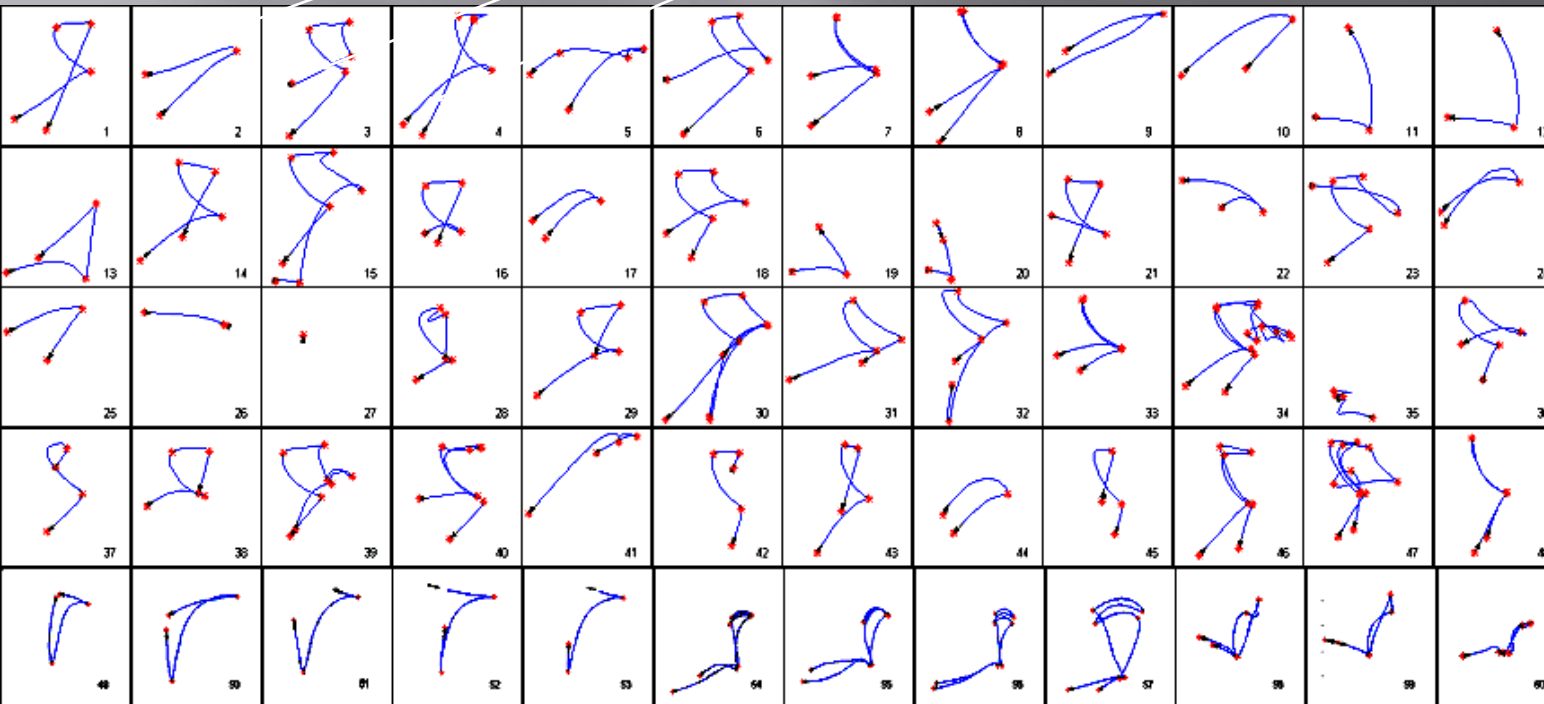
The 9th singular value of M, σ_9 , is the match measure.



Experimental Results

60 Action Trajectories

7 People



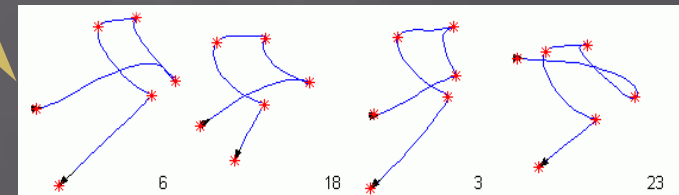
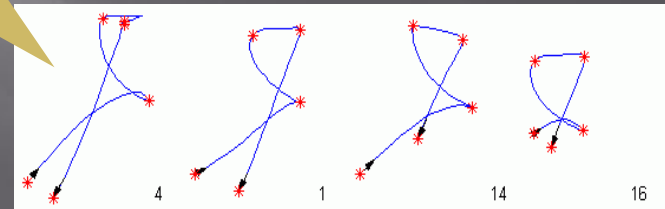
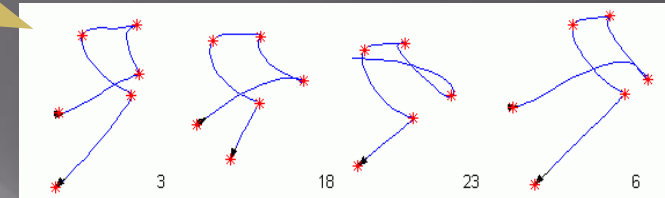
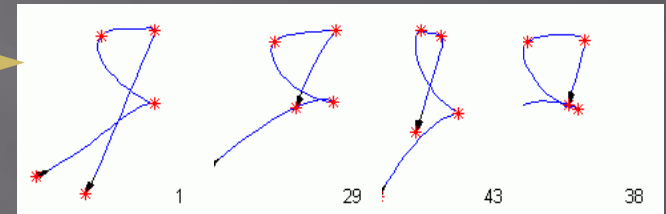
University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

Experimental Results

Actions	3 Best matches	Evaluation & comments
1	29 43 38	Correct
2	Pick up	Correct
3	18 23 6	Correct
4	1 14 16	One wrong
5		Unique action
6	18 3 23	Correct
7	48 33 8	correct
8	48 33 7	One wrong
9	Pick up	Correct
10	Put down	Correct
11	Pick up	Correct
12	Put down	Correct
13		Unique action
14	43 16 1	Correct
15		Unique action
16	14 29 1	Correct
17	Pick up	Incorrect, object hidden
18	6 3 23	Correct
19	Pick up	Correct
20		Unique random motion



Temporal Alignment Results

Before Temporal Alignment



After Temporal Alignment

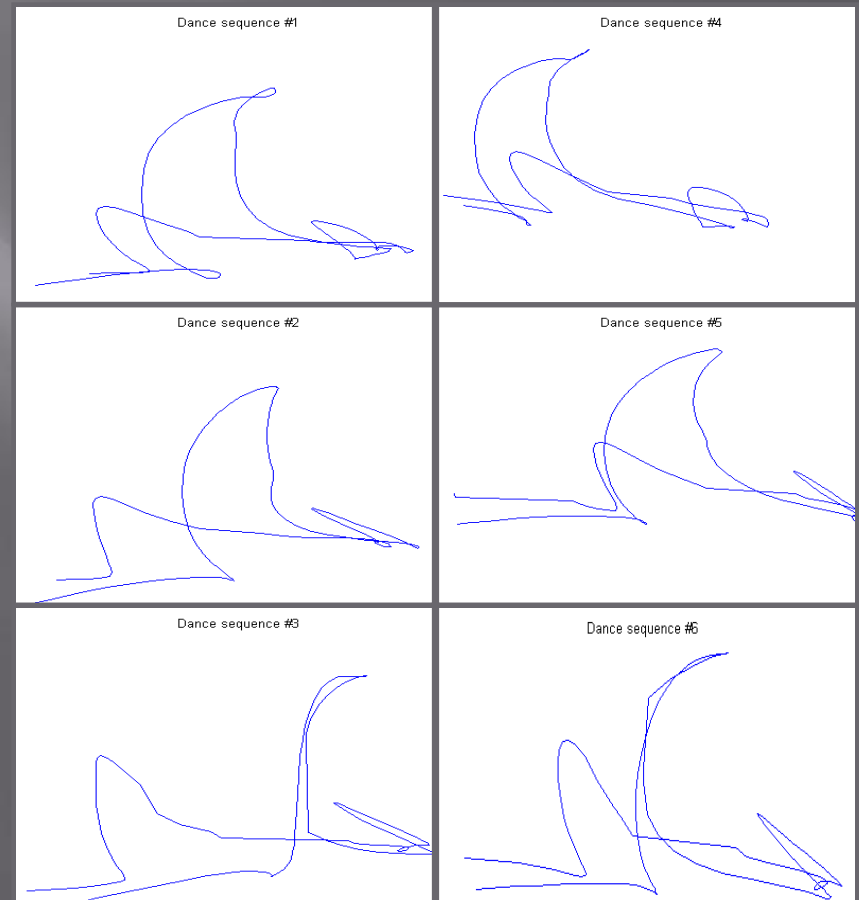


Temporal Alignment of Videos

Input videos:



Trajectories of the right foot:



Temporal Alignment Results

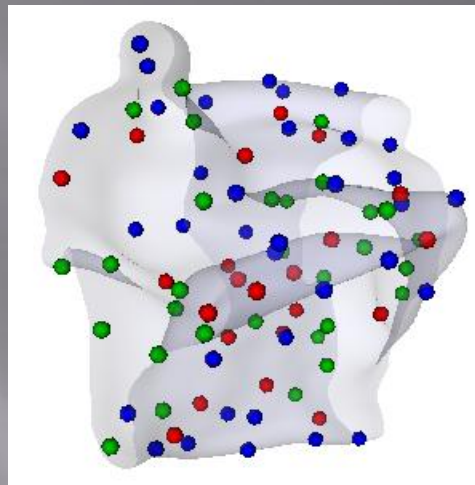
Synchronized videos:



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF



ACTION AS OBJECTS

Alper Yilmaz

- A. Yilmaz and M. Shah "Actions Sketch: A Novel Action Representation," IEEE Conf. on Computer Vision and Pattern Recognition (CVPR), 2005.
- A. Yilmaz and M. Shah "Representing Actions Using Differential Geometry," Computer Vision and Image Understanding (CVIU), submitted 2006.



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Actions As Objects

When something moves it develops a shape.

Santiago Calatrava

(Sculpture into architecture)

Milwaukee Museum of Art

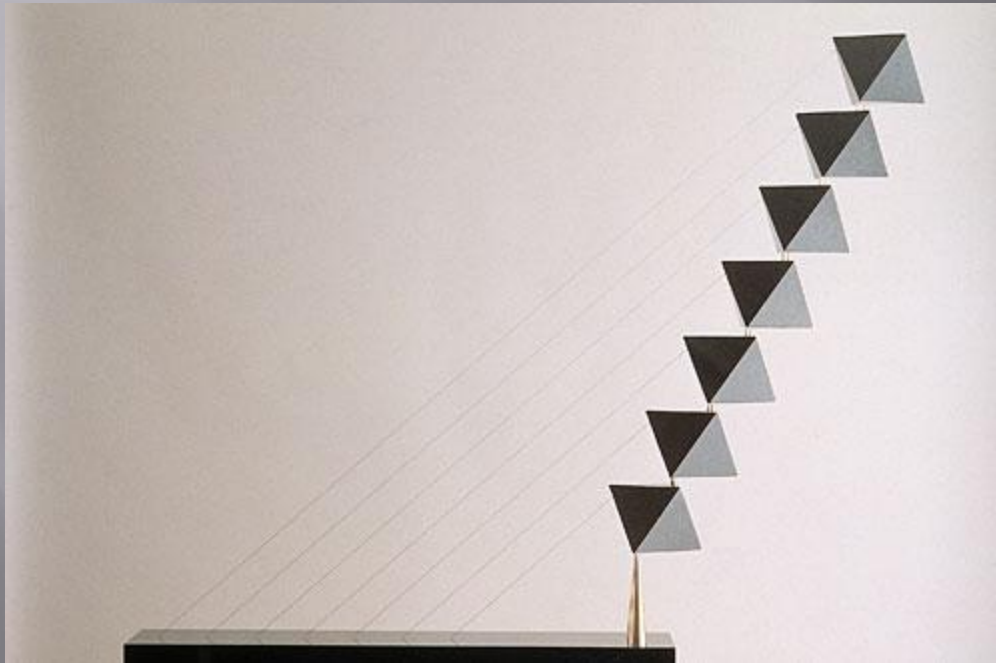


University of
Central Florida

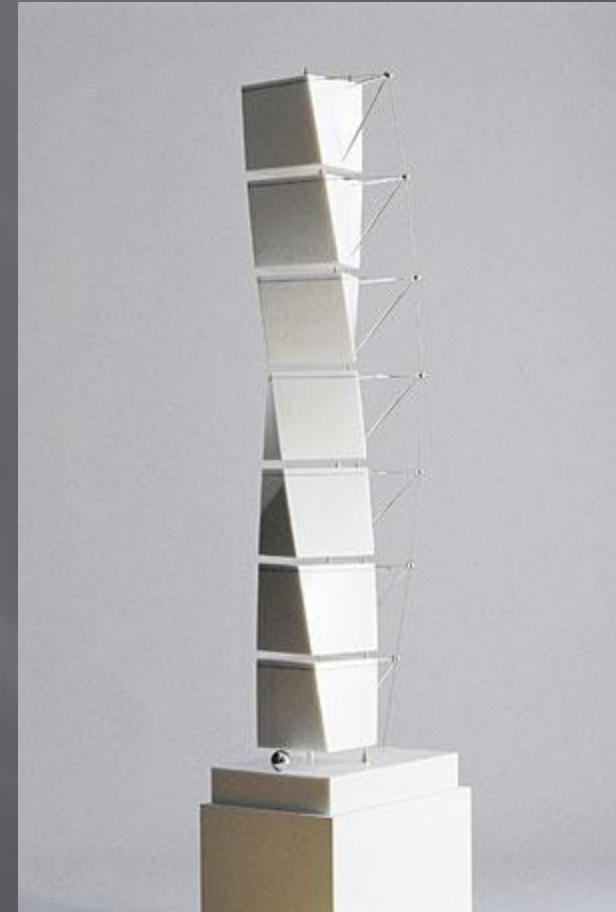
VISION

Copyrights Mubarak Shah, UCF

Actions As Objects



Musical Star

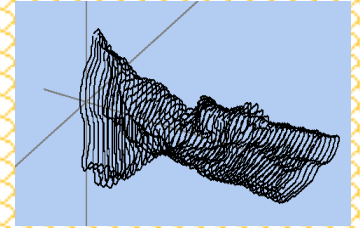


Turning Torso

Flow diagram

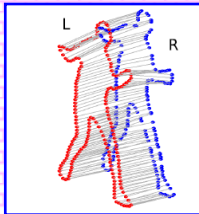
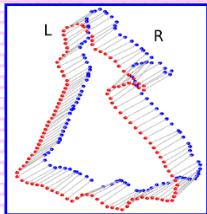


Contour Extraction



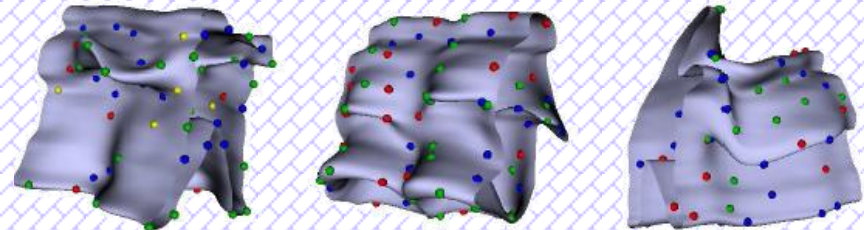
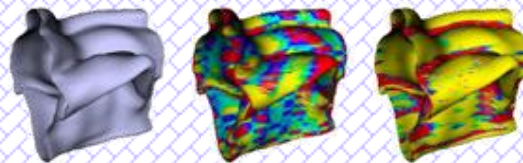
Action Volume Generation

- Graph theoretic volume generation
- volume smoothing



Feature Extraction & Recognition

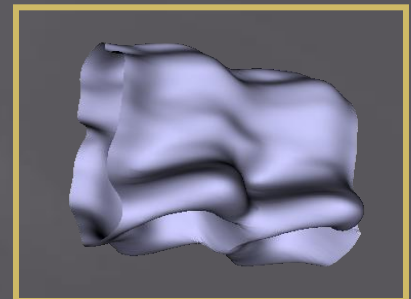
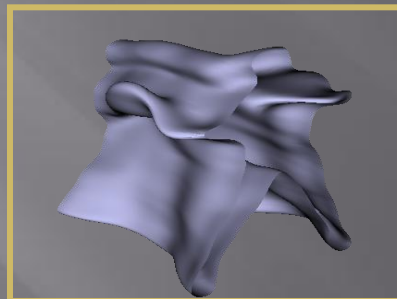
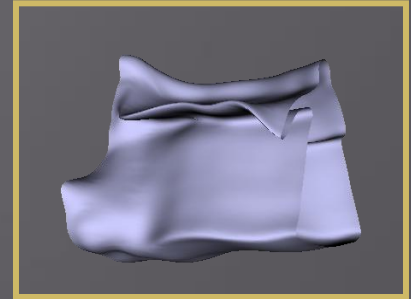
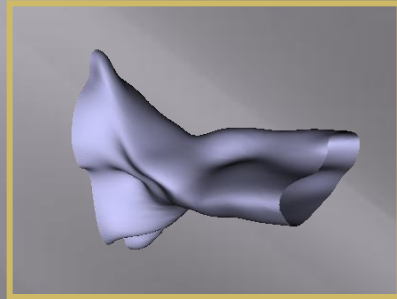
- Differential geometry
- Epi-polar geometry



	$K > 0$	$K = 0$	$K < 0$
$H < 0$	peak 	ridge 	saddle ridge
$H = 0$	none 	flat 	minimal
$H > 0$	pit 	valley 	saddle valley



Resulting Volume



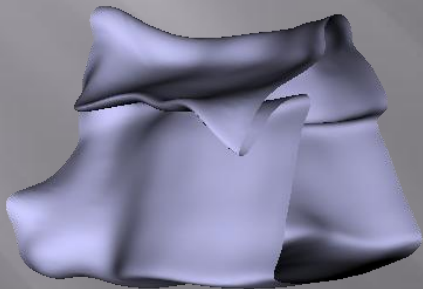
University of
Central Florida

VISION

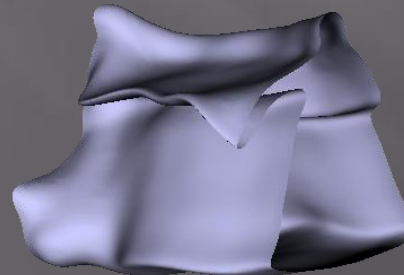
Copyrights Mubarak Shah, UCF

Properties of the Action Volume

- ▣ Space-time (3D) object
- ▣ Encodes shape and motion
- ▣ Uses complete object contours instead of a single point on the object.
- ▣ Suitable for fine action analysis
- ▣ Continuous representation
 - Same volume for same action of different lengths



40 frames



20 random
selected frames



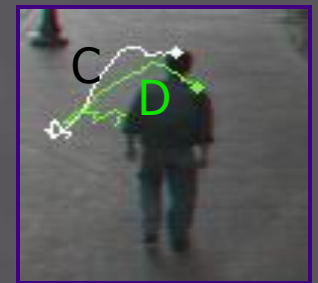
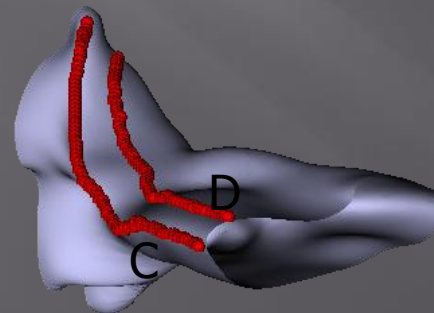
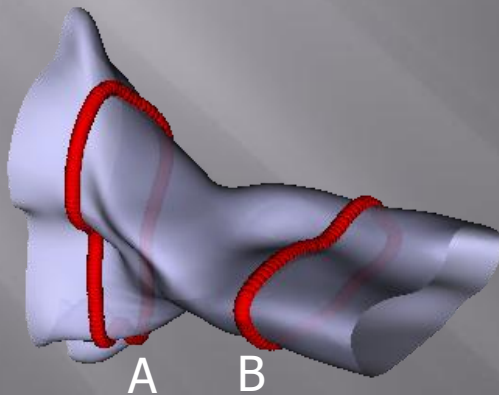
University of
Central Florida

VISION

Copyright Mubarak Shah, UCF
Copyrights Mubarak Shah, UCF

Properties of the Action Volume

- Can be represented in 2D
 - Arc length and time
- Can regenerate contour at time t
- Can provide spatial trajectory of contour points



What is the Action Sketch?

- ▣ Important action descriptors
 - Unique shape and motion characteristics
- ▣ Related to differential geometric properties of action volume
 - 1st and 2nd fundamental forms
 - ▣ Gaussian and mean curvatures
 - ▣ Fundamental surface types



Computing Gaussian (K) and Mean (H) Curvatures

- ▣ K and H are two algebraic invariants of Weingarten mapping S .

$$K = \det(S)$$

$$H = \frac{1}{2} \text{tra } S$$

$$\mathbf{g} = \begin{bmatrix} f_s f_s & f_s f_t \\ f_s f_t & f_t f_t \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} f_s \vec{n} & f_t \vec{n} \\ f_t \vec{n} & f_{tt} \vec{n} \end{bmatrix}$$

$$S = \mathbf{g}^{-1} \mathbf{b}$$

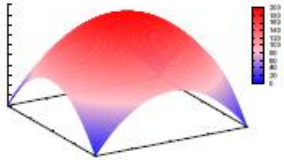
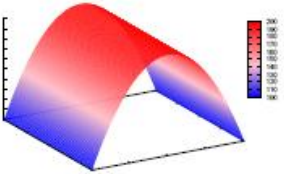
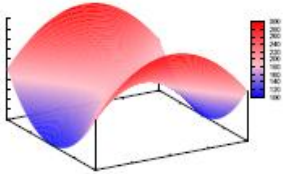
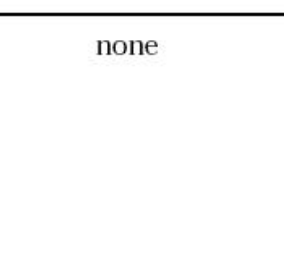
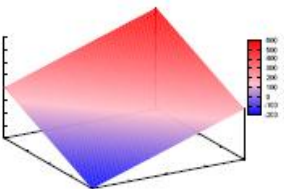
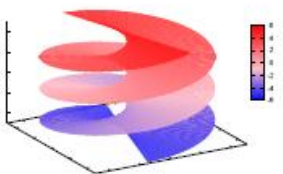
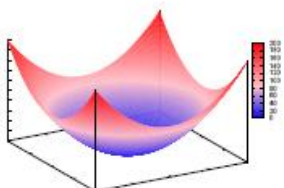
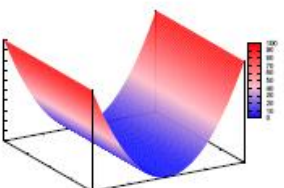
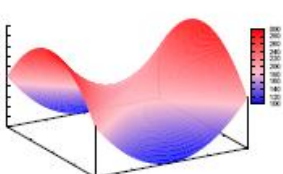
from 1st fundamental form

from 2nd fundamental form

where $f(s,t)$ is a point on the volume, n is normal at f



Fundamental Surface Types

	$K > 0$	$K = 0$	$K < 0$
$H < 0$	peak 	ridge 	saddle ridge 
$H = 0$	none 	flat 	minimal 
$H > 0$	pit 	valley 	saddle valley 



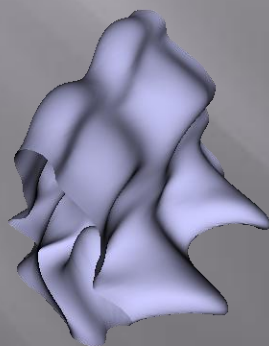
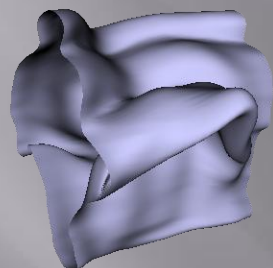
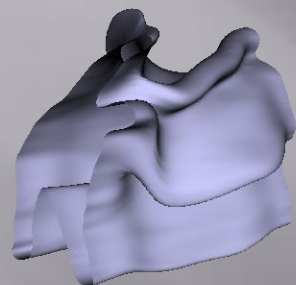
Properties of Surface Types

- ▣ Rotation and translation invariant in spatio-temporal space.
- ▣ Encodes intrinsic properties of surface.
 - Defines the convexity or concavity of surface.
- ▣ Related to speed and acceleration.

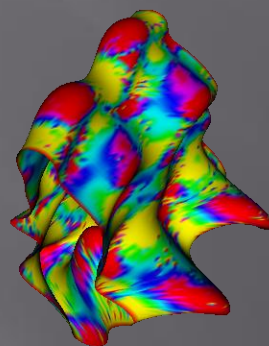
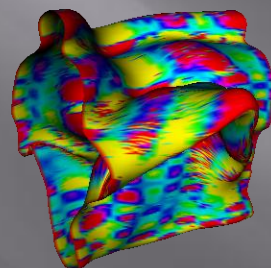
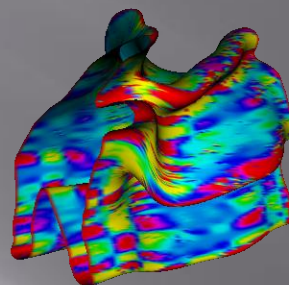


Differential Geometric Surface

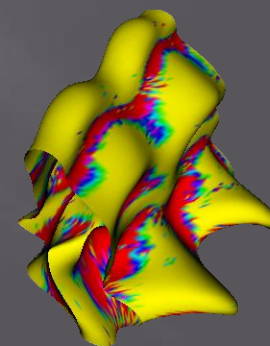
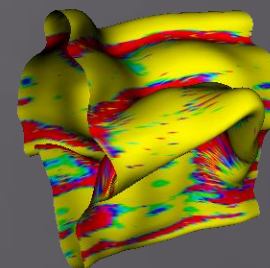
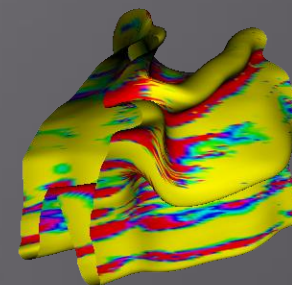
Action Volume



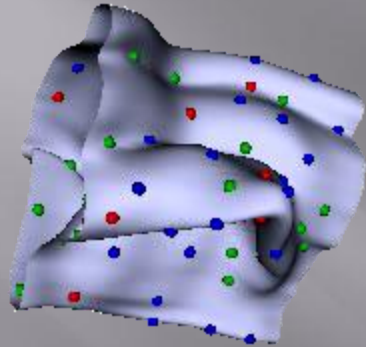
Gaussian Curvature



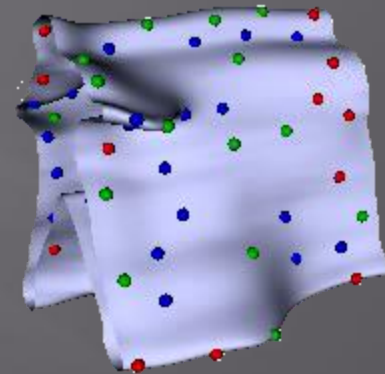
Mean Curvature



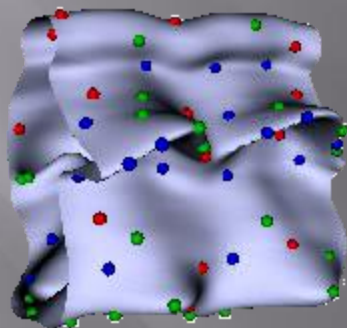
Examples



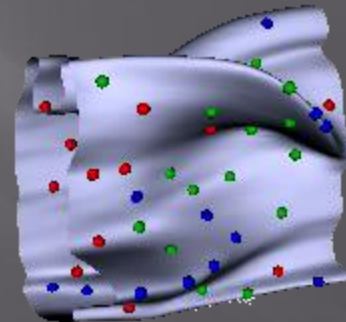
kicking



dance



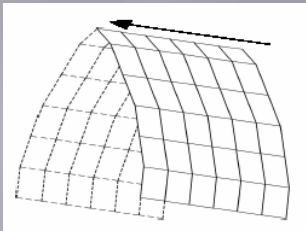
walking



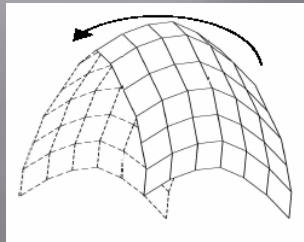
surrender



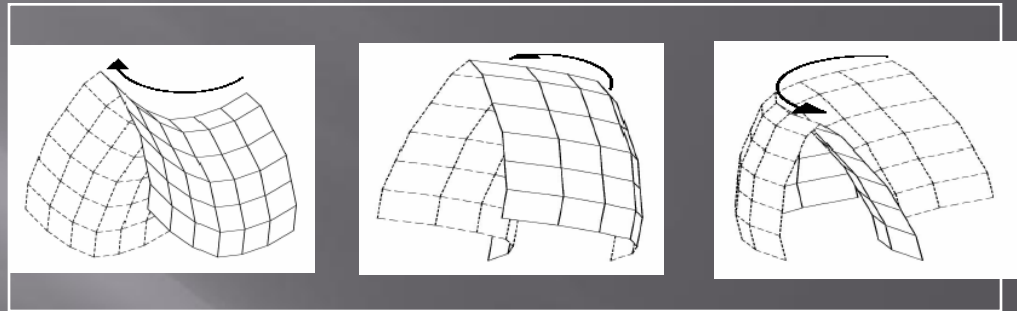
Surface patches & their relation to the object motion



ridge



peak

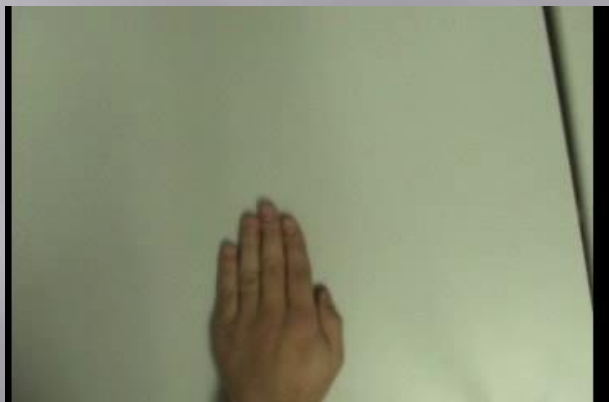


saddle ridge



Action Descriptors Relation to Object Motion

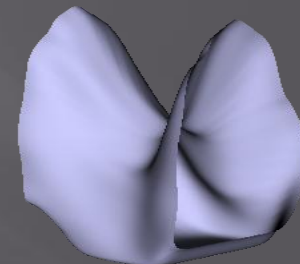
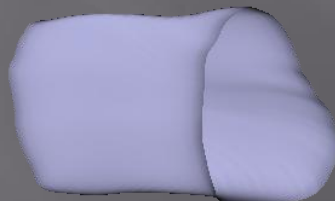
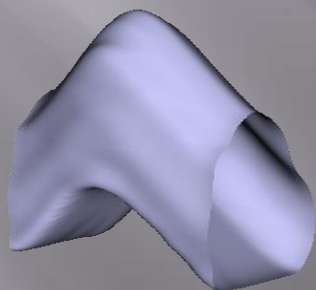
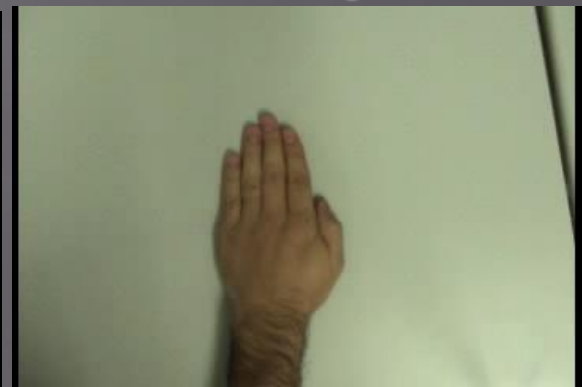
peak



ridge

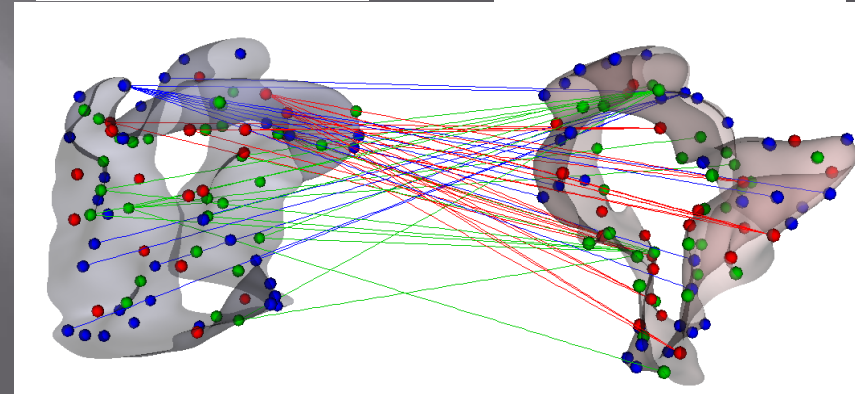
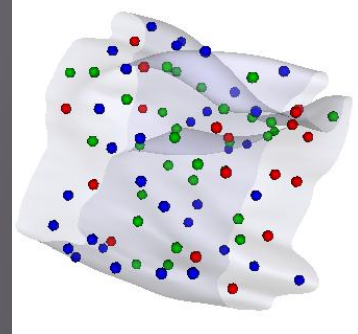
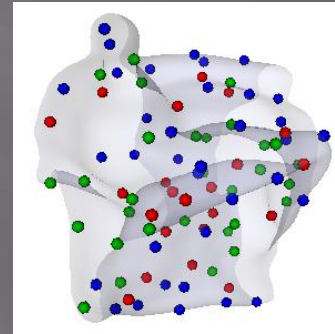


saddle ridge

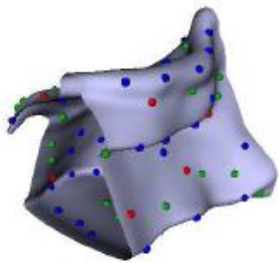


Matching Volumes: Establishing Correspondence

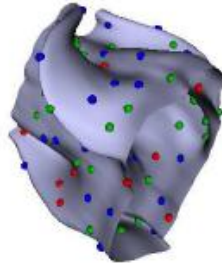
- Generate bipartite action graphs
 - Space-time proximity
 - Shape similarity
- Find Maximum Matching



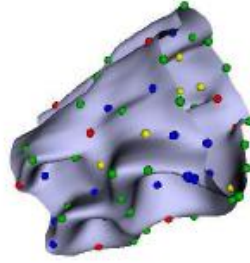
Action Volumes



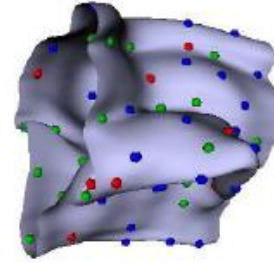
1) dance



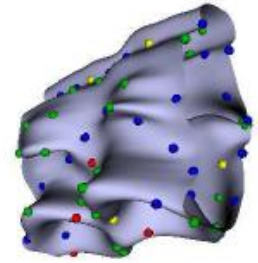
2) hand down



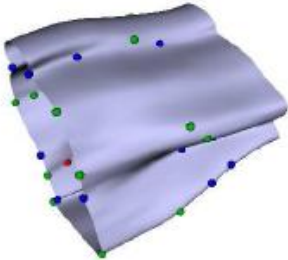
3) walk



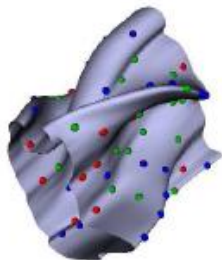
4) kick



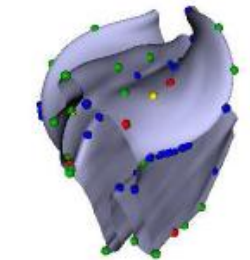
5) walk



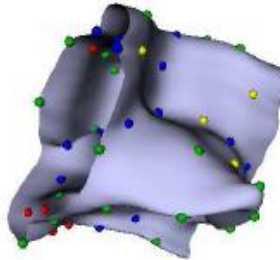
6) stand up



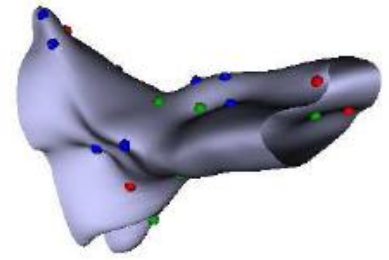
7) surrender



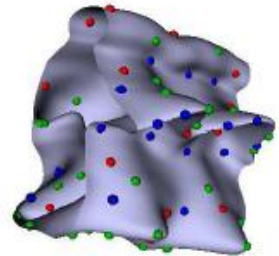
8) hand down



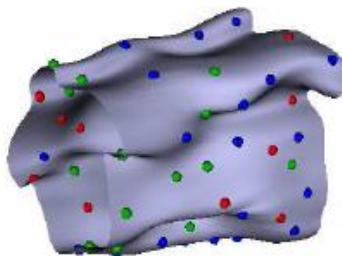
9) kick



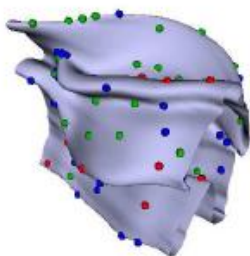
10) fall



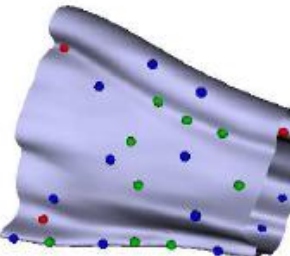
11) walk



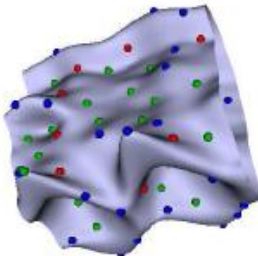
12) walk



13) aerobic 1

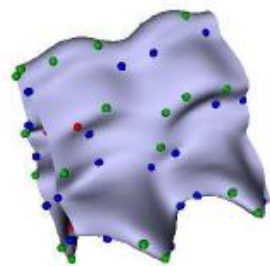


14) sit down

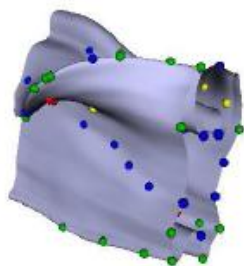


15) walk

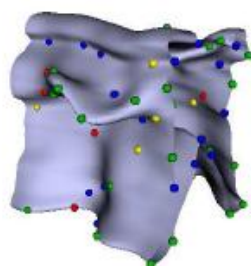
Action Volumes



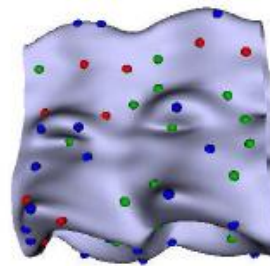
16) running



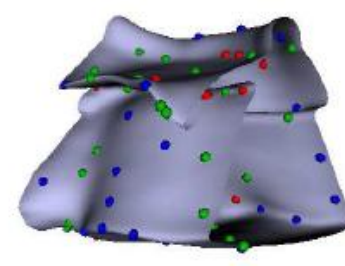
17) surrender



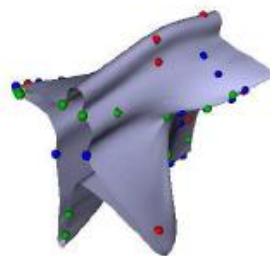
18) stroke



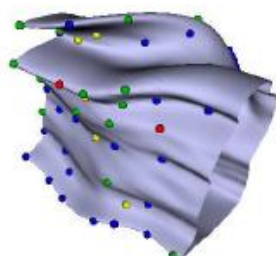
19) walk



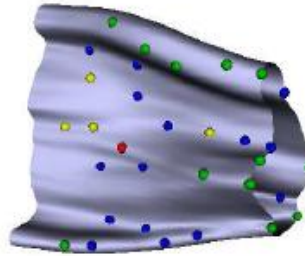
20) dance



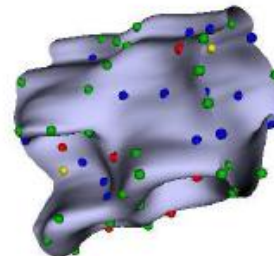
21) aerobic 2



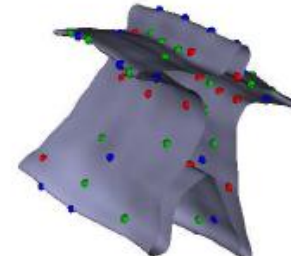
22) aerobic 3



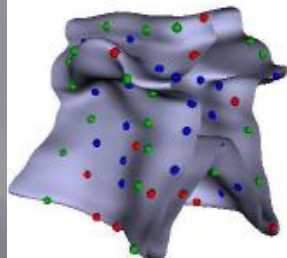
23) sit down



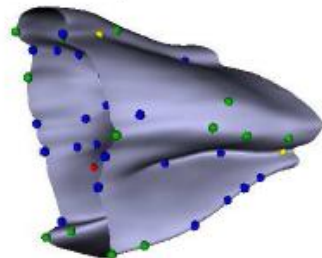
24) walk



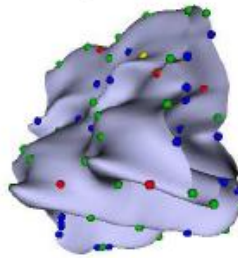
25) aerobic 4



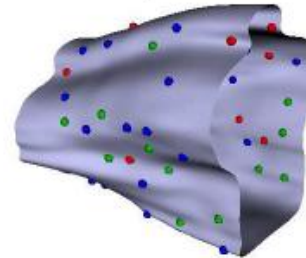
26) stroke



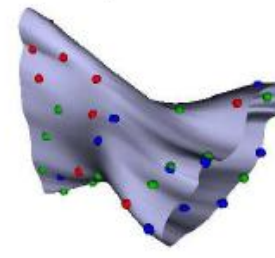
27) stand up



28) running



29) stand up



30) falling

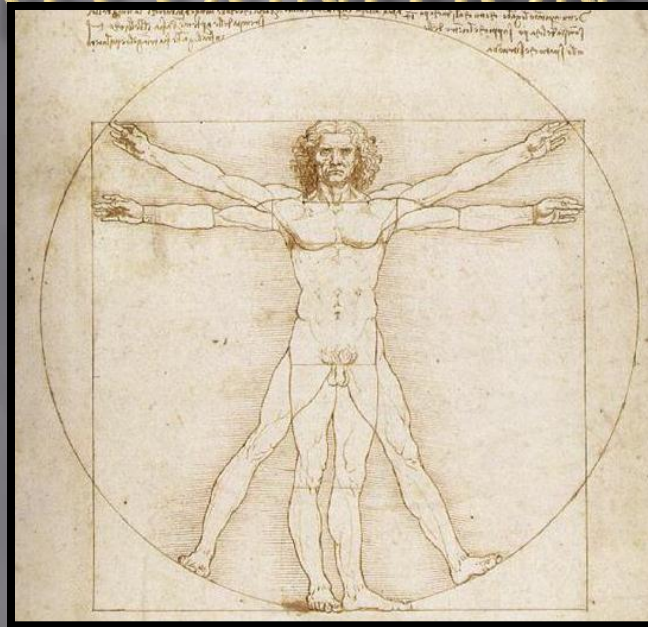
Recognition Results

Input Action	#	Matching action	#
Dance	1	<i>Dance</i>	20
Hand down	2	Stand up	29
Walking	3	<i>Walking</i>	11
Kicking	4	<i>Kicking</i>	9
Walking	5	<i>Walking</i>	11
Stand up	6	<i>Stand up</i>	29
Surrender	7	<i>Surrender</i>	17
Hands down	8	<i>Hands down</i>	82
Kicking	9	<i>Kicking</i>	4
Falling	10	<i>Falling</i>	30
Walking	11	<i>Walking</i>	11
Walking	12	Sit down	23
Sit down	14	<i>Sit down</i>	23

Video	#	Matching action	#
Walking	15	<i>Walking</i>	11
Running	16	<i>Running</i>	28
Surrender	17	<i>Surrender</i>	17
Tennis stroke	18	<i>Tennis stroke</i>	26
Walking	19	<i>Walking</i>	11
Dance	20	<i>Dance</i>	1
Sit down	23	<i>Sit down</i>	23
Walking	24	<i>Walking</i>	11
Tennis stroke	26	<i>Tennis stroke</i>	18
Stand up	27	<i>Stand up</i>	29
Running	28	<i>Running</i>	16
Stand up	29	Hands down	8
Falling	30	<i>Falling</i>	10



ANTHROPOMETRIC REPRESENTATION FOR INVARIANT ACTION RECOGNITION

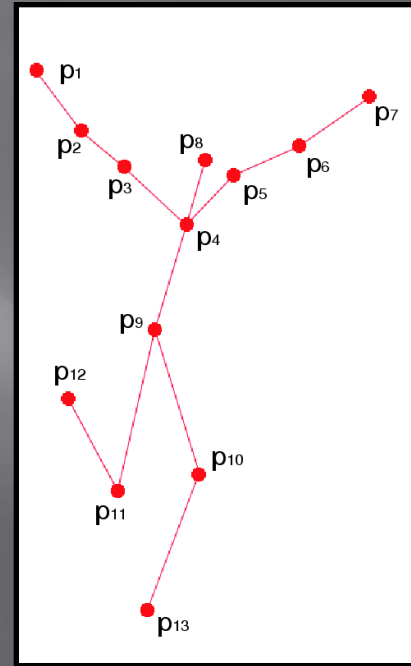
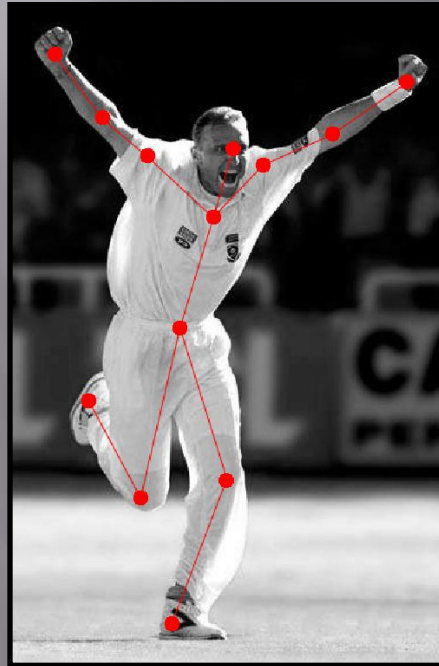


University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Representation of Actors



- Point-based model contains sufficient description for the recognition of human actions, [1].

[1] G. Johansson. Visual perception of biological motion and a model for its analysis. Perception and Psychophysics, 14(2): 201 – 211, 1993.

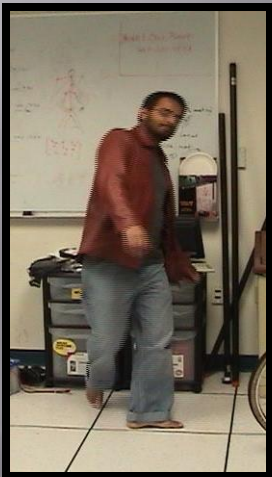
Anthropometry

- ▣ \An`thro*pom"e*try\, *n*. Measurement of the height and other dimensions of human beings, especially at different ages, or in different races, occupations, etc.
- ▣ Variability in human proportion is not *arbitrary*.
- ▣ Action Recognition must address this variation.



Pose and Posture

- Posture: The stance an actor has at a time instant
- Pose: The global orientation and position of an actor



Different Poses, Same Posture



Different Postures, Same Pose

Anthropometric Constraint

- ▣ **Conjecture:** The relationship between points of two actors X and Y in the same posture can be described by a matrix M

$$X_i = M Y_i$$

where $i = 1, 2 \dots n$, M is a 4×4 non-singular matrix, X_i and Y_i are sets of points describing two actors.

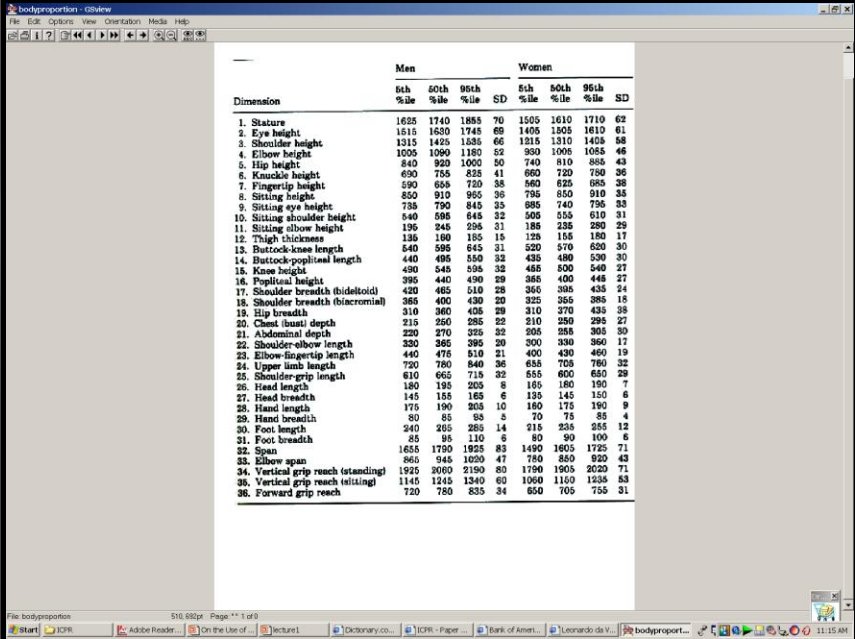
- ▣ This transformation simultaneously captures:
 - the different poses
 - difference in size/proportions.



Anthropometric Constraint

- This was verified empirically between the 5th percentile woman and 95th percentile man.

- Mean error of
 - 227.3 mm before the transformation,
 - 23.87 mm after the transformation.



Dimension	Men				Women			
	5th %ile	50th %ile	95th %ile	SD	5th %ile	50th %ile	95th %ile	SD
1. Stature	1625	1740	1865	70	1555	1610	1710	62
2. Eye height	1515	1620	1745	65	1405	1505	1610	61
3. Shoulder height	1315	1425	1535	66	1215	1310	1405	58
4. Elbow height	1005	1090	1180	52	950	1005	1055	46
5. Hip height	840	920	1000	50	740	810	865	43
6. Kneecap height	690	755	825	41	660	720	780	36
7. Fingertip height	590	655	720	36	560	625	685	28
8. Sitting height	850	910	965	36	795	850	910	35
9. Sitting eye height	735	790	845	35	685	740	795	33
10. Sitting shoulder height	540	595	645	32	505	555	610	31
11. Sitting elbow height	195	245	295	31	185	235	280	29
12. Thigh thickness	135	160	185	15	125	155	180	17
13. Buttock-knee length	540	595	645	31	520	570	620	30
14. Buttock-popliteal length	440	495	550	32	435	480	530	30
15. Knee height	430	545	595	32	465	500	540	27
16. Popliteal height	395	440	490	29	365	400	445	27
17. Shoulder breadth (biacromial)	420	465	510	28	365	395	435	24
18. Shoulder breadth (biacromial)	365	400	430	20	325	355	385	18
19. Hip breadth	310	360	405	29	310	370	435	38
20. Chest (bust) depth	215	250	285	22	210	250	295	27
21. Abdominal depth	220	270	325	32	205	255	305	30
22. Shoulder-elbow length	320	365	395	20	300	330	360	17
23. Elbow-fingertip length	440	475	510	21	400	430	460	19
24. Upper limb length	720	780	840	36	655	705	760	32
25. Shoulder-grip length	610	665	715	32	555	600	650	29
26. Head length	180	195	205	6	165	180	190	7
27. Head breadth	145	155	165	6	135	145	150	6
28. Hand length	175	190	205	10	160	175	190	9
29. Hand breadth	80	85	90	5	75	80	85	4
30. Foot length	240	265	285	14	215	235	255	12
31. Foot breadth	85	95	110	6	80	90	100	6
32. Span	1655	1750	1825	83	1490	1605	1725	71
33. Elbow span	865	945	1050	47	780	880	920	43
34. Vertical grip reach (standing)	1925	2060	2190	80	1790	1905	2020	71
35. Vertical grip reach (sitting)	1145	1245	1340	60	1060	1160	1255	53
36. Forward grip reach	720	780	835	34	650	705	755	31

R. Bridger. *Human Performance Engineering: A Guide for system designers*, Prentice Hall, 1982



Postural Constraint

- ▣ **Proposition 1:** If x_t and y_t describe the imaged posture of two actors at time t , a Fundamental Matrix can be uniquely associated with (x_t, y_t) if the two actors are in the same posture.

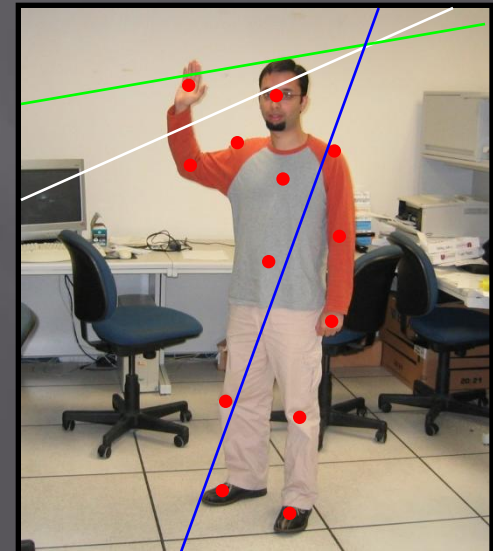
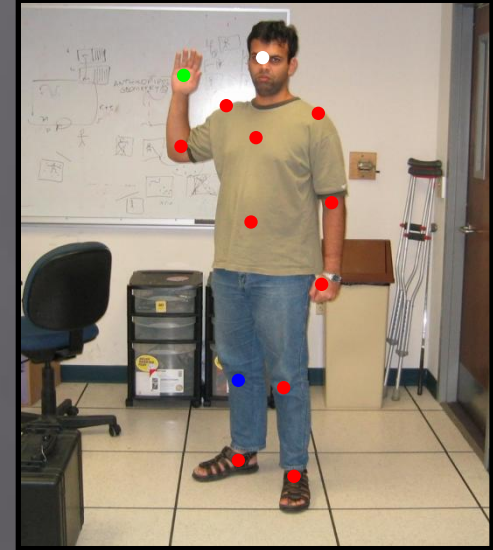
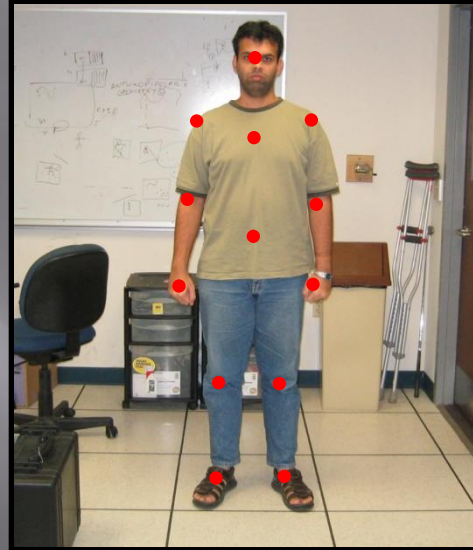
$$x_t^T F y_t = 0$$

- ▣ Two actors performing the action instead of two views.
- ▣ This is valid for a single time instance.



Capturing View Variance

- ▣ The fundamental matrix captures the variability in proportion as well as the change in view.



Postural Constraint

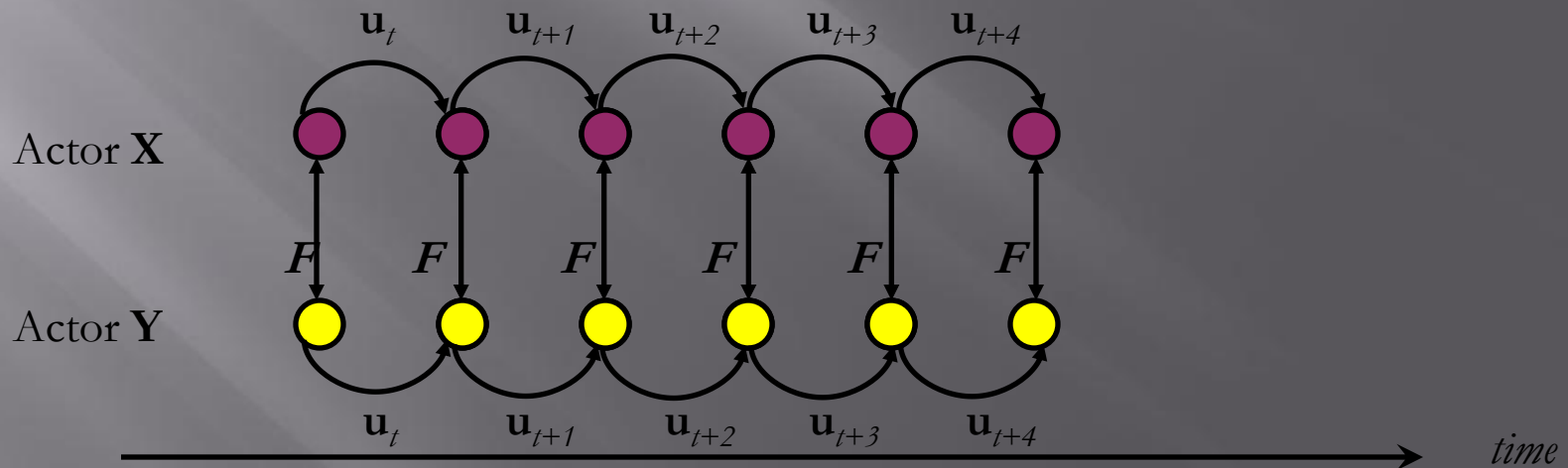
- ▣ The similarity of *posture* between two actors can be measured using the **ninth singular** value of a measurement matrix A , where $Af = 0$.

$$\begin{bmatrix} x'_1 x_1 & \dots & x'_n x_n \\ x'_1 y_1 & \dots & x'_n y_n \\ x'_1 & \dots & x'_n \\ y'_1 x_1 & \dots & y'_n x_n \\ y'_1 y_1 & \dots & y'_n y_n \\ y'_1 & \dots & y'_n \\ x_1 & \dots & x_n \\ y_1 & \dots & y_n \\ 1 & \dots & 1 \end{bmatrix}^T \begin{bmatrix} F_{11} \\ F_{12} \\ F_{13} \\ F_{21} \\ F_{22} \\ F_{23} \\ F_{31} \\ F_{32} \\ F_{33} \end{bmatrix} = Af = 0$$



Action Constraint

- ▣ **Proposition 2:** For an action element \mathbf{u}_t , the fundamental matrices associated with $(\mathbf{x}_t, \mathbf{y}_t)$ and $(\mathbf{x}_{t+1}, \mathbf{y}_{t+1})$ are the same if both actors perform the action element defined by \mathbf{u}_t .



Measuring Action Similarity

- Since all the F s are the same:

$$A_1 f = 0$$

$$A_2 f = 0$$

$$\vdots$$

$$A_k f = 0$$

- Thus the ninth singular value of

$$A = [A_1, A_2 \dots A_k]$$

can be used as a view invariant measure.



Experimental Results

- ▣ We performed a diverse set of experiments
 - Action Detection
 - ▣ Analyzing periodicity
 - ▣ Multiple view multiple people
 - Action Synchronization
 - ▣ Following the leader
 - Odd one out

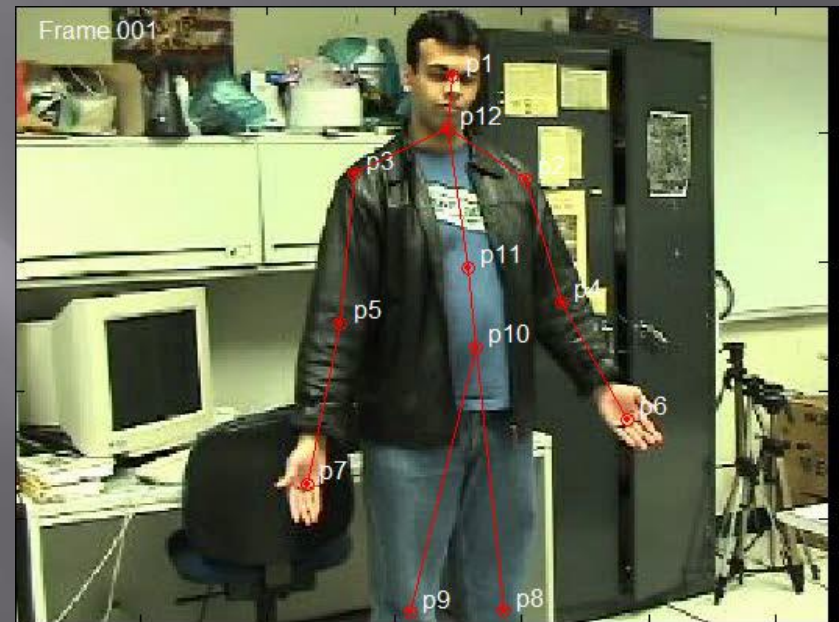


Action Detection

Analyzing Periodicity



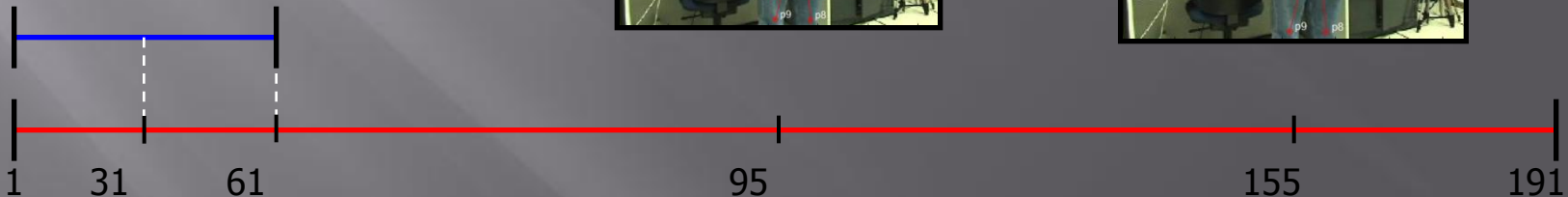
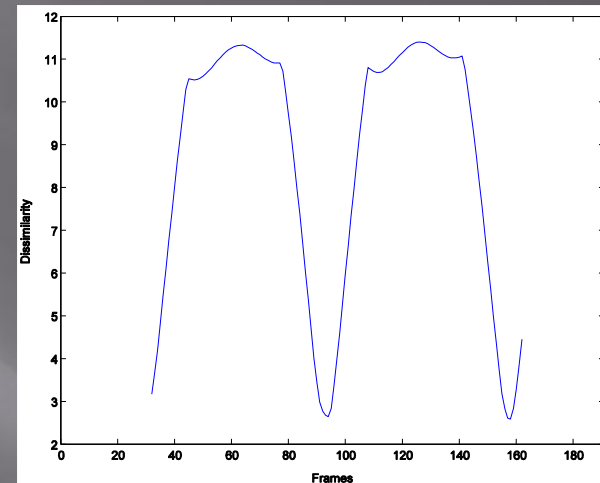
Reference Pattern



Test Sequence

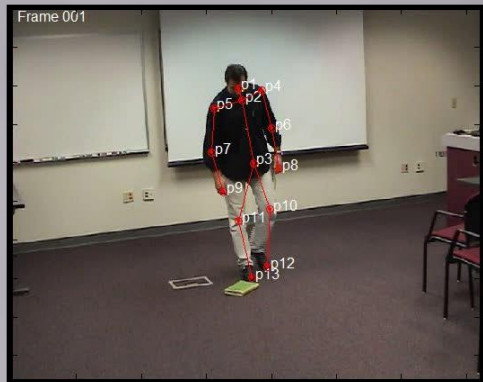
Action Detection

Analyzing Periodicity

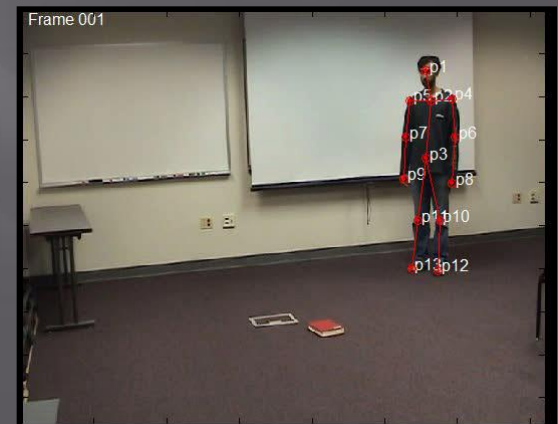
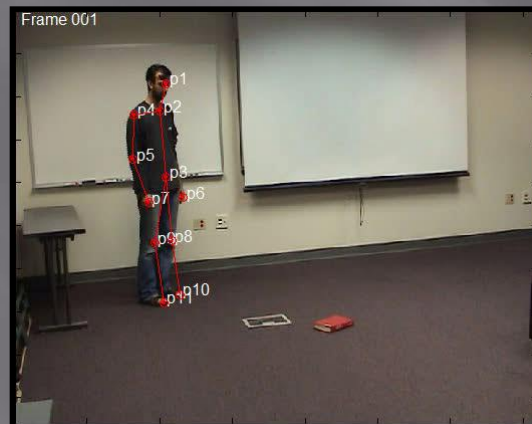
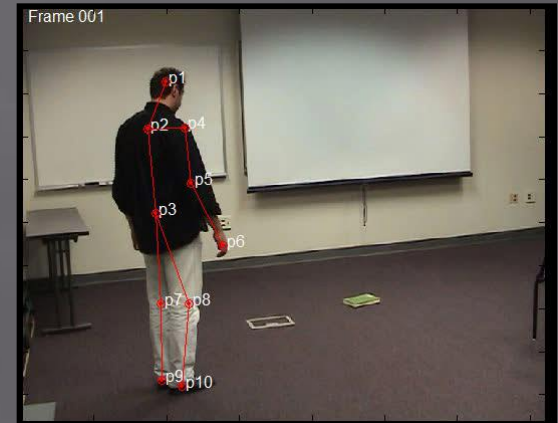
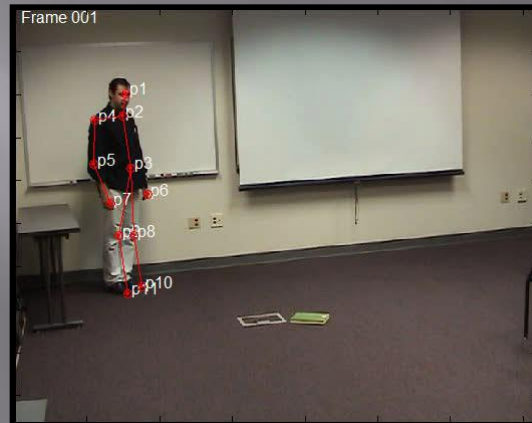


Action Detection:

Different approaches, different people, the same action



ReferencePattern

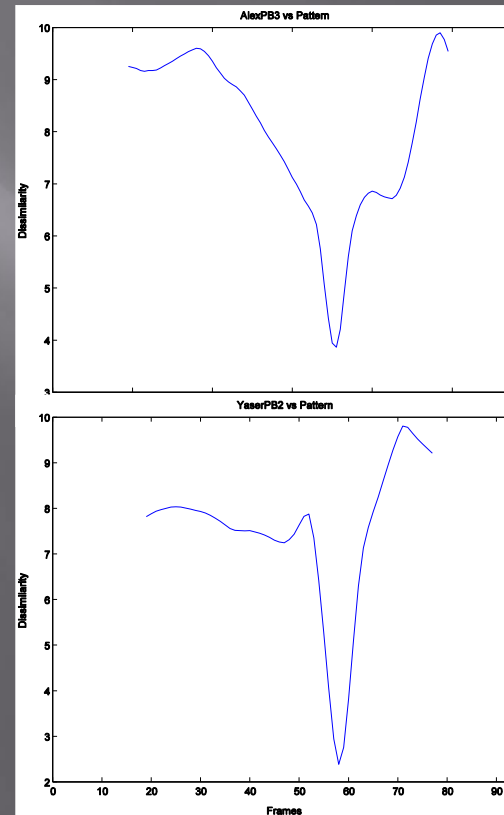
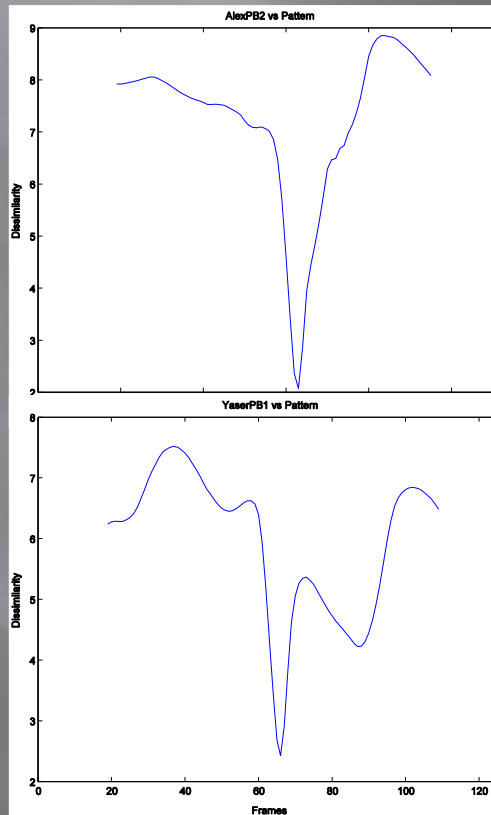


Test Sequences



Action Detection:

Different approaches, different people, the same action

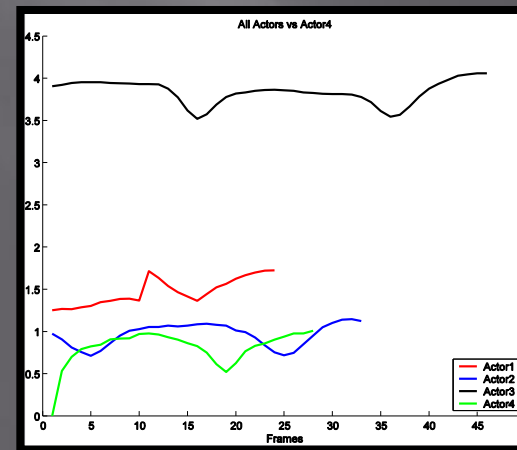
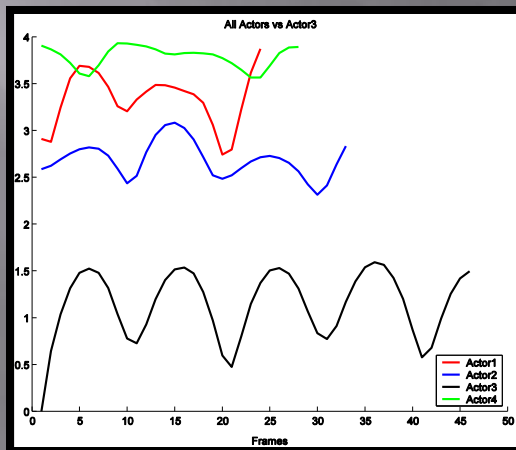
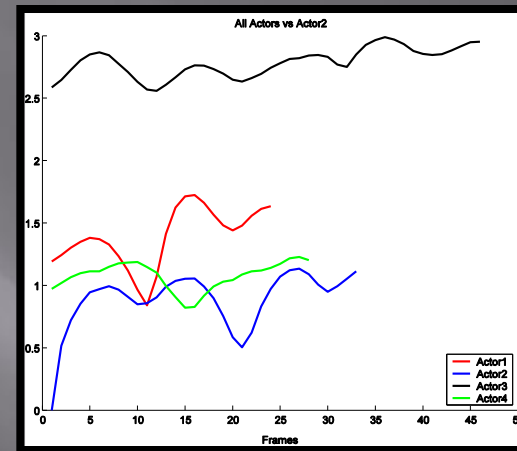
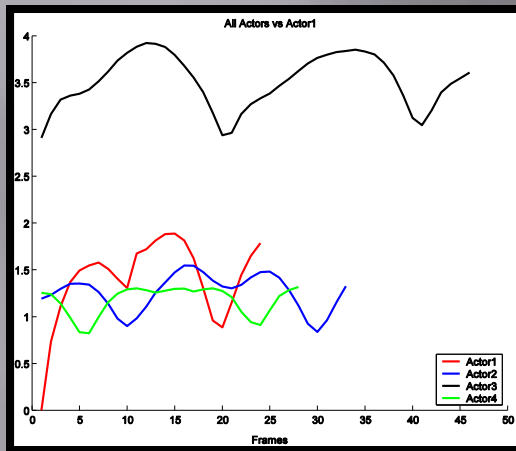


Analyzing Actions

Odd One Out

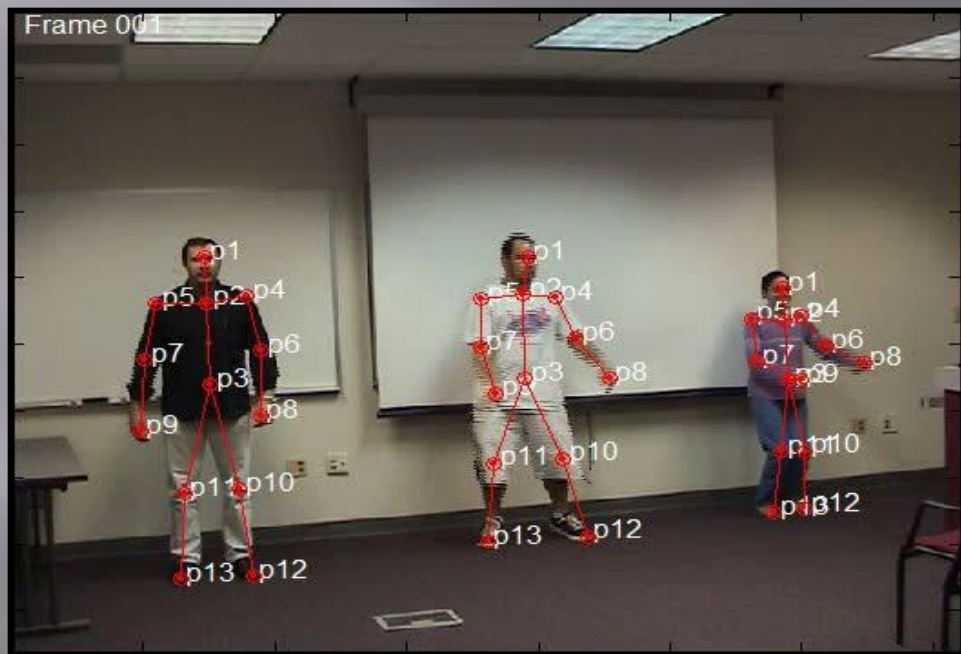


‘Odd One Out’



Action Synchronization

Following the Leader



Action Synchronization

Following the Leader



University of
Central Florida

VISION

Copyright Mubarak Shah, UCF

SPACE-TIME PROJECTION FOR UNIFORMLY MOVING CAMERAS

Yaser Sheikh, Alex Gritai and
Mubarak Shah
CVPR2007



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF



Space-time Projection Model

□ Notation

- Cartesian space-time world coordinate $\mathbf{X} = [T \ X \ Y \ Z]^T$
- Homogeneous space-time world coordinate $\mathbf{U} = [T \ X \ Y \ Z]^T$
- Inhomogeneous camera space-time coordinate $\mathbf{x} = [t \ u \ v]^T$
- Homogeneous camera space-time coordinate $\mathbf{u} = [t \ w \ u \ v]^T$

□ Projection

$$\mathbf{x} = \mathbf{K} \mathbf{u}$$

□ f is the focal length

□ α_t is the reciprocal of the frame-rate of the camera

□ (p_u, p_v) are the principal point offset



Fundamental Constraint Between Galilean Cameras

$$(u' \ v' \ t')^T \Gamma (u \ v \ t)^T = 0$$

$$\begin{pmatrix} u' \\ v' \\ t' \\ u \\ v \\ 1 \end{pmatrix}^T \begin{pmatrix} 0 & 0 & 0 & f_1 & f_2 & f_3 \\ 0 & 0 & 0 & f_4 & f_5 & f_6 \\ 0 & 0 & 0 & f_7 & f_8 & f_9 \\ f_{10} & f_{11} & f_{12} & f_{13} & f_{14} & f_{15} \\ f_{16} & f_{17} & f_{18} & f_{19} & f_{20} & f_{21} \\ f_{22} & f_{23} & f_{24} & f_{25} & f_{26} & f_{27} \end{pmatrix} \begin{pmatrix} u \\ v \\ t \\ u \\ v \\ 1 \end{pmatrix} = 0$$

$$\Gamma = \begin{pmatrix} 0 & \Delta F \\ \Delta F^T & F_{00} \end{pmatrix}$$



Epipolar Geometry

- Epipolar surface in one camera corresponding to a point in the other camera is defined by the

$$(u_k, v_k) \in \mathcal{I}_k \Rightarrow \begin{bmatrix} u_k & v_k & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = 0$$

where $\mathbf{s} = [s_1, \dots, s_6]$ is computed from given

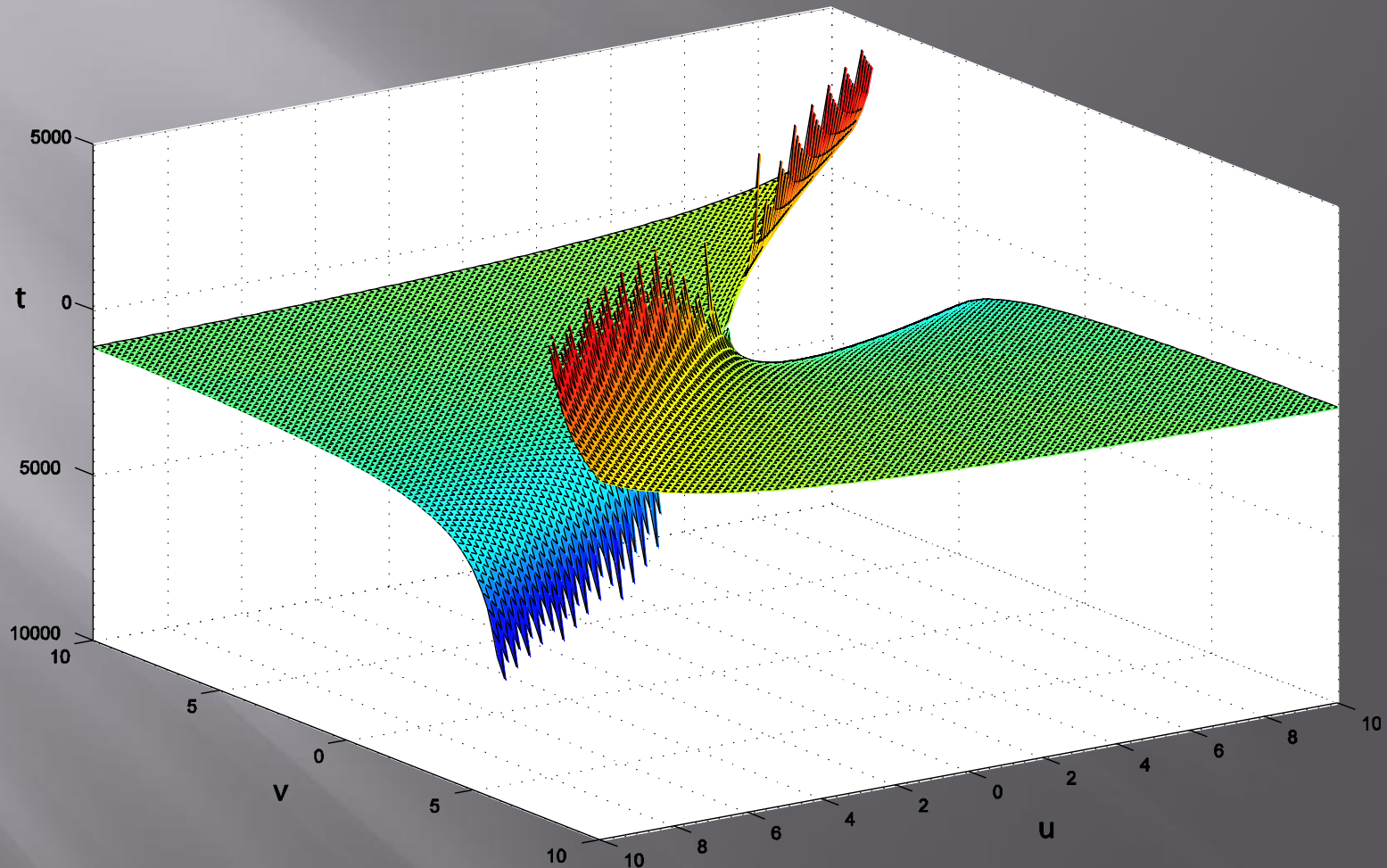
$$\mathbf{s} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \text{ as}$$

$$\begin{bmatrix} u_k & v_k & 1 \end{bmatrix}$$



Epipolar Surface

$$(0.049932u + 0.039102v + 0.012682) / (0.000034u + 0.000047v + 0.000011)$$

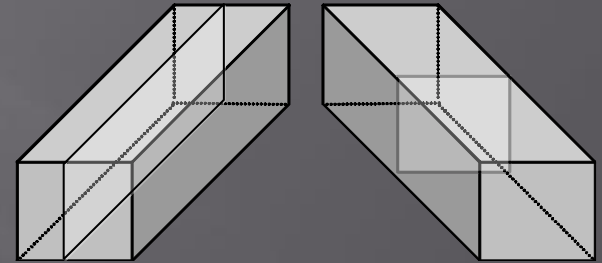
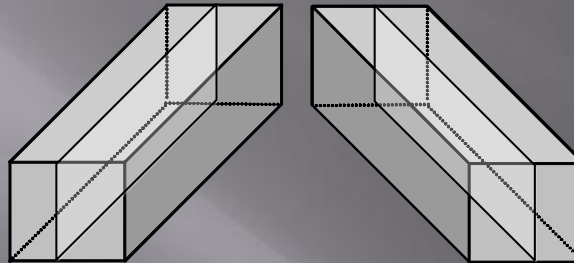
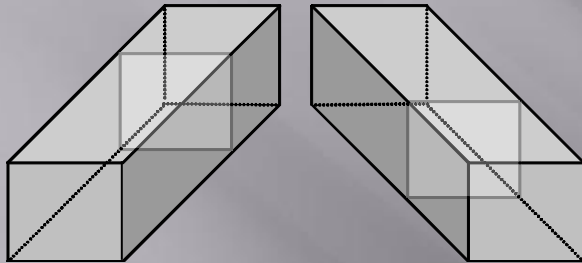


Specializations

Pushbroom and Perspective
images

Perspective images

Pushbroom images



Hartley and Faugeras

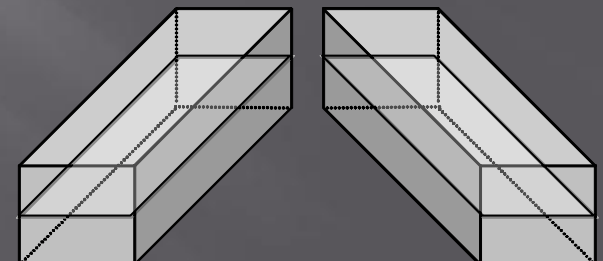
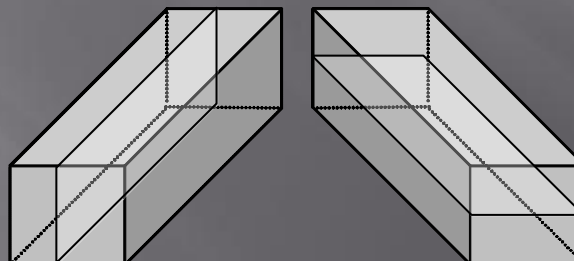
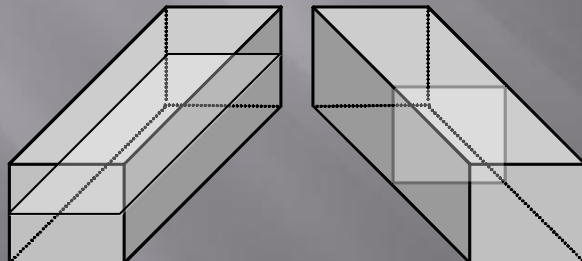
Gupta and Hartley

Khan, Rafi and Shah

EPI and Perspective
images

Pushbroom and EPI
images

EPI images



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF



APPLICATION

Action Recognition In Two Moving Cameras



Application to Action Recognition (two moving cameras)

- Corresponding points in two videos should satisfy:

$$(u'K' + t) = (uK + T)T$$

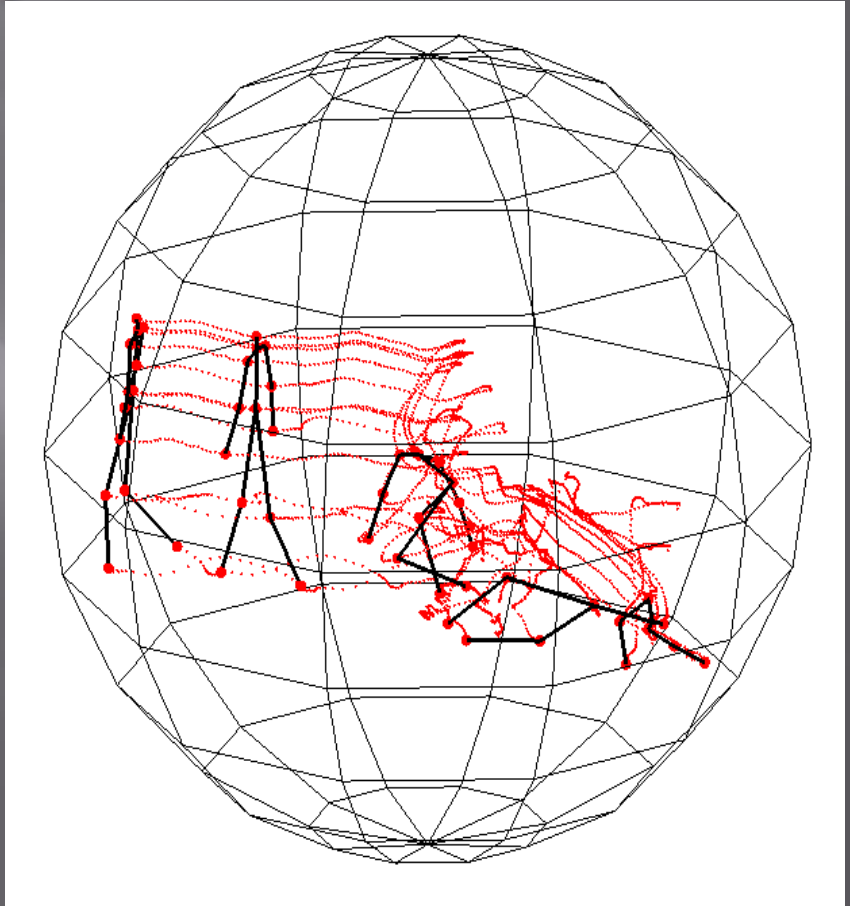
$$\begin{bmatrix} x'_1 & \dots & x'_n \\ y'_1 & \dots & y'_n \\ x'_1 & \dots & x'_n \\ y'_1 & \dots & y'_n \\ y'_1 & \dots & y'_n \\ y'_1 & \dots & y'_n \\ x_1 & \dots & x_n \\ y_1 & \dots & y_n \\ 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} E_{11} \\ E_{12} \\ E_{13} \\ E_{21} \\ E_{22} \\ E_{23} \\ E_{31} \\ E_{32} \\ E_{33} \end{bmatrix} = Af_T = C$$

$$K = \frac{\sigma_{27}}{\sigma_1}$$

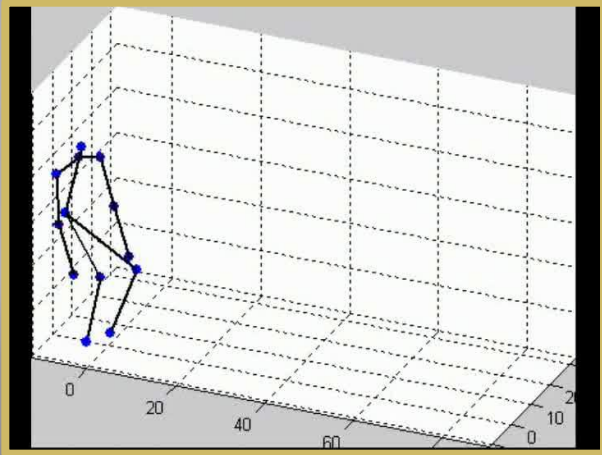


Experimental Results based on Motion Capture Data

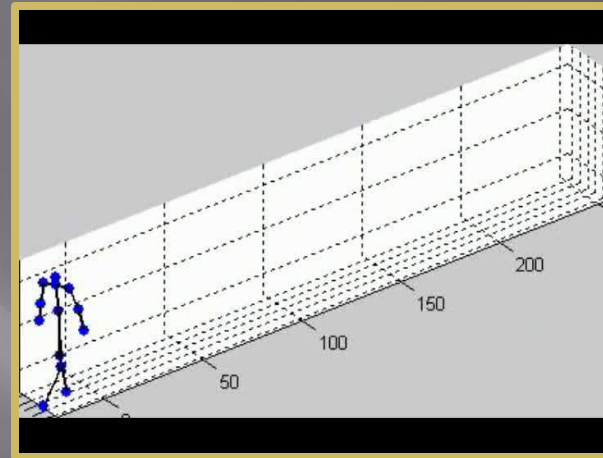
- ▣ Types of Transformations
 - Viewpoint
 - Anthropometric
 - Time
 - All three together



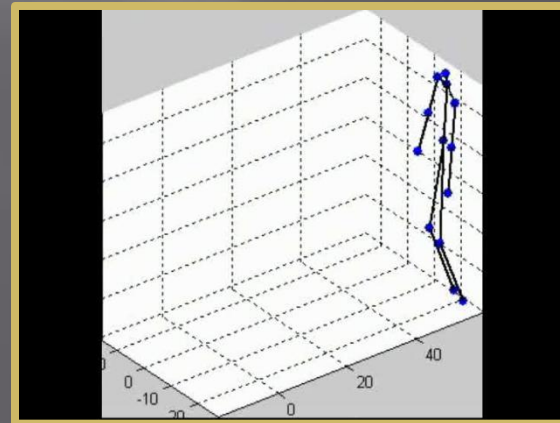
Robustness to View Point Transformations: Different Actions



Standing up



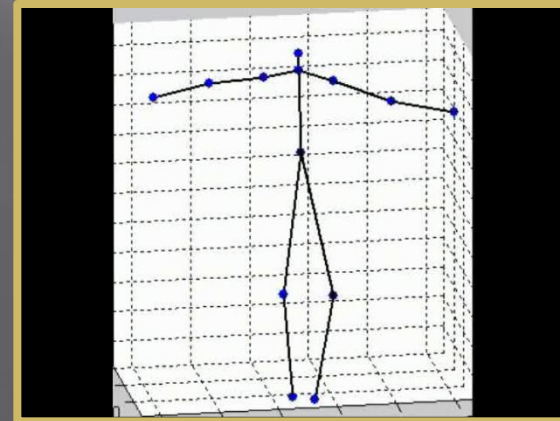
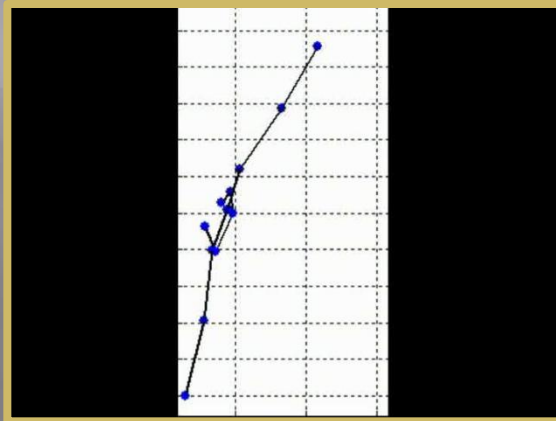
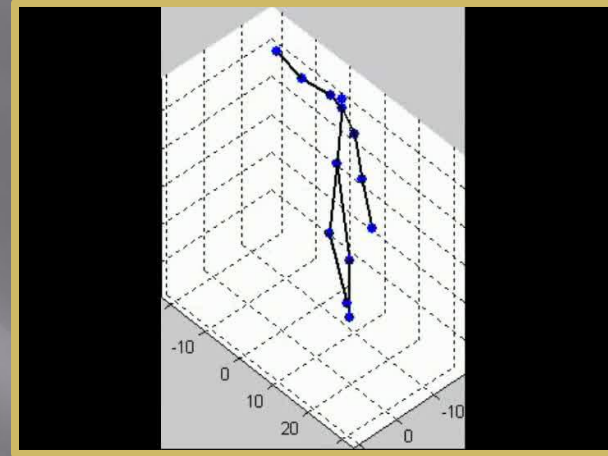
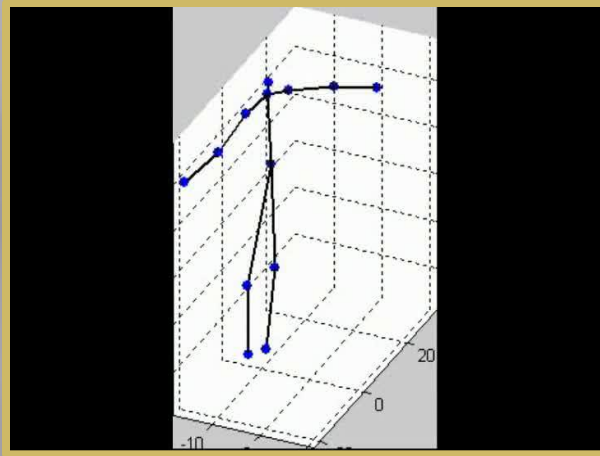
walking



sitting



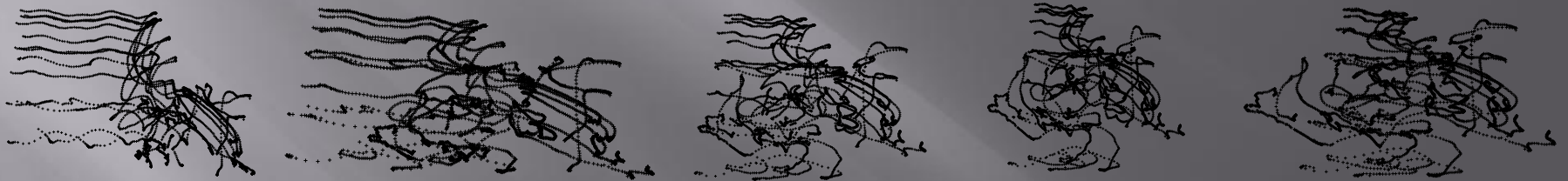
Robustness to View Point Transformations: Same Action



Human Actions Under Viewpoint, Anthropometric and Camera Motion Transformations (getting up)



Viewpoint



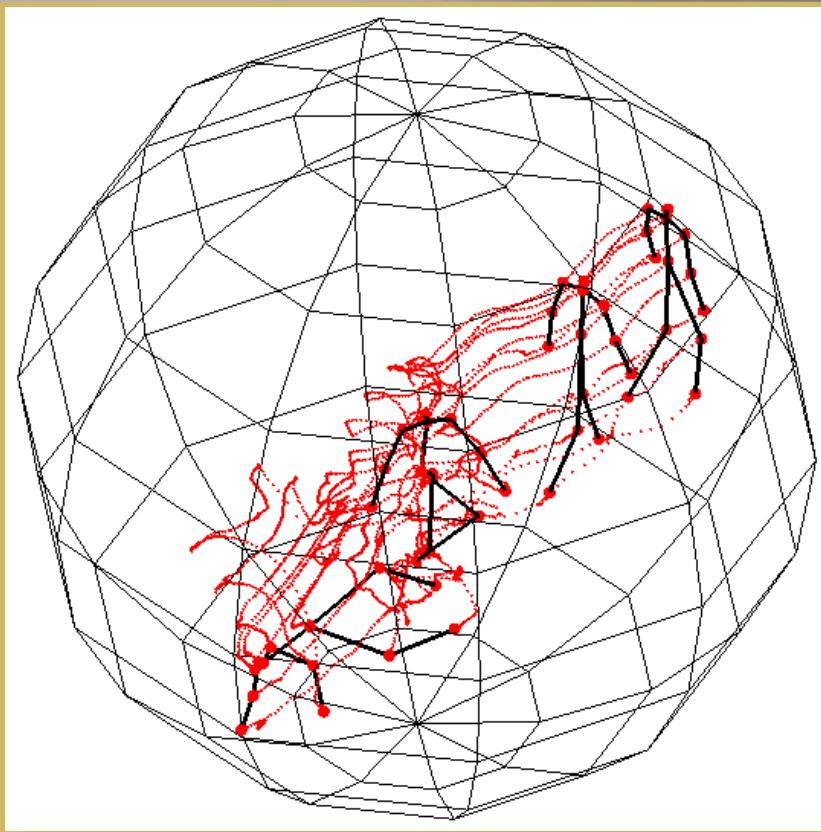
Anthropometric



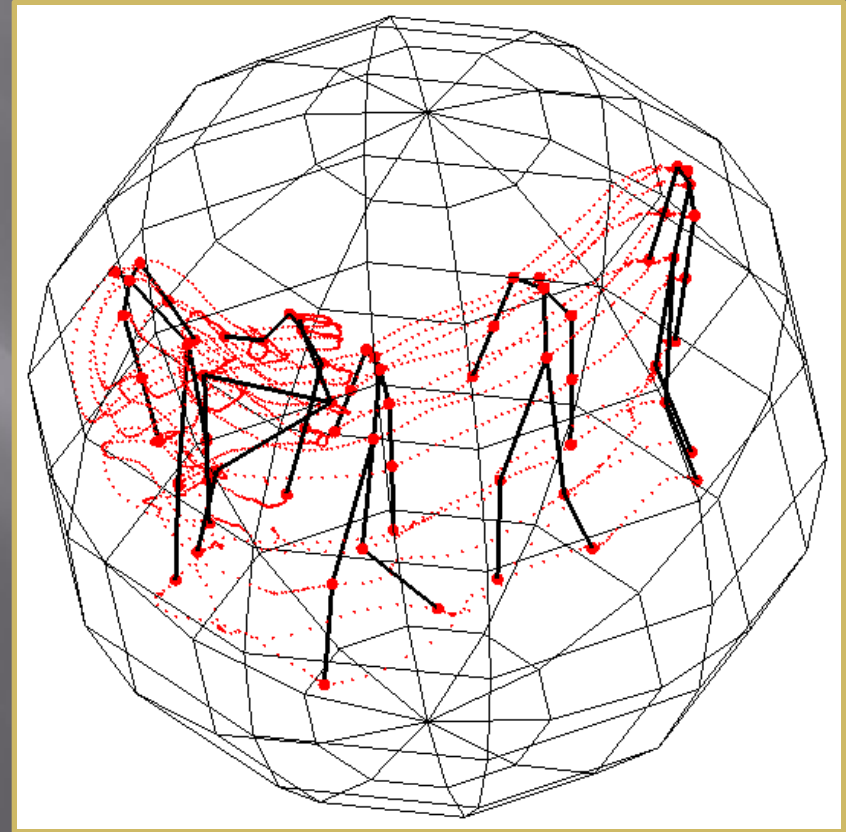
Camera Motion



Viewing Sphere (Azimuth 0-350, Elevation 0-90 (interval of 10))



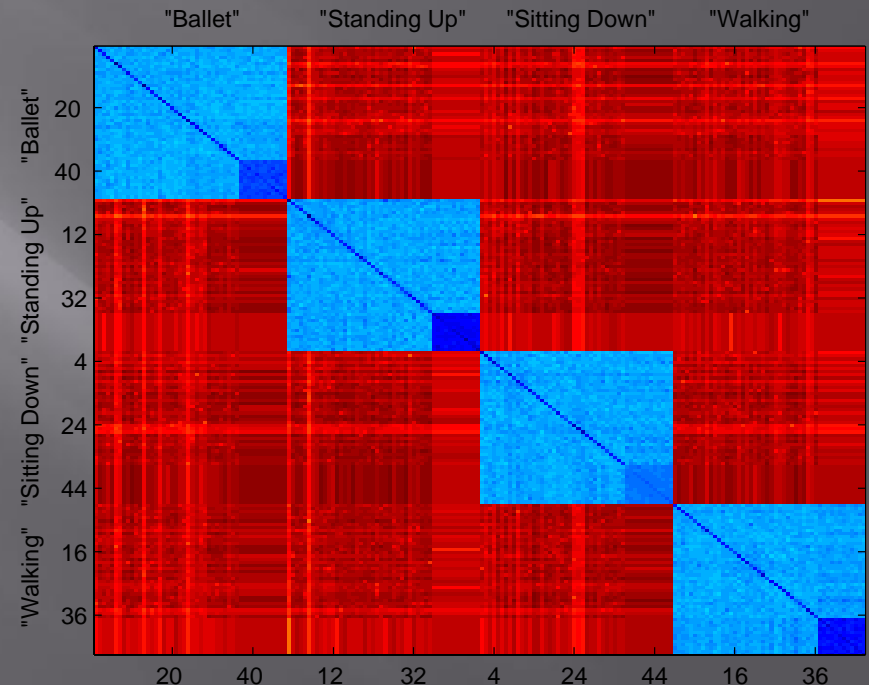
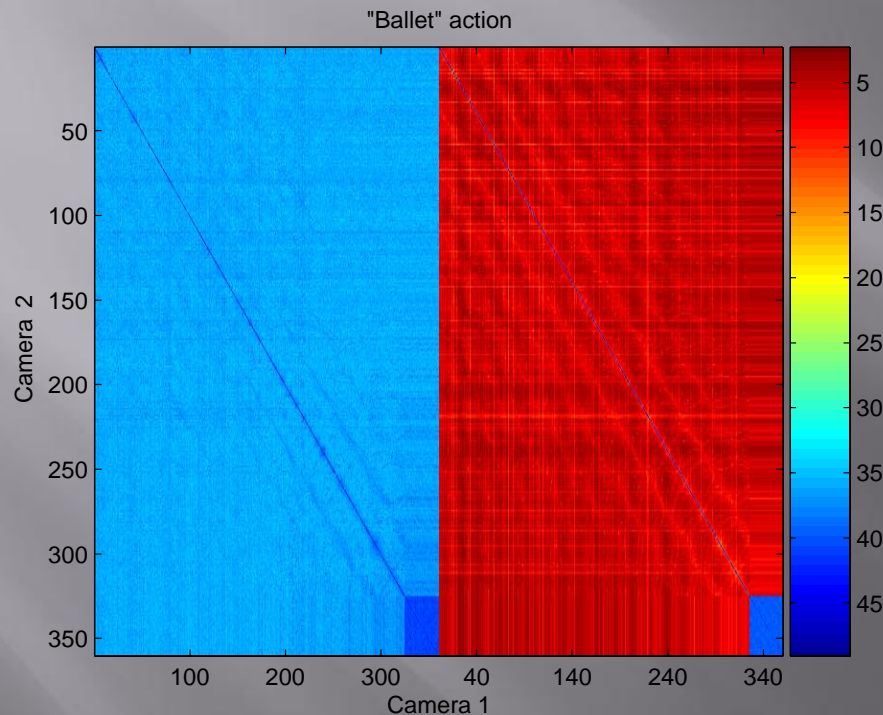
Getting Up



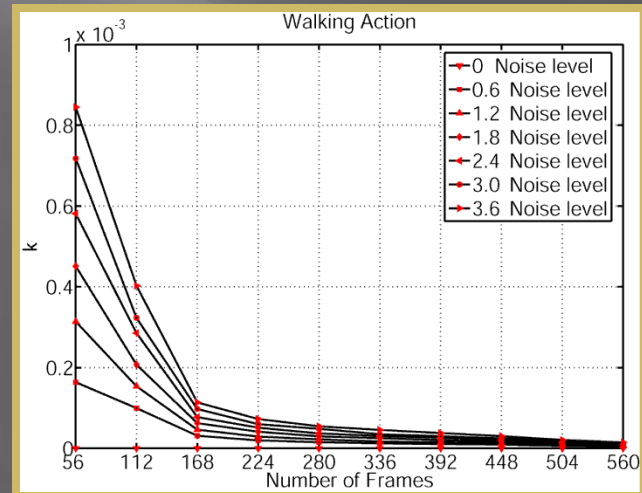
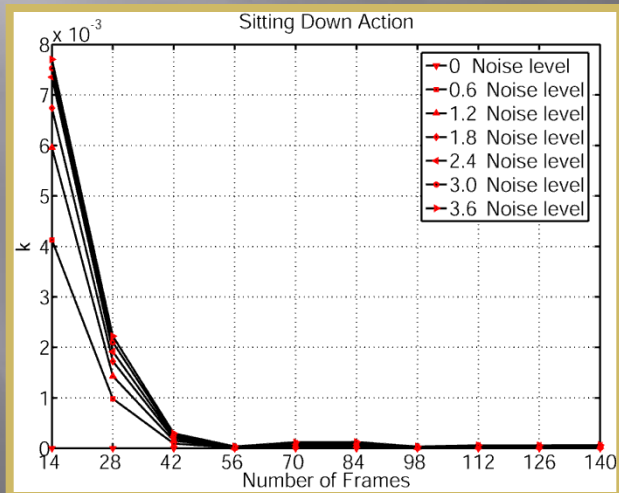
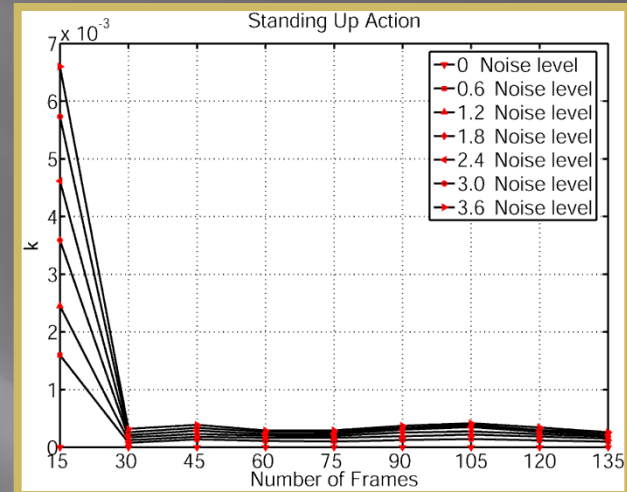
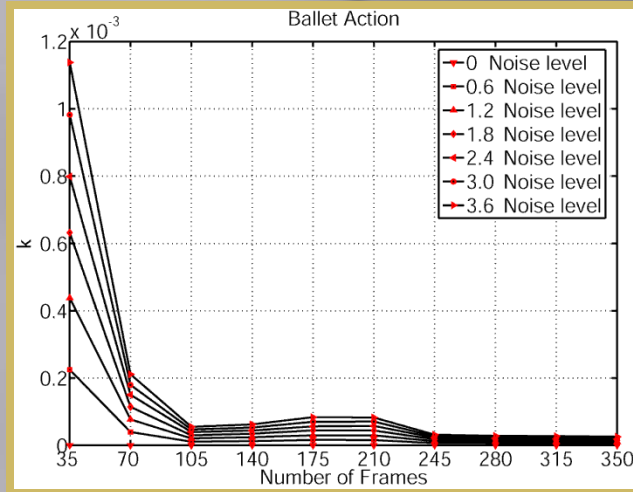
Sitting Down



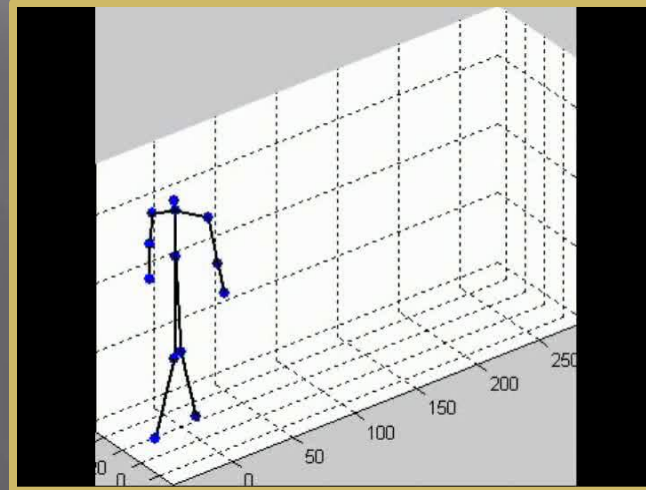
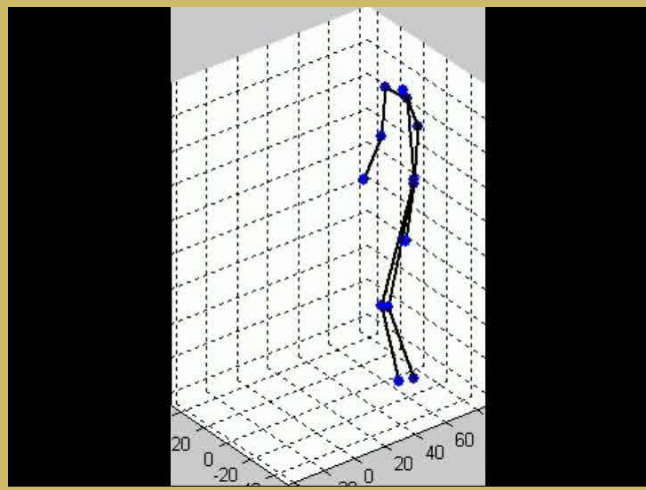
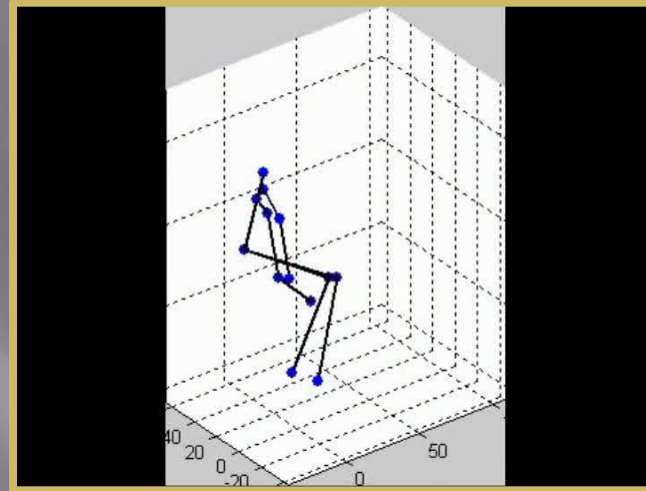
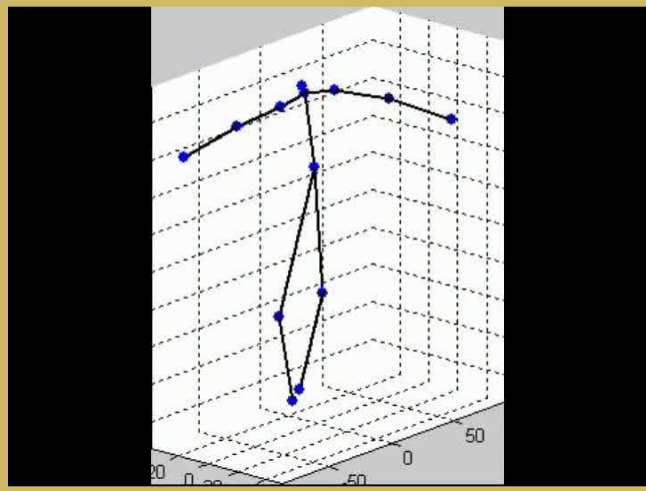
Robustness to View Point Transformations



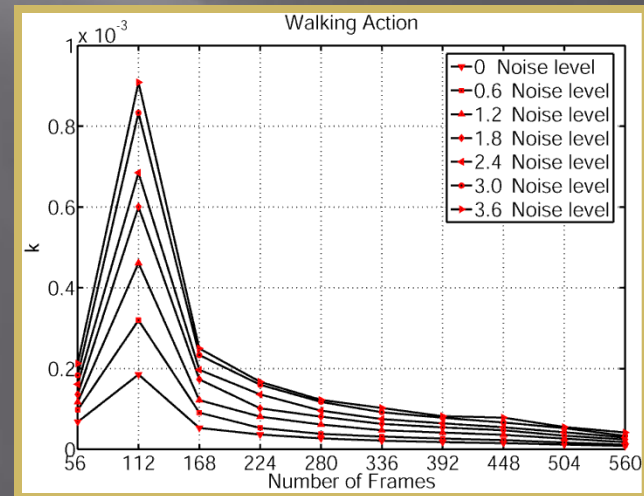
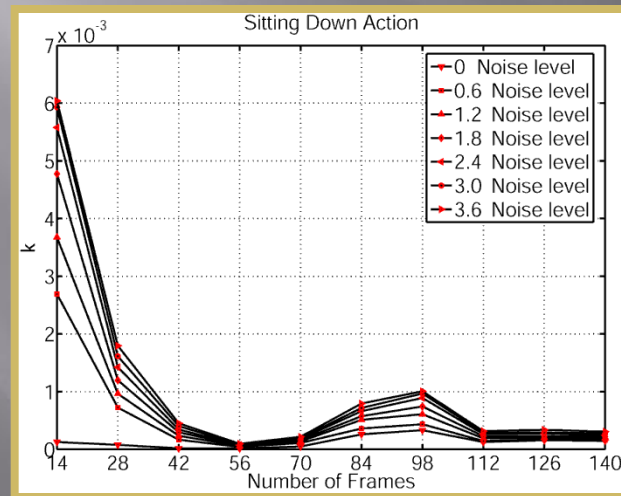
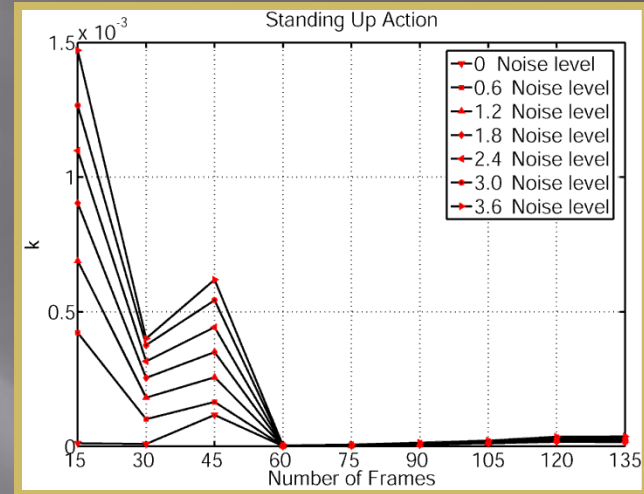
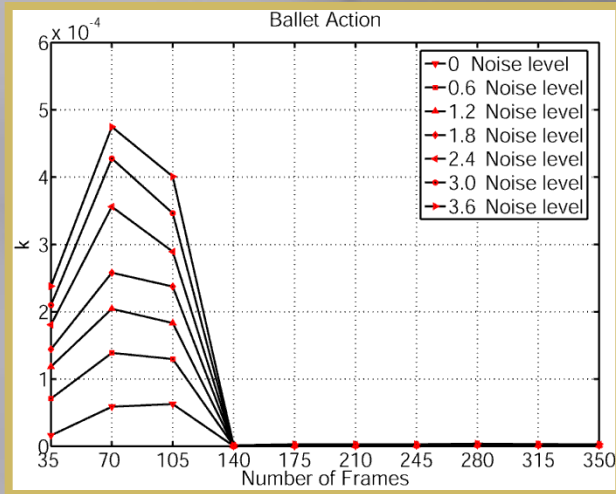
Viewpoint Transformations



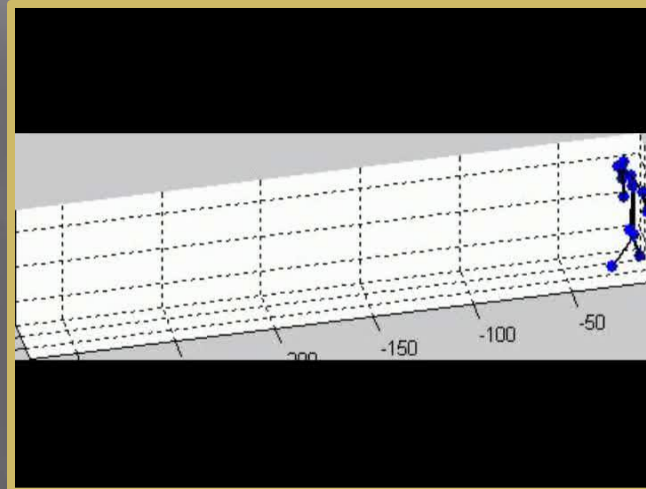
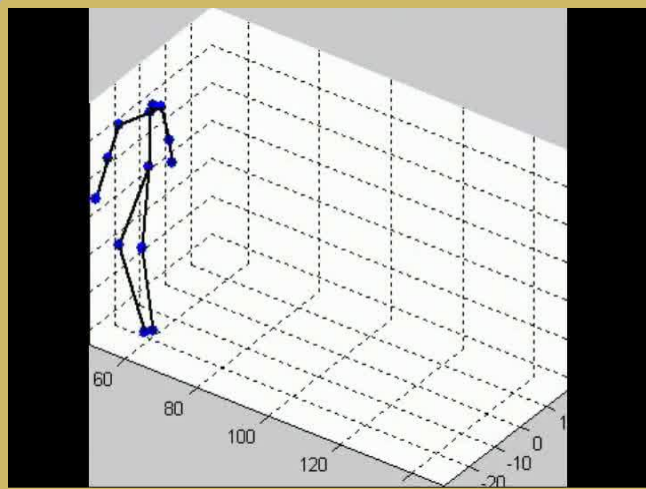
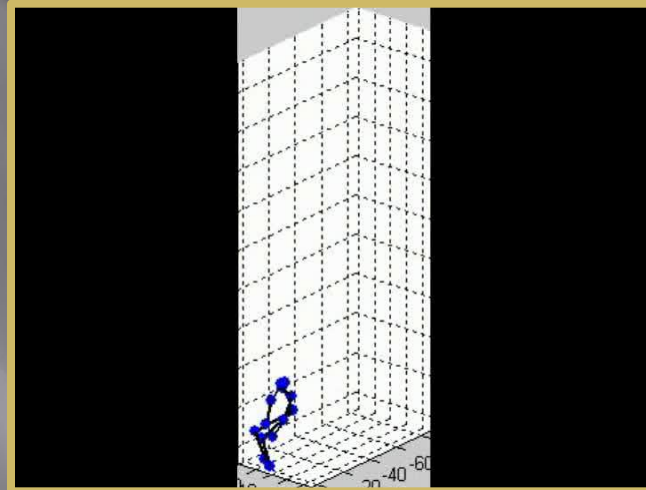
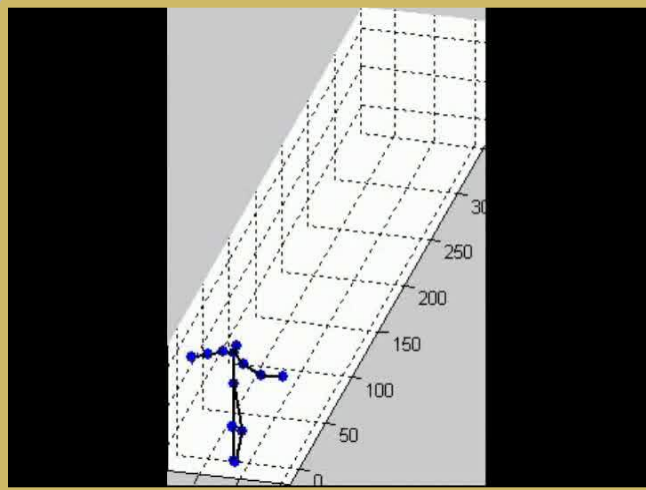
Anthropometric Transformations



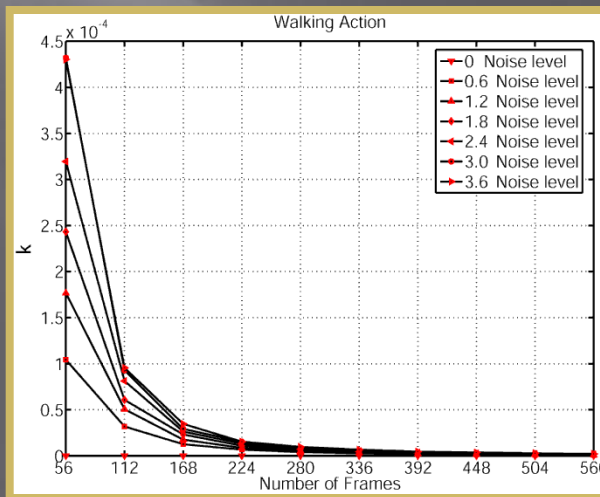
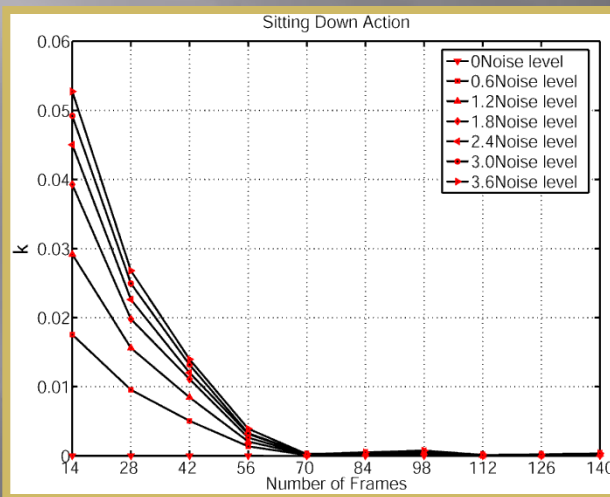
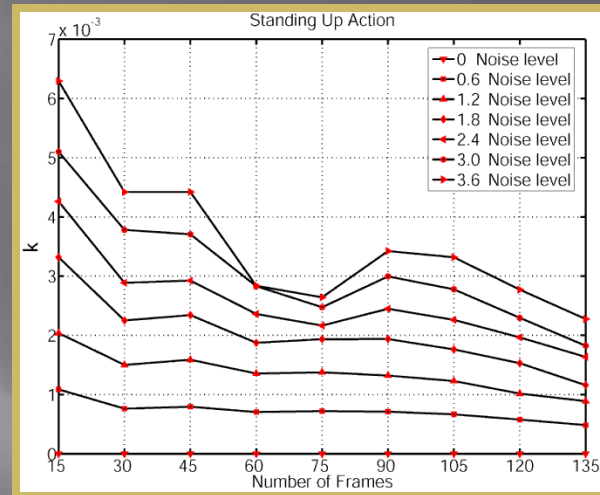
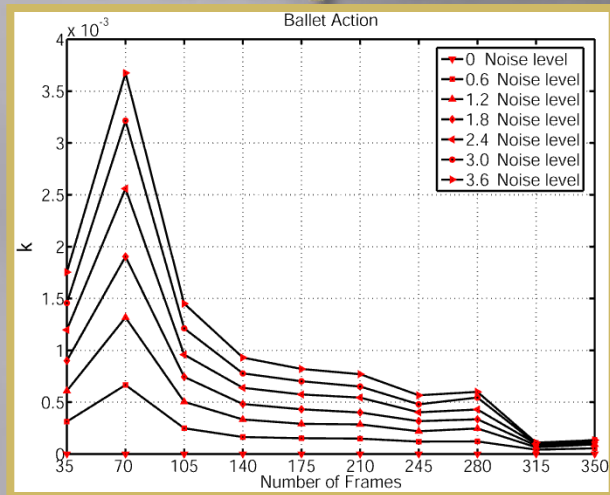
Anthropometric Transformations



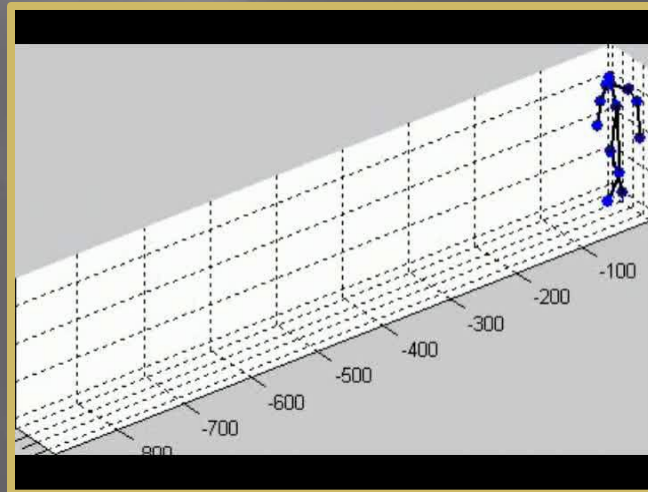
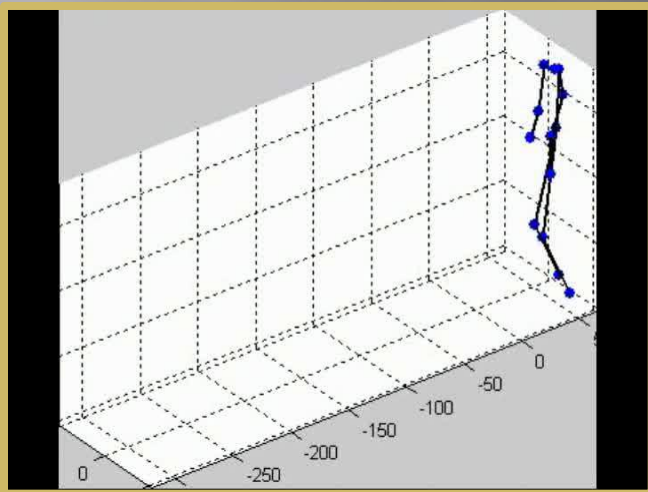
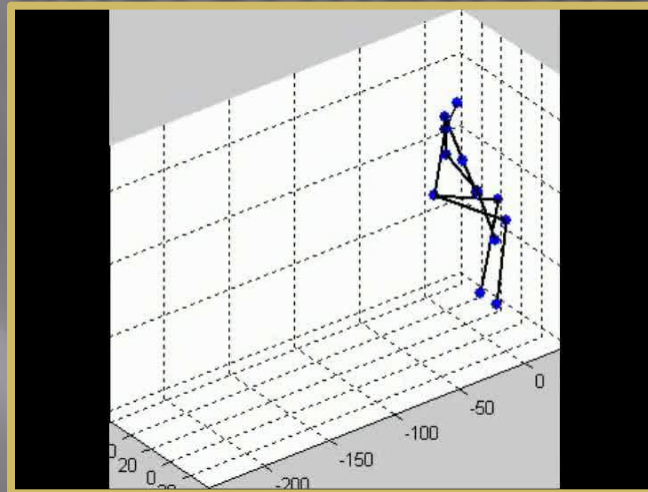
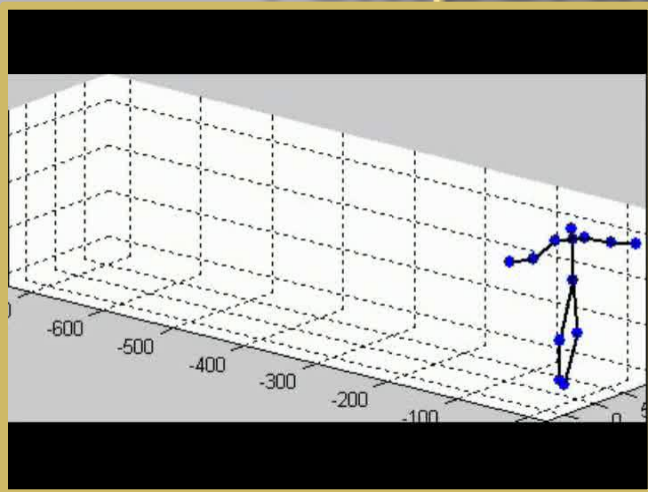
Temporal Transformations



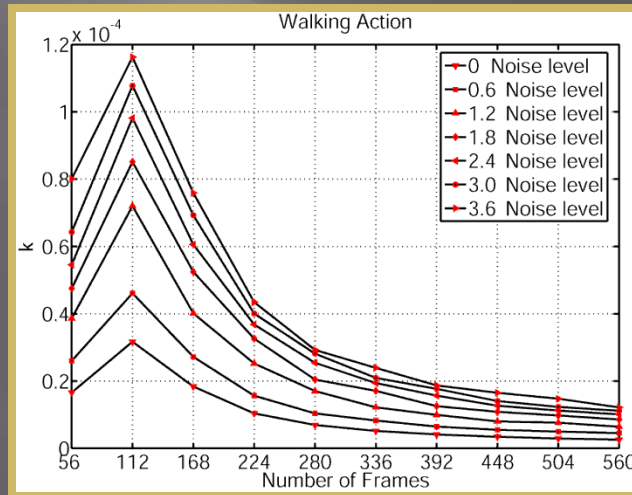
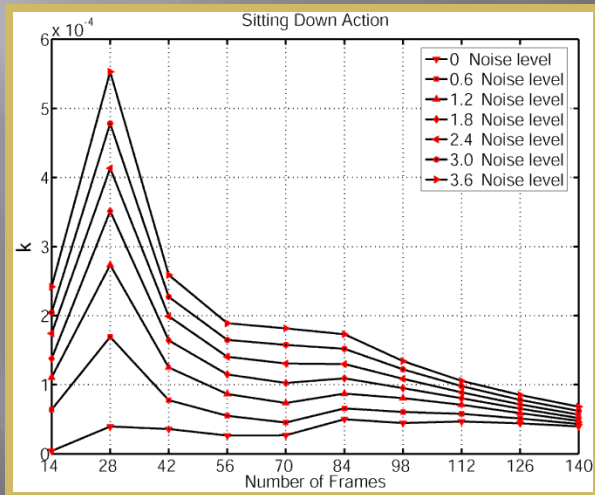
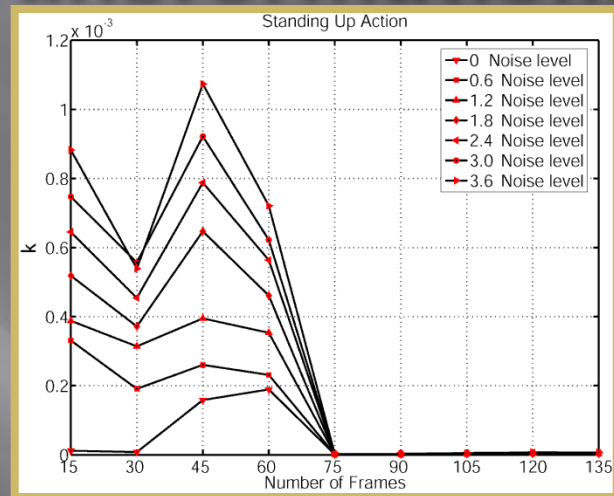
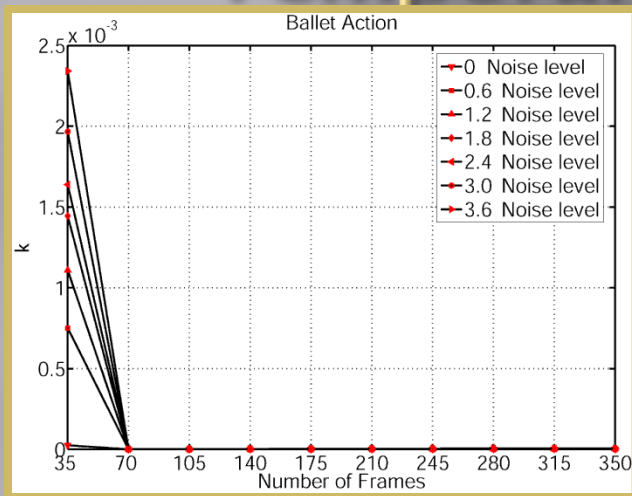
Temporal Transformations



Viewpoint, Anthropometric and Temporal Transformations



Viewpoint, Anthropometric and Temporal Transformations



CHAOTIC INVARIANTS FOR HUMAN ACTION RECOGNITION

Saad Ali, Arslan Basharat, Mubarak Shah
ICCV 2007



University of
Central Florida

VISION

Copyrights Mubarak Shah, UCF

Proposed Idea

$$f(\theta_1^t, \theta_2^t, \dots, \theta_N^t)$$

$$f(\theta_1^{t+1}, \theta_2^{t+1}, \dots, \theta_N^{t+1})$$

$$f(\theta_1^{t+2}, \theta_2^{t+2}, \dots, \theta_N^{t+2})$$

$$f(\theta_1^{t+3}, \theta_2^{t+3}, \dots, \theta_N^{t+3})$$



f

$$(\theta_1^t, \theta_2^t, \dots, \theta_N^t)$$



University of
Central Florida

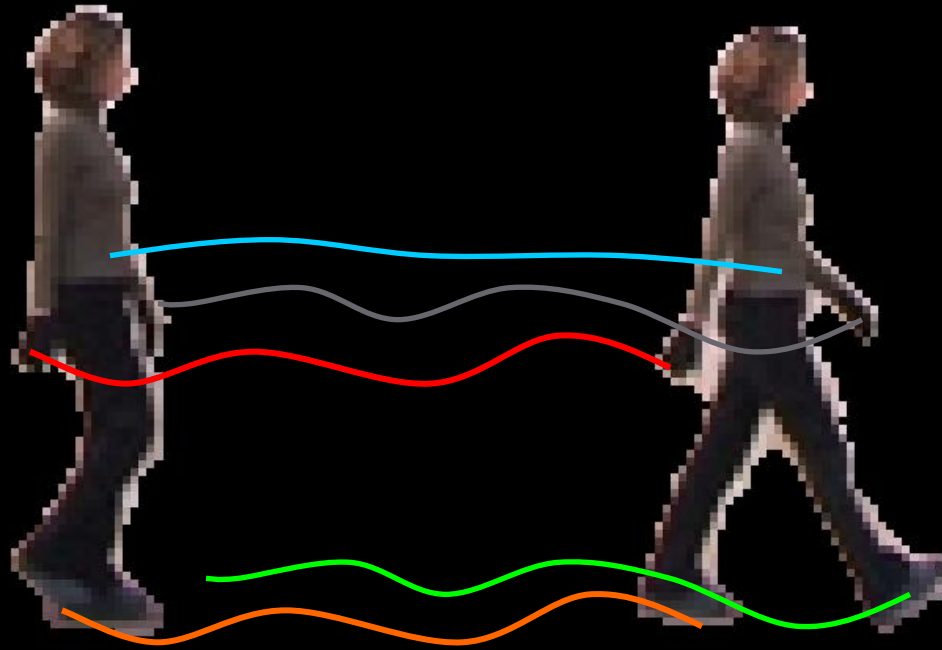
State

Space Variables

Copyrights Mubarak Shah, UCF

We have the access to the data generated by the dynamical system controlling this action !

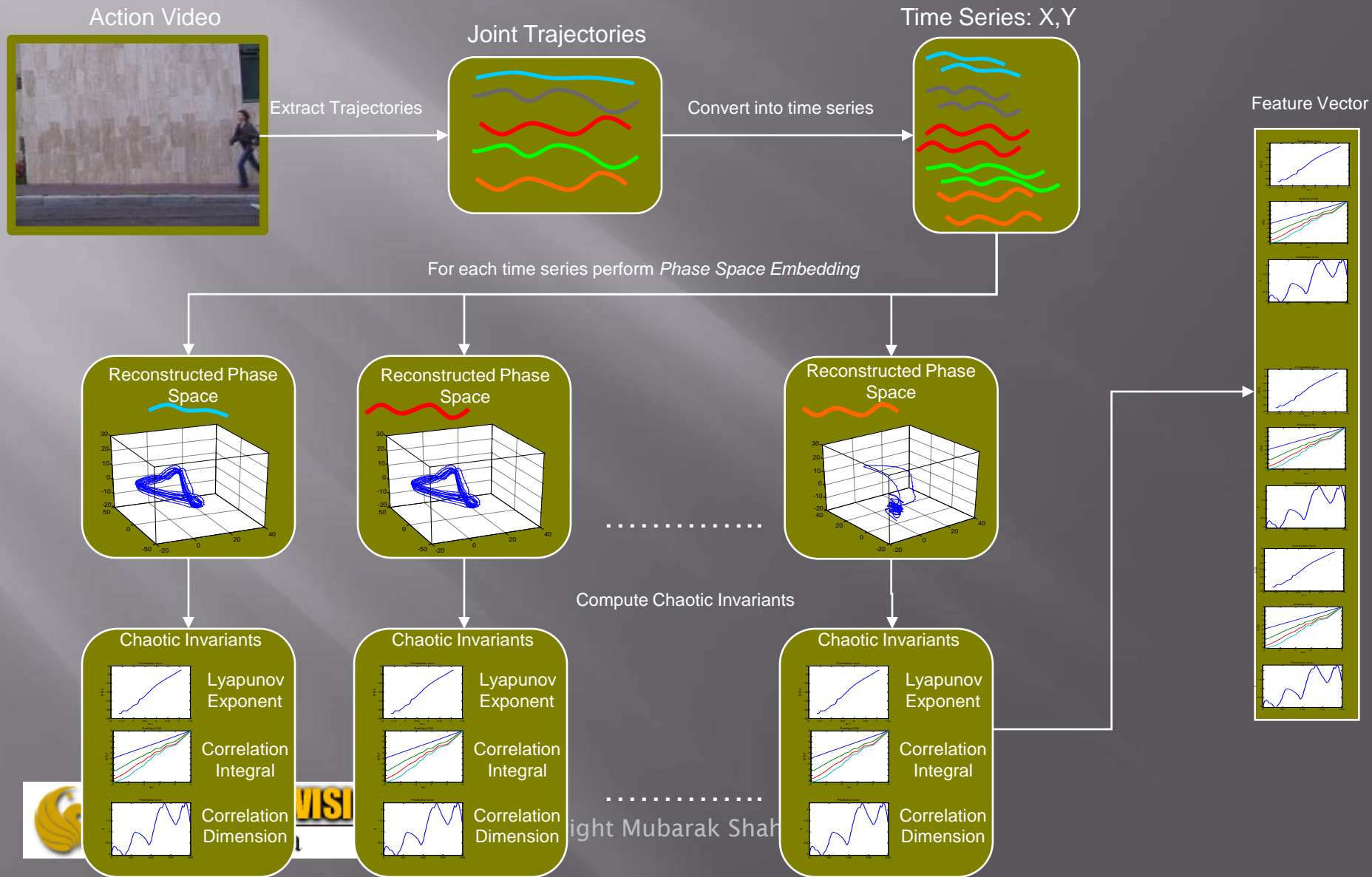
Proposed Idea



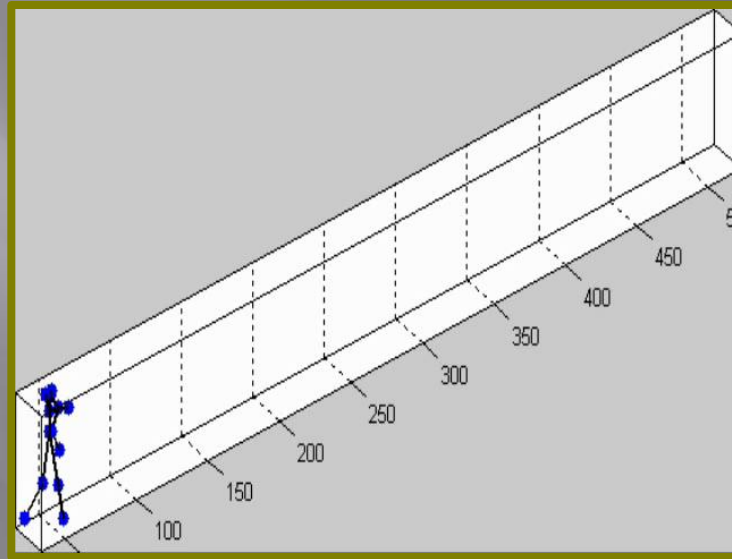
- ❑ Experimental Data: Trajectories of body joints
- ❑ From this data construct the phase space (and get periodic strange attractors) corresponding to the dynamical system responsible for generating the data.
- ❑ That is: Let the data speak to you and tell you what mechanisms are generating chaotic behavior.



Algorithmic Overview



Action Representation

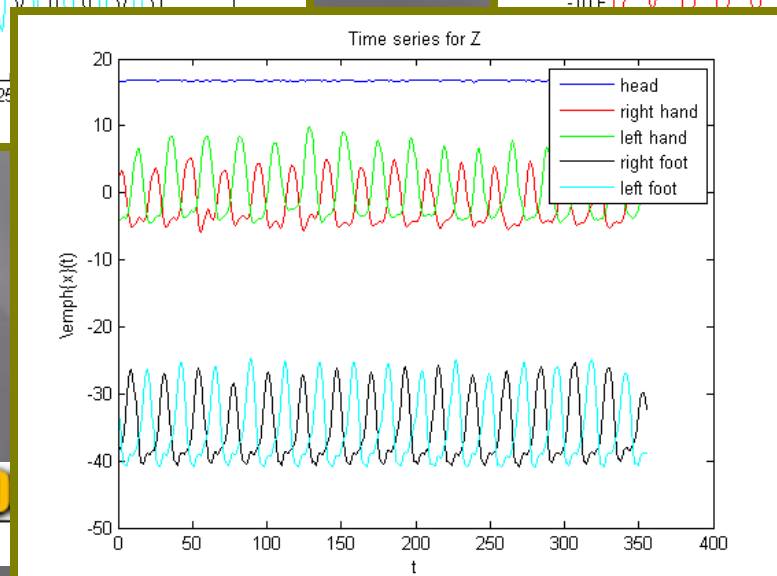
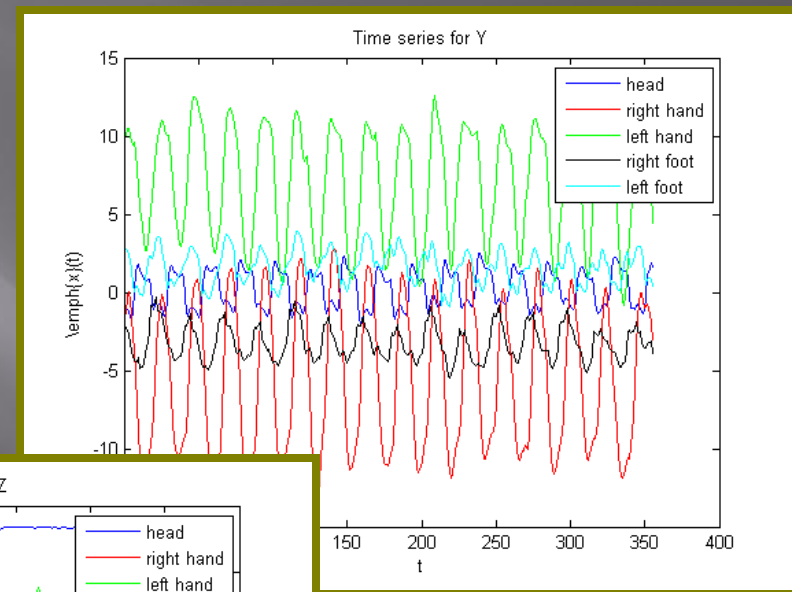
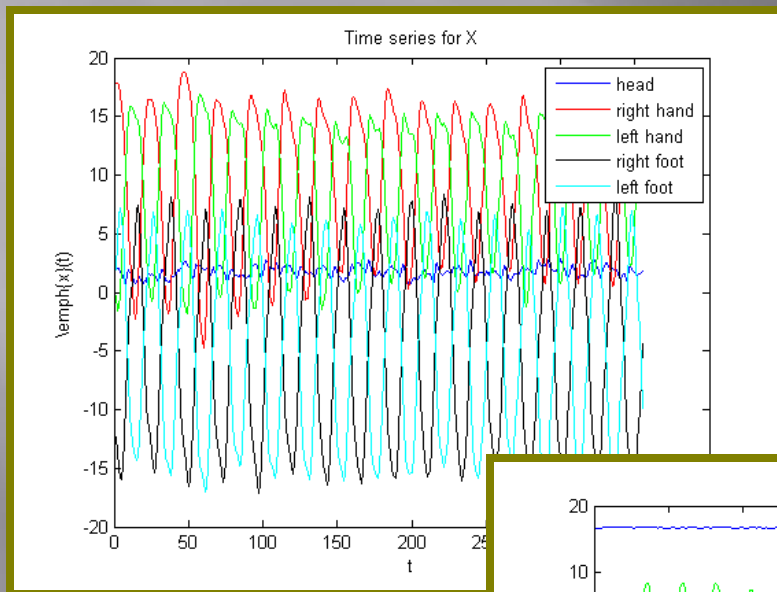


- Six Body Joints
 - Two Hands, Two Feet, Head, Belly.
- Normalized with respect to the belly point.
- Results in 5 trajectories per action.



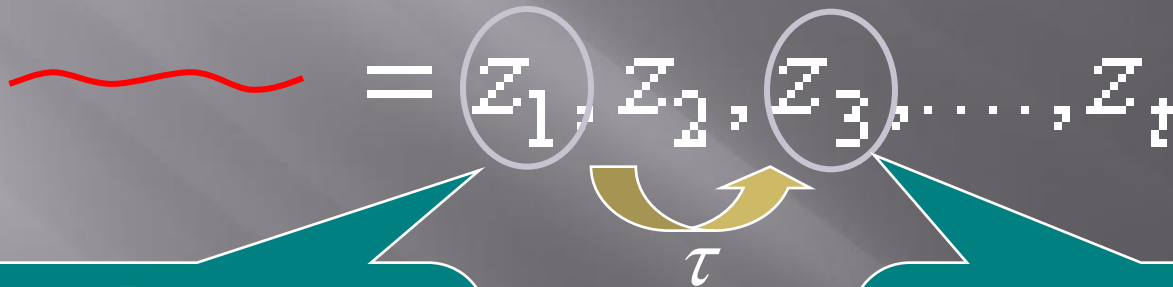
Action Representation

- Each dimension of the trajectory is considered as a univariate time series



Phase Space Embedding

- ▣ **Underlying Idea:** All the variables of the dynamical system influence each other.



Every point z_i of the series results from the intricate combination of influences of all the true state variables.

$$(\theta_1, \theta_2, \dots, \theta_N)$$

Therefore, $z_{i+\tau}$ can be considered as a second substitute variable which carries the influence of all the systems variables during time interval τ .

$$z_{i+2\tau}, z_{i+3\tau}, \dots, z_{i+m\tau}$$

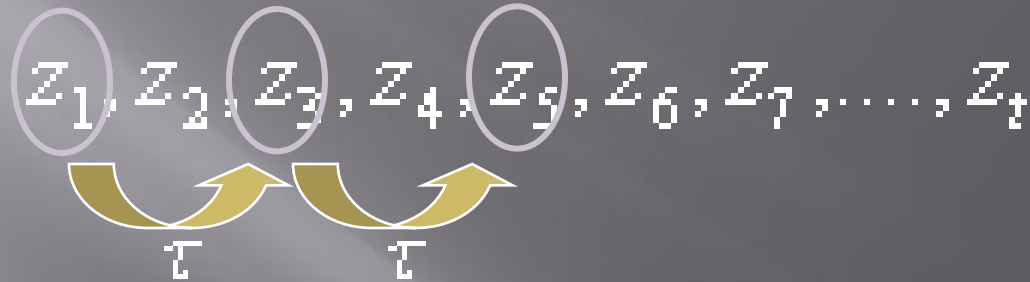
Using this reasoning, introduce a series of substitute variables and obtain the whole m -dimensional space.



Reconstructed Phase Space

$$m = 3$$

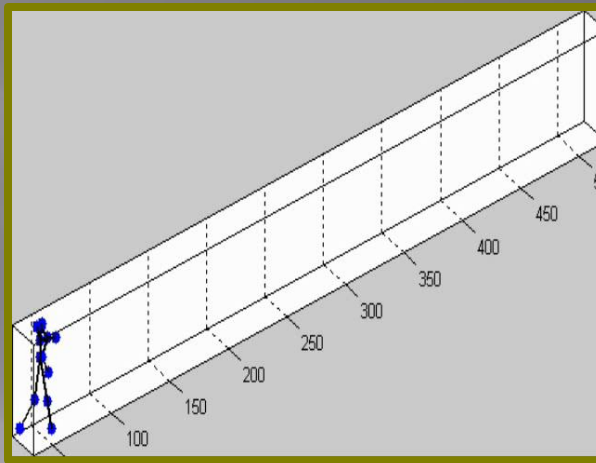
$$\tau = 2$$



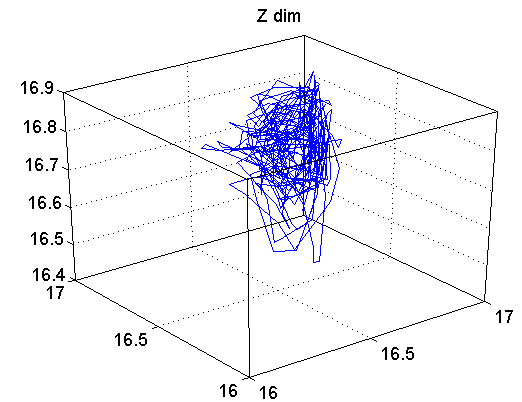
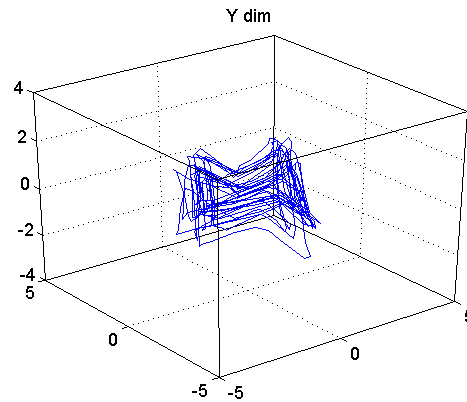
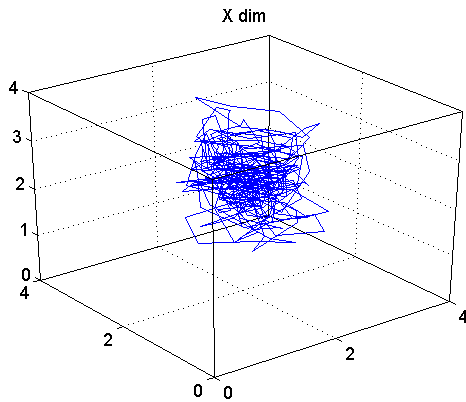
$$X = \begin{bmatrix} z_1 & z_3 & z_5 \\ z_2 & z_4 & z_6 \\ \vdots & \vdots & \vdots \end{bmatrix}$$

Each row is a point in a m-dimensional phase space.

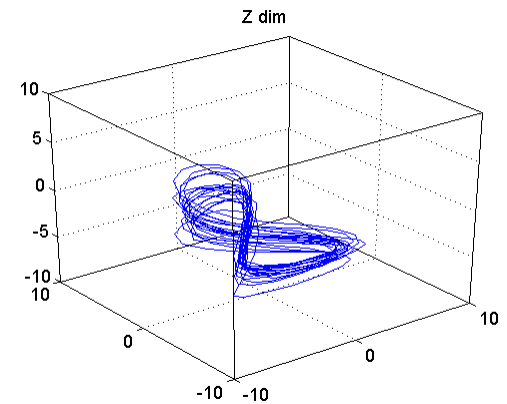
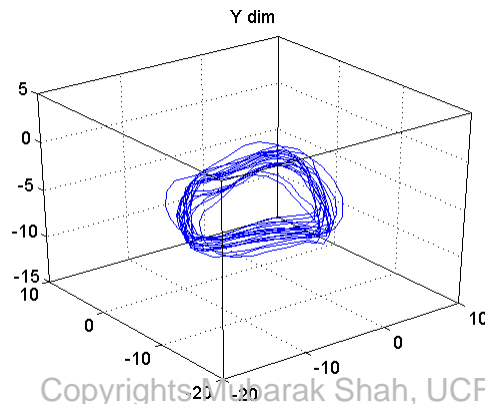
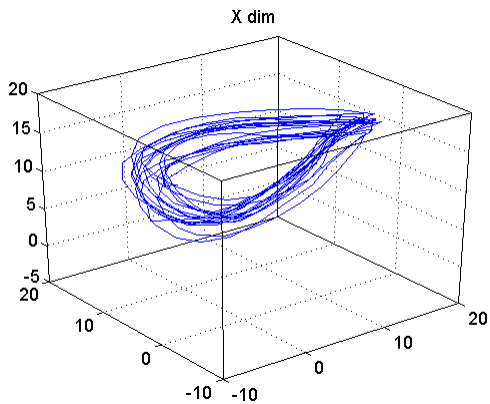
m-dimensional reconstructed phase space



Phase Spaces

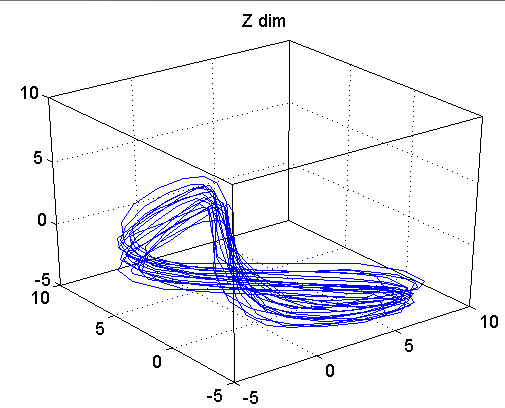
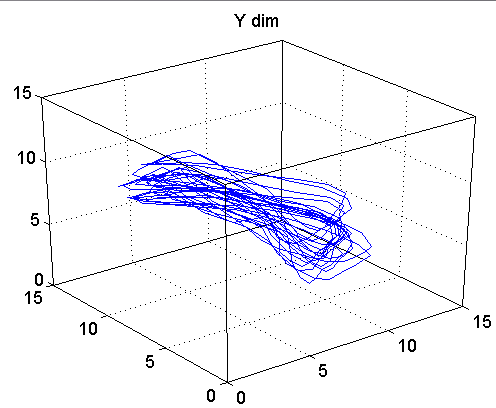
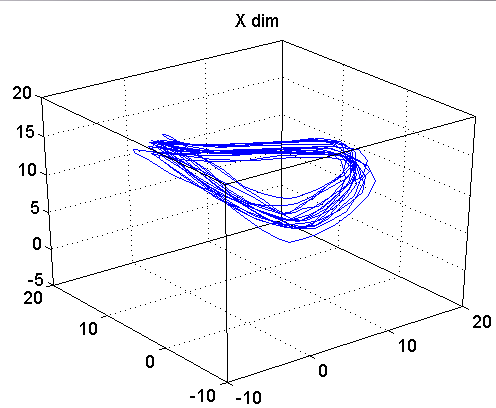


Head

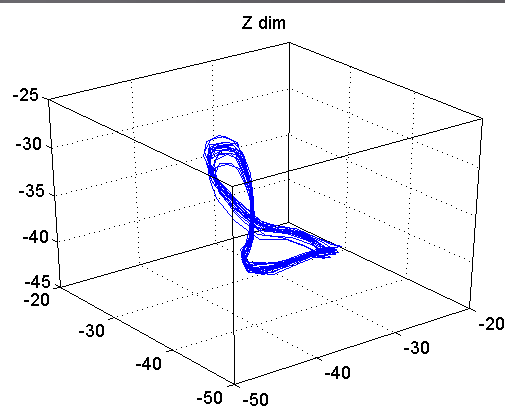
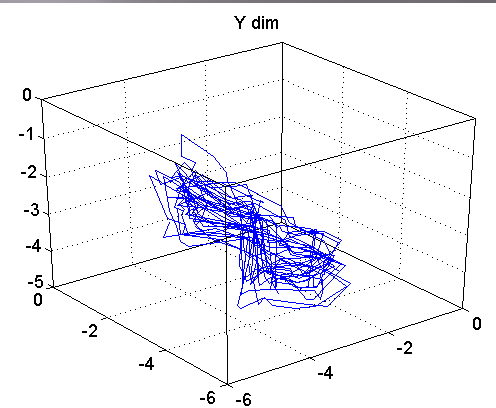
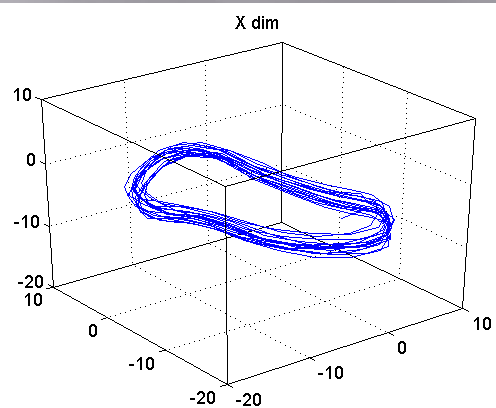


Right Hand

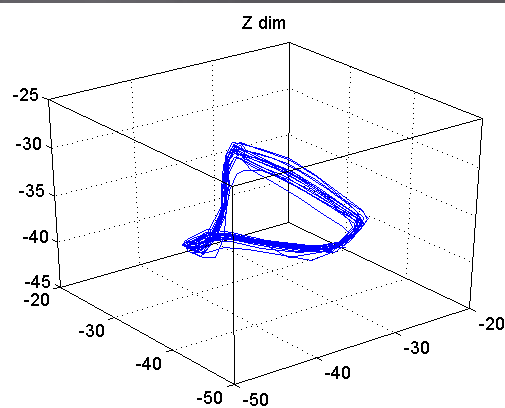
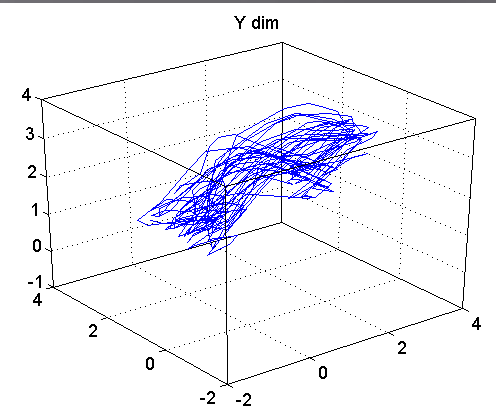
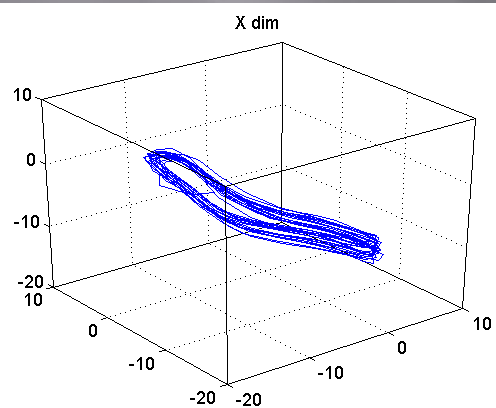
Left Hand



Right Foot



Left Foot



Invariant Features

- ▣ Maximal Lyapunov Exponent
- ▣ Correlation Integral
- ▣ Correlation Dimension



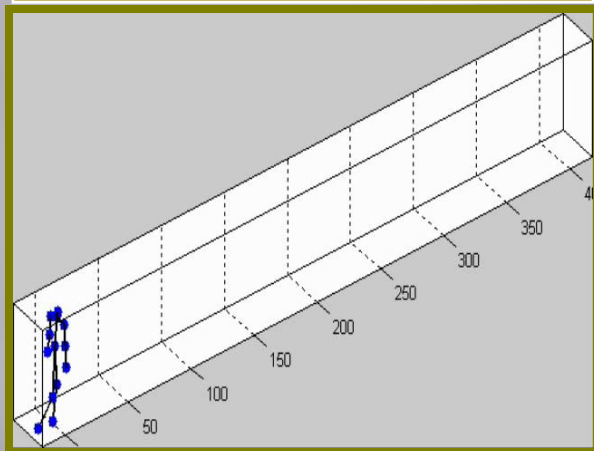
Experiment-I

- ▣ Motion capture data
- ▣ Dataset size
 - Dance : 19
 - Run : 26
 - Walk : 46
 - Sit: 14
 - Jump: 33
- ▣ Leave-One-Out Cross validation using K-means classifier.

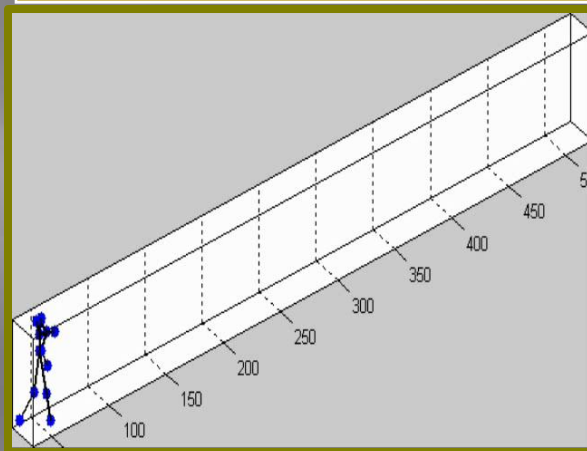


Experiments Motion Capture Data

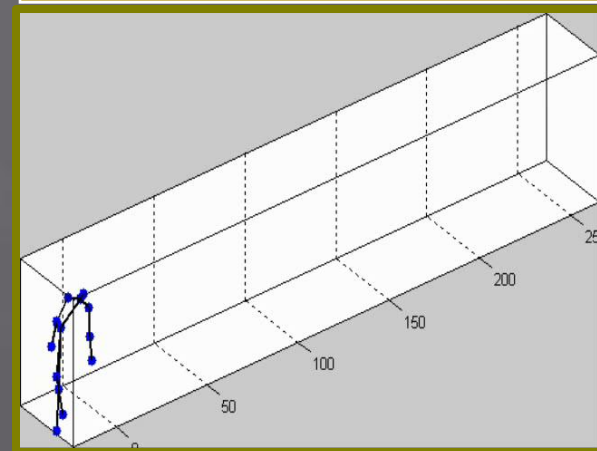
Walking



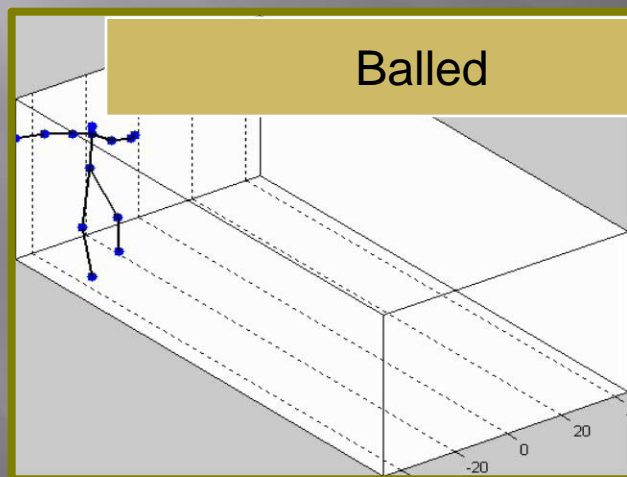
Running



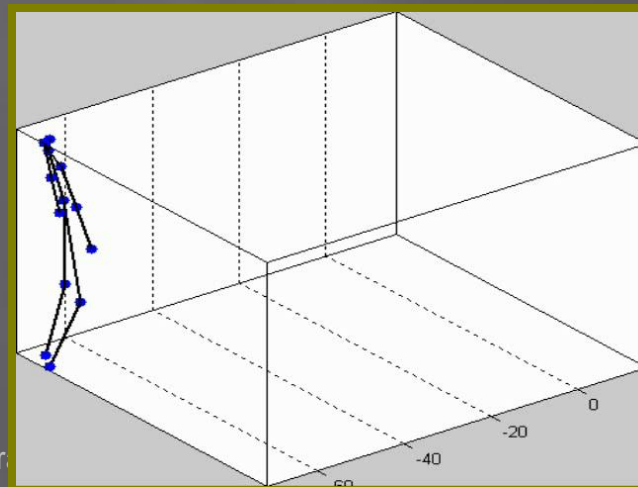
Jumping



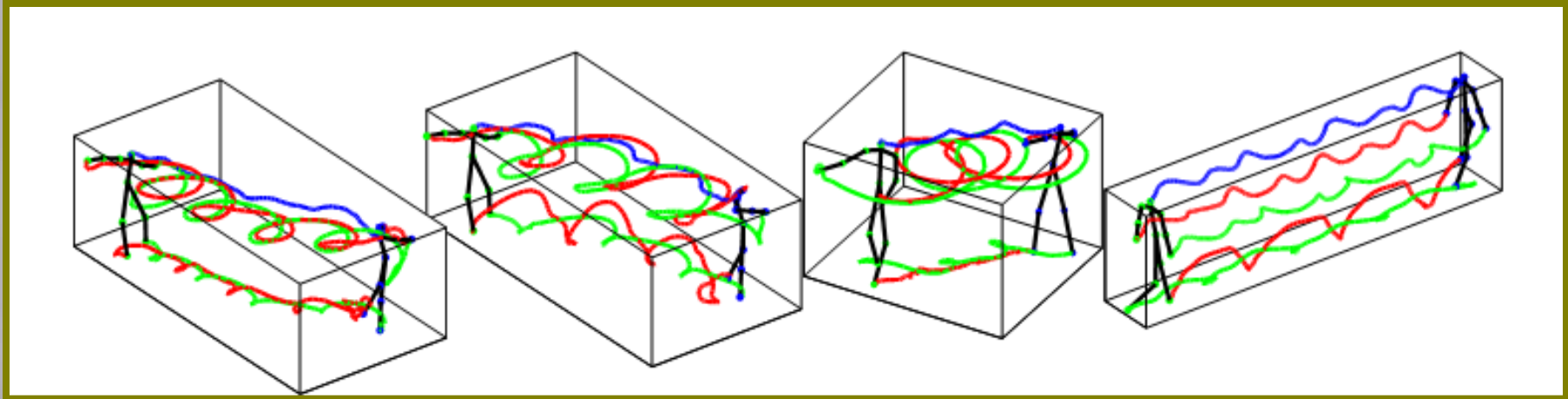
Balled



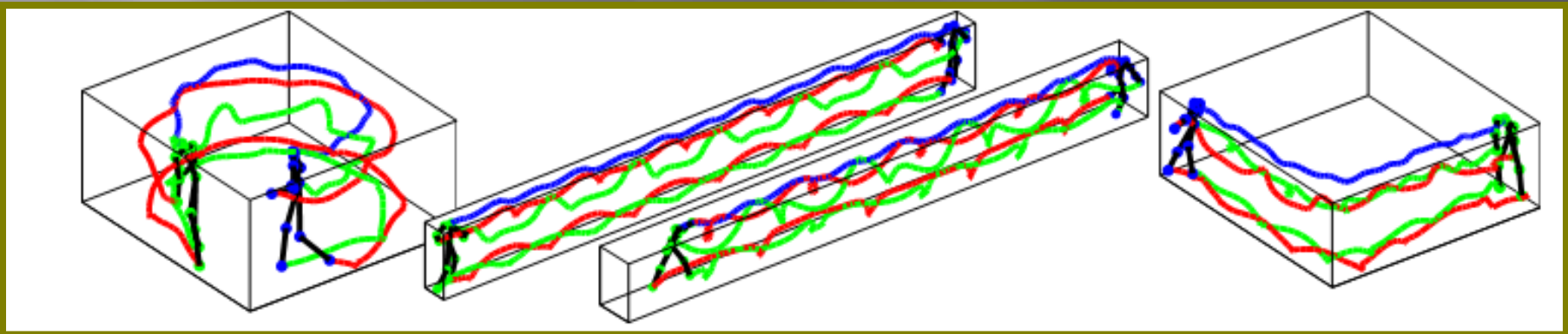
Sitting



Experiments Motion Capture Data

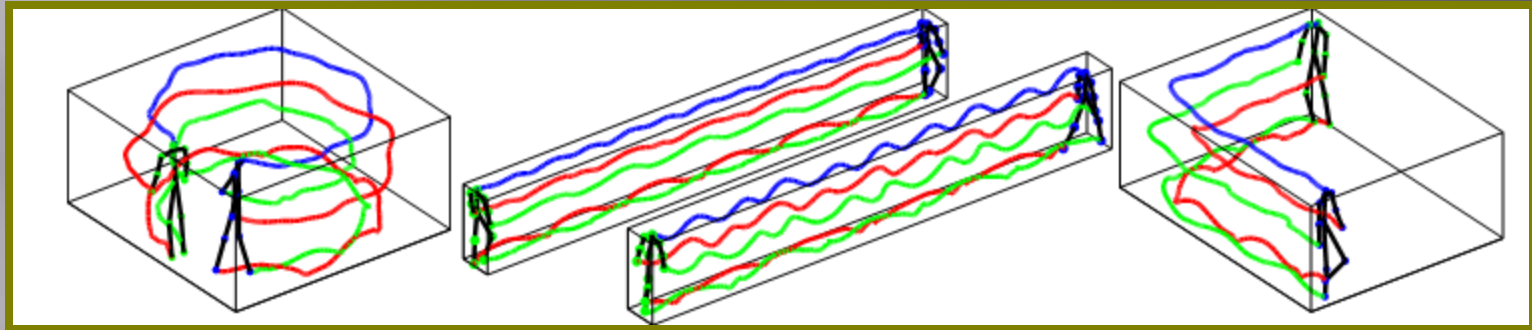


Dance

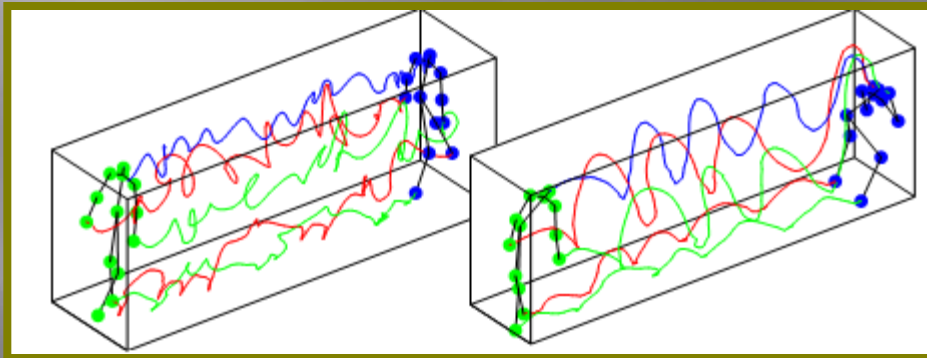


Run

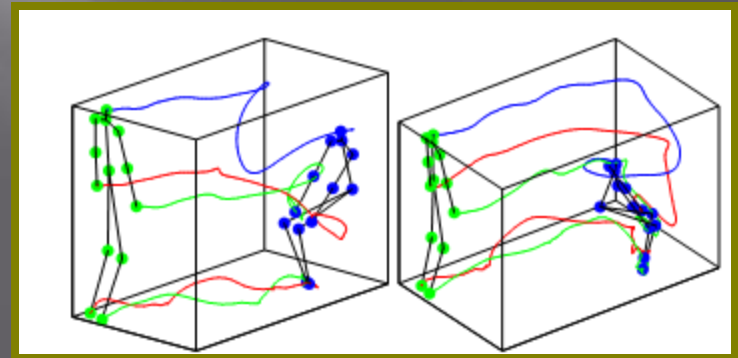
Experiments Motion Capture Data



Walk



Jump



Sit



Experiments Motion Capture Data

	Dance	Jump	Run	Sit	Walk
Dance	28				2
Jump		13			1
Run	2	1	22	1	4
Sit				33	
Walk	3		2		43

Mean Accuracy: 89.7%

Experiment -II

- ▣ Wizemann Action Data Set
- ▣ Nine actions performed by nine different actors:
 - Bend, Jumping Jack, Jump Forward, Jump in Place, Run, Side Gallop, Walk, Wave One Hand, Wave Two Hands
- ▣ 81 videos



Experiment-II



Experiments Wizemann Action Data Set

	Bend	Jumping Jack	Jumping Forward	Jumping In Place	Run	Side Gallop	Walk	Wave1	Wave2
Bend	9								
Jumping Jack		9							
Jump Forward			5	2	2				
Jump In Place				9					
Run					8		1		
Side Gallop					1	8			
Walk							9		
Wave1								9	
Wave2									9



Papers

- ▣ Cen Rao, Alper Yilmaz, and Mubarak Shah, View-Invariant Representation And Recognition of Actions, International Journal of Computer Vision, Vol.50, Issue 2, 2002.
(<http://www.cs.ucf.edu/~vision/papers/ijcv2002.pdf>).
- ▣ Alper Yilmaz and Mubarak Shah, Actions As Objects: A Novel Action Representation, IEEE CVPR 2005, San Diego, June 20-26.
(http://www.cs.ucf.edu/~vision/papers/yilmaz_cvpr_2005.pdf)
- ▣ Alexei Gritai, Yaser Sheikh, and Mubarak Shah, On the Invariant Analysis of Human Actions, 17th conference of the International Conference on Pattern Recognition, 2004.
(http://www.cs.ucf.edu/~vision/papers/gritai_icpr_2004.pdf)



Papers

- ▣
Yaser Sheikh, Alexei Gritai, and Mubarak Shah On the Spacetime Geometry of Galilean Cameras, IEEE Conference on Computer Vision and Pattern Recognition, Minneapolis, USA 2007.
(<http://server.cs.ucf.edu/~vision/papers/758.pdf>)

- ▣ Saad Ali, Arslan Basharat, and Mubarak Shah, Chaotic Invariants for Human Action Recognition, ICCV 2007, Rio de Janeiro, Brazil.
(<http://server.cs.ucf.edu/~vision/papers/SaadICCV07.pdf>)

