

Action Recognition

Copyright Mubarak Shah 2003

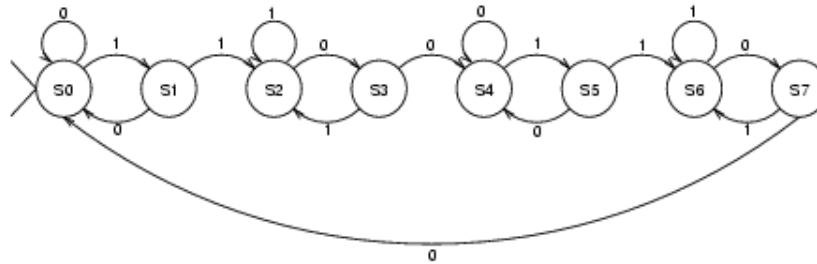
Approaches

- FSA
- HMMs/NNs
- Rule-based
- ---

- Representation is important

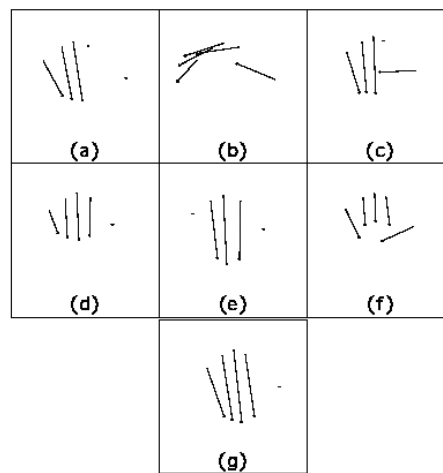
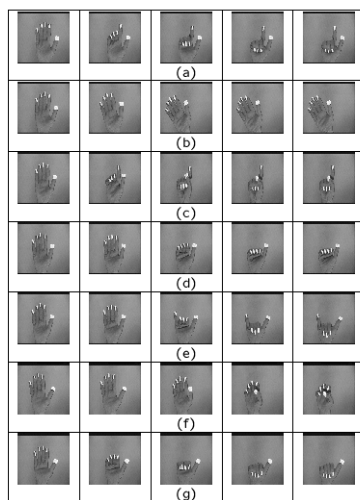
Copyright Mubarak Shah 2003

FSA: Hand Gesture Recognition



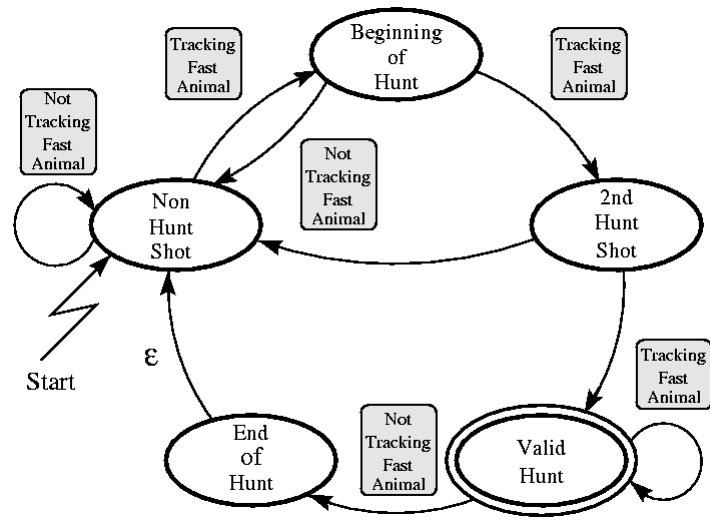
Copyright Mubarak Shah 2003

FSA: Hand Gesture Recognition

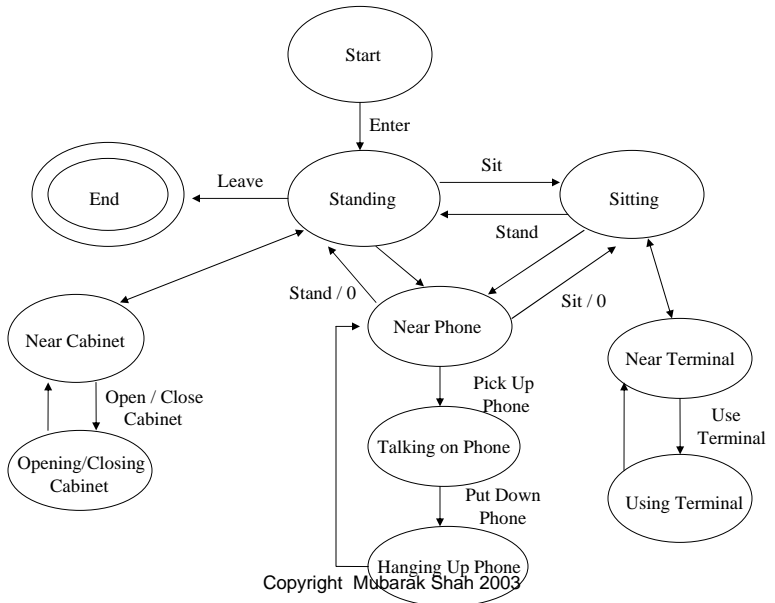


Mubarak Shah 2003

Hunt events



FSA: Recognizing Human Behavior in Office Environment



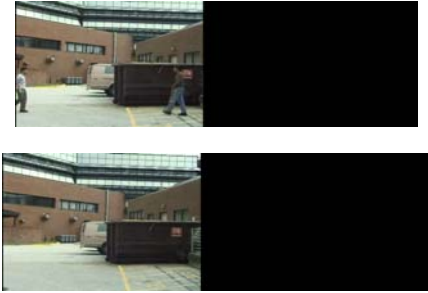
Copyright Mubarak Shah 2003

Rule-Based: Detecting Violence



Copyright Mubarak Shah 2003

Rule-Based: Detecting Violence



Copyright Mubarak Shah 2003

Rule-Based: Recognizing Outdoor Activities



Copyright Mubarak Shah 2003

Limitations

- A priori knowledge
- Extensive training
- No explanation
- No learning
- Representation
- View invariance

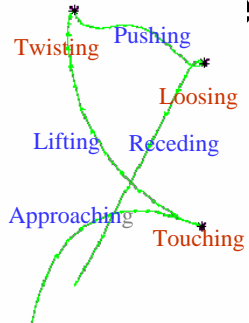
Copyright Mubarak Shah 2003

View-invariant Representation and Recognition of Human Action

Hand Actions Recognition

- hand generates a 3-D trajectory with respect to time.
- analyze 2-D projection of this 3-D trajectory.
- View invariance issues.

Instant- presentation



Copyright Mubarak Shah 2003

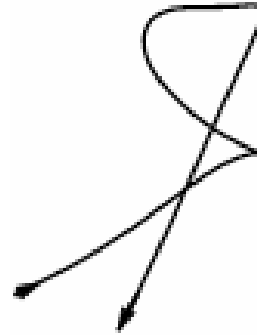
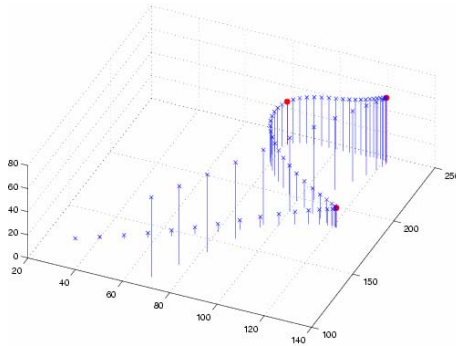
Hand Actions



Copyright Mubarak Shah 2003

Spatiotemporal Curve

$$r_{st} = [x(t) \quad y(t) \quad t]$$



Copyright Mubarak Shah 2003

Spatiotemporal Curvature

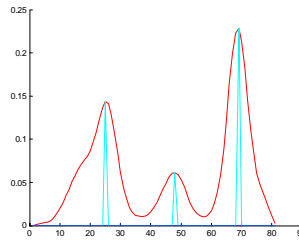
$$k = \frac{\sqrt{A^2 + B^2 + C^2}}{\left((x')^2 + (y')^2 + (t')^2\right)^{\frac{3}{2}}} \quad A = \begin{vmatrix} y' & t' \\ y'' & t'' \end{vmatrix}, B = \begin{vmatrix} t' & x' \\ t'' & x'' \end{vmatrix}, C = \begin{vmatrix} x' & y' \\ x'' & y'' \end{vmatrix}$$

Spatiotemporal curvature captures both the **speed** and **direction** changes in **one quantity**.

Copyright Mubarak Shah 2003

Representation of Actions

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



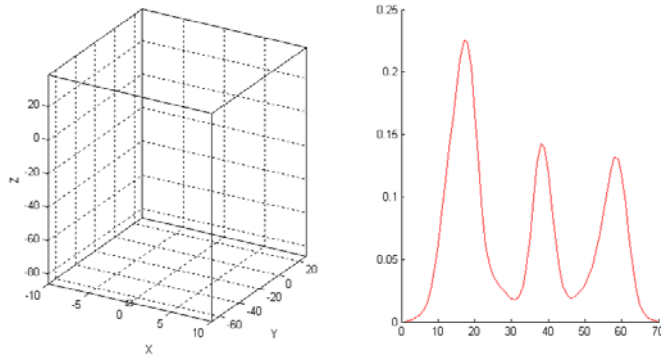
Copyright Mubarak Shah 2003

Action Units

- Psychology research shows:
 - People tend to divide an action into atomic units at the places where the motion characteristics change the most.
 - Stops, starts, pauses, dynamic instants

Copyright Mubarak Shah 2003

Viewing Directions and Spatiotemporal Curvature

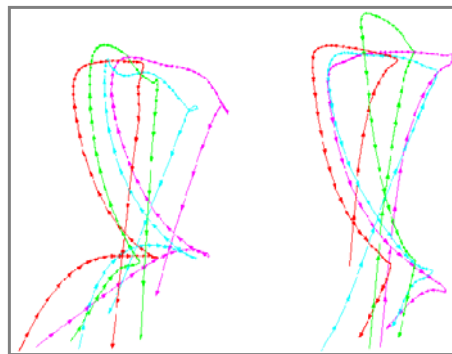


Although the viewing directions are quite different, the peak locations are consistent.

Copyright Mubarak Shah 2003

View-invariant Action Representation

- Opening and closing an overhead cabinet



Opening Closing
Copyright Mubarak Shah 2003

View-invariant Matching

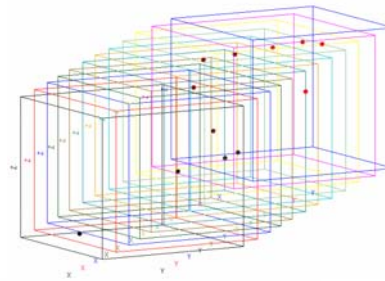
- Consider 3D trajectories as 3D objects.

Copyright Mubarak Shah 2003

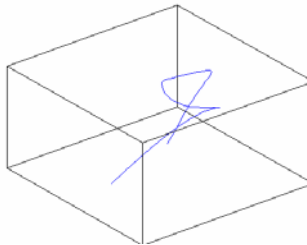
Action Trajectory in 4D



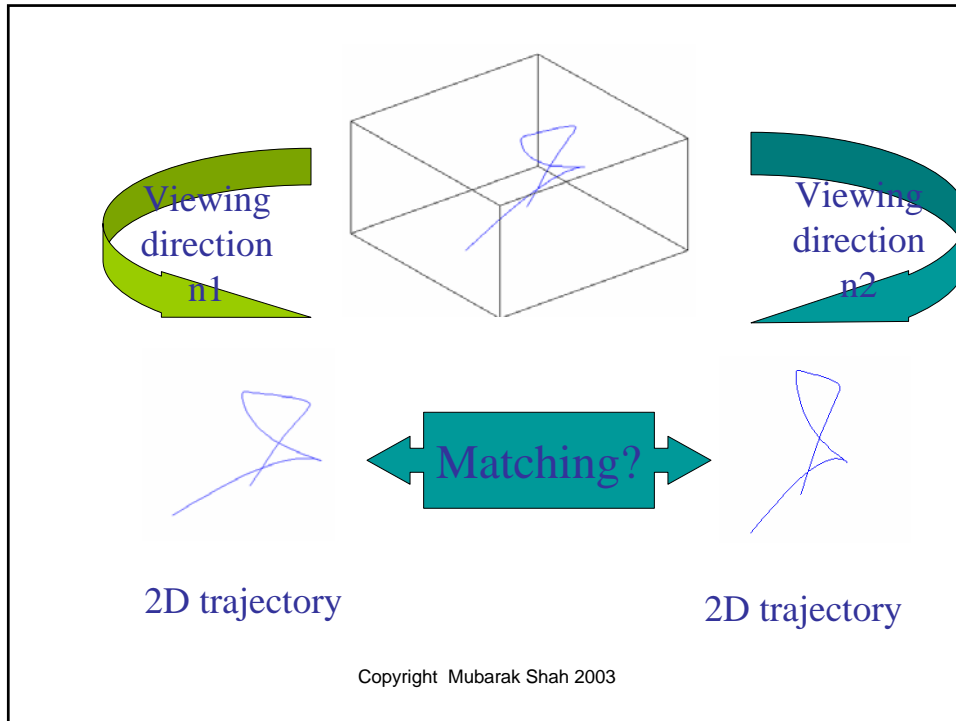
Sampling in time



Ignore the time index



Copyright Mubarak Shah 2003



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:

$$M = \begin{bmatrix} \mu_1^{v_1} & \mu_2^{v_1} & \dots & \mu_n^{v_1} \\ v_1^{v_1} & v_2^{v_1} & \dots & v_n^{v_1} \\ \mu_1^{v_2} & \mu_2^{v_2} & \dots & \mu_n^{v_2} \\ v_1^{v_2} & v_2^{v_2} & \dots & v_n^{v_2} \end{bmatrix} = P \cdot S = \begin{bmatrix} \Pi_{v_1} \\ \Pi_{v_2} \end{bmatrix} \cdot \begin{bmatrix} X_1 & X_2 & \dots & X_n \\ Y_1 & Y_2 & \dots & Y_n \\ Z_1 & Z_2 & \dots & Z_n \end{bmatrix}$$

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

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

Copyright Mubarak Shah 2003

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:

$$M = \begin{bmatrix} \mu_1^i & \mu_2^i & \dots & \mu_n^i \\ \nu_1^i & \nu_2^i & \dots & \nu_n^i \\ \mu_1^j & \mu_2^j & \dots & \mu_n^j \\ \nu_1^j & \nu_2^j & \dots & \nu_n^j \end{bmatrix} \quad \text{dist}_{i,j} = |\sigma_4|$$

Copyright Mubarak Shah 2003

Perspective View-Invariant Matching

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

Copyright Mubarak Shah 2003

Perspective View-invariant Measure

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

$$\begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix}^T F \begin{bmatrix} u'_i \\ v'_i \\ 1 \end{bmatrix} = 0, \quad Mf = \begin{bmatrix} u_1 u'_1 & u_1 v'_1 & u'_1 & v_1 u'_1 & v_1 v'_1 & v'_1 & u_1 & v_1 & 1 \\ u_2 u'_2 & u_2 v'_2 & u'_2 & v_2 u'_2 & v_2 v'_2 & v'_2 & u_2 & v_2 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ u_n u'_n & u_n v'_n & u'_n & v_n u'_n & v_n v'_n & v'_n & u_n & v_n & 1 \end{bmatrix} f = 0$$

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.

Copyright Mubarak Shah 2003

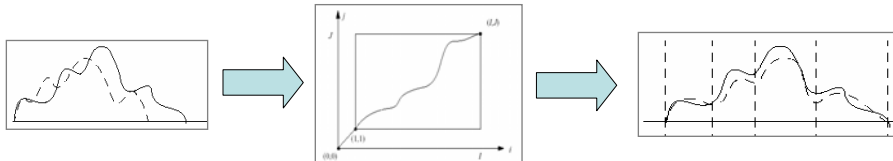
Instant-Interval Representation

- We have not used the interval information
- View Invariance
- Temporal Invariance

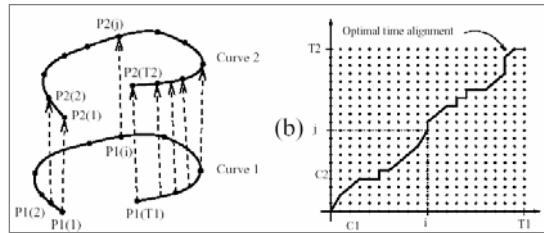
Copyright Mubarak Shah 2003

DTW and Temporal Signals

- Match two 1D temporal signals:



- Match two 2D temporal curves:



Warping

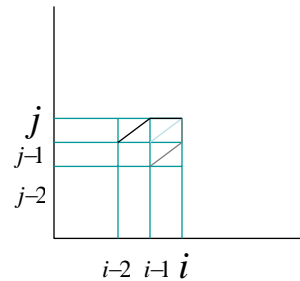
$$A = [a_1, a_2, \dots, a_i, a_I]$$

$$B = [b_1, b_2, \dots, b_j, b_J]$$

$$d_{ij} = |a_i - b_j|$$

$$g_{11} = 2d_{11}$$

$$g(i, j) = \min \begin{bmatrix} g(i-1, j-2) + 2d(i, j-1) + d(i, j) \\ g(i-1, j-1) + 2d(i, j) \\ g(i-2, j-1) + 2d(i-1, j) + d(i, j) \end{bmatrix}$$



Copyright Mubarak Shah 2003

Affine View-invariant DTW

Step 1: Pick up 4 instants from trajectories I and I' , (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , (x_4, y_4) and (x'_1, y'_1) , (x'_2, y'_2) , (x'_3, y'_3) , (x'_4, y'_4) are image coordinates.

Step 2: Apply Dynamic Time Warping

- the similarity between i^{th} (u_i, v_i) and j^{th} (u'_j, v'_j) point is:

σ_4 is the 4th singular value of matrix M

Copyright Mubarak Shah 2003

Affine View-invariant DTW (Con.)

Step 3: if there are more than 4 pairs of instants in the trajectories, go back to step 1 and try other combinations of 4 instants.

Step 4: pick the minimal matching error as the similarity measurement and get the correspondence result.

Copyright Mubarak Shah 2003

Matching using View-invariant Dynamic Time Warping

Step 1) Pick up 4 instants from each of trajectory

Step 2) DTW using the view invariant similarity measure, σ_4 , for intervals

3) More Instants left?

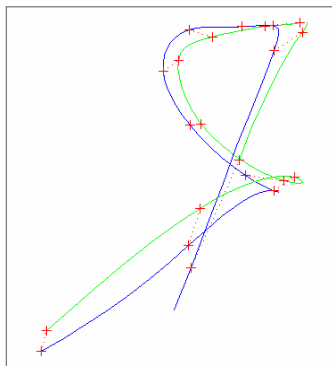
Yes

No

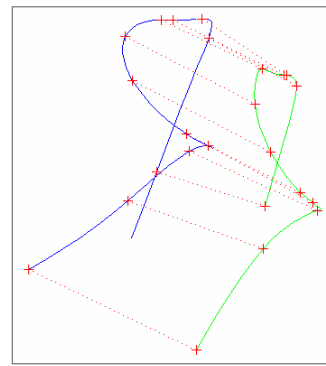
Step 4) get the minimum similarity measurement and correspondence result

Copyright Mubarak Shah 2003

Action Recognition Results



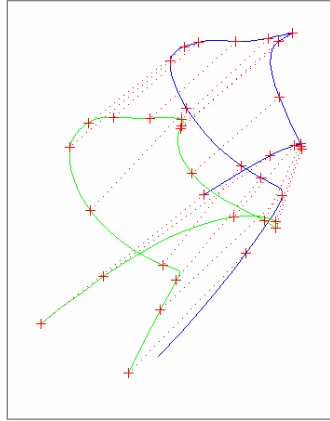
Difference = 2



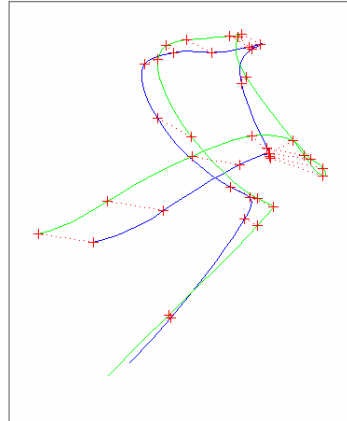
Difference = 2.3

Copyright Mubarak Shah 2003

Action Recognition Results



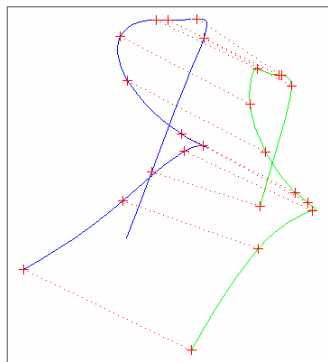
Difference = 2.5



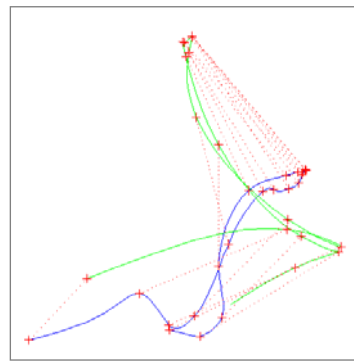
Difference = 3.2

Copyright Mubarak Shah 2003

View-invariant DTW Results



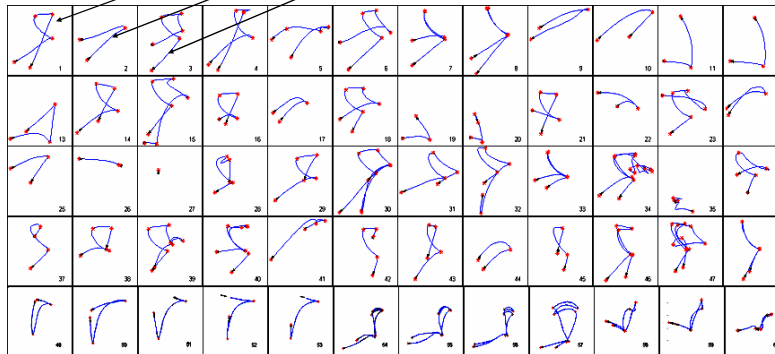
Difference = 2.3



Difference = 71

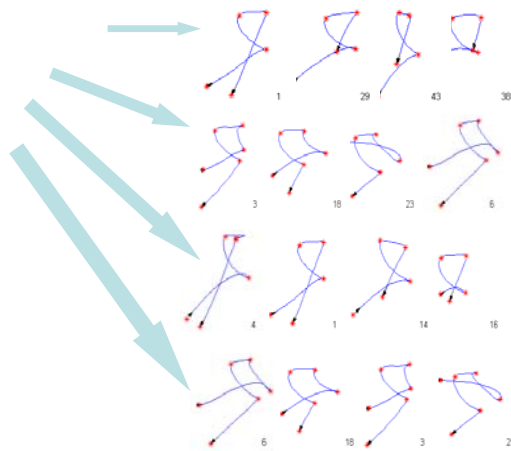
Copyright Mubarak Shah 2003

Experimental Results
60 Action Trajectories
7 People



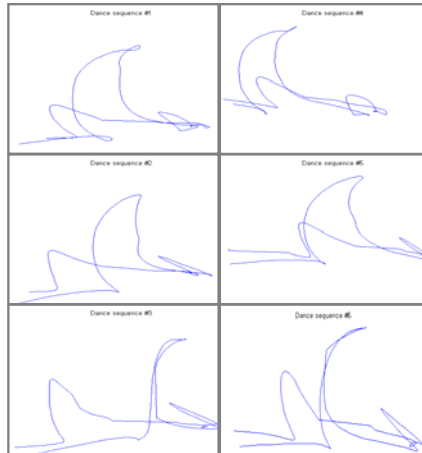
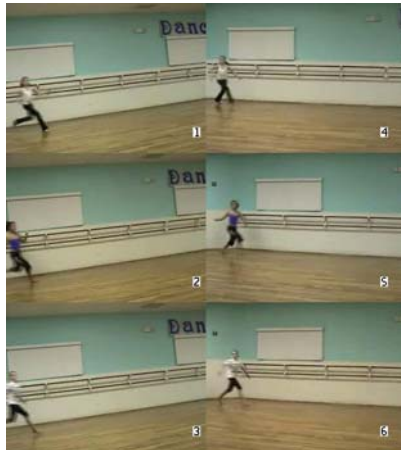
Copyright Mubarak Shah 2003

Experimental Results



< Shah 2003

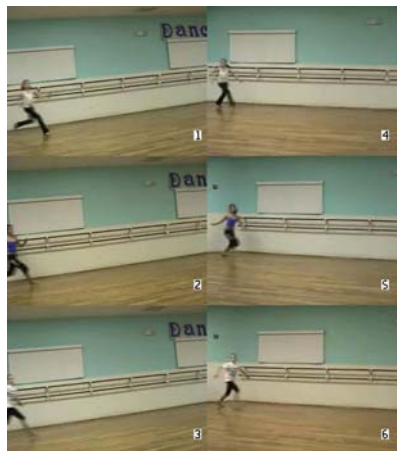
Temporal Alignment of Videos



Input videos: Copyright Mubarak Shah 2003 Trajectories of the right foot:

Temporal Alignment Results

Synchronized videos:



Copyright Mubarak Shah 2003

Temporal Alignment Results

Before Temporal Alignment



After Temporal Alignment



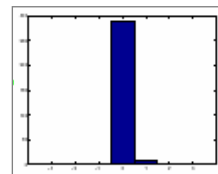
Temporal Alignment Results

Non-overlapping video temporal alignment:



Left Seq

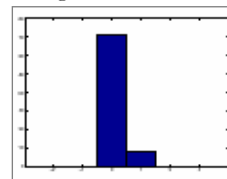
Right Seq



Histogram of Misalignment



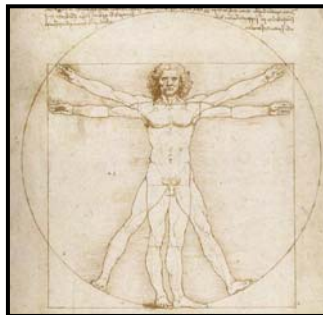
Copyright Mubarak Shah 2003



- Cen Rao, Alexei Gritai, Mubarak Shah, [View-invariant Alignment and Matching of Video Sequences](#). The Ninth IEEE International Conference on Computer Vision, Nice, France, 2003.
[Project web page](#).
- Cen Rao, Alper Yilmaz, Mubarak Shah. [View-Invariant Representation And Recognition of Actions](#), International Journal of Computer Vision, Vol. 50, Issue 2, 2002.

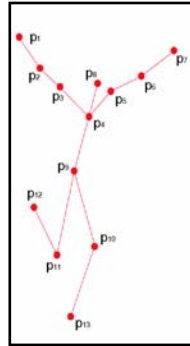
Copyright Mubarak Shah 2003

Anthropometric Representation for Invariant Action Recognition



Copyright Mubarak Shah 2003

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.

Copyright Mubarak Shah 2003

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.

Copyright Mubarak Shah 2003

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

Copyright Mubarak Shah 2003

Anthropometric Constraint

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

$$\mathbf{X}_i = M \mathbf{Y}_i$$

where $i = 1, 2 \dots n$, M is a 4×4 non-singular matrix, \mathbf{X}_i and \mathbf{Y}_i are sets of points describing two actors.

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

Copyright Mubarak Shah 2003

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.

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

Dimension	Men					Women				
	5th %ile	50th %ile	90th %ile	95th %ile	SD	5th %ile	50th %ile	90th %ile	95th %ile	SD
1. Stature	1625	1740	1855	70	1505	1610	1710	62		
2. Eye height	1515	1630	1745	69	1405	1505	1610	61		
3. Shoulder height	1315	1425	1535	66	1215	1310	1405	58		
4. Elbow height	1055	1090	1180	52	930	1000	1085	46		
5. Hip height	840	920	1000	50	740	810	885	43		
6. Knuckle height	690	755	825	41	605	720	790	36		
7. Fingertip height	590	655	720	35	540	625	695	30		
8. Sitting height	850	915	965	36	790	850	910	35		
9. Sitting eye height	735	790	845	35	685	740	795	33		
10. Sitting shoulder height	540	595	645	32	500	555	610	31		
11. Sitting elbow height	395	445	505	31	385	435	490	29		
12. Thigh thickness	135	180	195	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	430	480	530	30		
15. Knee height	490	545	595	32	485	530	580	27		
16. Popliteal height	395	440	490	29	385	430	480	27		
17. Shoulder breadth (biacromial)	430	485	510	28	365	395	425	24		
18. Shoulder breadth (biacromial)	385	400	430	20	325	355	385	18		
19. Hip breadth	310	360	405	29	310	370	435	28		
20. Chest (max) depth	215	260	285	22	210	260	295	27		
21. Abdominal depth	220	270	305	32	205	255	305	30		
22. Shoulder-elbow length	330	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	655	715	32	555	600	650	29		
26. Head length	180	195	205	8	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	70	75	80	4		
30. Foot length	240	265	285	14	215	235	255	12		
31. Foot breadth	85	88	110	6	80	90	100	6		
32. Span	1655	1750	1825	83	1490	1605	1725	71		
33. Elbow span	865	945	1020	47	785	865	920	43		
34. Vertical grip reach (standing)	1925	2050	2190	80	1790	1905	2020	71		
35. Vertical grip reach (sitting)	1145	1245	1340	60	1060	1160	1235	53		
36. Forward grip reach	720	780	835	34	650	705	755	31		

Copyright Mubarak Shah 2003

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.

Copyright Mubarak Shah 2003

Postural Constraint

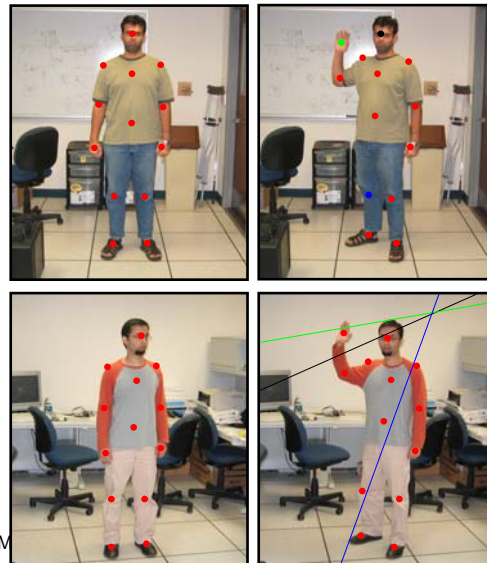
- The similarity of *posture* between two actors can be measured using the **ninth singular** value of a measurement matrix A , where $Af = 0$.

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

Copyright Mubarak Shah 2003

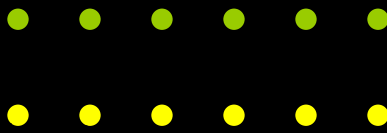
Capturing View Variance

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



Copyright M

- **Proposition 2:** For an action element \mathbf{u}_t , the fundamental matrices associated with (x_t, y_t) and (x_{t+1}, 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

$$\mathbf{A} = [A_1, A_2 \dots A_k]$$

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

Copyright Mubarak Shah 2003

Action Detection

Analyzing Periodicity



Reference Pattern

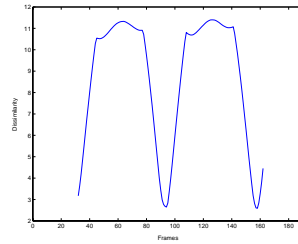


Test Sequence

Copyright Mubarak Shah 2003

Action Detection

Analyzing Periodicity



911
159
51

Copyright Mubarak Shah 2003

Action Detection: Different approaches, different people, the same action



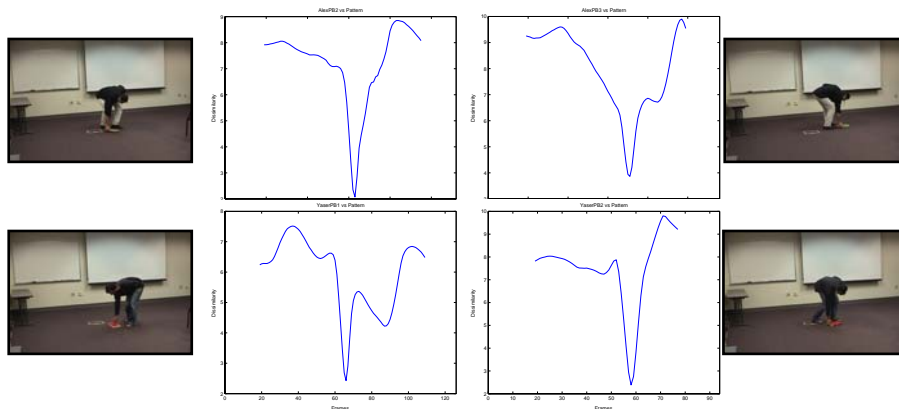
ReferencePattern



Test Sequences

Copyright Mubarak Shah 2003

Action Detection: Different approaches, different people, the same action



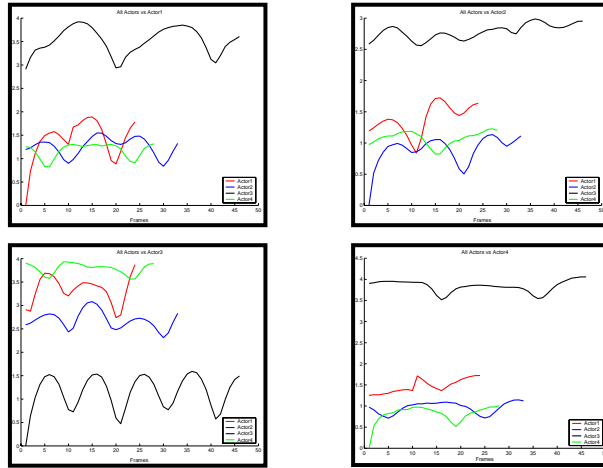
Copyright Mubarak Shah 2003

Analyzing Actions *Odd One Out*



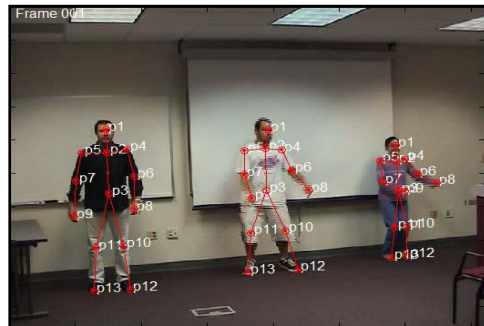
Copyright Mubarak Shah 2003

'Odd One Out'



Copyright Mubarak Shah 2003

Action Synchronization *Following the Leader*



Copyright Mubarak Shah 2003

Action Synchronization

Following the Leader



Copyright Mubarak Shah 2003