

+ Large volumes of video



- For YouTube alone
 - More than 1 billion unique user visits each month
 - Over 6 billion hours of video are watched each month
 - 100 hours of video are uploaded every minute

+ Available Datasets

Dataset
 KTH
 Weizmann
 IXMAS
 Hollywood
 UCF Sports
 Hollywood2
 UCF YouTube
 MSR
 Olympic
 UCF50
 HMDB51

#Classes

6
 9
 11
 8
 9
 12
 11
 3
 16
 50
 51



<http://serre-lab.clps.brown.edu/resource/hmdb-a-large-human-motion-database/>

+ Videos in the wild

- Unrestricted type of events with various activities



Harlem Shake : <http://www.youtube.com/watch?v=4hpEnLtgUDg>

+ Our attempts

- Videos as sequence of frames
 - Detect concepts in each frame
 - Utilize image search engines
- Discover important knowledge from videos itself
 - Discriminate parts
- Understand actions in videos
 - Simple but effective descriptors





Utilizing large volumes of weakly labeled images

Pinar Duygulu, January 13, 2014, CMU

+ Utilize image search results

Query : Paris

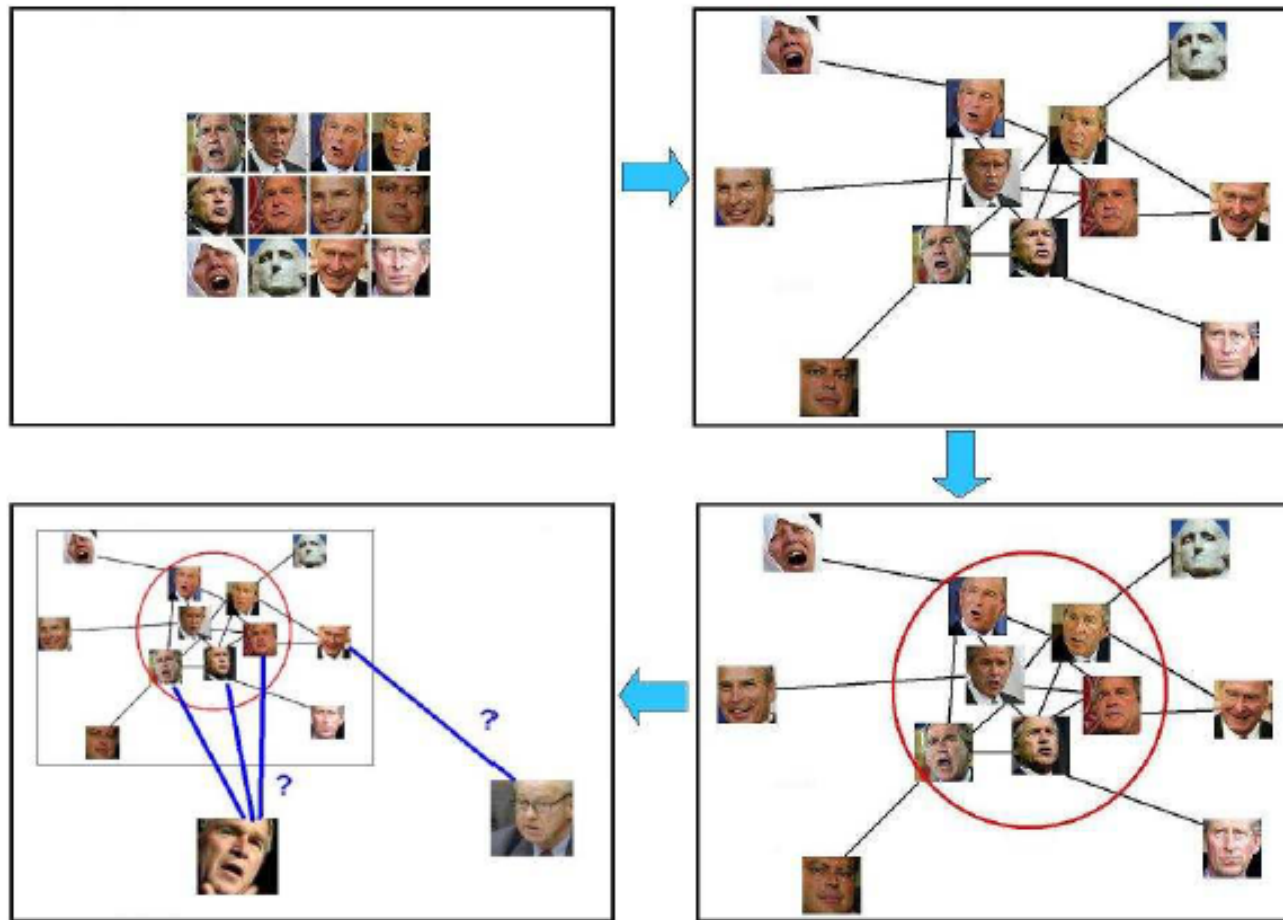


+ Single Dominant Category

Query : George W. Bush



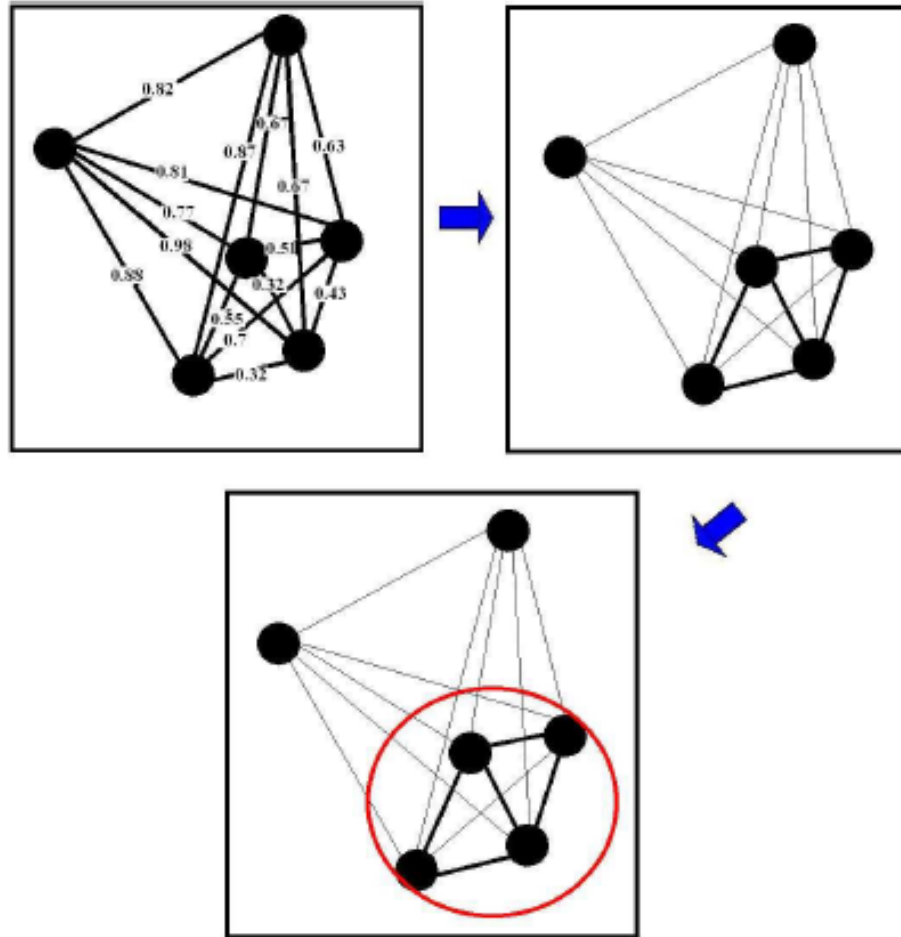
+ Naming faces



Among the faces associated with a name find the correct subset :
The most similar subset of faces

Ozkan, D., Duygulu, P., "Interesting Faces: A Graph Based Approach for Finding People in News", Pattern Recognition, 2010
 Ozkan, D., Duygulu, P., "A Graph Based Approach for Naming Faces in News Photos", CVPR, 2006
 Ozkan, D., Duygulu, P., "Finding People Frequently Appearing in News", CIVR, 2006

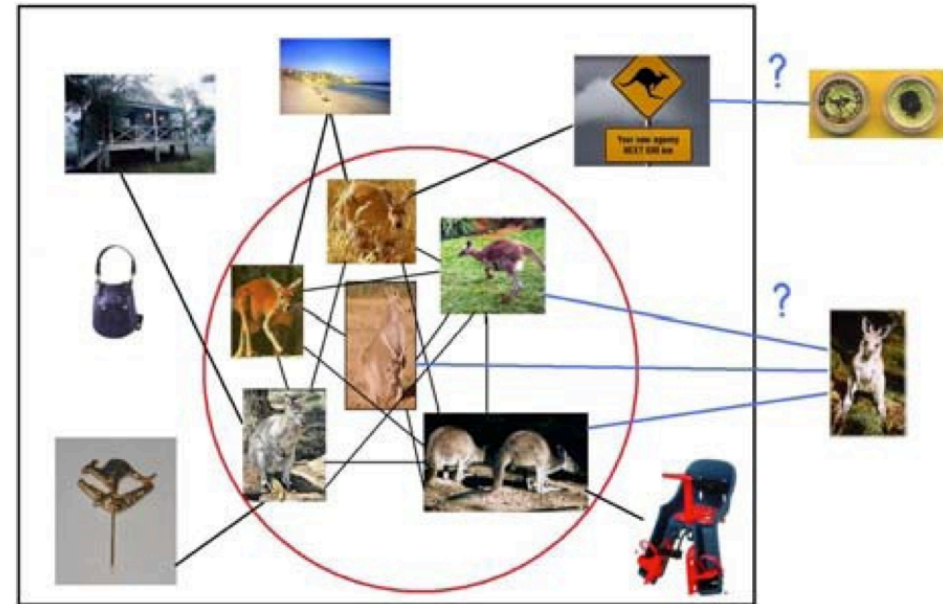
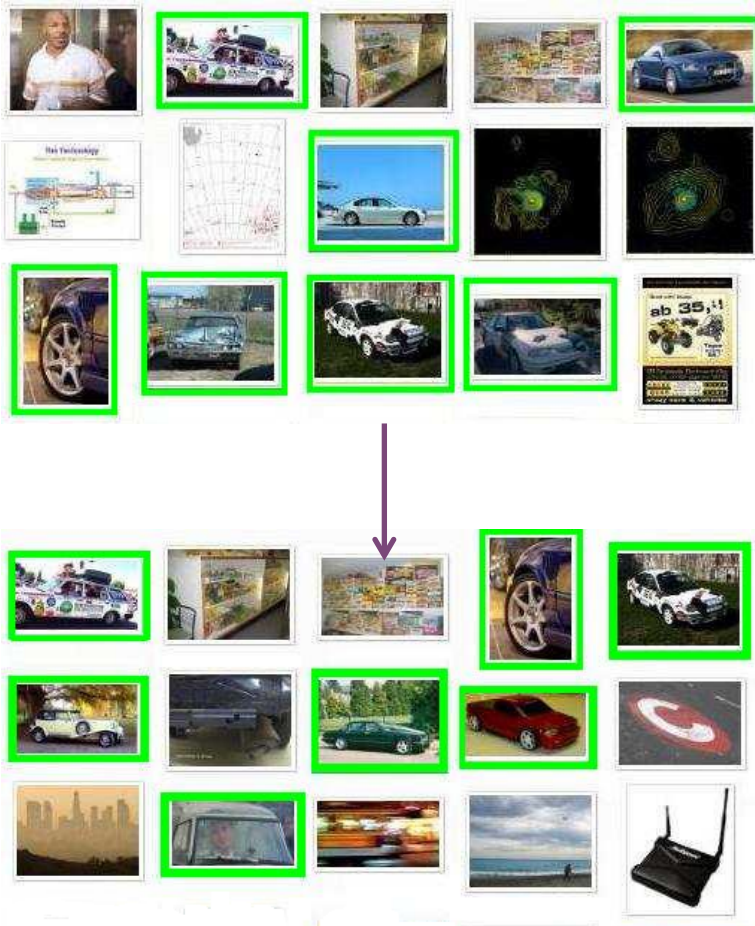
+ Finding Densest component



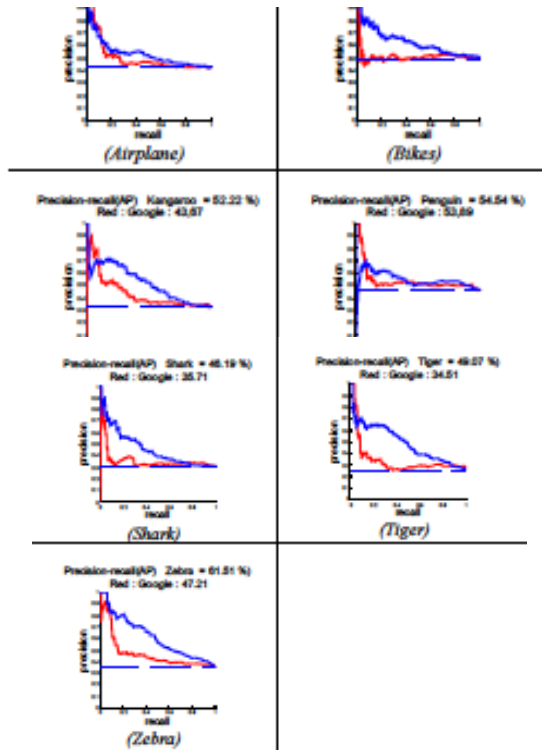
$$f(S) = \frac{|E(S)|}{|S|},$$

Node with the minimal degree is removed at each iteration (Charikar, 2000)

+ Image Re-ranking



+ Multiple Instance Learning for re-ranking



On the dataset by Schroff, F., ICCV 2007
 “Harvesting Image Databases from the Web”.

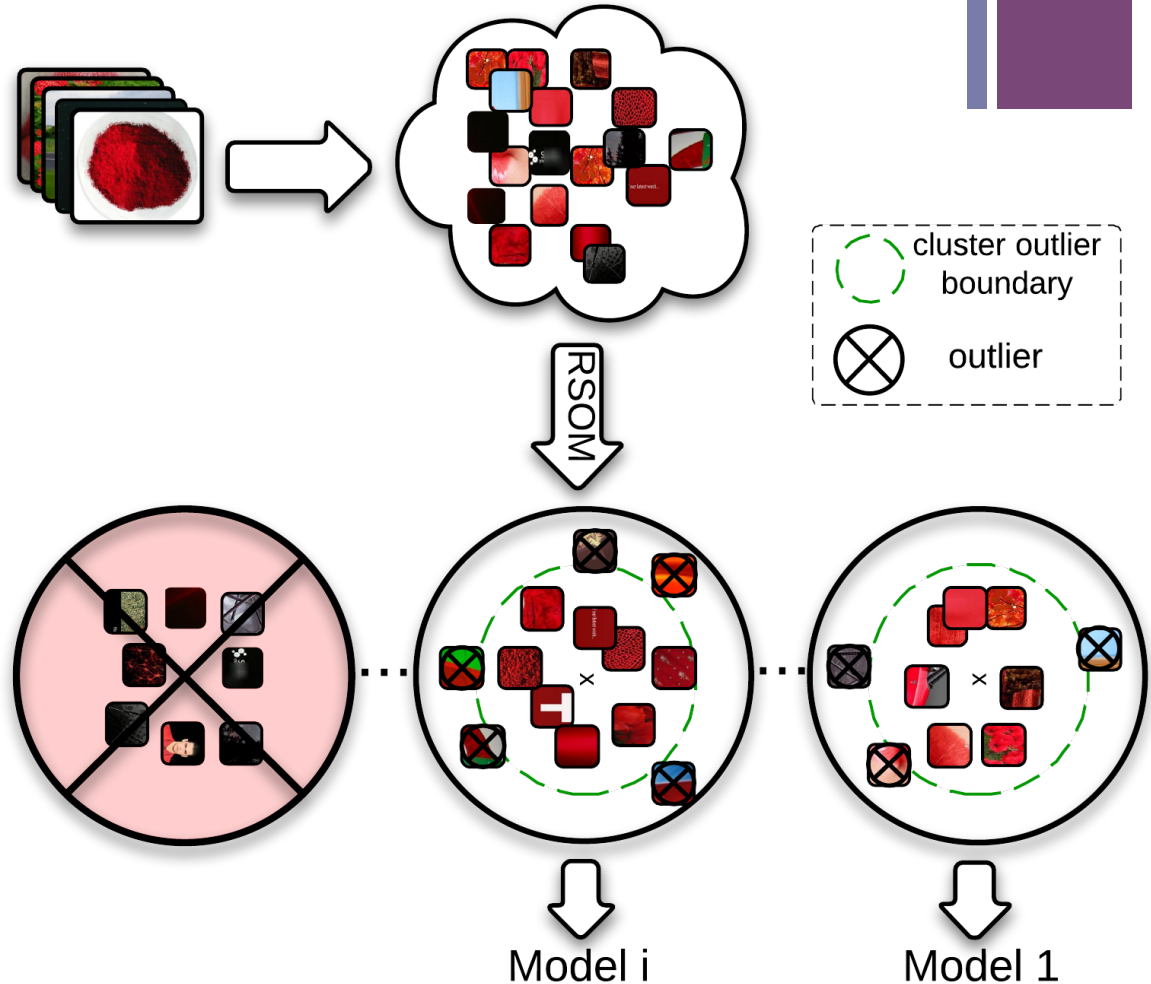
+ Multiple meanings/variations



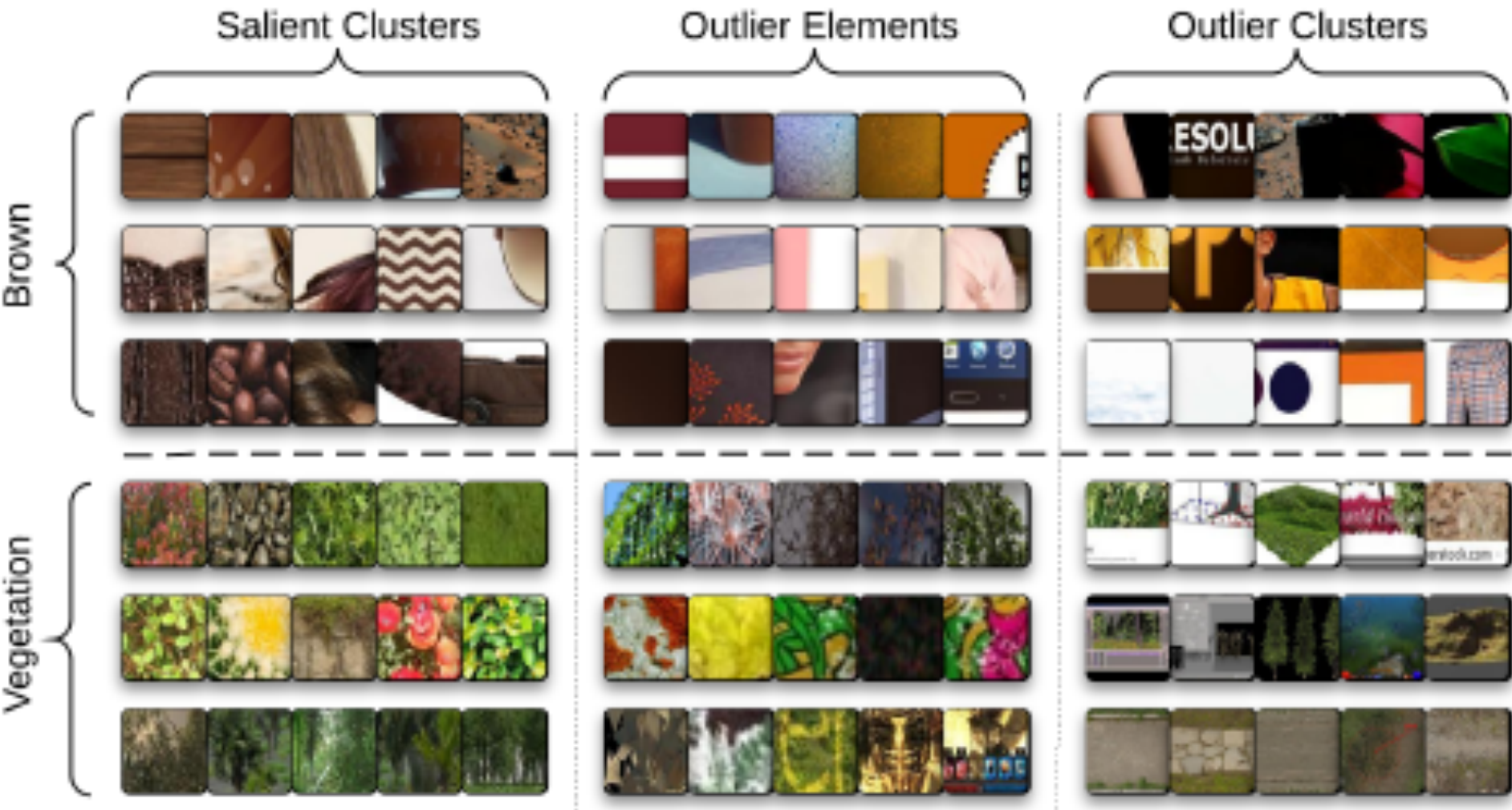
the attributes are observed in different forms and in small portions requiring grouping and non-attribute parts to be eliminated.

+ RSOM for Concept Learning

- Collect images from web for a keyword
- Clustering and outlier detection
- RSOM (Rectifying Self Organizing Maps)
- Learn a model for each cluster



+ Color and Texture Attributes



+ Scene Concepts



+ Attribute and Scene Learning

Attribute learning

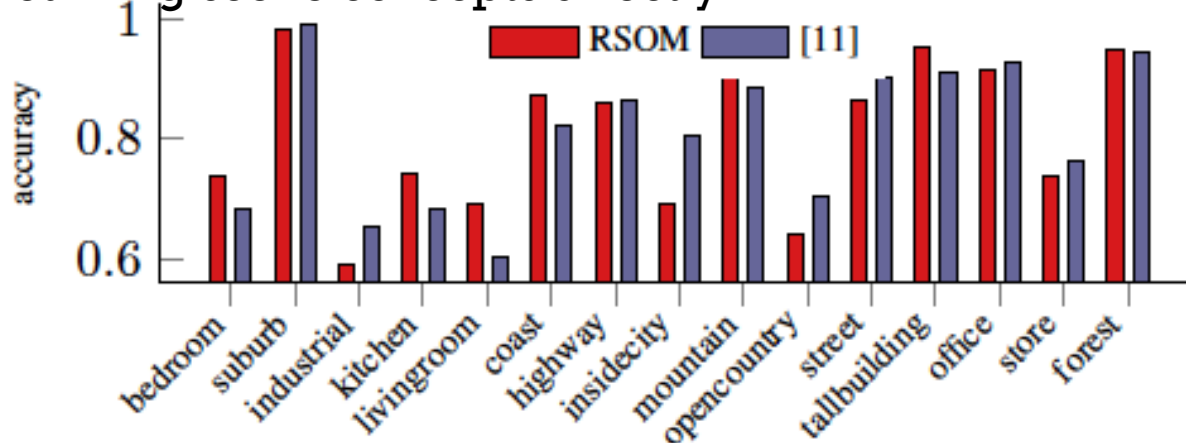
Method	RSOM-M	RSOM	PLSA-reg [22].
cars	0.97	0.92	0.93
shoes	1.0	0.97	0.99
dresses	1.0	1.0	0.99
pottery	0.98	0.92	0.94
overall	0.99	0.95	0.96

Attribute based scene recognition

Method	MIT-indoor [17]	Scene-15 [11]
RSOM-A	46.2%	82.7%
RSOM-S	-	80.7%
RSOM-S+HM	-	81.3%
Li <i>et al.</i> [12] VQ	47.6%	82.1%
Pandey <i>et al.</i> [16]	43.1%	-
Kwitt <i>et al.</i> [9]	44%	82.3%

On ImageNet: 37.4% (RSOM), 36.8% (Russakovsky & Fei-Fei, 2012)

Learning scene concepts directly



[17] Quattoni and Torralba, "Recognizing Indoor Scenes". 2009

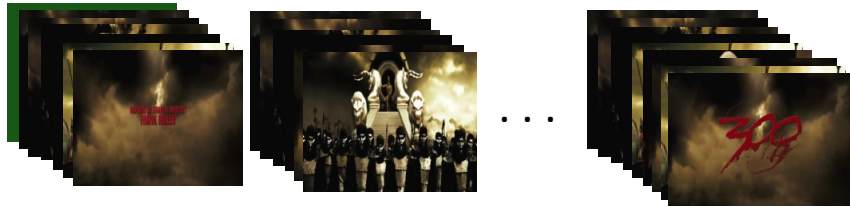
[11] Lazebnik, Schmid, Ponce, "Beyond Bags of features: Spatial pyramid matching for recognizing natural scene categories", CVPR 2006

[22] Van de Weijer, Schmid, Verbeek, Larlus, "Learning Color Names for Real-world Applications", 2009

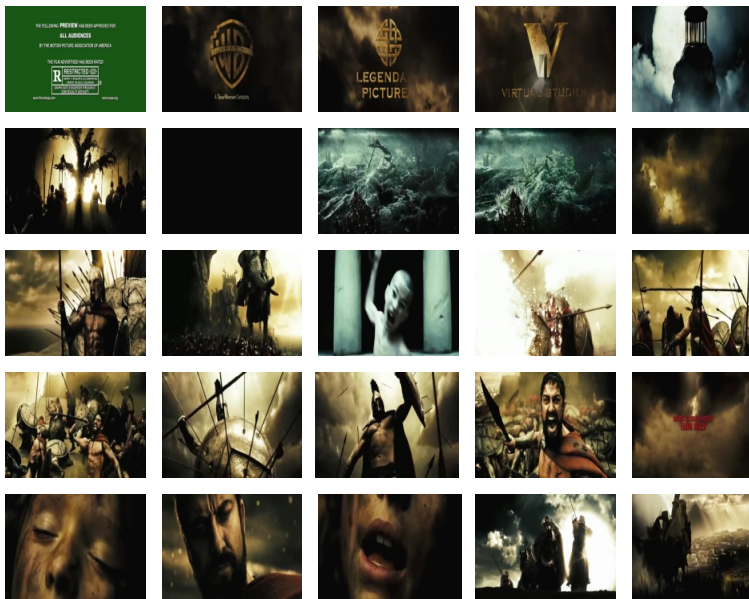


Utilizing videos

+ Movie Genre & MPAA Ratings



300 Spartans – Movie Trailer (represented as frames)

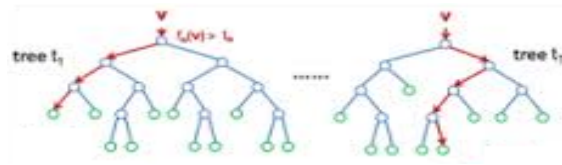


■ Statistical features :

- Key frame count, Black and White frame counts, Avg. in shot frame difference, Total number of frames, Avg. shot length, Avg. key frame HSV color histogram, Avg. key frame edge histogram

■ Visual features:

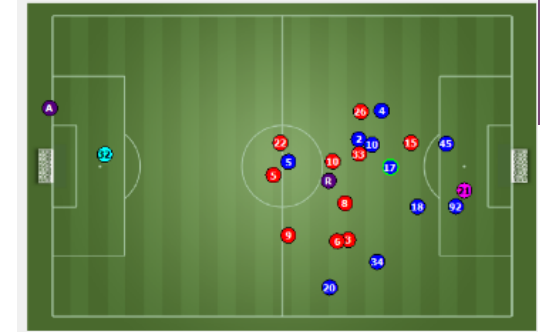
- SIFT, LAB histogram, Local Binary Patterns, GIST



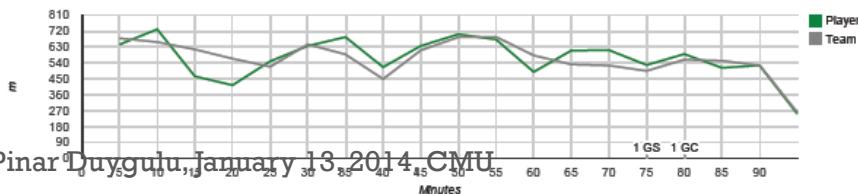
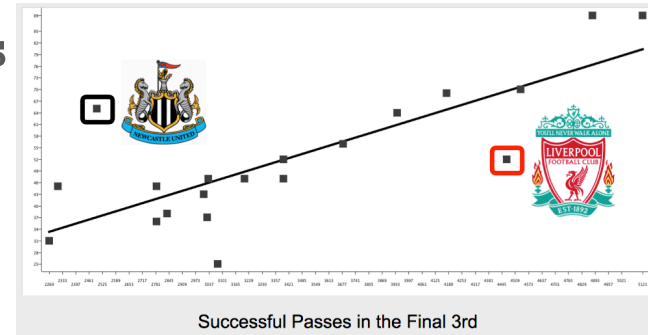
<p>GENERAL AUDIENCES</p> <p>G</p> <p>GENERAL AUDIENCES ALL AGES ADMITTED</p>	
<p>PARENTAL GUIDANCE SUGGESTED</p> <p>PG</p> <p>PARENTAL GUIDANCE SUGGESTED SOME MATERIAL MAY BE OFFENSIVE TO CHILDREN</p>	
<p>PARENTS STRONGLY CAUTIONED</p> <p>PG-13</p> <p>PG-13 PARENTS STRONGLY CAUTIONED SOME MATERIAL MAY BE OFFENSIVE TO CHILDREN</p>	
<p>RESTRICTED</p> <p>R</p> <p>RESTRICTED PARENTS STRONGLY CAUTIONED</p>	
<p>NO ONE 17 AND UNDER ADMITTED</p> <p>NC-17</p> <p>NO ONE 17 AND UNDER ADMITTED</p>	

Genre: action, horror, animation, comedy, drama

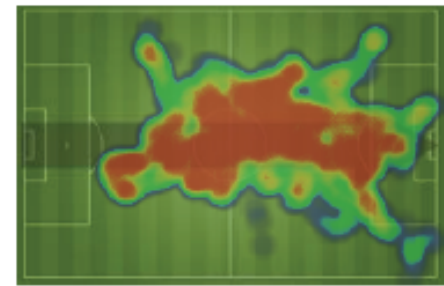
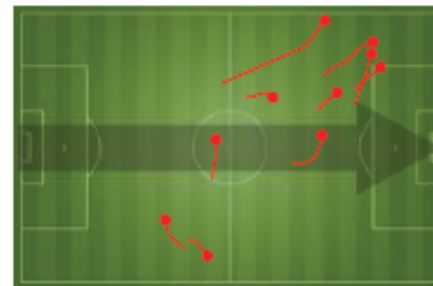
+ Video Analysis for Sports



- Track soccer players in real-time in all matches
 - Extract XY coordinates and event data
- Using the big data
 - Player/team performance analysis/comparison
 - Extract game features that is correlated with seasonal success
 - Performance indexing and market value evaluation
 - Fatigue and injury prediction



Pinar Duygulu, January 13, 2014, CMU



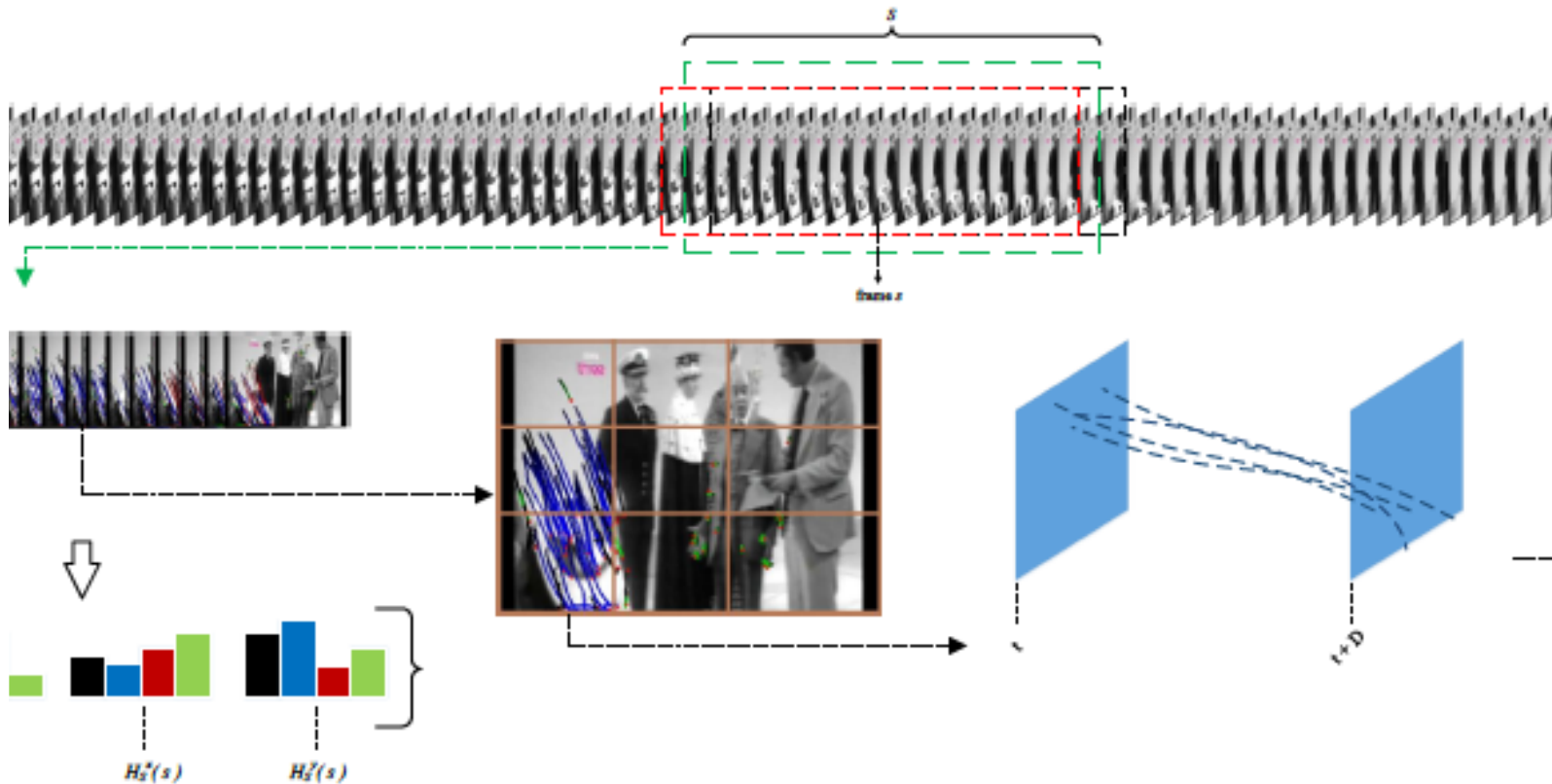
+ Utilize video data: Capture usualness in unusual videos



+ Usual versus unusual

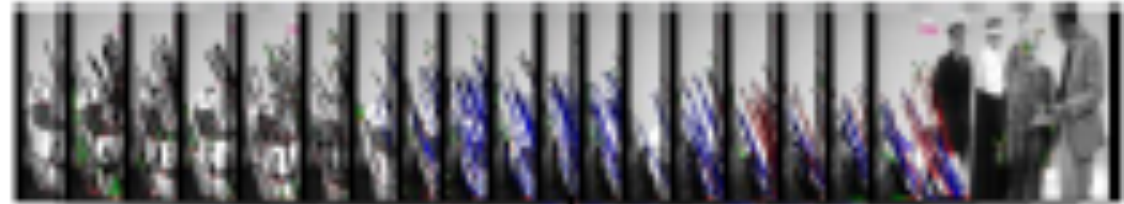
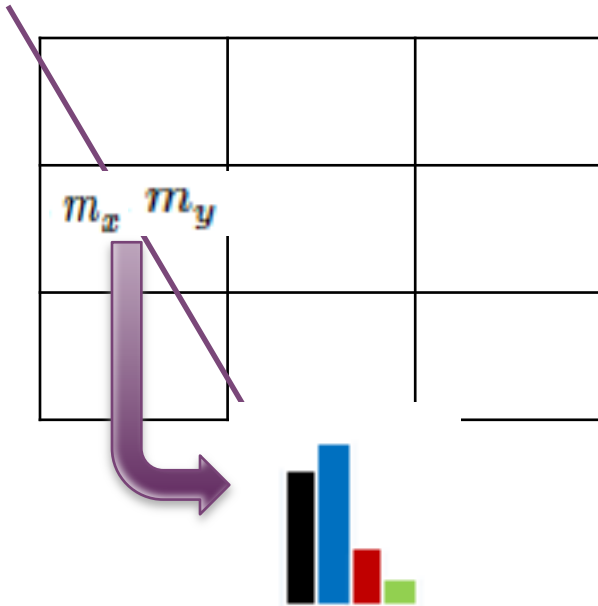


+ Trajectory Snippet Histograms



Iscen, A., Armagan, A., Duygulu, P., "What is usual in unusual videos? Trajectory snippet histograms for discovering unusualness", submitted to CVPR 2014

+ Representation



$$H_S^l = \sum_{t=s-(\|S\|/2)}^{s+(\|S\|/2)} H_S^l(t) \quad H_S = (H_S^l, H_S^x, H_S^y)$$

Velocity and spatial extension of the motion

$$H_S^l(t) = (H_S^l(t)_{[1,1]}, \dots, H_S^l(t)_{[1,N]}, \dots, H_S^l(t)_{[N,N]})$$



$$\vec{T} = (P_t, \dots, P_{t+D-1}) \quad \vec{P}_t = (x_t, y_t)$$

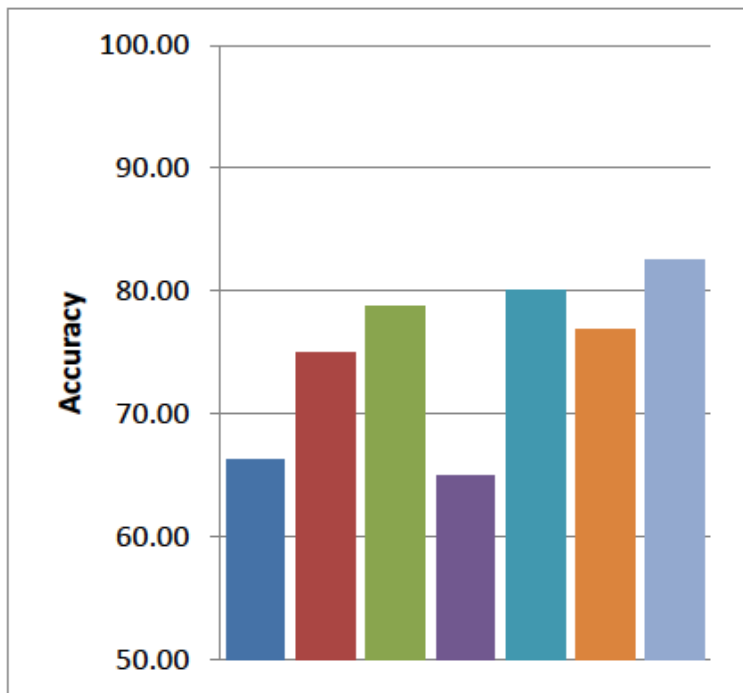
$$m_x = \frac{1}{D} \sum_t^{t+D-1} x_t, v_x = \frac{1}{D} \sum_t^{t+D-1} (x_t - m_x)^2$$

$$m_y = \frac{1}{D} \sum_t^{t+D-1} y_t, v_y = \frac{1}{D} \sum_t^{t+D-1} (y_t - m_y)^2,$$

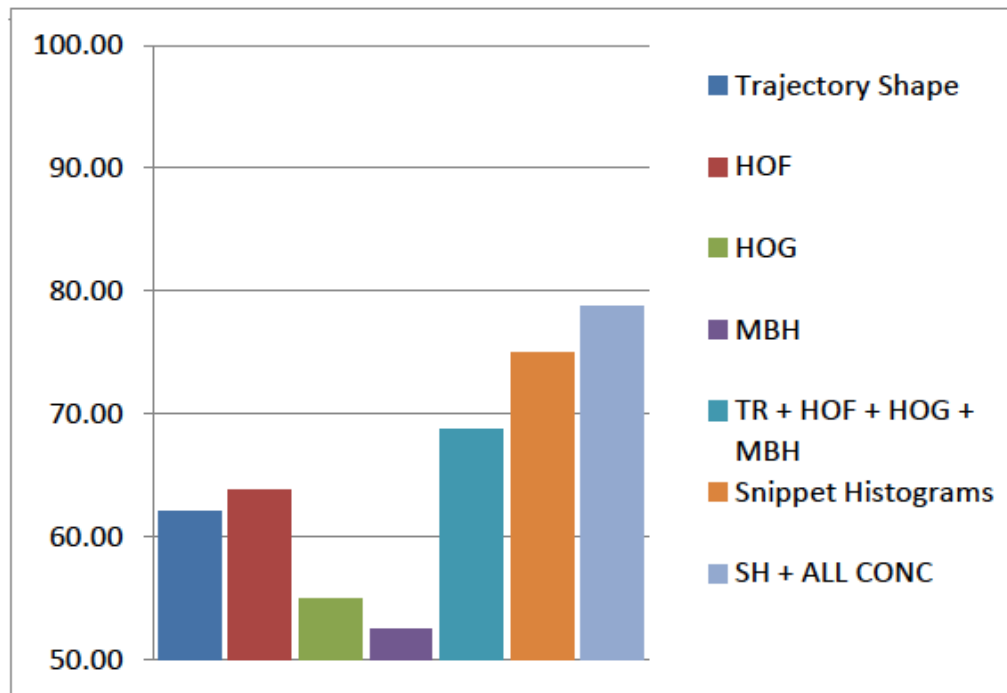
$$l = \sum_t^{t+D-1} \sqrt{(x_{t+1} - x_t)^2 + (y_{t+1} - y_t)^2}$$



Classification



People Falling



Funny videos

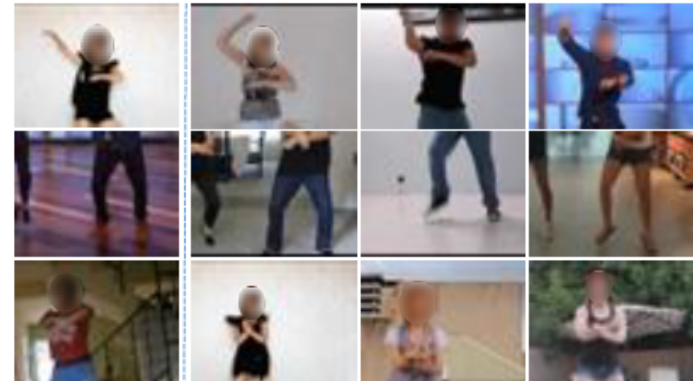
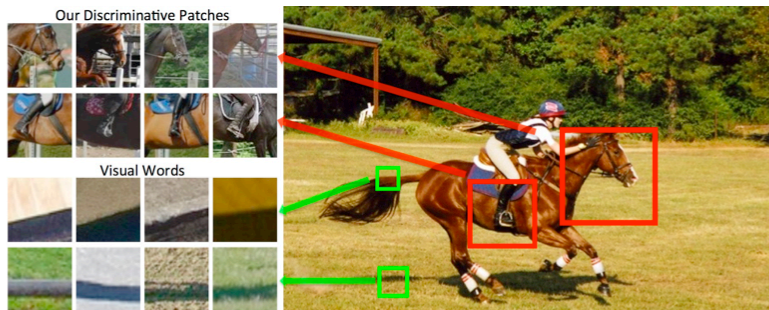
+ Snapshot Discovery

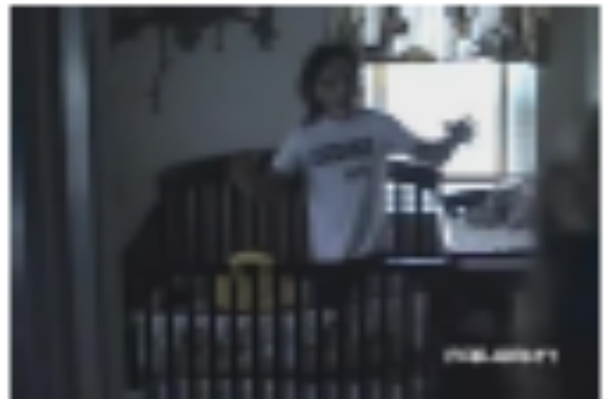
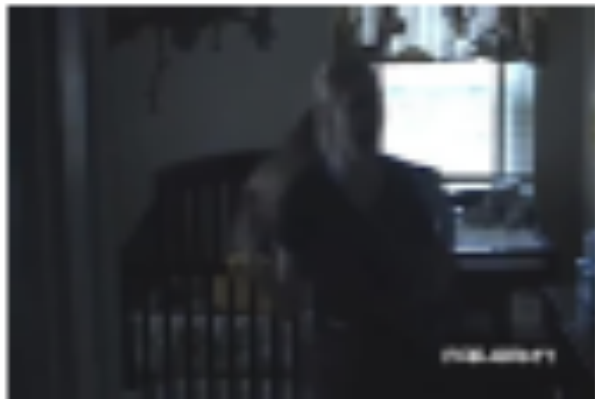
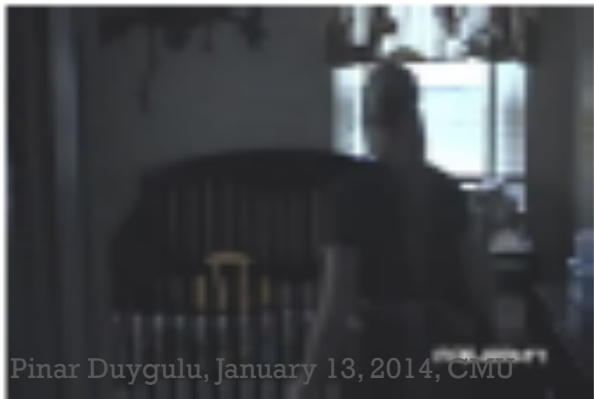
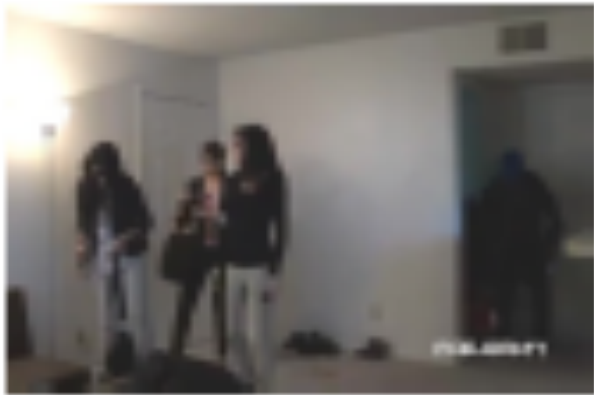
- Discriminative video patch idea over snippets (short video sequences)

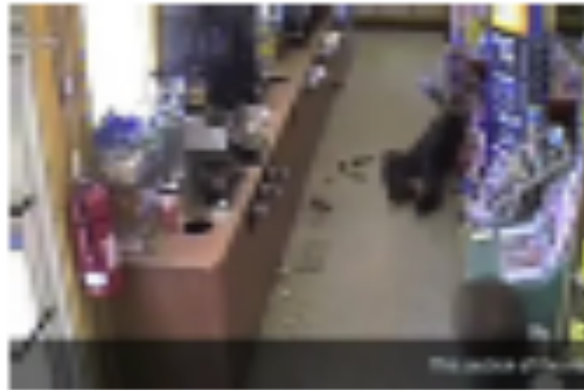
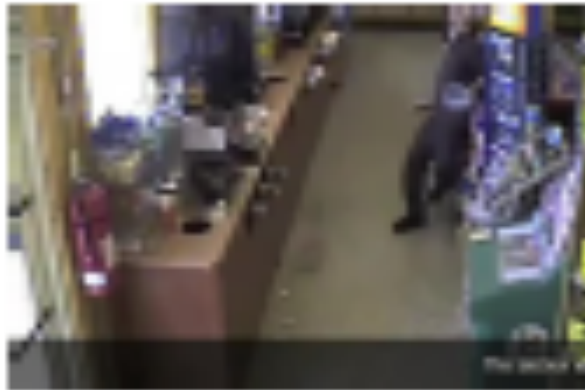
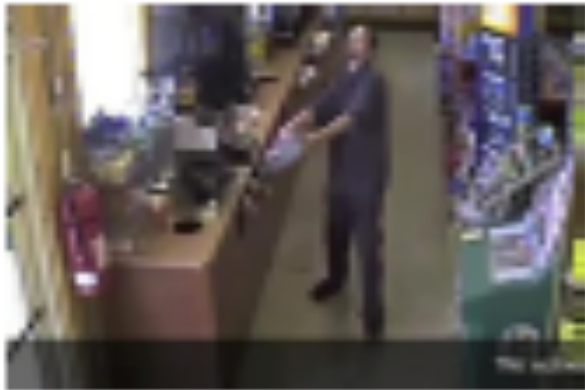
Singh ECCV 2012

Jain CVPR 2013

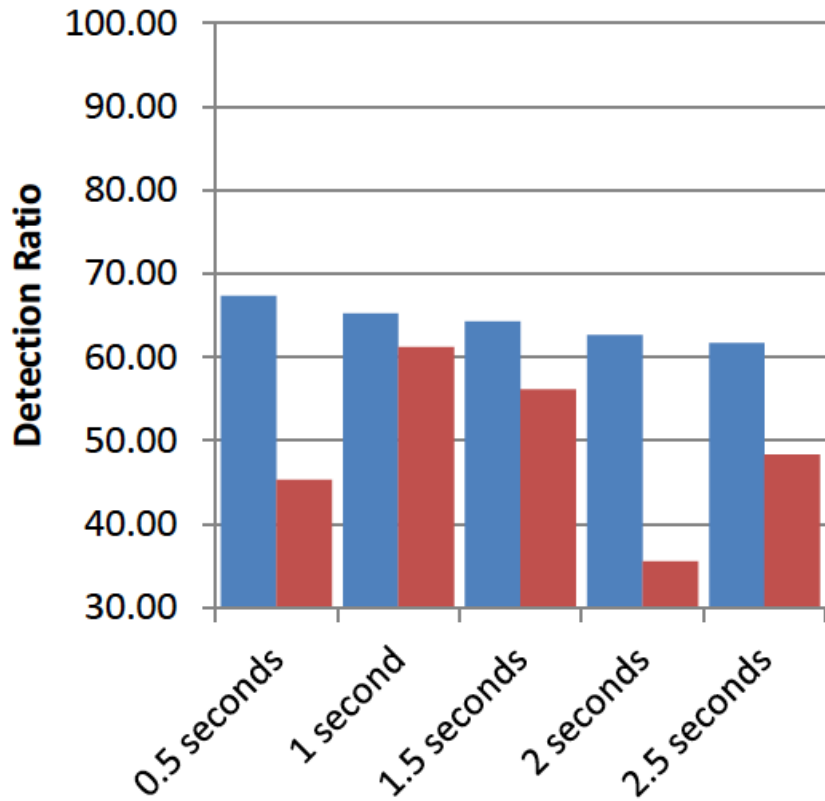
Unsupervised Discovery of Mid-Level Discriminative Patches



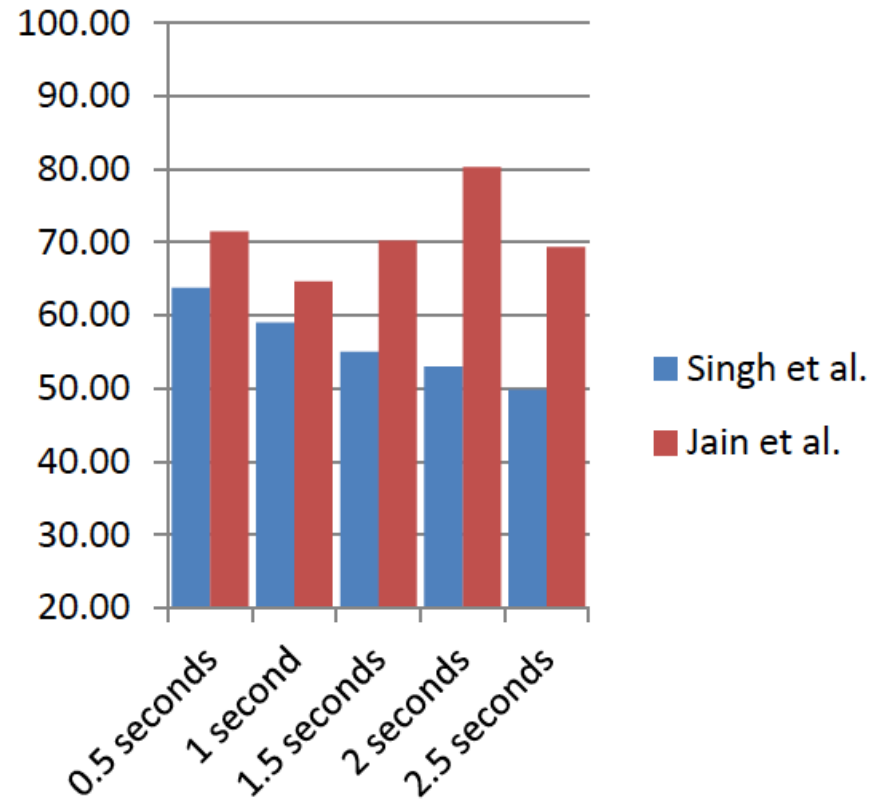




+ Snapshot discovery



People Falling



Funny videos

HOG3D --- people falling 25 %, funny videos 32 %

+ Birthday event



Flowing candles





Human Activity Analysis

Pinar Duygulu, January 13, 2014, CMU

30



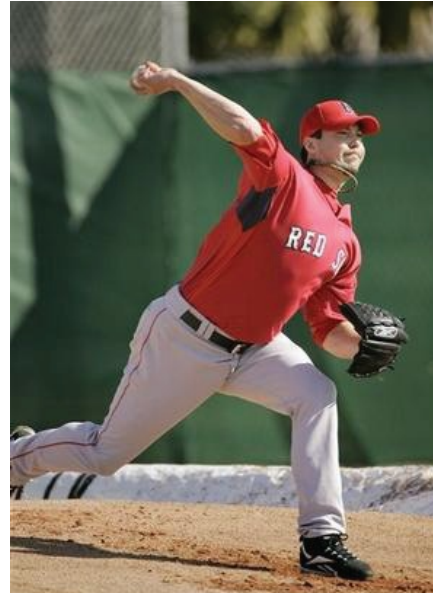
+ What do these people do?



running



walking



throwing

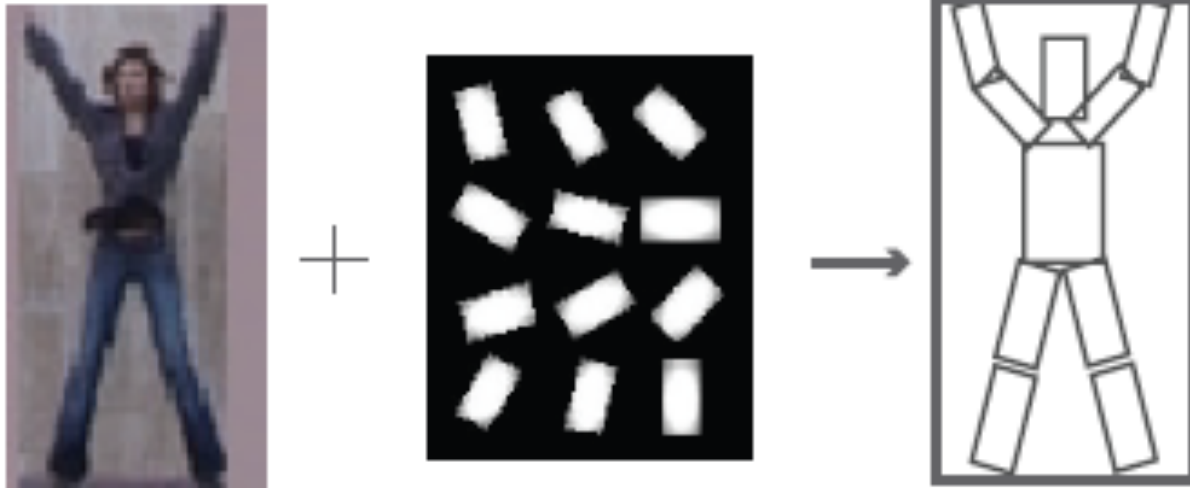


crouching

- Pose tells a lot about the actions.
- How can we describe the pose?

+ Pose as a Collection of Rectangles

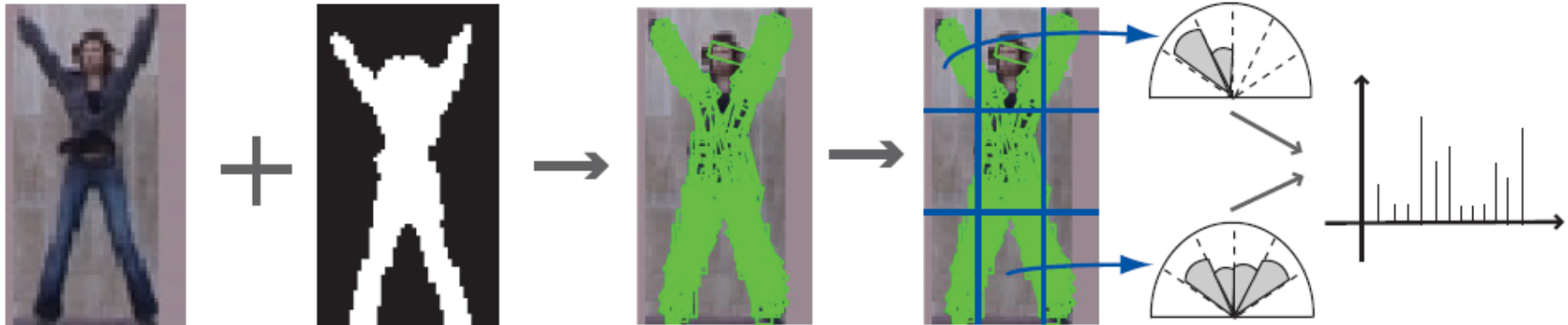
- Human body is composed of cylindrical parts.
- The projection of a cylinder on 2D is a rectangle.
- Body can be thought as a collection of rectangular regions
- We can represent the pose based on the orientation of these rectangles



Ikizler, N. Duygulu, P. "Human Action Recognition Using Distribution of Oriented Rectangular Patches", Proc. 2nd Workshop on Human Motion: Understanding, Modeling, Capture and Animation, In conjunction with ICCV2007

Ikizler, N. Duygulu, P. "Histogram of Oriented Rectangles: A New Pose descriptor for Human Action Recognition", Image and Vision Computing, volume 27, Issue 10, pages 1515-1526, September 2009

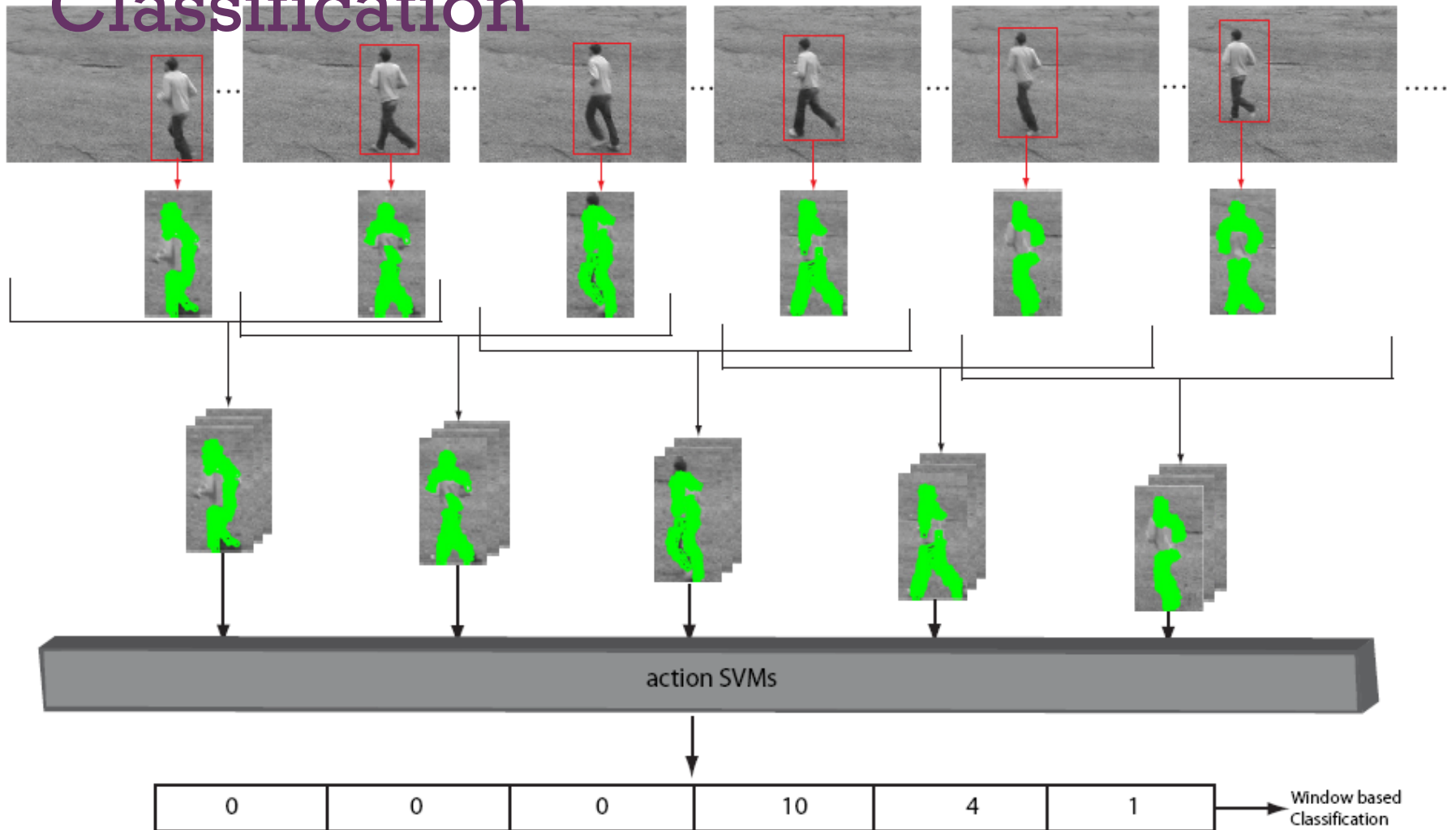
+ Histogram of Oriented Rectangles (HOR)



- Rectangular regions are extracted over silhouettes using convolution of a zero-padded rectangular 2D Gaussian on different orientations and scales
 - 12 angles 15° apart



Classification

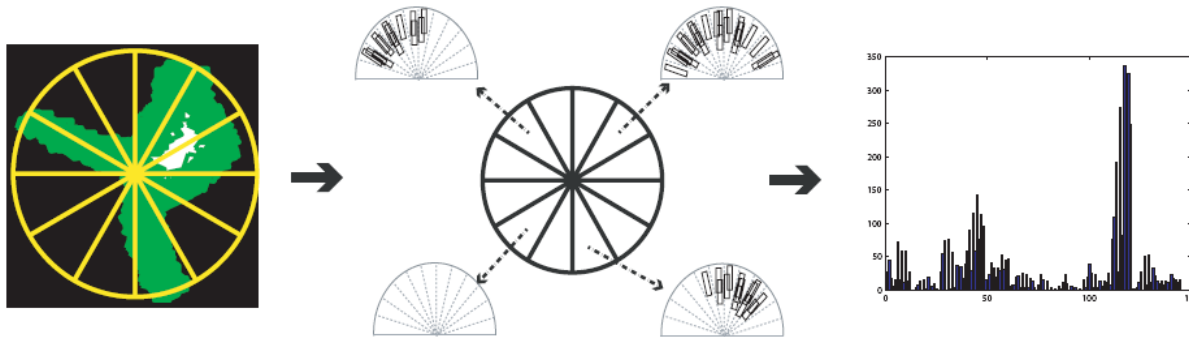


- Use snippets of frames and form histogram of oriented rectangles over a window (HORW)

+ Action Recognition in Still Images



- Pose estimation by Ramanan's method based on CRFs.
- Form Circular HORs (CHORs)
- Classification based on LDA+SVM



Ikizler, N., Cinbis, R. G., Pehlivan, S., Duygulu, P., "Recognizing actions from still images", Proc. 19th International Conference on Pattern Recognition (ICPR 2008)

+ Still Image Results



running



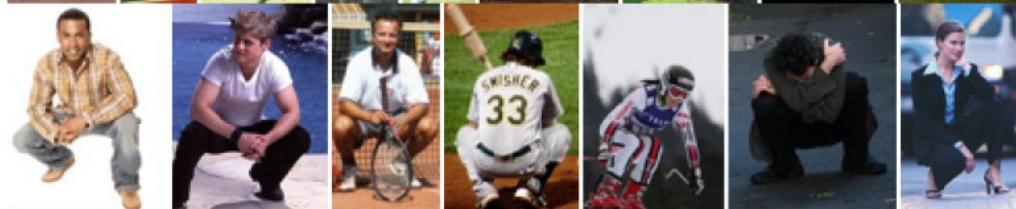
walking



throwing



catching



crouching



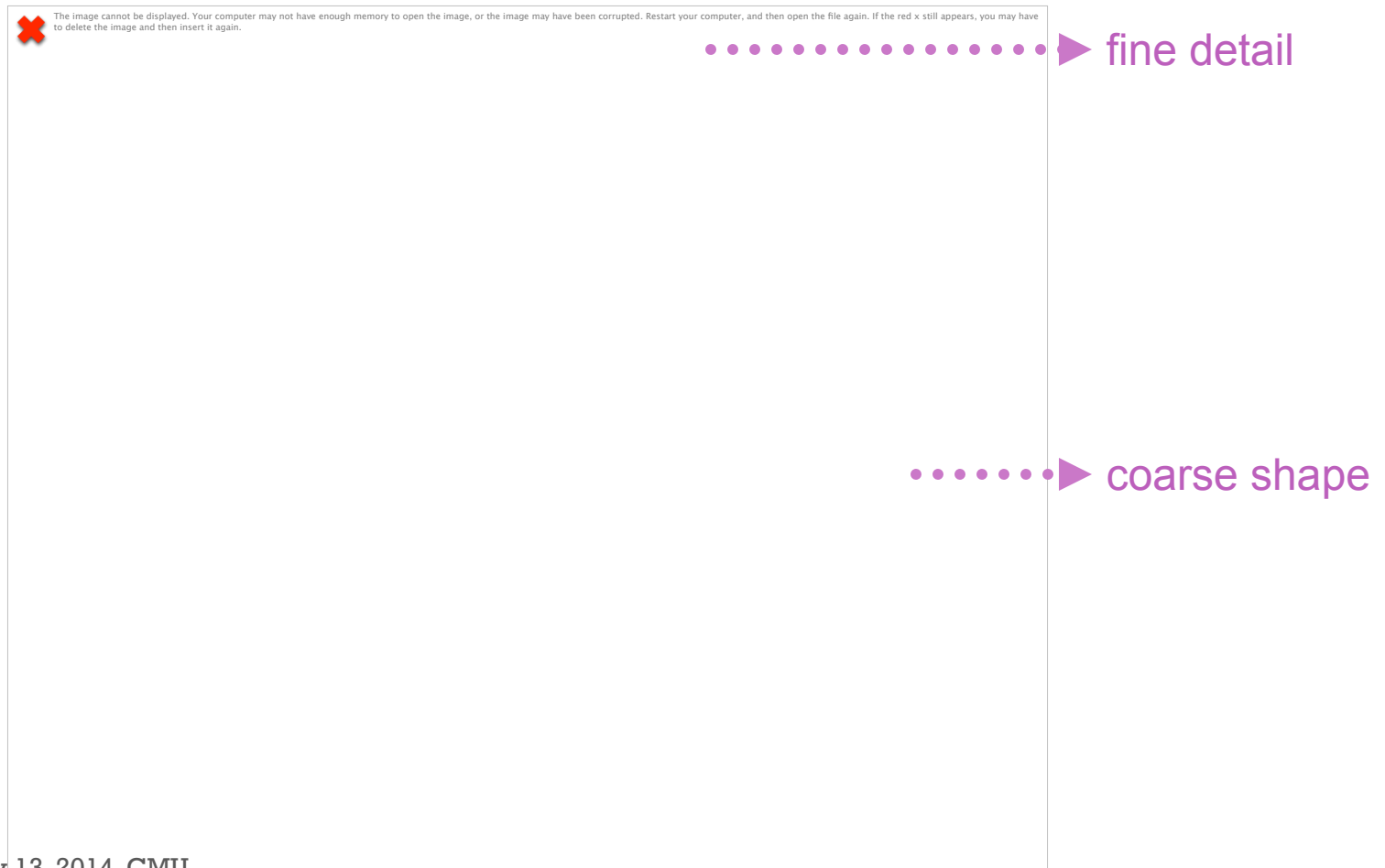
kicking

ActionWeb dataset -
467 images collected
from the web

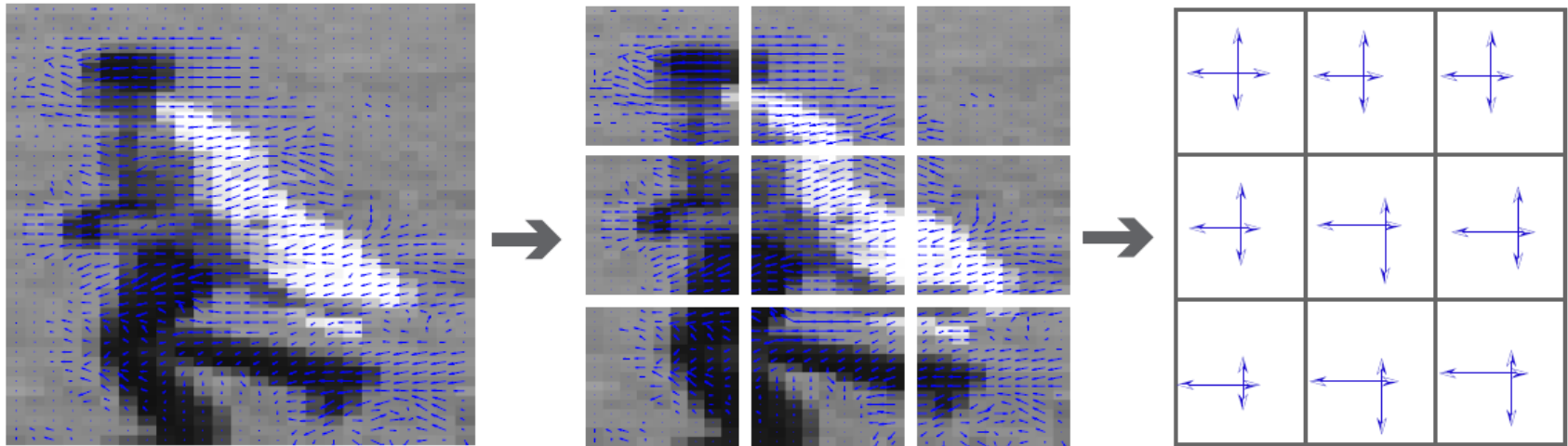
Correctly
classified
action images

+ Boundary-fitted Lines

- In the absence of silhouettes, we can use lines fitted to the boundaries (Pb) (Martin PAMI2004) of human figures

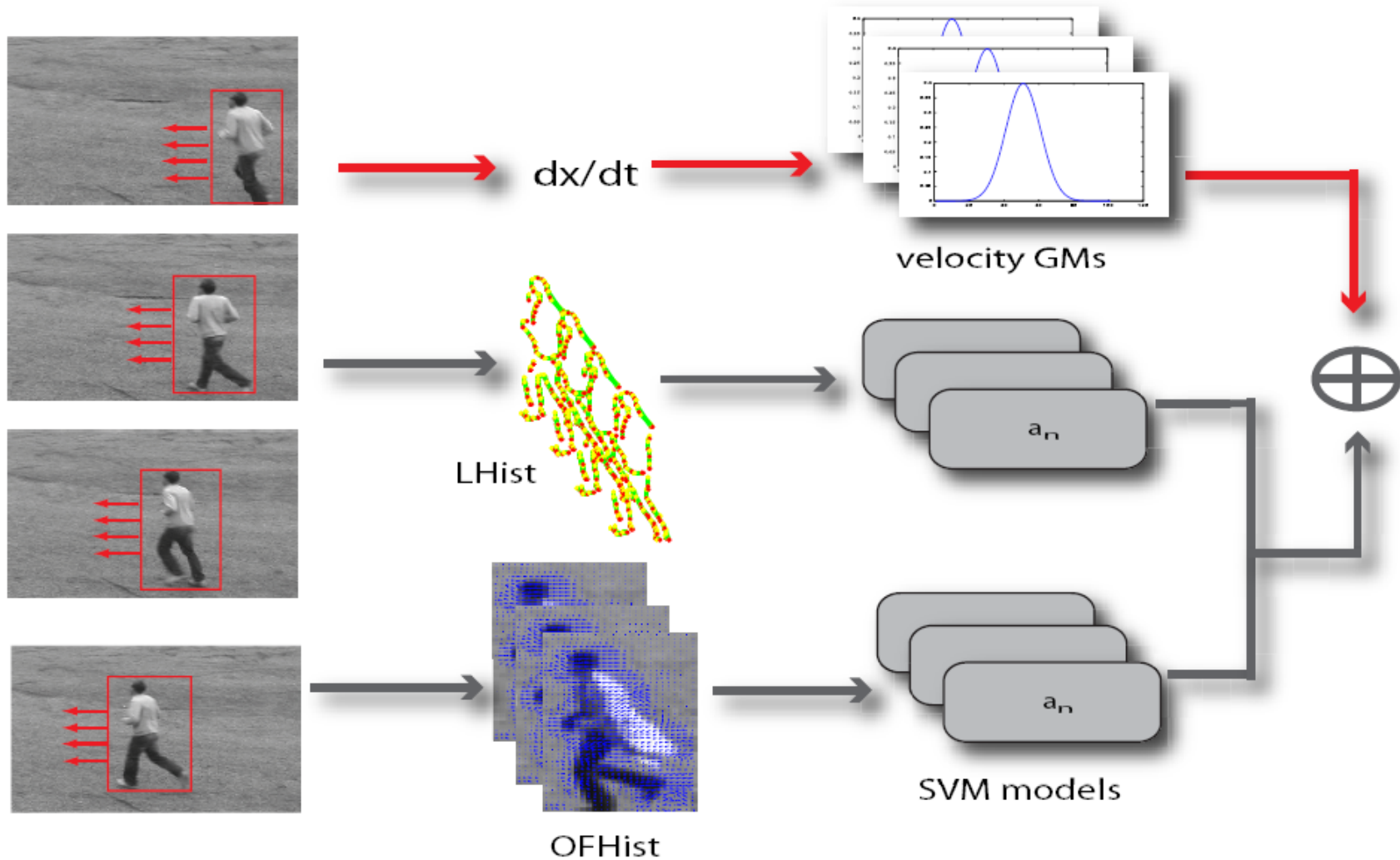


+ ..and Optical Flow

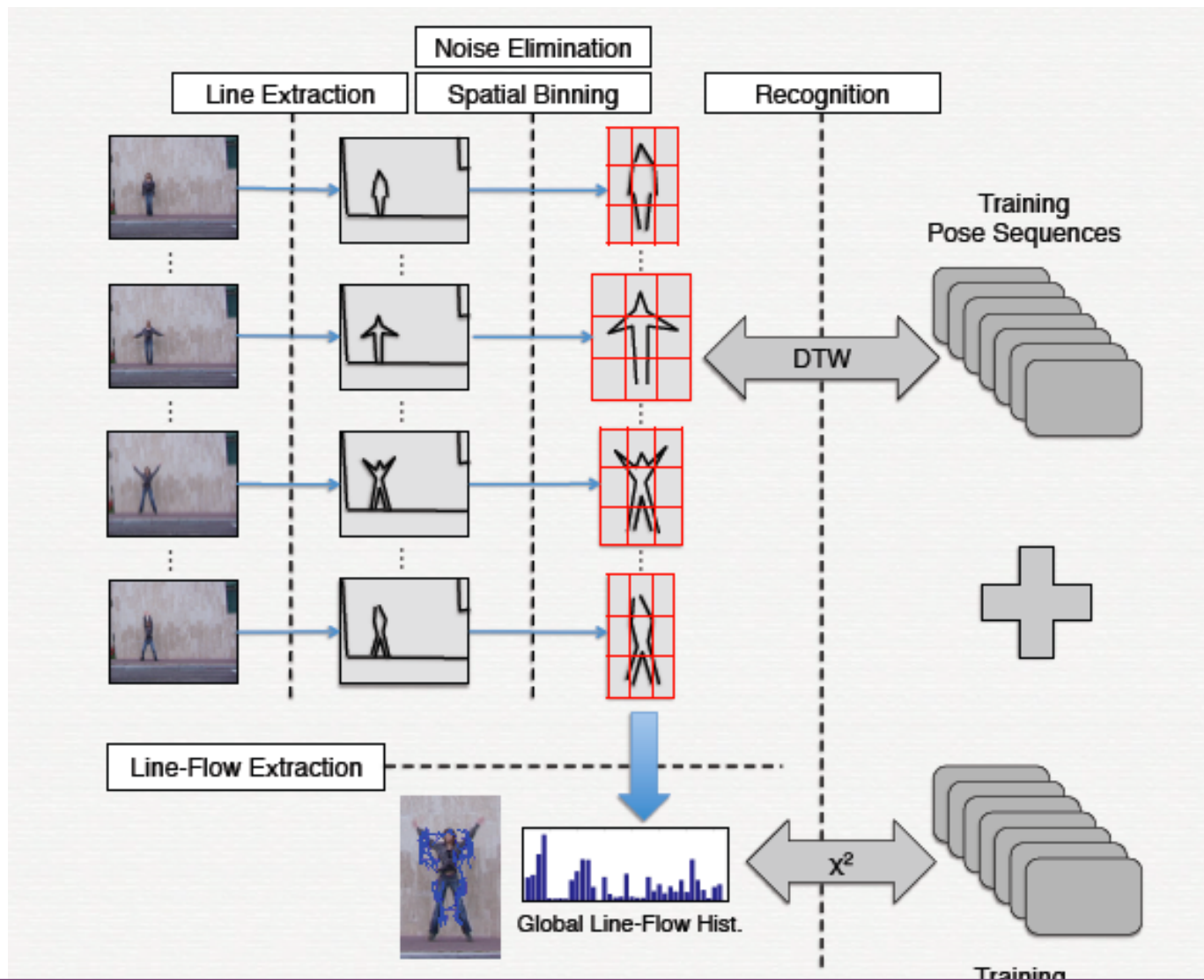


- Dense block-based optical flow calculation
 - L_1 block distance
 - 5x5 template size with a window size of 3

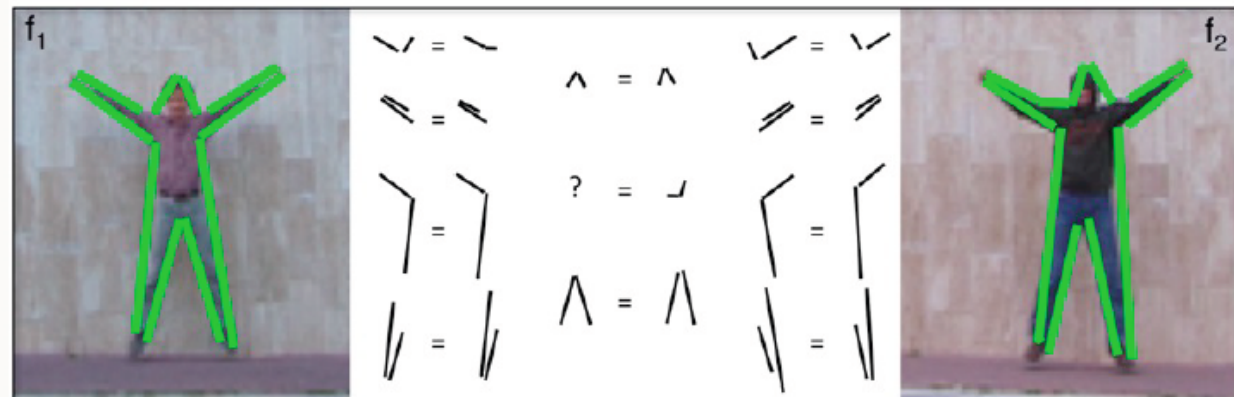
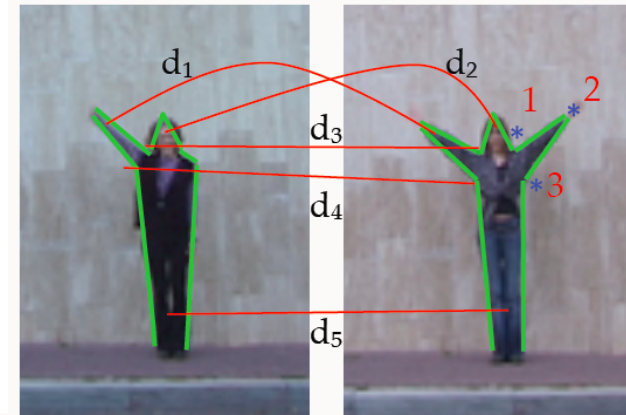
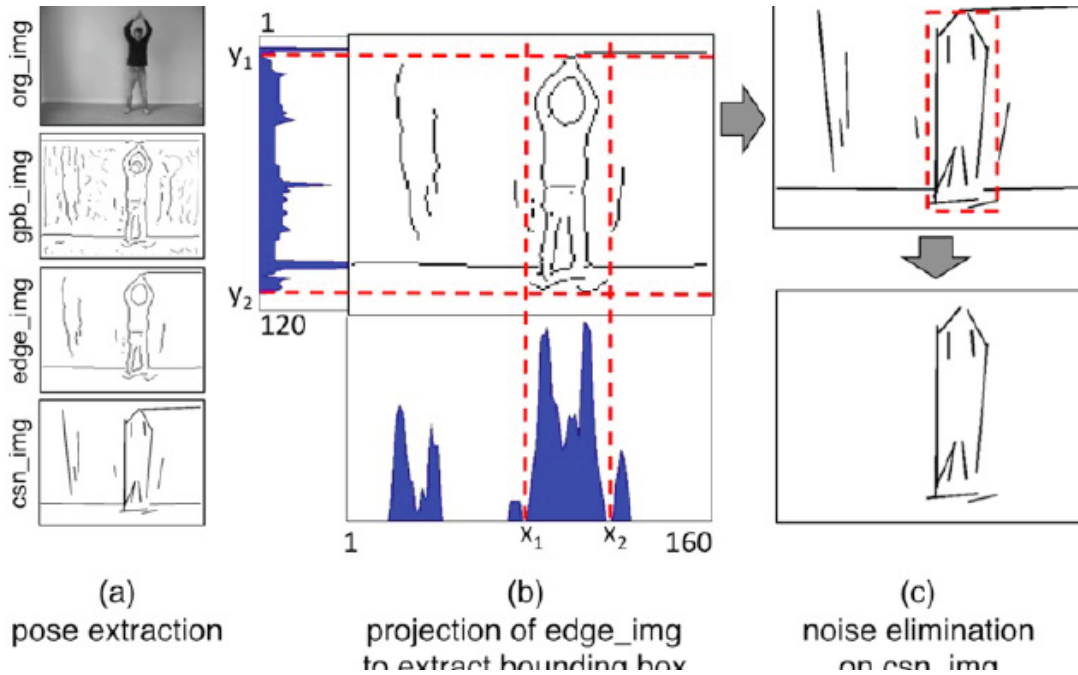
+ Recognition with LHist and OFHist



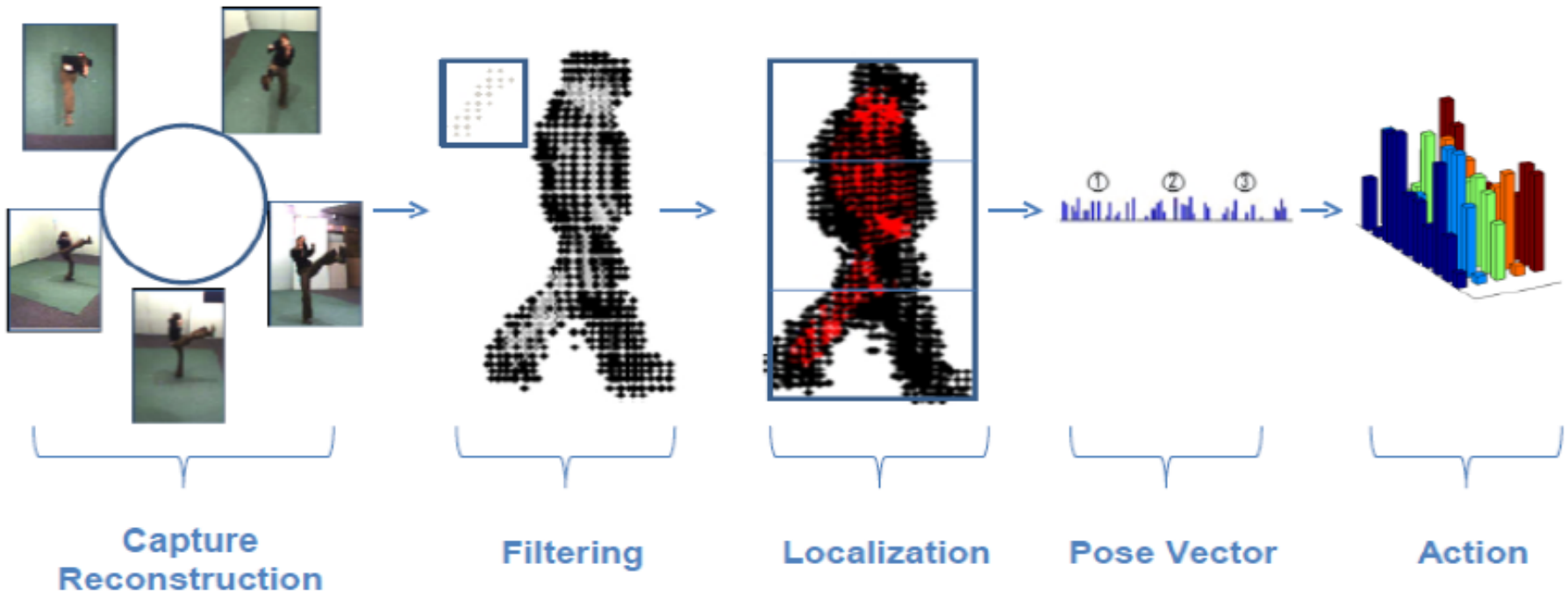
+ Pose as line segments



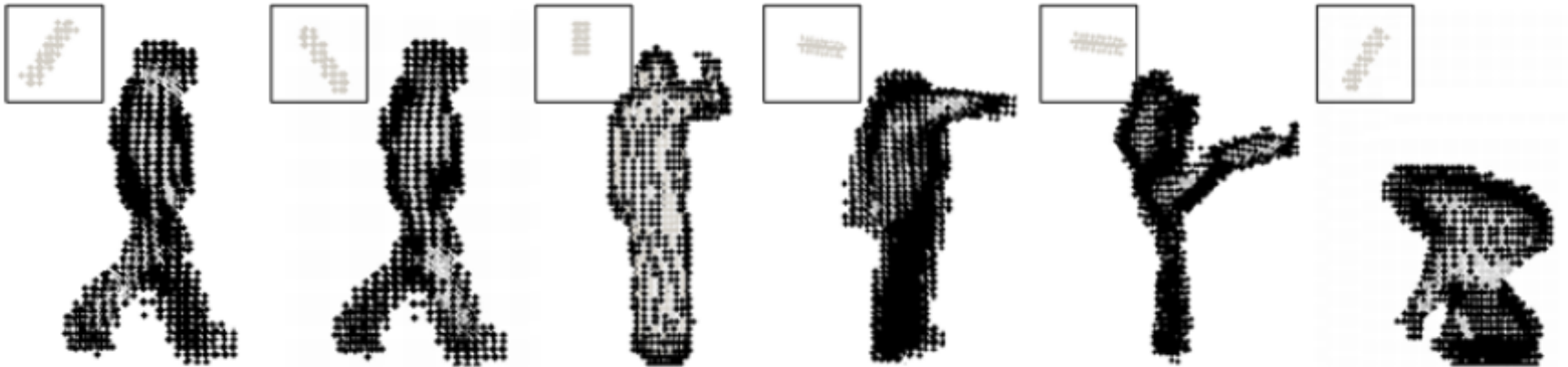
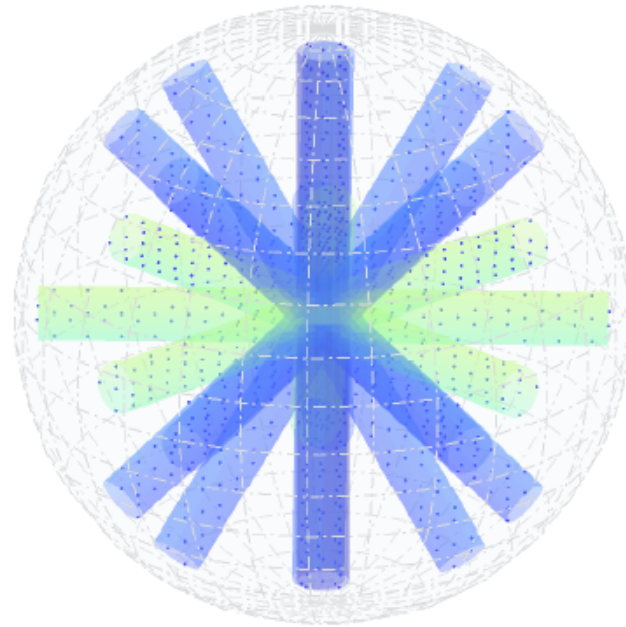
+ Line pairs



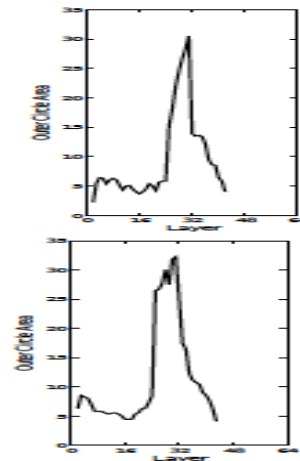
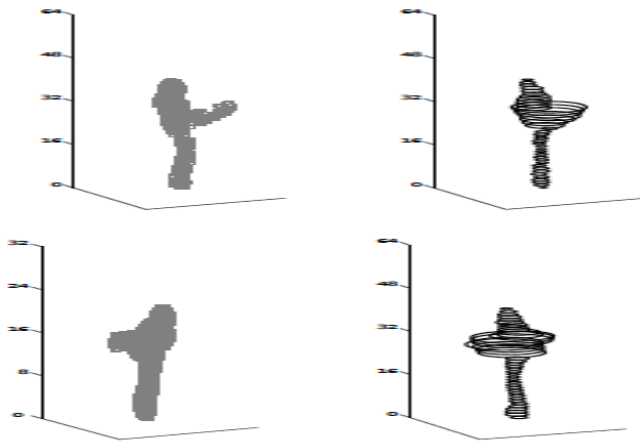
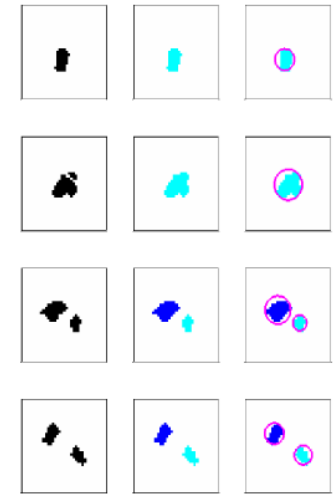
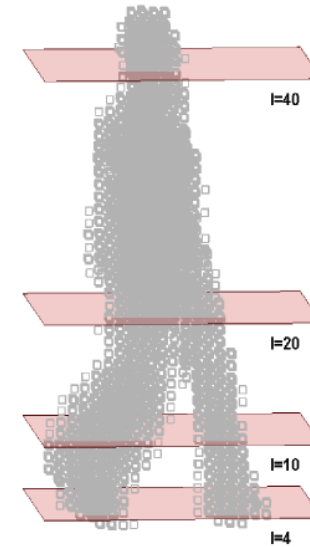
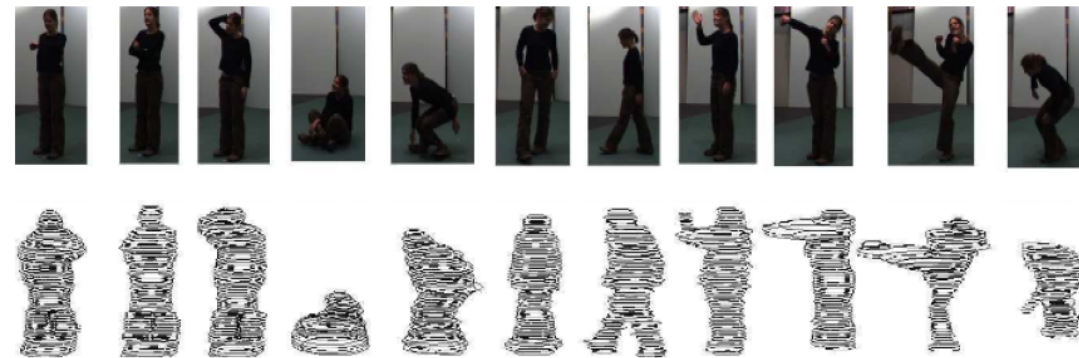
+ Multiple camera views



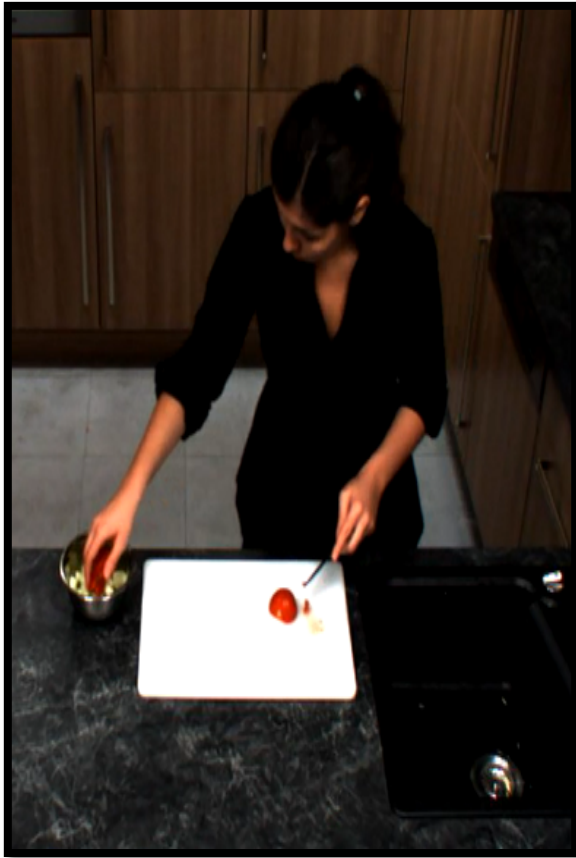
+ Oriented cylinders



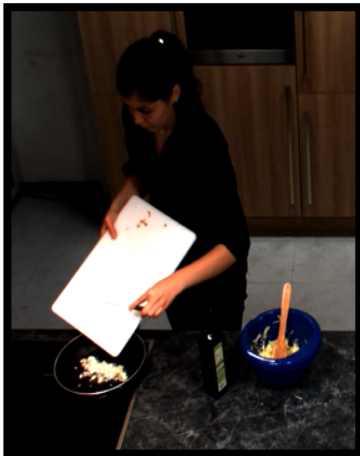
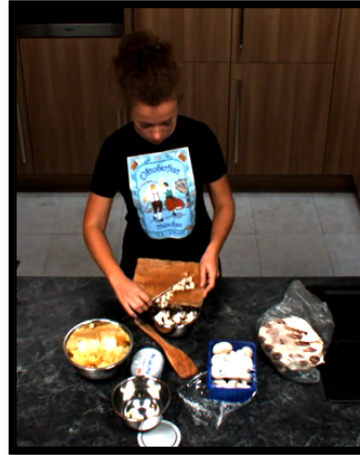
+ Projections as circles



+ Cooking Activities: High Intra-class Variance



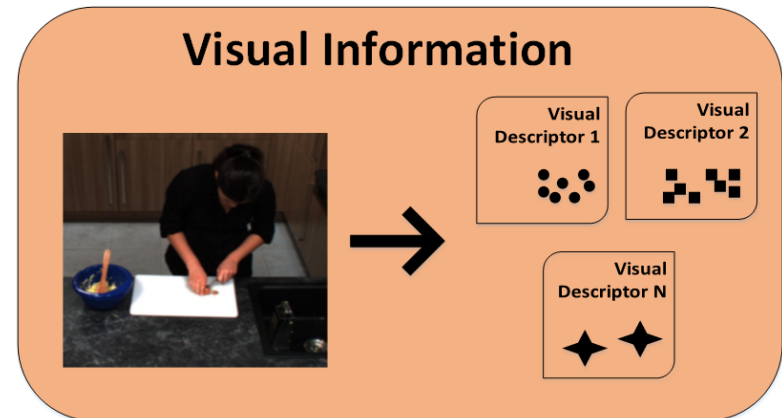
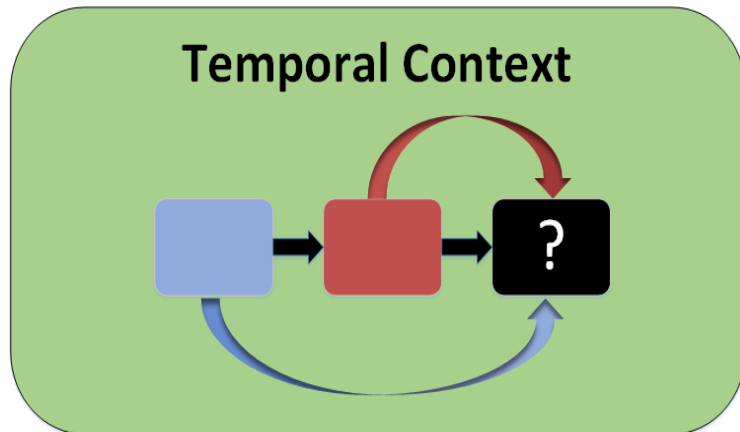
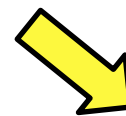
+ Low Inter-class Variance



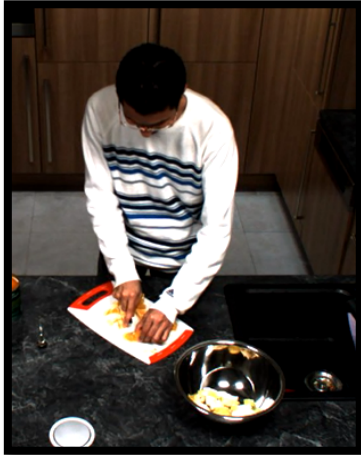
+ Solution

$$y = \underset{i}{\operatorname{argmax}} P(c_i | x)$$

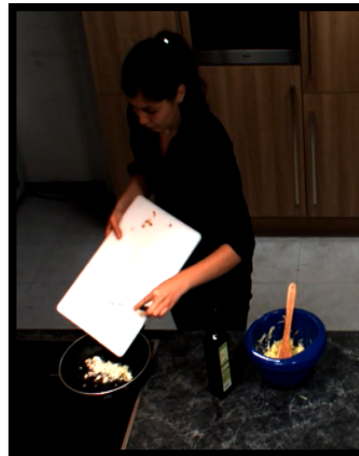
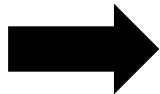
$$P(c_i | x) = T(c_i) \cdot A(c_i, x)$$



+ Put in Pan or Put in Bowl?



**P("put in bowl" | "cut dice") >
P("put in pan" | "cut dice")**



**P("put in pan" | "spread") >
P("put in bowl" | "spread")**

+ Asthma Inhaler use

Correct steps to use an inhaler:

Shake the inhaler (for 5 second)

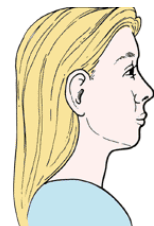
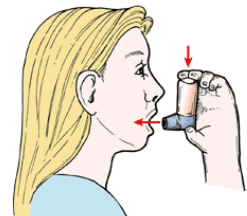
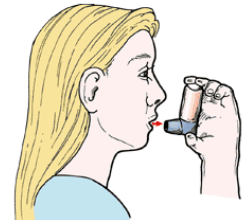
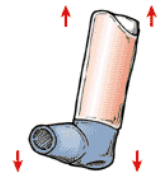
Breathe out

Put the inhaler about 2 inches in front of your mouth

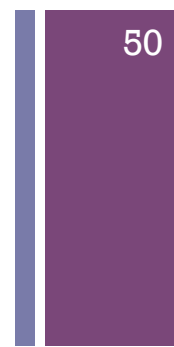
Breathe in and push down the button at the same time

Hold your breath for 10 seconds

Breathe out slowly

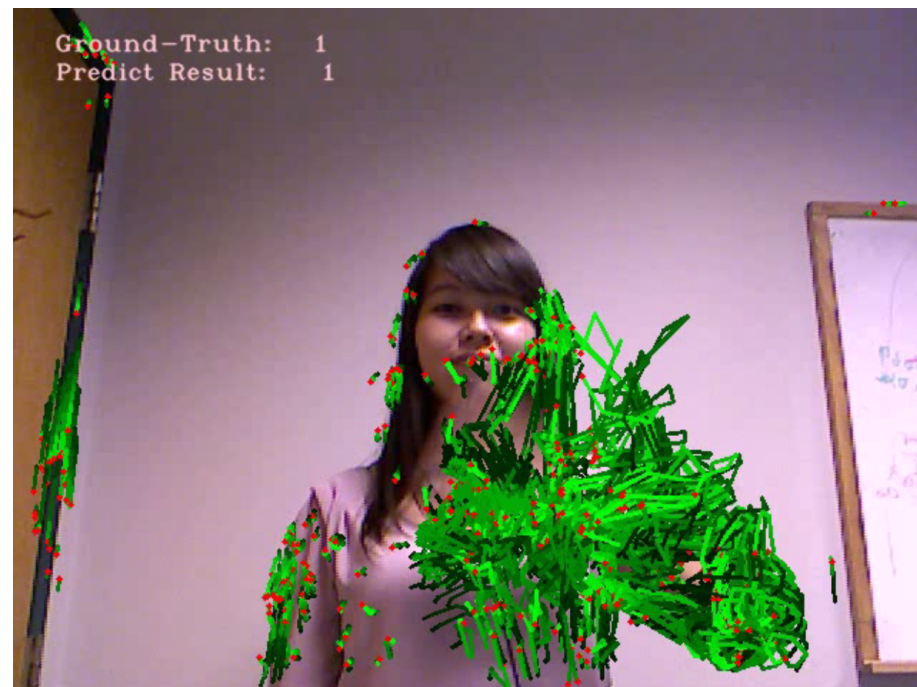


Ground-Truth: 0
Predict Result: 0



Reaching mouth

Ground-Truth: 1
Predict Result: 1



shaking

+ Contributors

- Ahmet Iscen
- Eren Golge
- Anil Armagan
- Sermetcan Baysal
- Fadime Sener
- Hilal Zitouni
- Sare Gul Sevil
- Selen Pehlivan
- Gokberk Cinbis
- Derya Ozkan
- Nazli Ikizler





Experimental Evaluation of HOR

Method	Accuracy
HOR	100%
Blank et al. [12]	99.64%
Jhuang et al. [48]	98.8%
Wang et al. [96]	97.78%
Niebles et al. [63]	72.8%

Comparison to other methods on the Weizzman dataset

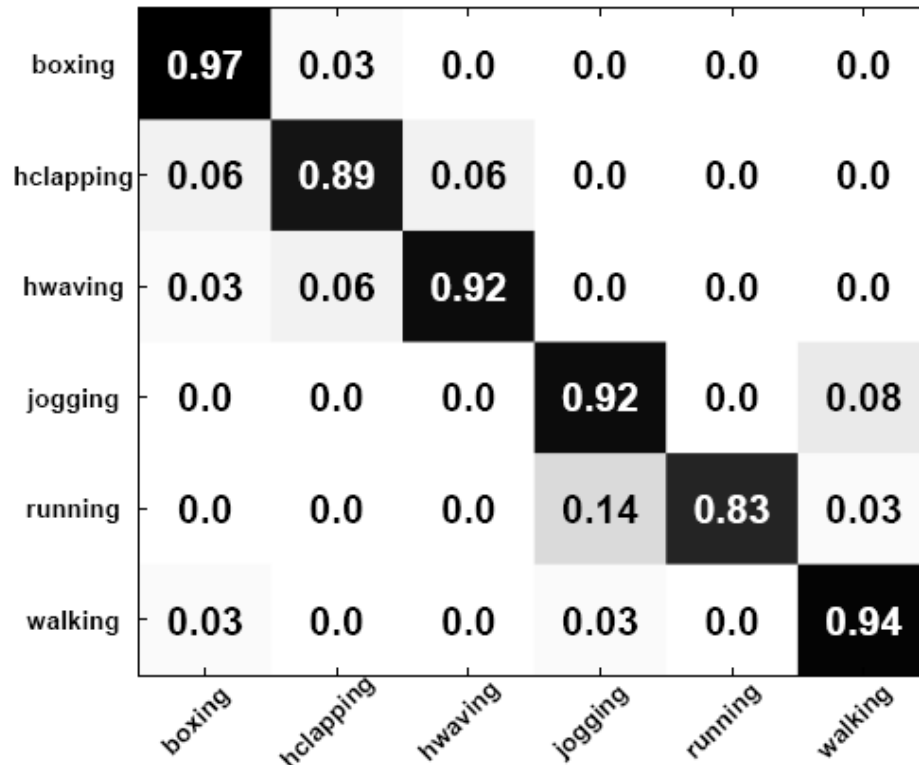
Method	Accuracy
Jhuang et al. [48]	91.7%
Wong et al. [100]	91.6%
HORW	89.4%
Niebles et al. [64]	81.5%
Dollár et al. [24]	81.2%
Ke et al. [50]	80.9%
Schuldt et al. [84]	71.7%

Comparison to other methods on the KTH dataset

Comparison to HOGs on the KTH

	HOG	HOR	HORW
SVM	76.85%	77.31%	85.65%
DTW	67.59%	74.54%	78.24%
v+SVM	82.41%	81.48%	89.35%

+ Line and Flow Results



- Shape and flow are complimentary to each other.
- Again, this depends on the nature of the actions in mention.

Method	Accuracy
LFV	94.0%
Jhuang [48]	91.7%
Wong [100]	91.6%
Niebles [64]	81.5%
Dollár [24]	81.2%
Ke [50]	80.9%
Schuldt [84]	71.7%

Condition	LFV	Jhuang [48]
s1	98.2%	96.0%
s2	90.7%	86.1%
s3	88.9%	89.8%
s4	98.2%	94.8%

Still Image Results



(a) catch, walk, catch, throw

(b) run, run, run, kick



(c) catch, kick, walk, crouch



(d) run, throw, run, run



(e) kick, walk, walk, catch



(f) throw, walk, run, throw

running	0.83	0.04	0.04	0.05	0.04	0.0
walking	0.04	0.94	0.0	0.0	0.01	0.01
throwing	0.0	0.07	0.85	0.01	0.03	0.04
catching	0.15	0.04	0.04	0.72	0.0	0.06
crouching	0.04	0.03	0.01	0.01	0.89	0.01
kicking	0.03	0.03	0.04	0.03	0.0	0.87
	running	walking	throwing	catching	crouching	kicking

Total accuracy 85.1%

Misclassified action images

+ Multi-view

Method	Accuracy (over 11 actions)	Accuracy (over 13 actions)
Weinland et al. [61]	93.33%	-
<i>Our method</i>	90.91%	88.63%
Liu et al. [32]	-	82.8%
Weinland et al. [59]	81.27%	-
Lv et al. [35]	-	80.6%
Yan et al. [62]	78.0%	-