

Action Recognition



Computer Vision

Jia-Bin Huang, Virginia Tech

This section: advanced topics

- Convolutional neural networks in vision
- Action recognition
- 3D Scenes and Context

What is an action?



Action: a transition from one state to another

- Who is the actor?
- How is the state of the actor changing?
- What (if anything) is being acted on?
- How is that thing changing?
- What is the purpose of the action (if any)?

How do we represent actions?

Categories

Walking, hammering, dancing, skiing, sitting down, standing up, jumping

Poses



Nouns and Predicates

<man, swings, hammer>

<man, hits, nail, w/ hammer>

What is the purpose of action recognition?

- To describe

<https://www.youtube.com/watch?v=bcgXAQcvxdc>

- To predict

<http://www.youtube.com/watch?v=LQm25nW6aZw>

How can we identify actions?

Motion



Pose



Held
Objects



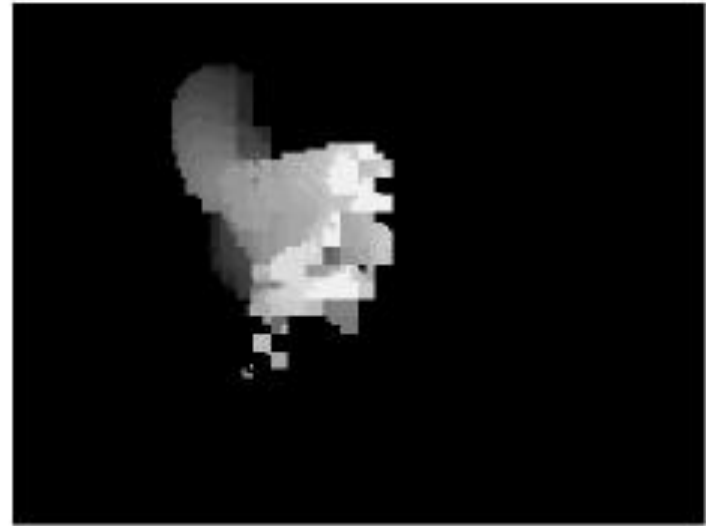
Nearby
Objects

Representing Motion

Optical Flow with Motion History



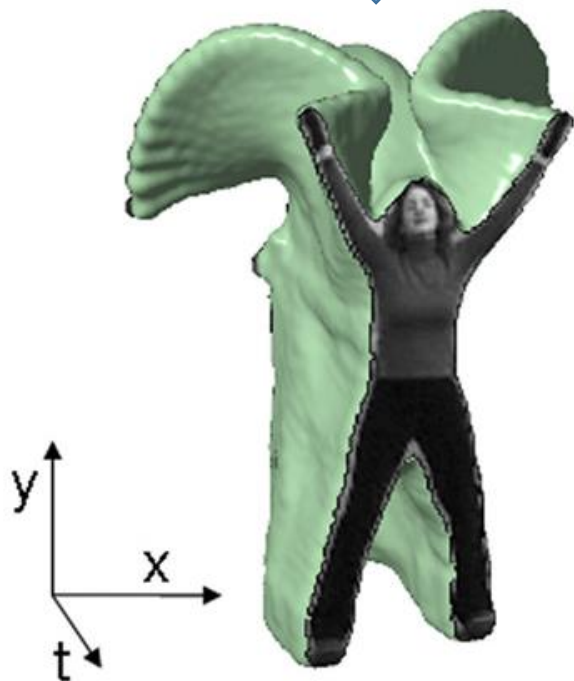
sit-down



sit-down MHI

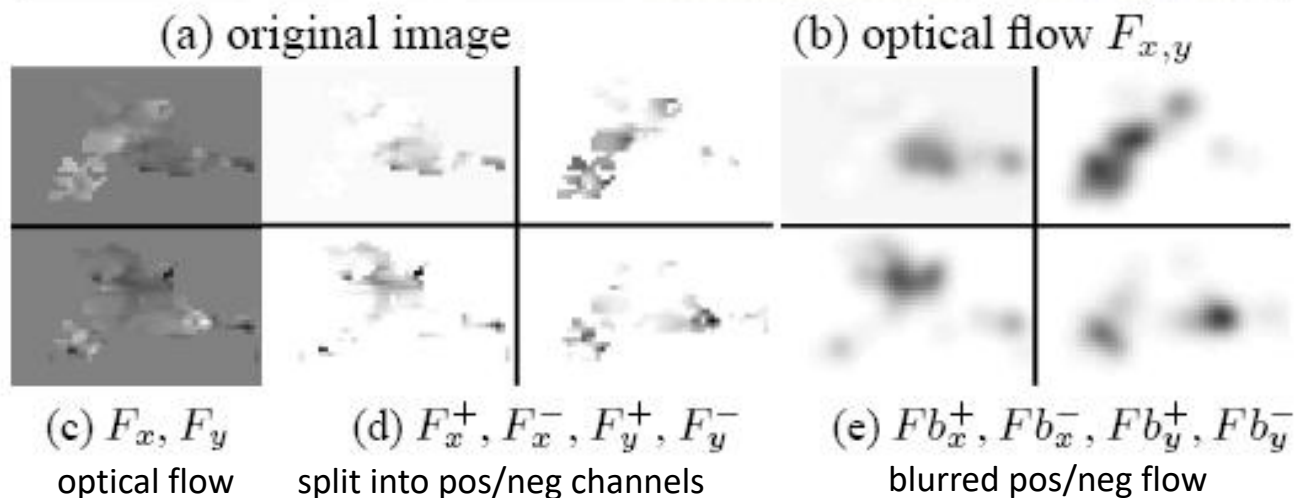
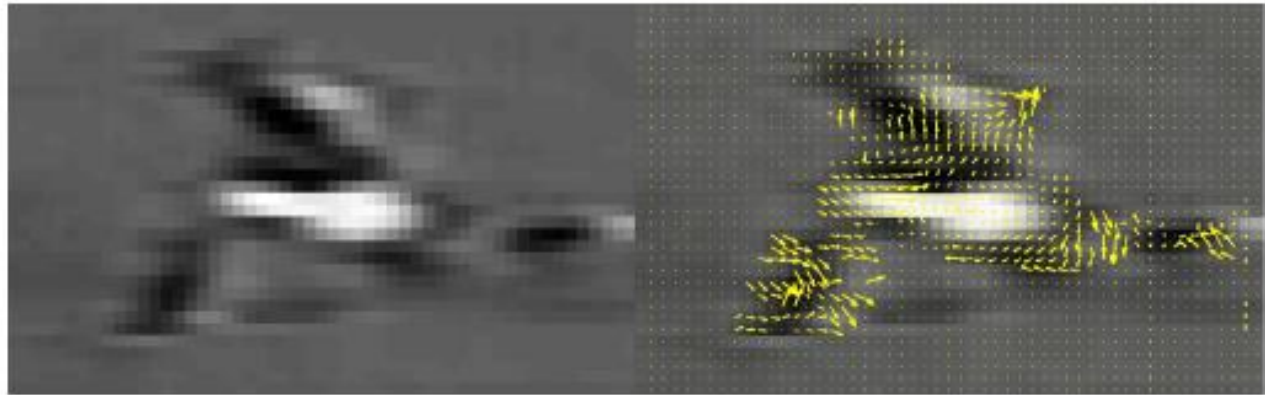
Representing Motion

Space-Time Volumes



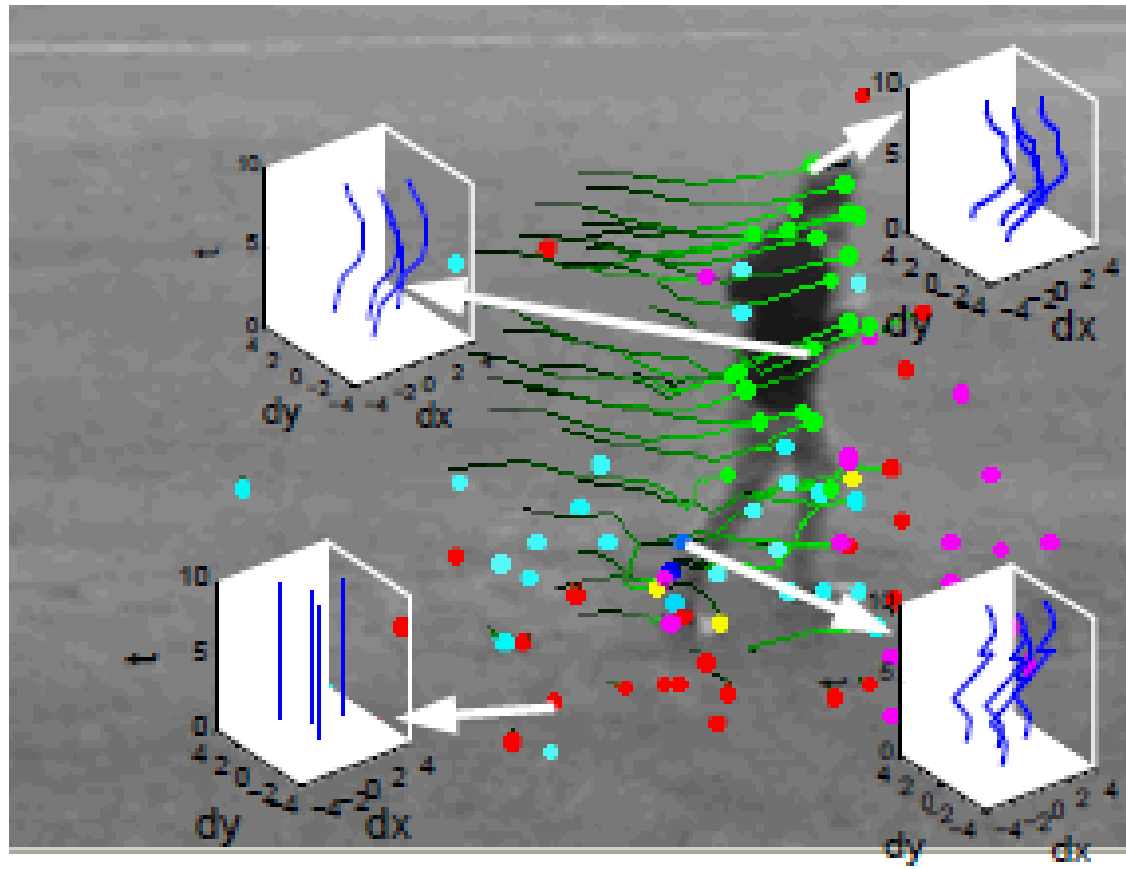
Representing Motion

Optical Flow with Split Channels



Representing Motion

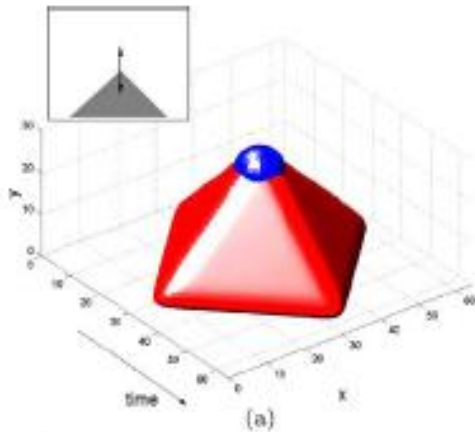
Tracked Points



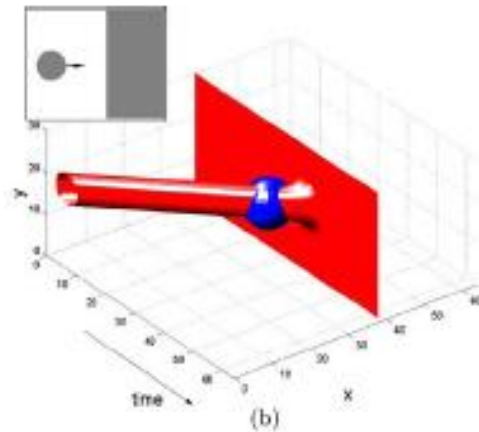
Representing Motion

Space-Time Interest Points

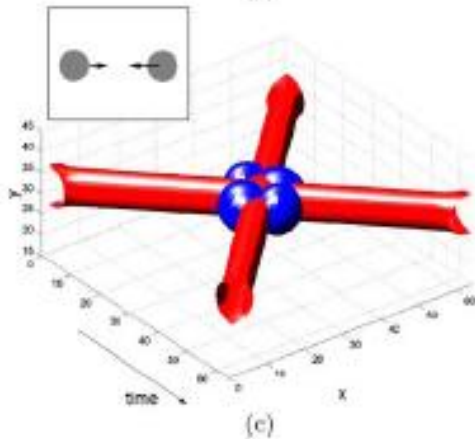
Moving corner



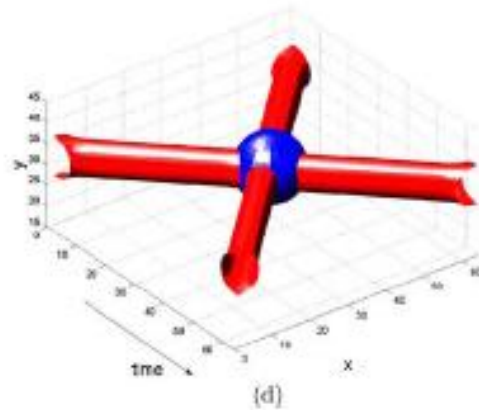
Ball hits wall



Corner detectors in
space-time



Balls collide

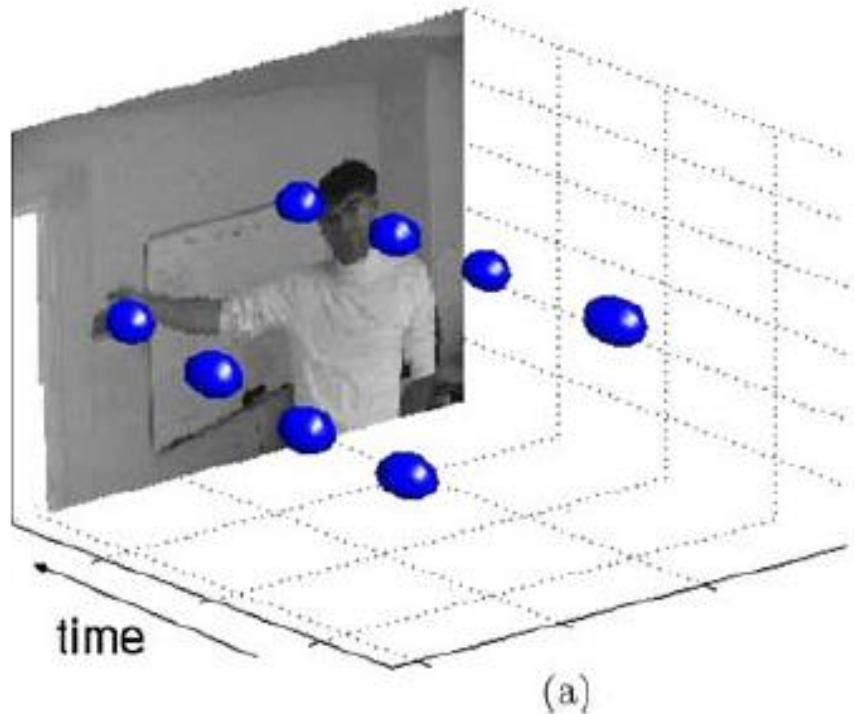


Balls collide (different scale)

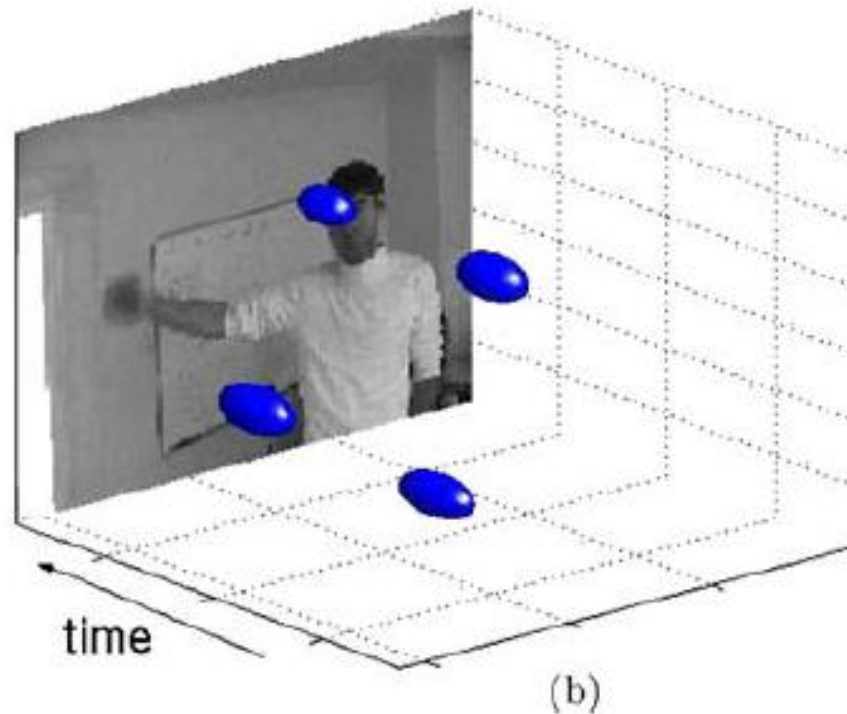
Representing Motion

Space-Time Interest Points

Hand waves with high frequency



Hand waves with low frequency



Examples of Action Recognition Systems

- Feature-based classification
- Recognition using pose and objects

Action recognition as classification

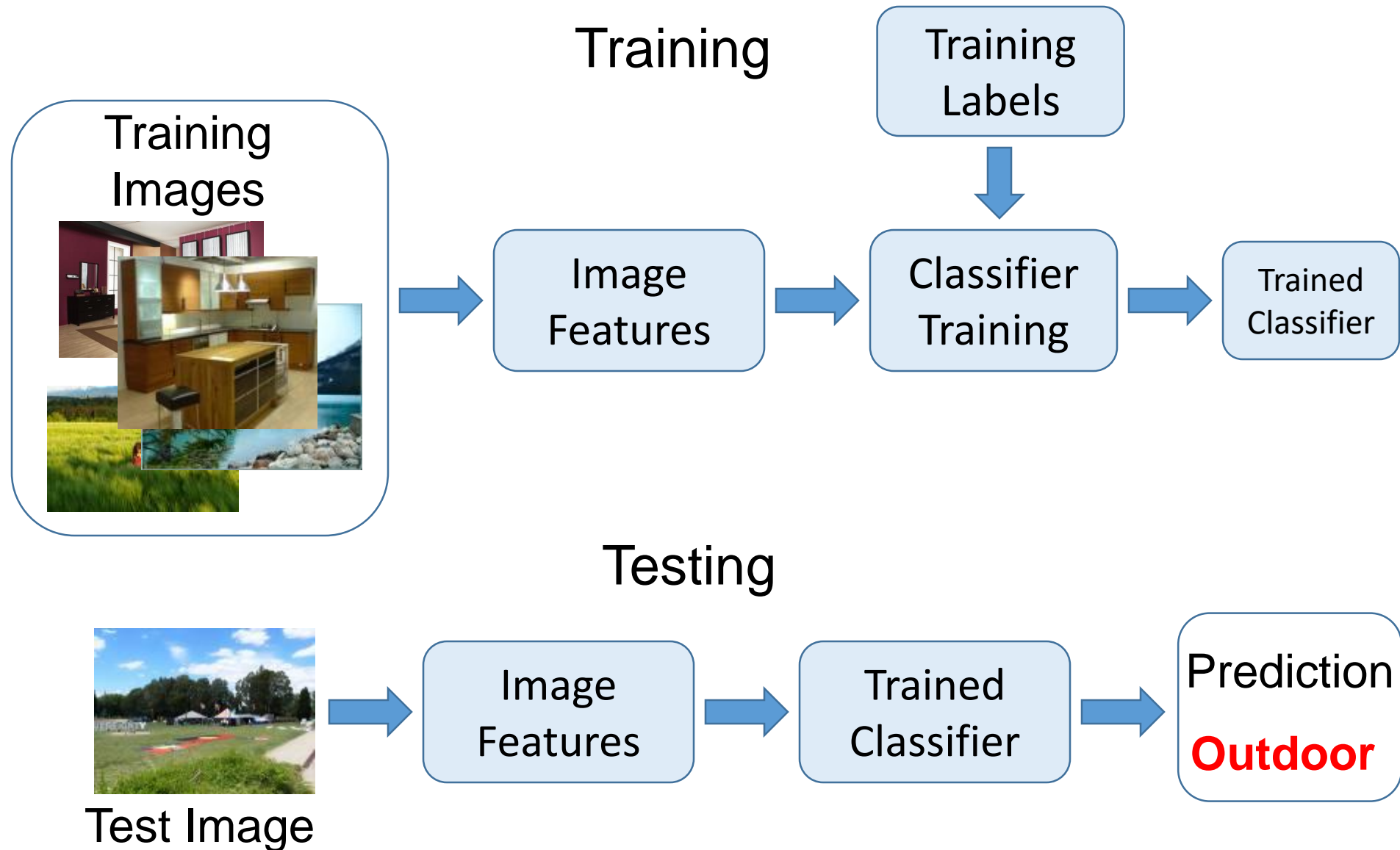
training samples



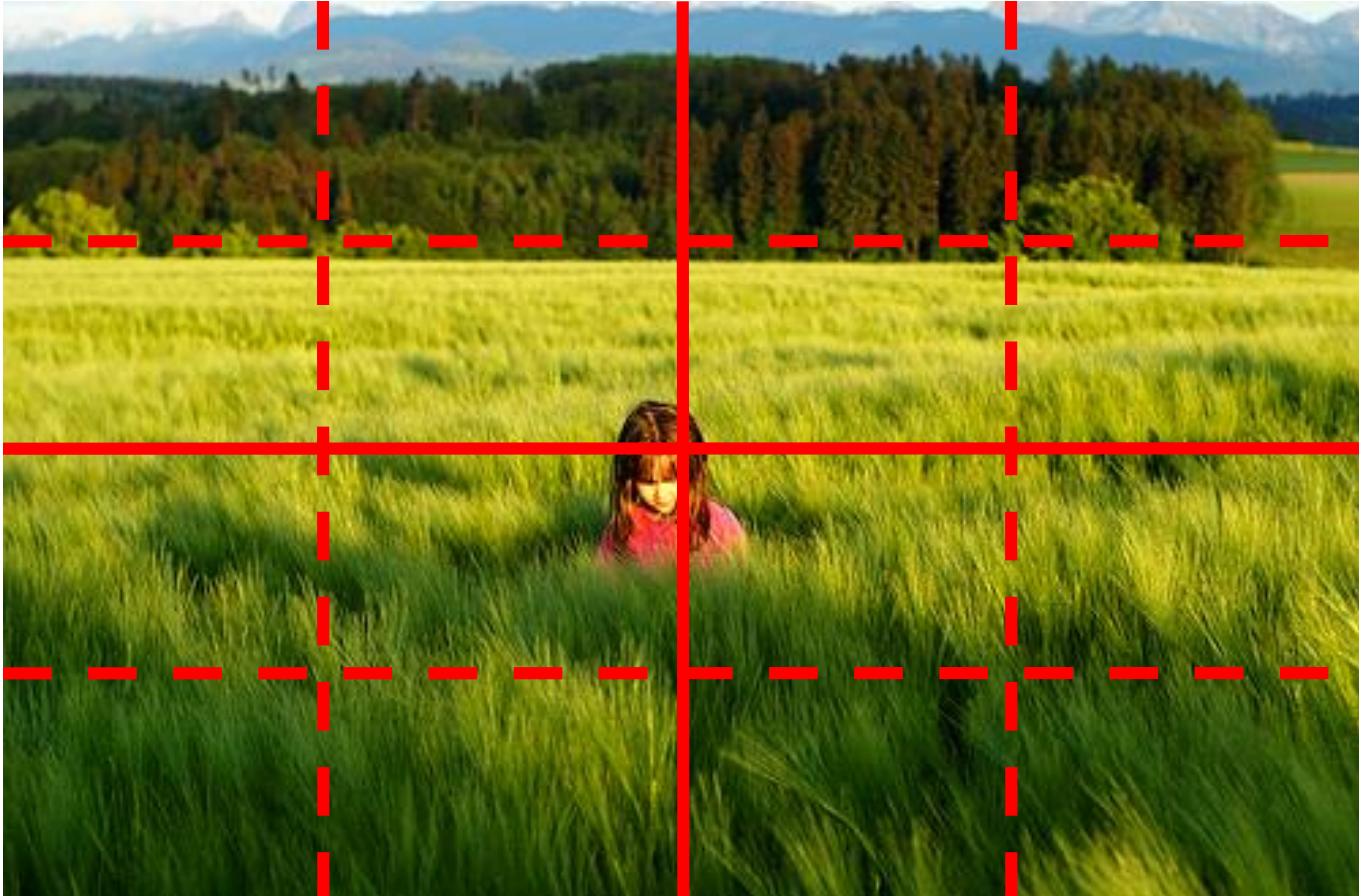
test samples



Remember image categorization...



Remember spatial pyramids....

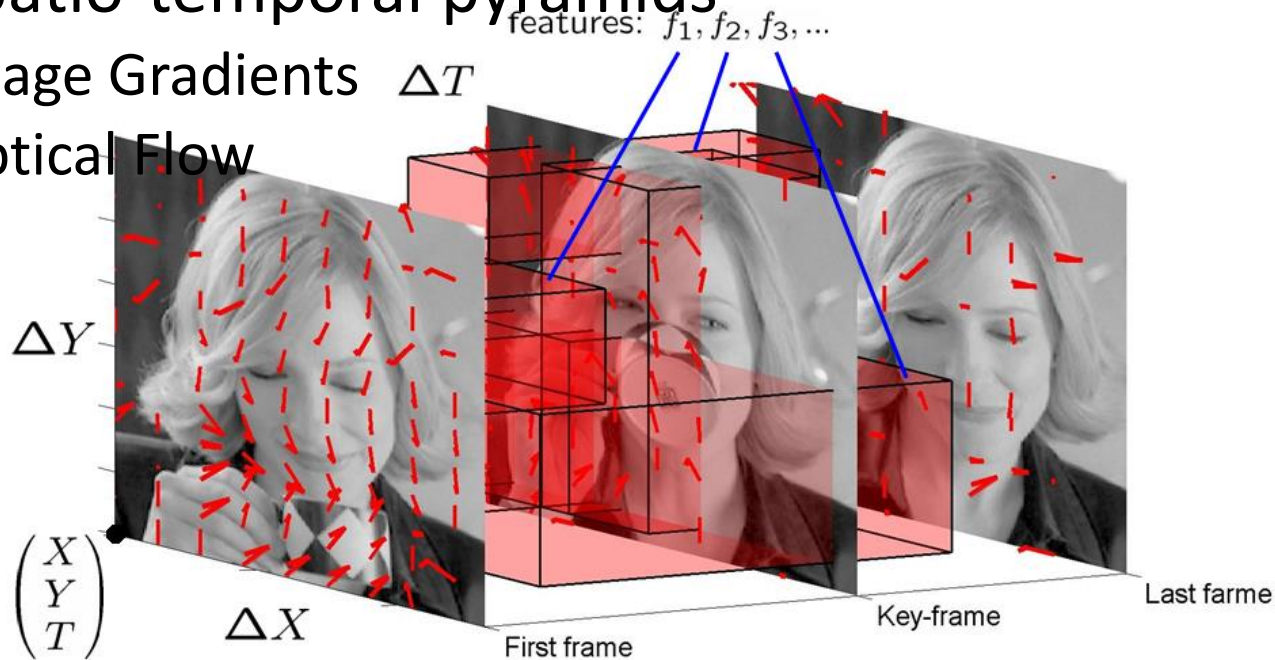


Compute histogram in each spatial bin

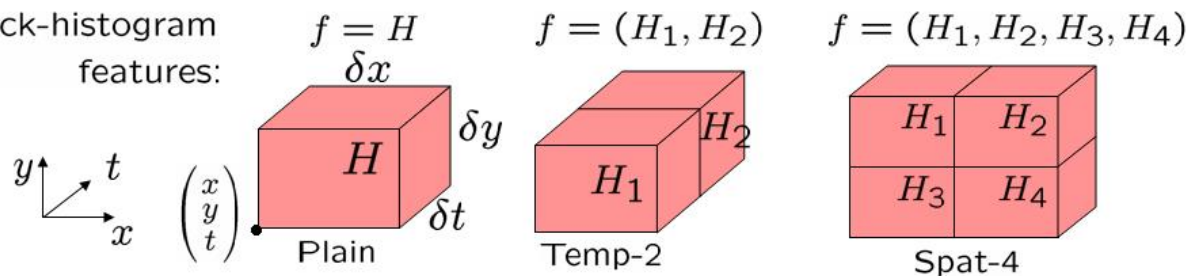
Features for Classifying Actions

1. Spatio-temporal pyramids

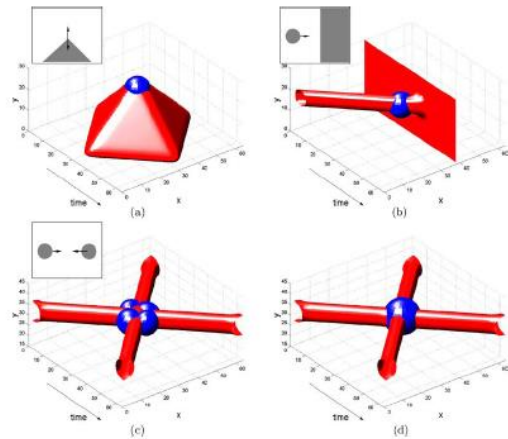
- Image Gradients ΔT
- Optical Flow



block-histogram
features:

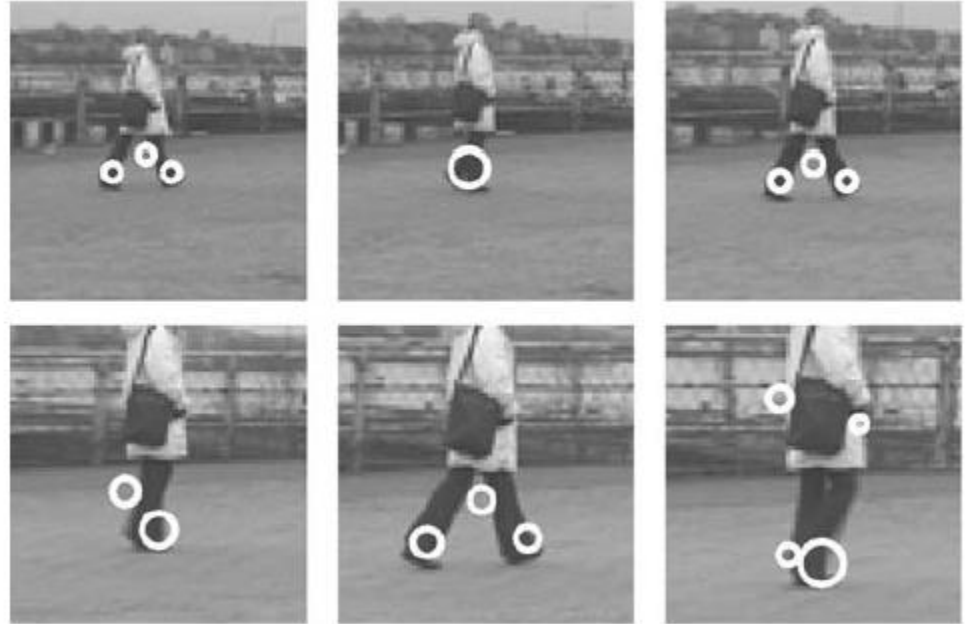


Features for Classifying Actions



Corner detectors in
space-time

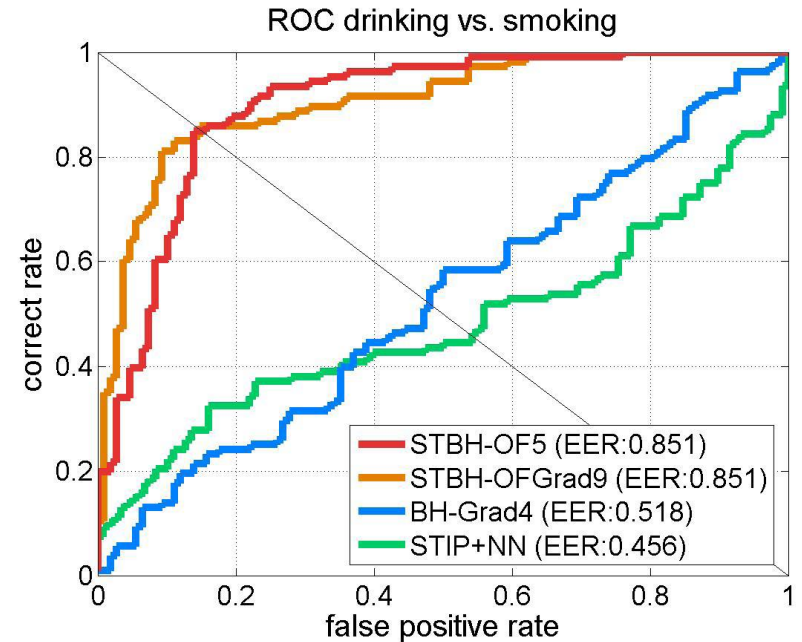
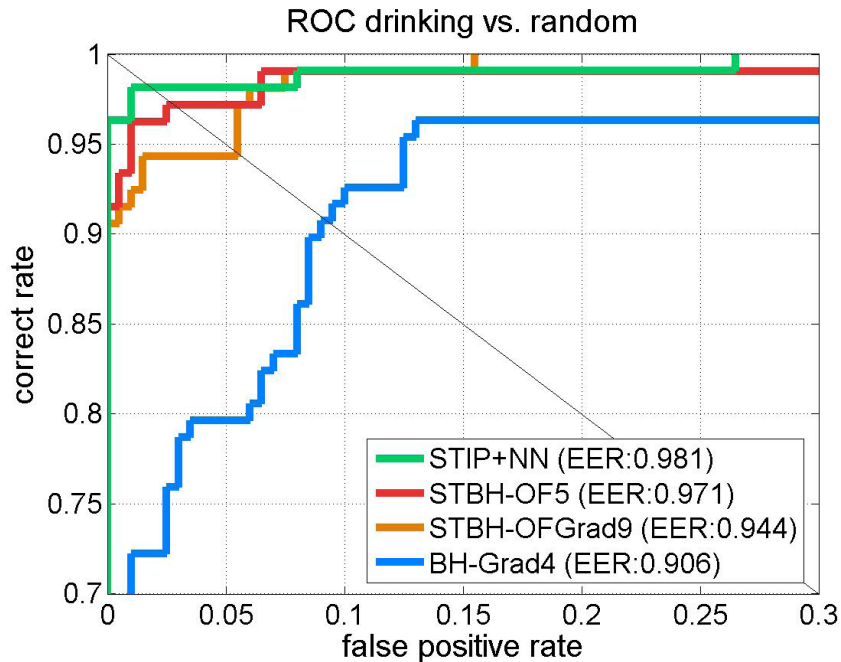
I in



Descriptors based on Gaussian derivative filters over x, y, time

Classification

- Boosted stubs for pyramids of optical flow, gradient
- Nearest neighbor for STIP



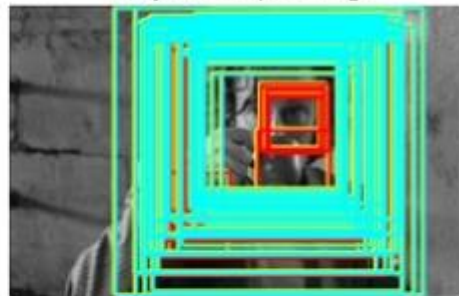
Searching the video for an action

1. Detect keyframes using a trained HOG detector in each frame
2. Classify detected keyframes as positive (e.g., “drinking”) or negative (“other”)

Test frame samples



Keyframe priming

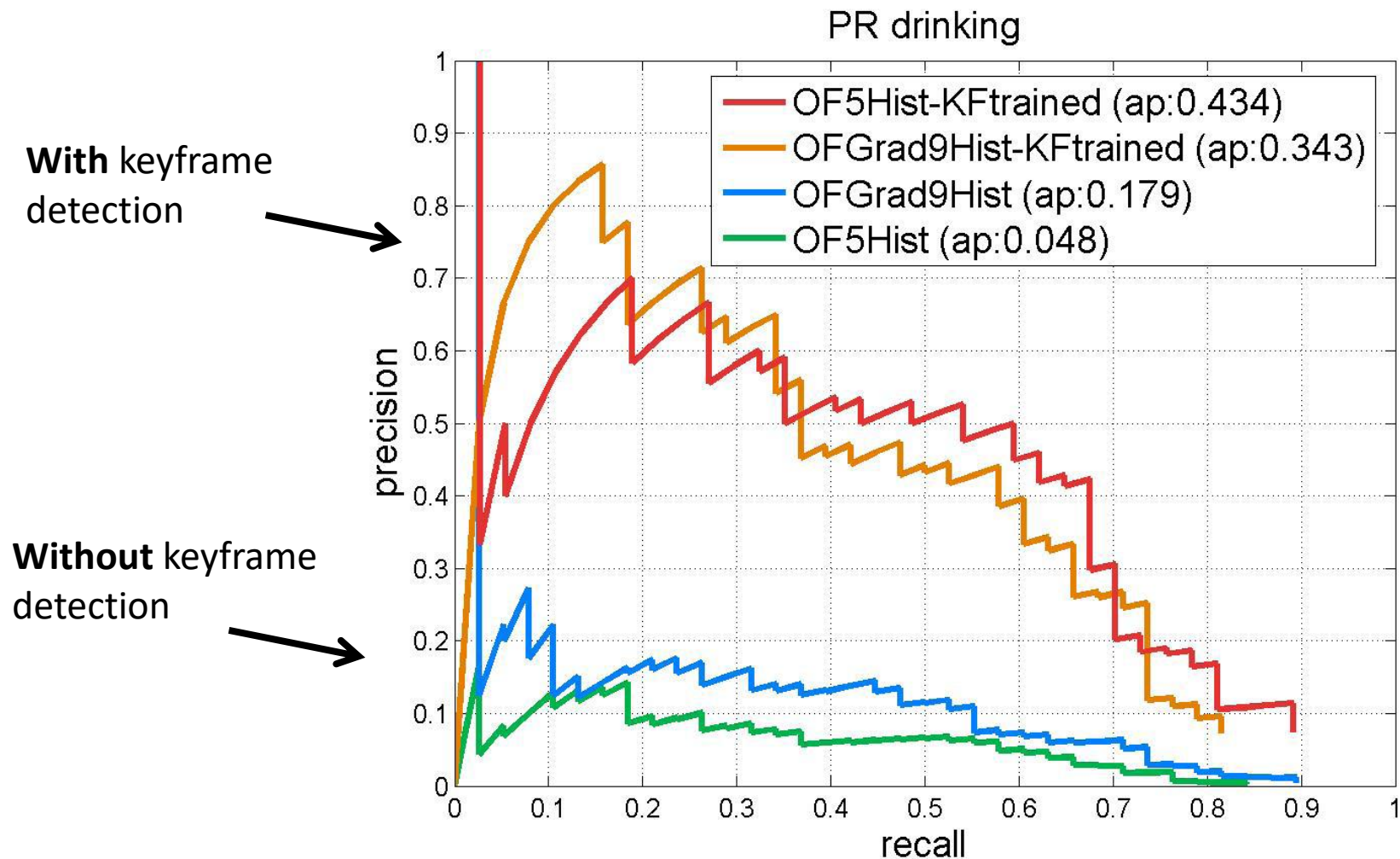


Keyframe-primed
event detection



Keyframe
detections

Accuracy in searching video





“Talk on phone”



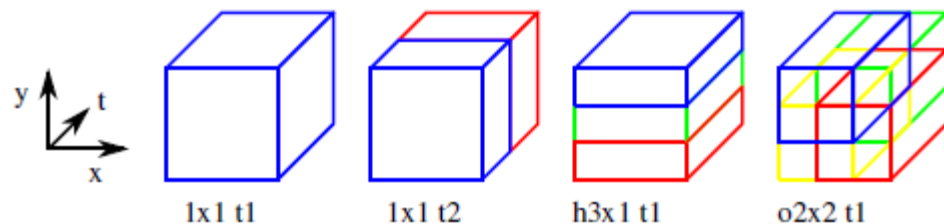
“Get out of car”

Approach

- Space-time interest point detectors
- Descriptors
 - HOG, HOF
- Pyramid histograms (3x3x2)
- SVMs with Chi-Squared Kernel

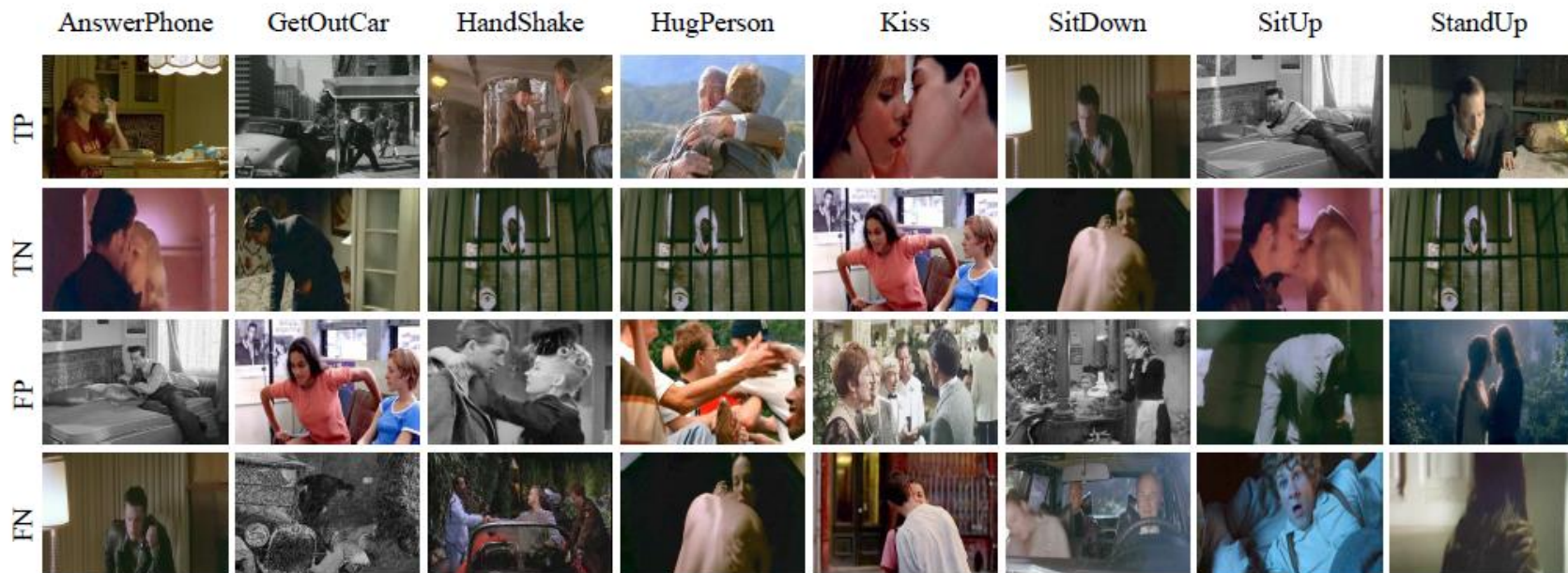


Interest Points



Spatio-Temporal Binning

Results



Task	HoG BoF	HoF BoF	Best channel	Best combination
KTH multi-class	81.6%	89.7%	91.1% (hof h3x1 t3)	91.8% (hof 1 t2, hog 1 t3)
Action AnswerPhone	13.4%	24.6%	26.7% (hof h3x1 t3)	32.1% (hof o2x2 t1, hof h3x1 t3)
Action GetOutCar	21.9%	14.9%	22.5% (hof o2x2 1)	41.5% (hof o2x2 t1, hog h3x1 t1)
Action HandShake	18.6%	12.1%	23.7% (hog h3x1 1)	32.3% (hog h3x1 t1, hog o2x2 t3)
Action HugPerson	29.1%	17.4%	34.9% (hog h3x1 t2)	40.6% (hog 1 t2, hog o2x2 t2, hog h3x1 t2)
Action Kiss	52.0%	36.5%	52.0% (hog 1 1)	53.3% (hog 1 t1, hof 1 t1, hof o2x2 t1)
Action SitDown	29.1%	20.7%	37.8% (hog 1 t2)	38.6% (hog 1 t2, hog 1 t3)
Action SitUp	6.5%	5.7%	15.2% (hog h3x1 t2)	18.2% (hog o2x2 t1, hog o2x2 t2, hog h3x1 t2)
Action StandUp	45.4%	40.0%	45.4% (hog 1 1)	50.5% (hog 1 t1, hof 1 t2)

Action Recognition using Pose and Objects



[Modeling Mutual Context of Object and Human Pose in Human-Object Interaction Activities](#), B. Yao and Li Fei-Fei, 2010

Human-Object Interaction

Holistic image based classification



Integrated reasoning

- **Human pose estimation**



Human-Object Interaction

Holistic image based classification



Integrated reasoning

- Human pose estimation
- **Object detection**



Human-Object Interaction

Holistic image based classification



Integrated reasoning

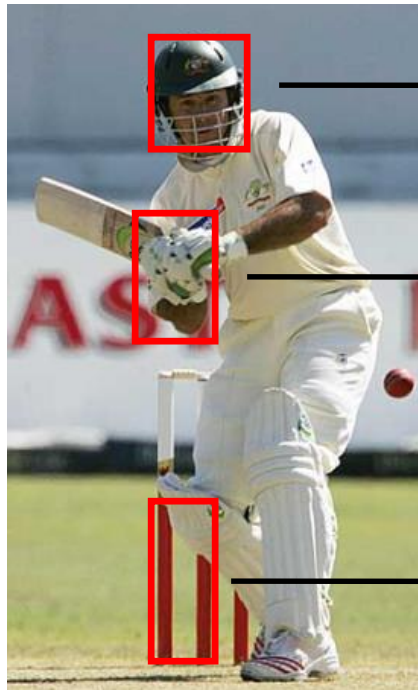
- **Human pose estimation**
- **Object detection**
- **Action categorization**



Activity: Tennis Forehand

Human pose estimation & Object detection

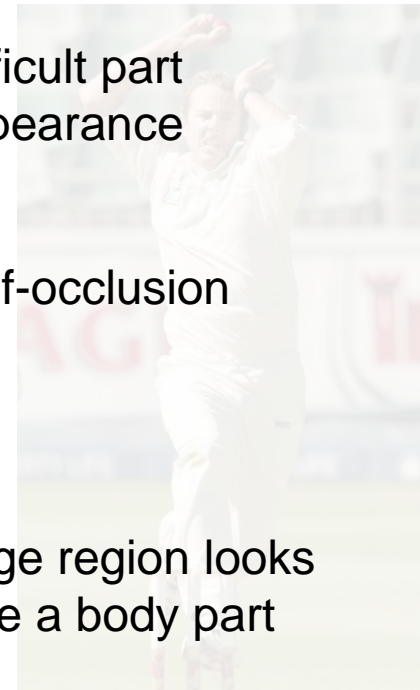
Human pose estimation is challenging.



Difficult part appearance

Self-occlusion

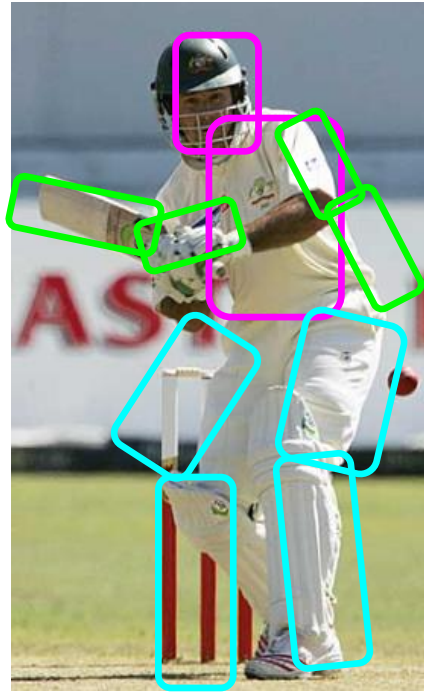
Image region looks like a body part



- Felzenszwalb & Huttenlocher, 2005
- Ren et al, 2005
- Ramanan, 2006
- Ferrari et al, 2008
- Yang & Mori, 2008
- Andriluka et al, 2009
- Eichner & Ferrari, 2009

Human pose estimation & Object detection

Human pose estimation is challenging.

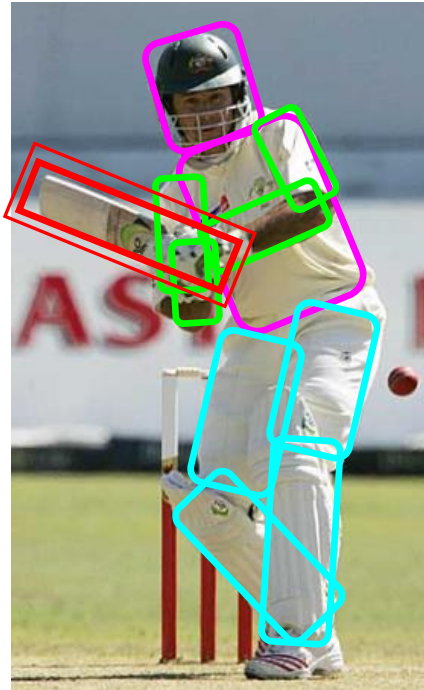


- Felzenszwalb & Huttenlocher, 2005
- Ren et al, 2005
- Ramanan, 2006
- Ferrari et al, 2008
- Yang & Mori, 2008
- Andriluka et al, 2009
- Eichner & Ferrari, 2009

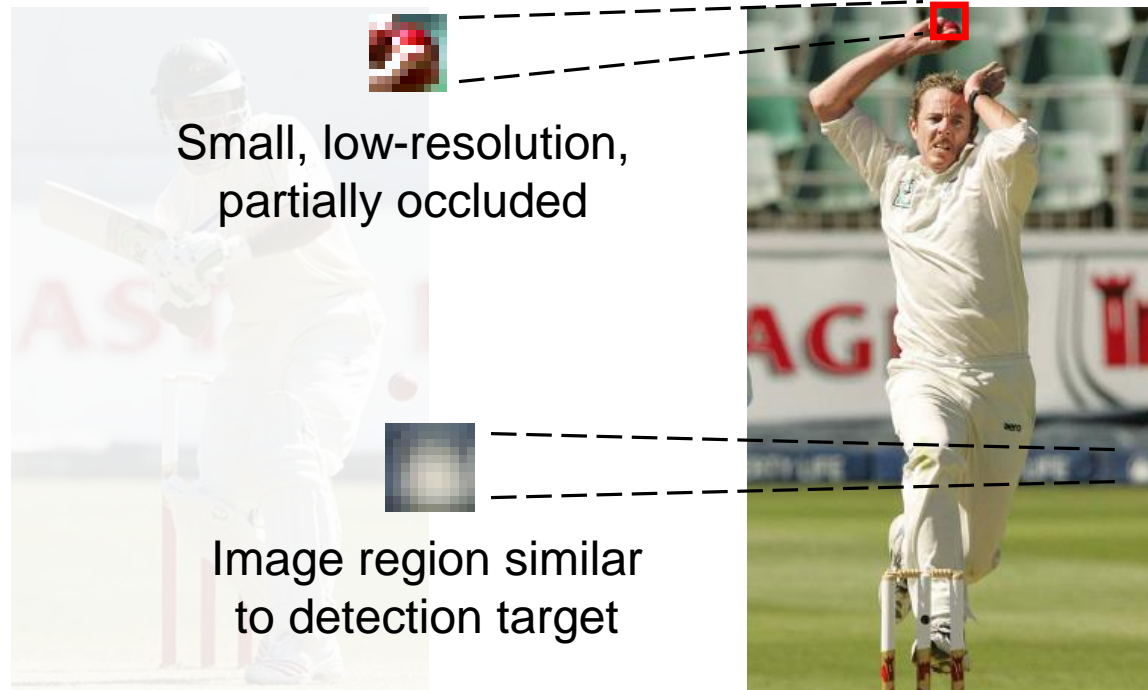
Human pose estimation & Object detection

Facilitate

Given the
object is
detected.



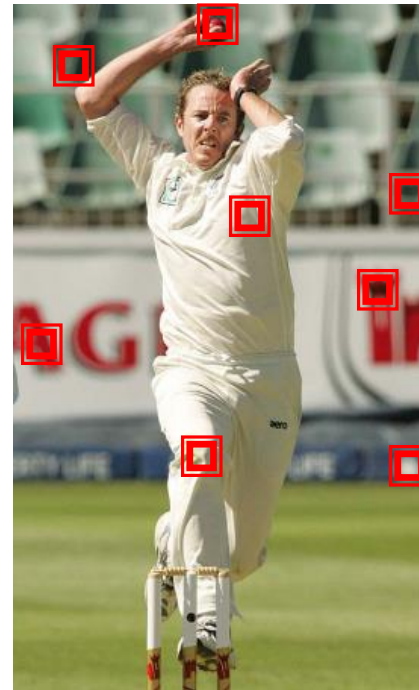
Human pose estimation & Object detection



Object
detection is
challenging

- Viola & Jones, 2001
- Lampert et al, 2008
- Divvala et al, 2009
- Vedaldi et al, 2009

Human pose estimation & Object detection

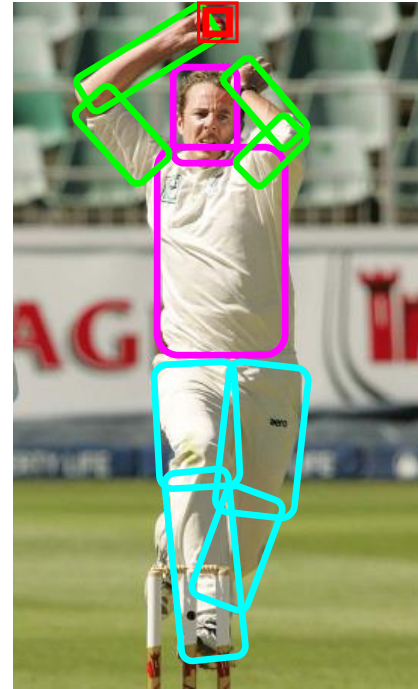


Object
detection is
challenging

- Viola & Jones, 2001
- Lampert et al, 2008
- Divvala et al, 2009
- Vedaldi et al, 2009

Human pose estimation & Object detection

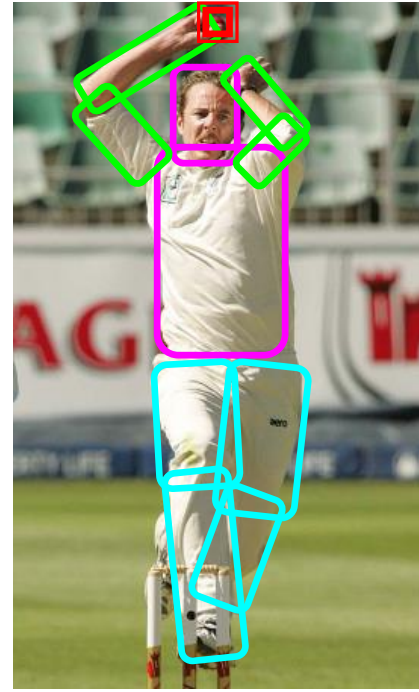
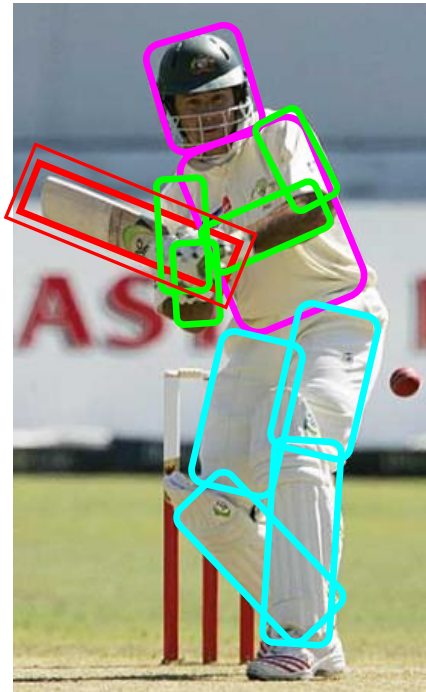
Facilitate



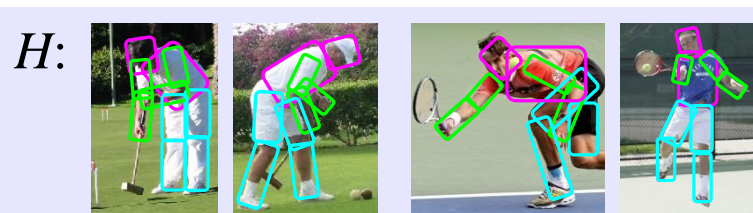
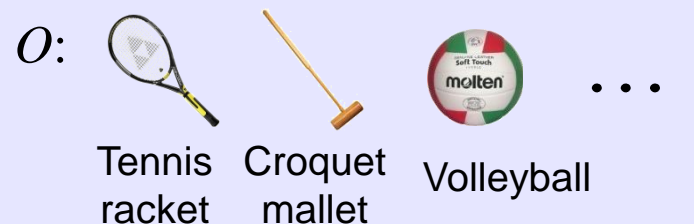
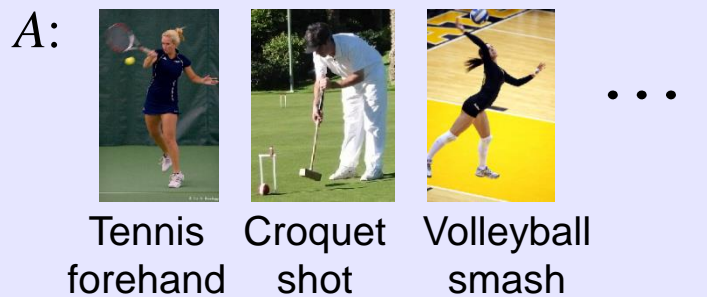
Given the pose is estimated.

Human pose estimation & Object detection

Mutual Context



Mutual Context Model Representation

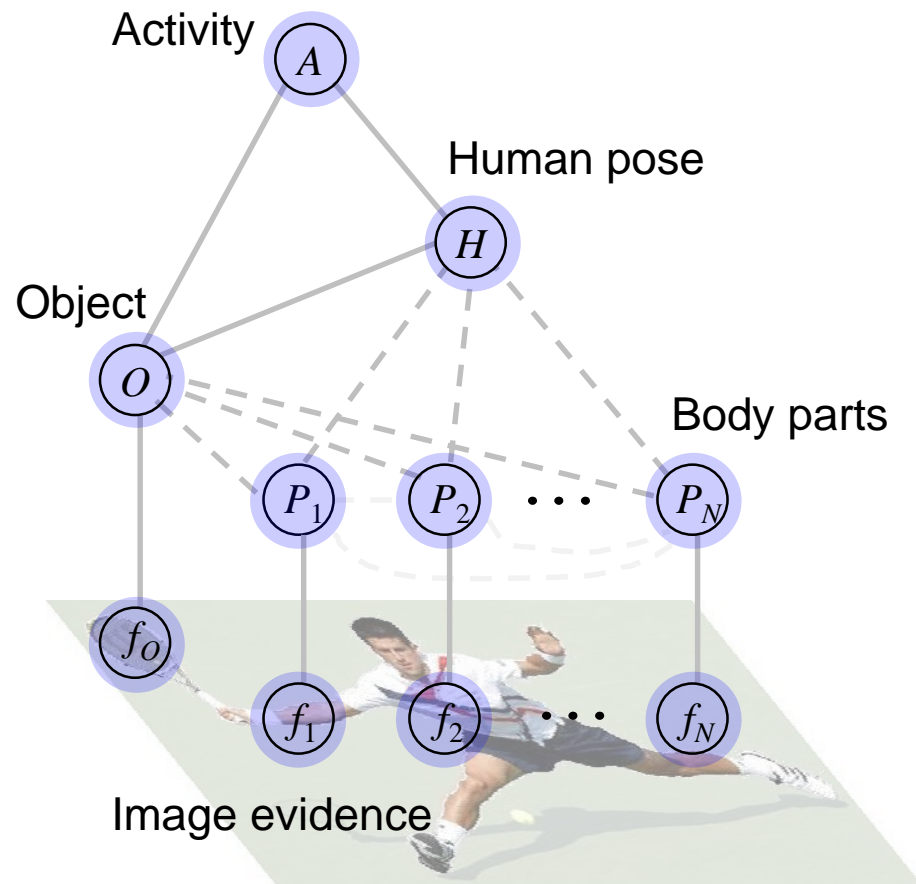


Intra-class variations

- More than one H for each A ;
- **Unobserved** during training.

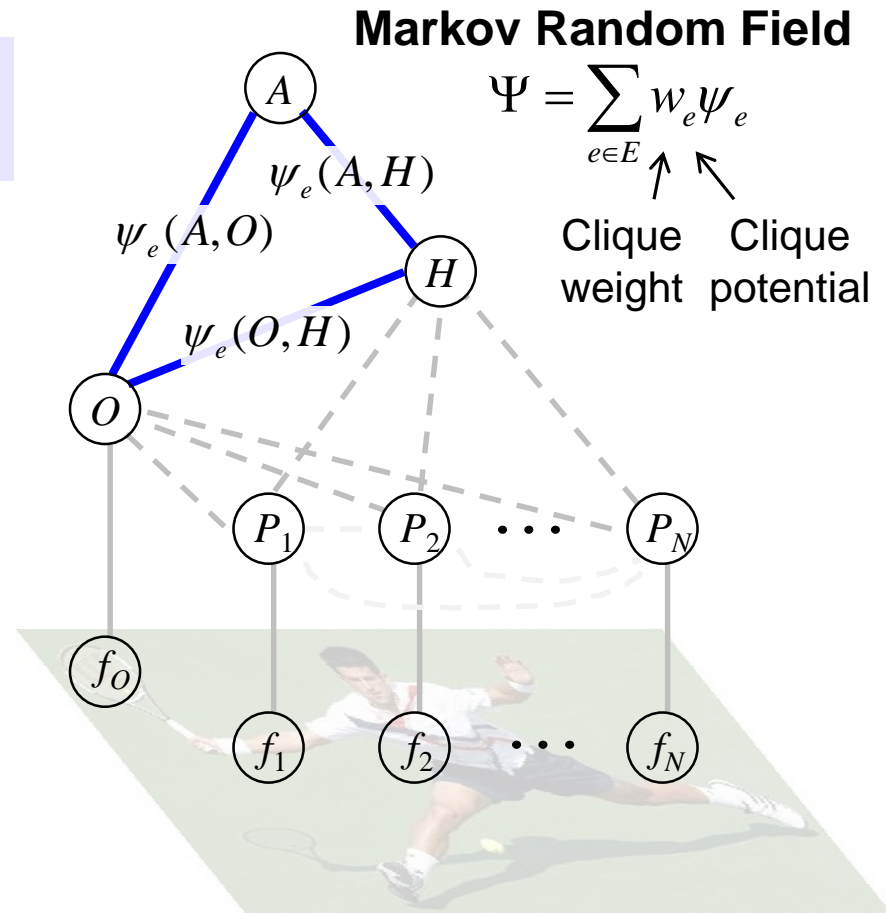
P : l_p : location; θ_p : orientation; s_p : scale.

f : Shape context. [Belongie et al, 2002]



Mutual Context Model Representation

- $\psi_e(A, O), \psi_e(A, H), \psi_e(O, H)$: Frequency of **co-occurrence** between A , O , and H .

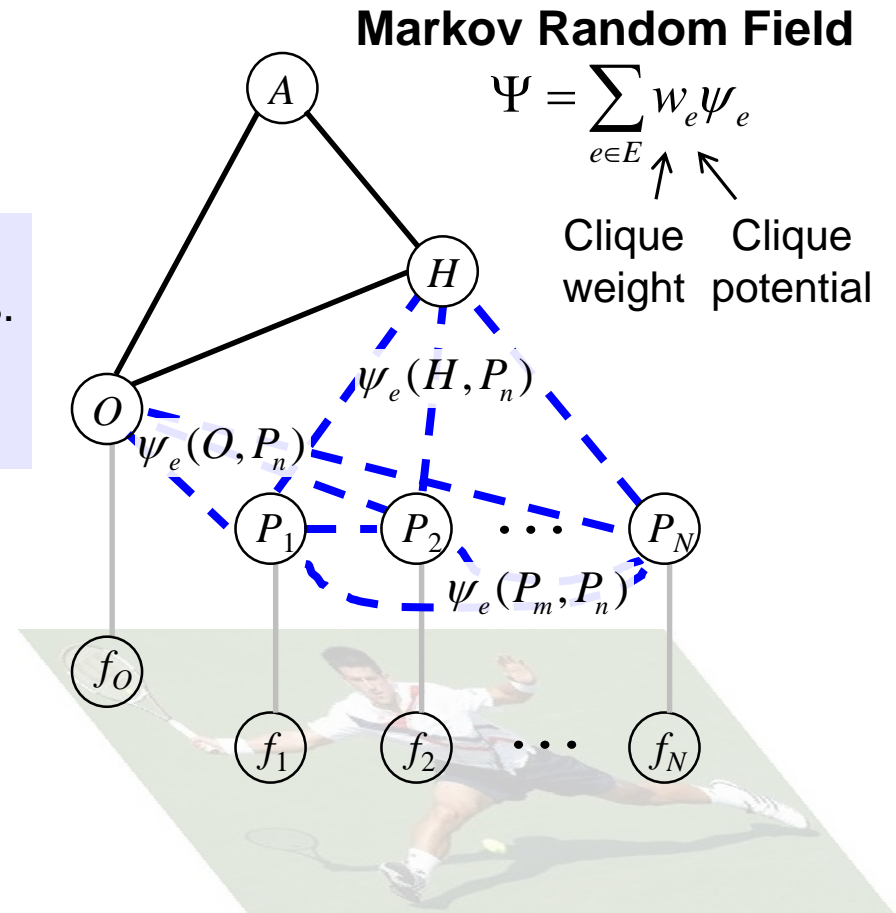


Mutual Context Model Representation

- $\psi_e(A, O), \psi_e(A, H), \psi_e(O, H)$: Frequency of **co-occurrence** between A , O , and H .

- $\psi_e(O, P_n), \psi_e(H, P_n), \psi_e(P_m, P_n)$: **Spatial relationship** among object and body parts.

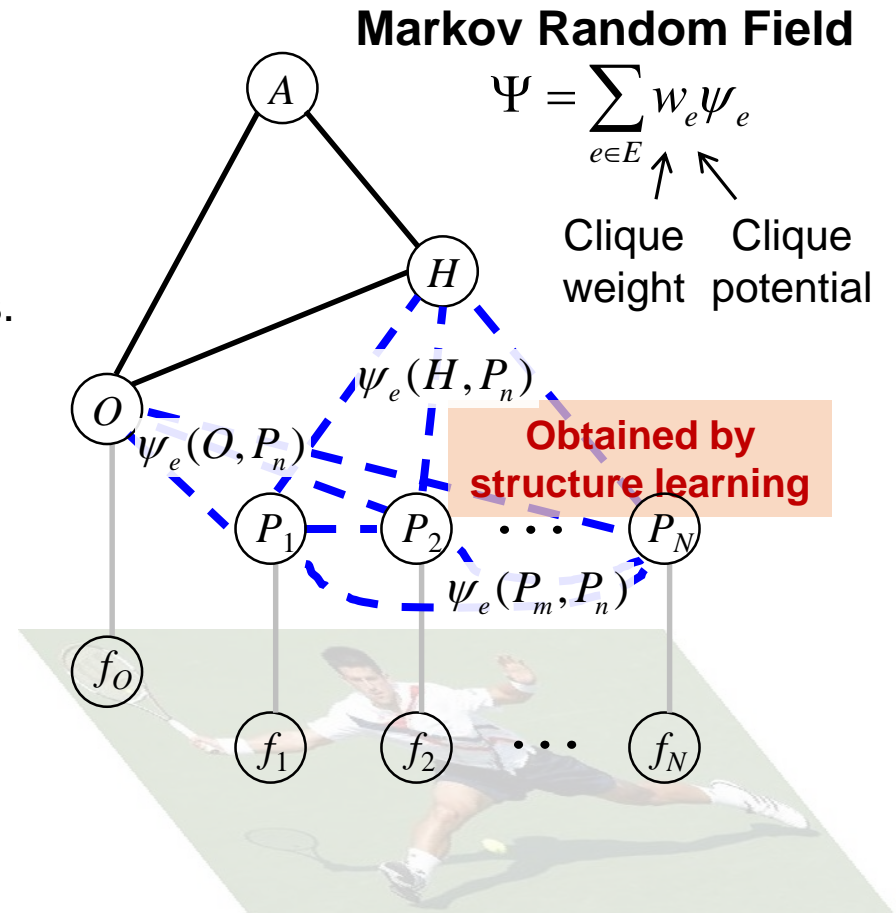
$$\underbrace{\text{bin}(l_O - l_{P_n})}_{\text{location}} \cdot \underbrace{\text{bin}(\theta_O - \theta_{P_n})}_{\text{orientation}} \cdot \underbrace{N(s_O / s_{P_n})}_{\text{size}}$$



Mutual Context Model Representation

- $\psi_e(A, O), \psi_e(A, H), \psi_e(O, H)$: Frequency of **co-occurrence** between A , O , and H .
- $\psi_e(O, P_n), \psi_e(H, P_n), \psi_e(P_m, P_n)$: **Spatial relationship** among object and body parts.

$$\underbrace{\text{bin}(l_O - l_{P_n})}_{\text{location}} \cdot \underbrace{\text{bin}(\theta_O - \theta_{P_n})}_{\text{orientation}} \cdot \underbrace{N(s_O / s_{P_n})}_{\text{size}}$$
- **Learn structural connectivity** among the body parts and the object.



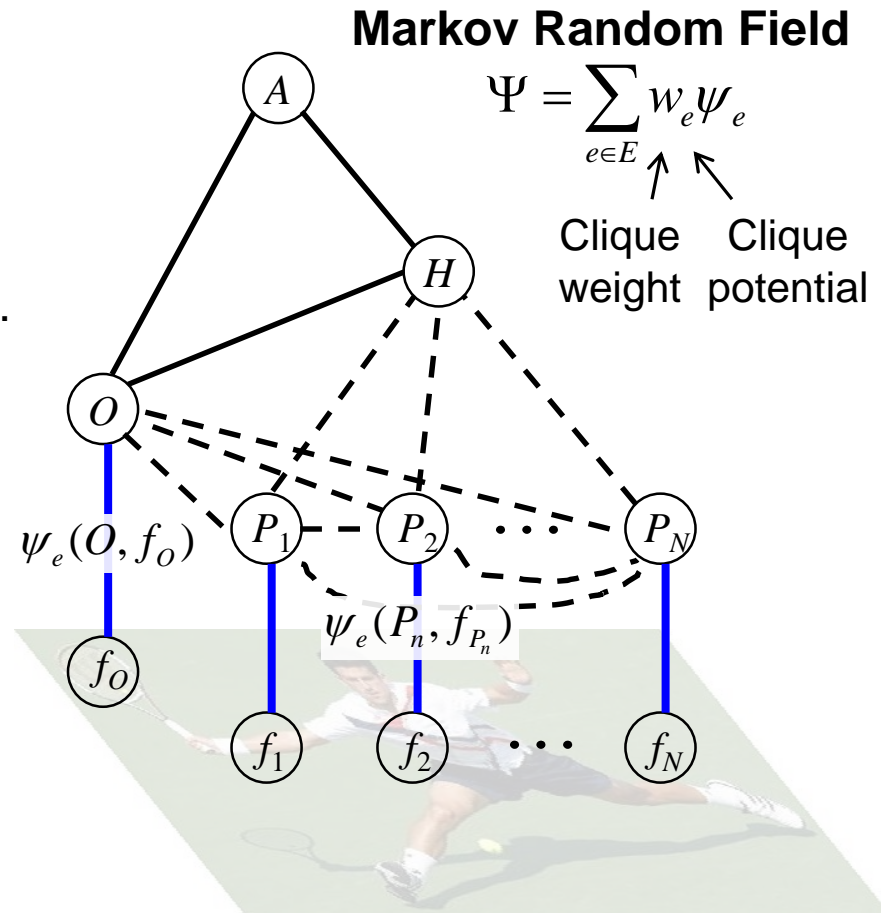
Mutual Context Model Representation

- $\psi_e(A, O), \psi_e(A, H), \psi_e(O, H)$: Frequency of **co-occurrence** between A , O , and H .
- $\psi_e(O, P_n), \psi_e(H, P_n), \psi_e(P_m, P_n)$: **Spatial relationship** among object and body parts.

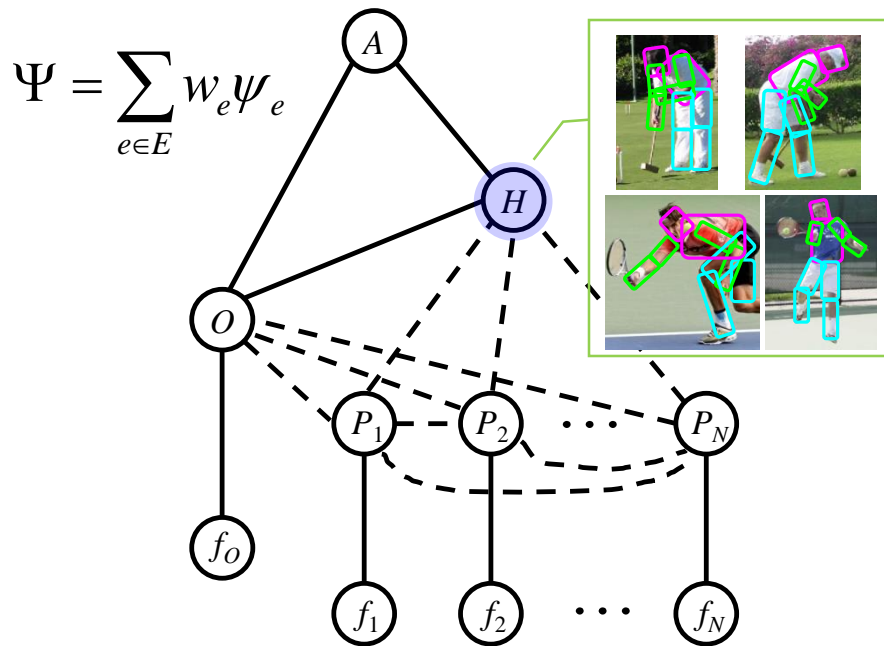
$$\underbrace{\text{bin}(l_O - l_{P_n})}_{\text{location}} \cdot \underbrace{\text{bin}(\theta_O - \theta_{P_n})}_{\text{orientation}} \cdot \underbrace{N(s_O / s_{P_n})}_{\text{size}}$$
- **Learn structural connectivity** among the body parts and the object.
- $\psi_e(O, f_O)$ and $\psi_e(P_n, f_{P_n})$: **Discriminative part detection** scores.

Shape context + AdaBoost

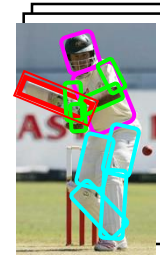
[Andriluka et al, 2009]
 [Belongie et al, 2002]
 [Viola & Jones, 2001]



Model Learning



Input:



cricket
shot



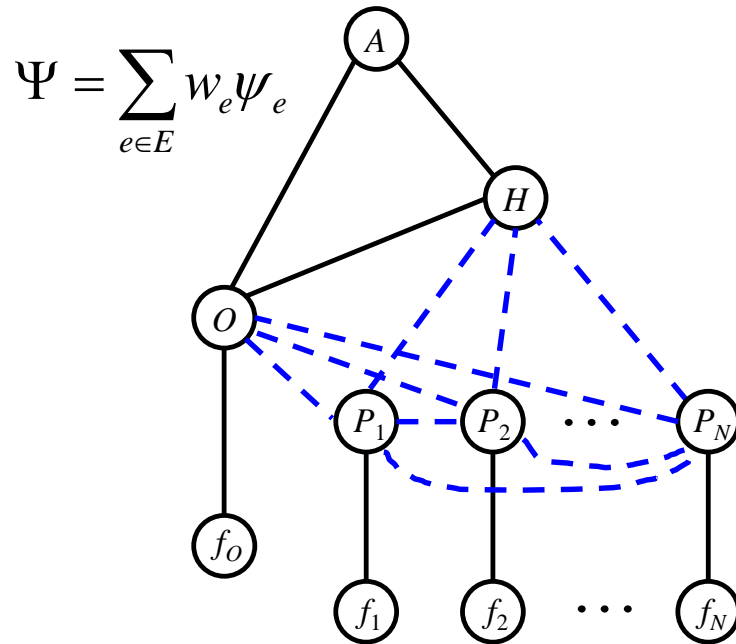
cricket
bowling

...

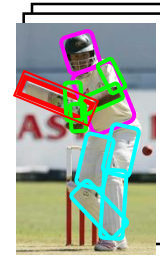
Goals:

Hidden human poses

Model Learning



Input:



cricket
shot



cricket
bowling

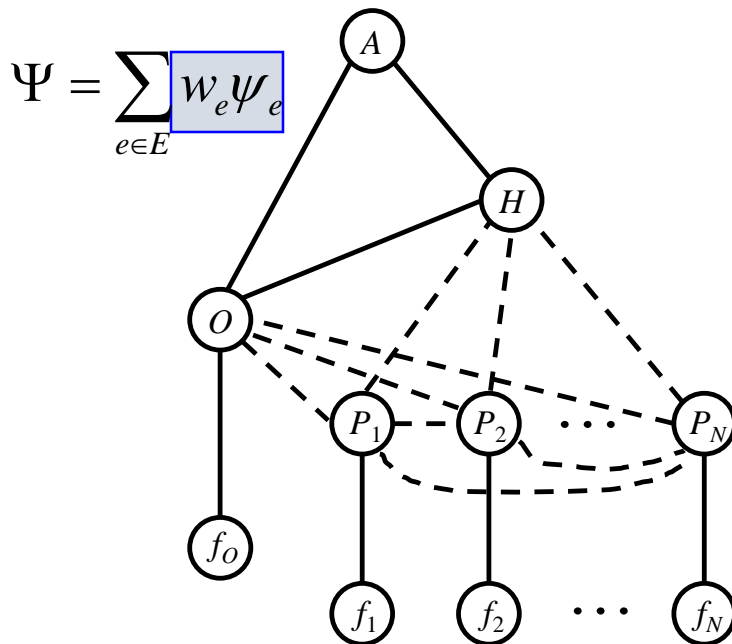
...

Goals:

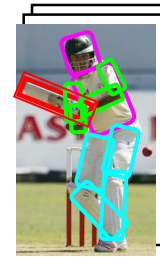
Hidden human poses

Structural connectivity

Model Learning



Input:



cricket
shot



cricket
bowling

...

Goals:

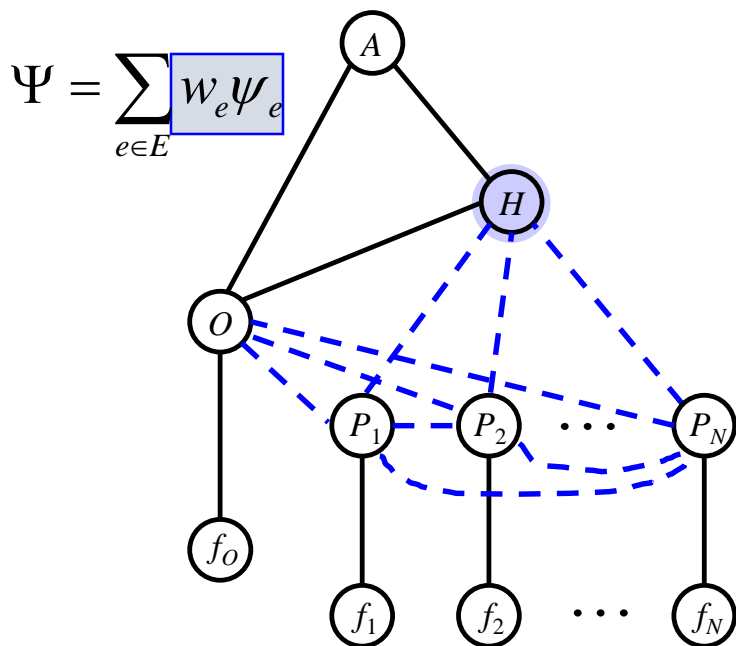
Hidden human poses

Structural connectivity

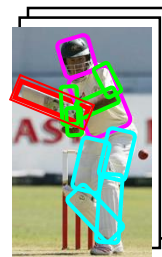
Potential parameters

Potential weights

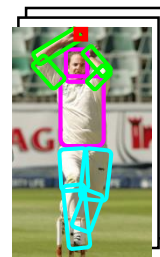
Model Learning



Input:



cricket
shot



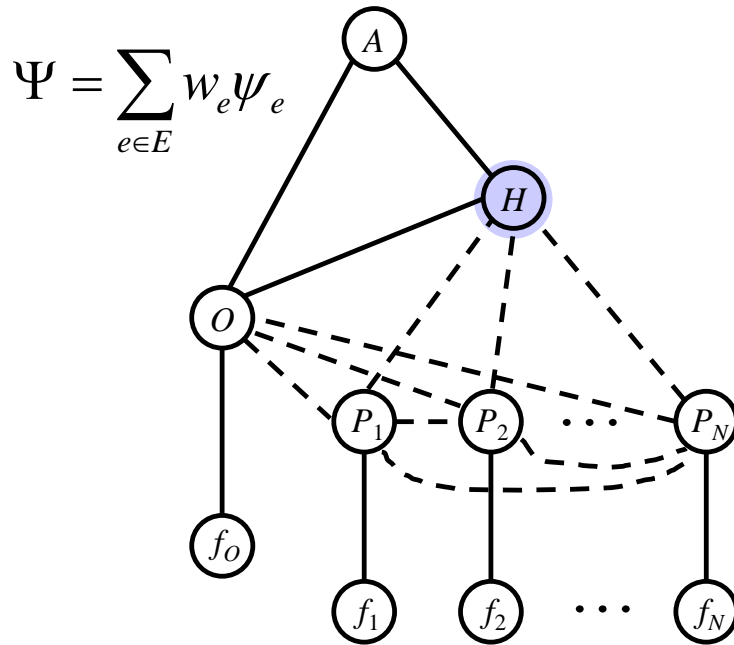
cricket
bowling

...

Goals:

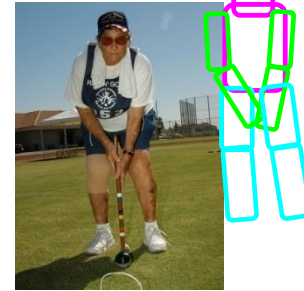
- Hidden human poses → **Hidden variables**
 - Structural connectivity → **Structure learning**
 - Potential parameters
 - Potential weights
- } **Parameter estimation**

Model Learning



Approach:

croquet shot



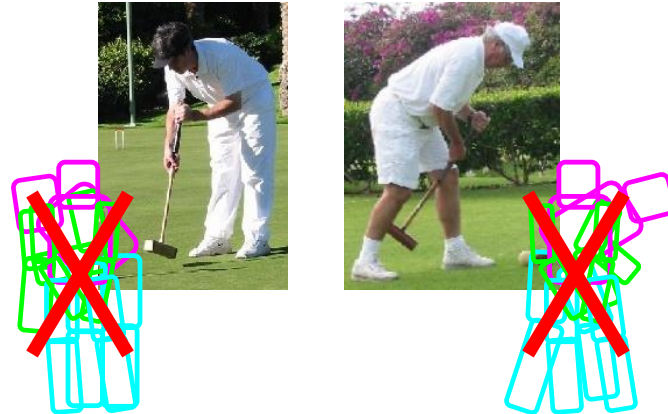
Goals:

Hidden human poses

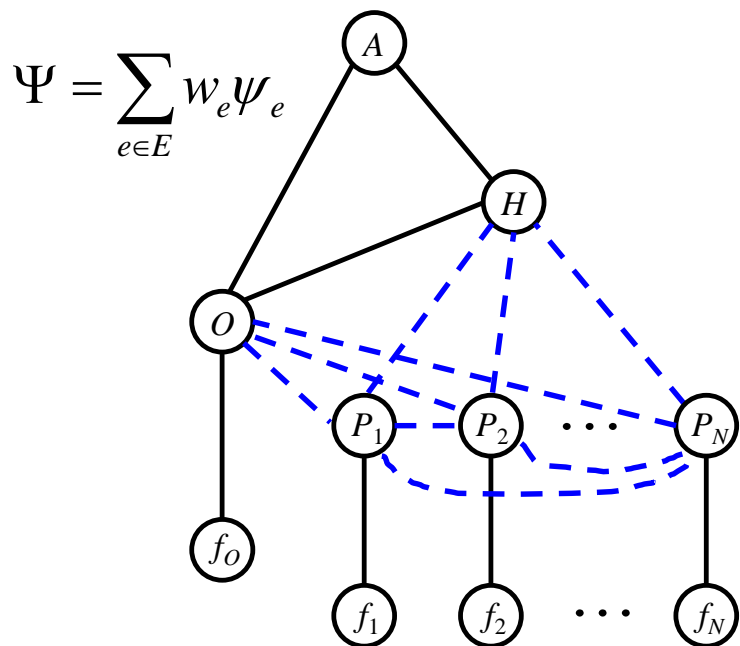
Structural connectivity

Potential parameters

Potential weights



Model Learning



Approach:

Hill-climbing $\max_{E=\{e\}} \left\{ \underbrace{\sum_e w_e \psi_e}_{\text{Joint density of the model}} - \frac{(|E| - \mu)^2}{\underbrace{2\sigma^2}_{\text{Gaussian priori of the edge number}}} \right\}$

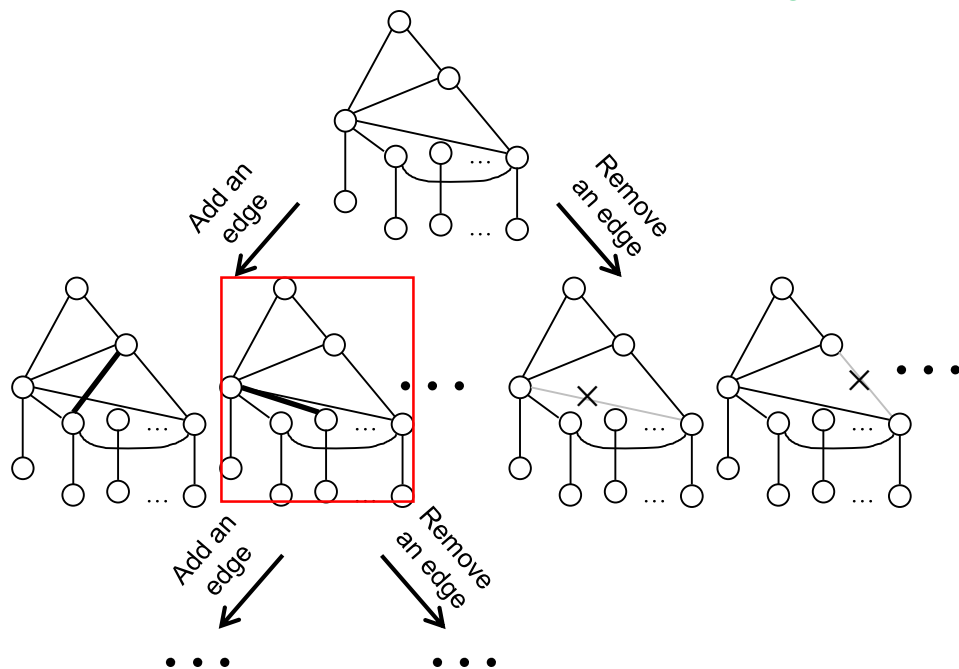
Goals:

Hidden human poses

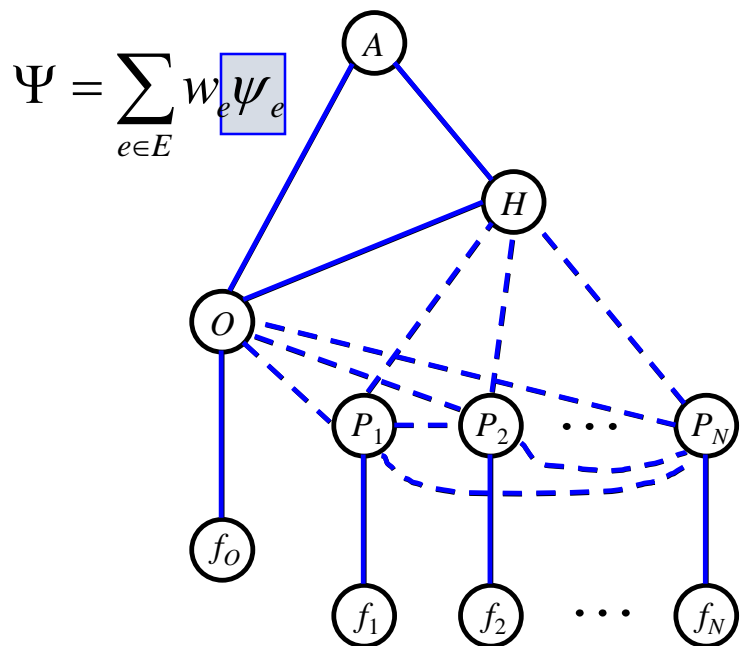
Structural connectivity

Potential parameters

Potential weights



Model Learning



Approach:

- Maximum likelihood

$$\begin{aligned} &\psi_e(A, O) \quad \psi_e(A, H) \quad \psi_e(O, H) \\ &\psi_e(H, P_n) \quad \psi_e(O, P_n) \quad \psi_e(P_m, P_n) \end{aligned}$$

- Standard AdaBoost

$$\psi_e(O, f_o) \quad \psi_e(P_n, f_{P_n})$$

Goals:

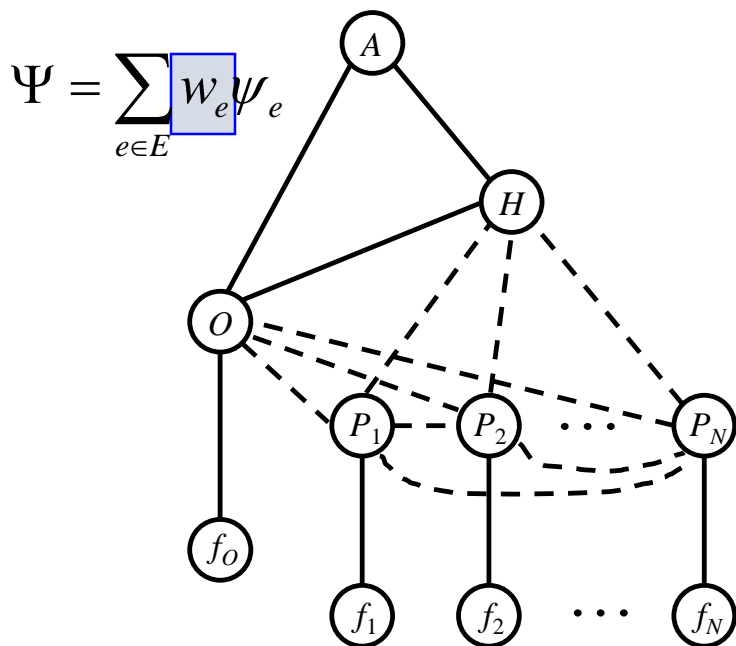
Hidden human poses

Structural connectivity

Potential parameters

Potential weights

Model Learning



Goals:

Hidden human poses
Structural connectivity
Potential parameters
Potential weights

Approach:

Max-margin learning

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \sum_r \|\mathbf{w}_r\|_2^2 + \beta \sum_i \xi_i$$

s.t. $\forall i, r$ where $y(r) \neq y(c_i)$,

$$\mathbf{w}_{c_i} \cdot \mathbf{x}_i - \mathbf{w}_r \cdot \mathbf{x}_i \geq 1 - \xi_i$$

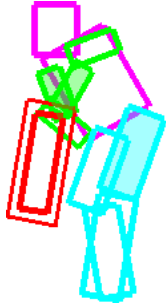
$$\forall i, \xi_i \geq 0$$

Notations

- \mathbf{x}_i : Potential values of the i -th image.
- \mathbf{w}_r : Potential weights of the r -th pose.
- $y(r)$: Activity of the r -th pose.
- ξ_i : A slack variable for the i -th image.

Learning Results

Cricket
defensive
shot



Cricket
bowling

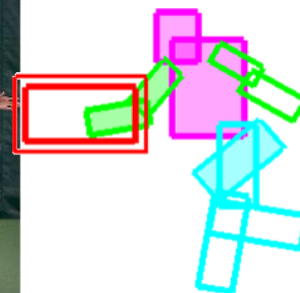


Croquet
shot



Learning Results

Tennis
forehand



Tennis
serve



Volleyball
smash

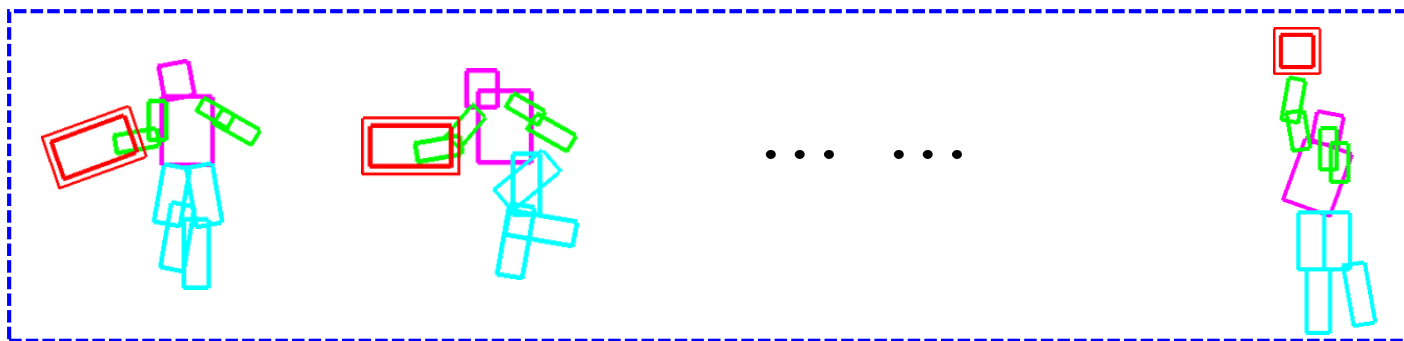


Model Inference

I



The learned models

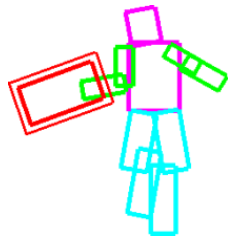


Model Inference

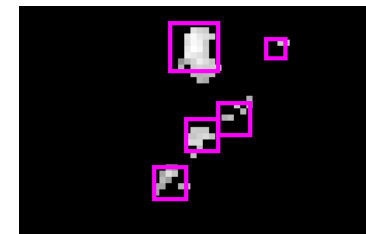
I



The learned models



...



Head detection



Torso detection

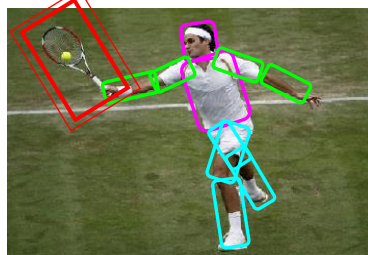
⋮



Tennis racket detection

Compositional Inference

[Chen et al, 2007]



$$\Psi\left(A_1, H_1, O_1^*, \{P_{1,n}^*\}_n\right)$$

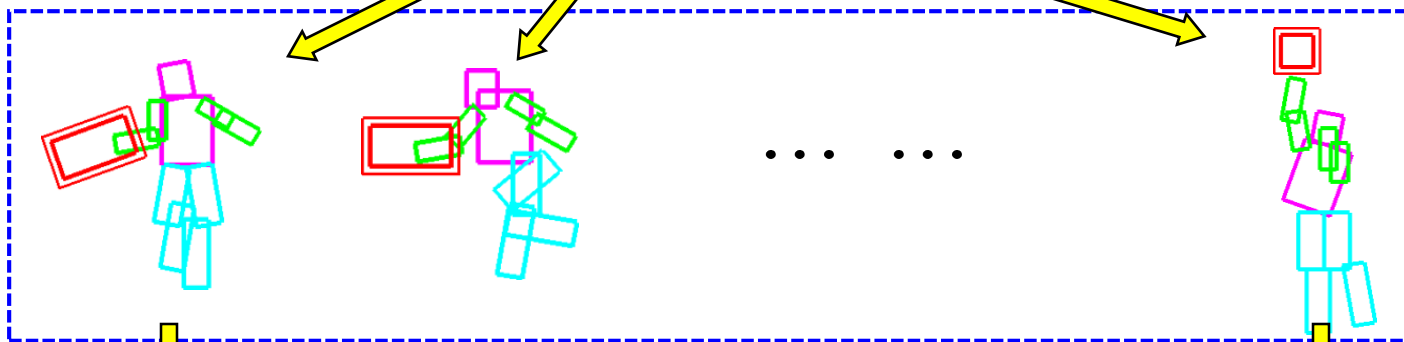
Layout of the **object** and **body parts**.

Model Inference

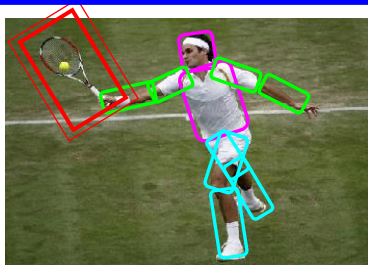
I



The learned models

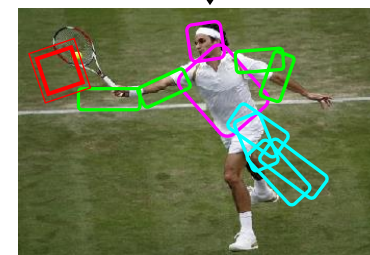


Output



$$\Psi\left(A_1, H_1, O_1^*, \{P_{1,n}^*\}_n\right)$$

...



$$\Psi\left(A_K, H_K, O_K^*, \{P_{K,n}^*\}_n\right)$$

Dataset and Experiment Setup

Sport data set: 6 classes

180 training (supervised with object and part locations) & 120 testing images



Cricket
defensive shot



Cricket
bowling



Croquet
shot



Tennis
forehand



Tennis
serve



Volleyball
smash

Tasks:

- Object detection;
- Pose estimation;
- Activity classification.

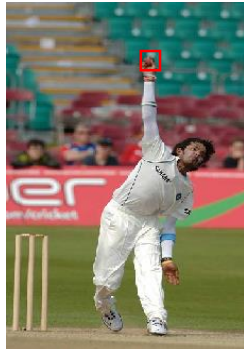
Dataset and Experiment Setup

Sport data set: 6 classes

180 training (supervised with object and part locations) & 120 testing images



Cricket
defensive shot



Cricket
bowling



Croquet
shot



Tennis
forehand



Tennis
serve



Volleyball
smash

Tasks:

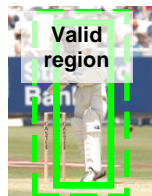
- **Object detection;**
- Pose estimation;
- Activity classification.

Object Detection Results



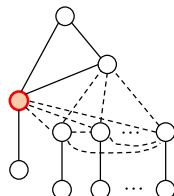
Sliding window

[Andriluka et al, 2009]



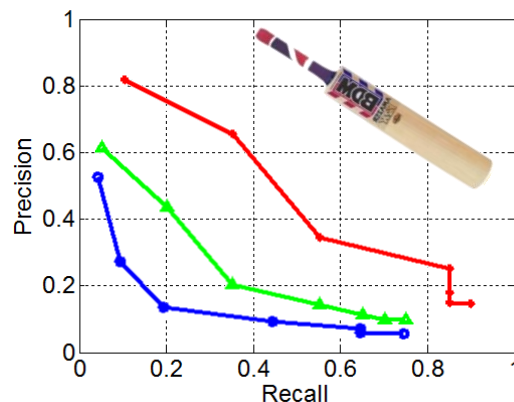
Pedestrian context

[Dalal & Triggs, 2006]

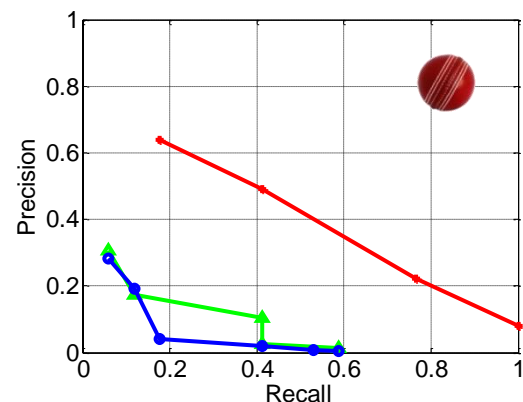


Our Method

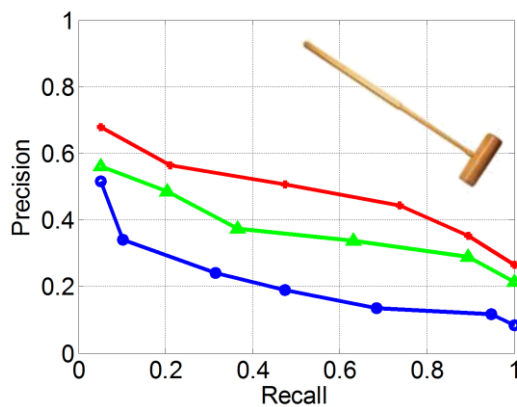
Cricket bat



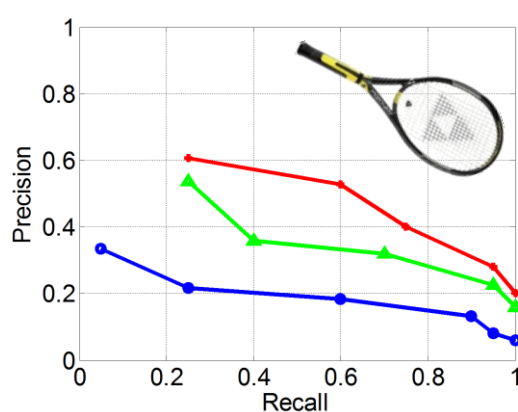
Cricket ball



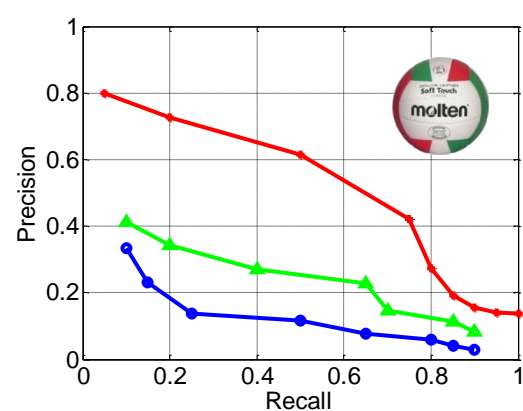
Croquet mallet



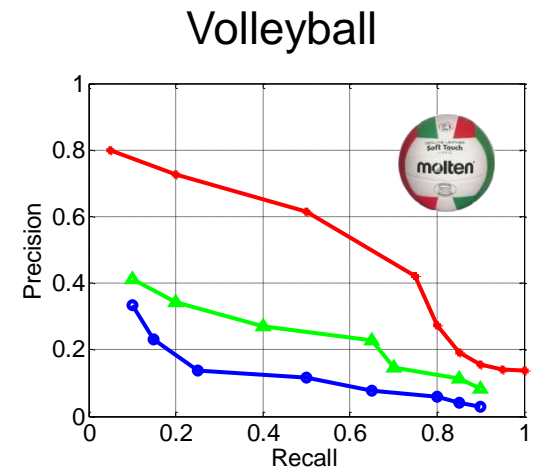
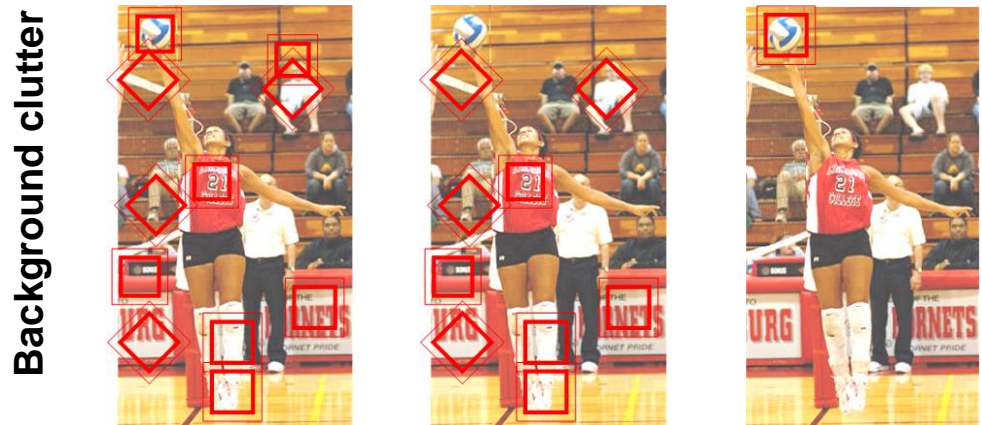
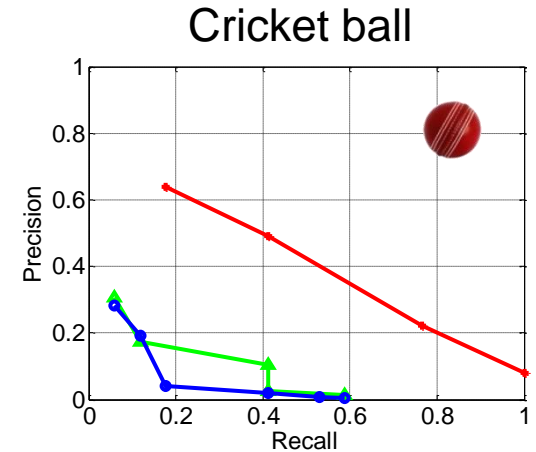
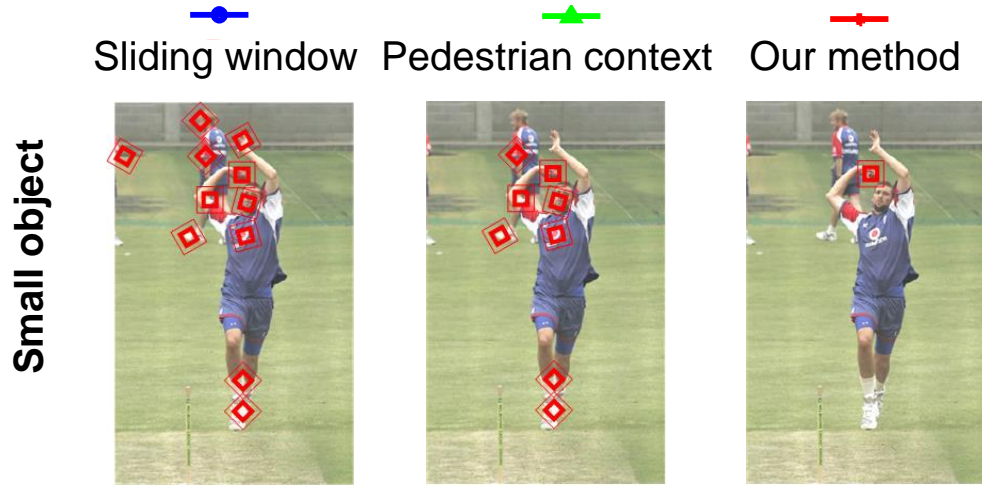
Tennis racket



Volleyball



Object Detection Results



Dataset and Experiment Setup

Sport data set: 6 classes

180 training & 120 testing images



Cricket
defensive shot



Cricket
bowling



Croquet
shot



Tennis
forehand



Tennis
serve



Volleyball
smash

Tasks:

- Object detection;
- **Pose estimation;**
- Activity classification.

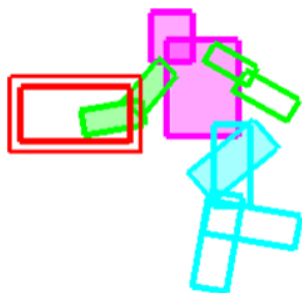
[Gupta et al, 2009]

Human Pose Estimation Results

Method	Torso	Upper Leg		Lower Leg		Upper Arm		Lower Arm		Head
Ramanan, 2006	.52	.22	.22	.21	.28	.24	.28	.17	.14	.42
Andriluka et al, 2009	.50	.31	.30	.31	.27	.18	.19	.11	.11	.45
Our full model	.66	.43	.39	.44	.34	.44	.40	.27	.29	.58

Human Pose Estimation Results

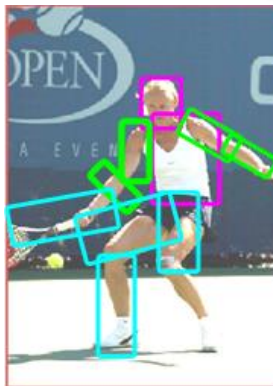
Method	Torso	Upper Leg		Lower Leg		Upper Arm		Lower Arm		Head
Ramanan, 2006	.52	.22	.22	.21	.28	.24	.28	.17	.14	.42
Andriluka et al, 2009	.50	.31	.30	.31	.27	.18	.19	.11	.11	.45
Our full model	.66	.43	.39	.44	.34	.44	.40	.27	.29	.58



Tennis serve model



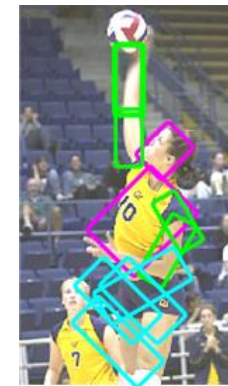
Our estimation result



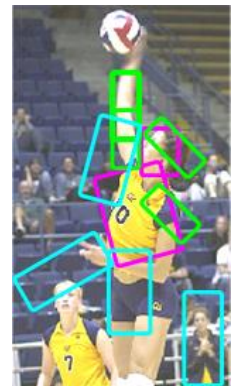
Andriluka et al, 2009



Volleyball smash model



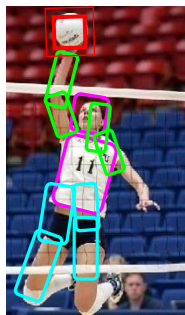
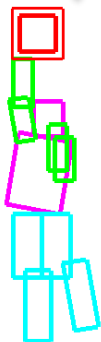
Our estimation result



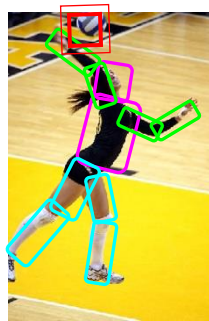
Andriluka et al, 2009

Human Pose Estimation Results

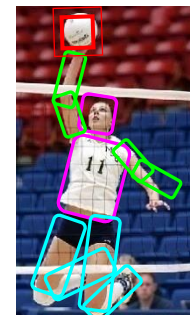
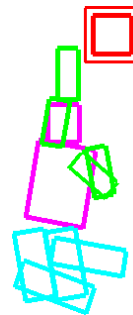
Method	Torso	Upper Leg		Lower Leg		Upper Arm		Lower Arm		Head
Ramanan, 2006	.52	.22	.22	.21	.28	.24	.28	.17	.14	.42
Andriluka et al, 2009	.50	.31	.30	.31	.27	.18	.19	.11	.11	.45
Our full model	.66	.43	.39	.44	.34	.44	.40	.27	.29	.58
One pose per class	.63	.40	.36	.41	.31	.38	.35	.21	.23	.52



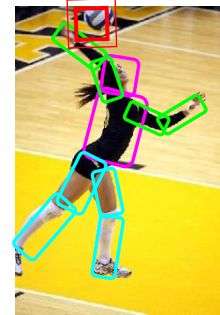
Estimation
result



Estimation
result



Estimation
result

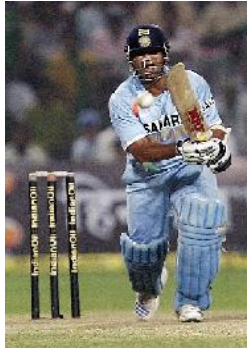


Estimation
result

Dataset and Experiment Setup

Sport data set: 6 classes

180 training & 120 testing images



Cricket
defensive shot



Cricket
bowling



Croquet
shot



Tennis
forehand



Tennis
serve



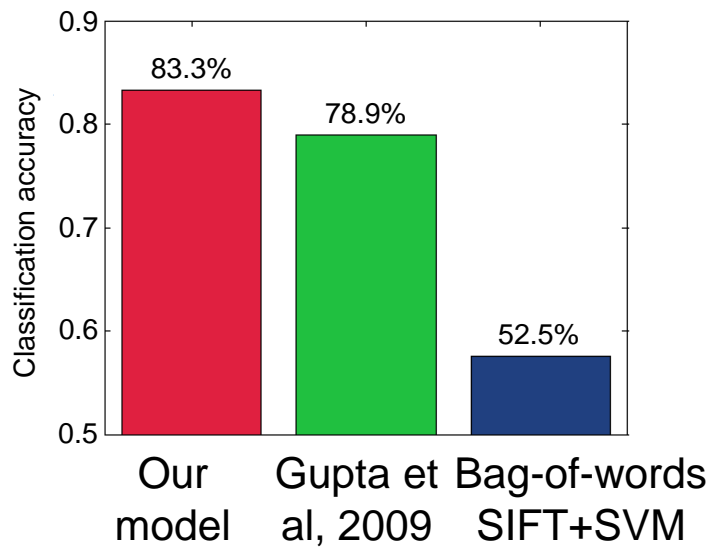
Volleyball
smash

Tasks:

- Object detection;
- Pose estimation;
- **Activity classification.**

[Gupta et al, 2009]

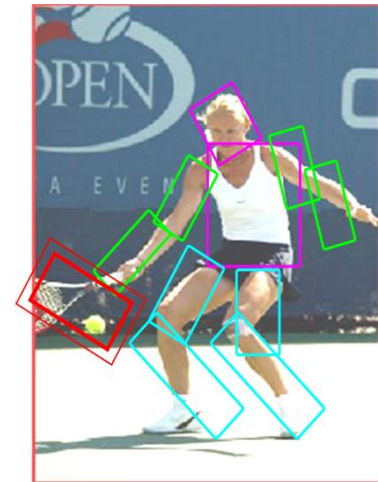
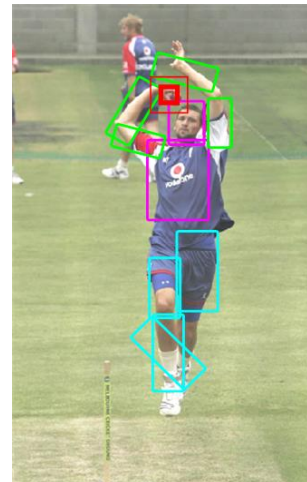
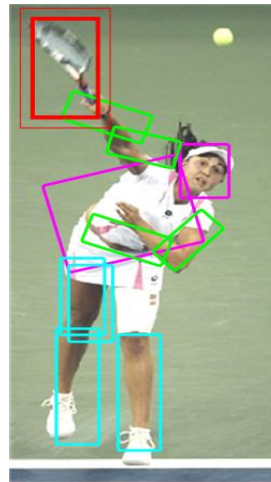
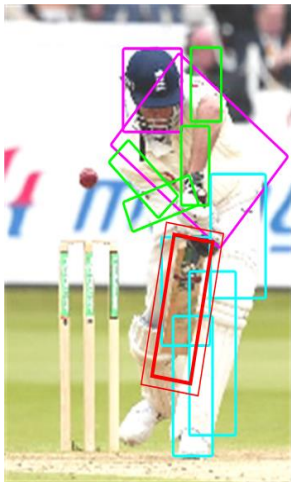
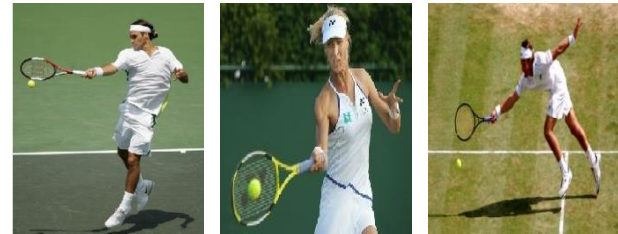
Activity Classification Results



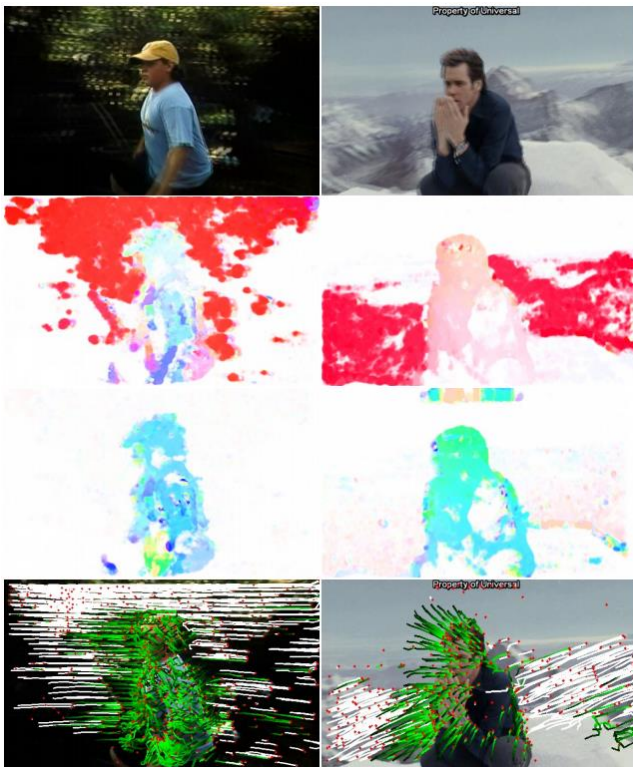
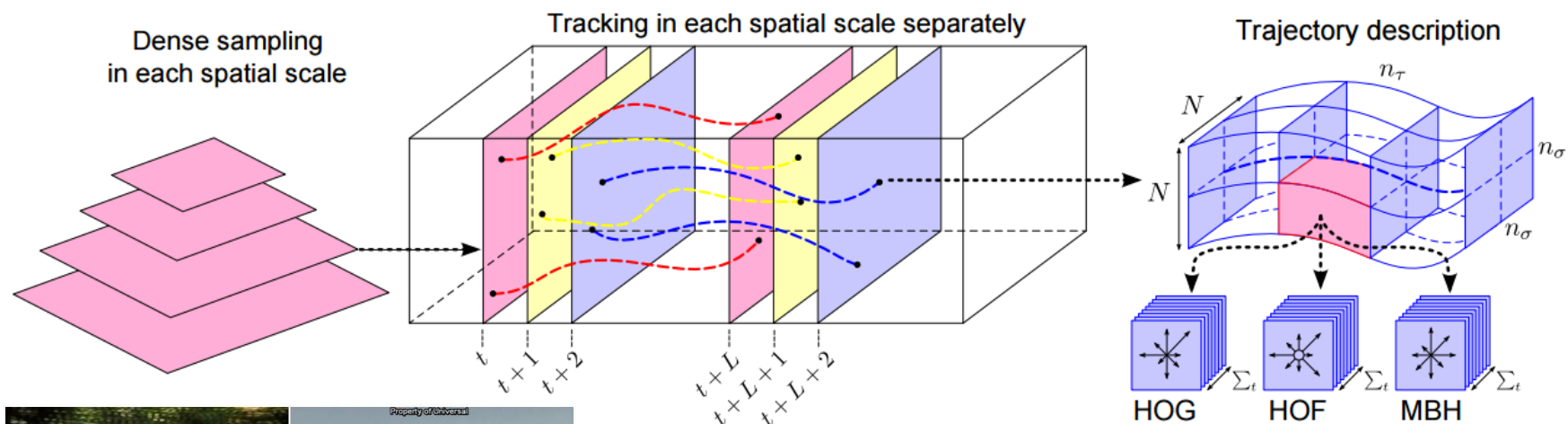
Cricket shot



Tennis forehand

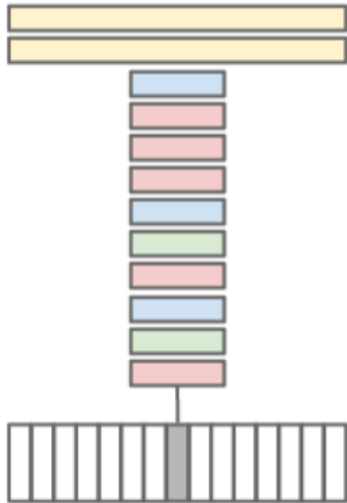


Motion features – Dense Trajectory

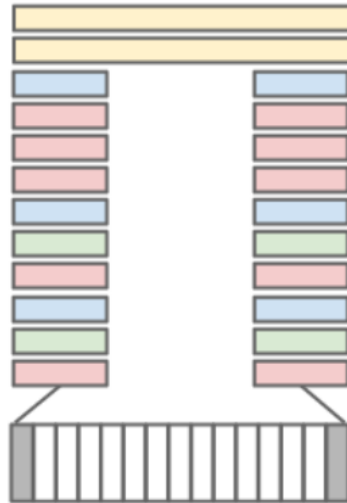


Video classification with CNNs

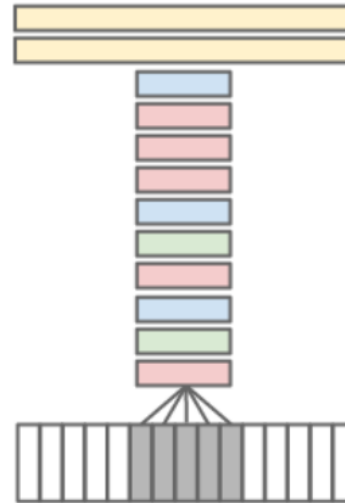
Single Frame



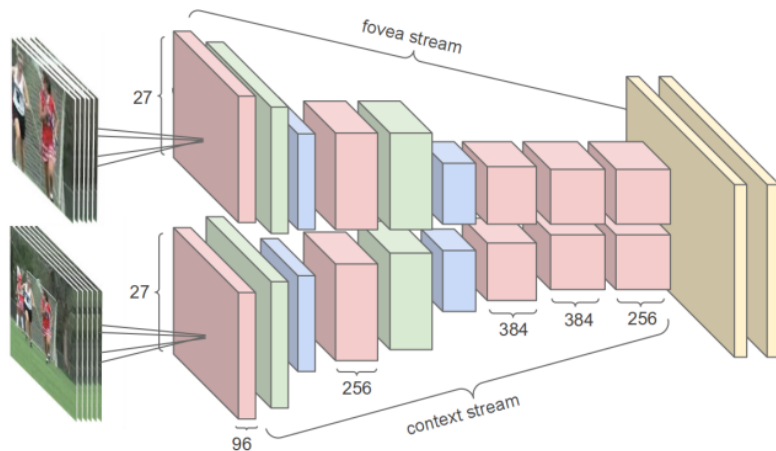
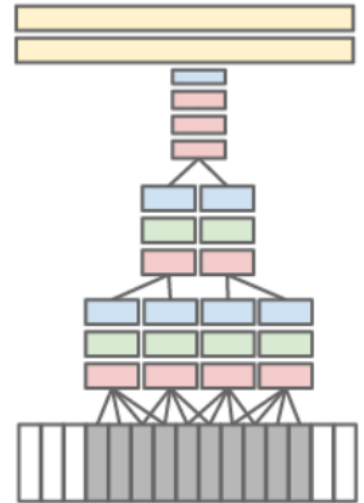
Late Fusion



Early Fusion



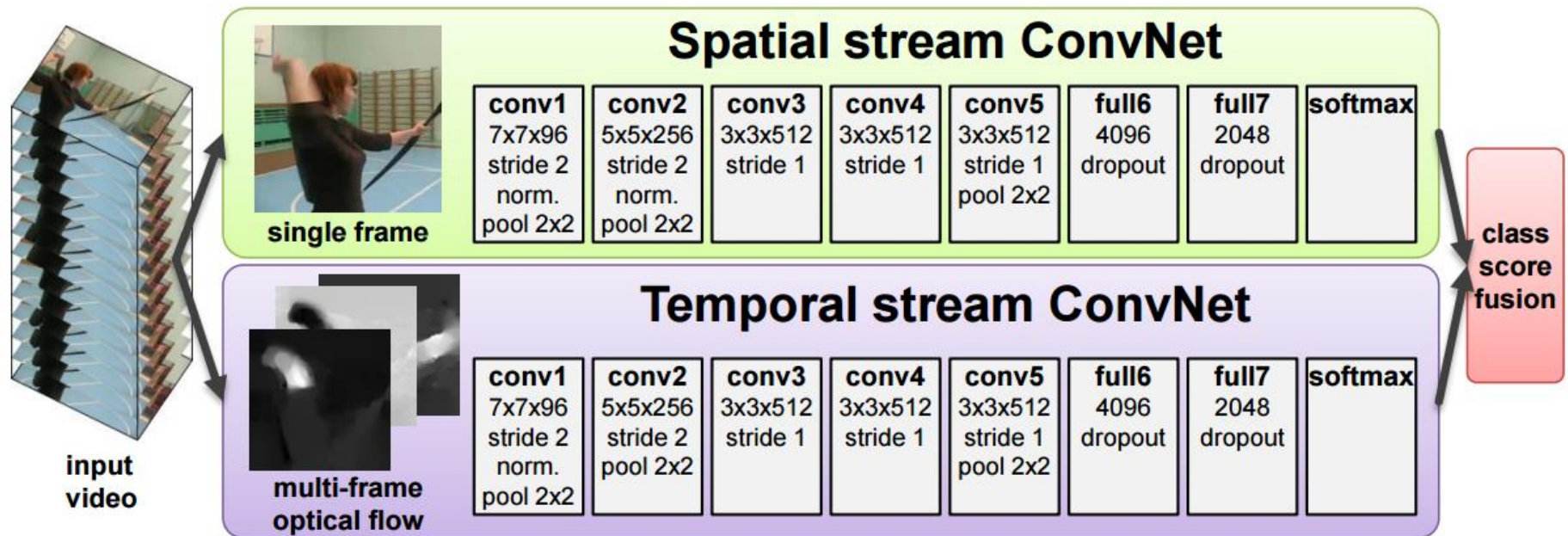
Slow Fusion



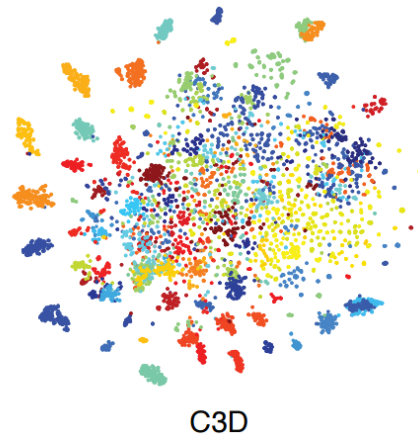
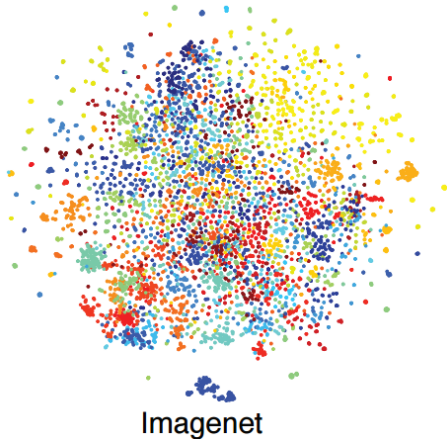
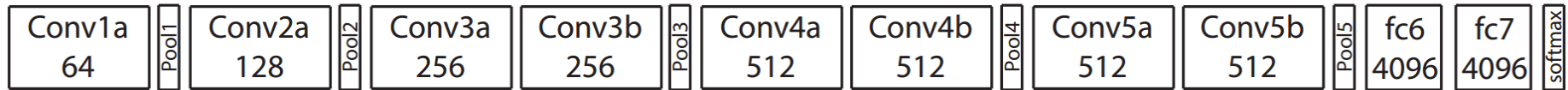
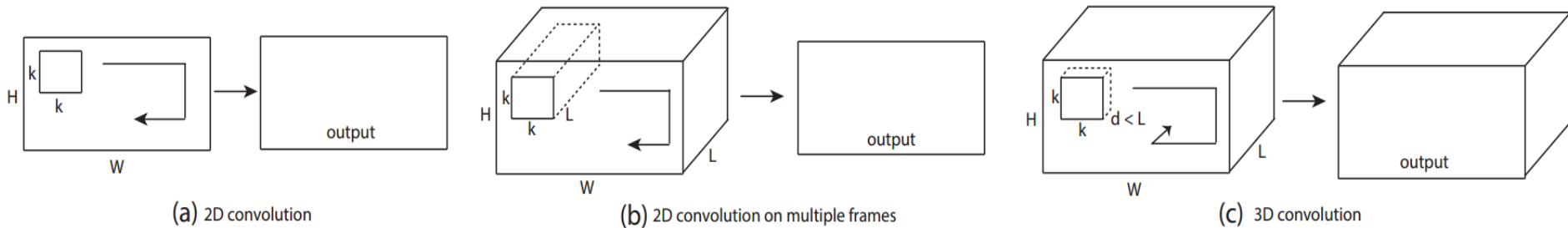
Video classification with CNNs

Sports Video Classification

Two-stream CNN



3D Convolutional Networks



Take-home messages

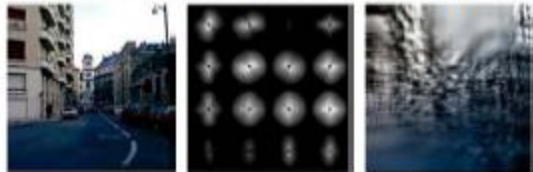
- Action recognition is an open problem.
 - How to define actions?
 - How to infer them?
 - What are good visual cues?
 - How do we incorporate higher level reasoning?

Take-home messages

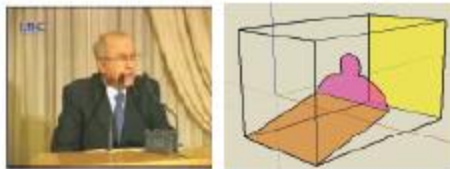
- Some work done, but it is just the beginning of exploring the problem. So far...
 - Actions are mainly categorical (could be framed in terms of effect or intent)
 - Most approaches are classification using simple features (spatial-temporal histograms of gradients or flow, s-t interest points, SIFT in images)
 - Just a couple works on how to incorporate pose and objects
 - Not much idea of how to reason about long-term activities or to describe video sequences

Next class: 3D Scenes and Context

Scene-Level Geometric Description



a) Gist, Spatial Envelope



b) Stages

Retinotopic Maps



c) Geometric Context



d) Depth Maps

Highly Structured 3D Models



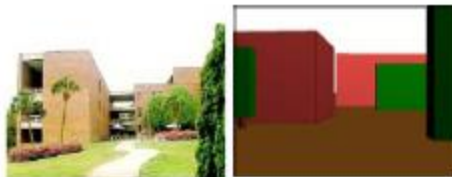
e) Ground Plane



f) Ground Plane with Billboards



g) Ground Plane with Walls



h) Blocks World



i) 3D Box Model

