

Emerging Topics in Human Activity Recognition

Michael Ryoo NASA Jet Propulsion Laboratory

Ivan Laptev INRIA

Greg Mori Simon Fraser University

Sangmin Oh Kitware

CVPR tutorial on 2014/06/23



Group Activity Recognition

Greg Mori
Simon Fraser University

CVPR tutorial on 2014/06/23



2012/03/26 PM07:37:04
GardenNVR -

What does activity recognition involve?



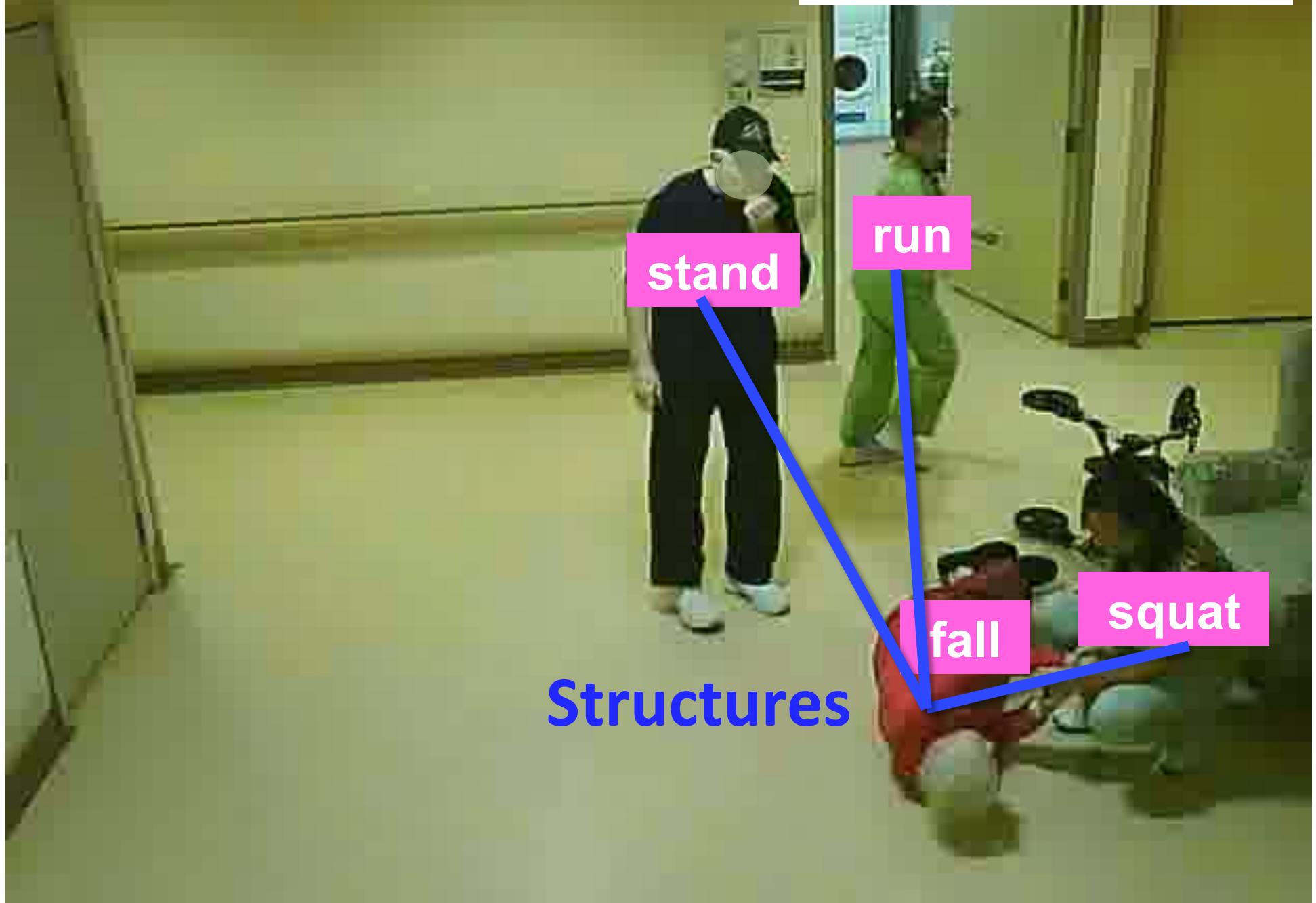
2012/03/26 PM07:37:04
GardenNVR -

Detection: are there people?



2012/03/26 PM07:37:04
GardenNVR -

Action recognition



2012/03/26 PM07:37:04
GardenNVR -

Group activity recognition



help the fallen
person

2012/03/26 PM07:37:04
GardenNVR -

Intention/social role





- Intertwined problems
- Build models that jointly solve these problems
 - Ability to focus on sub-problems, leverage information from others
 - Principled learning algorithms

Surveillance

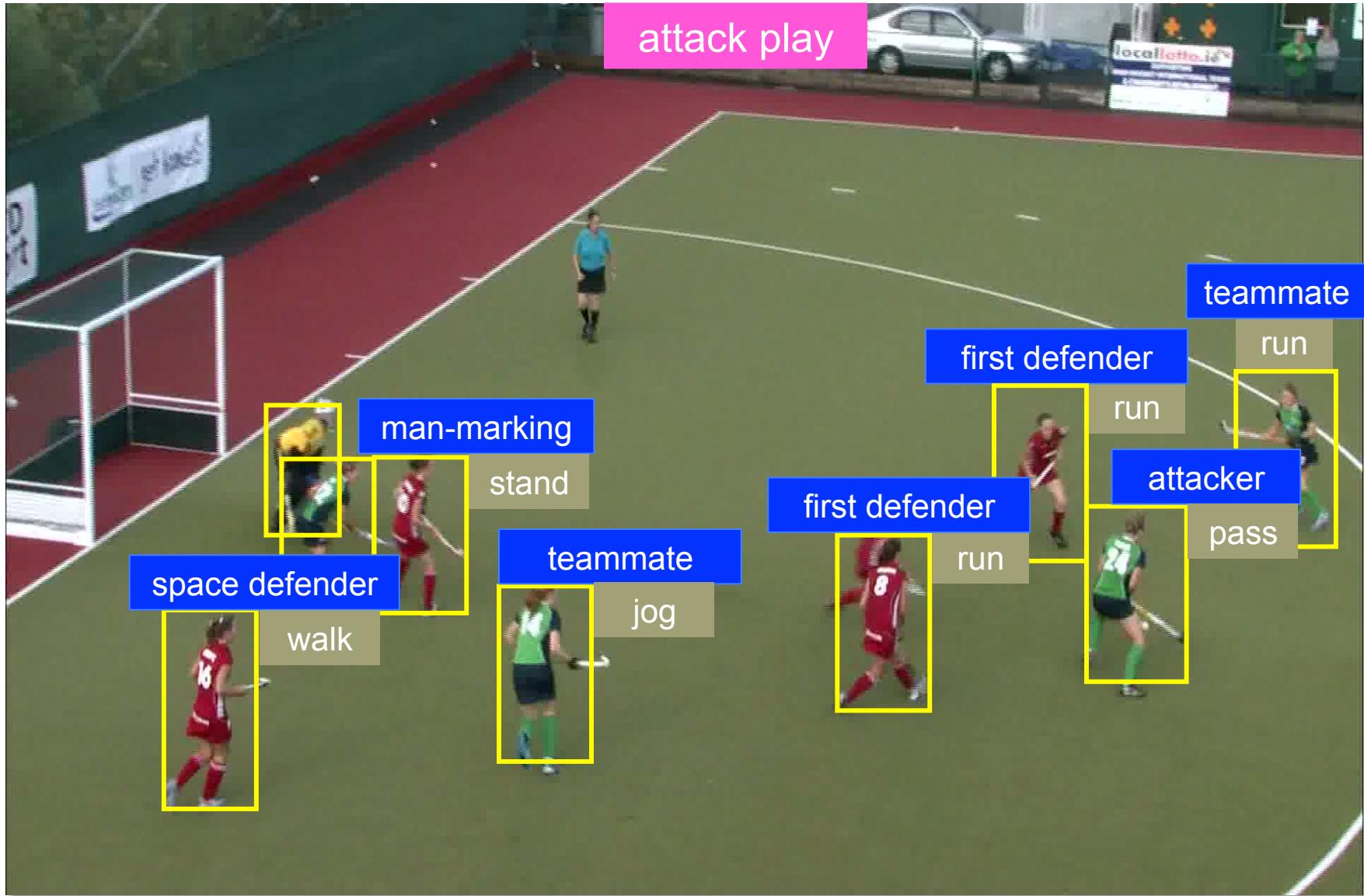


VIRAT dataset: Oh et al. CVPR 2011

Surveillance – Road Safety

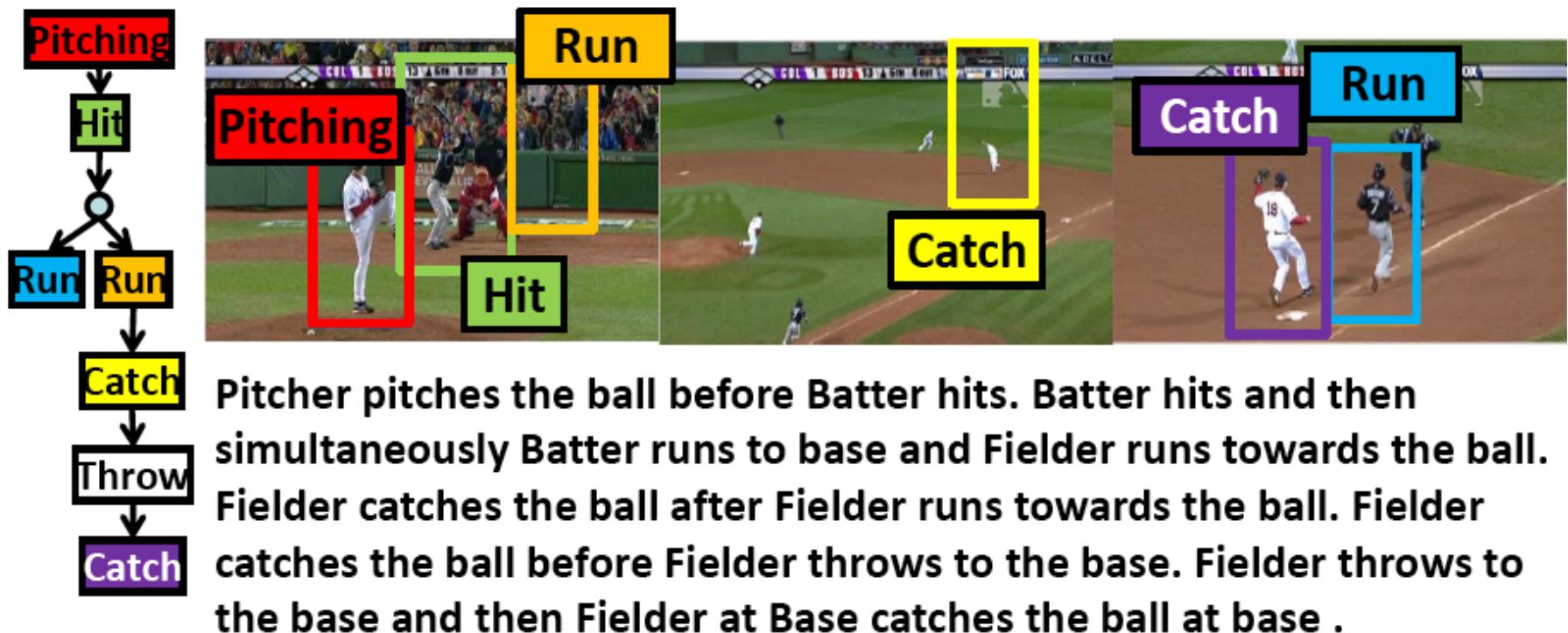


Sports Analysis

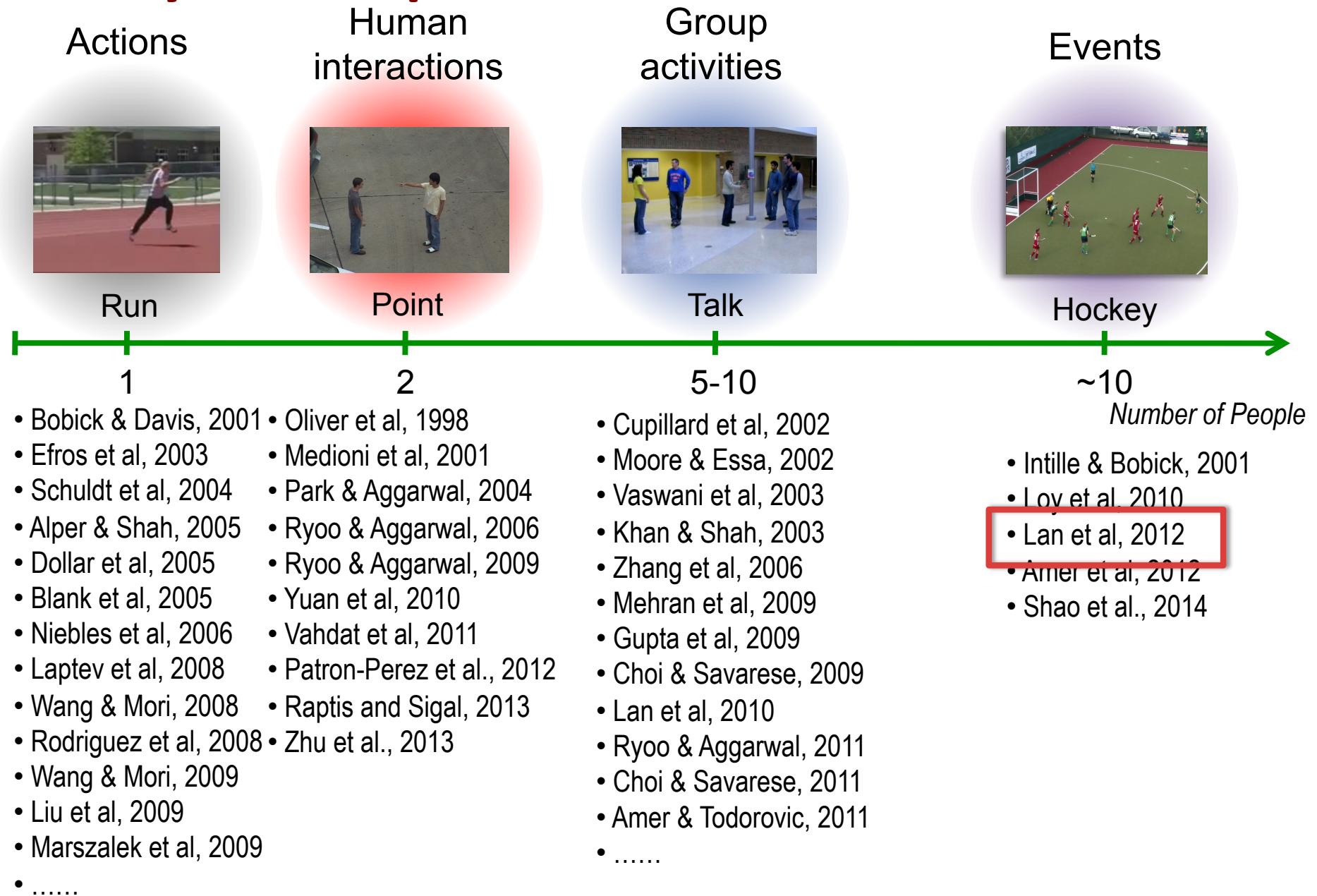


Lan, Sigal, Mori CVPR 2012

Text Generation



Activity landscape



Lan, Sigal, and Mori CVPR 2012

Role of Context in Actions



Is this a fallen person?

Role of Context in Actions



Who has the puck?

Semantic Descriptions of Videos

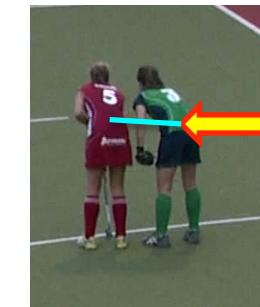
actions	social roles	event
walk run jog bend shoot dribble pass	attacker first defenders man-marking defend-space teammate	corner hit free hit attack play



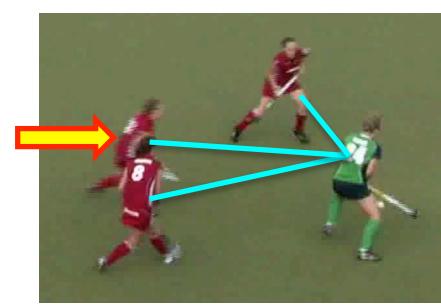
Social Roles

- Mid-level semantics that describe individual/group behaviors in the context of social interactions.

man-marking

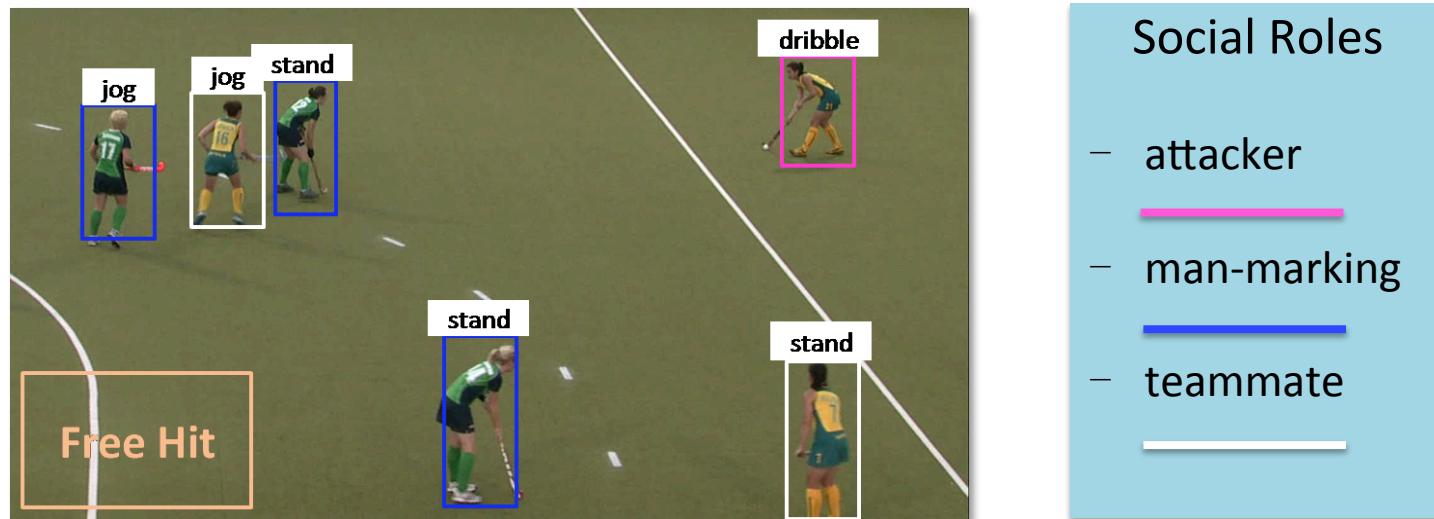


first defenders



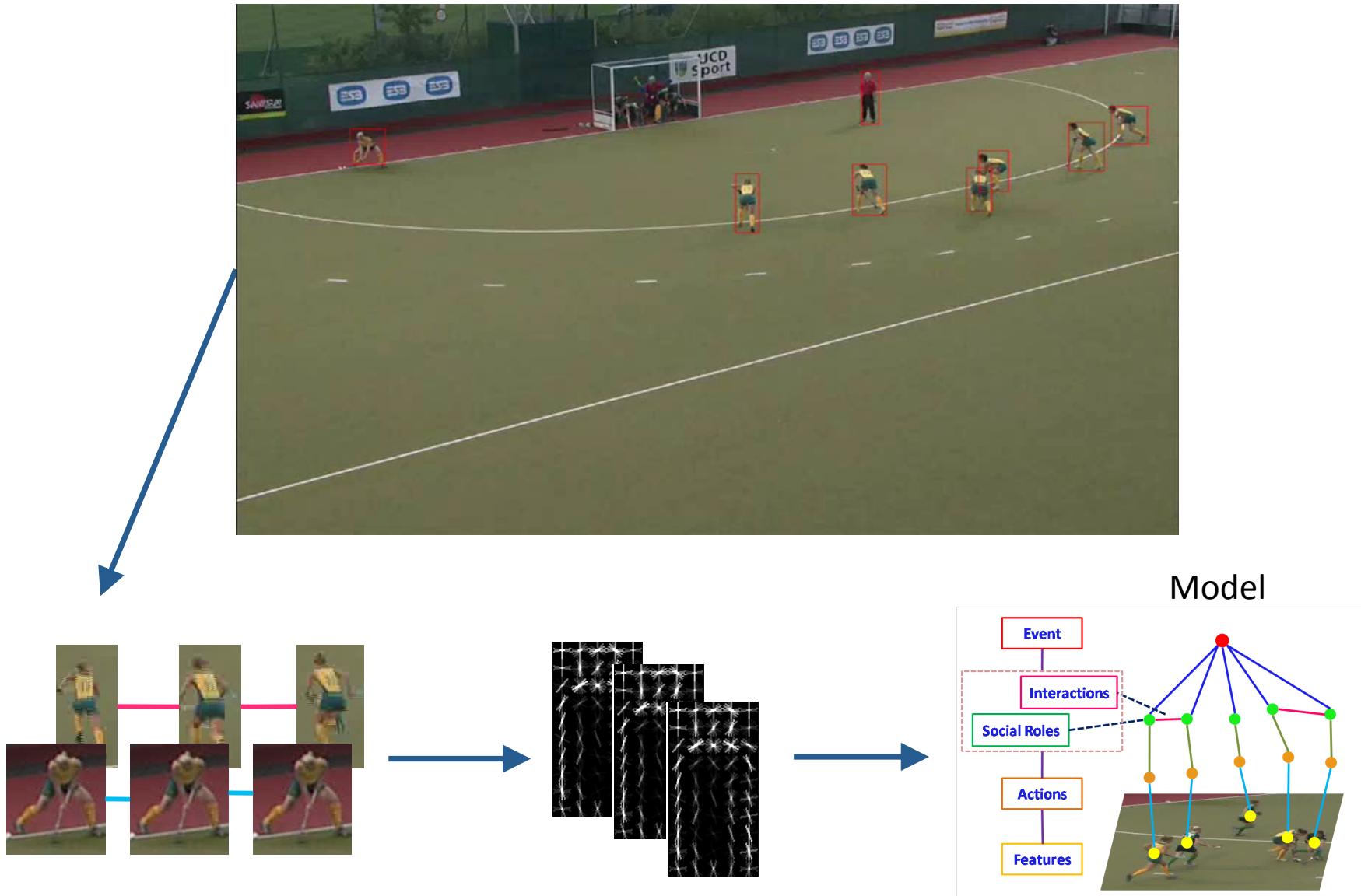
Goal

- Label all individuals' actions, social roles and the scene-level events.

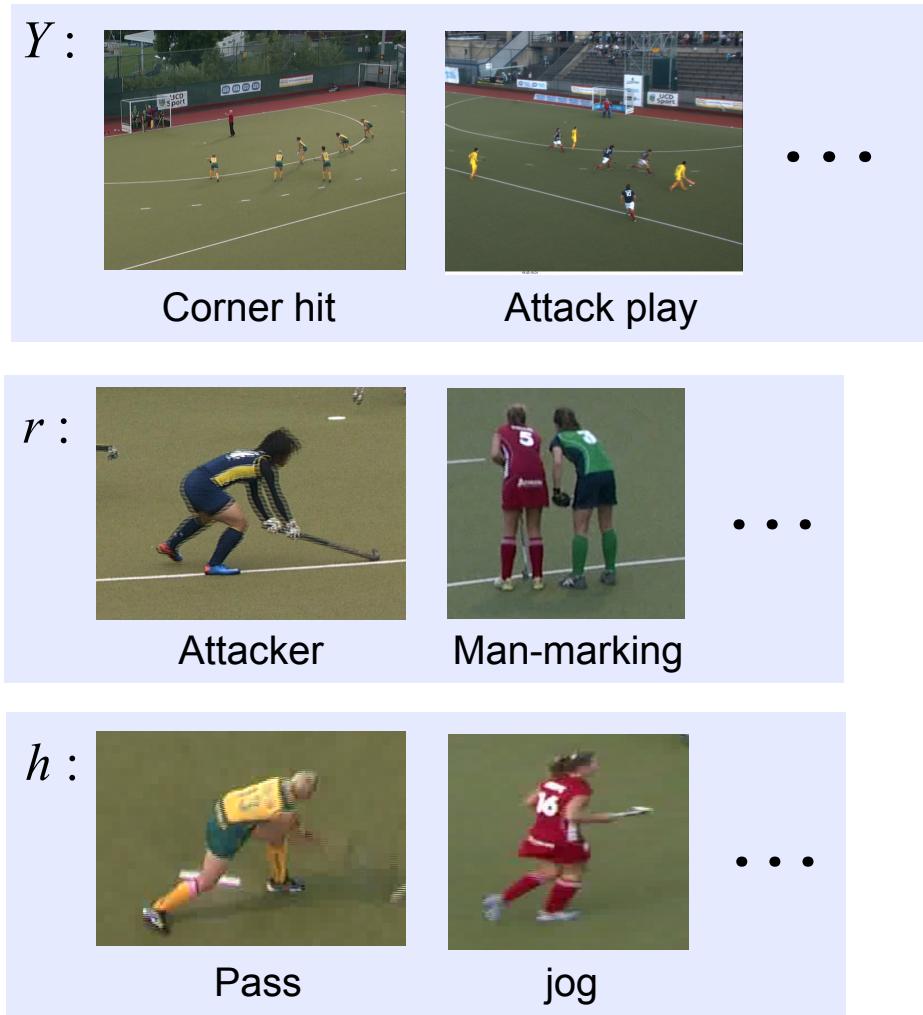


- Search for event/social role/action of interest
 - Who is the attacker? What's the overall game situation?

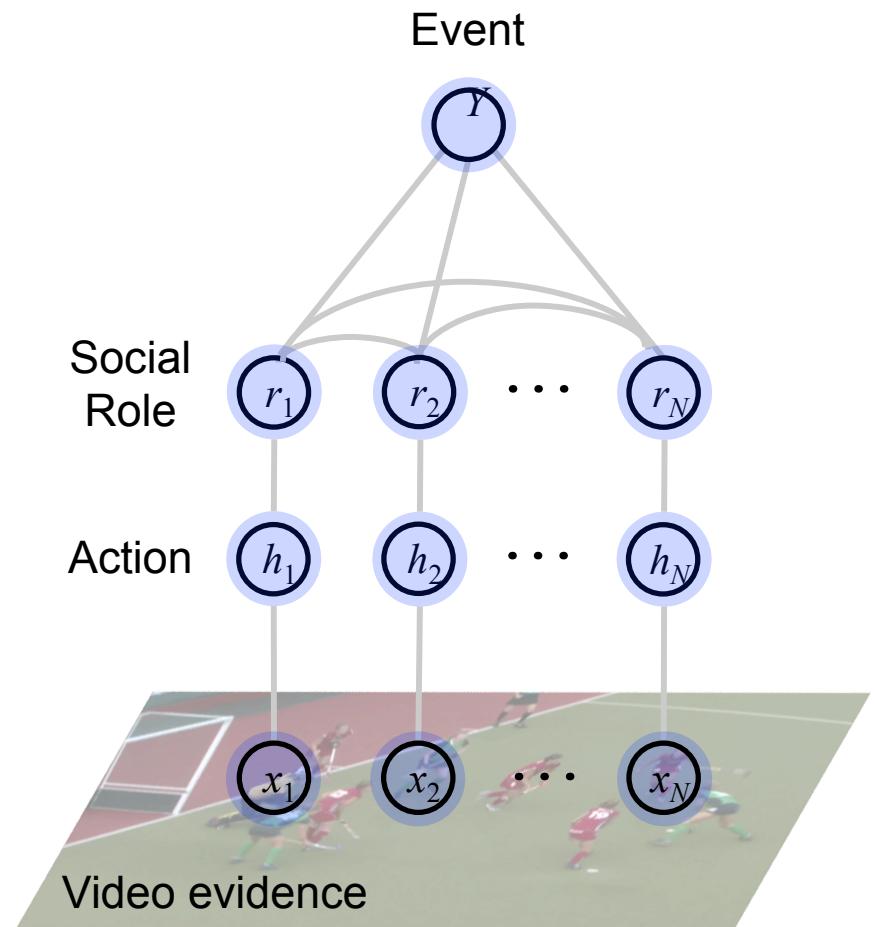
System Overview



Activity Hierarchy Model Representation

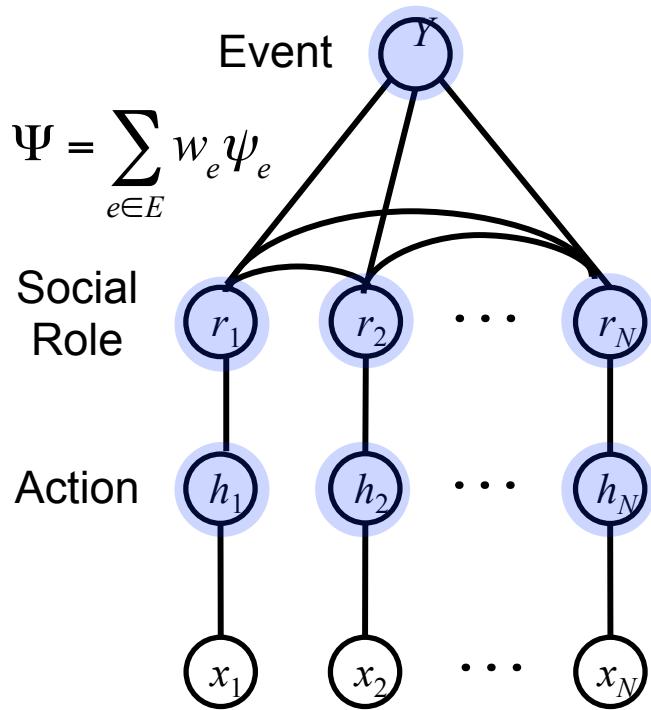


x : Concatenated HOG [Dalal & Triggs, 2005]



Lan, Sigal, and Mori CVPR 2012

Model Learning



Query for event: $loss = \Delta(y, y_i)$

$$\Delta(y, y_i) = \begin{cases} 1 & \text{if } y \neq y_i \\ 0 & \text{otherwise} \end{cases}$$

Query for social roles: $loss = \Delta(r, r_i)$

Query for actions: $loss = \Delta(h, h_i)$

Scene labeling: $loss = \Delta(y, y_i) + \Delta(r, r_i) + \Delta(h, h_i)$

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

s.t. $\forall i, y, r, h$

$$\mathbf{w}_{y_i r_i h_i} \cdot \psi_i - \mathbf{w}_{y r h} \cdot \psi_i \geq loss - \xi_i$$

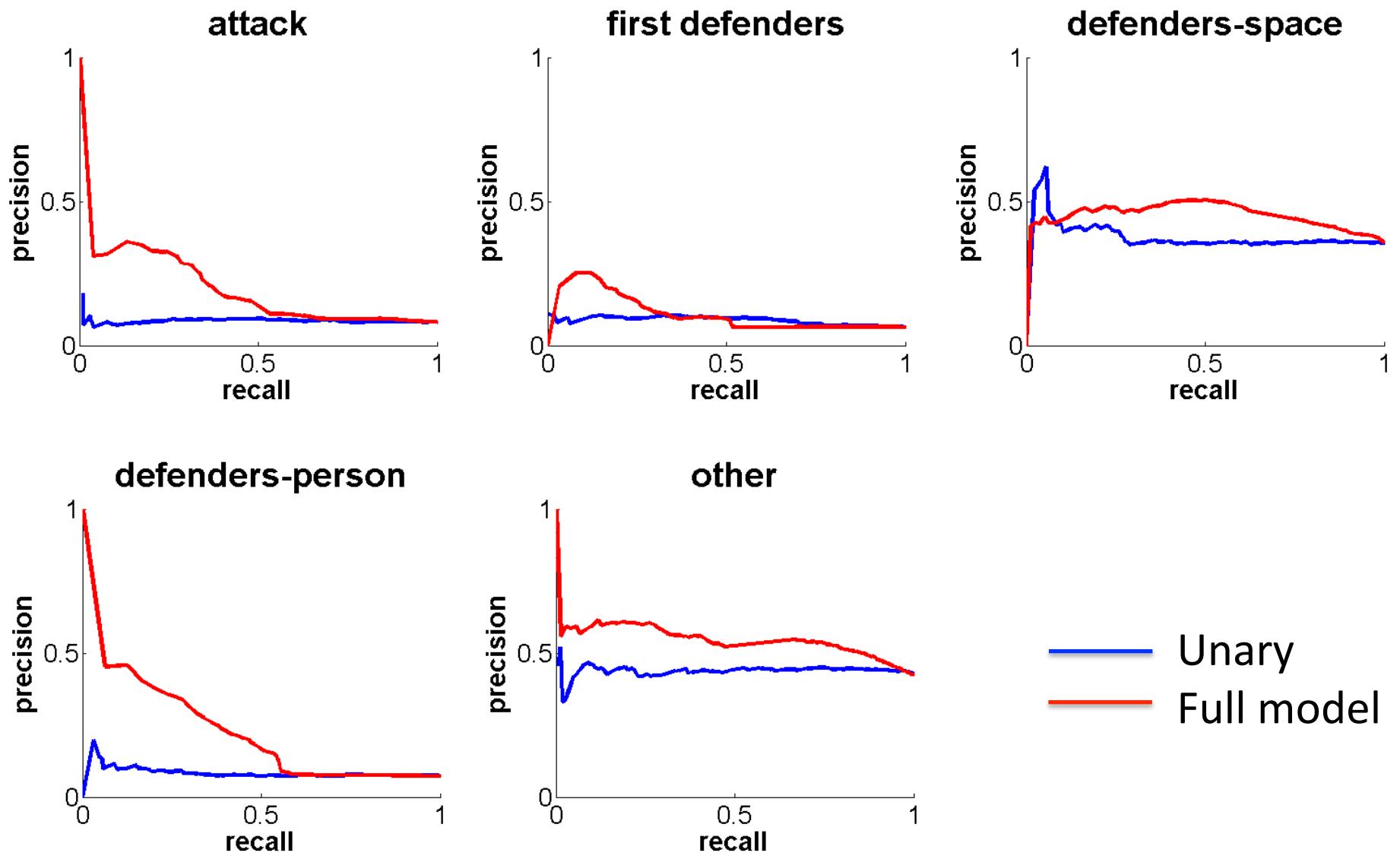
$$\forall i, \xi_i \geq 0$$

ESPN Broadcast Field Hockey Data



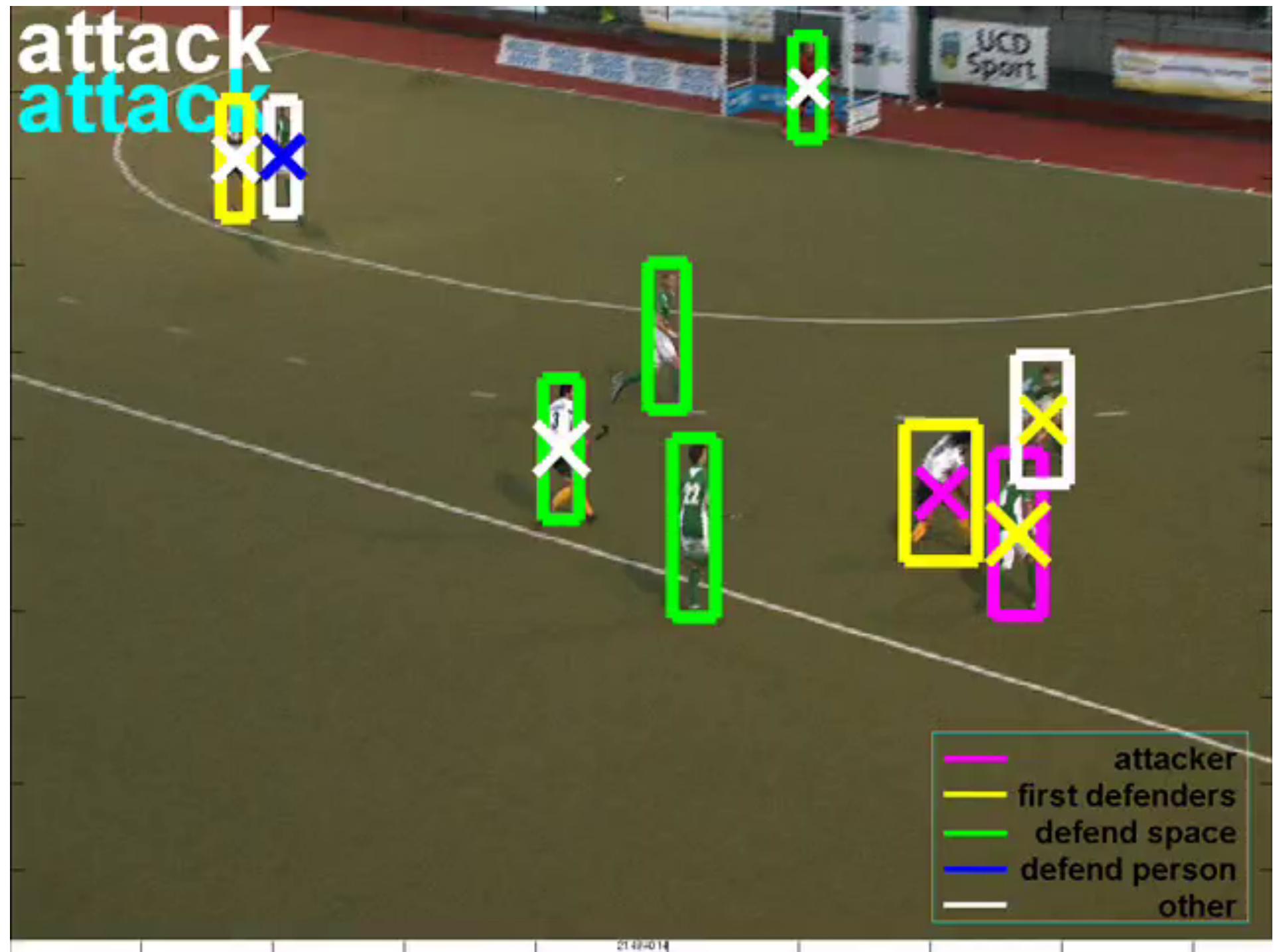
- 58 videos, 11 actions, 5 social roles, 3 scene-level events

Results – Query for Social Roles

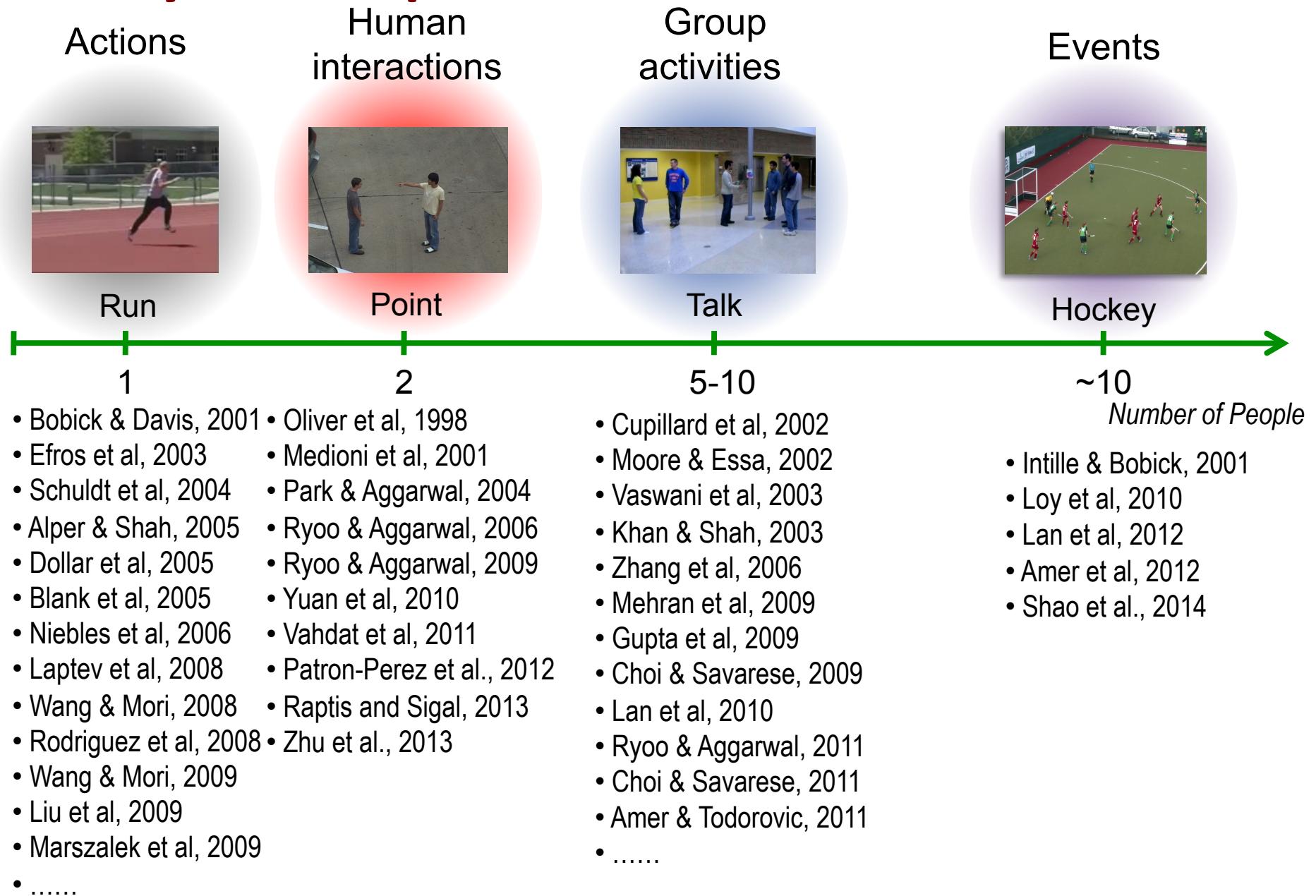


attack

attack



Activity landscape



Interaction Models I: Rule-based System

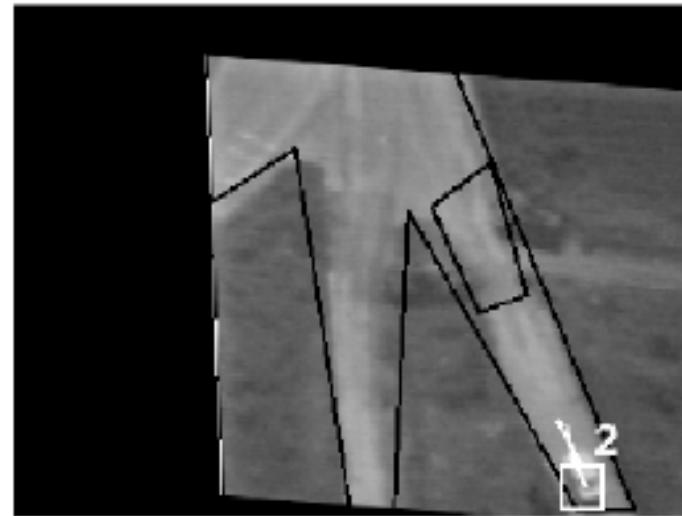
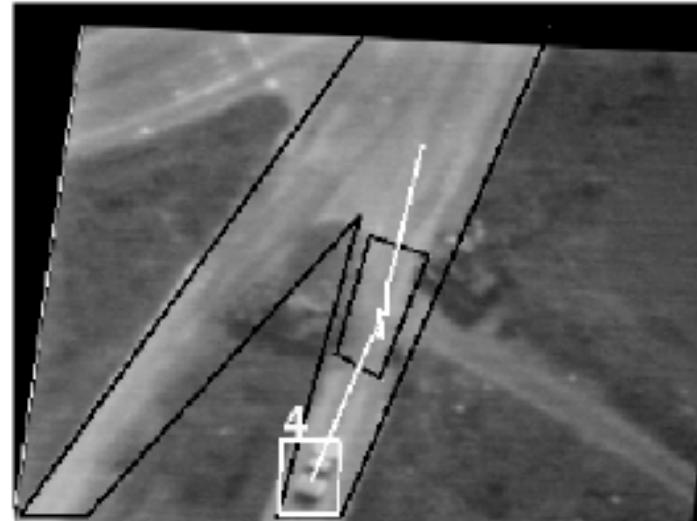
Detect and track moving objects

Manually identify key regions in scene

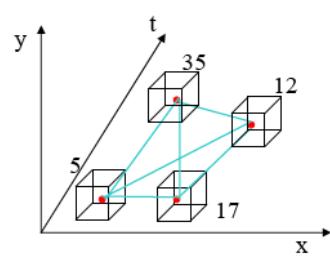
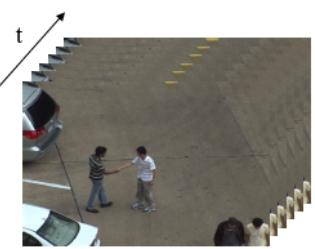
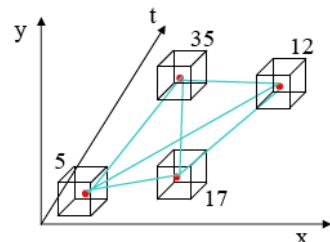
- E.g. road, checkpoint

Scenarios describe relative arrangements of objects in scene

- E.g. proximity of car to checkpoint
- Notions of scene **context**

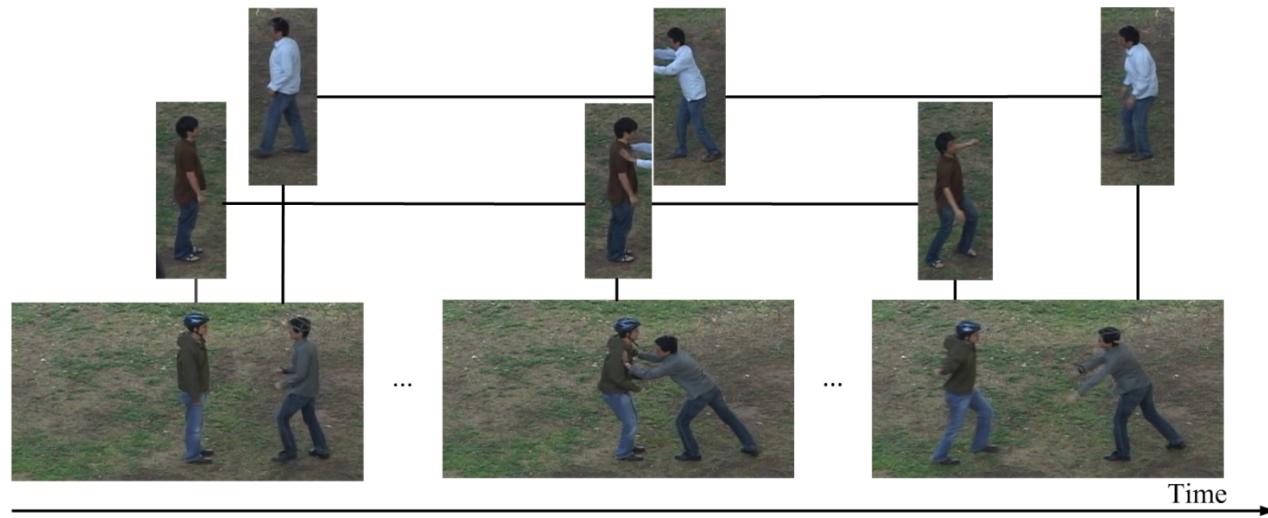


Interaction Models II: Local Feature Relationships



- UT Interaction dataset
- Focus on interactions between people
- Local feature approach
 - Define spatio-temporal relationships between points
 - Novel kernel for comparing sets

Interaction Models III: Key poses



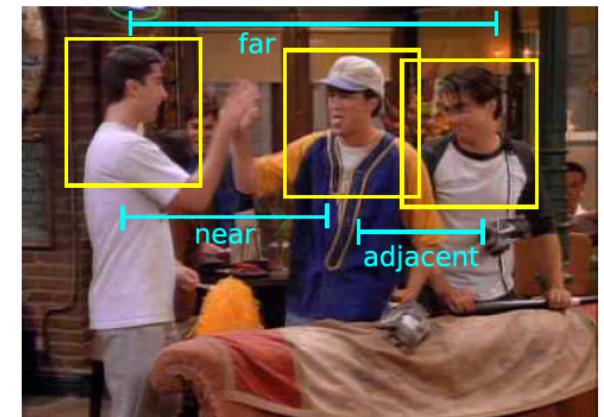
28

- Identify key points in an interaction
 - Poses of people
 - Relative positions of people
- Latent SVM formulation

Interaction Models IV: Upper-body Position / App.



- Human-human interactions in TV shows (dataset)
 - Upper-body detector
 - Estimate head orientation
 - Model combination of relative positions, facing, appearance



Interaction Models V: Human-object interactions

opening trunk



getting out of vehicle



entering a facility

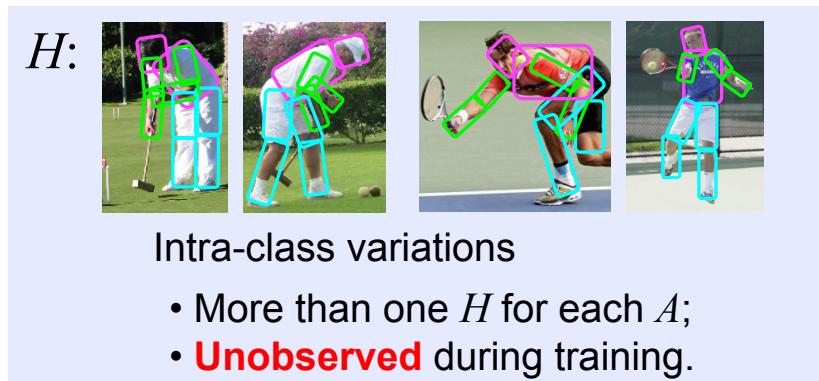
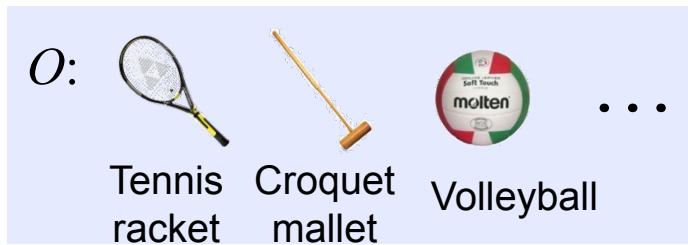
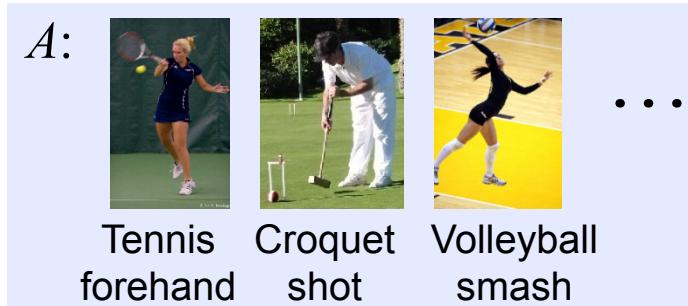


exiting a facility



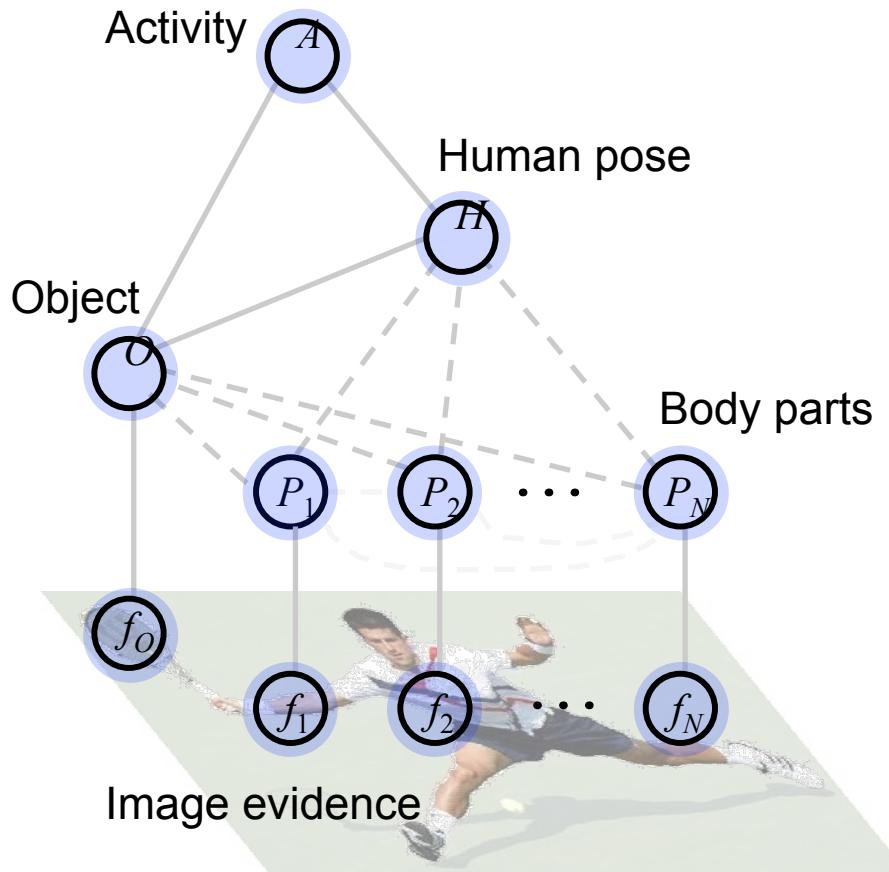
- Human-object interactions in surveillance video (VIRAT)
 - Detect moving objects in video
 - Features: velocity, object type, proximity to vehicle trunk, person is carrying object, person is in parking lot, ...
 - Structural learning for parameters

Interaction Models VI: Objects and Human Pose



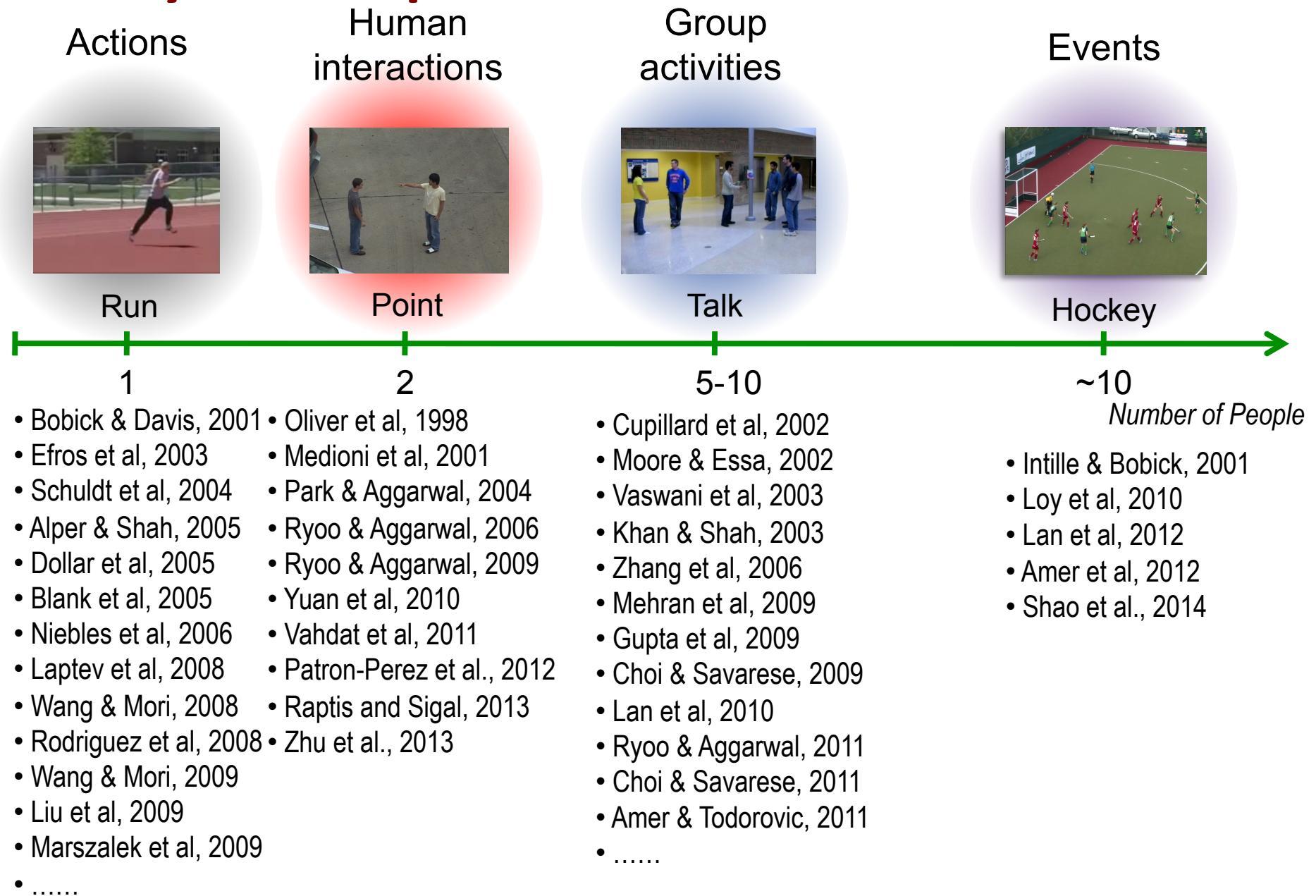
P : l_P : location; θ_P : orientation; s_P : scale.

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

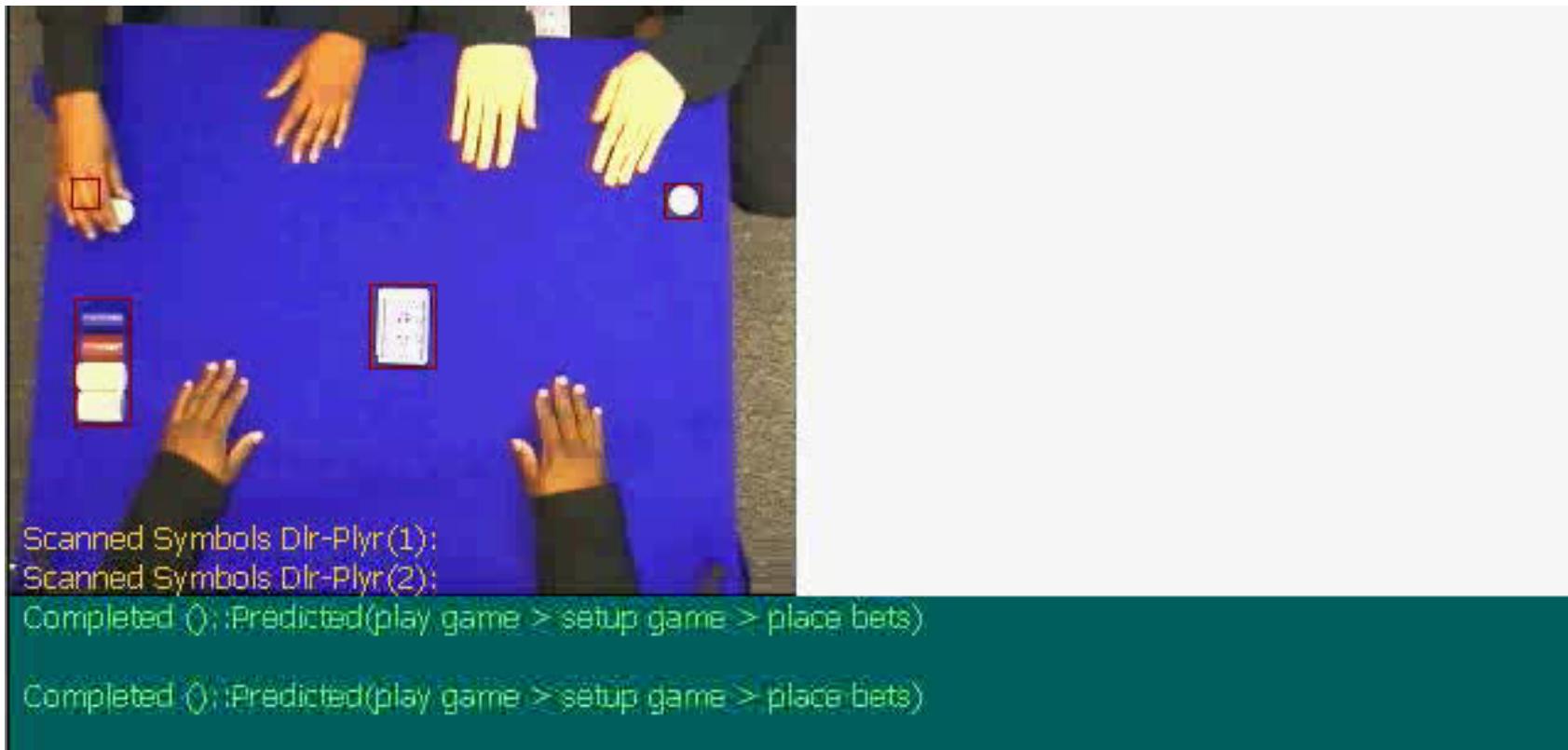


Yao and Fei-Fei, CVPR 2010

Activity landscape

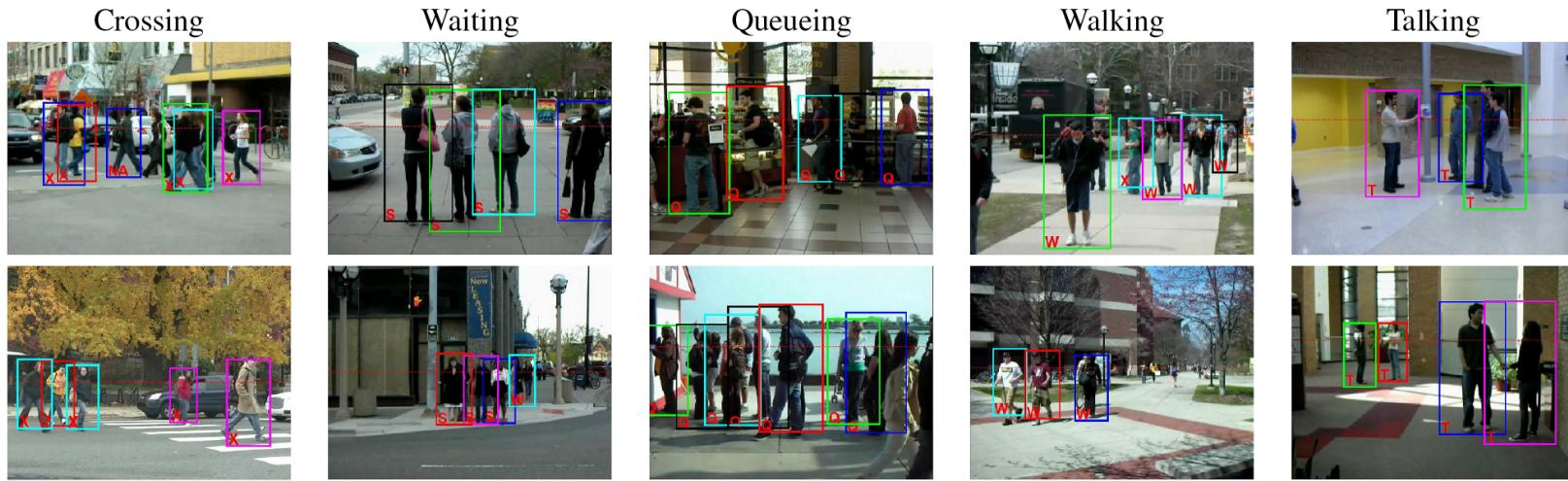


Group Models I: Stochastic Grammars

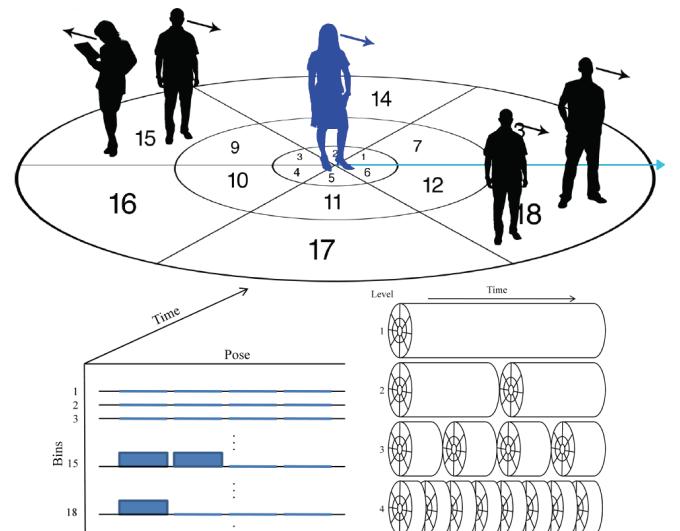


- Probabilistic grammar for describing domain
- Hand tracking, model object interactions, interactions between players in game
- Detect player strategies, roles, low-level actions

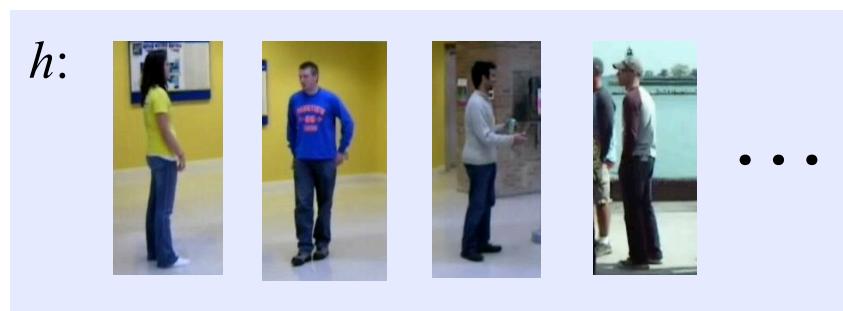
Group Models II: Person context



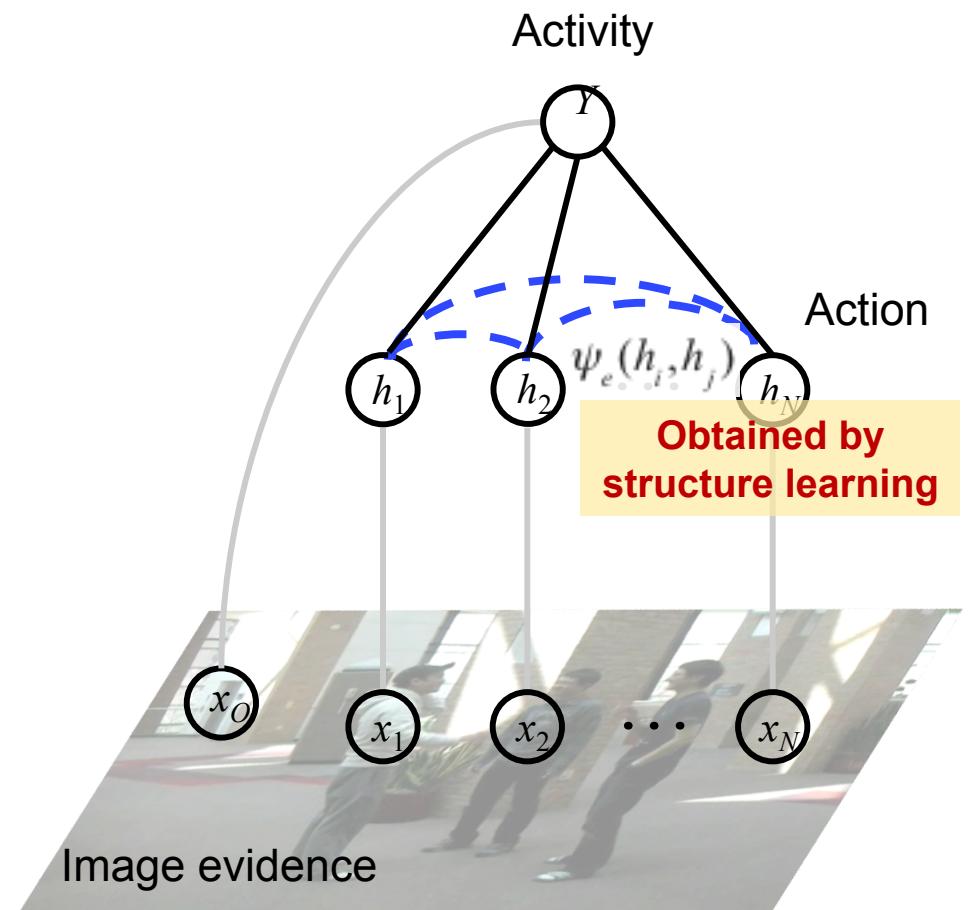
- Collective activity dataset
 - Pedestrian behaviours
 - Activities are context-dependent
 - Actions of others in the scene
 - Spatio-temporal context containing pose of other people in scene



Group Models III: Adaptive Structures

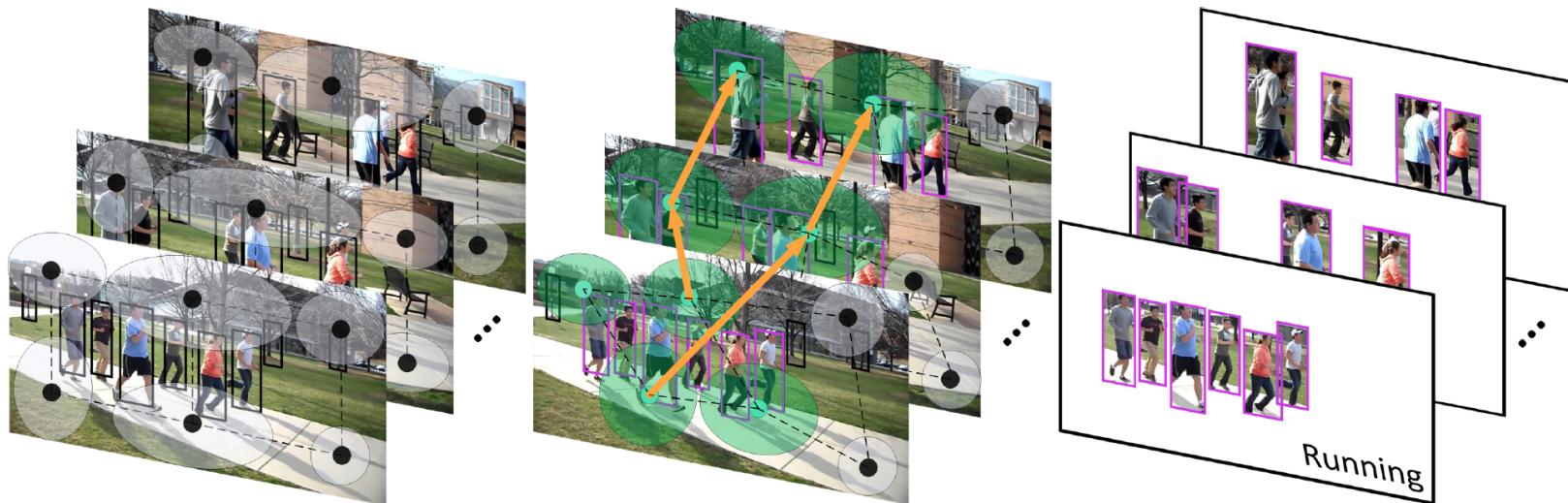


$x:$ HOG [Dalal & Triggs, 2005]



Lan, Wang, Yang, Mori, NIPS 2010

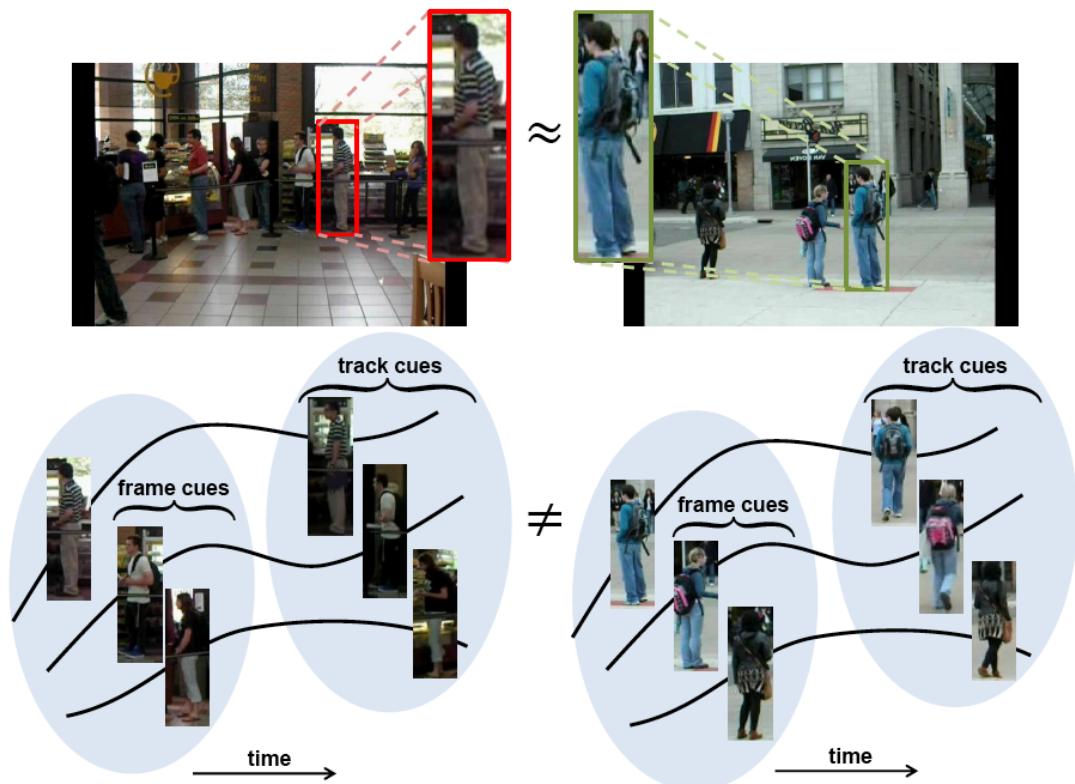
Group Models IV: Temporal extent



- “Chains model”
 - Find people whose actions are relevant to the group activity
 - Reason about temporal links, data association between noisy detections, temporal extent of a group activity

Group Models V: Linking tracklets

- Collective activity dataset
 - Features on individual detections
 - Features on tracks
 - Short tracklets
 - Data association terms
- Efficient inference



Group Models VI: Storyline Model

Captioned baseball videos in training

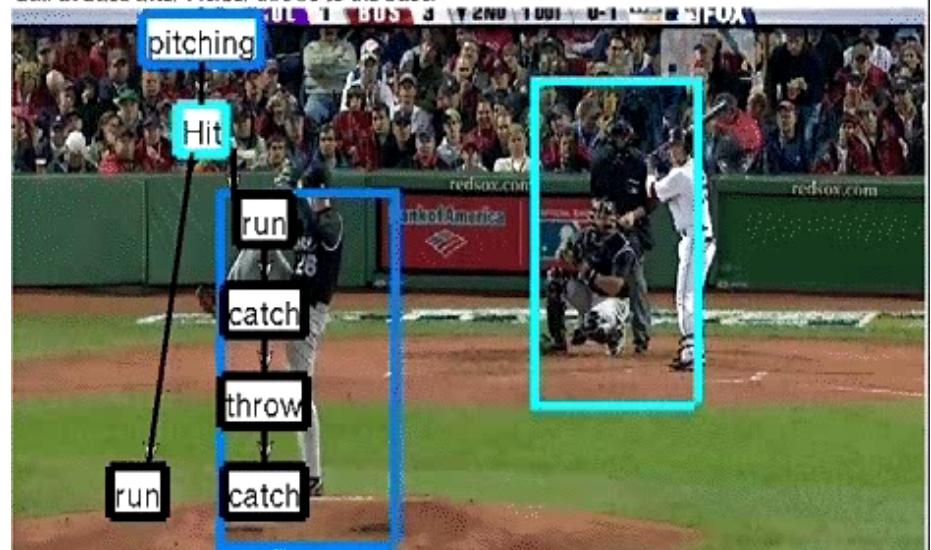
Build AND-OR graph representation of activities

- AND specifies elements of an activity that must occur
- OR allows variation in how an element appears

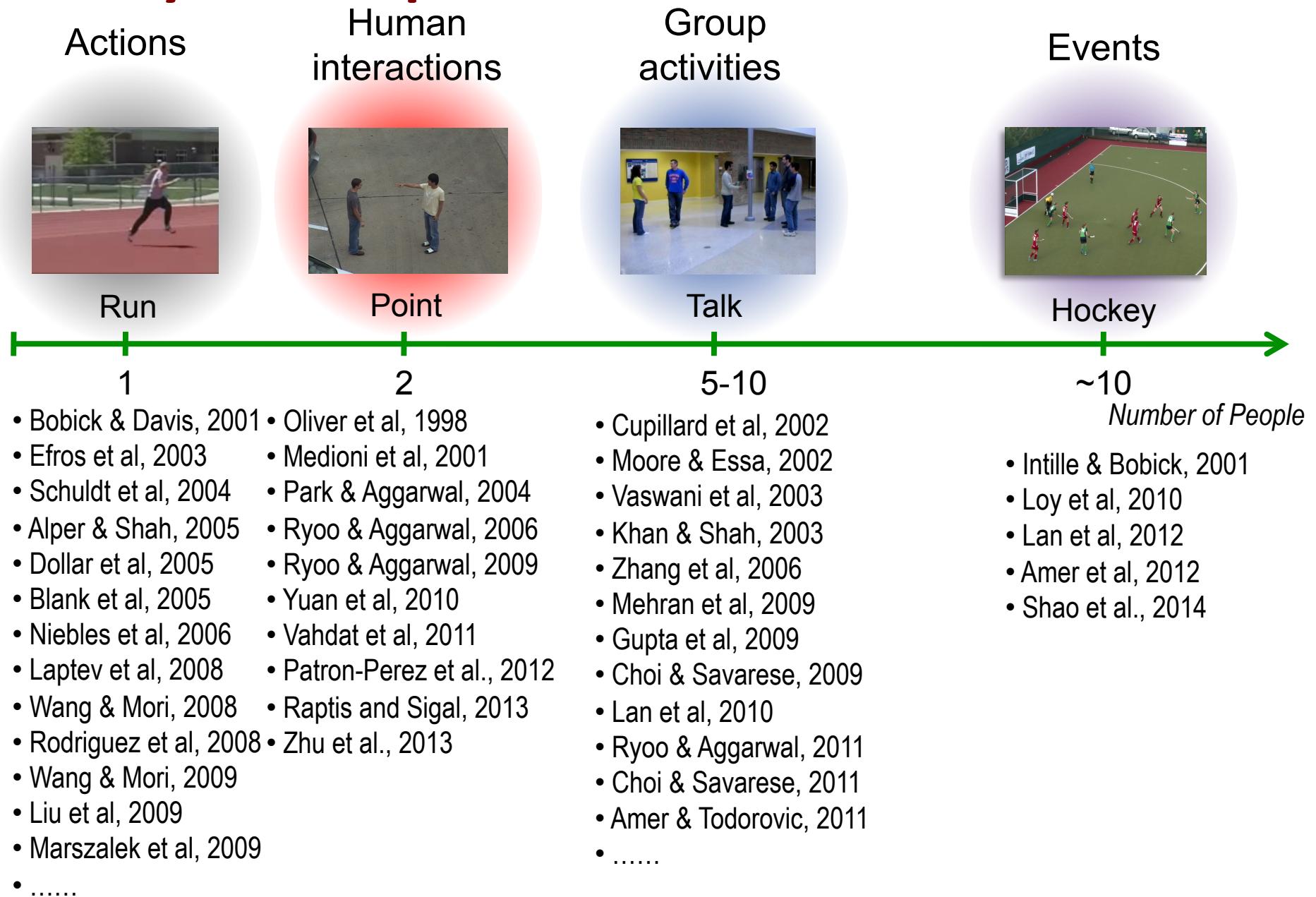
Describe low-level tracks using STIPs

Match tracks to actions in AND-OR graph

Pitcher pitches the ball before Batter hits. Batter hits and then simultaneously Batter runs to base and Fielder runs towards the ball. Fielder runs towards the ball and then Fielder catches the ball. Fielder catches the ball and then Fielder throws to the base. Fielder at Base catches the ball at base after Fielder throws to the base.

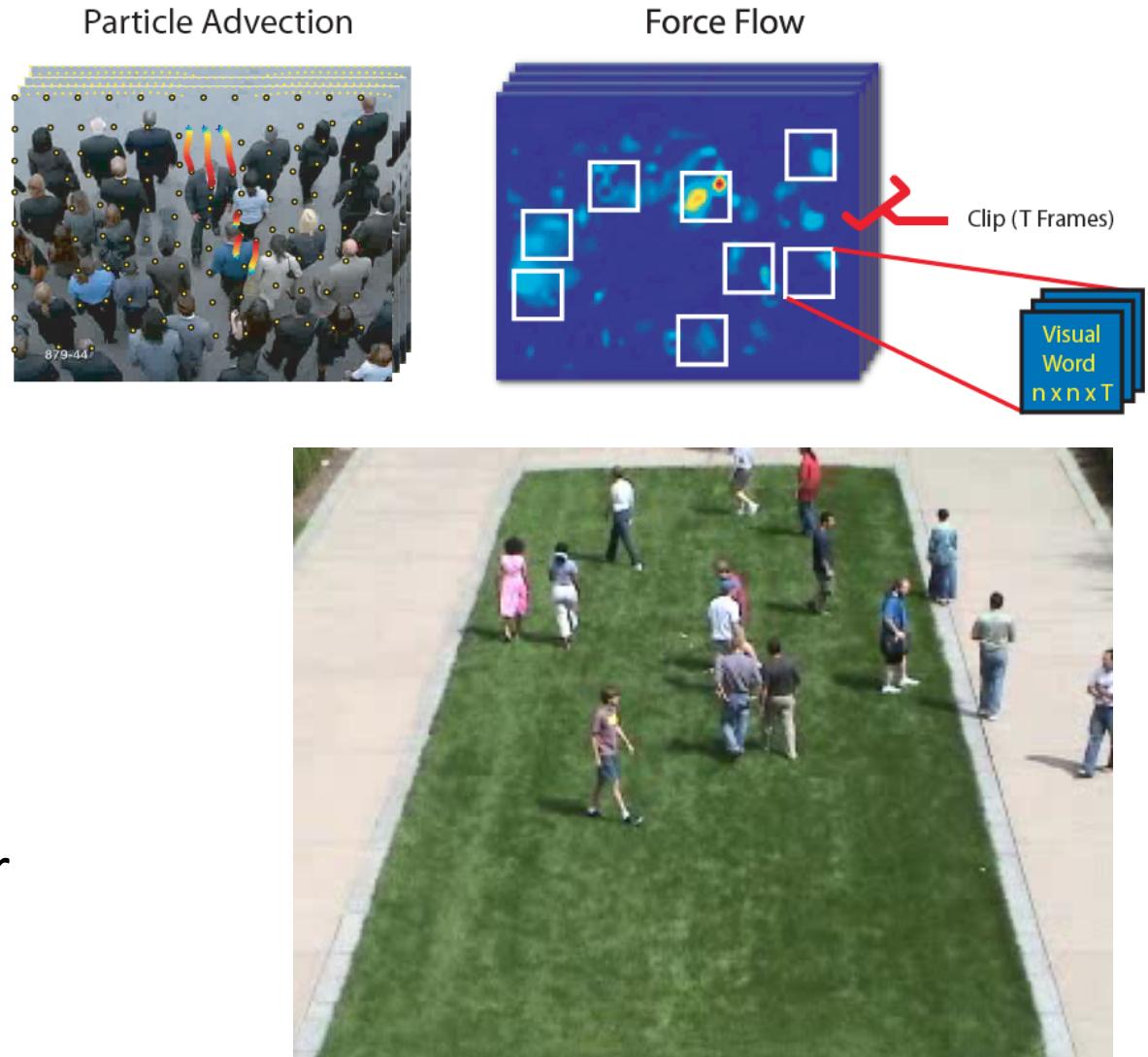


Activity landscape



Event Models I: Holistic flow models

- Global model of crowd
 - No tracking of individuals
 - Compute flow on dense grid
 - Interaction forces between particles estimated in local regions
 - Quantized into bag of words representation over the video



Mehran, Oyama, Shah, CVPR 2009

Event Models II: Group Flow



- Model group motion using Markov chain parameters – probabilistic spatial transition model (from tracklets of individuals)
- Define intra-group and inter-group measurements
 - Collectiveness: fit of individual members to group parameters
 - Stability: maintenance of nearest neighbours within a group, ...
 - ...

Shao, Loy, Wang, CVPR 2014

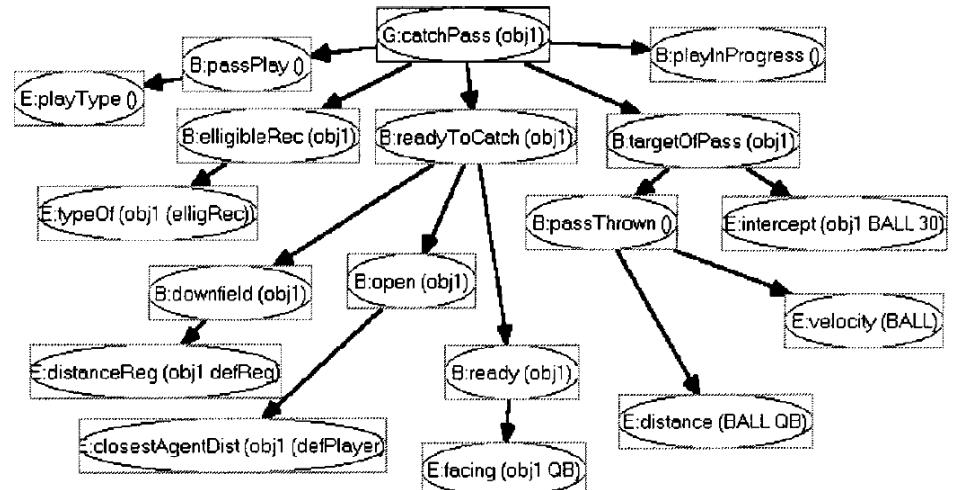
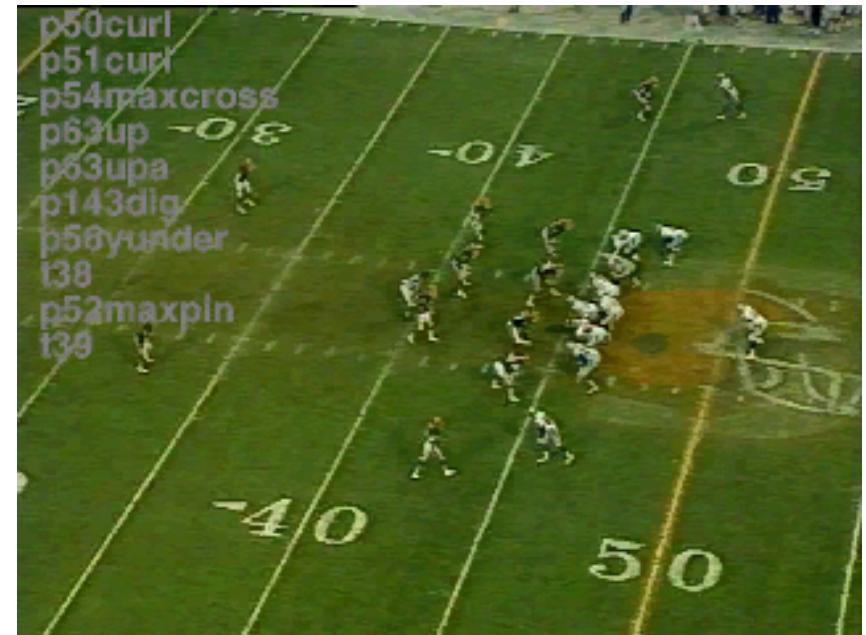
Event Models III: Bayes Nets

Detect and track players, ball

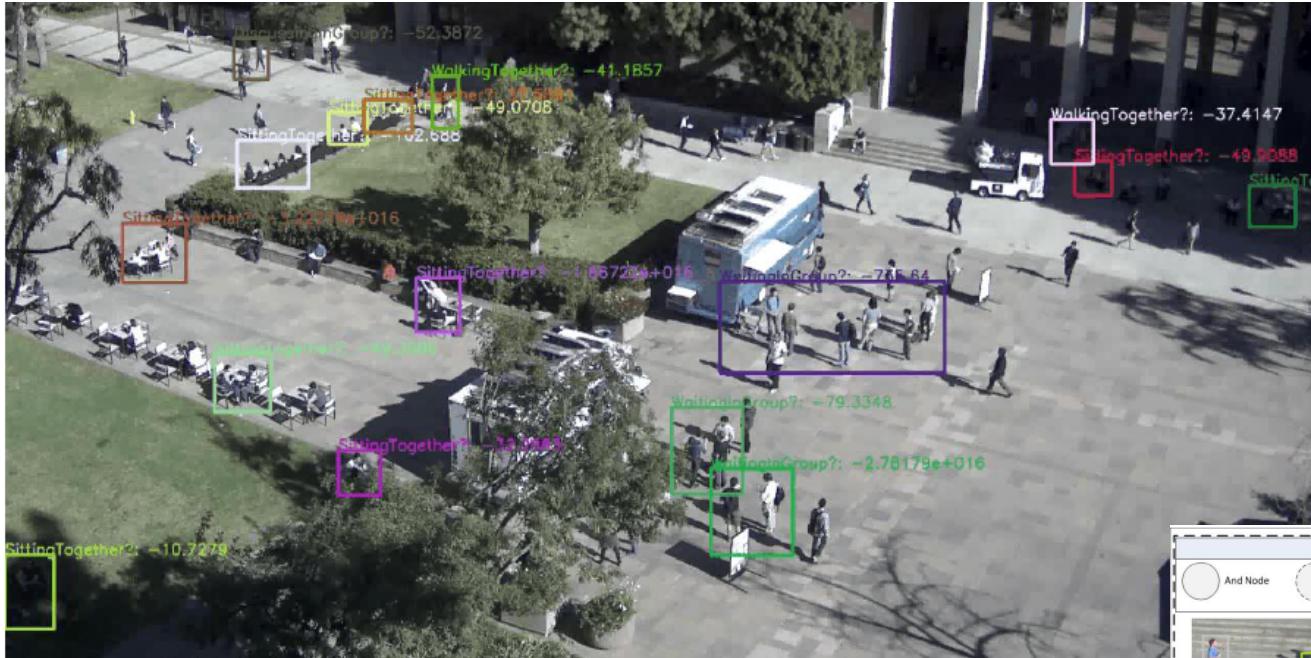
Low-level action detectors for individual players

Hand-constructed Bayes net for each activity

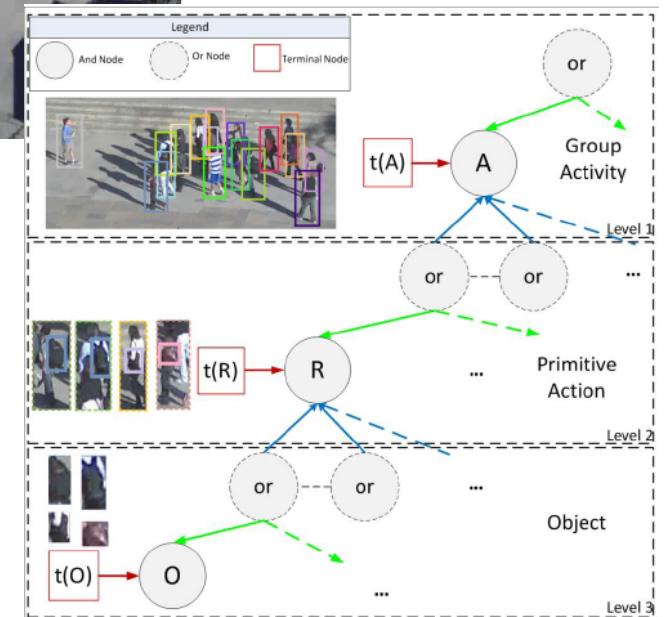
- Spatial and temporal relations between low-level actions



Event Models IV: Hierarchical models



- UCLA Courtyard dataset: actions, interactions, objects
- AND-OR graph to model relationships between group activity, primitive actions, objects
 - DPM detectors for people (poses) and objects
- Efficient inference



Conclusion

Broader scope

- Pairs of people, groups, levels of interaction



Types of models

- Probabilistic models
- Grammars
- Discriminative models
(LSVM)

Level of detail

- Individual people: actions, social roles, interactions
- Overall crowd flow

