**Scene dynamism** (vertical axis: Static scene → Dynamic scene)

**Number of group members** (horizontal axis: Dyadic interaction → Crowd interaction)

Rehg, CVPR13
Prabhaker, ECCV12
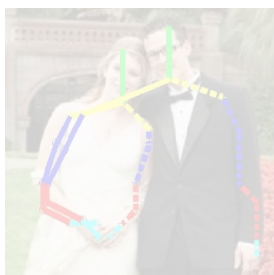Prabhakar, CVPR12
Patron-Perez, BMVC10

Lan, CVPR12
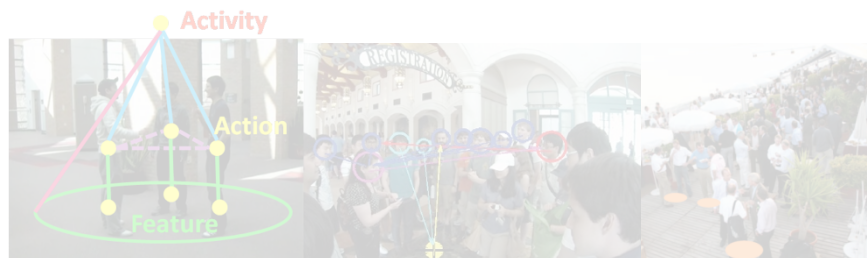Ramanathan, CVPR13
Antic, ECCV14

Ding, ECCV10
Choi, ECCV12, CVPR14
Direkoglu, ECCV12

Rodriguez, ICCV11a, ICCVb
Mehran, CVPR09
Alahi, CVPR14

Yang, CVPR12
Hoai, CVPR14

Fathi, CVPR12
Choi, ECCV14
Park, NIPS12, ICCV13

Cristani, BMVC11
Park, CVPR15
Arev, SIGGRAPH14

Wang, ECCV10
Gallagher, CVPR09

# Data-driven Approaches in Social Dynamics

**Scene dynamism**

Dynamic scene

Static scene

**Number of group members**

Dyadic interaction

Crowd interaction

Rehg, CVPR13
Prabhakar, ECCV12
Prabhaker, CVPR13
Patron, ECCV14

Rodriguez, ICCV11a, ICCVb
Mehran, CVPR09
Kim, CVPR14

Lan, CVPR12
Ramanathan, CVPR13
Antic, ECCV14

Ding, ECCV10
Choi, ECCV12, CVPR14

Activity

Action

Fathi, CVPR12
Choi, ECCV14
Park, NIPS12, ICCV13

Cristani, BMVC11
Park, CVPR15
Arev, SIGGRAPH14
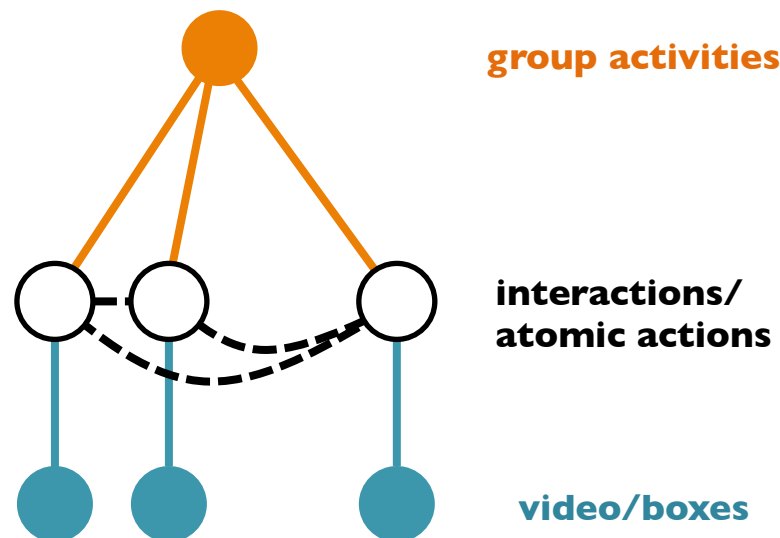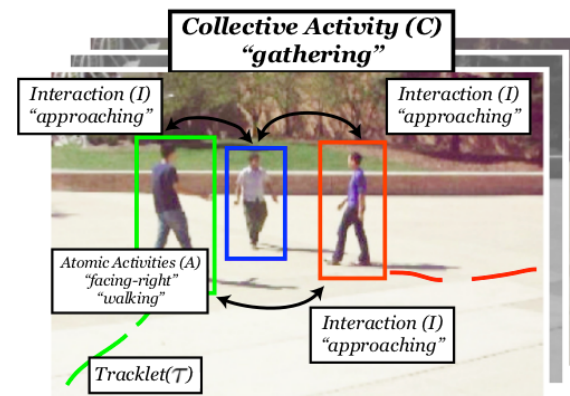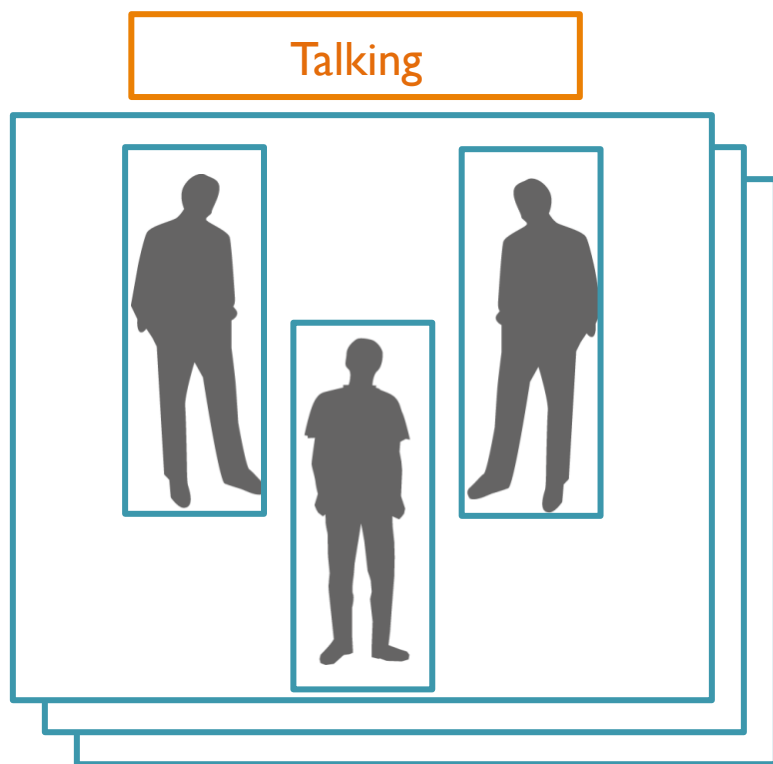
Wang, ECCV10
Gallagher, CVPR09

Yang, CVPR12
Hoai, CVPR14

Gathering

# Group Activities in Videos

[Choi VSWS09, CVPR11, ECCV12, Lan WSGA09, NIPS10, Khamis CVPR12, ECCV12]

Input: video and box tracks
Output: group activity labels in time



Collective Activity (C) "gathering"

Interaction (I) "approaching"

Interaction (I) "approaching"

Atomic Activities (A) "facing-right" "walking"

Interaction (I) "approaching"

Tracklet(T)

Talking



group activities

interactions/ atomic actions

video/boxes

# Collective Activities

Definition:

Activities that are defined or reinforced by the existence of a coherent behavior of a group of individuals in time and space.

**Waiting**

**Queuing**

**Talking**

# Individual Appearance



- Individual appearance/motion does not provide an important social signal to recognize collective activities.

# Crowd Context



- Contextual relationship among people is the key signal to recognize collective activities.

# Challenges

- Large <span style="color:red">intra-class</span> variation.
  - View point variation.
  - Number of group participants.
- <span style="color:red">Multiple groups</span> in the scene.
- Activities changes over time.

# Key Modules

- Crowd Context and the Representation.
  - Individual posture representation.
  - Encoding the context with posture representation.

- Exploit spatial-temporal correlations.
  - Utilize the structure in group activities.

# Key Modules

- **Crowd Context and the Representation.**
  - Individual posture representation.
  - Encoding the context with posture representation.

- Exploit spatial-temporal correlations.
  - Utilize the structure in group activities.

# Crowd Context: How-to

- Given trajectories (boxes over time) of people.
- Represent each individual box with a posture description:
  - Combination of view point and velocity.
  - Finite set of activity labels.
  - A bag of mid-level discriminative parts.
- Encode the context using a spatio-temporal descriptor.



Crowd Context 1

Crowd Context 2

Crowd Context 3

# Individual Posture Representation

- Extract a feature (e.g. HoG) in a bounding box.
- Represent the box in a posture space (e.g. SVM).



Dalal and Triggs, CVPR 2005

Front

**Simple and well-studied.**
**Require the definition and annotation of the posture space.**

# Encoding the Context

- Given a person of interest (anchor), aggregate the posture information of the others around the anchor person.

- Common ideas:
  - Define spatio-temporal support regions.
  - Pull the features in the space.

# Encoding Context: Action Context



(a)

(b) Focal person + Context → action

(c) action

$$C_i = \left[ \max_{j \in \mathcal{N}_1(i)} S_{1j}, \ldots, \max_{j \in \mathcal{N}_1(i)} S_{Kj}, \ldots, \max_{j \in \mathcal{N}_M(i)} S_{1j}, \ldots, \max_{j \in \mathcal{N}_M(i)} S_{Kj} \right]$$

# Encoding Context: Spatio-Temporal Local (STL) Descriptor



**e.g. SVM**

# Encoding Context: Learning the Contextual Model

# Learning the Contextual Model



(a) Waiting

(b) Talking



Anchors are looking upward.

Red: Facing forward,
Blue: Facing down,
Green: Facing right

# Collective Activity Dataset

**Crossing**  **Waiting**  **Queuing**  **Walking**  **Talking**



- 44 videos with multiple people.
- Crossing, Waiting, Queuing, Walking, Talking.
- Leave-One-Video-Out.

# Atomic Activity Feature v.s. Crowd Context

## Classification Accuracy

■ STIP (Baseline)



31.8

# Qualitative Examples



X: Crossing, S: Waiting, Q: Queuing, W: Walking, T: Talking

# Key Modules

- Crowd Context and the Representation.
  - Individual posture representation.
  - Encoding the context with posture representation.


- Exploit spatial-temporal correlations.
  - Utilize the structure in group activities.

# Hierarchical Activity Model

Input: video with tracklets

# Hierarchical Activity Model

# Hierarchical Activity Model



$$\Psi(C, I, A, O, f) =$$
$$\Psi(A, O) + \Psi(I, A, f) + \Psi(C, I) + \Psi(C, O) +$$
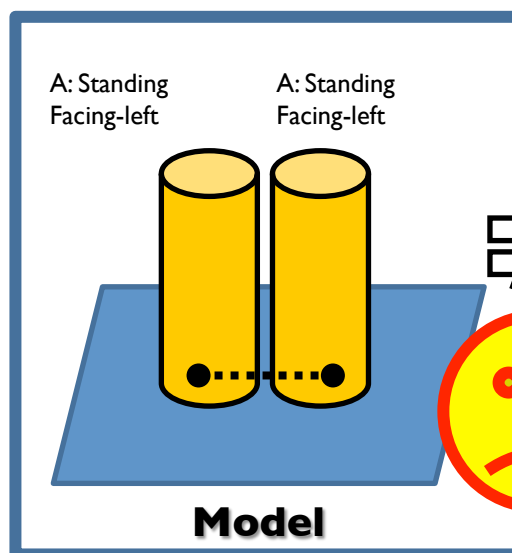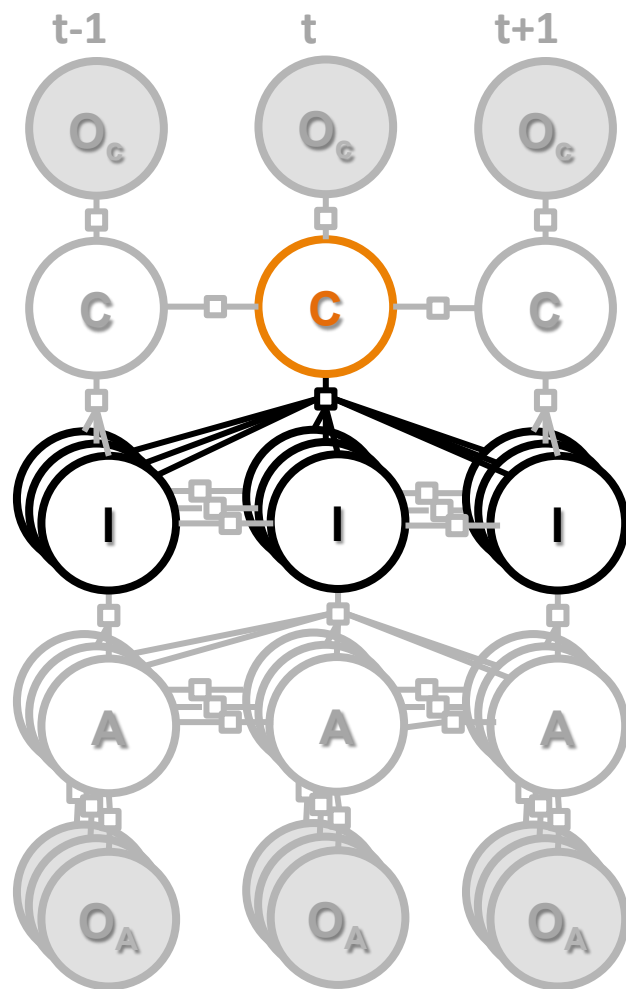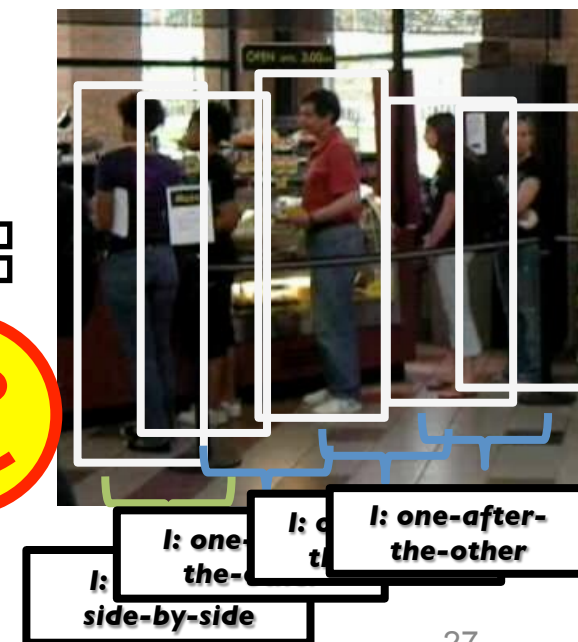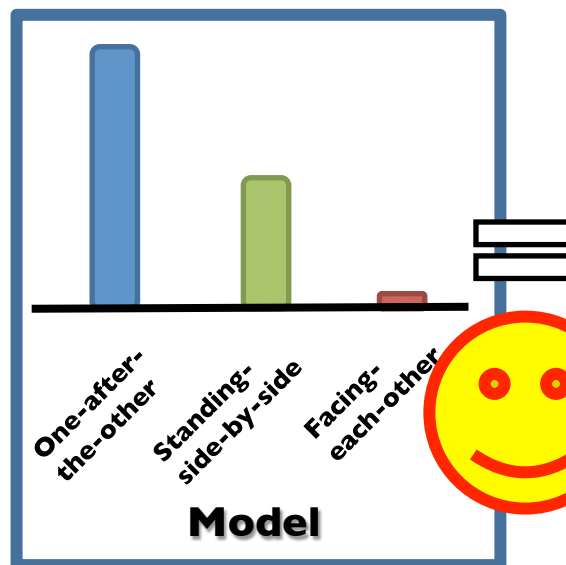$$\Psi(C) + \Psi(I) + \Psi(A) - c^{T}f, \ f \in S$$

# Atomic-Observation Potential

**t-1**   **t**   **t+1**

$$\Psi(C, I, A, O, f) =$$
$$\underline{\Psi(A, O)} + \Psi(I, A, f) + \Psi(C, I) + \Psi(C, O) +$$
$$\Psi(C) + \Psi(I) + \Psi(A) - c^T f, \ f \in S$$

Atomic Activity Models
- Action: BoW with STIP
- Pose: HoG

Dollar et al, 06; Niebles et al, 07

Dalal and Triggs, 05

# Interaction-Atomic Potential

$$\Psi(C, I, A, O, f) =$$
$$\Psi(A, O) + \Psi(I, A, f) + \Psi(C, I) + \Psi(C, O) +$$
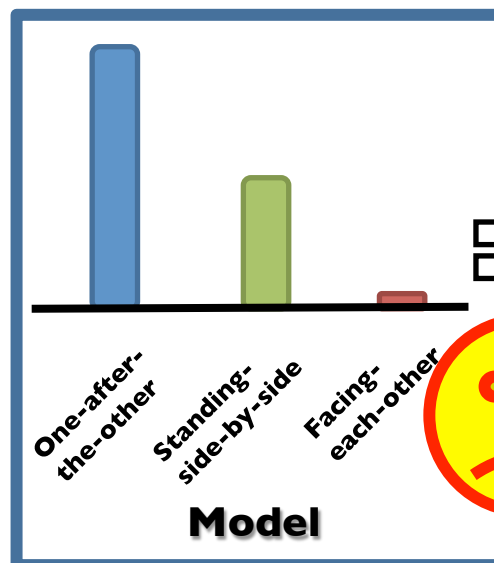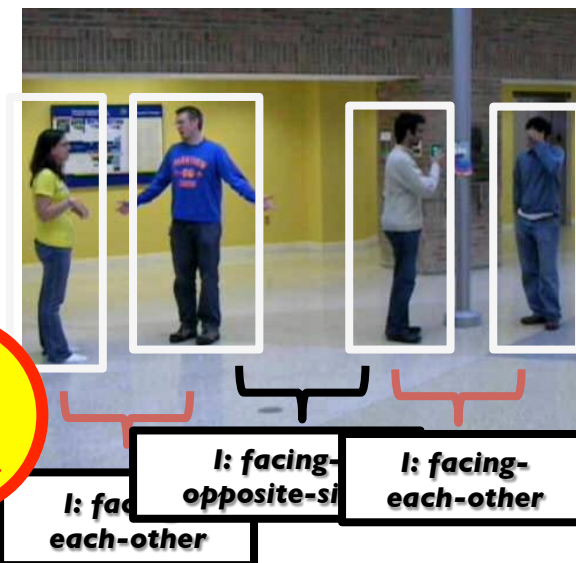$$\Psi(C) + \Psi(I) + \Psi(A) - c^T f, \ f \in S$$

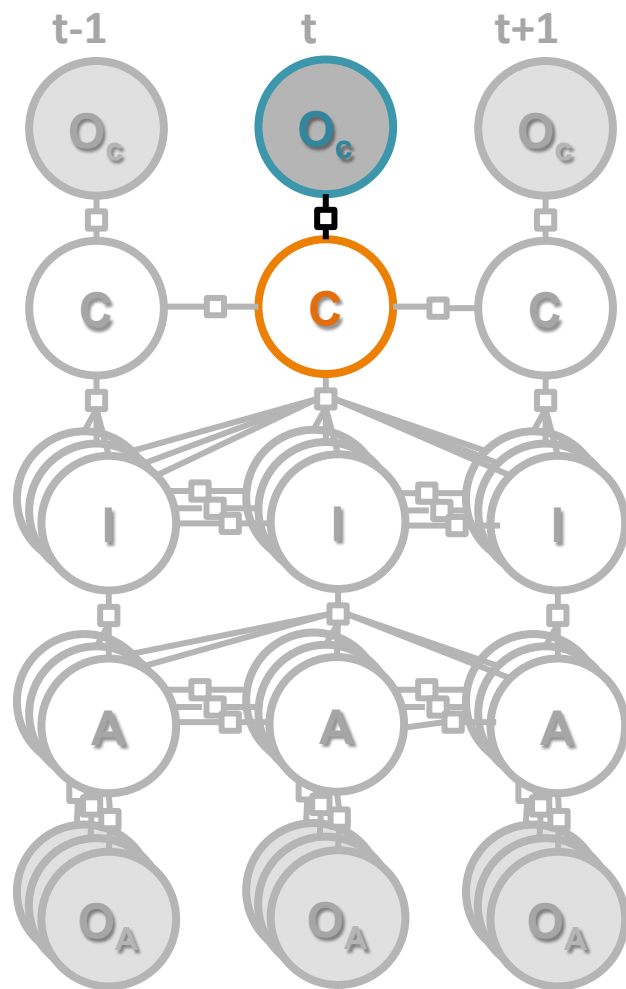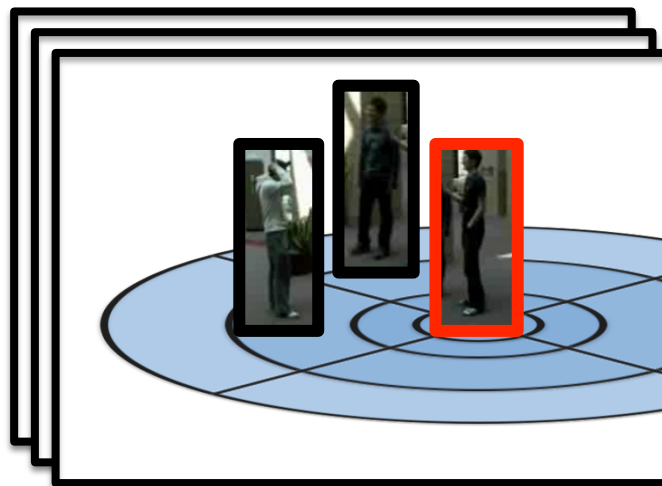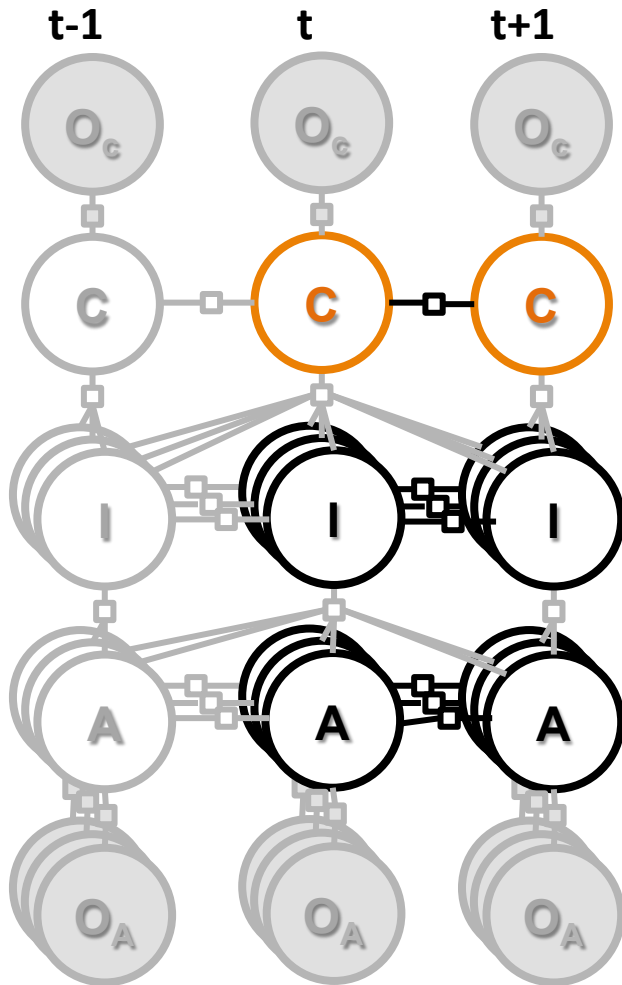**I: Standing-in-a-line**

A: Standing Facing-left

A: Standing Facing-left

**Model**

A: Standing Facing-left

A: Standing Facing-left
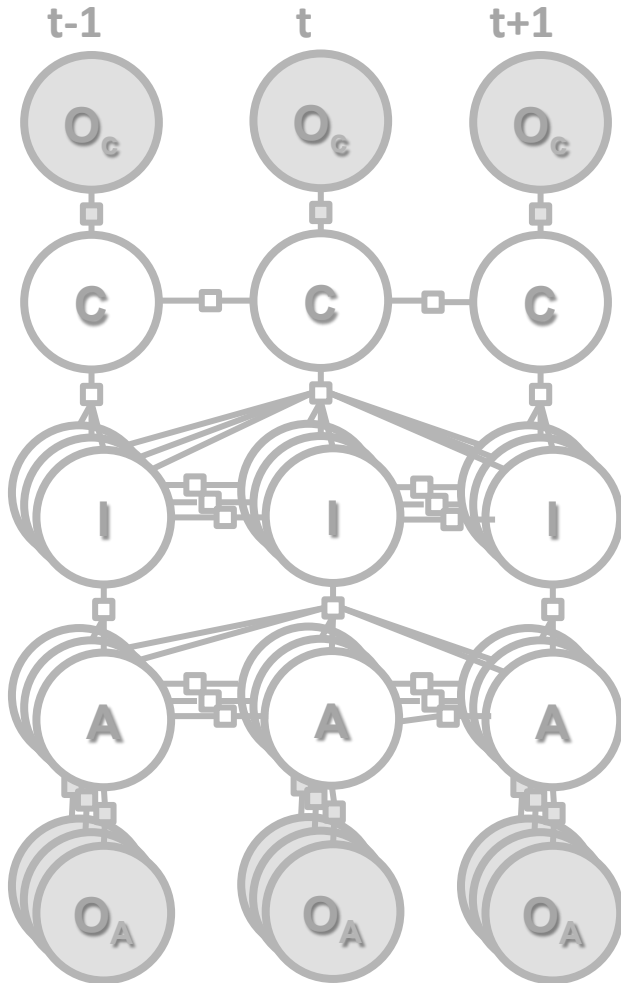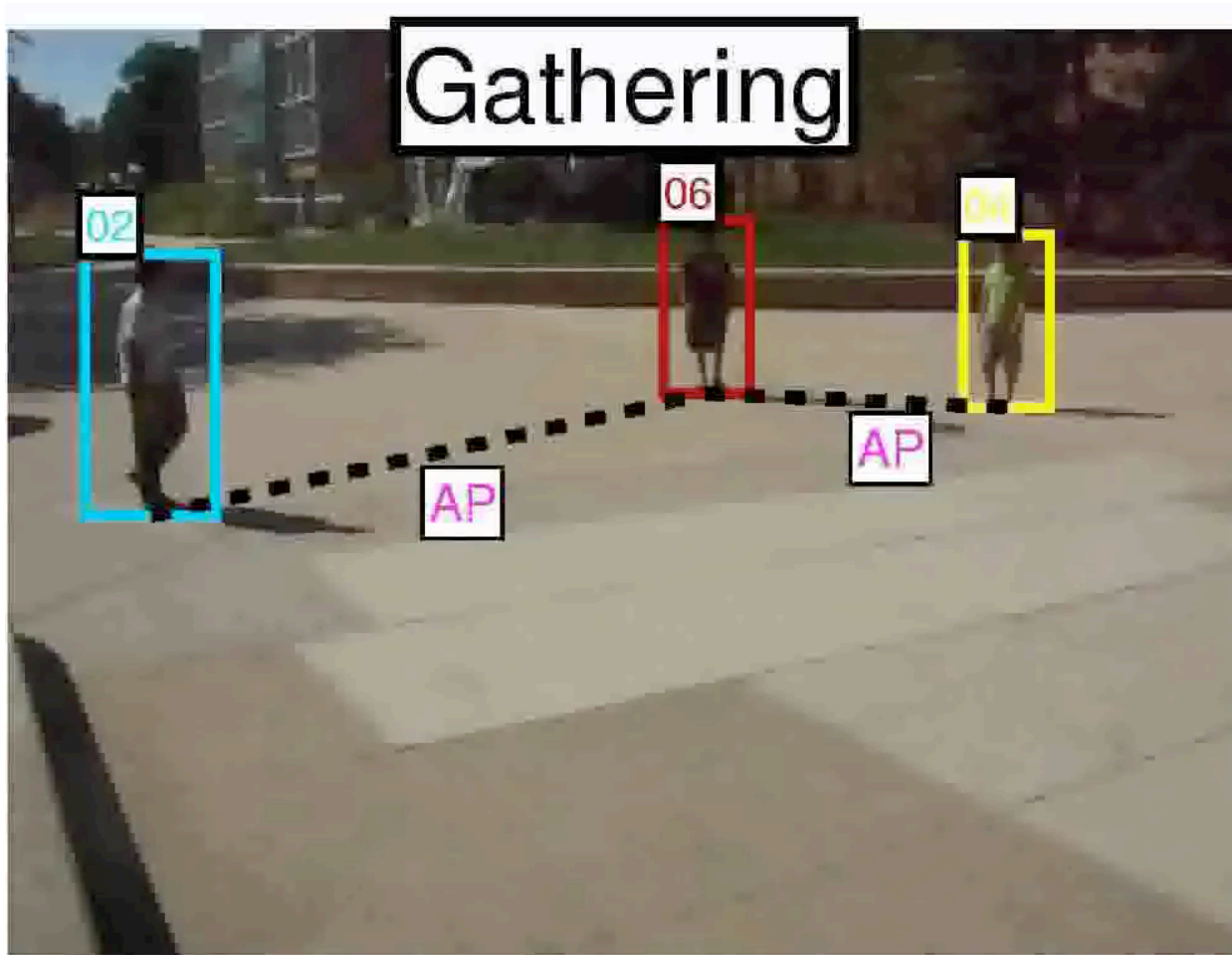
t-1   t   t+1

25

# Interaction-Atomic Potential

$$\Psi(C, I, A, O, f) =$$
$$\Psi(A, O) + \underline{\Psi(I, A, f)} + \Psi(C, I) + \Psi(C, O) +$$
$$\Psi(C) + \Psi(I) + \Psi(A) - c^T f, \ f \in S$$

**I: Standing-in-a-line**



A: Standing Facing-left    A: Standing Facing-left

**Model**

A: Standing Facing-right    A: Standing Facing-left

# Collective-Interaction Potential

$$\Psi(C,I,A,O,f) =$$
$$\Psi(A,O) + \Psi(I,A,f) + \underline{\Psi(C,I)} + \Psi(C,O) +$$
$$\Psi(C) + \Psi(I) + \Psi(A) - c^T f, \; f \in S$$

**C: Queuing**

**Model**

One-after-the-other

Standing-side-by-side

Facing-each-other

t-1   t   t+1

$O_C$   $O_C$   $O_C$

C   C   C

I   I   I

A   A   A

$O_A$   $O_A$   $O_A$

I: one-after-the-other

I: side-by-side

27

# Collective-Interaction Potential

$$\Psi(C, I, A, O, f) =$$
$$\Psi(A, O) + \Psi(I, A, f) + \underline{\Psi(C, I)} + \Psi(C, O) +$$
$$\Psi(C) + \Psi(I) + \Psi(A) - c^T f, \; f \in S$$

**C: Queuing**



One-after-the-other

Standing-side-by-side

Facing-each-other

**Model**

$\neq$

I: facing-each-other

I: facing-opposite-si

I: facing-each-other

# Collective-Observation Potential

**t-1**    **t**    **t+1**

$$\Psi(C,I,A,O,f) =$$
$$\Psi(A,O)+\Psi(I,A,f)+\Psi(C,I)+\underline{\Psi(C,O)}+$$
$$\Psi(C)+\Psi(I)+\Psi(A)-c^{T}f, \ f \in S$$

Collective Activity
- STL of all targets



Choi et al, 09

# Activity Transition Potential

$$\Psi(C, I, A, O, f) =$$
$$\Psi(A, O) + \Psi(I, A, f) + \Psi(C, I) + \Psi(C, O) +$$
$$\underline{\Psi(C) + \Psi(I) + \Psi(A)} - c^T f, \ f \in S$$

Smooth activity transition

# Trajectory Estimation



**t-1**    **t**    **t+1**

$\Psi(C, I, A, O, f) =$
$\quad \Psi(A, O) + \Psi(I, A, f) + \Psi(C, I) + \Psi(C, O) +$
$\quad \Psi(C) + \Psi(I) + \Psi(A) - c^T f, \ f \in S$

# Training the Graphical Model

- Model weights can be learned in a Max-Margin framework using Structural SVM.

$$\min_{\mathbf{w}, \boldsymbol{\xi}} \quad \frac{1}{2}\|\mathbf{w}\|^2 + \frac{C}{n}\sum_{i=1}^{n}\xi_i, \;\; \text{s.t.} \;\; \forall i, \xi_i \geq 0$$

$$\forall i, \; \forall \mathbf{y} \in \mathcal{Y} \setminus \mathbf{y}_i : \; \langle \mathbf{w}, \delta\Psi_i(\mathbf{y})\rangle \geq \triangle(\mathbf{y}_i, \mathbf{y}) - \xi_i$$

Tsochantaridis et al, 2004

# Example Classification Result



**Interaction labels**
AP: approaching
FE: facing-each-other
SR: standing-in-a-row
...

# Example Classification Result



**Atomic Activities**
Action:
W - walking
S – standing

Pose (8 directions)
L - left
LF– left/front
F – front
RF- right/front
etc.

# Example Classification Result



**Pair-Interactions**
- AP: approaching
- .....
- FE: facing-each-other
- SS: standing-side-by-side
- SQ: standing-in-a-queue

# Example Classification Result

# Learning Latent Constituents for Group Activity Recognition (Antic ECCV14)

Input: video and box tracks
Output: group activity labels in time



Talking

group activities

latent constituents

video/boxes

# Meaningful Parts of Group Behavior

# Meaningful Parts of Group Behavior

# Meaningful Parts of Group Behavior

# Less Meaningful Parts



Less meaningful parts

# Learning Mid-level Constituents



$$\Delta(\mathbf{f}_i, \mathbf{f}_v) := d(\mathbf{f}_i, \mathbf{f}_v) + \lambda_x \|\mathbf{x}_i - \mathbf{x}_v\| + \lambda_s |s_i - s_v|.$$

# Encoding Social Signal with Latent Constituents

# Behavior Recognition with Constituents

# Behavior Recognition with Constituents



crossing    waiting    queueing    talking    dancing    jogging

# Behavior Recognition with Constituents

# Behavior Recognition with Constituents



crossing  waiting  queueing  talking  dancing  jogging

# Quantitative Evaluation



crossing  waiting  queueing  talking  dancing  jogging

**Holistic approach (full b-boxes):**

70.4%
83.3%

**Latent constit's (functional grouping):**

75.1%
90.1%

■ 5 Activities Dataset  ■ 6 Activities Dataset

# Social Role Discovery (Ramanathan CVPR13)

Input: videos with event labels, box tracks
Output: groups with activity label



BIRTHDAY
b'day
person
parents
friends
guests

WEDDING
bride
groom
priest
grooms men
brides maid

event

social roles

video/boxes

Wedding

Role1
Role2
Role3

Wedding

Role1
Role2
Role3

# Humans in Social Setting



priest claps

bride gets ring
from bridesmaid

groom gets ring
from groomsman

bride and
groom kiss

bride and groom
exchange rings

# Goal: Identify social roles

# Problem setup



parent

b'day boy

friends

guest

person-specific role features    inter-role interaction features

# Social role model

# Social role model

## Person-specific features

- Role specific **actions**
  - **HOG3D** from person tube
  - **Trajectory** of person
- **Color** and **Gender** features
- Object interaction features

# Social role model

Person-specific features

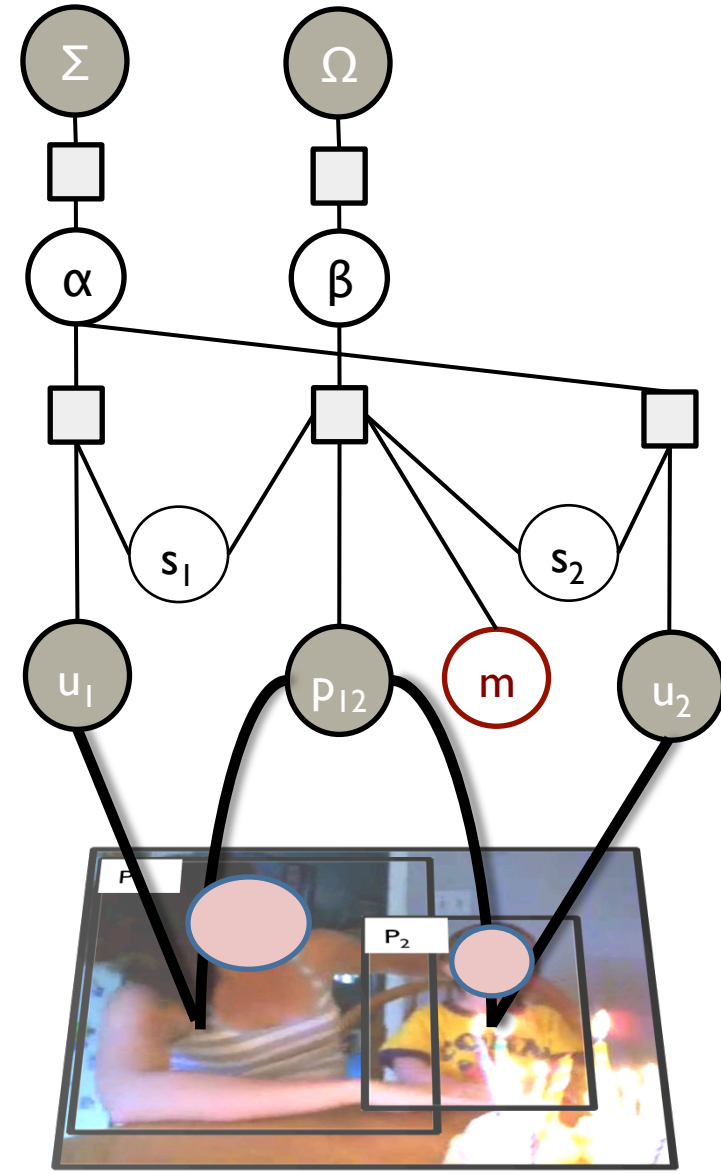## Inter-role features

- Spatio-temporal features
- Proxemic features

# Social role model

Person-specific features

Inter-role features

## Reference role

only interactions with reference
considered for tractable inference

# Social role model

Person-specific features

Inter-role features

Reference role

## Model parameters

$\alpha$    is the person-specific feature weight

$\beta$    is the inter-role feature weight

# Social role model

Person-specific features
Inter-role features
Reference role

## Model parameters

$\alpha$   is the person-specific feature weight

$\beta$   is the inter-role feature weight

$\Sigma$

$\Omega$   Gaussian priors for regularization

# Social role model

**Variational Inference**

*Jointly*  *Learn model parameters. Assign social roles.*

# Dataset

- Youtube Social Role Dataset
  - Available at
    http://vision.stanford.edu/vigneshr_release_data/
    youtube_CVPR13_social.tar.gz
  - Only the event type is provided.
  - Social roles are discovered in unsupervised fashion.

# Results: role clusters

# Results: role clusters

# Results: role clusters
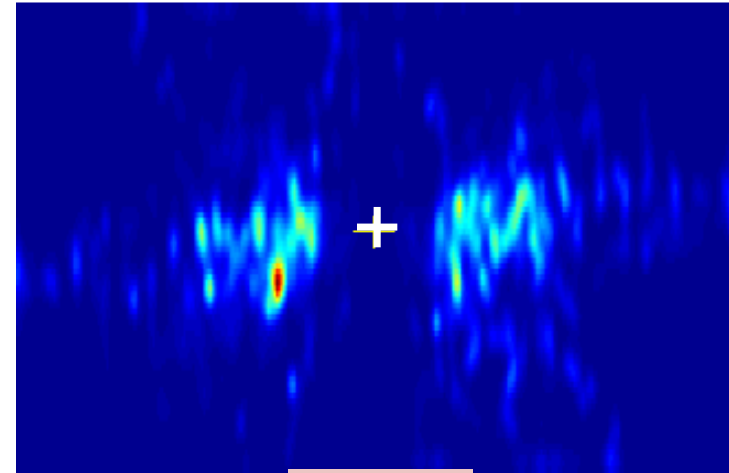


Bride

Brides maid

reference
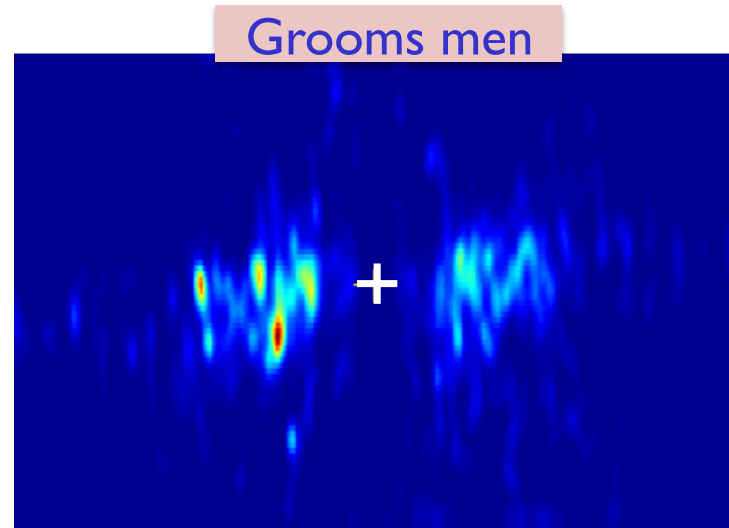
Groom

Priest

Grooms men

# Spatial Relations

# Spatial Relations



Bride

reference

Priest

Brides maid

Grooms men

Groom

# Results – Role clusters



bride
groom
priest
grooms men
brides maid

b'day person
parent
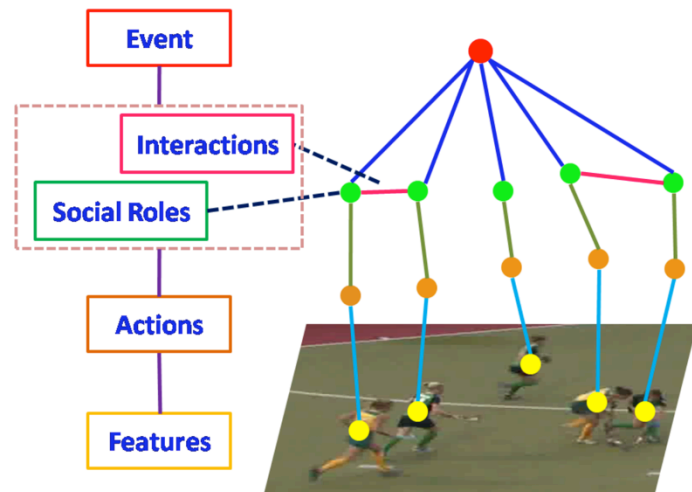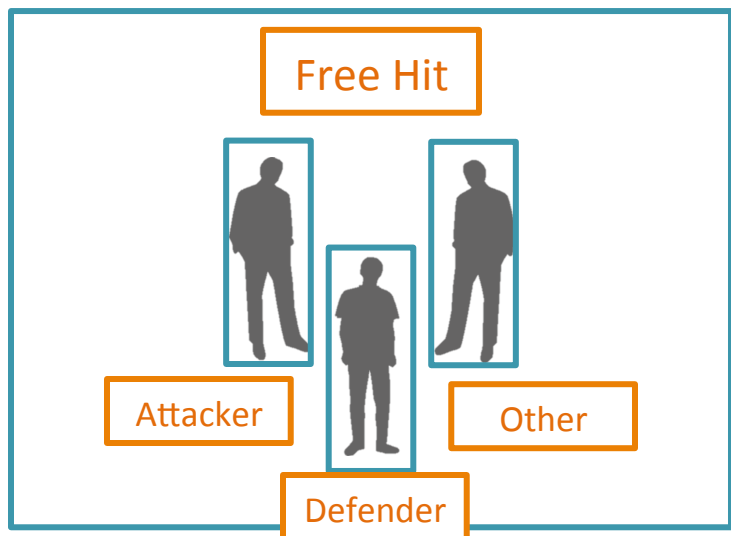friends
guest

presenter
recipient
host
distributor

instructor
presenter

# Social Role Discovery (Lan CVPR12)

Input: a video with box tracks

Output: social role and activity labels

# Semantic Description of Videos



| actions | social roles | event |
|---------|--------------|-------|
| walk | attacker | corner-hit |
| run | first defender | free-hit |
| jog | man-marking | attack play |
| bend | defend-space | … |
| shoot | Teammate | |
| dribble | … | |
| … | | |

**Social roles**
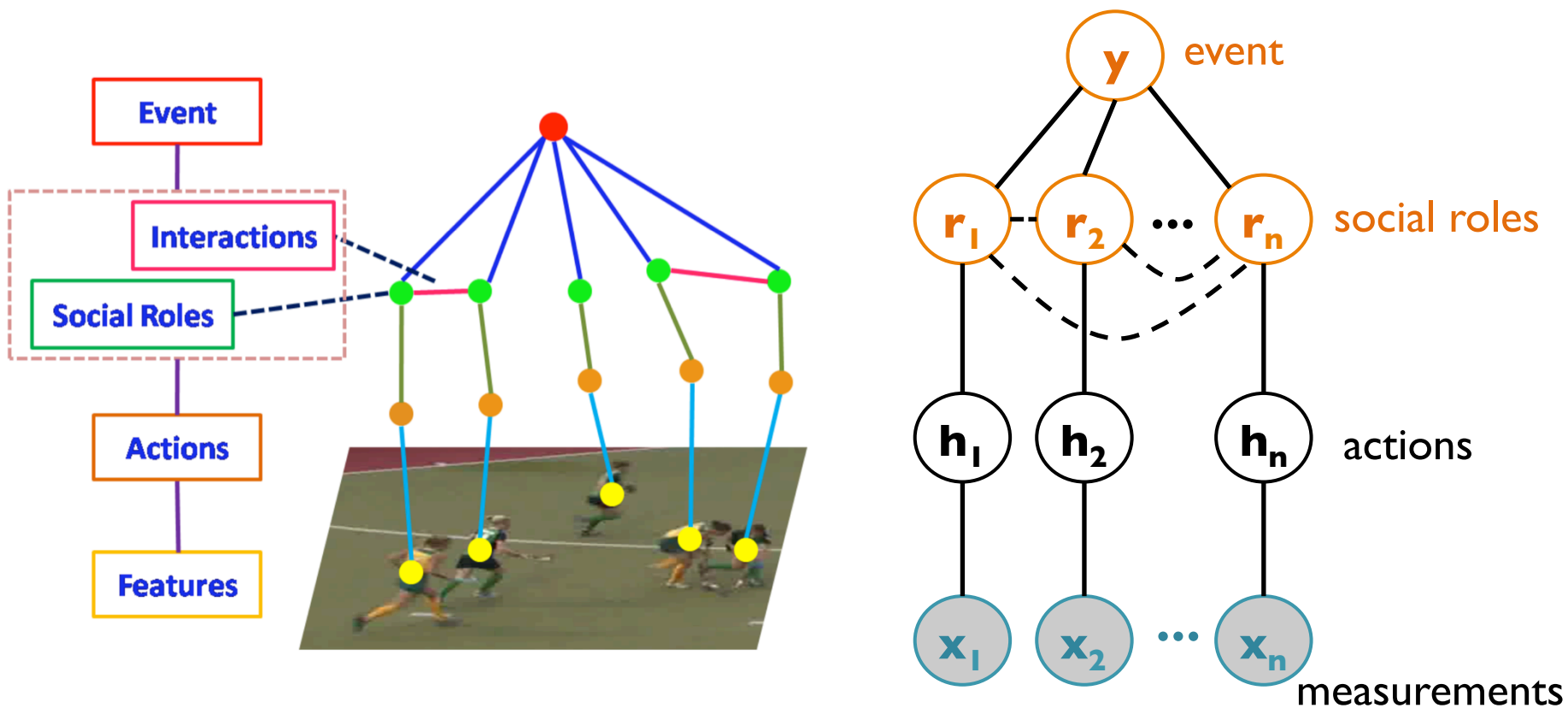
Mid-Level semantics that describe individual/group behaviors in the context of social interactions.
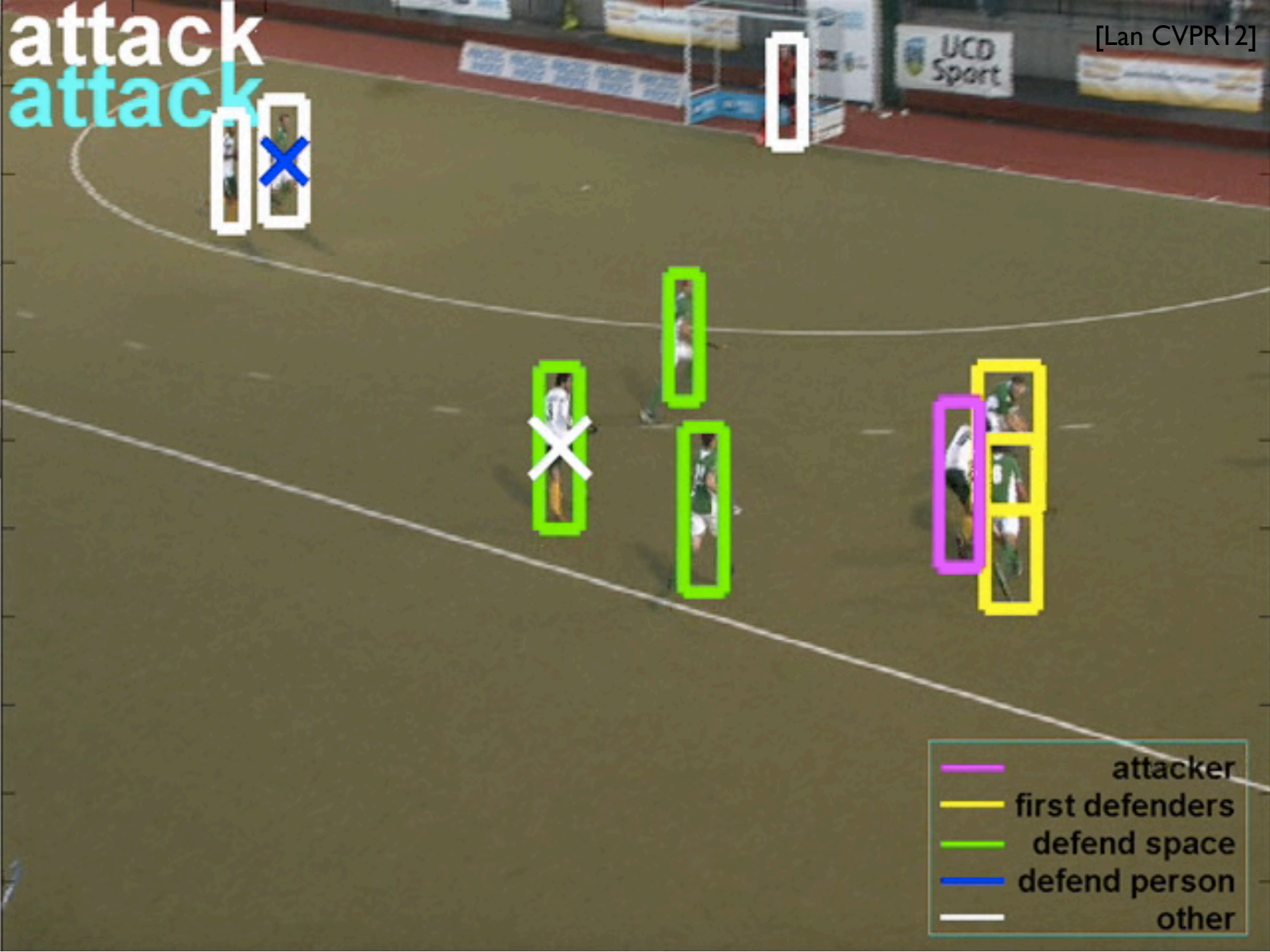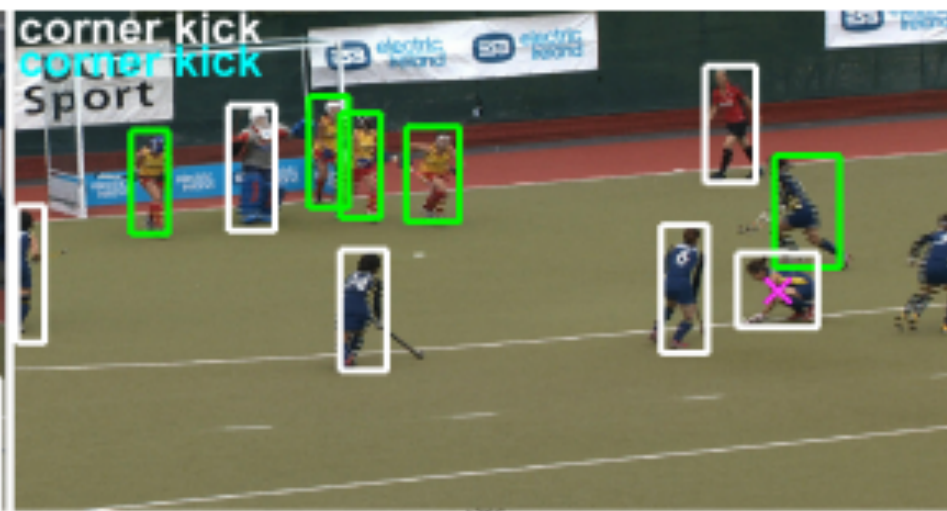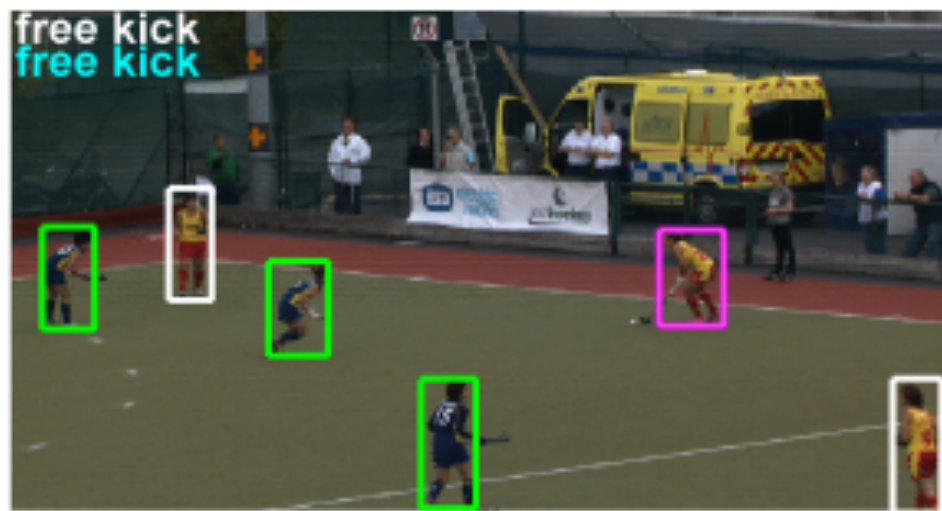
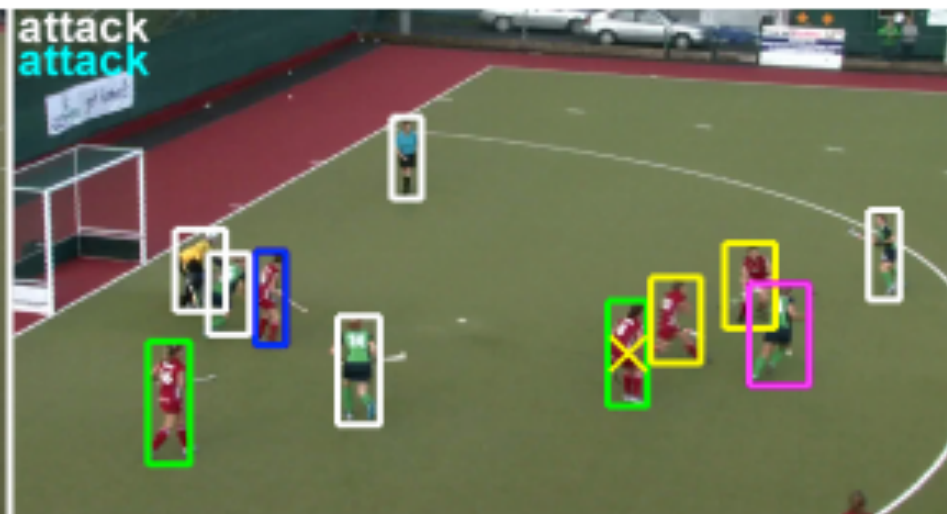man-marking

first defenders

# Hierarchical Model

attack
attack

[Lan CVPR12]

UCD
Sport

attacker
first defenders
defend space
defend person
other
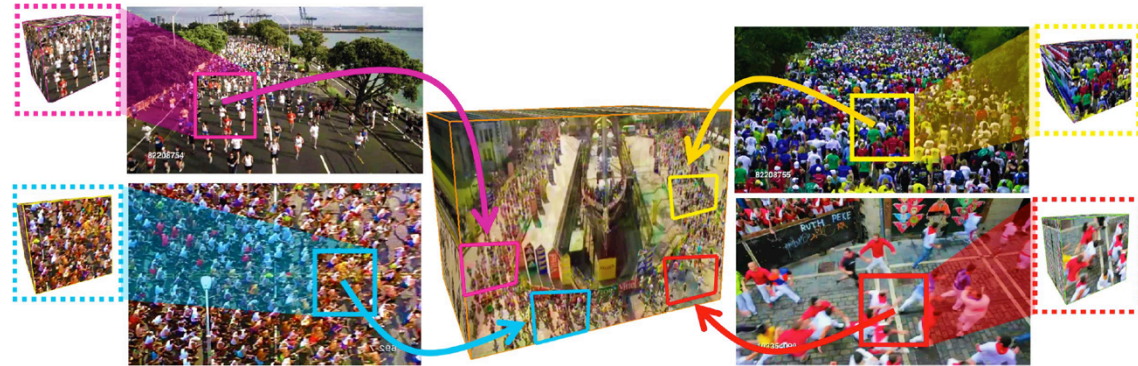
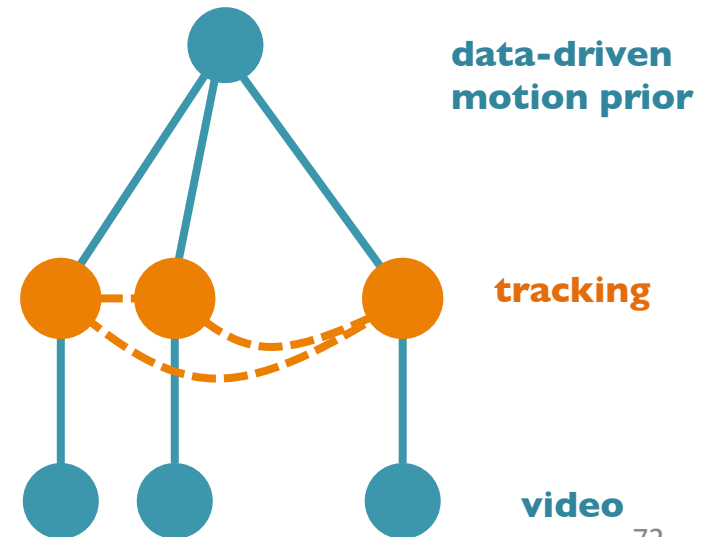attacker, first defender, defend space, defend person, other
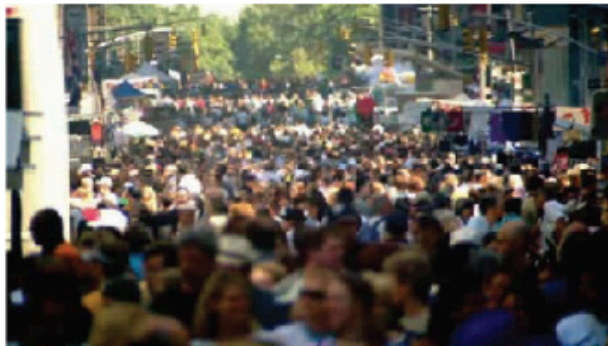
# Data-driven Crowd Analysis (Rodriguez ICCV11)
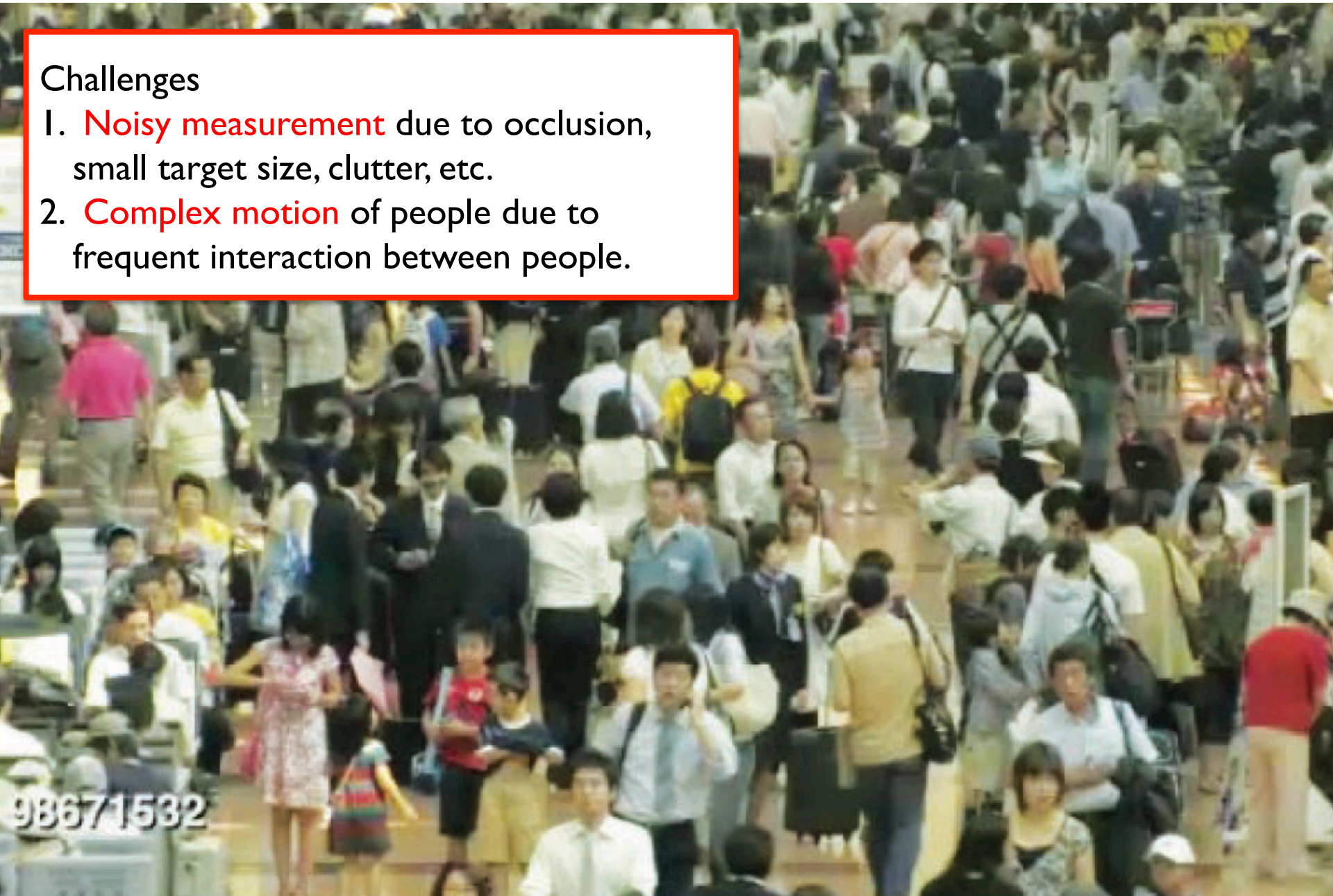
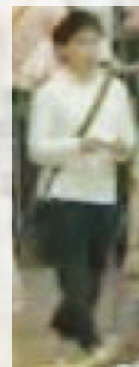Input: a crowd video
Output: individual tracking



Crowd video



data-driven
motion prior

tracking

video

[Rodriguez ICCV11]

Challenges
1. Noisy measurement due to occlusion, small target size, clutter, etc.
2. Complex motion of people due to frequent interaction between people.

Challenges
1. Noisy measurement due to occlusion, small target size, clutter, etc.
2. Complex motion of people due to frequent interaction between people.

Challenges
1. Noisy measurement due to occlusion, small target size, clutter, etc.
2. Complex motion of people due to frequent interaction between people.
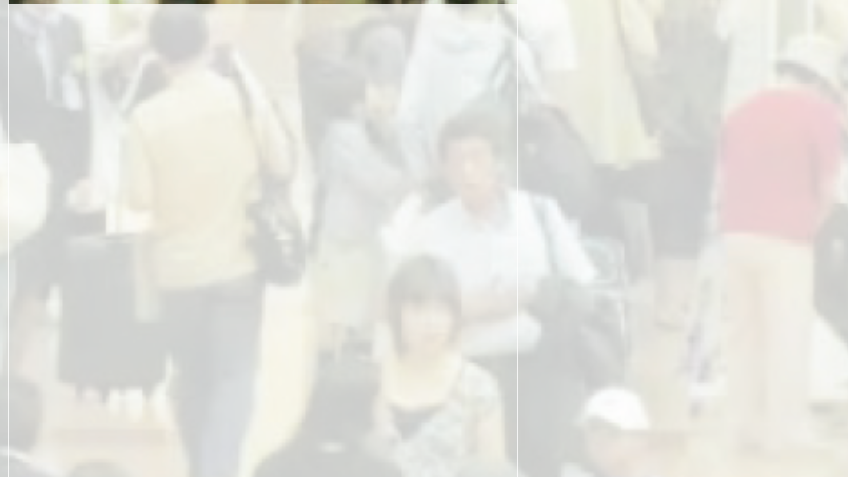
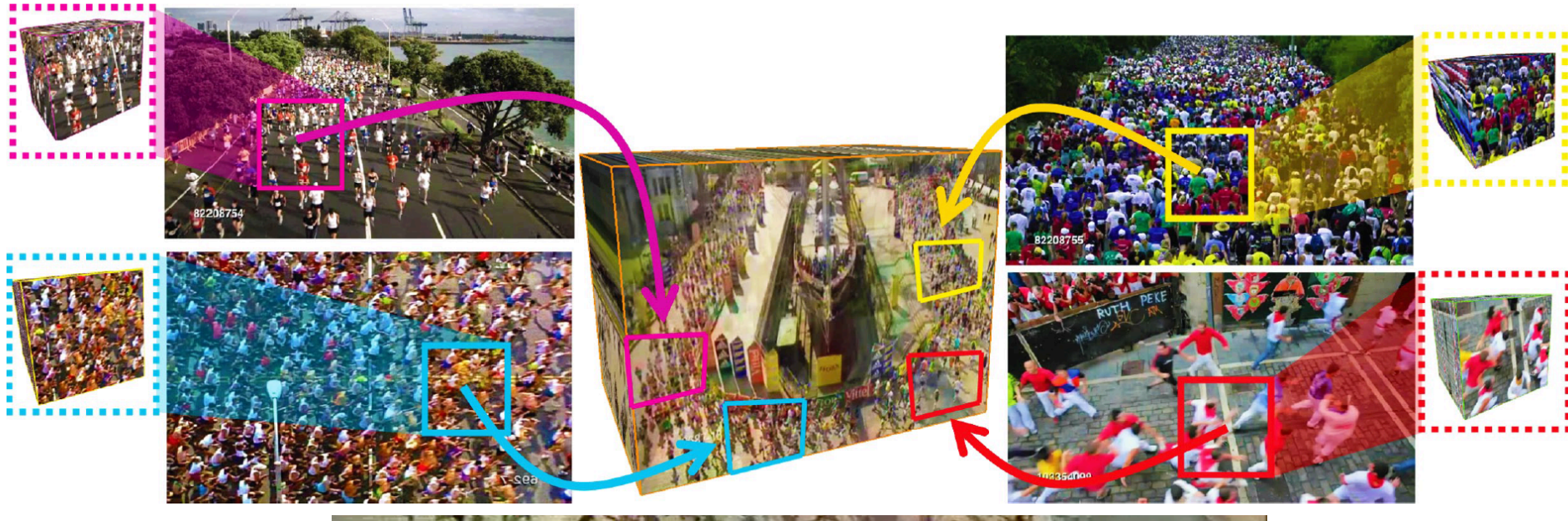Solution
1. See larger are to encode collective signal.

Challenges
1. Noisy measurement due to occlusion, small target size, clutter, etc.
2. Complex motion of people due to frequent interaction between people.

Solution
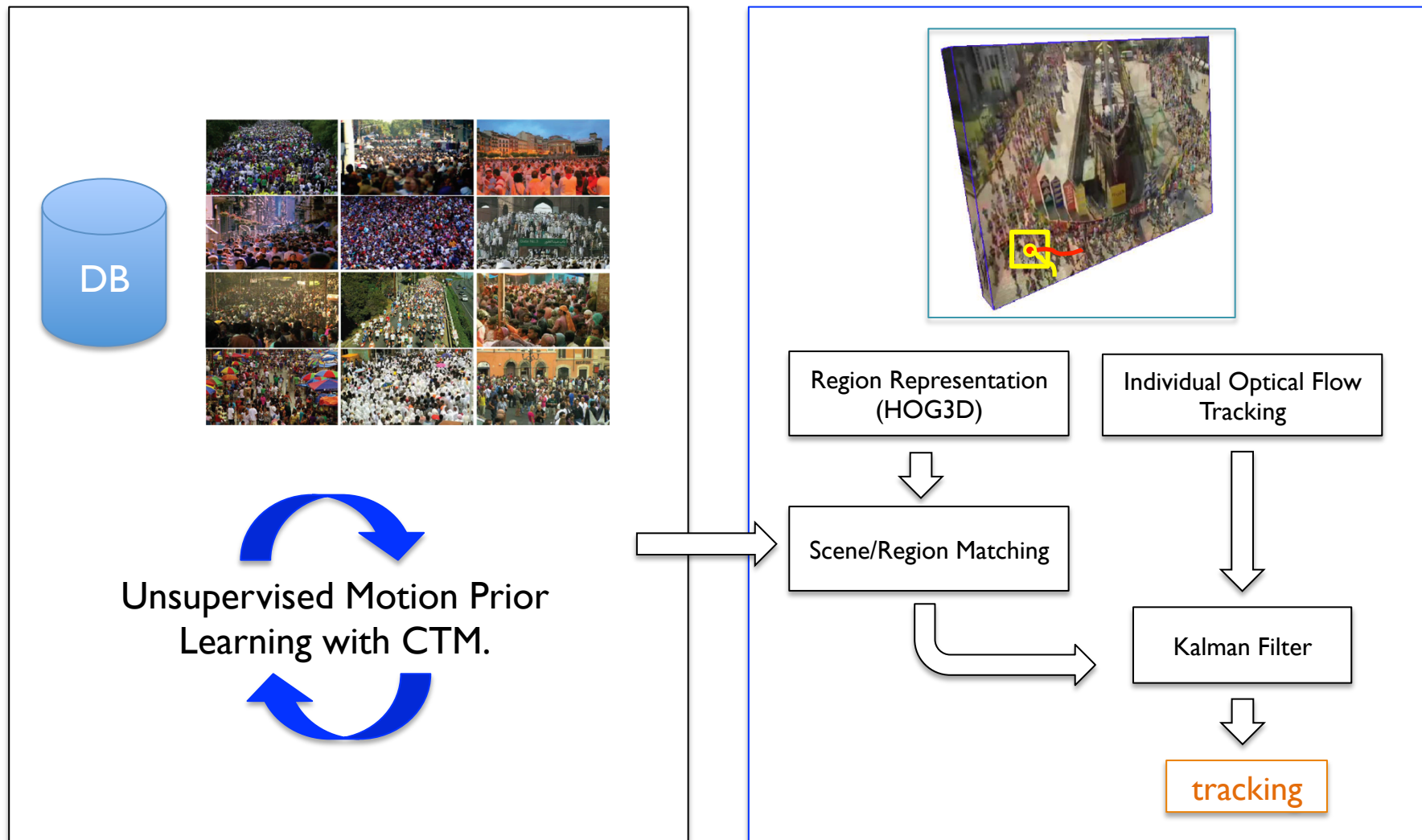1. See larger are to encode collective signal.
2. Transfer the motion prior of data with similar collective signal.

DB

# Mixture of Other Videos

# Framework



DB

Unsupervised Motion Prior
Learning with CTM.

Region Representation
(HOG3D)

Individual Optical Flow
Tracking

Scene/Region Matching

Kalman Filter

tracking

# Tracking in Crowd Videos



Ground Truth, Tracking Results

# Tracking in Crowd Videos



Ground Truth, Tracking Results
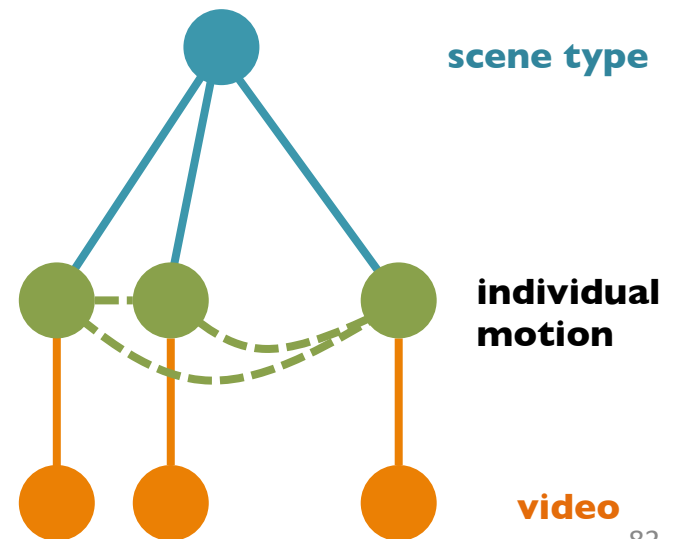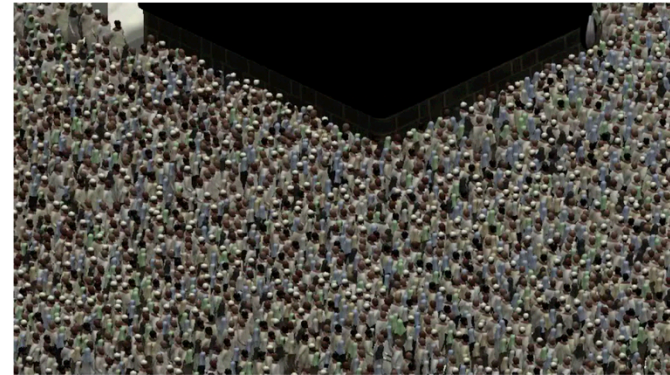
# Tracking in Rare Event

Ground-truth, Base-line, Data-driven racking Results
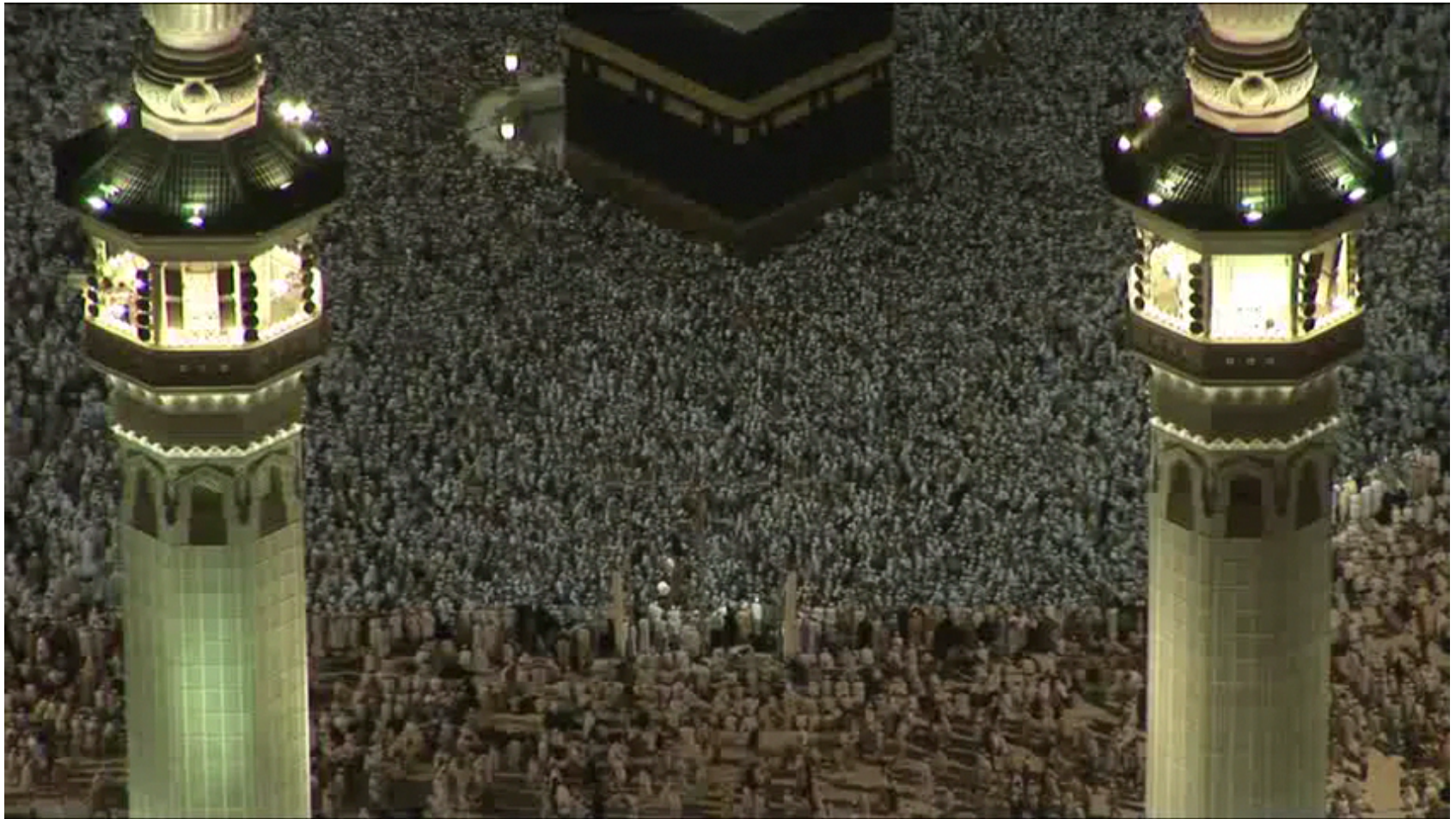
# Application: Crowd Simulation [Curtis LC11]
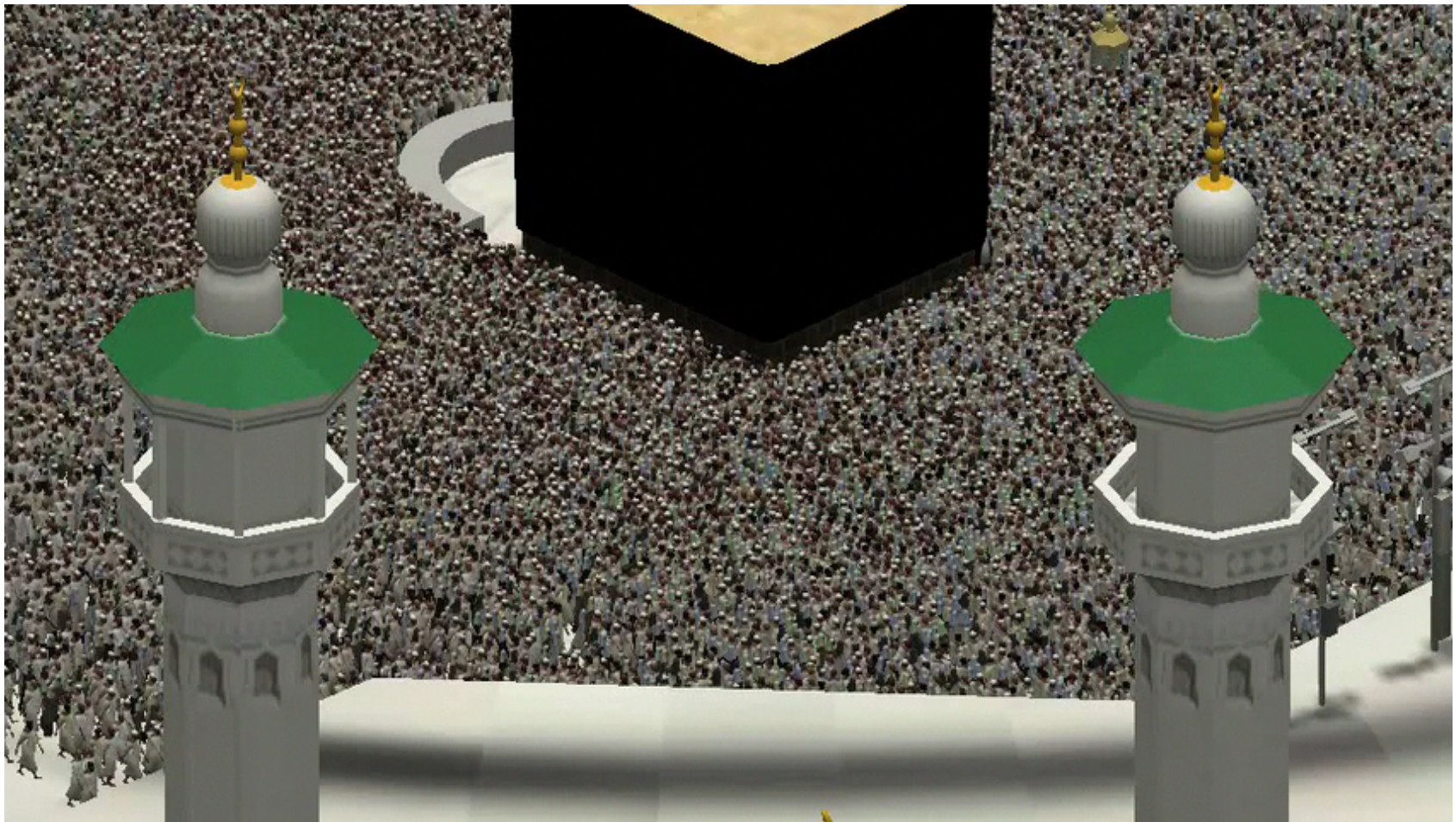
Input: scene type
Output: video of a crowd
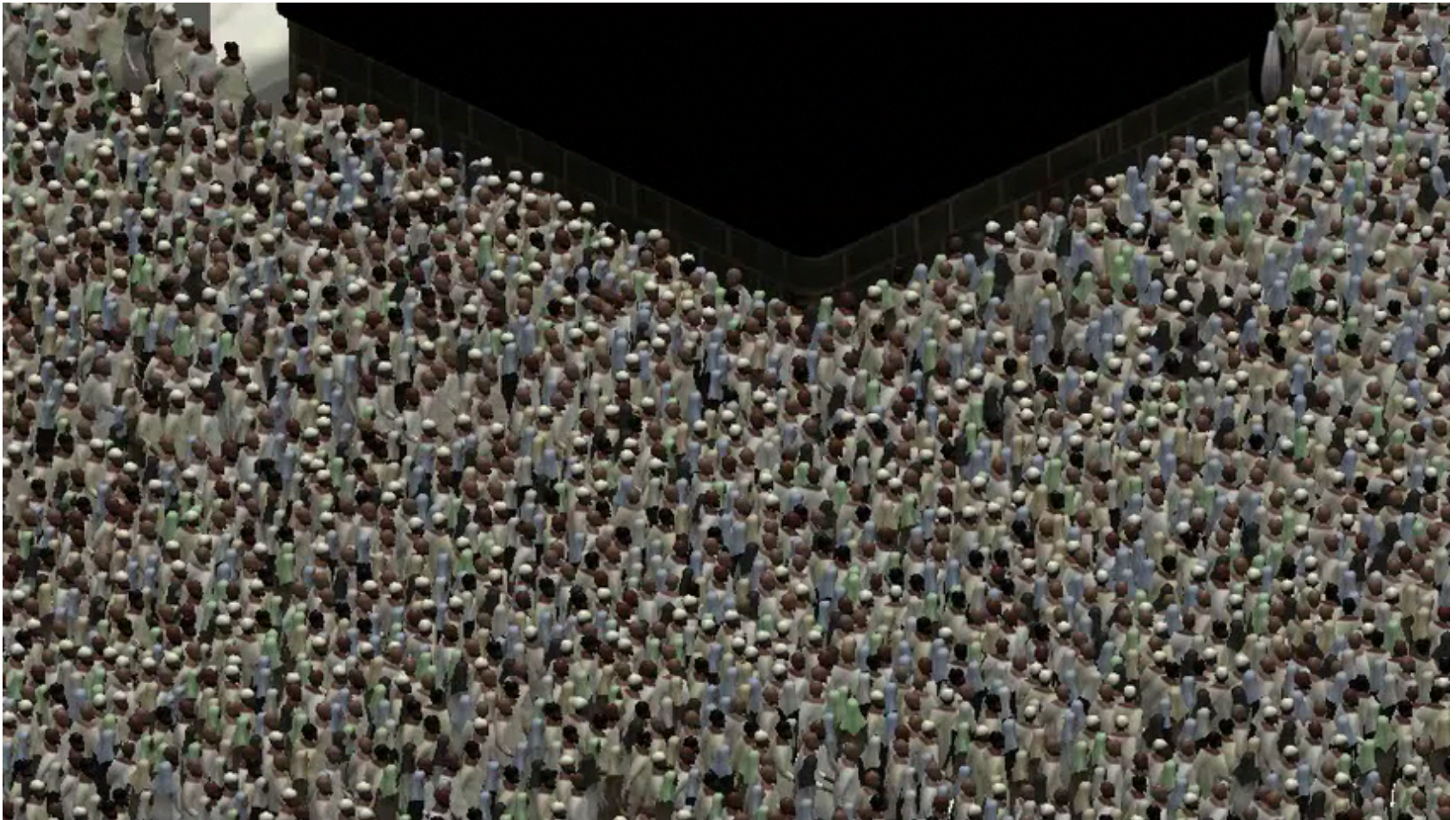Model: Social force + FSM



scene type

individual
motion
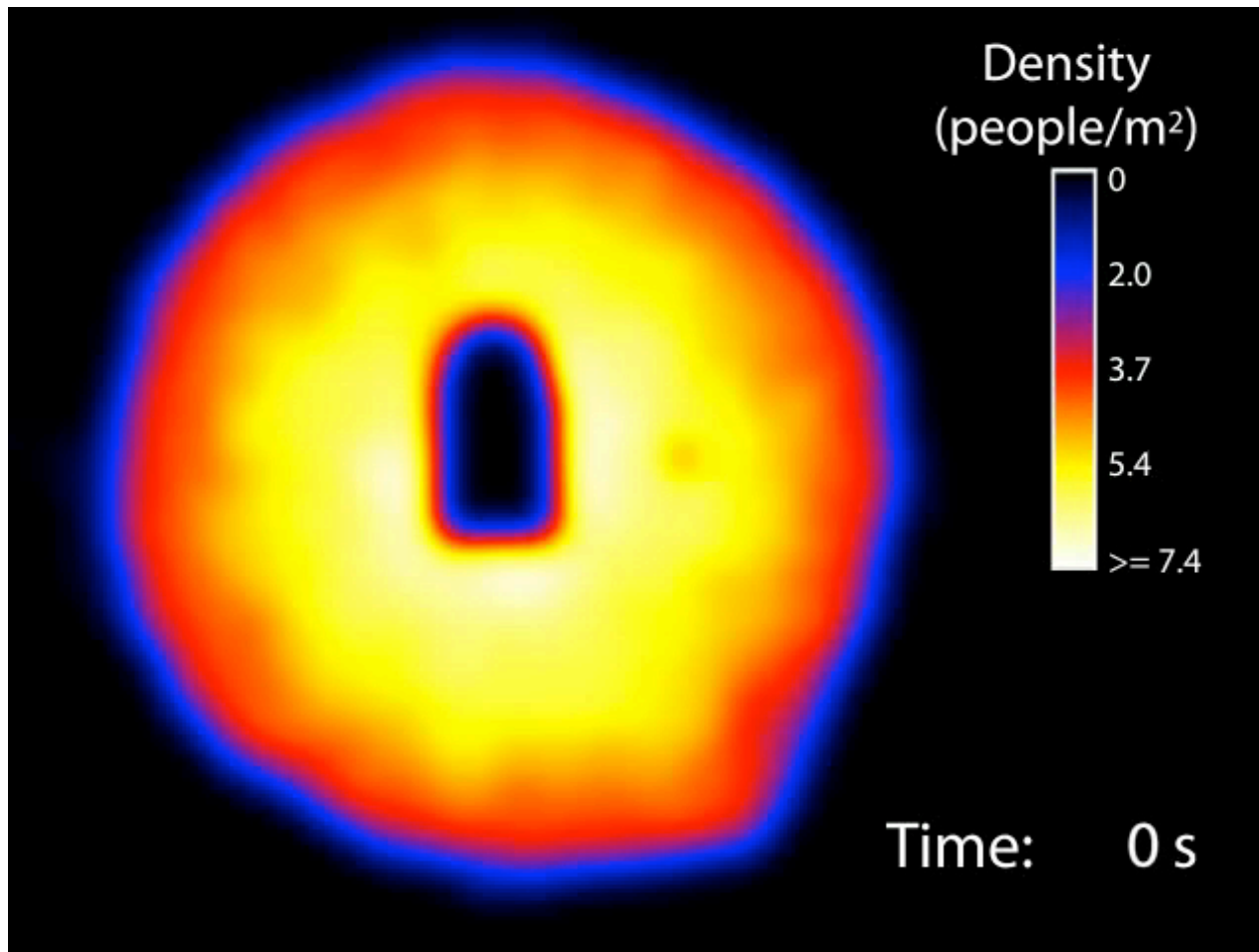
video

# Real Video from Kaaba during Hajj

# Simulation Results

# Simulation Result - Zoomed

# Density of People

# Speed of People