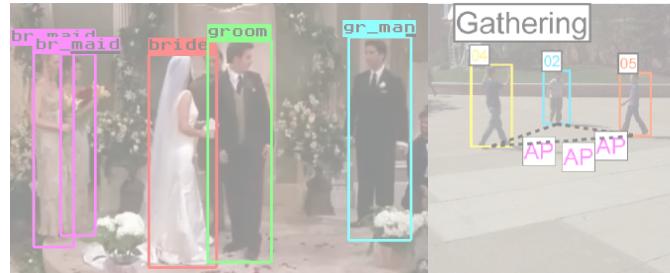


Scene dynamism

Dynamic scene
 Static scene



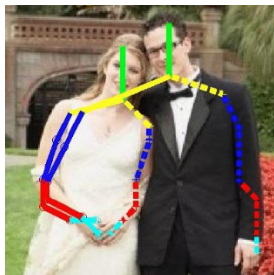
Rehg, CVPR13
 Prabhakar, ECCV12
 Prabhakar, CVPR12
 Patron-Perez, BMVC10



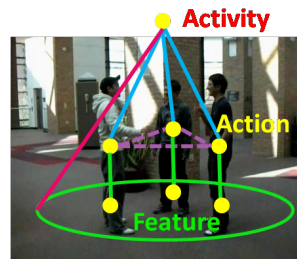
Lan, CVPR12
 Ramanathan, CVPR13
 Antic, ECCV14
 Ding, ECCV10
 Choi, ECCV12, CVPR14
 Direkoglu, ECCV12



Rodriguez, ICCV11a, ICCV1b
 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

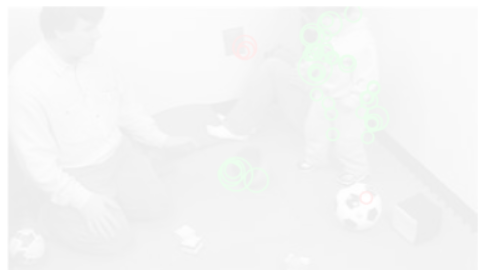
Dyadic interaction

Crowd interaction

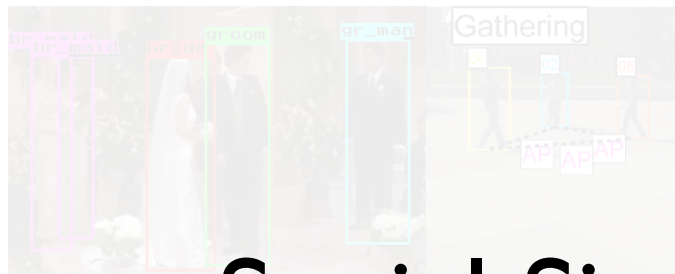
Number of group members

Scene dynamism

Dynamic scene
Static scene



Rehg, CVPR13
Prabhakar, ECCV12
Prabhakar, CVPR12
Patino, Perez, BMVC10



Antic, ECCV14
Direkoglu, ECCV12



Rodriguez, ICCV11a, ICCV11b
Mehran, CVPR09
Alahi, CVPR14

Body Pose as Social Signals



Yang, CVPR12
Hoai, CVPR14



Fathi, CVPR12
Choi, ECCV14
Park, NIPS12, ICCV13

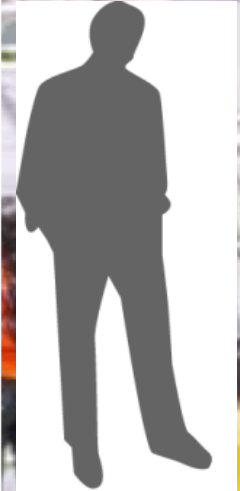


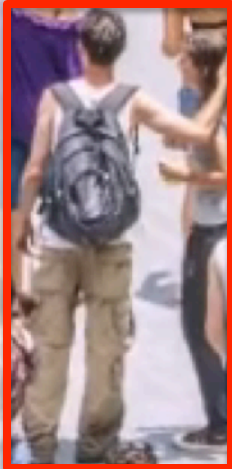
Cristani, BMVC11
Park, CVPR15
Arev, SIGGRAPH14
Wang, ECCV10
Gallagher, CVPR09

Dyadic interaction

Crowd interaction

Number of group members





Pose as a skeleton



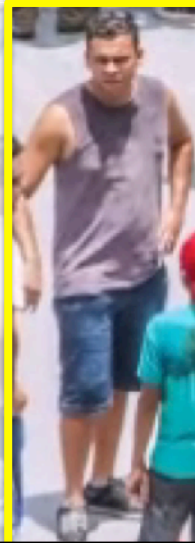
+ Almost complete representation.
+ Provide fine details of the person.

- Hard to estimate in complex scene.
- Large amount of variation.



Pose as a finite set of semantic classes

Standing
Sitting
Running
Facing-Front
Facing-Right
...



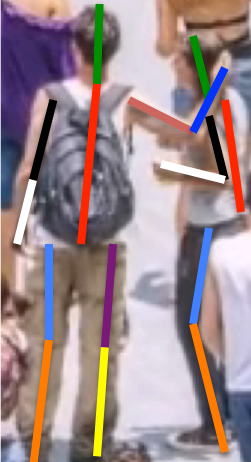
+ Compact representation.
+ Higher accuracy.

- Require definition of classes.
- Less details.

Standing and Facing-forward



Standing and Facing-forward



The left touches the right's head.

Proxemics

Group Activity Recognition



Conversation

Standing and Facing-front

Standing and Facing-right

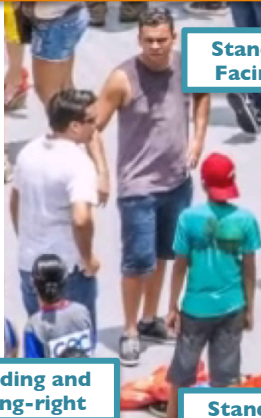
Standing and Facing-back

Exploiting Structures in Group Behavior



The left touches the right's head.

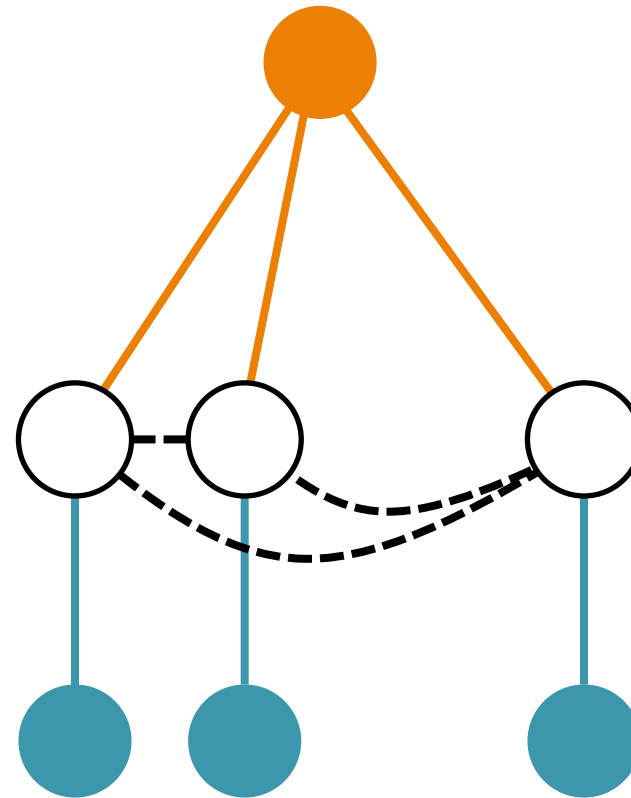
Conversation



Standing and Facing-front

Standing and Facing-right

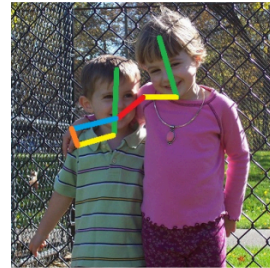
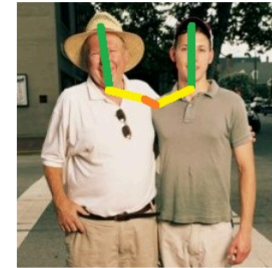
Standing and Facing-back



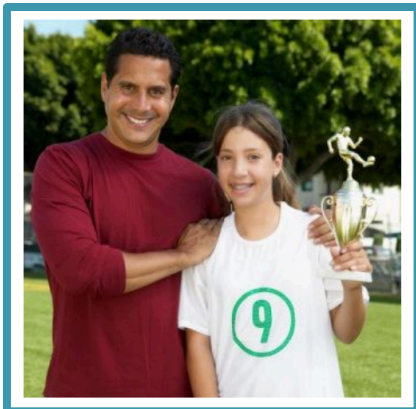
Proxemics in Personal Photos (Yang CVPR 12)

Input: image

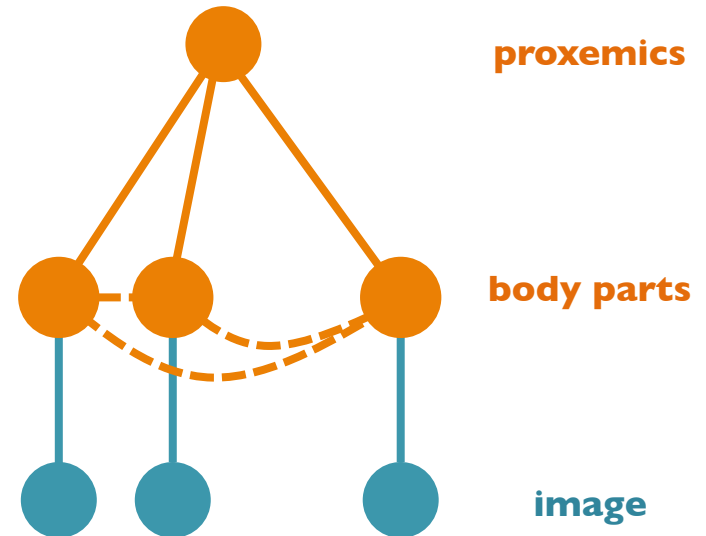
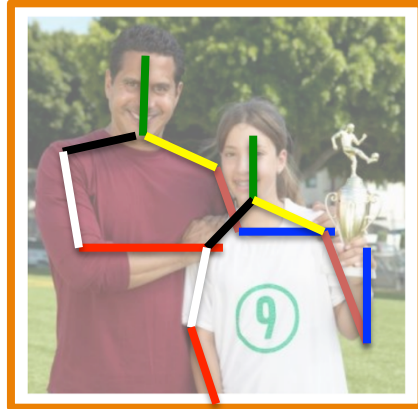
Output: proxemics label
and skeletons



Image



Hand-Shoulder



Proxemics

Proxemics: the study of spatial arrangement of people as they interact

- anthropologist [Edward T. Hall](#) in 1963

Proxemics in Personal Photos

Hand touch Hand



Hand touch Shoulder



Shoulder touch Shoulder

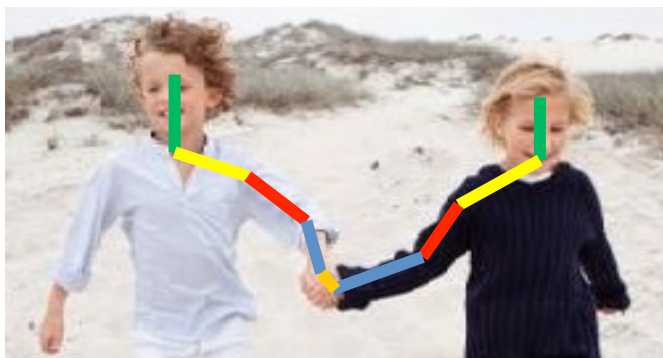


Arm touch Torso

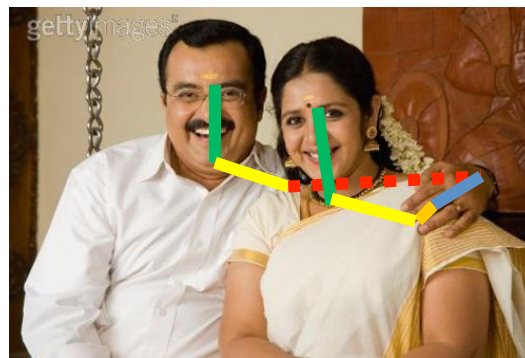


Proxemics in Personal Photos

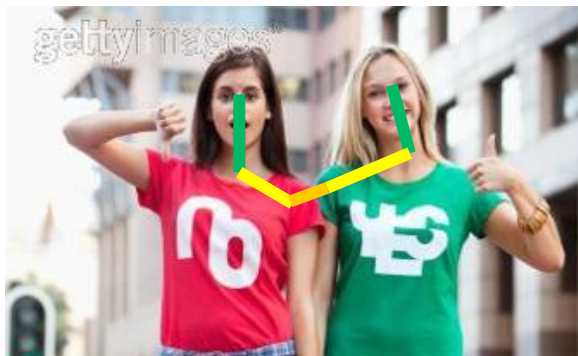
Hand touch Hand



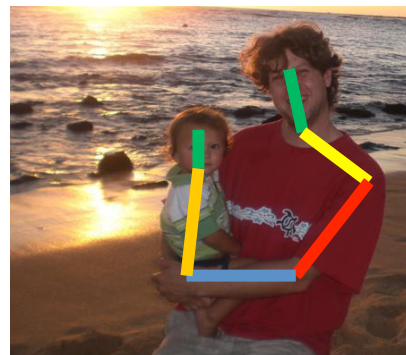
Hand touch Shoulder



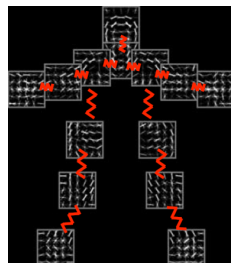
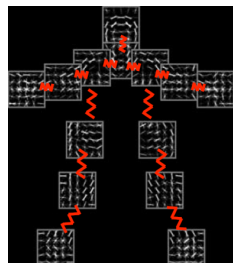
Shoulder touch Shoulder



Arm touch Torso



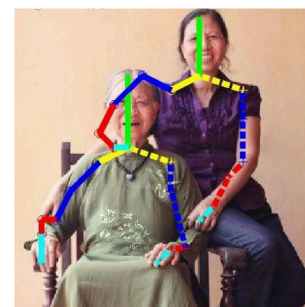
Naïve Solution



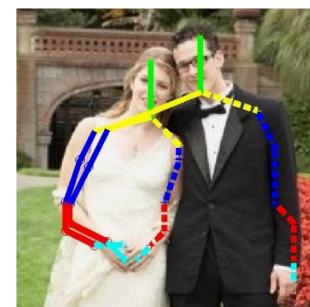
Hand touch Hand



(a) Extreme Pose



(b) Occlusion



(c) Part Ambiguity

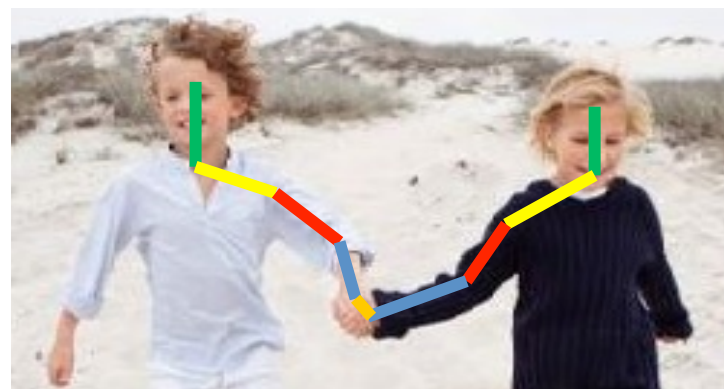
Joint Pose Estimation



Connected



Touch code

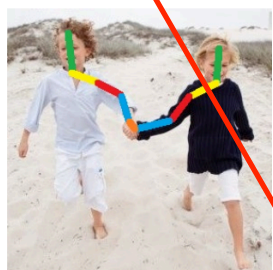
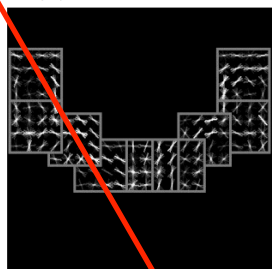


Hand touch Hand

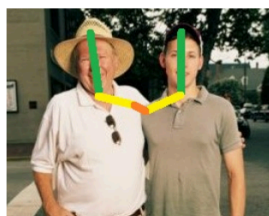
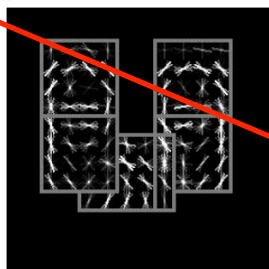
Tree structure!
Efficient Inference.

Proxemic Dependent Pose Models

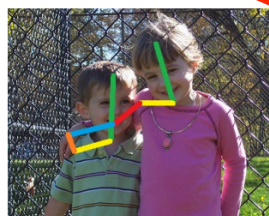
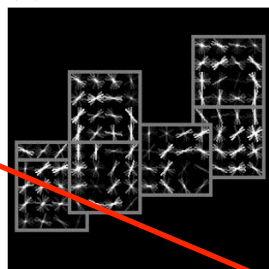
(a) Hand-Hand



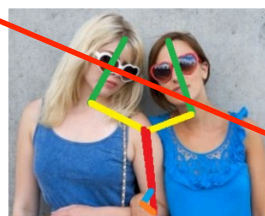
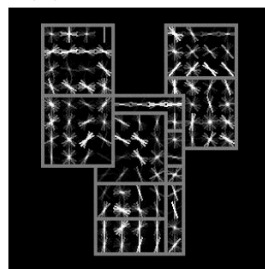
(b) Shoulder-Shoulder



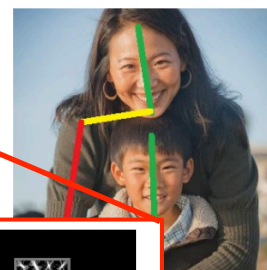
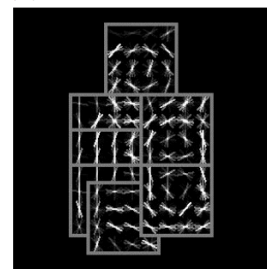
(c) Hand-Shoulder



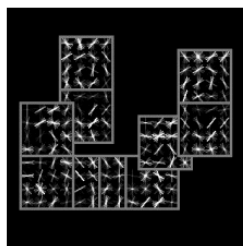
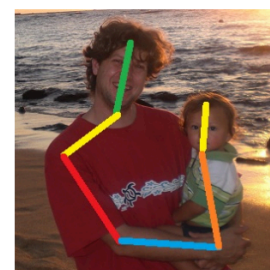
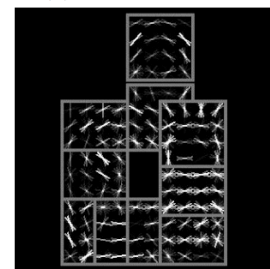
(d) Hand-Elbow



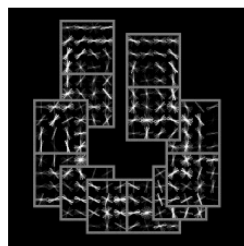
(e) Elbow-Shoulder



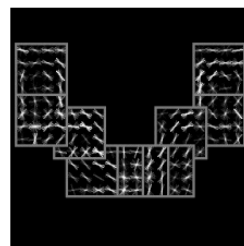
(f) Hand-Torso



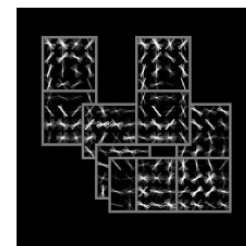
(a) Left left



(b) Left right



(c) Right left



(d) Right right

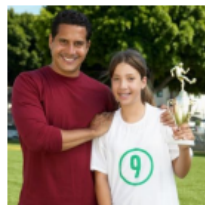
Experimental Setup

- Proxemic Dataset

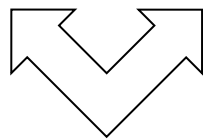
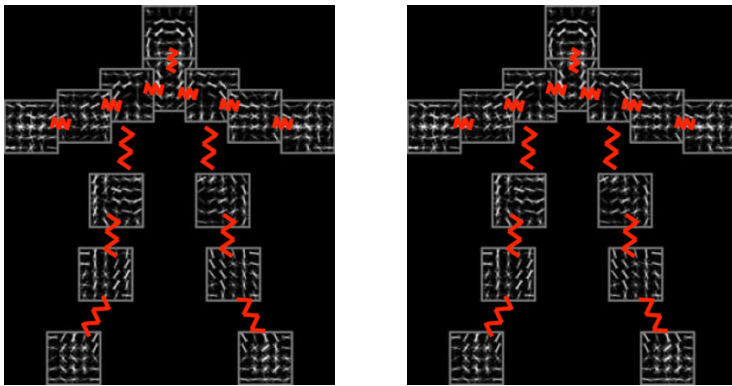
- Available at:

- <http://www.ics.uci.edu/~yyang8/research/proxemics/index.html>

- 589 images, 1207 people, 1332 pairs.

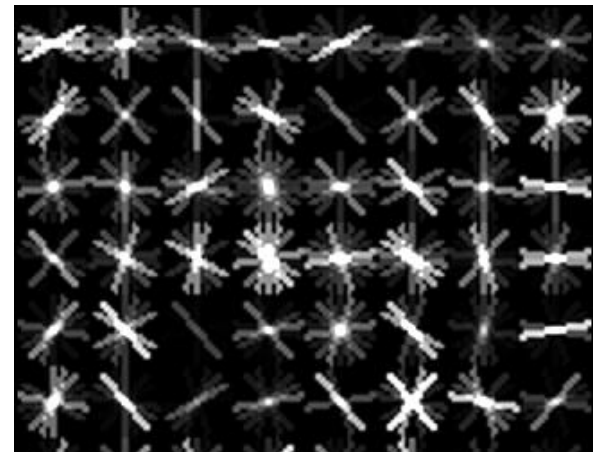


Baselines



Distance?

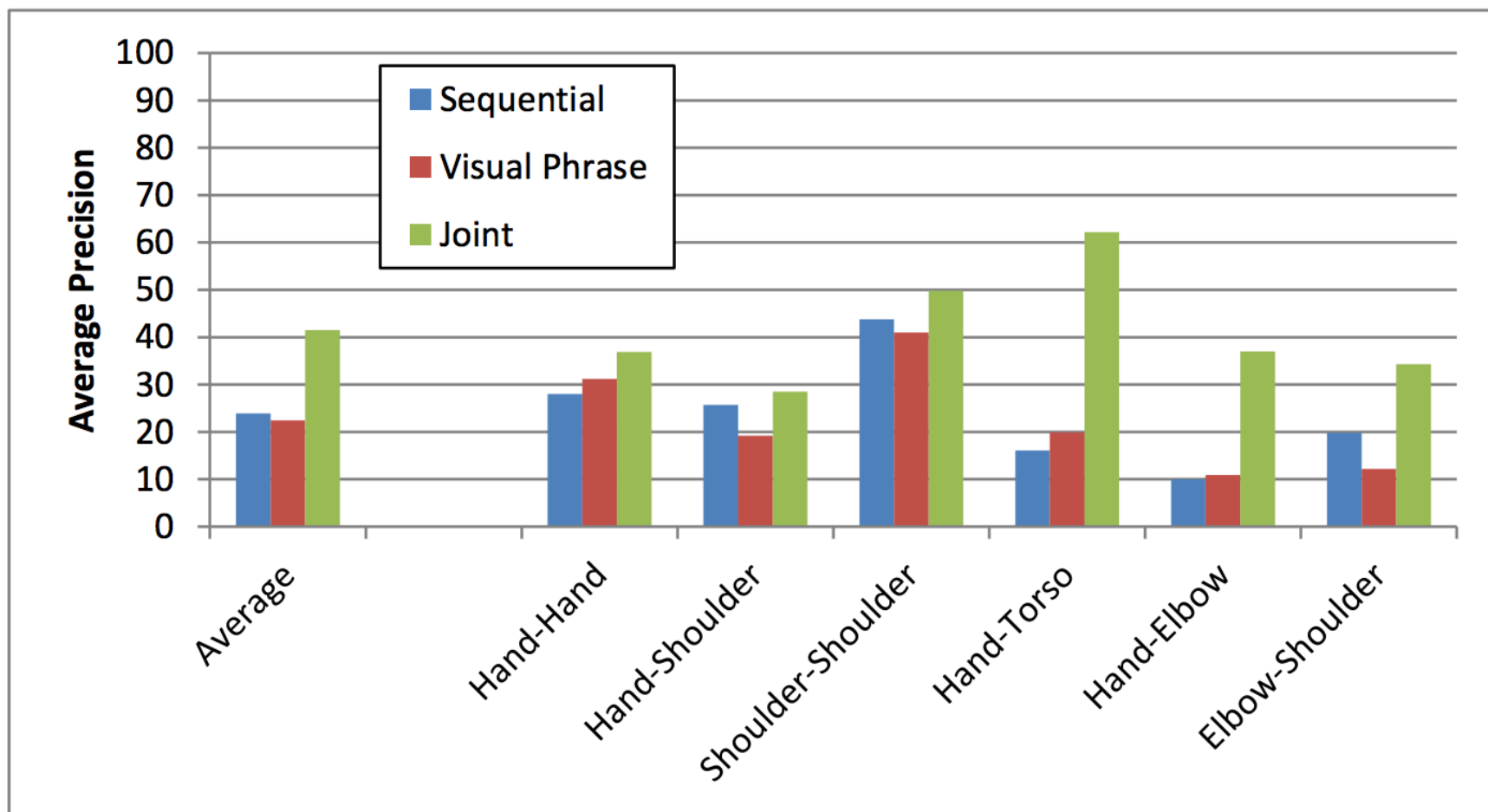
Baseline 1: Sequential



HoG Template

Baseline 2: Phrases

Quantitative Evaluation

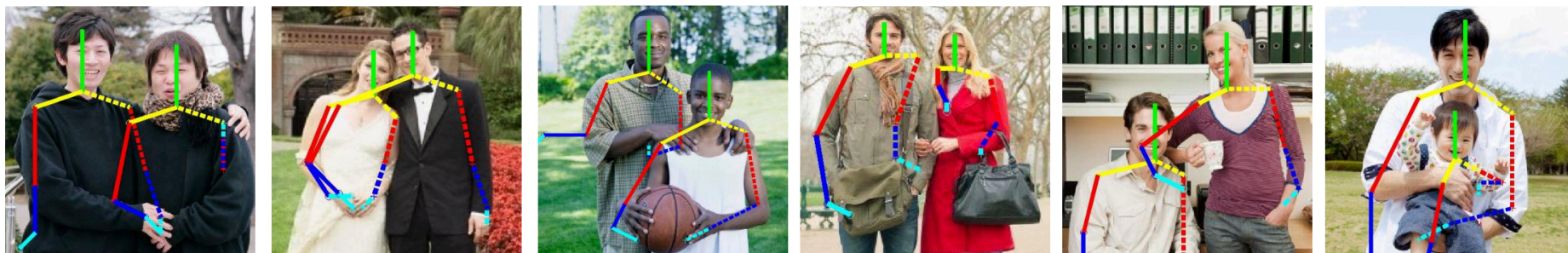


Quantitative Evaluation

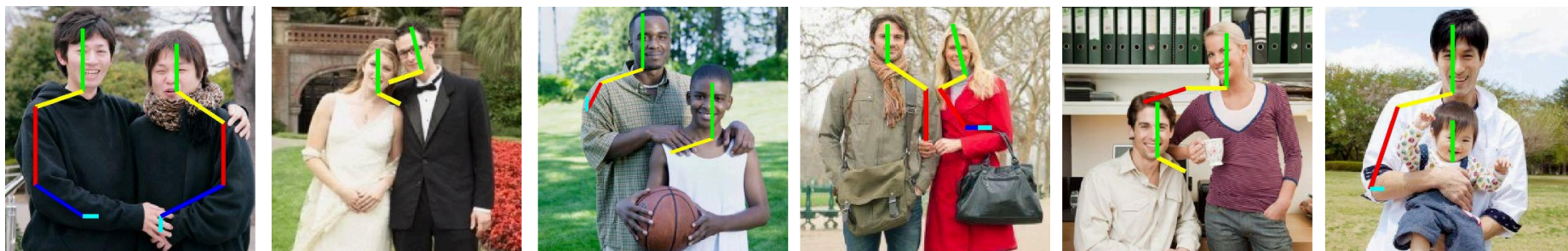
- Pose Accuracy:
 - Sequential Estimation: 47.5%
 - Joint Estimation: 73.6%

Qualitative Examples

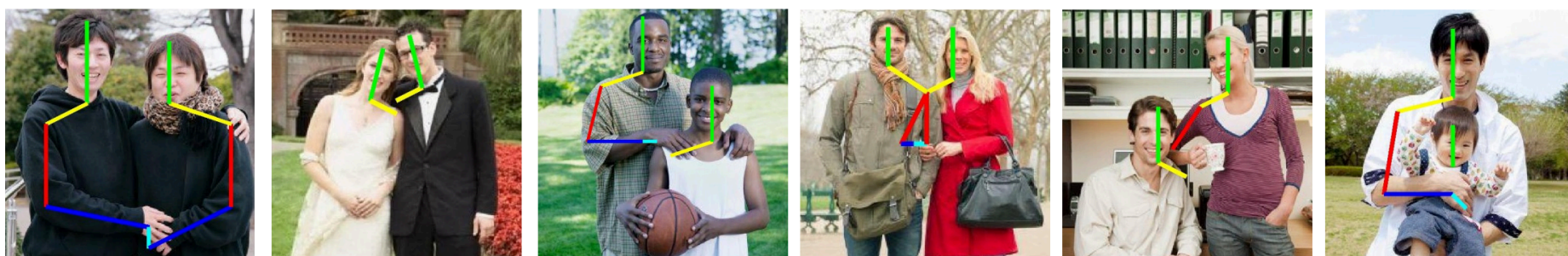
Sequential



No Spring



Joint



(a) Hand hand

(b) Shoulder shoulder

(a) Hand shoulder

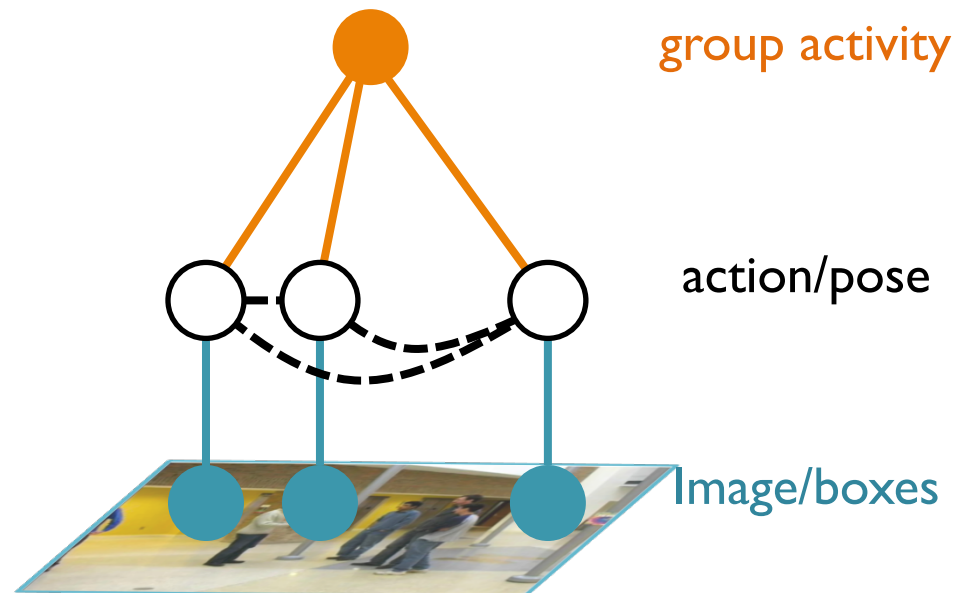
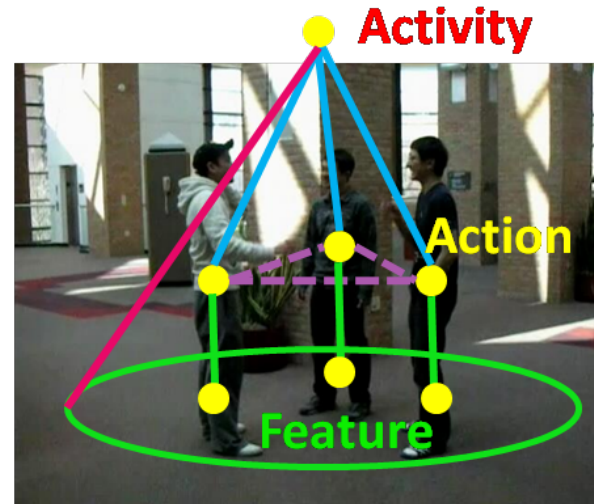
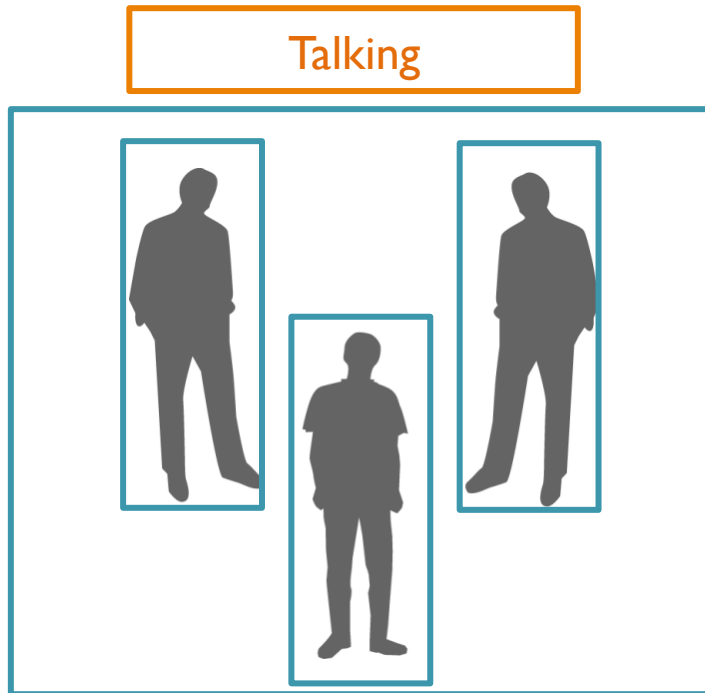
(d) Hand elbow

(c) Elbow shoulder

(e) Hand torso

Group Activity Recognition (Lan NIPS10)

Input: image and bounding boxes
Output: group activity label



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.

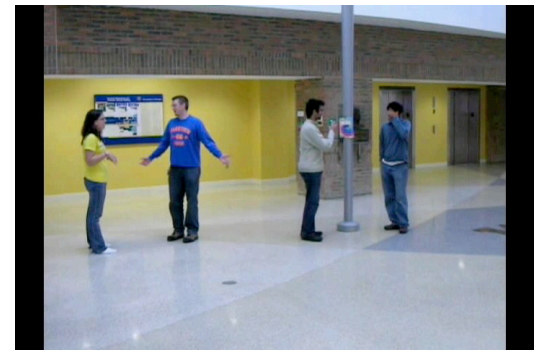
Waiting



Queuing



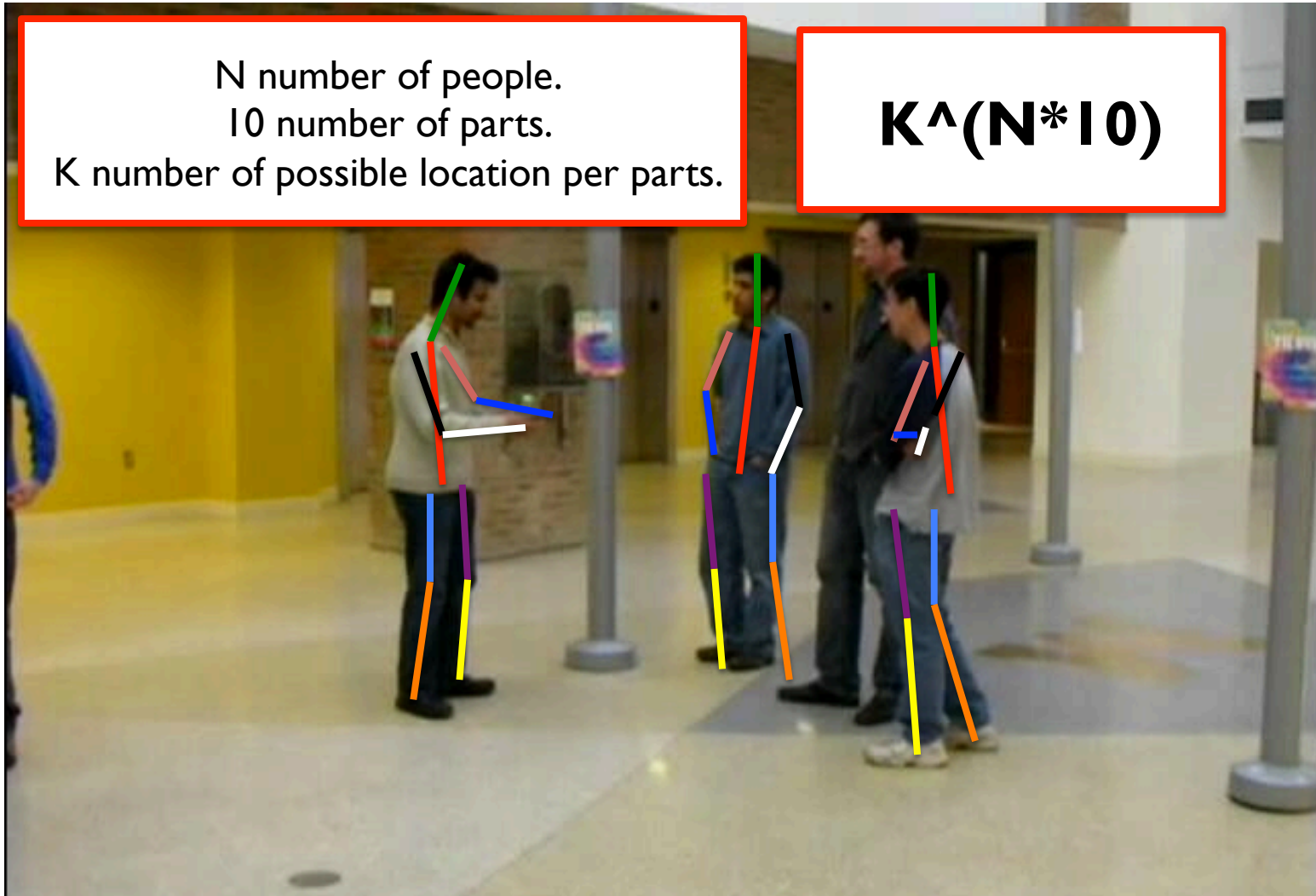
Talking



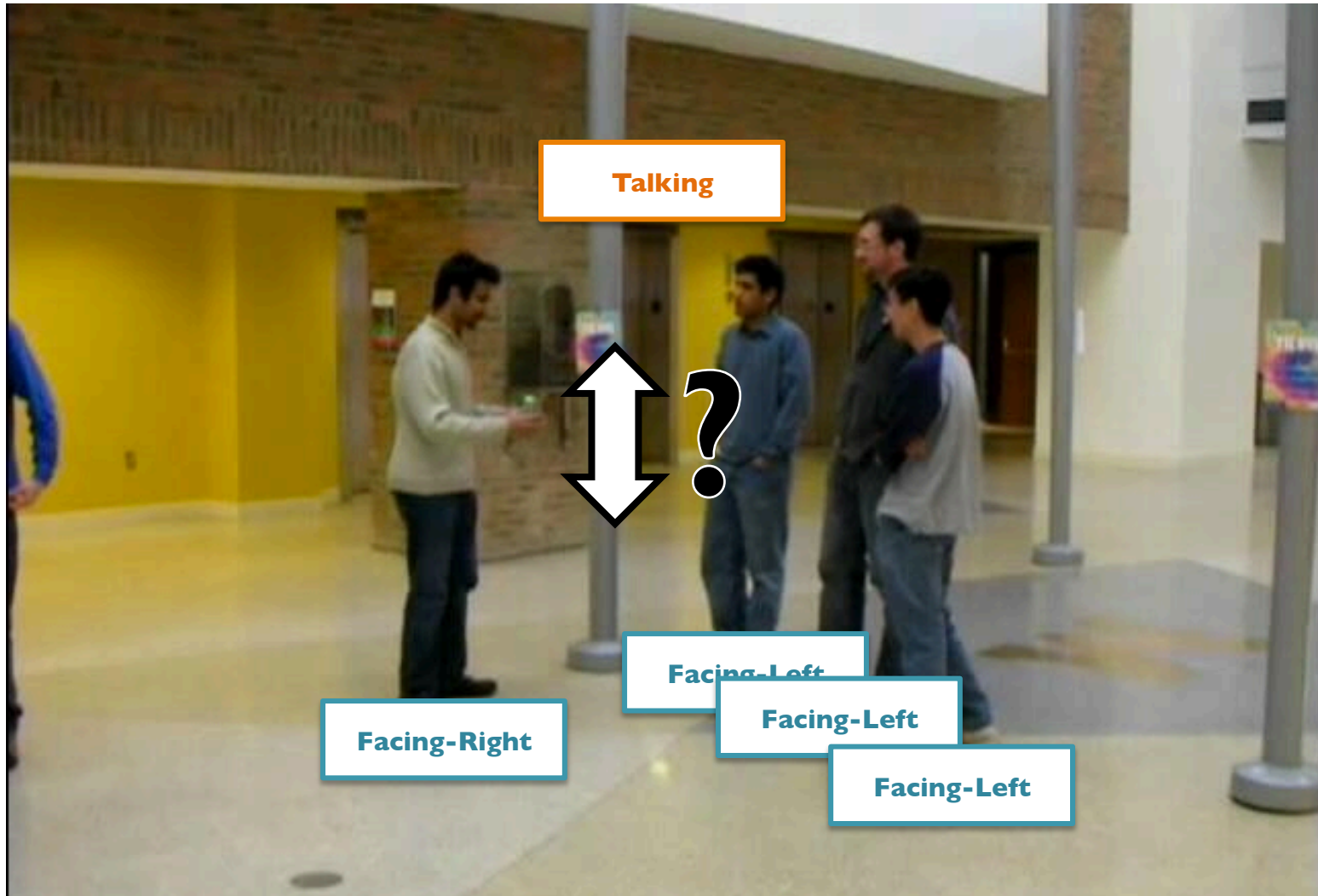
Pose Representation for Groups

N number of people.
10 number of parts.
K number of possible location per parts.

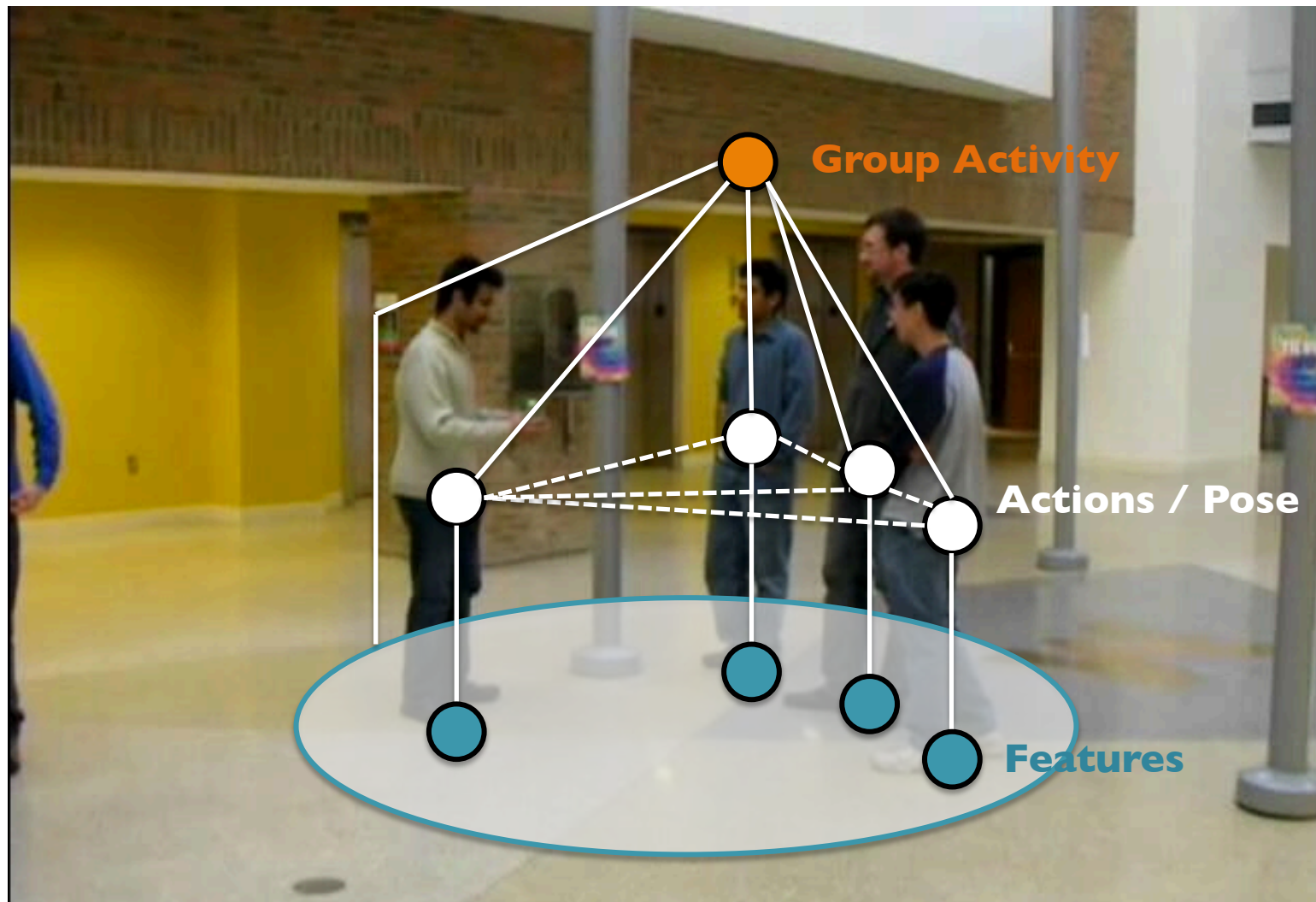
$$K^{(N*10)}$$



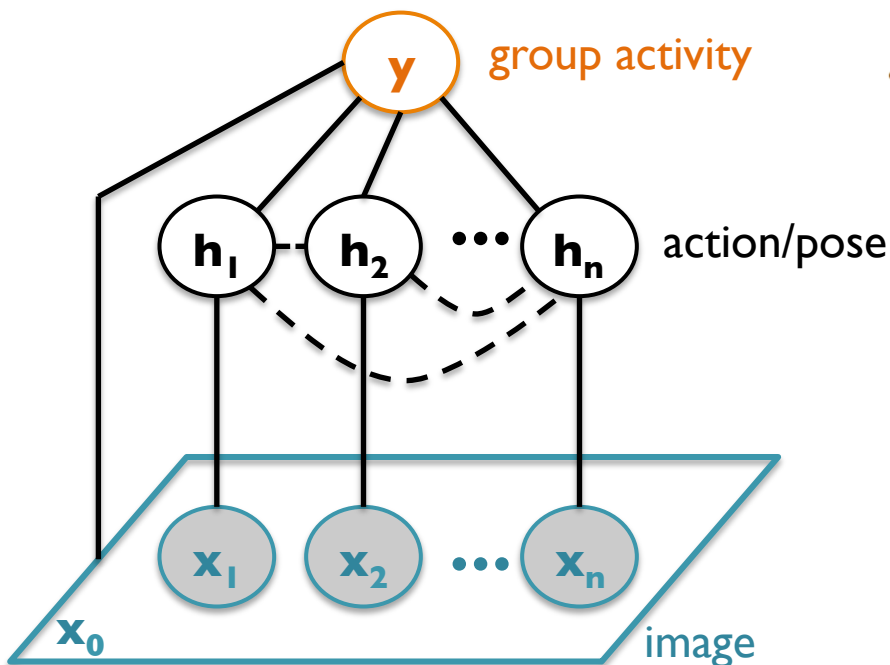
Pose to Group Behavior?



Hierarchical Model For Group Activity Recognition

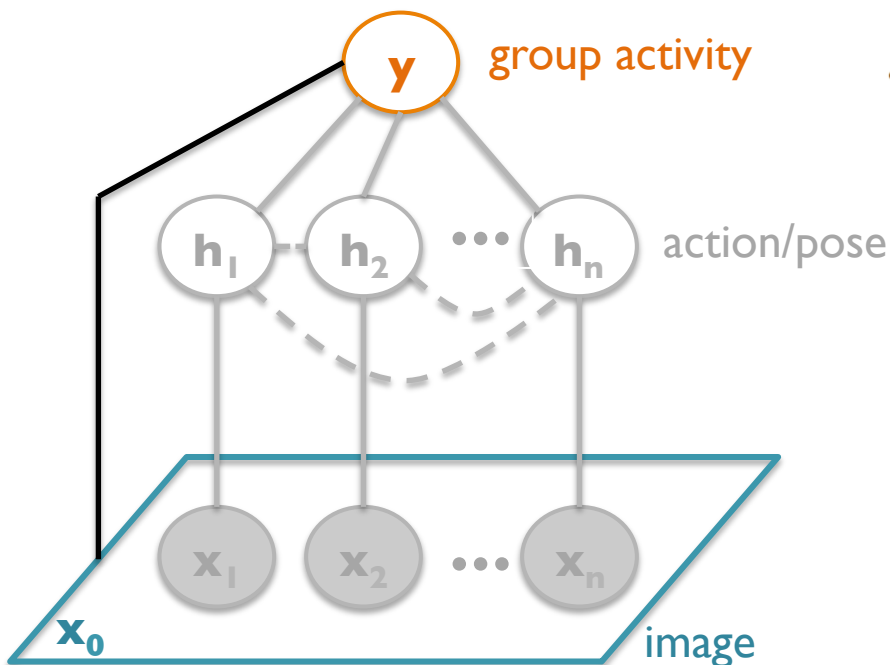


Hierarchical Model For Group Activity Recognition



$$\begin{aligned}
 f_w(\mathbf{x}, \mathbf{h}, y; \mathcal{G}) &= w^\top \Psi(y, \mathbf{h}, \mathbf{x}; \mathcal{G}) \\
 &= w_0^\top \phi_0(y, x_0) \\
 &+ \sum_{j \in \mathcal{V}} w_1^\top \phi_1(x_j, h_j) \\
 &+ \sum_{j \in \mathcal{V}} w_2^\top \phi_2(y, h_j) \\
 &+ \sum_{j, k \in \mathcal{E}} w_3^\top \phi_3(y, h_j, h_k)
 \end{aligned}$$

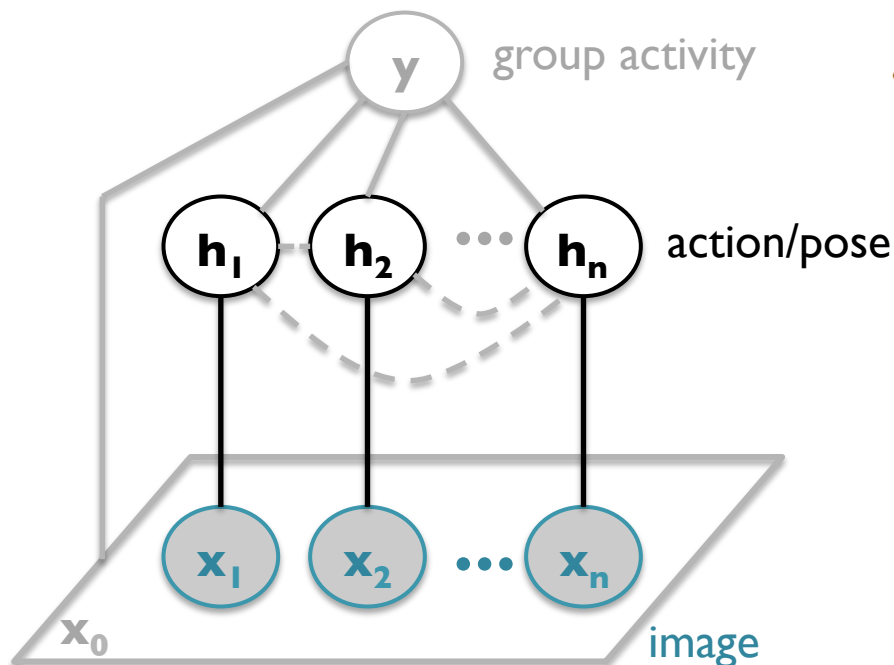
Hierarchical Model For Group Activity Recognition



$$\begin{aligned}
 f_w(\mathbf{x}, \mathbf{h}, y; \mathcal{G}) &= w^\top \Psi(y, \mathbf{h}, \mathbf{x}; \mathcal{G}) \\
 &= w_0^\top \phi_0(y, x_0) \\
 &+ \sum_{j \in \mathcal{V}} w_1^\top \phi_1(x_j, h_j) \\
 &+ \sum_{j \in \mathcal{V}} w_2^\top \phi_2(y, h_j) \\
 &+ \sum_{j, k \in \mathcal{E}} w_3^\top \phi_3(y, h_j, h_k)
 \end{aligned}$$

Direct Measurement for Groups.

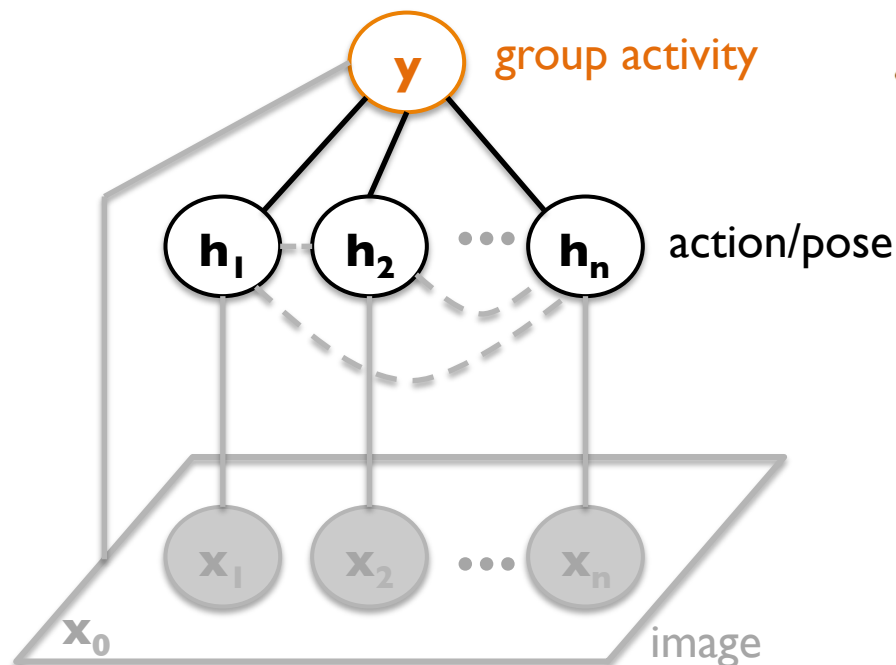
Hierarchical Model For Group Activity Recognition



$$\begin{aligned}
 f_w(\mathbf{x}, \mathbf{h}, y; \mathcal{G}) &= w^\top \Psi(y, \mathbf{h}, \mathbf{x}; \mathcal{G}) \\
 &= w_0^\top \phi_0(y, x_0) \\
 &+ \sum_{j \in \mathcal{V}} w_1^\top \phi_1(x_j, h_j) \\
 &+ \sum_{j \in \mathcal{V}} w_2^\top \phi_2(y, h_j) \\
 &+ \sum_{j, k \in \mathcal{E}} w_3^\top \phi_3(y, h_j, h_k)
 \end{aligned}$$

Measurements for Action/Pose.

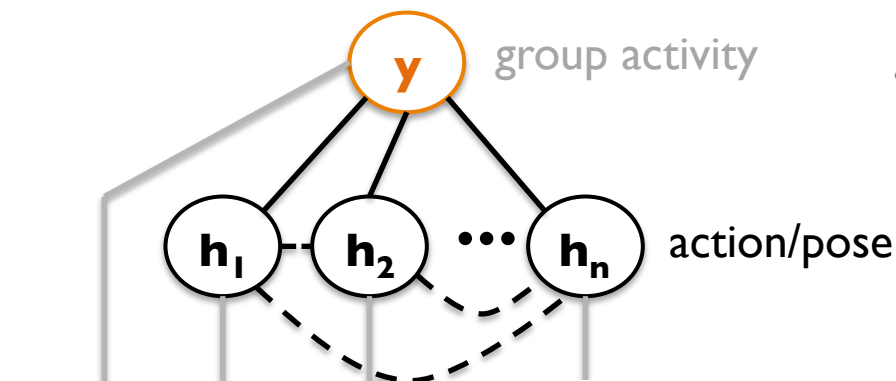
Hierarchical Model For Group Activity Recognition



$$\begin{aligned}
 f_w(\mathbf{x}, \mathbf{h}, y; \mathcal{G}) &= w^\top \Psi(y, \mathbf{h}, \mathbf{x}; \mathcal{G}) \\
 &= w_0^\top \phi_0(y, x_0) \\
 &+ \sum_{j \in \mathcal{V}} w_1^\top \phi_1(x_j, h_j) \\
 &+ \sum_{j \in \mathcal{V}} w_2^\top \phi_2(y, h_j) \\
 &+ \sum_{j, k \in \mathcal{E}} w_3^\top \phi_3(y, h_j, h_k)
 \end{aligned}$$

Group-Individual Relationship
via co-occurrence.

Hierarchical Model For Group Activity Recognition



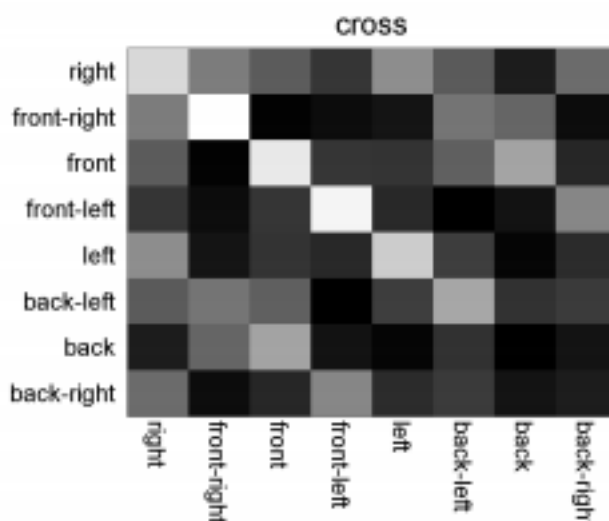
$$f_w(\mathbf{x}, \mathbf{h}, y; \mathcal{G}) = w^\top \Psi(y, \mathbf{h}, \mathbf{x}; \mathcal{G})$$

$$= w_0^\top \phi_0(y, x_0)$$

$$+ \sum_{j \in \mathcal{V}} w_1^\top \phi_1(x_j, h_j)$$

$$+ \sum_{j \in \mathcal{V}} w_2^\top \phi_2(y, h_j)$$

$$+ \sum_{j, k \in \mathcal{E}} w_3^\top \phi_3(y, h_j, h_k)$$



Group Dependent
Pairwise Individual Relationship
via co-occurrence.

Dataset

- Collective Activity Dataset
 - Available at <http://wwwweb.eecs.umich.edu/vision/activity-dataset.html>
 - 44 videos with multiple people
 - Crossing, Waiting, Queuing, Walking, Talking



Crossing



Waiting



Queuing



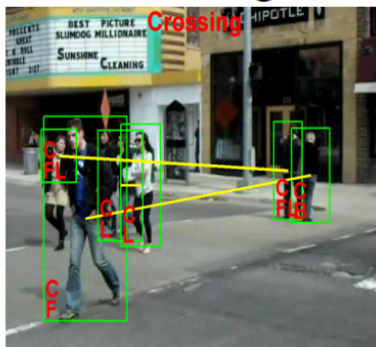
Walking



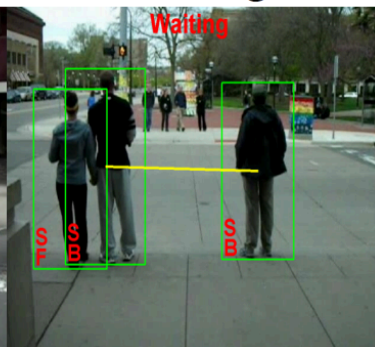
Talking

Qualitative Examples

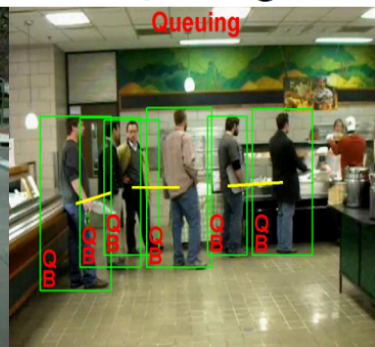
Crossing



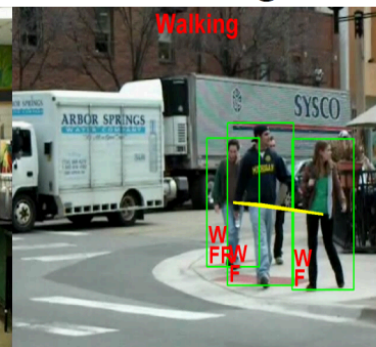
Waiting



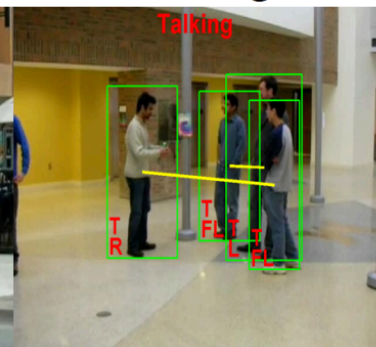
Queuing



Walking



Talking



Qualitative Examples

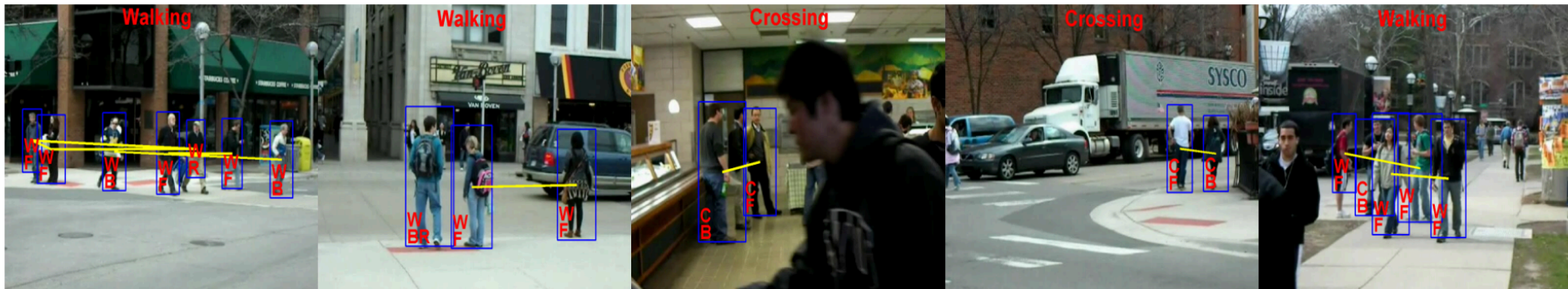
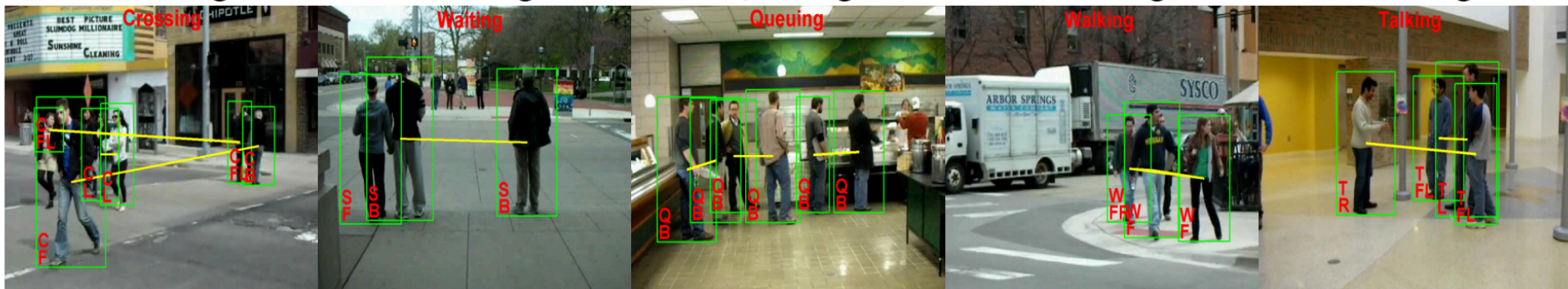
Crossing

Waiting

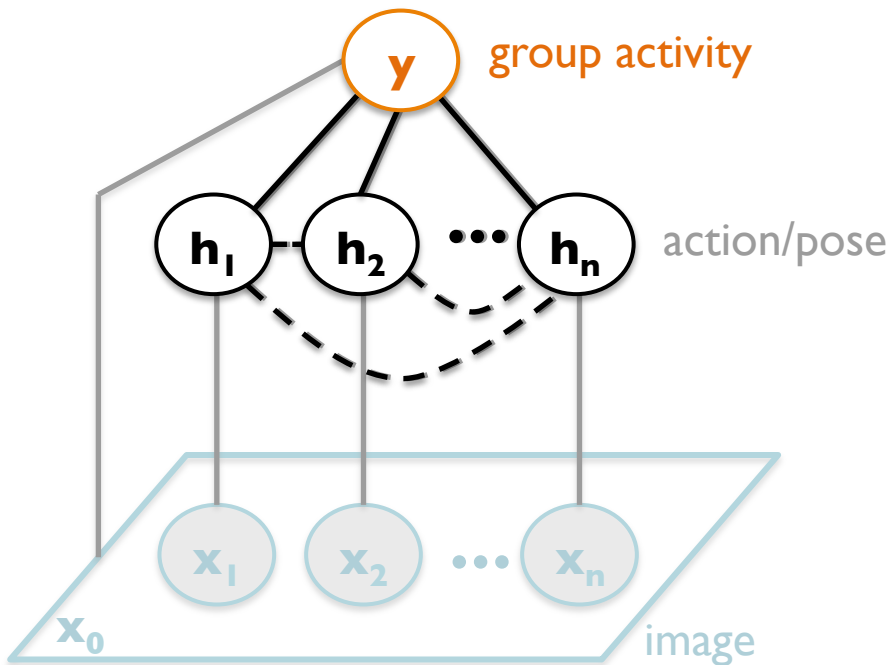
Queuing

Walking

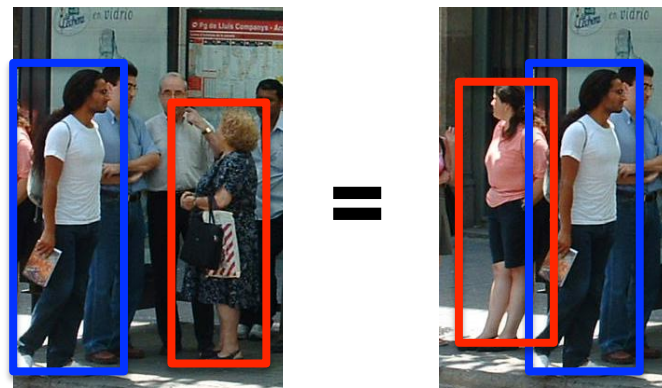
Talking



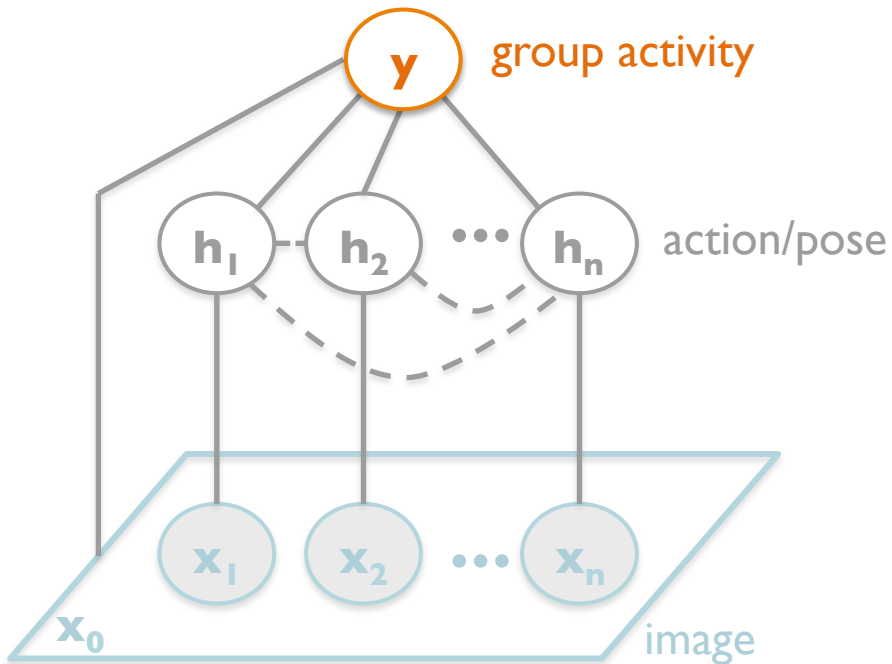
Limitations



No spatial relationship



Limitations



No spatial relationship

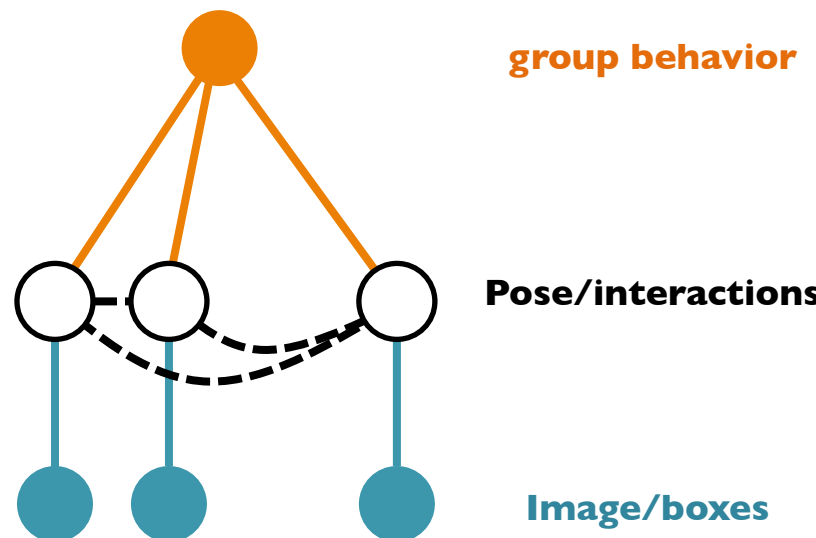
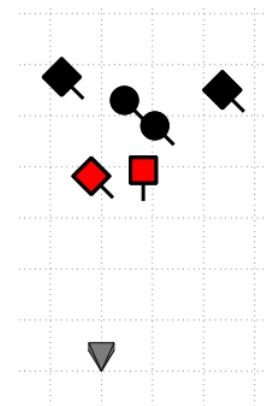
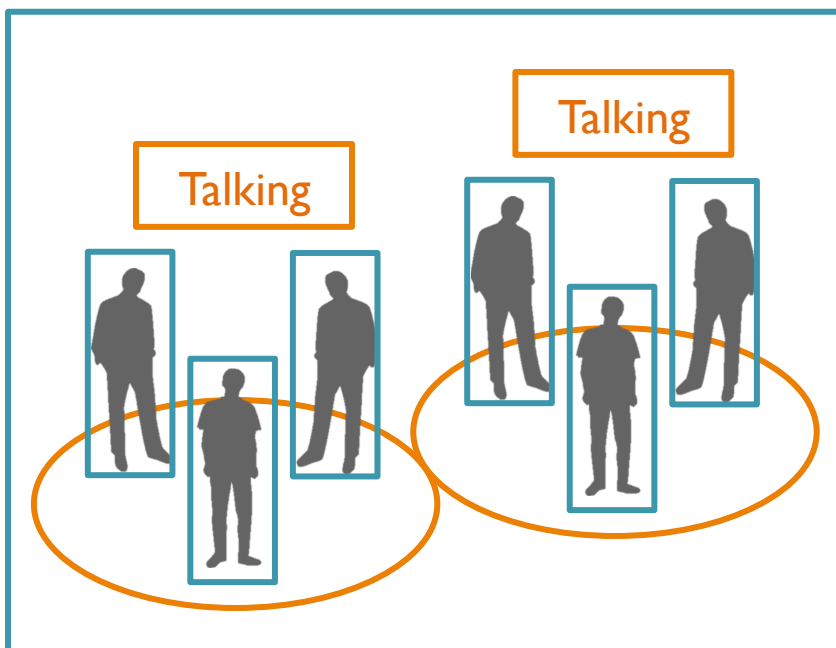
One group assumption



Group Discovery (Choi ECCV14)

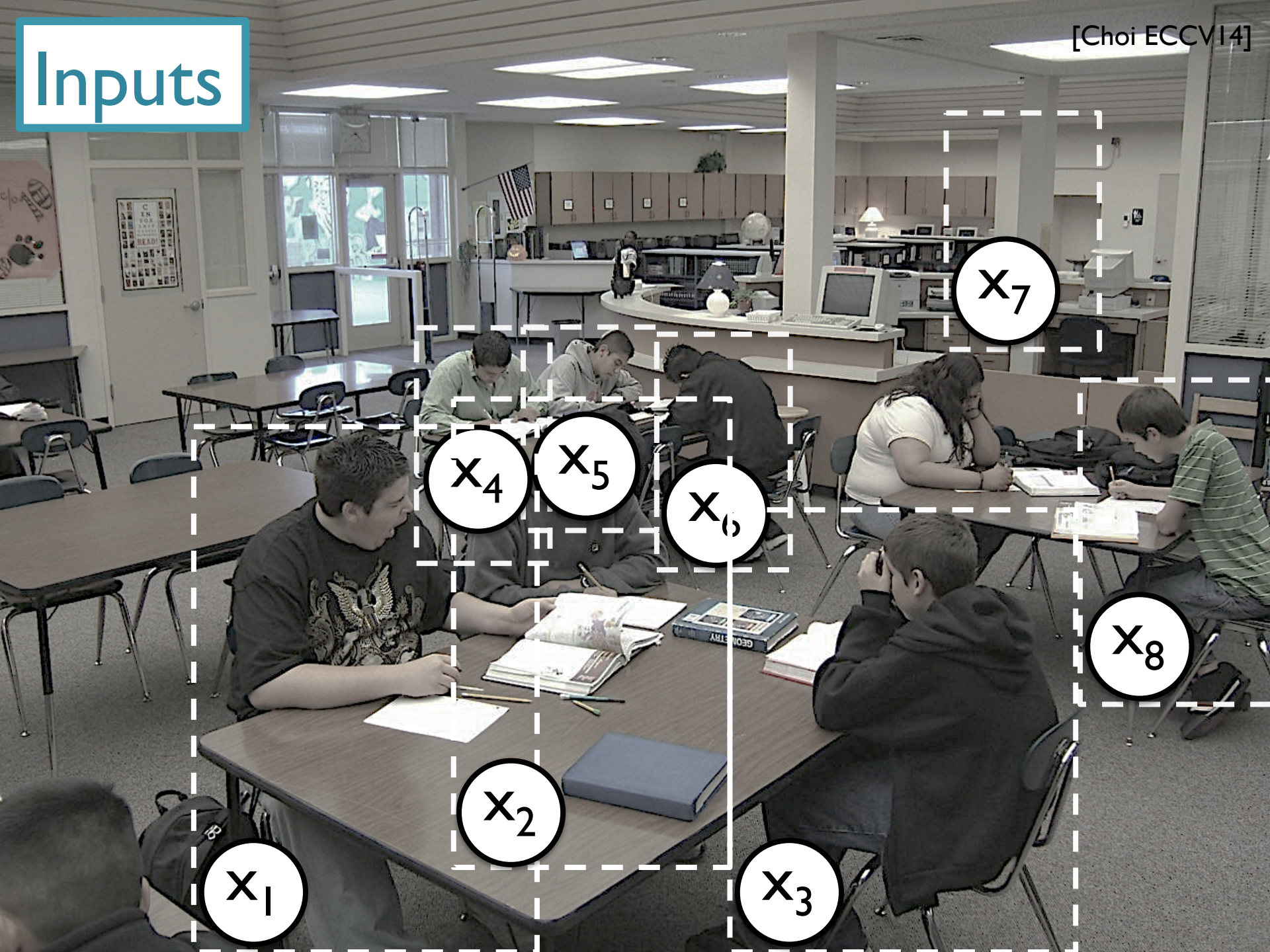
Input: image and bounding boxes

Output: clustered groups
with activity label





Inputs



X₁

X₂

X₄

X₅

X₆

X₃

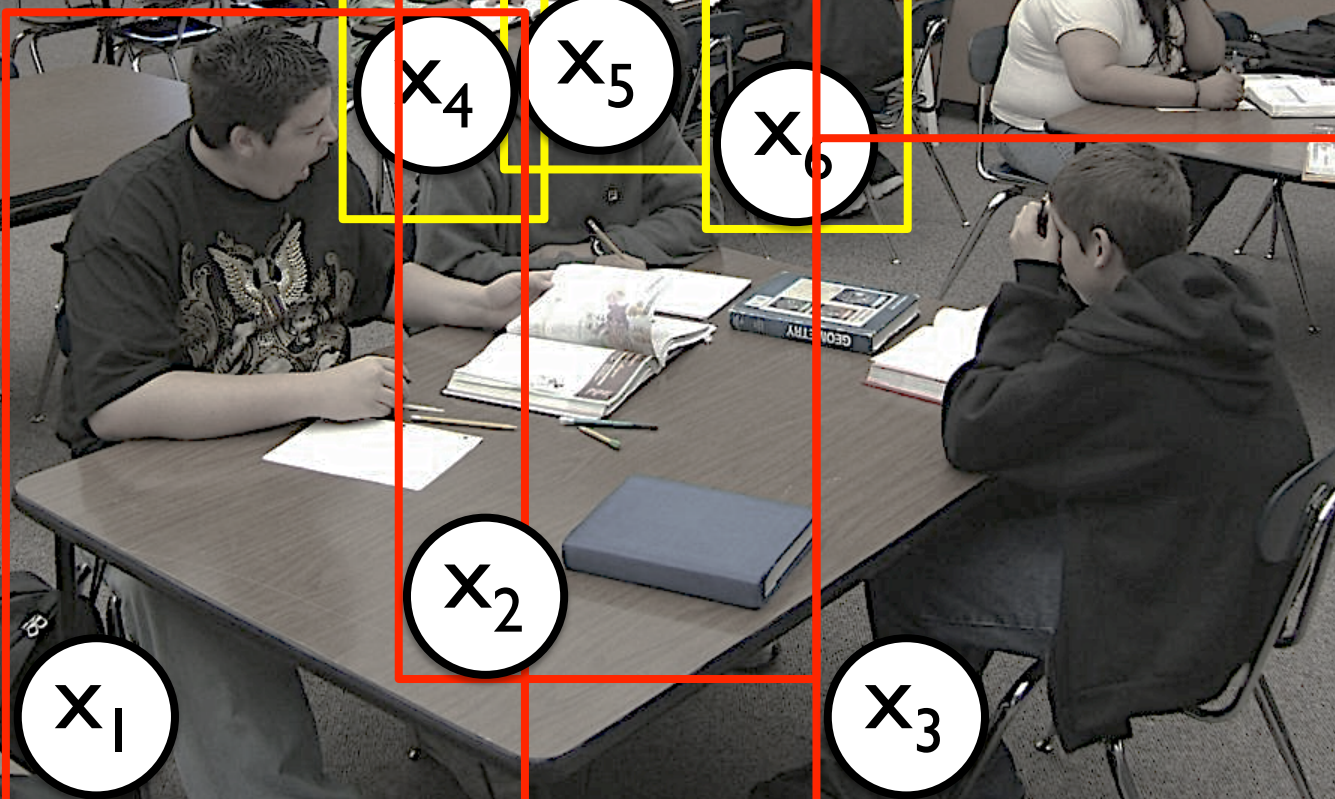
X₇

X₈

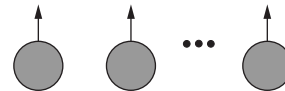
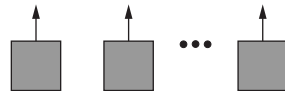
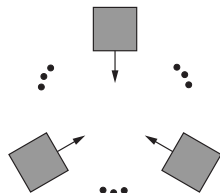
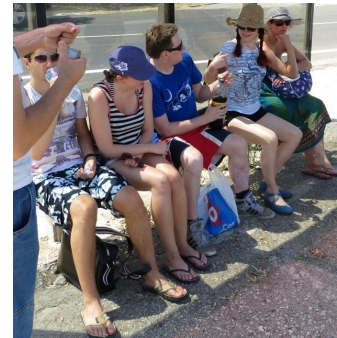
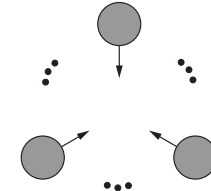
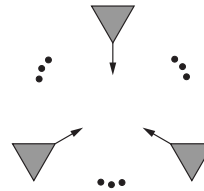
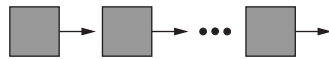
Outputs

c_1 = facing-each-other-sitting

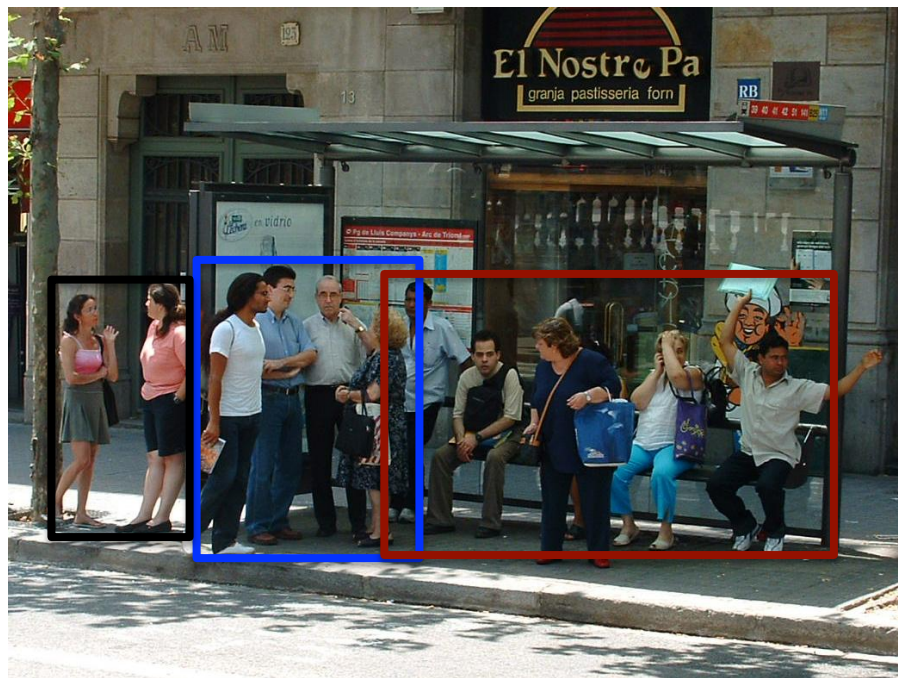
c_2 = facing-each-other-sitting



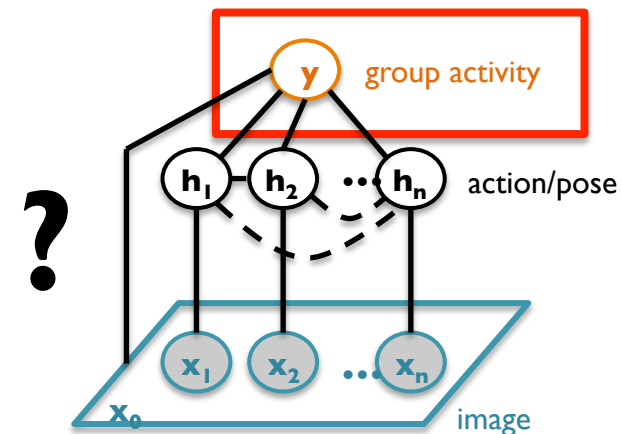
Structured Groups



Challenges



- Unknown number of groups.



Challenges



- Unknown number of groups.
- Large intra-class variation.

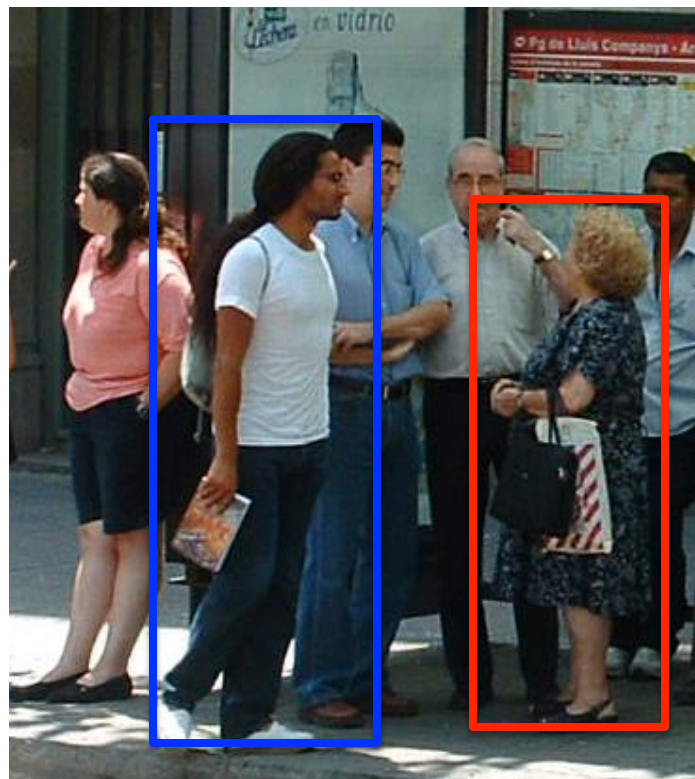
Interactions as Key Social Signal



**Standing
Facing-right**

- Individual Pose
 - **Weak** social signal.

Interactions as Key Social Signal



**Standing
Facing-right**

**Standing
Facing-left**

- Individual Pose
 - **Weak** social signal.
- Pair Interaction
 - **Strong** social signal within a group.

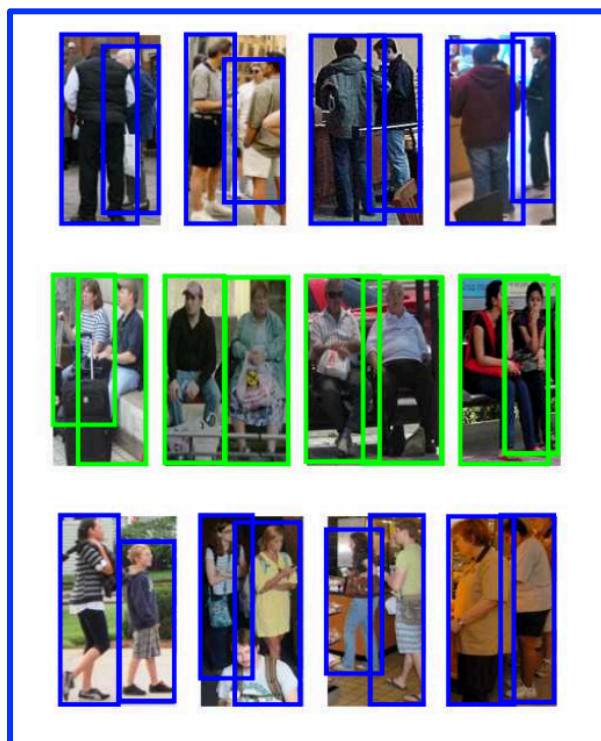
Interactions as Key Social Signal



- Individual Pose
 - **Weak** social signal.
- Pair Interaction
 - **Strong** social signal within a group.
 - **Repulsive** signal in different groups.

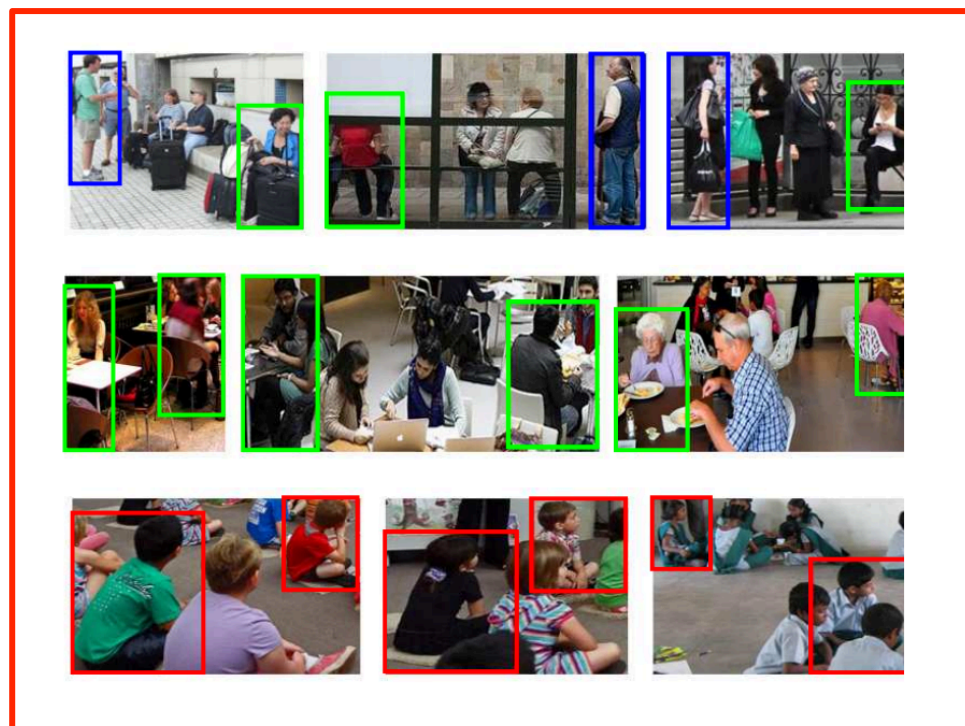
Interactions as Key Social Signal

Intra-group Interactions

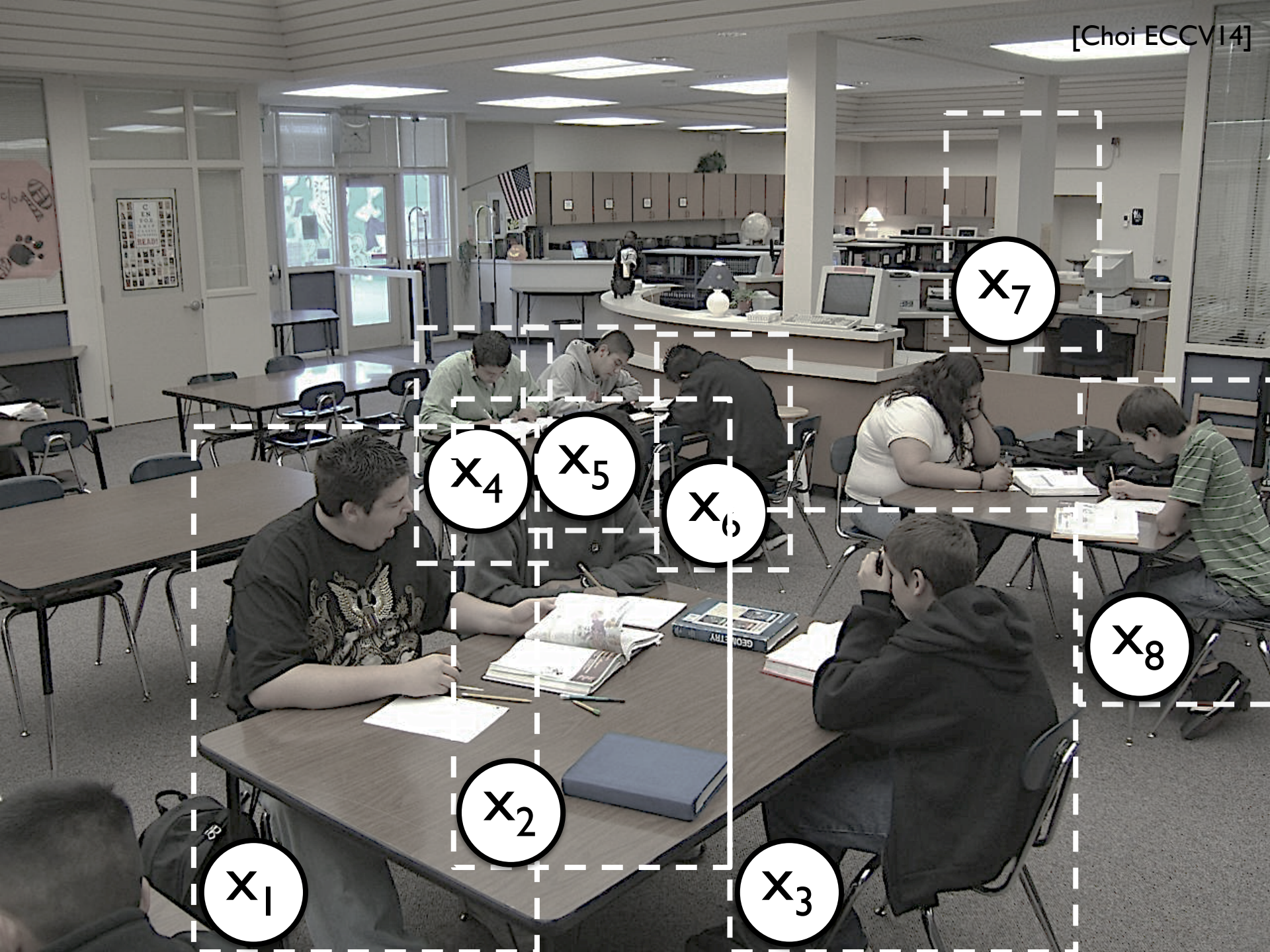


Attractive Potential for a Group

Inter-group Interactions



Repulsive Potential for a Group



X₁

X₂

X₄

X₅

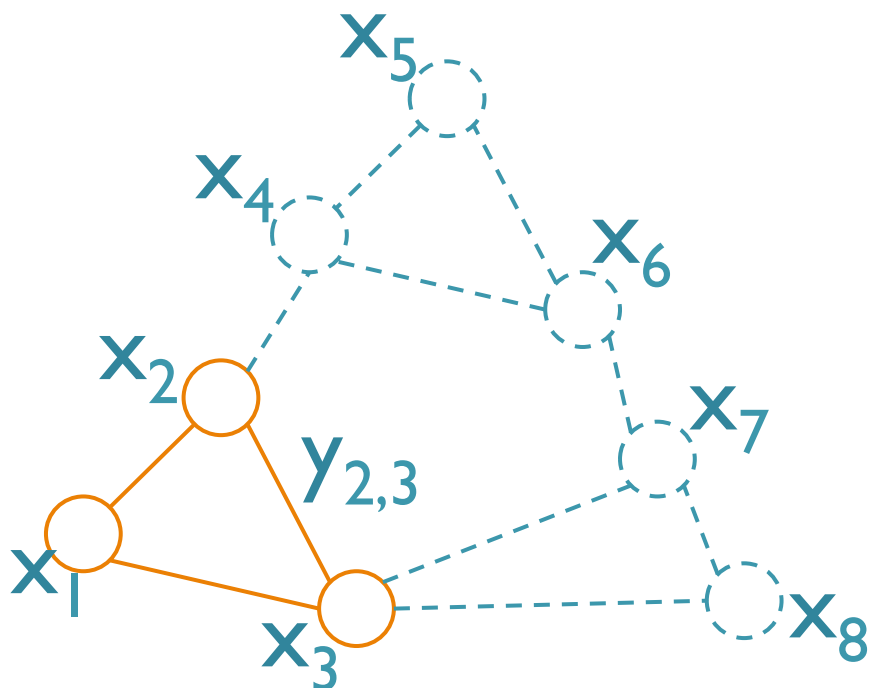
X₆

X₃

X₇

X₈

Group Discovery Model



$$C_1 = \{c_1 = \text{facing-each-other-sitting}, \\ H_1 = [1, 1, 1, \dots, 0]\}$$

Nodes

- Individual likelihood to belong in a group C .

Edges

- Likelihood of interaction in a group C .

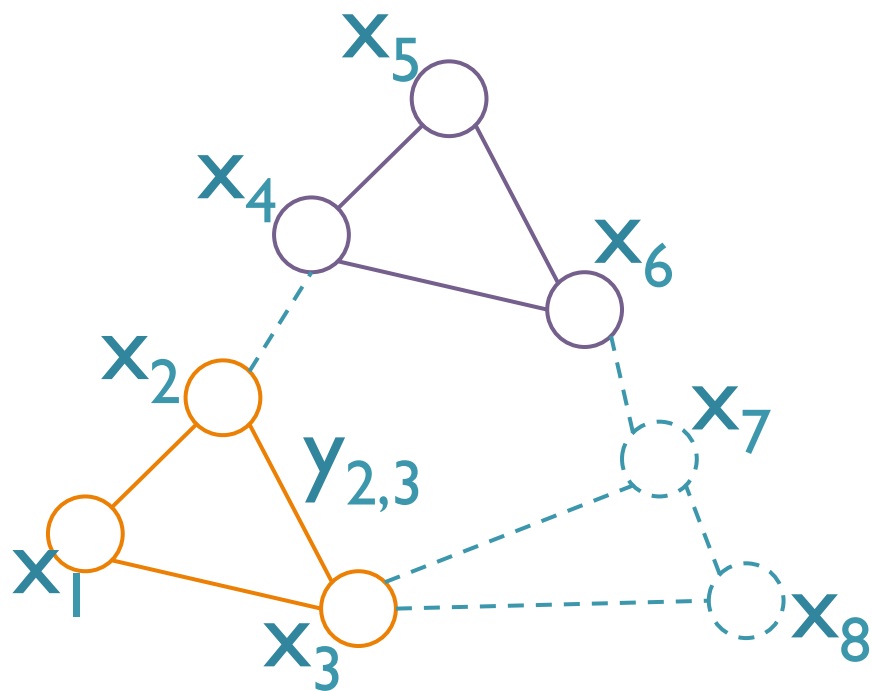
Solid lines

- Members belonging to the group.

Dashed lines

- Not belonging to the group.

Group Discovery Model



$$C_1 = \{c_1 = \text{facing-each-other-sitting}, \\ H_1 = [1, 1, 1, \dots, 0]\}$$

Nodes

- Individual likelihood to belong in a group C .

Edges

- Likelihood of interaction in a group C .

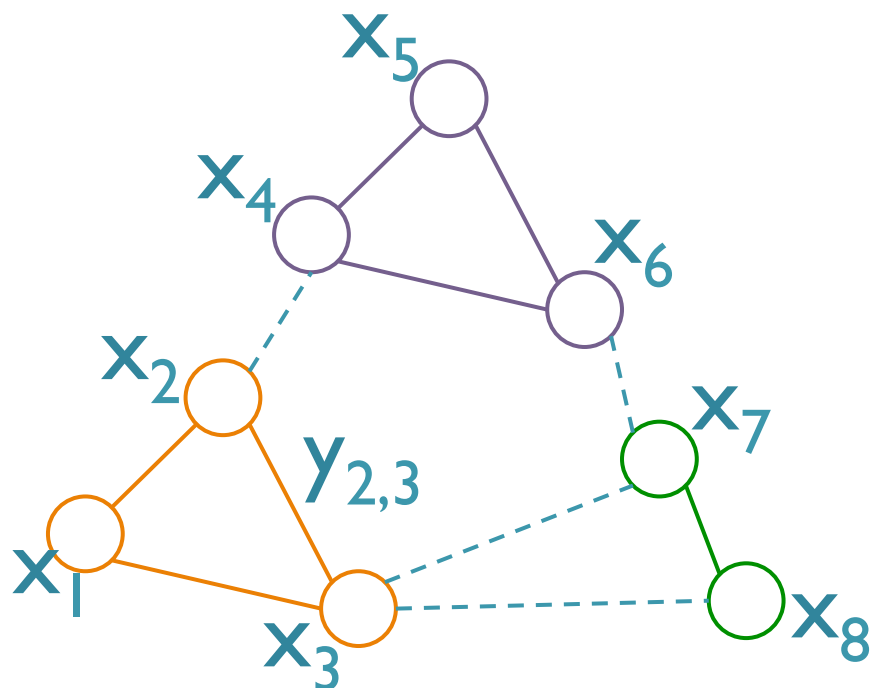
Solid lines

- Members belonging to the group.

Dashed lines

- Not belonging to the group.

Group Discovery Model



$$C_1 = \{c_1 = \text{facing-each-other-sitting}, \\ H_1 = [1, 1, 1, \dots, 0]\}$$

Nodes

- Individual likelihood to belong in a group C .

Edges

- Likelihood of interaction in a group C .

Solid lines

- Members belonging to the group.

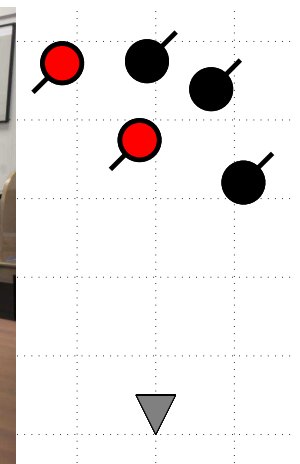
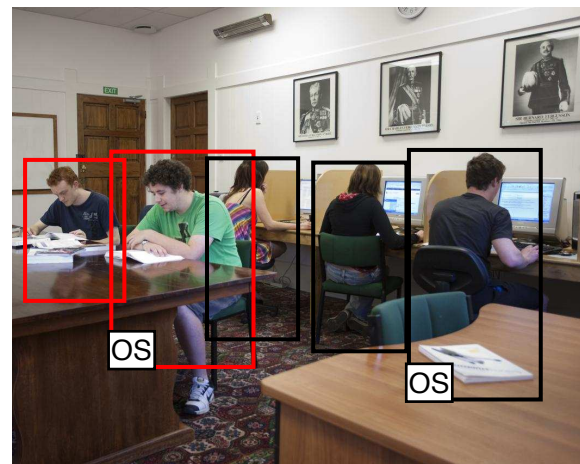
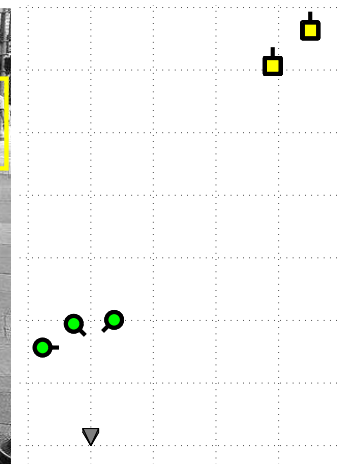
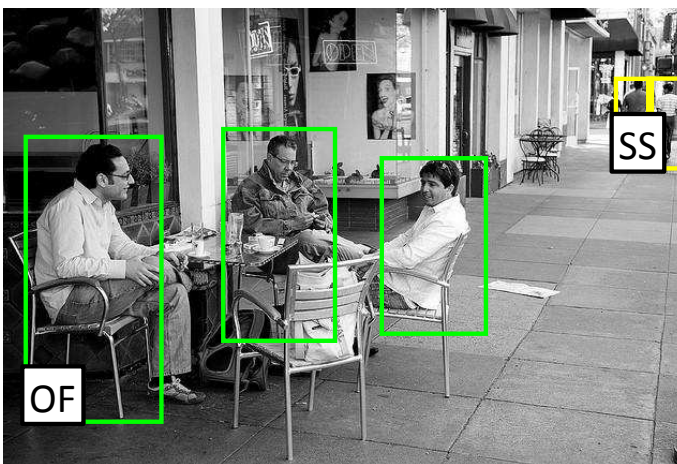
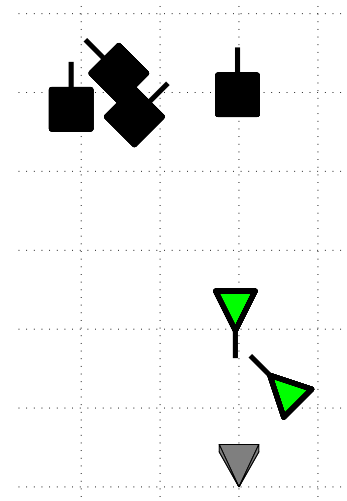
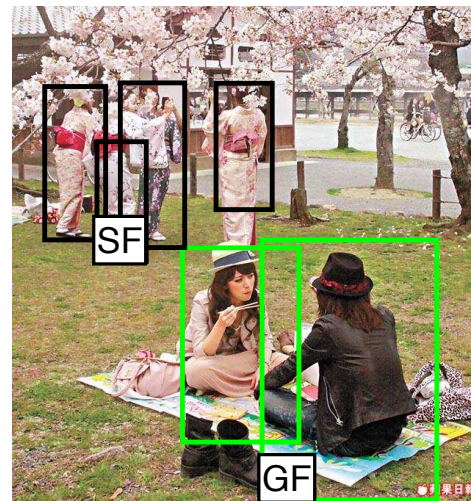
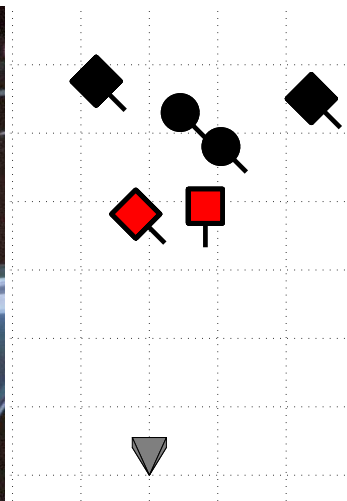
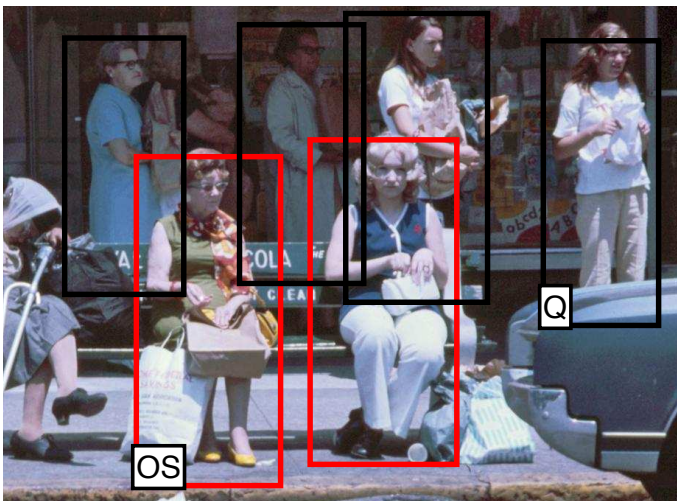
Dashed lines

- Not belonging to the group.

Datasets

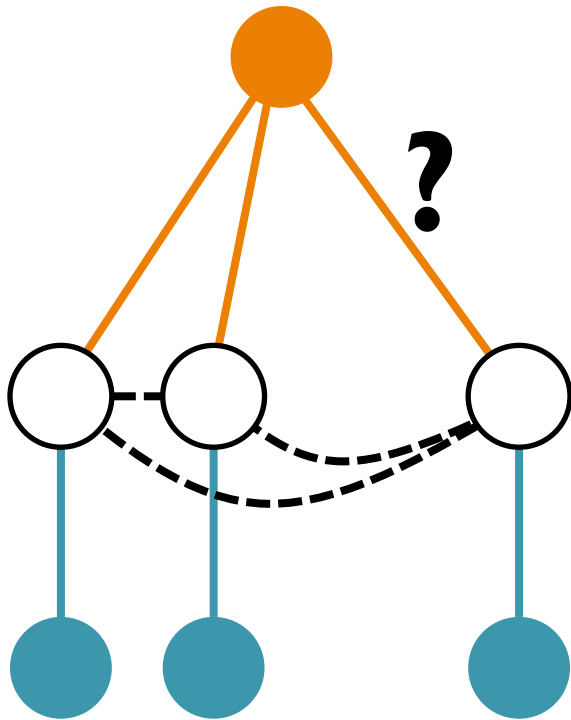
- Structured Group Dataset
 - Available at <http://cvgl.stanford.edu/projects/groupdiscovery/>
 - 588 images with 5,415 individuals and 1,719 groups.
 - 7 structured groups:
 - Queue, Standing-Facing, Sitting-Object-Facing, Sitting-Ground-Facing, Standing-Side, Sitting-Object-Side, Sitting-Ground-Side.

Group Discovery Results



Learning the Models

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

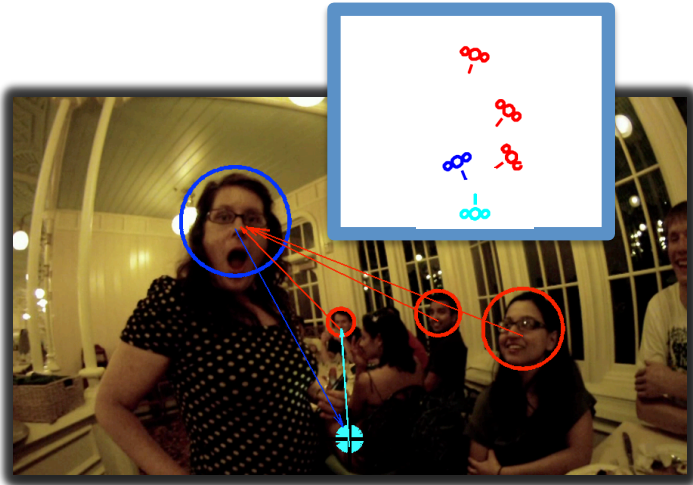


$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i, \quad \text{s.t. } \forall i, \xi_i \geq 0$$

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

Tsochantaridis et al, 2004

Signals for Social Statics



The left touches the right's head.



Conversation

