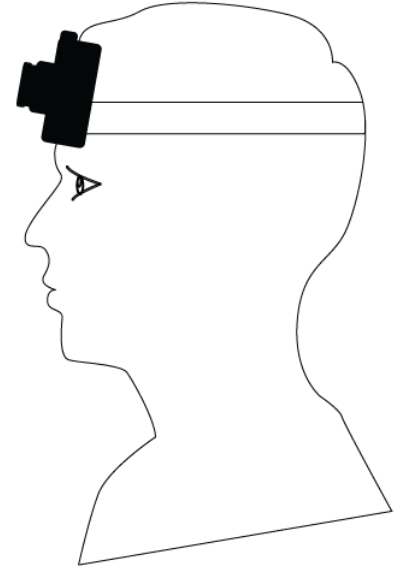# Social Attention
## :What can first person cameras tell us about our social interactions?
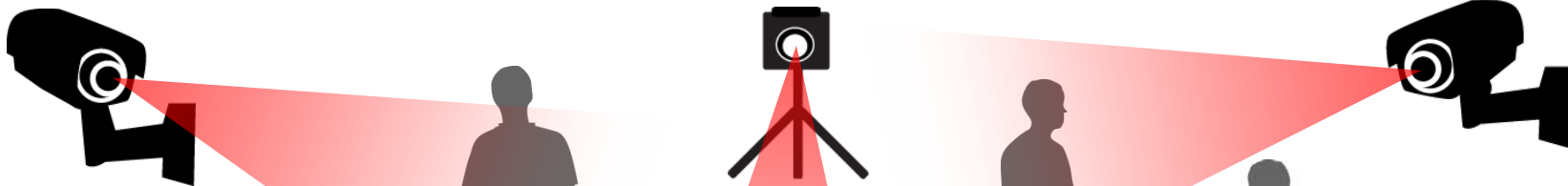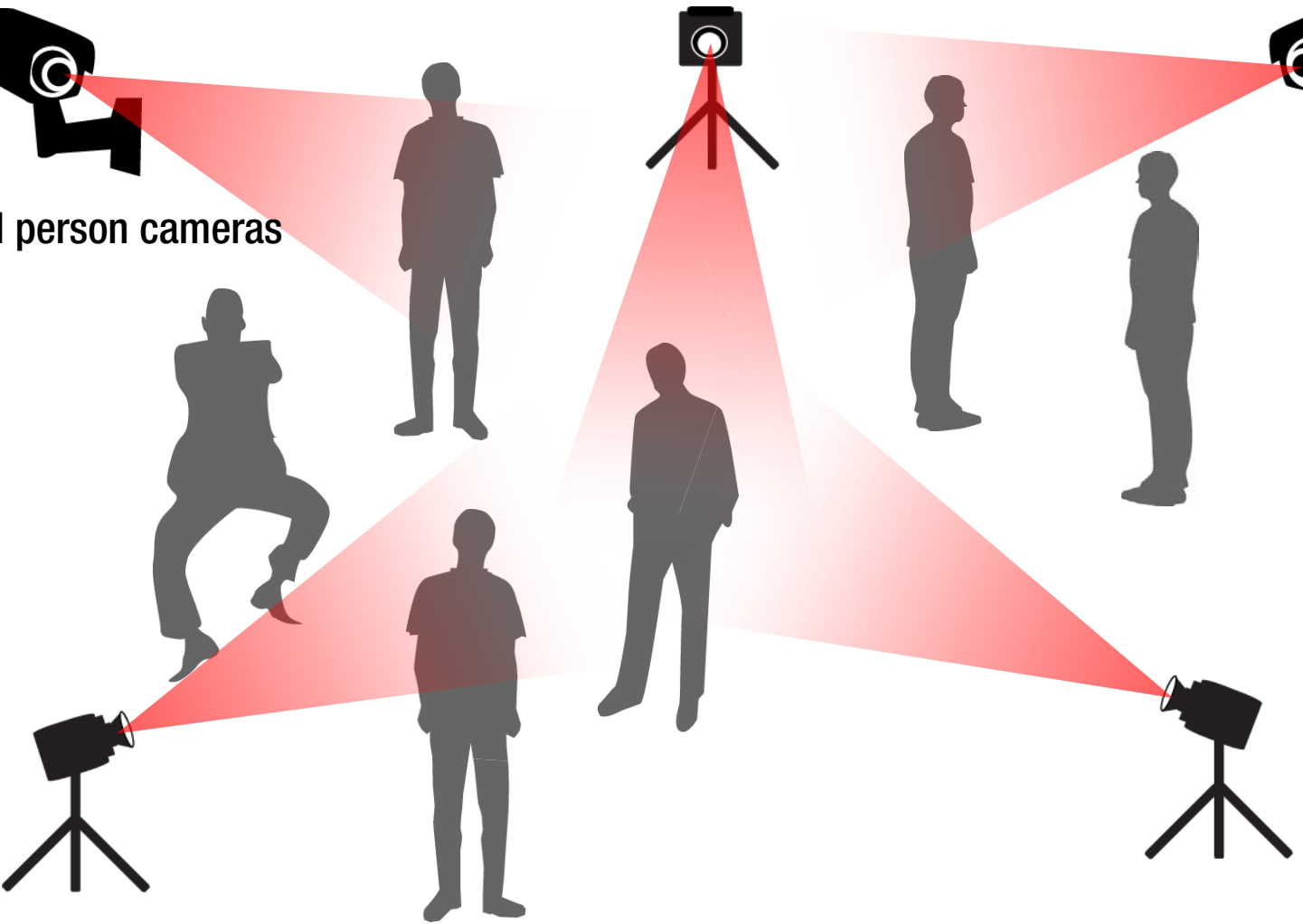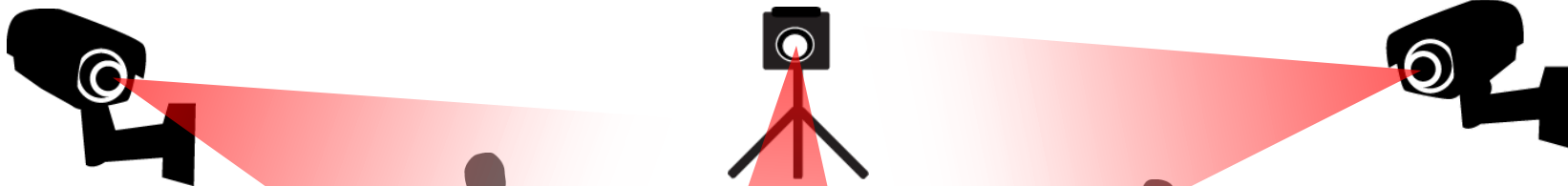
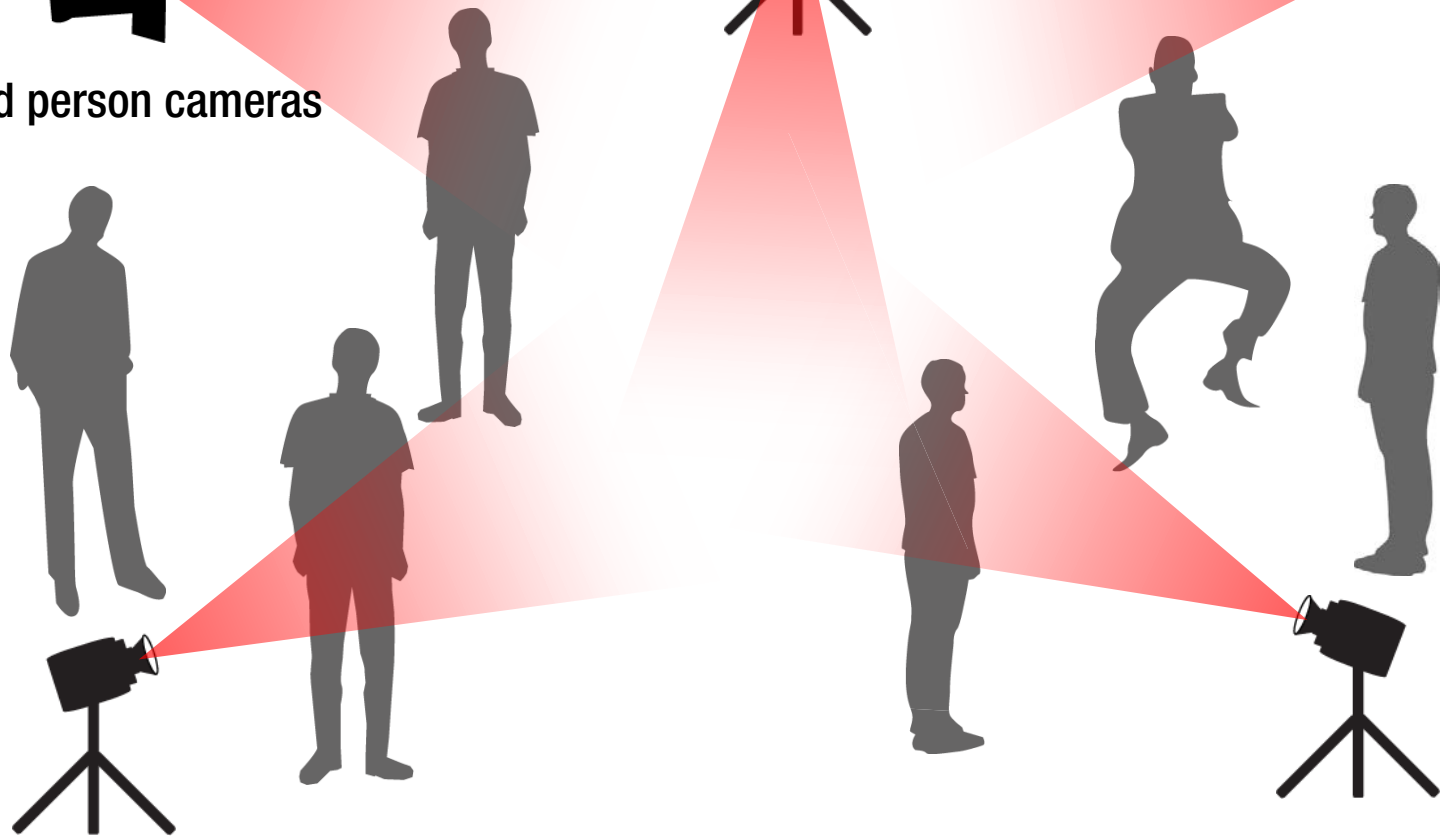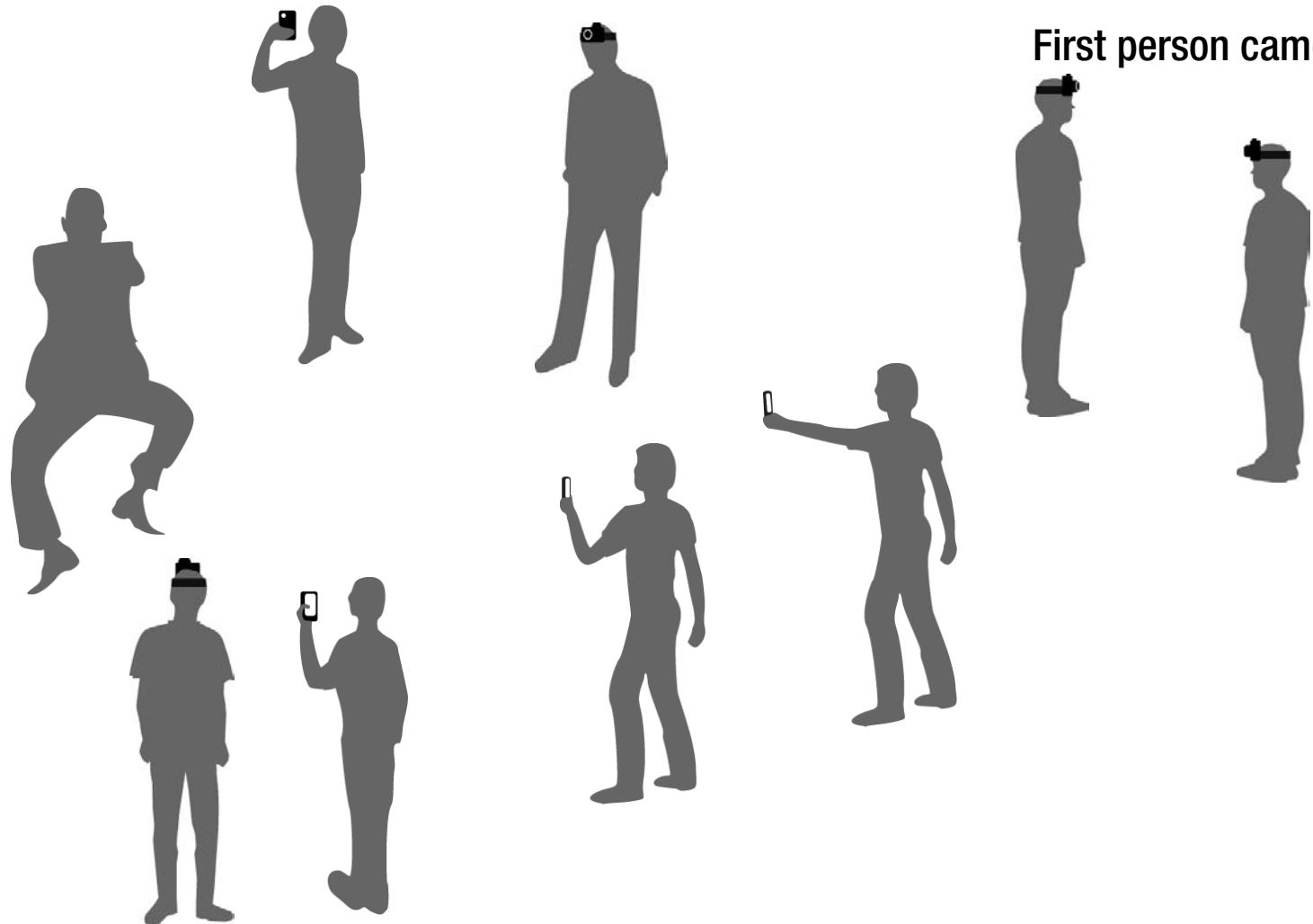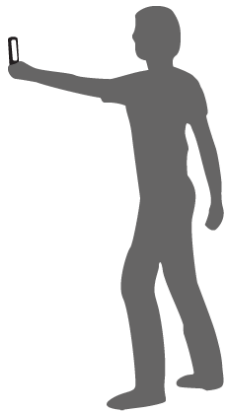# First person cameras are ideal sensors to measure social behaviors.
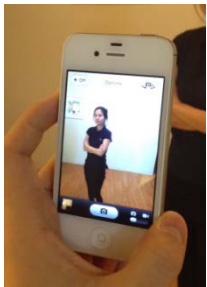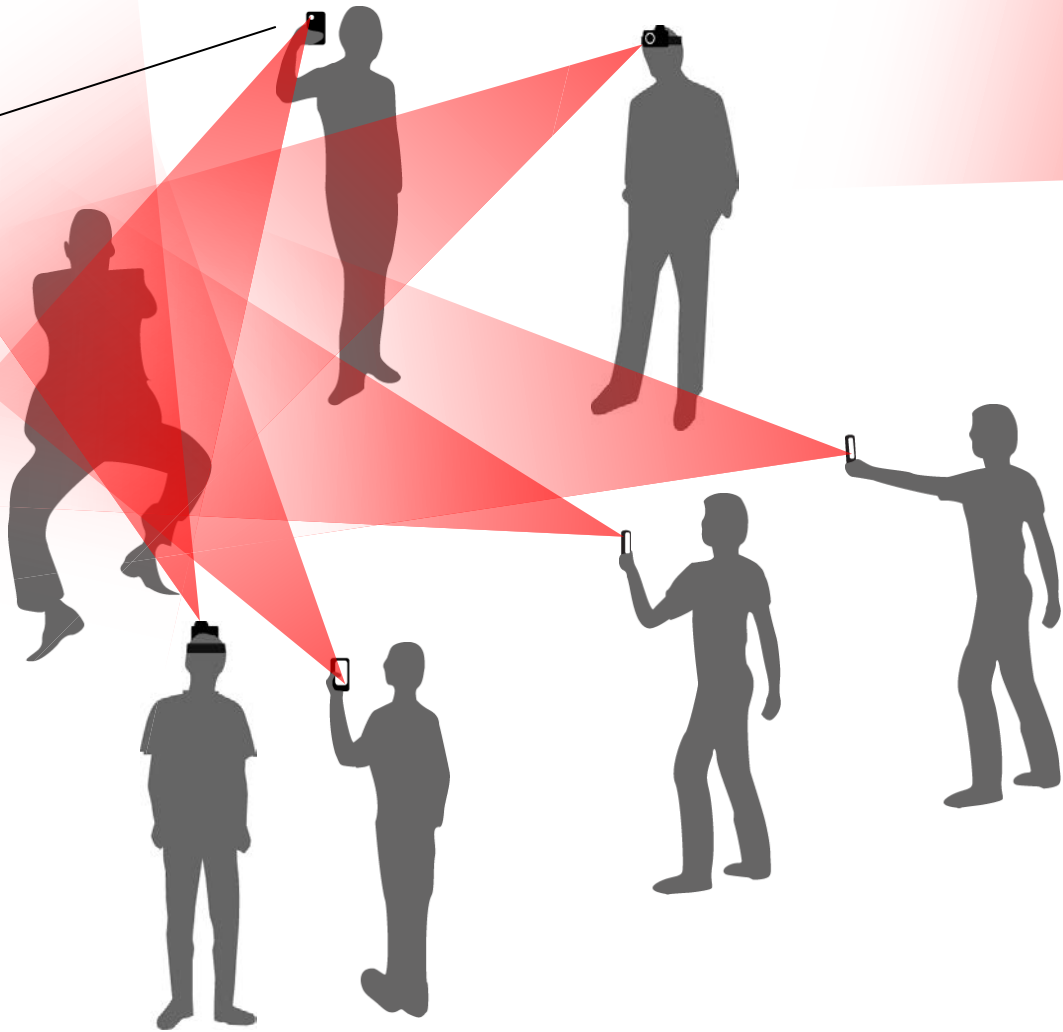
Third person cameras
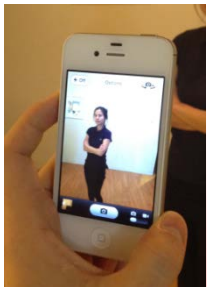
Third person cameras
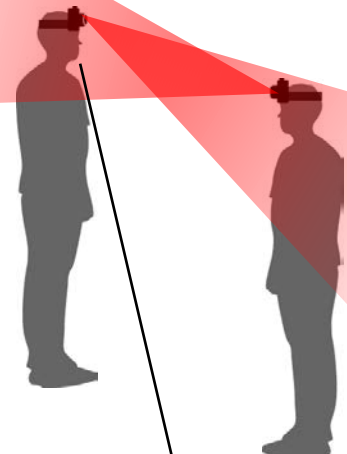
Third person cameras

First person cameras

First person cameras

First person camera

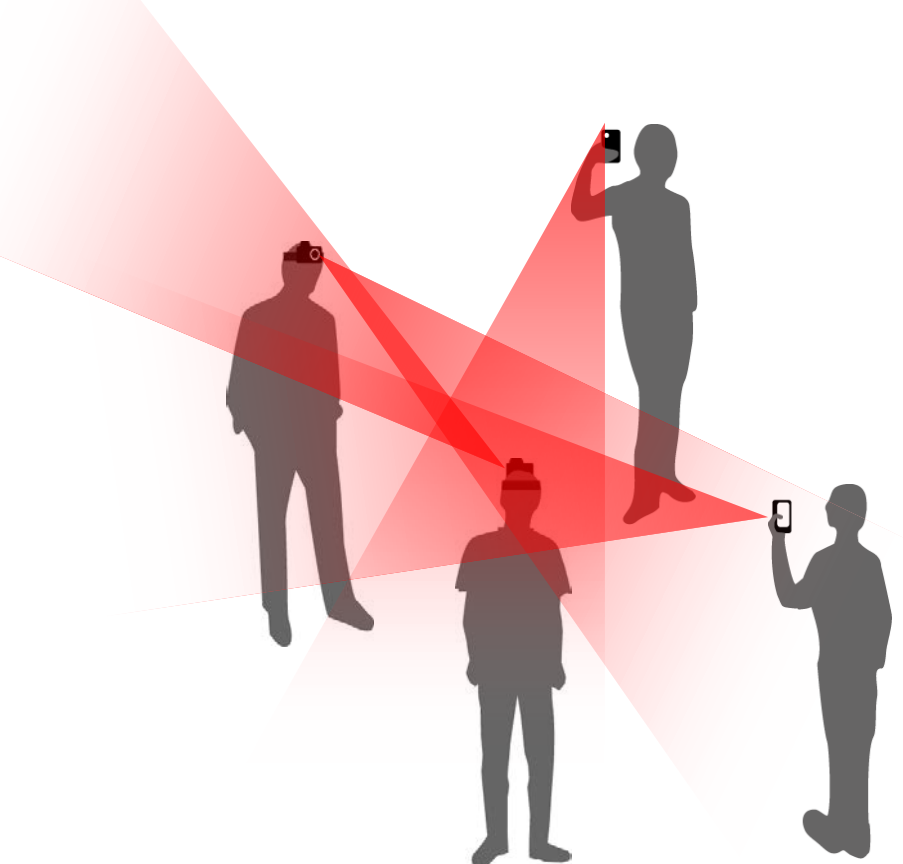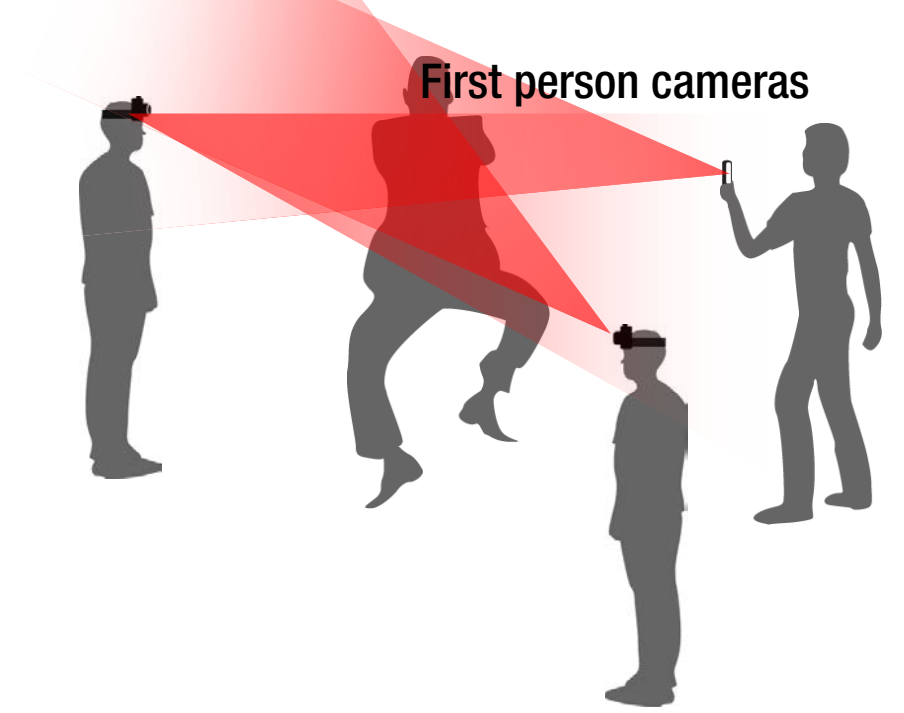First person cameras

First person camera

First person cameras

First person cameras are ideal sensors to measure <u>social behaviors</u> because they

1. secure the best views
2. produce more views of events of greater interest
3. follow social behaviors

Gaze direction
Camera direction (frustum)

Gaze direction

Camera direction (frustum)

Gaze direction

Camera direction (frustum)

Joint attention

**Joint Attention**
Autism spectrum disorder (ASD)
Attention deficit hyperactivity disorder (ADHD)

→ Gaze direction

Camera direction (frustum)

Joint attention

Primary gaze ray

# Eye-in-head motion

—Primary gaze direction
o Center of eyes

Primary gaze ray

$\mathcal{W}$

$\mathcal{C}$

Cone-shaped distribution
of the point of regard

Camera pose

Cone-shaped gaze distribution

Primary gaze ray

3D gaze model registration

Cone-shaped gaze model

Gaze direction

Eye center

Top view

# Attention density field:

$$f(\mathbf{x}) = k\left(\frac{\text{dist}(\mathbf{l}, \mathbf{x})}{h}\right)$$

where k is a kernel density function (  ).



dist(l, x)

x

l

Gaze direction

h

Eye center

Top view

## Attention density field:

$$f(\mathsf{x}) = \sum_i k\left(\frac{\mathrm{dist}(\mathsf{l}_i, \mathsf{x})}{h}\right)$$

where k is a kernel density function ( ).



- - - Primary gaze ray
○ Center of eyes

Gaze directions

Eye centers

Top view

Attention density field:

$$f(\mathbf{x}) = \sum_i k\left(\frac{\mathrm{dist}(\mathbf{l}_i, \mathbf{x})}{h}\right)$$

where k is a kernel density function (  ).



--- Primary gaze ray
○ Center of eyes

Gaze directions

Multiple social cliques

Eye centers

Top view

Attention density field:

$$f(\mathsf{x}) = \sum_i k\left(\frac{\mathrm{dist}(\mathsf{l}_i, \mathsf{x})}{h}\right)$$

where k is a kernel density function (  ).

3D joint attention

$\approx$ Modes of attention density field

# Mode-seeking: Gaze Concurrences

Two groups

Multiple groups

3D joint attention

Gaze directions

# Party Scene: 4 groups in a room



Couch

Pool table

4 gaze concurrences

1x speed

Dining area

Ping pong table

Poster

Whiteboard

Poster

Realtime Joint Attention

# Social Anomaly Detection

⟶ Predicted gaze direction

Can we predict <u>without</u> first person cameras?

Target

Target

Target

First person

Second person

Third person

3D estimation error < 5cm

Measurement accuracy

Learning

Prediction

Noninvasiveness

Gaze detection

True positive head detection

True positive gaze detection

Halloween show

Children

Social saliency

Top view

Social saliency: likelihood of joint attention

Halloween show

Children

Social saliency

Top view

⊖ : Ground truth joint attention

⊕ : Head location

Location of performer:
Ground truth joint attention

Children

$$g\left(\underbrace{\{⊕_i\}_{i=1}^N}_{\text{Social formation}}\right) = ⊖$$

where $N$ is the number of social members.

**: Ground truth joint attention**

**: Head location**

$$g\left(\underbrace{\left\{\oplus_i\right\}_{i=1}^N}_{\text{Social formation}}\right) = \ominus$$

where *N* is the number of social members.

cf. $g\left(\left\{\oplus \rightarrow_i\right\}_{i=1}^N\right) = \ominus$

Geometric localization: deterministic

: Ground truth joint attention

: Head location

$$g\left(\underbrace{\left\{ \oplus_i \right\}_{i=1}^N}_{\text{Social formation}}\right) = \ominus$$

where $N$ is the number of social members.

cf. $g\left(\left\{ \oplus \longrightarrow_i \right\}_{i=1}^N\right) = \ominus$

Geometric localization: deterministic

$$g\left(\underbrace{\{\oplus_i\}_{i=1}^{N}}_{\text{Social formation}}\right) = \ominus$$

where *N* is the number of social members.

## Scale variation

$\ominus$ : Ground truth joint attention

$\oplus$ : Head location
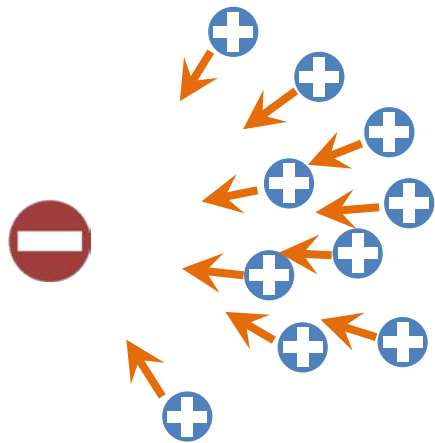
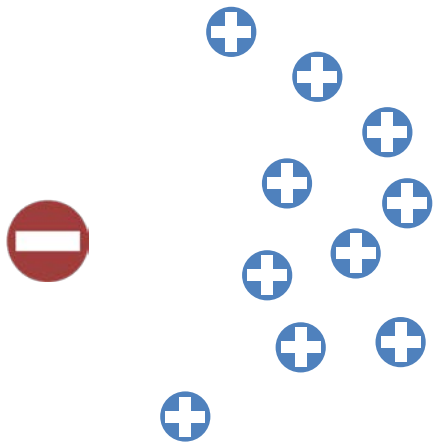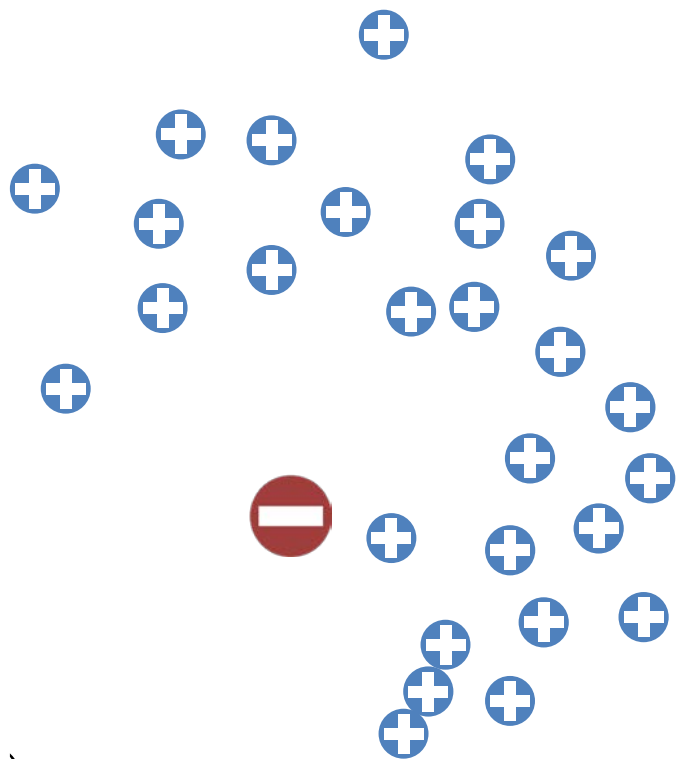$$g\left(\underbrace{\{\oplus_i\}_{i=1}^{N}}_{\text{Social formation}}\right) = \ominus$$

where *N* is the number of social members.

Scale variation
Orientation variation

⊖ : Ground truth joint attention

⊕ : Head location

**Representation:
Social Dipole Moment**

Water molecule, H₂O



$$\mathbf{q}_e = \sum_i^N (\mathbf{e} - \mathbf{p}_i)$$

Electric dipole moment

Water molecule, H₂O

**O** −

**H** +

104.5°

**H** +

$$\mathbf{q}_e = \sum_{i}^{N} (\mathbf{e} - \mathbf{p}_i)$$

Electric dipole moment

Water molecule, H₂O

Water molecule, H₂O

Water molecule, H$_2$O

**Social dipole moment**

$$\mathbf{q} = \mathbf{s} - \frac{1}{N}\sum_{i}^{N}\mathbf{p}_i = \mathbf{s} - \mathbf{c}$$

**Electric dipole moment**

$$\mathbf{q}_e = \sum_{i}^{N}(\mathbf{e} - \mathbf{p}_i)$$

**Legend:**
- ⊕ Social member
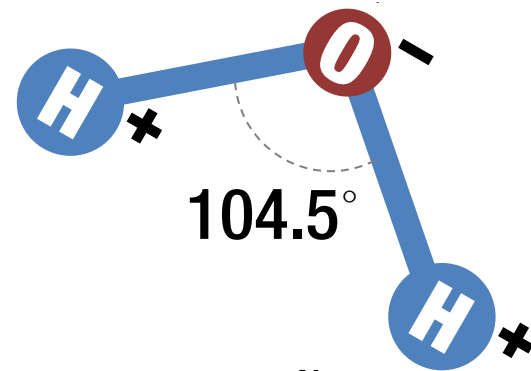- ⊖ Joint attention
- ◑ Center of mass
- → Social dipole moment

$$\mathbf{c} = \frac{1}{N}\sum_{i}^{N}\mathbf{p}_i$$

**Social dipole moment**

$$\mathbf{q} = \mathbf{s} - \frac{1}{N}\sum_{i}^{N}\mathbf{p}_i = \mathbf{s} - \mathbf{c}$$

**Electric dipole moment**

$$\mathbf{q}_e = \sum_{i}^{N}(\mathbf{e} - \mathbf{p}_i)$$

$104.5°$

Social dipole moment

$$\mathbf{q} = \mathbf{s} - \frac{1}{N}\sum_{i}^{N}\mathbf{p}_i = \mathbf{s} - \mathbf{c}$$

Electric dipole moment

$$\mathbf{q}_e = \sum_{i}^{N}(\mathbf{e} - \mathbf{p}_i)$$

Legend:
- ⊕ Social member
- ⊖ Joint attention
- ◑ Center of mass
- → Social dipole moment

Social dipole moment

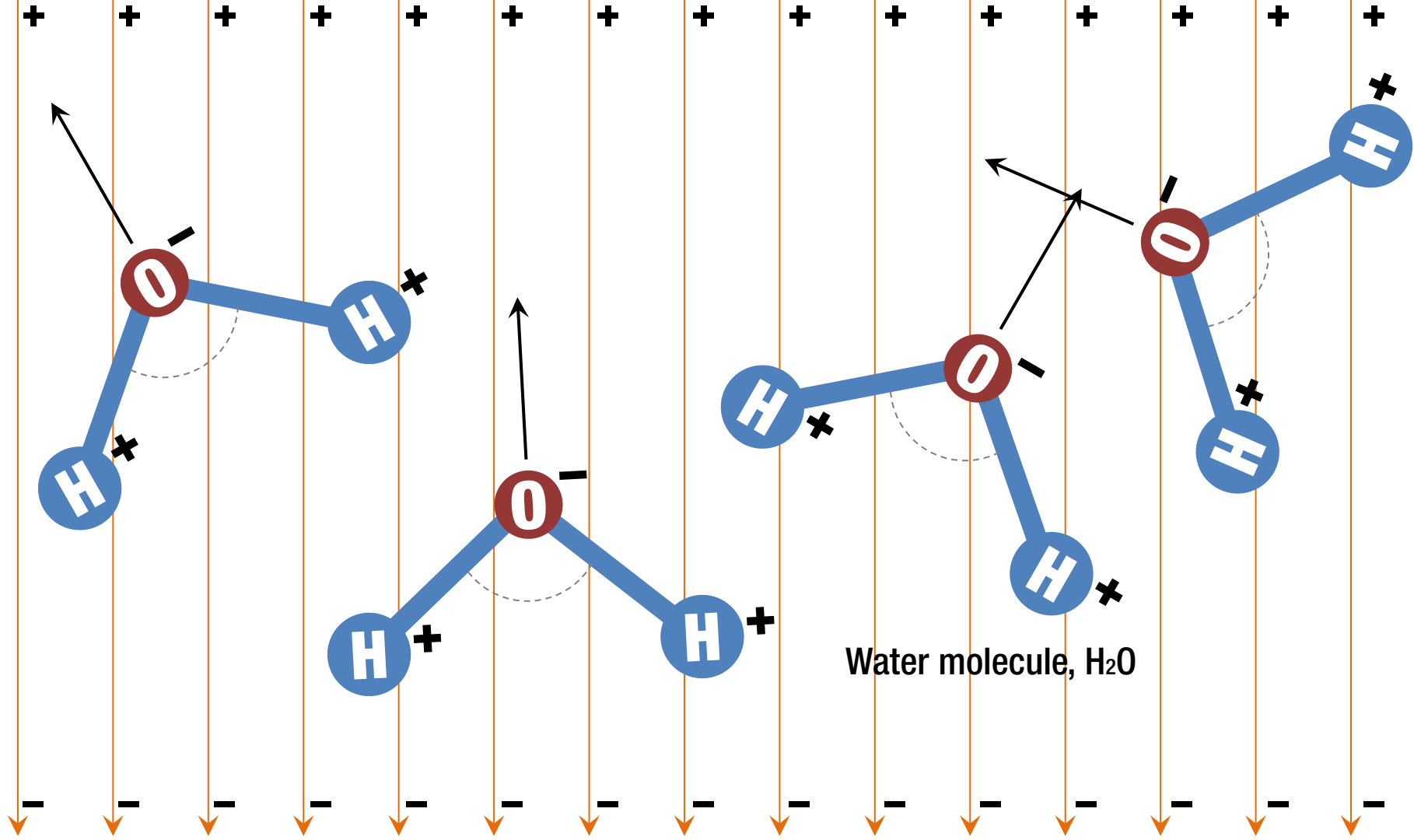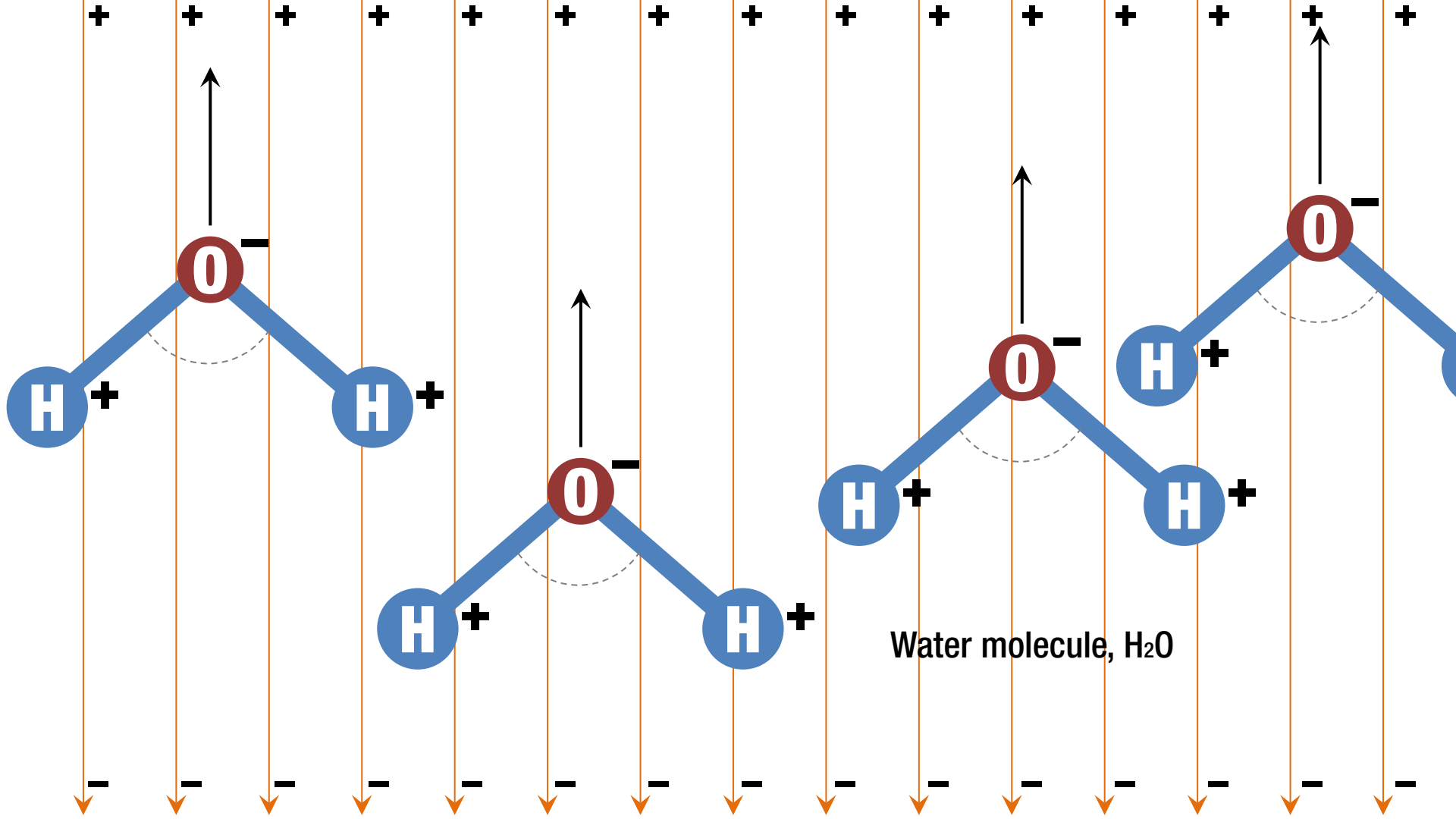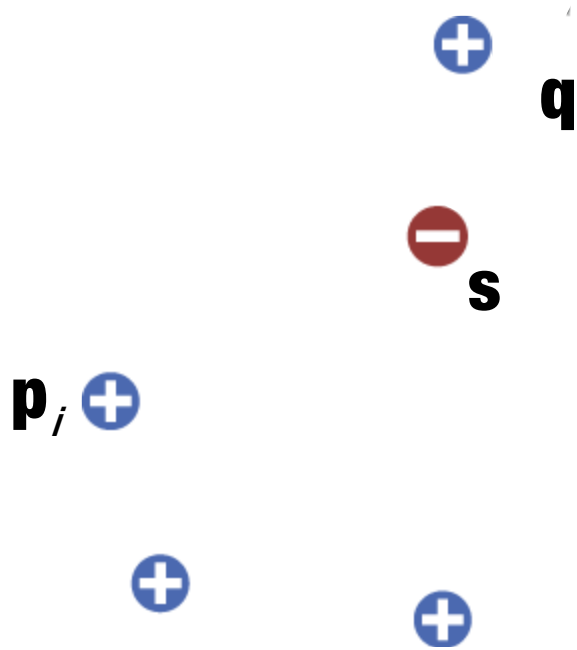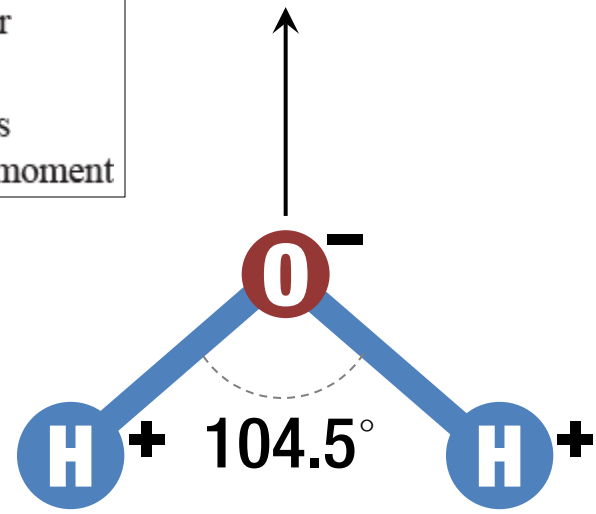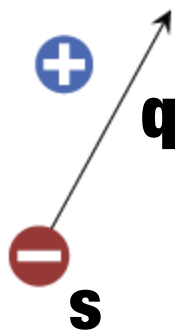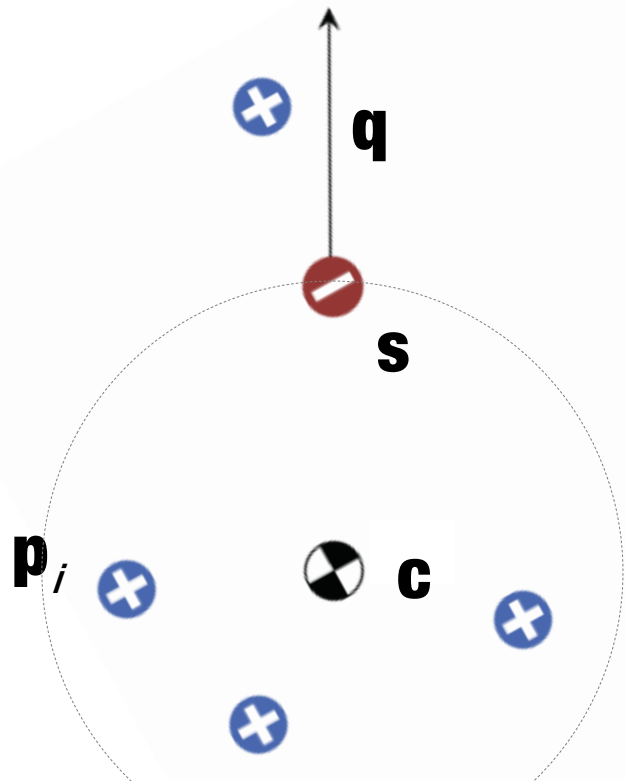$$\mathbf{q} = \mathbf{s} - \frac{1}{N}\sum_{i}^{N}\mathbf{p}_i = \mathbf{s} - \mathbf{c}$$
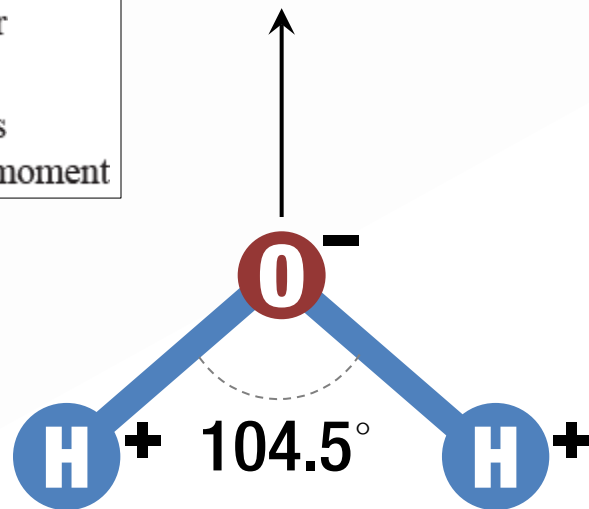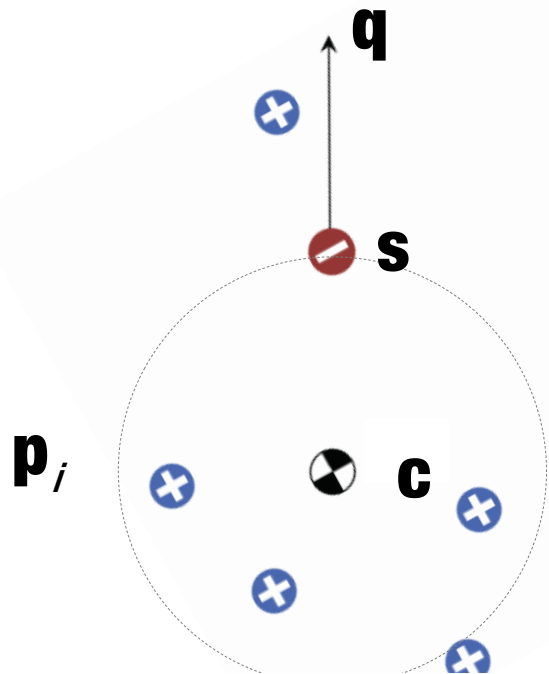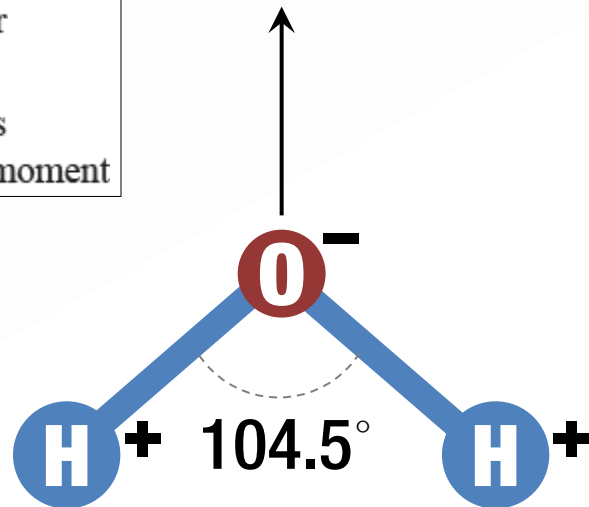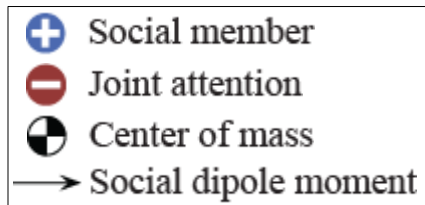
Electric dipole moment

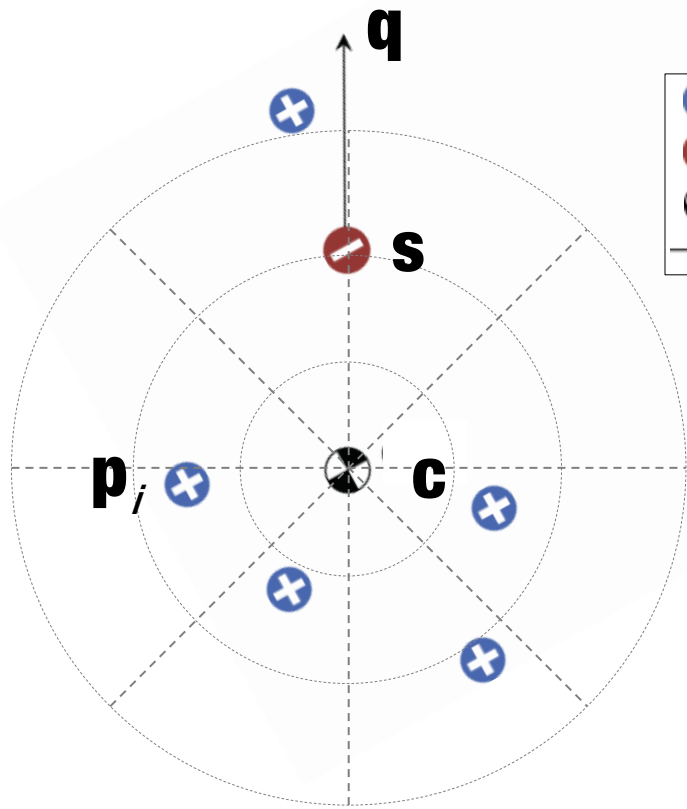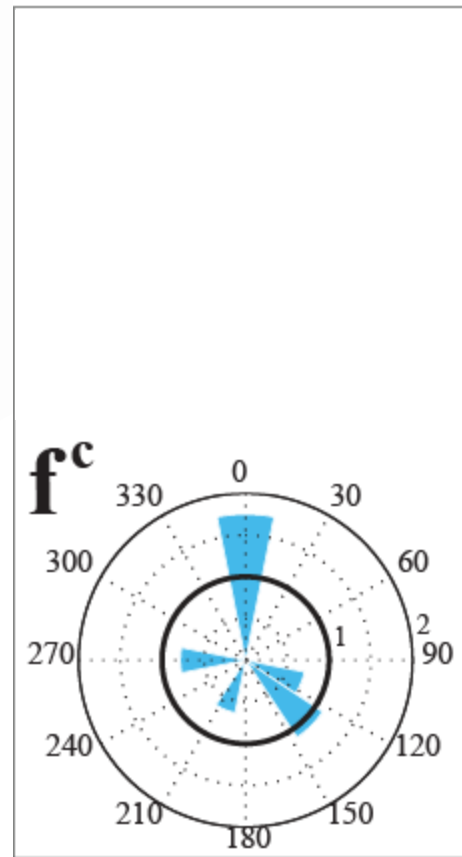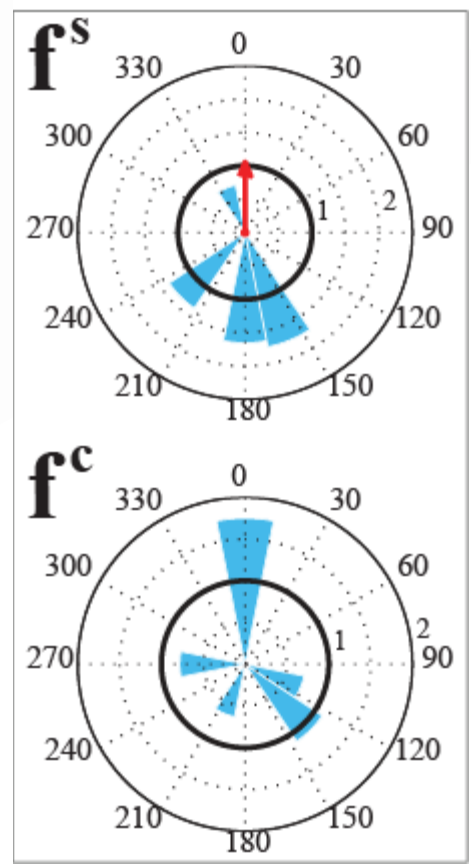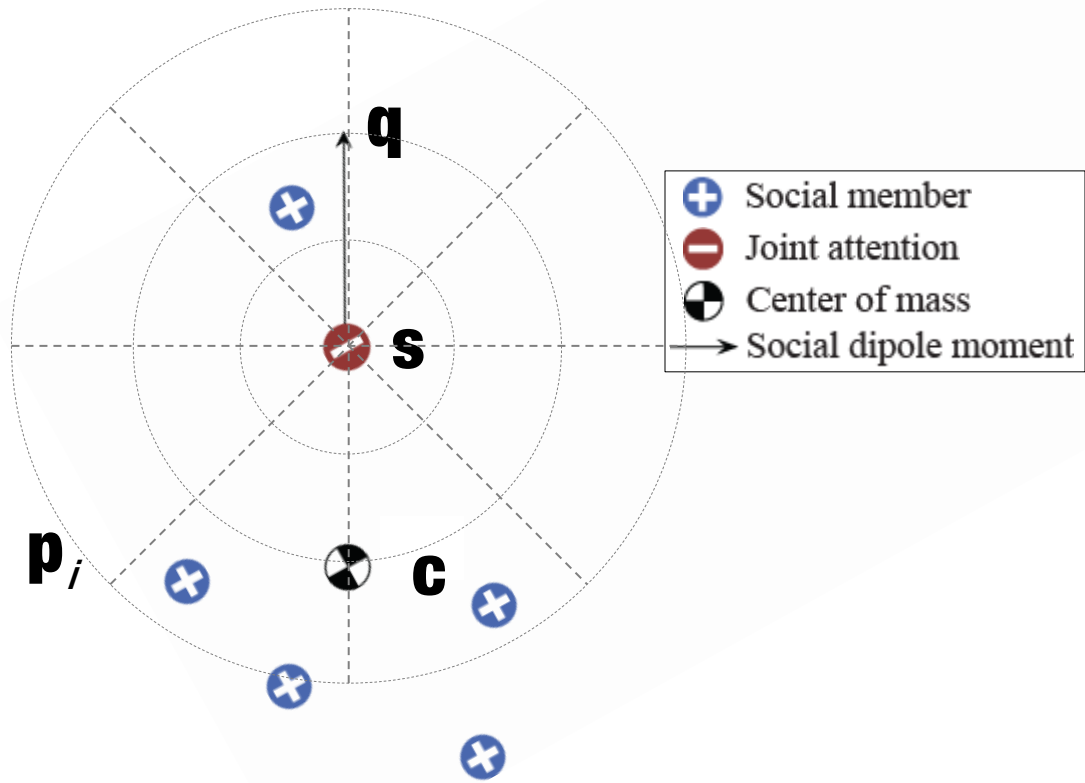$$\mathbf{q}_e = \sum_{i}^{N}(\mathbf{e} - \mathbf{p}_i)$$

Social member
Joint attention
Center of mass
Social dipole moment

$f^c$

Social formation feature
Social dipole moment

| Scene | N | T(sec) | F |
|---|---|---|---|
| B-boy I | 18 | 105 | 317 |
| B-boy II | 18 | 450 | 1351 |
| B-boy III | 18 | 160 | 528 |
| B-boy IV | 18 | 50 | 180 |
| Surprise party | 11 | 120 | 2227 |
| Class | 11 | 360 | 3590 |
| Croquet | 6 | 300 | 6000 |
| Busker I | 6 | 120 | 3566 |
| Busker II | 6 | 180 | 5394 |
| Card game | 3 | 180 | 768 |
| Hide and seek | 3 | 180 | 214 |
| Block building | 3 | 700 | 2702 |
| Social game | 8 | 450 | 2086 |
| Meeting I | 11 | 120 | 832 |
| Meeting II | 5 | 440 | 1120 |
| Picnic | 6 | 60 | 965 |
| Musical | 7 | 180 | 2184 |
| Dance | 6 | 180 | 5301 |
| 4 way party | 11 | 180 | 1909 |
| Snowman | 4 | 753 | 8256 |

*Total 49,490 social formations*

# Dyadic interaction



Center of mass with std.

Distribution of members



**Social formation theory (F-formation)***

Face-to-face

Side-to-side

Corner-to-corner

R - SPACE

P - SPACE

O - SPACE

*A. Kendon, "Conducting Interaction: Patterns of Behavior in Focused Encounters", Cambridge University Press, 1990.

Dyadic interaction

# Social Group Detection

q

s

x

$p_i$

c

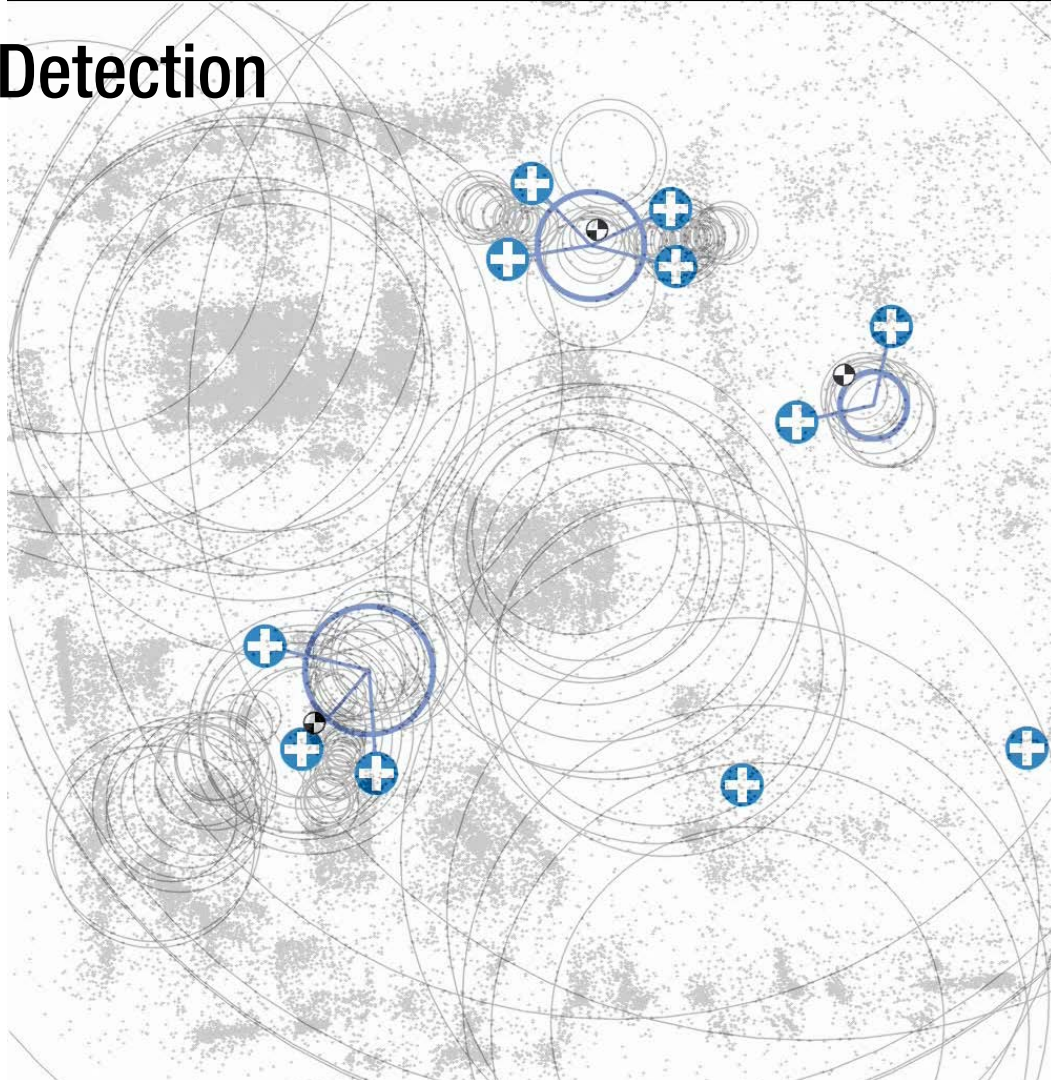| | | |
|---|---|---|
| ⊕ | Social member | |
| ⊖ | Joint attention | |
| ◑ | Center of mass | |
| → | Social dipole moment | |

$$\Phi\left(\mathbf{f^c}, \mathbf{f^s};\right) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{s} \end{cases}$$

$\mathbf{f^s}$

$\mathbf{f^c}$

| | |
|---|---|
| ◀ | Social formation feature |
| → | Social dipole moment |

$$\Phi\left(\mathbf{f^c}, \mathbf{f^s};\right) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{s} \\ 0 & \text{if } \mathbf{x} \neq \mathbf{s} \end{cases}$$

Social member
Joint attention
Center of mass

**x**

$\mathbf{p}_i$

**c**

$\mathbf{f^s}$

$\mathbf{f^c}$

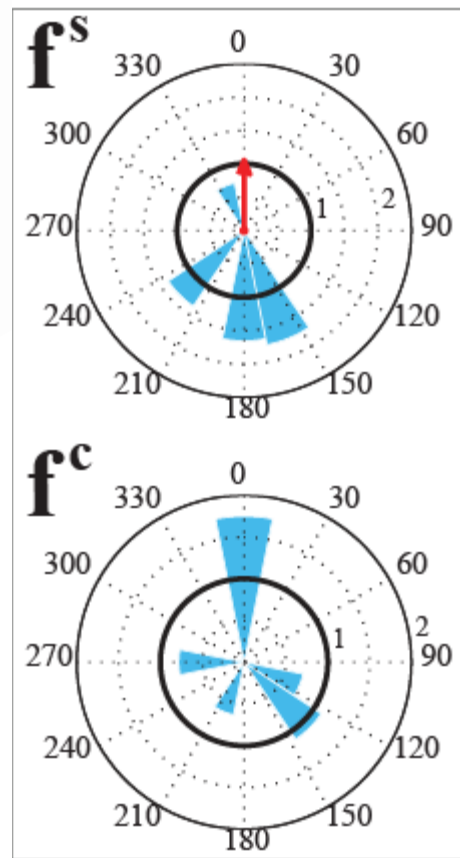Social formation feature
Social dipole moment

$$\Phi\left(\mathbf{f^c}, \mathbf{f^s};\right) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{s} \\ 0 & \text{if } \mathbf{x} \neq \mathbf{s} \end{cases}$$

**Legend:**
- Social member
- Joint attention
- Center of mass

$\mathbf{x}$

$\mathbf{f^s}$

$\mathbf{f^c}$

$\mathbf{p}_i$

$\mathbf{c}$

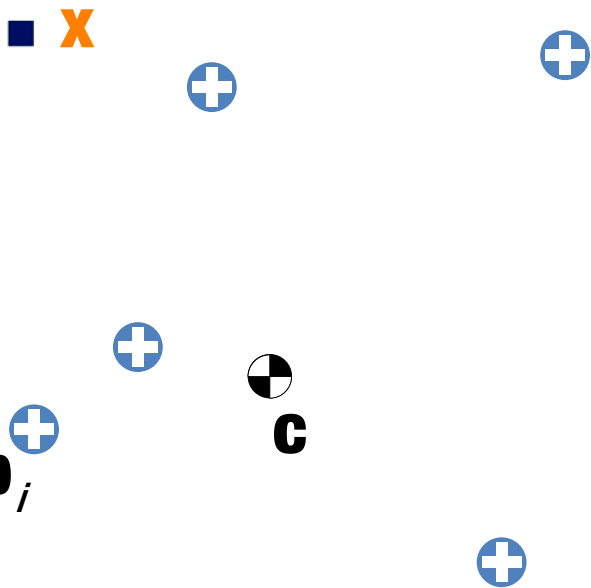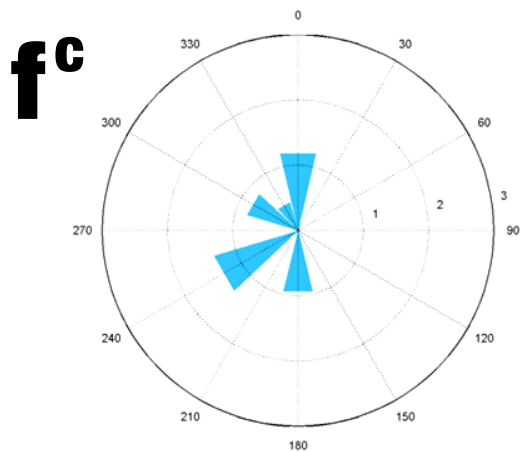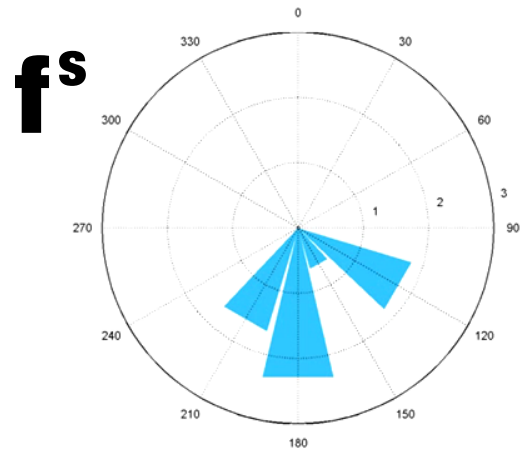Social saliency:
Likelihood of joint attention

$$\Phi\left(\mathbf{f^c},\mathbf{f^s};\right) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{s} \\ 0 & \text{if } \mathbf{x} \neq \mathbf{s} \end{cases}$$

- Social formation feature
- Social dipole moment

Social saliency:
Likelihood of joint attention

$$\Phi\left(\mathbf{f^c}, \mathbf{f^s};\right) = \begin{cases} 1 & \text{if } \mathbf{x} = \mathbf{s} \\ 0 & \text{if } \mathbf{x} \neq \mathbf{s} \end{cases}$$
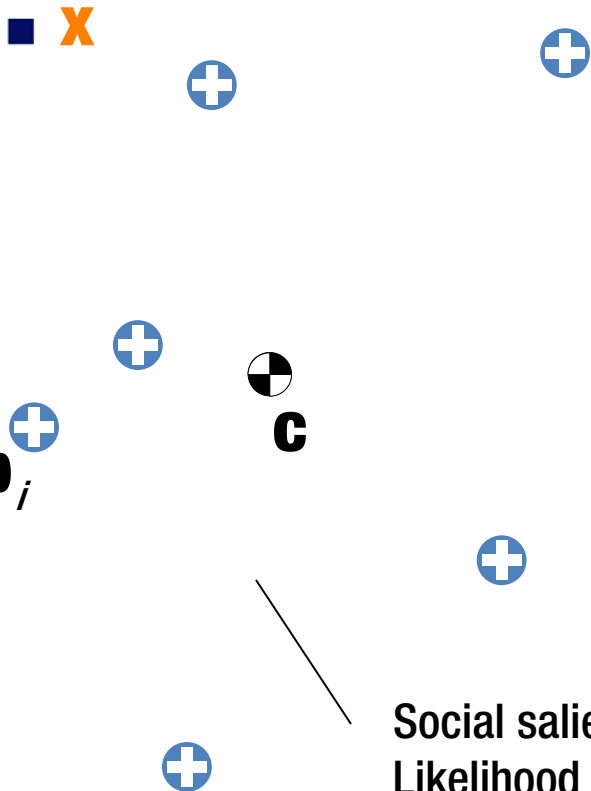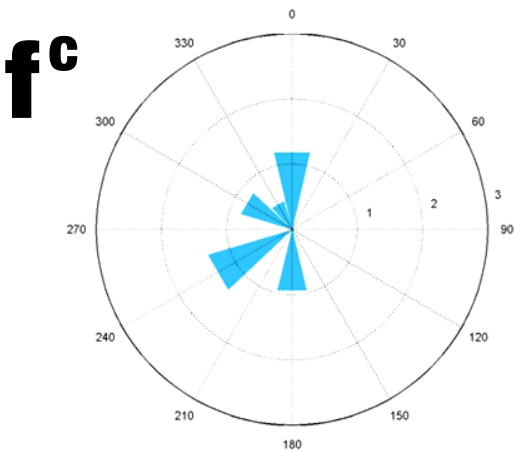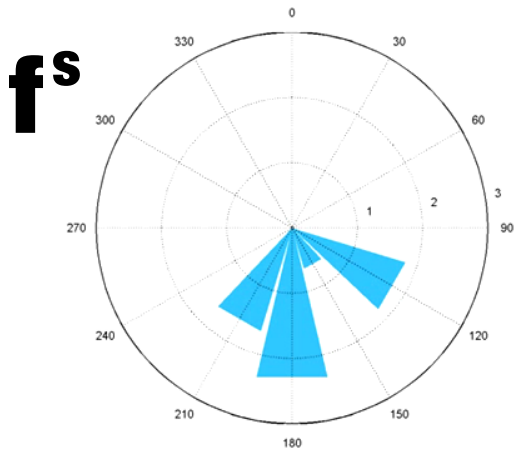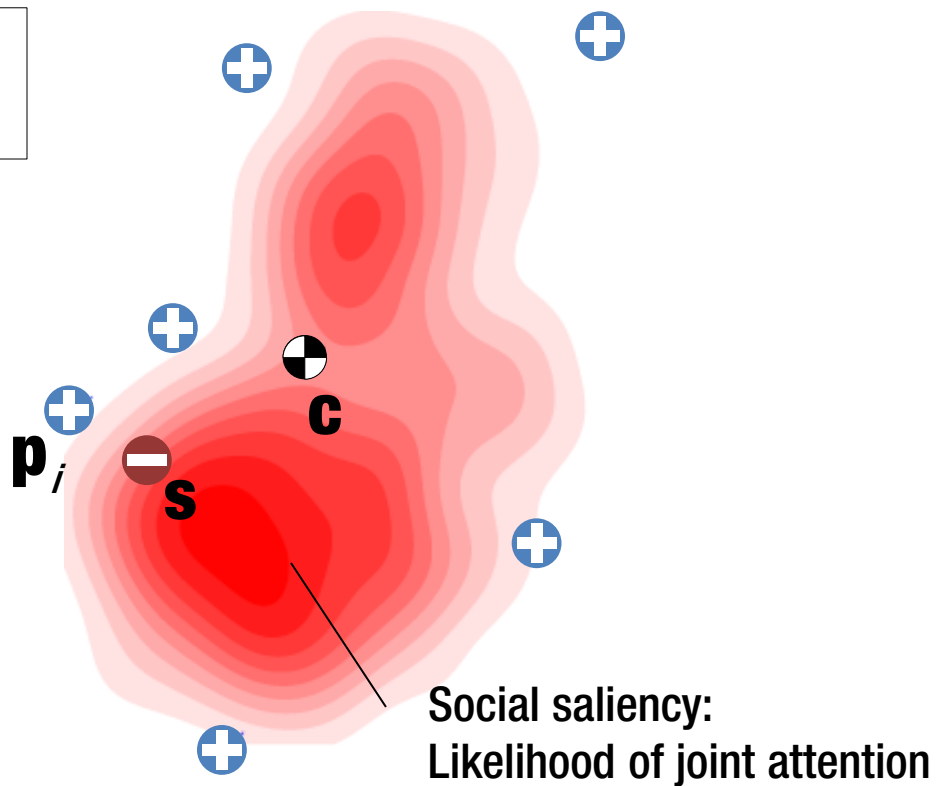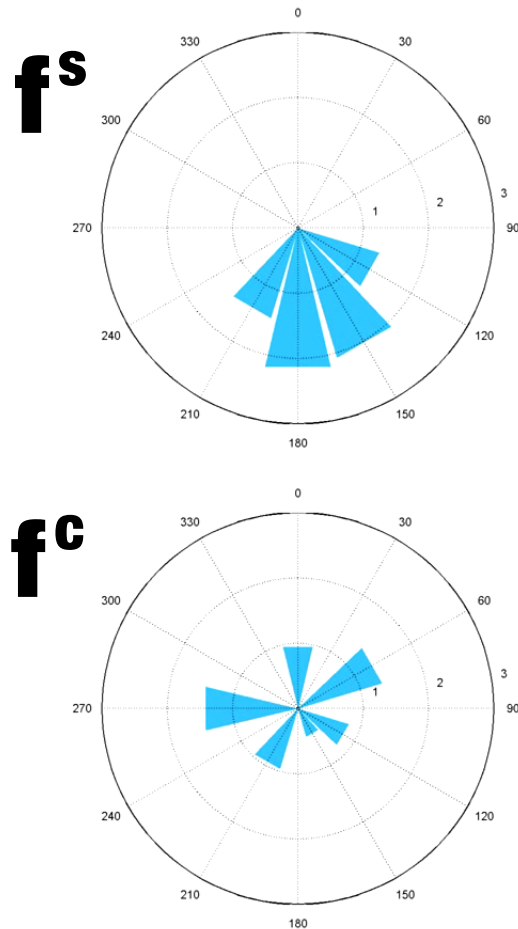
Social member location
Joint attention location

Social member location
Joint attention location
Center of mass (COM)

Social member location
Joint attention location
Center of mass (COM)
Center of circumcircle (CC)

**Group meeting**

**Street performance**

**Class interactions**

CF+Boosting

⊕ Social members
⊖ Ground truth joint attention
◑ Center of mass
⊗ Center of circumcircle
■ Social saliency

CF: Context feature (Lan et al., PAMI 2012)

⊕ Social member location
⊖ Joint attention location
◑ Center of mass (COM)
⊗ Center of circumcircle (CC)

# CF: Context feature (Lan et al., PAMI 2012)



⊕ Social member location
⊖ Joint attention location
◐ Center of mass (COM)
⊗ Center of circumcircle (CC)

# Mean average precision

| Scenes | SFF+Boosting | SFF+RF | CC | COM | CF |
|---|---|---|---|---|---|
| Dance | 0.2769 | 0.1381 | 0.3299 | 0.0419 | 0.0106 |
| Meeting I | 0.2941 | 0.3599 | 0.2418 | 0.2350 | 0.0649 |
| B-boy I | 0.7178 | 0.6907 | 0.2078 | 0.1232 | 0.1225 |
| Class | 0.7678 | 0.7386 | 0.1445 | 0.2757 | 0.1873 |
| Busker | 0.2919 | 0.2059 | 0.3432 | 0.1929 | 0.0103 |
| Picnic | 0.1364 | 0.1349 | 0.1115 | 0.1808 | 0.0244 |
| Social game | 0.5425 | 0.4419 | 0.3461 | 0.2463 | 0.0020 |

Group meeting / Street performance / Class interactions ROC curves

Legend:
- SFF+Boosting
- SFF+RF
- CC predictor
- COM predictor
- CF+Boosting

- ⊕ Social members
- ⊖ Ground truth joint attention
- Center of mass
- Center of circumcircle
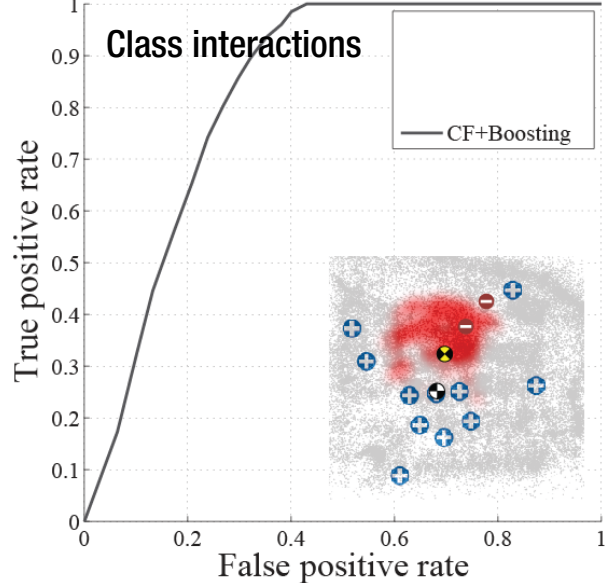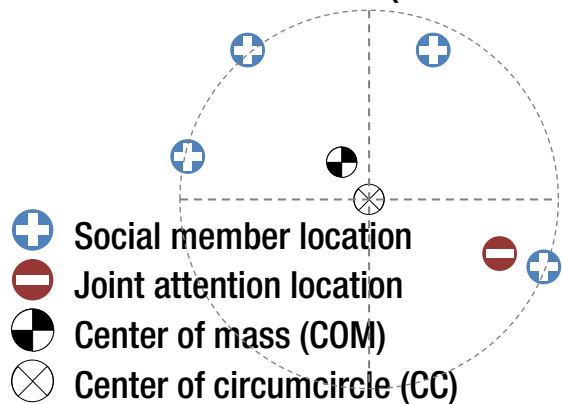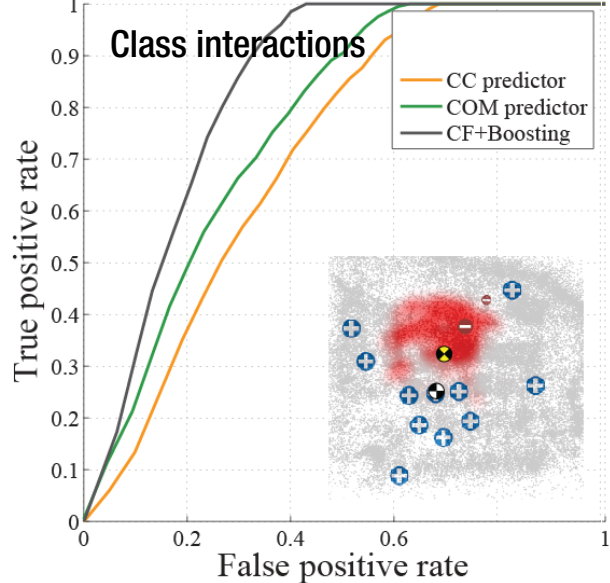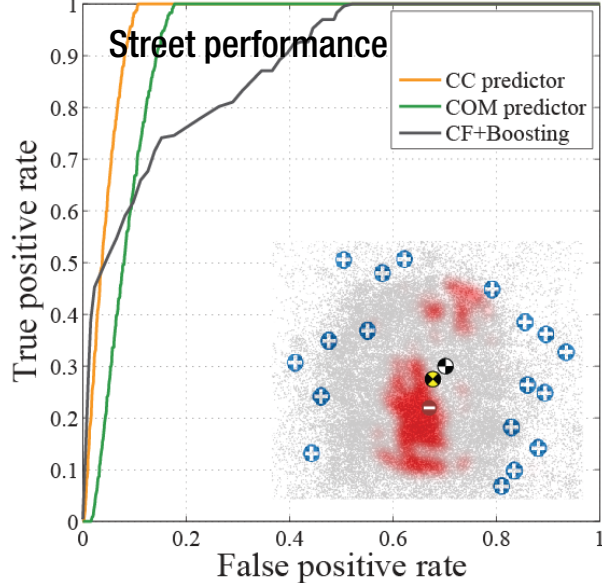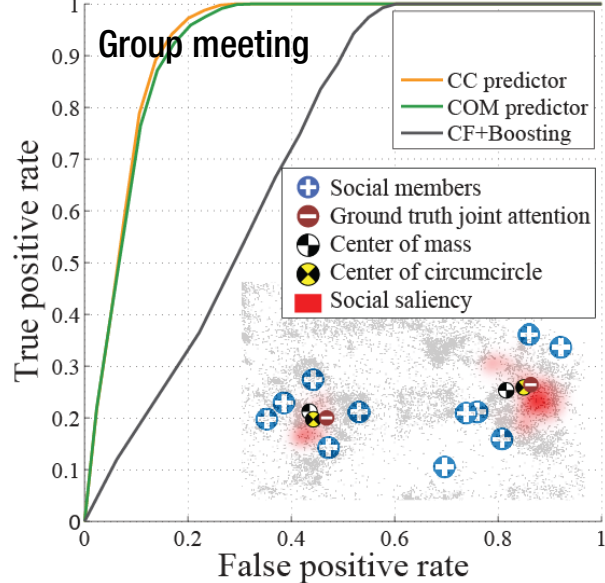- Social saliency

**CF: Context feature (Lan et al., PAMI 2012)**

- ⊕ Social member location
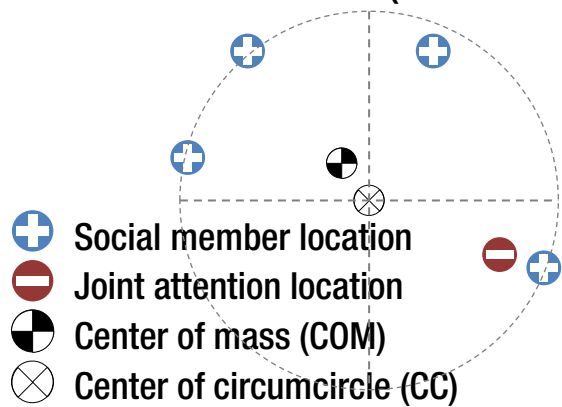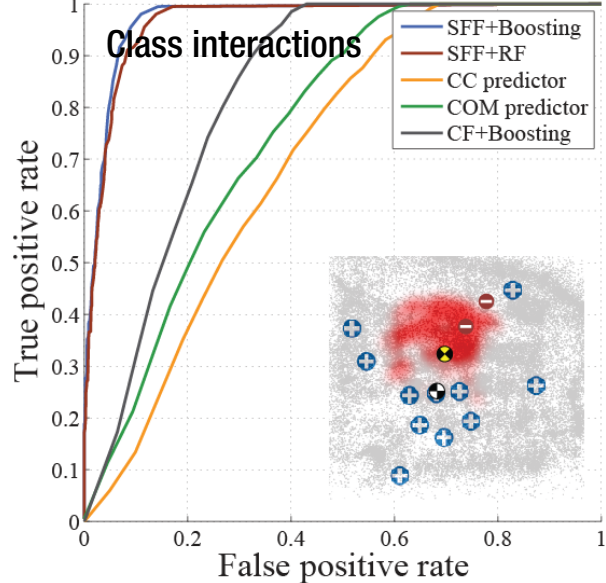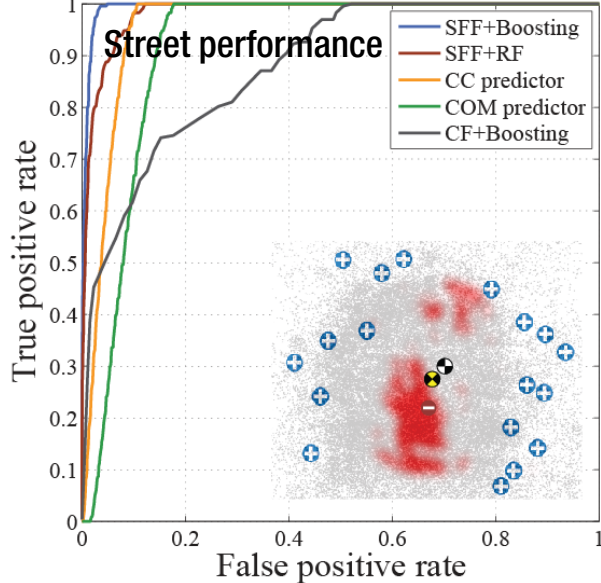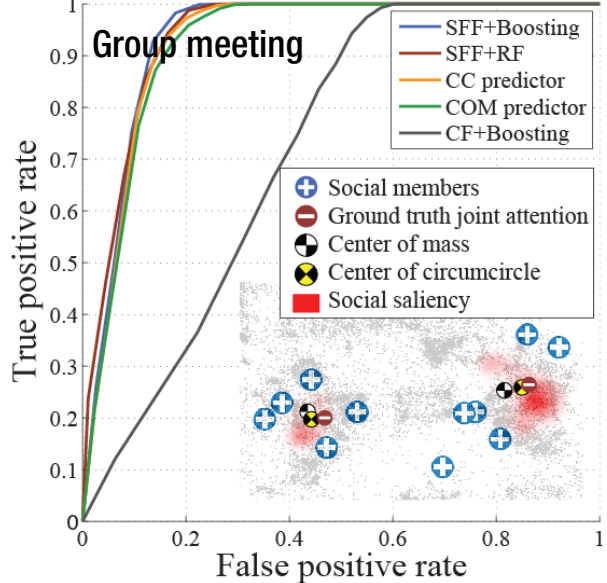- ⊖ Joint attention location
- Center of mass (COM)
- ⊗ Center of circumcircle (CC)

## Mean average precision

| Scenes | SFF+Boosting | SFF+RF | CC | COM | CF |
|---|---|---|---|---|---|
| Dance | 0.2769 | 0.1381 | 0.3299 | 0.0419 | 0.0106 |
| Meeting I | 0.2941 | 0.3599 | 0.2418 | 0.2350 | 0.0649 |
| B-boy I | 0.7178 | 0.6907 | 0.2078 | 0.1232 | 0.1225 |
| Class | 0.7678 | 0.7386 | 0.1445 | 0.2757 | 0.1873 |
| Busker | 0.2919 | 0.2059 | 0.3432 | 0.1929 | 0.0103 |
| Picnic | 0.1364 | 0.1349 | 0.1115 | 0.1808 | 0.0244 |
| Social game | 0.5425 | 0.4419 | 0.3461 | 0.2463 | 0.0020 |

Social saliency

Social saliency

Camera

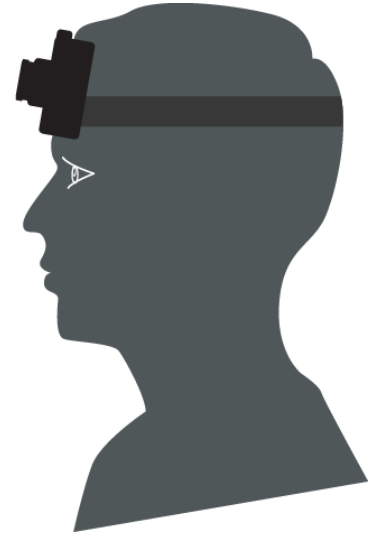Top view

Time Square

Top view

Source: https://www.youtube.com/watch?v=ezyrSKgcyJw
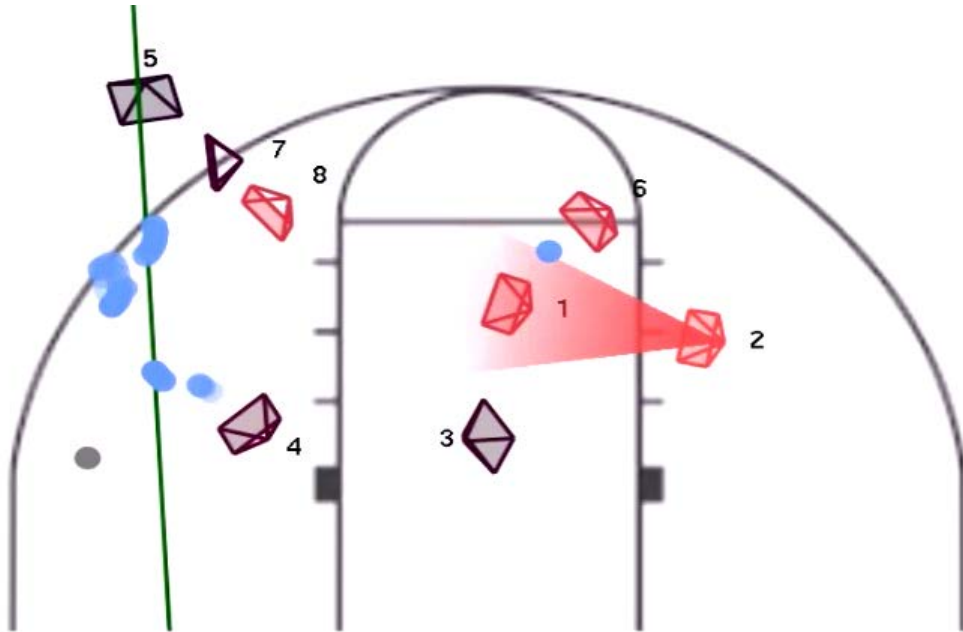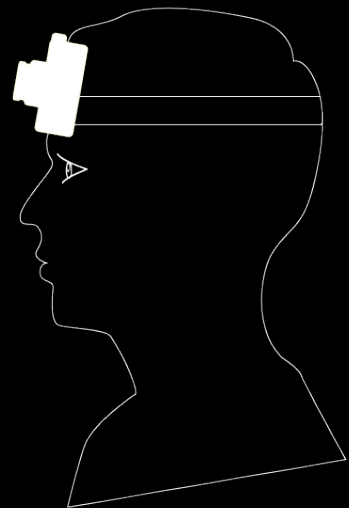
# How would this joint attention be useful for social event?

# Problems of videos taken by social cameras:

- Produce too much information to digest at once
- Are biased by an intimate and personal view
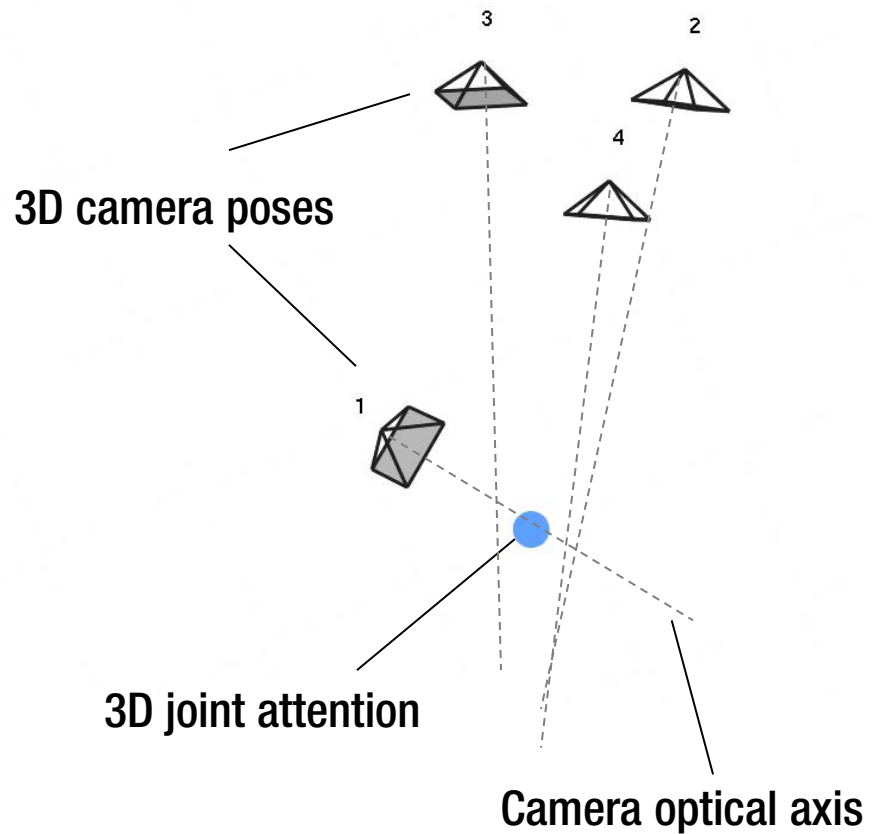
Input: Synchronized Social Footage

Output: Edited Video
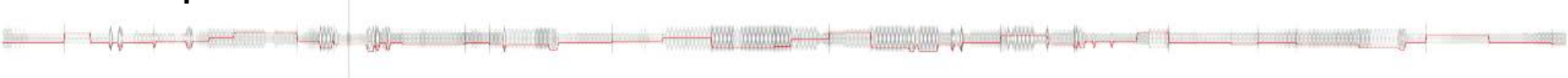
How to define the content of the social event?

3D camera poses

3D joint attention

Camera optical axis

Reprojection of joint attention

Joint attention ≈ Content of the event

**Input video feeds**

**Timeline**

Ours

Selected camera: 2

**Video feeds**        **3D geometry**        **Output video**
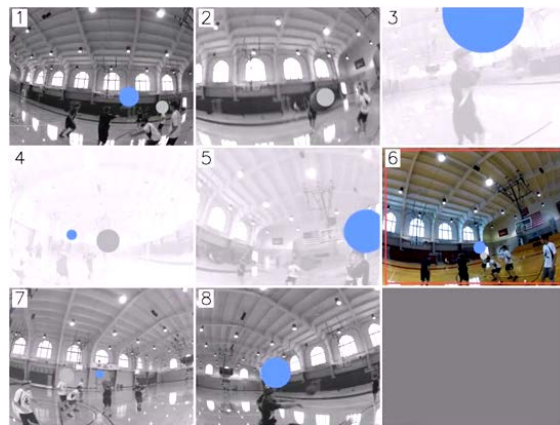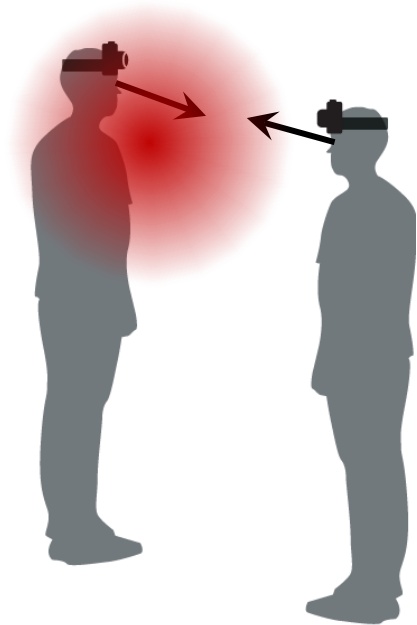
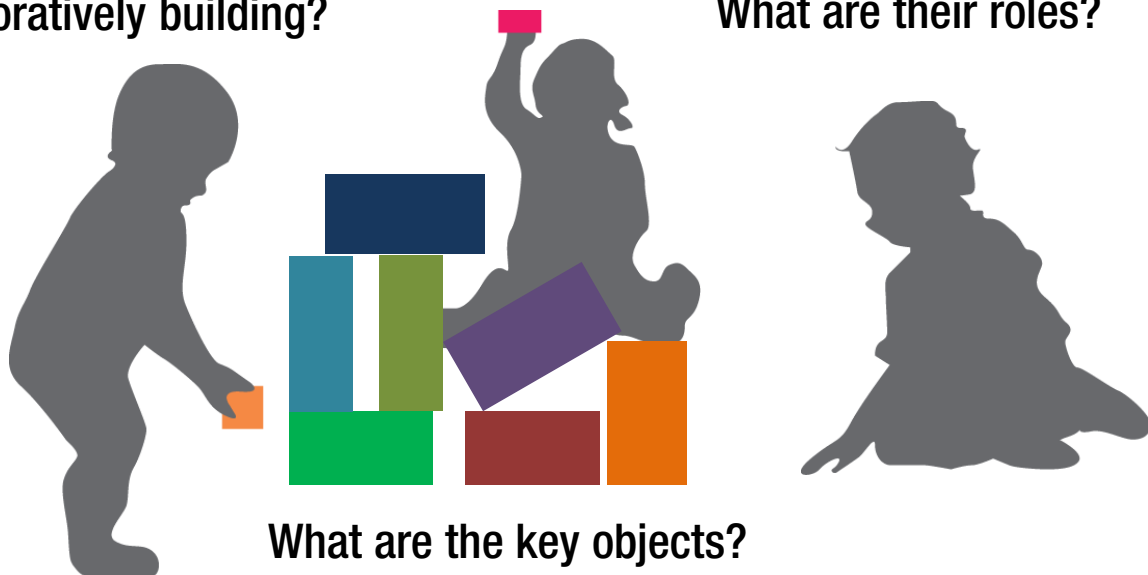Our method                    Professional Editor

# First person cameras are ideal sensors to measure social behaviors.

What are they collaboratively building?

What are their roles?

What are the key objects?

# Social Attention
## :What can first person cameras tell us about our social interactions?