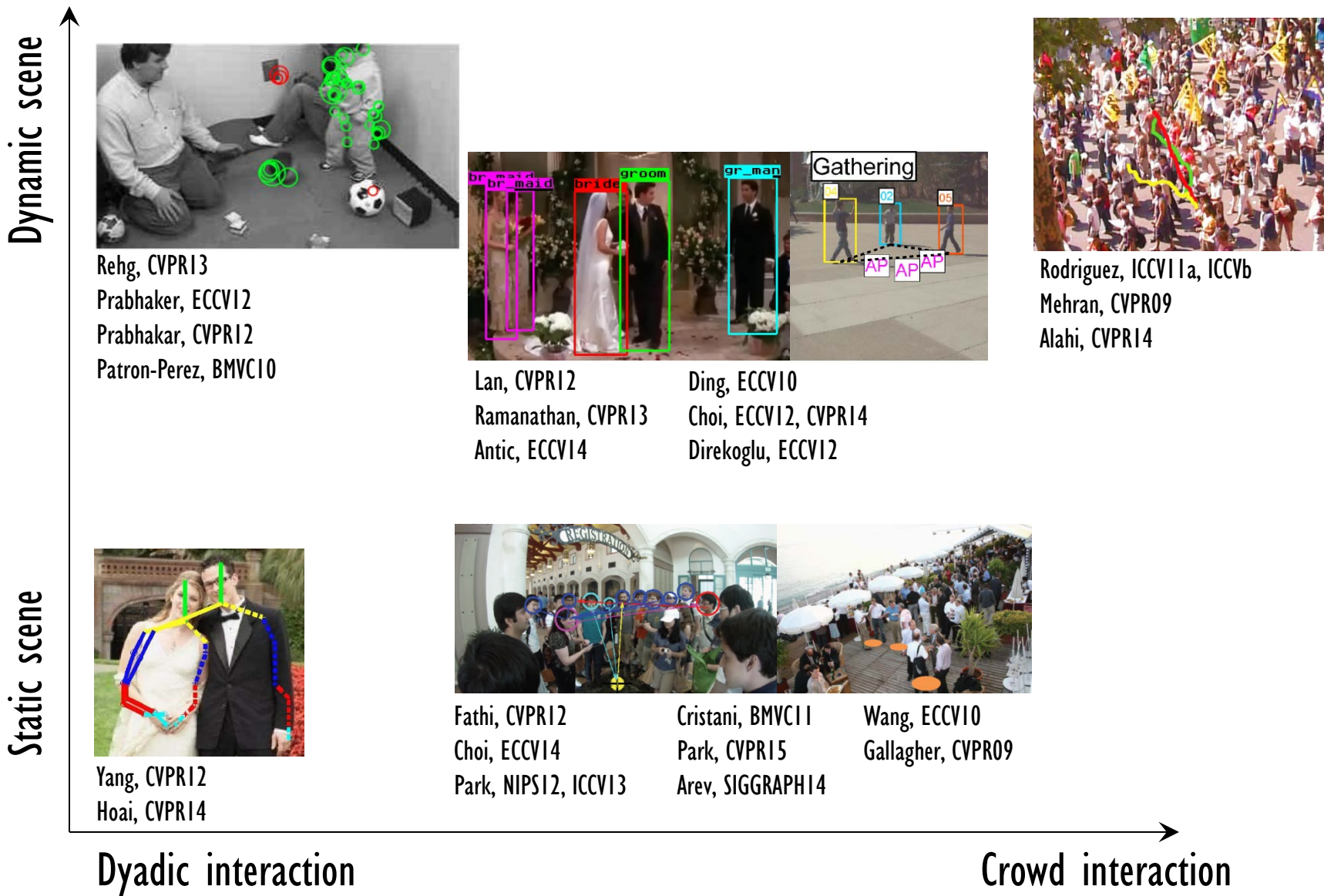


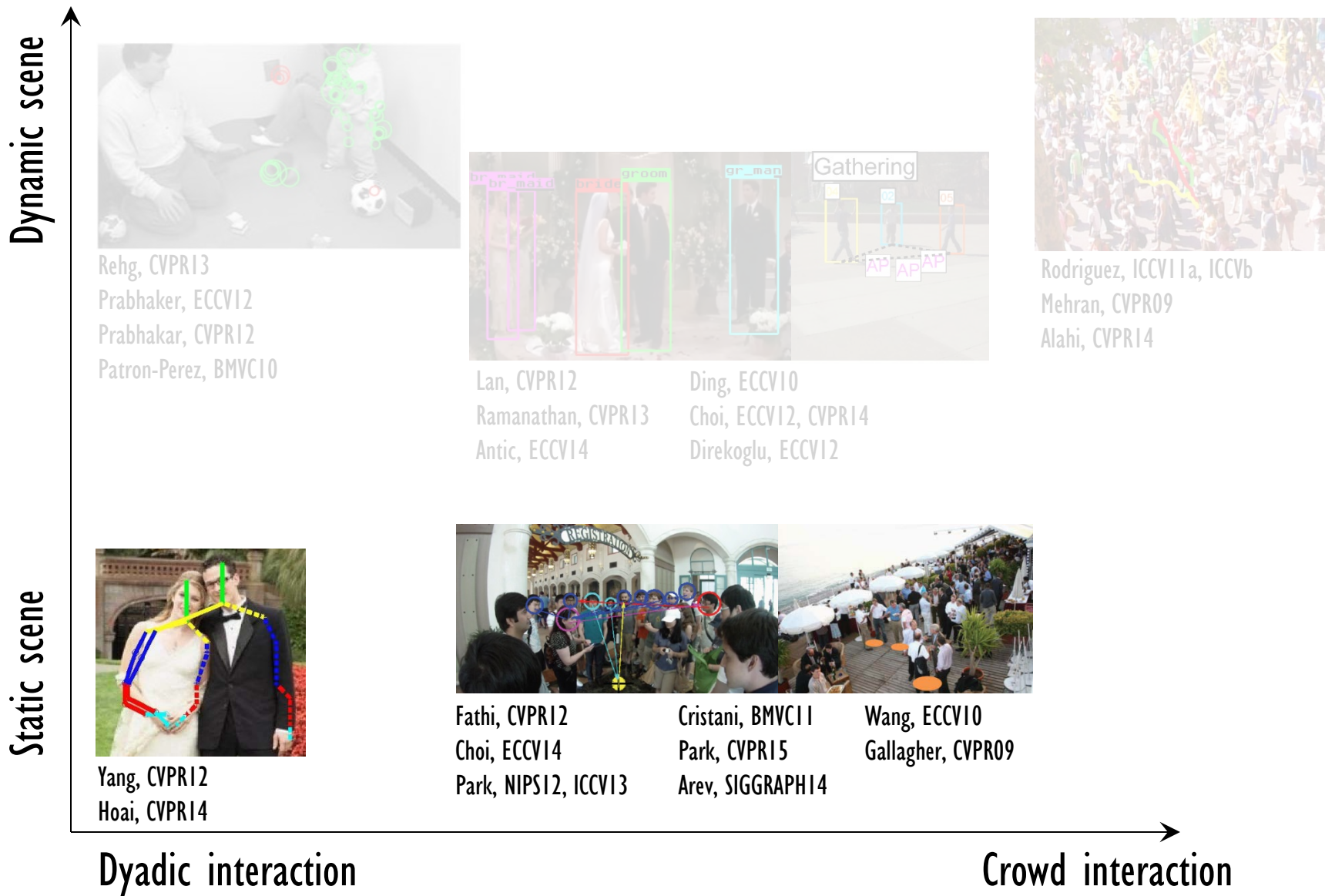
Social Statics



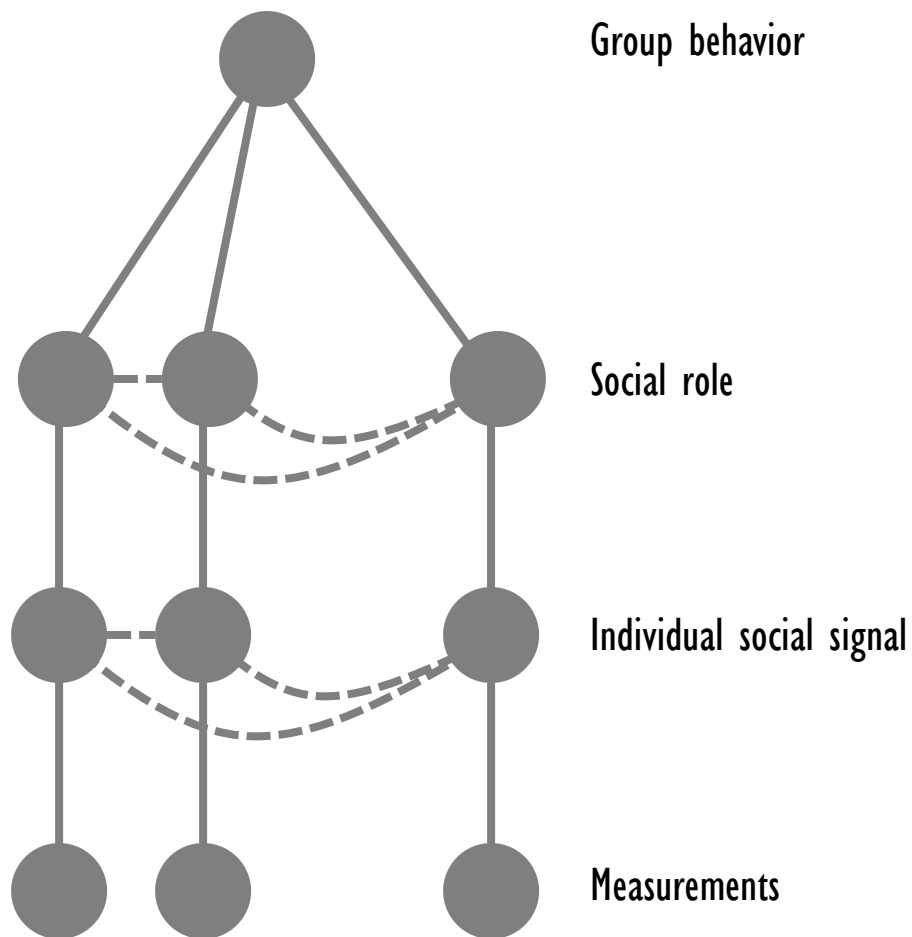
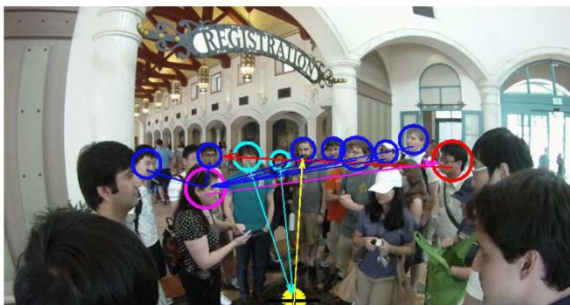
Scene dynamism

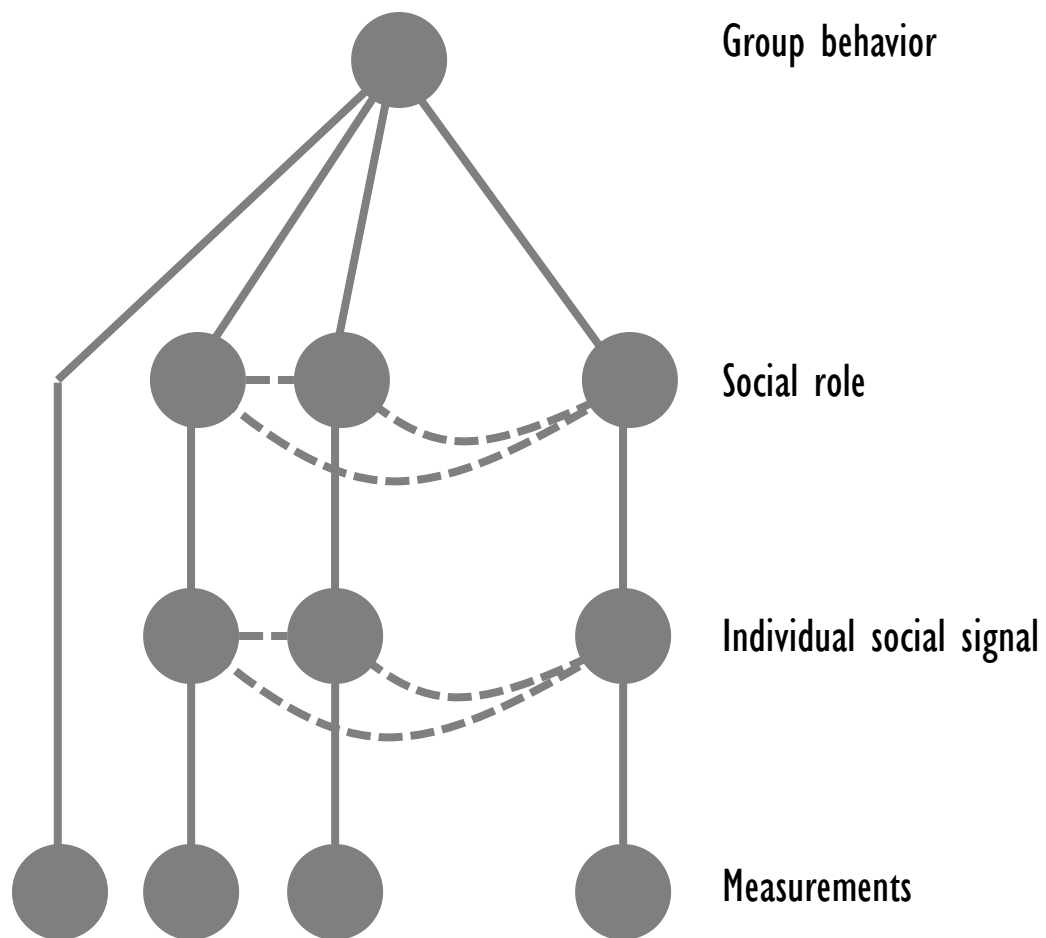


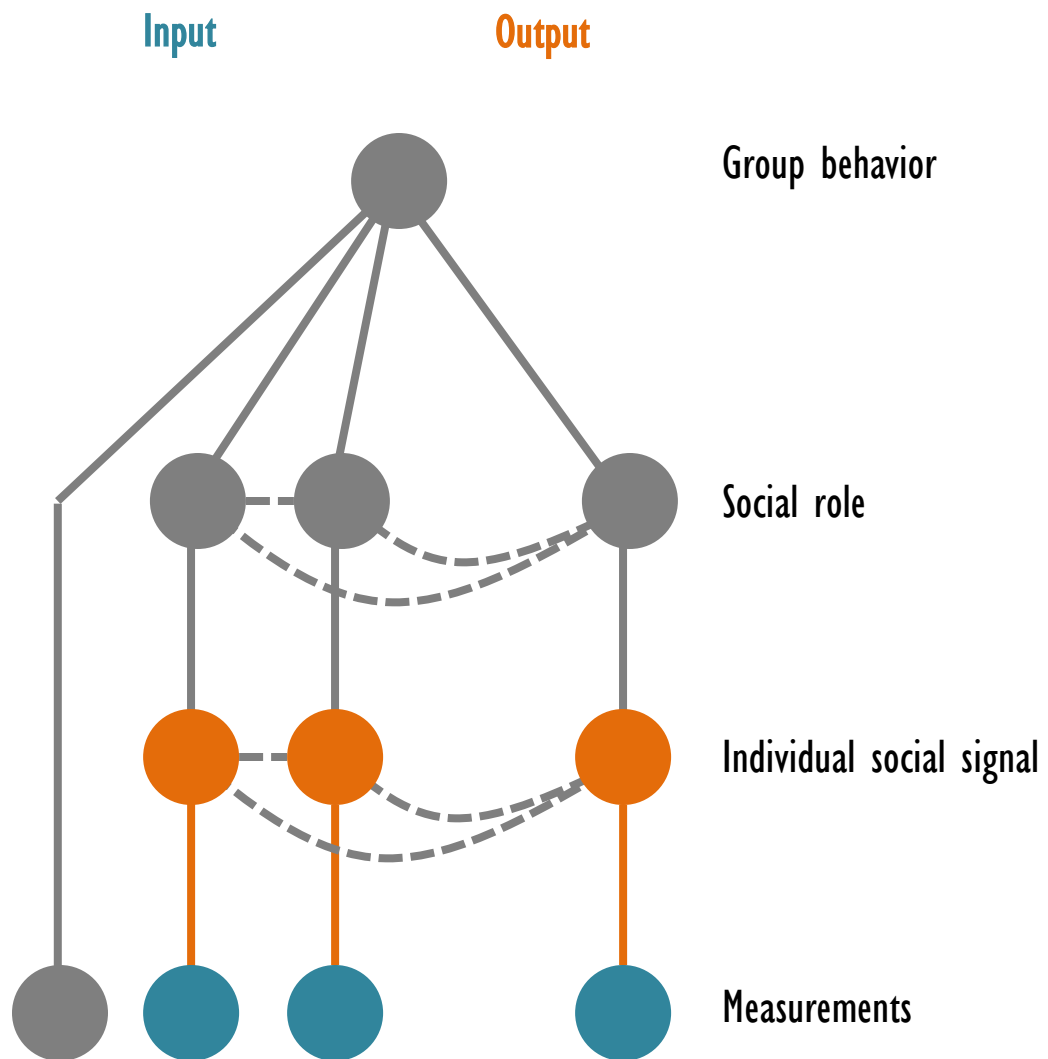
Scene dynamism

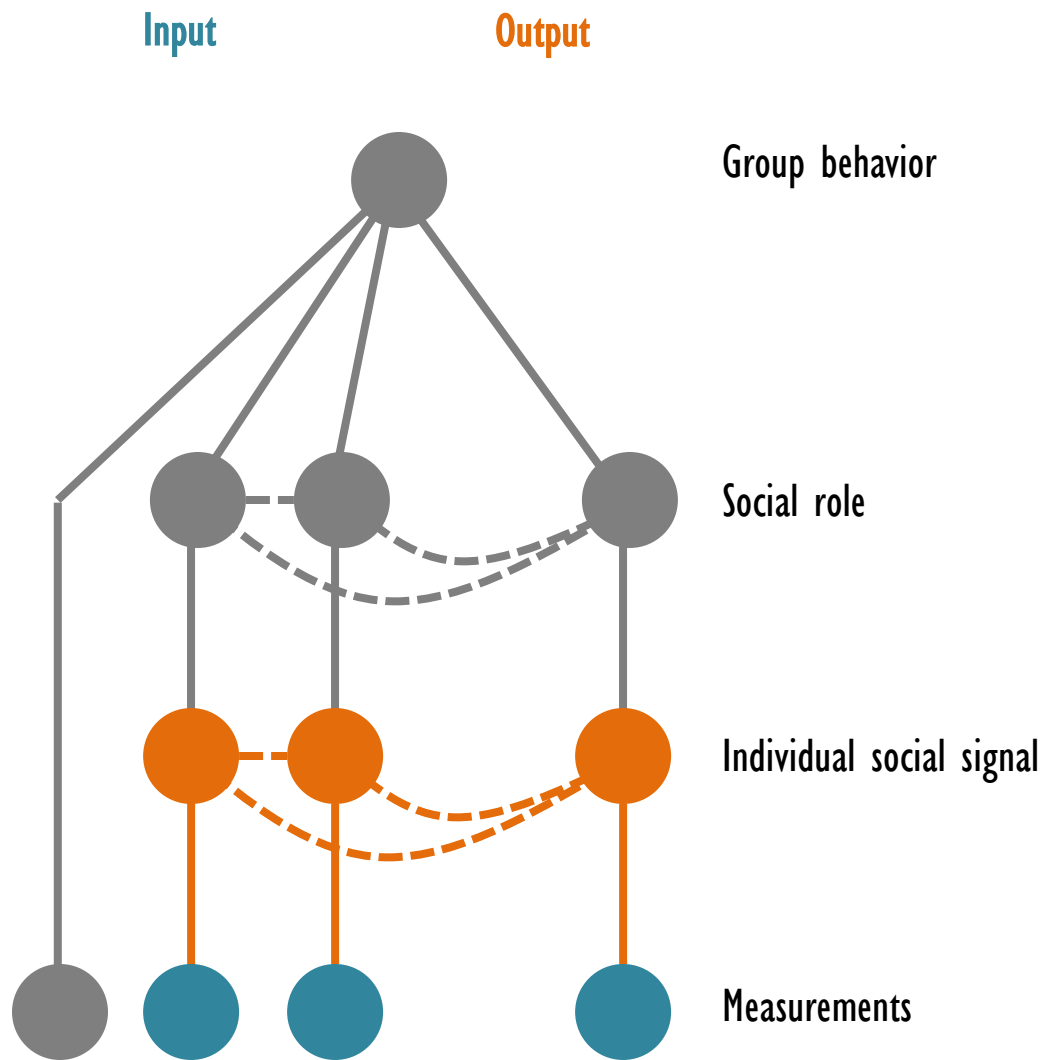


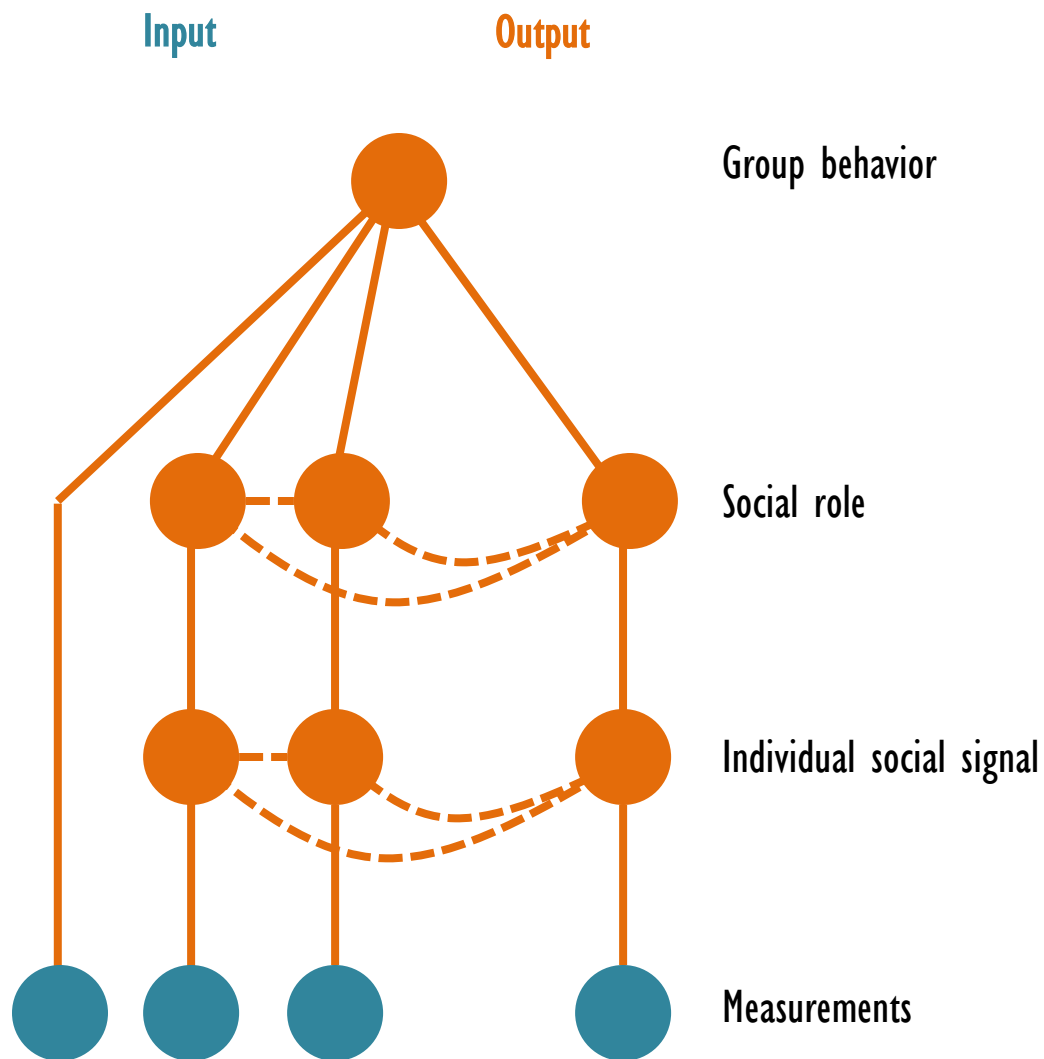
Number of group members

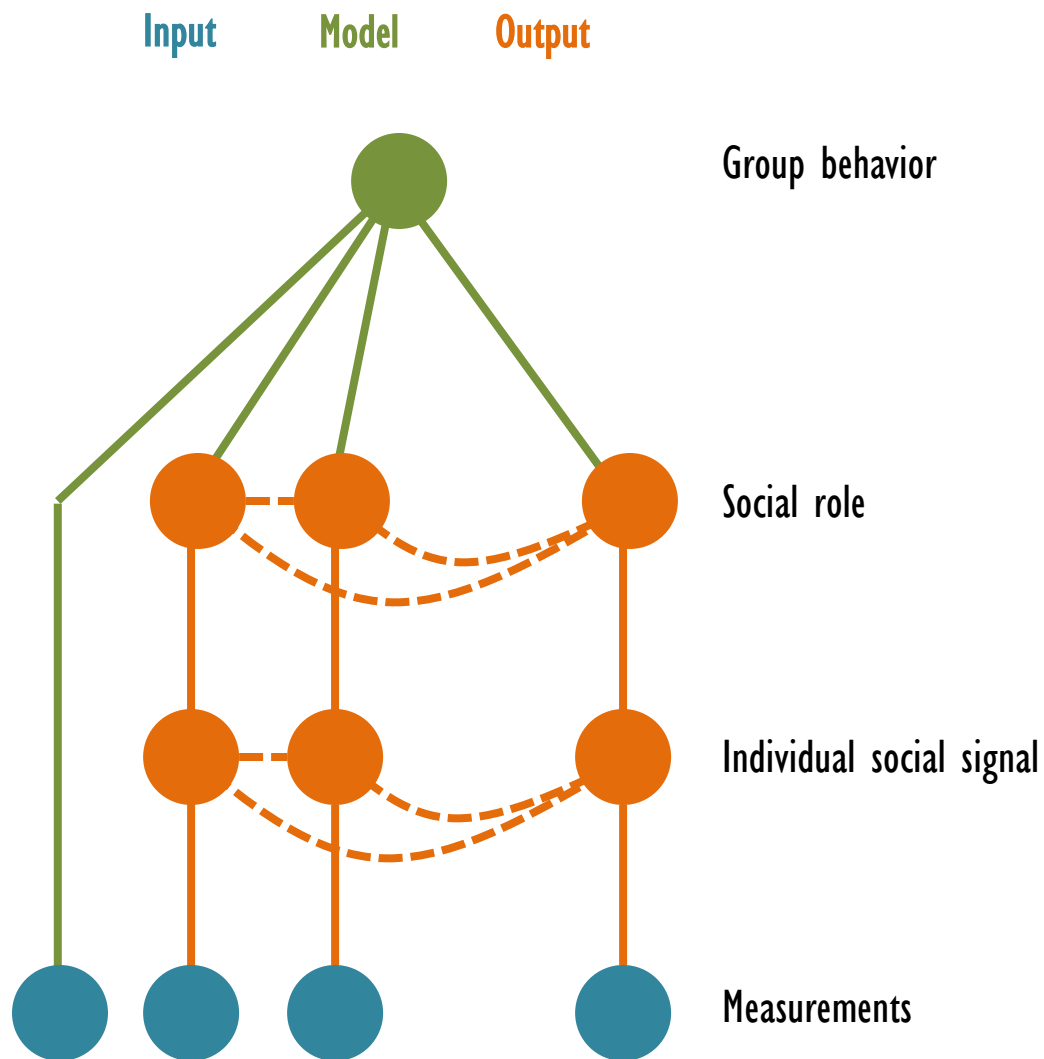


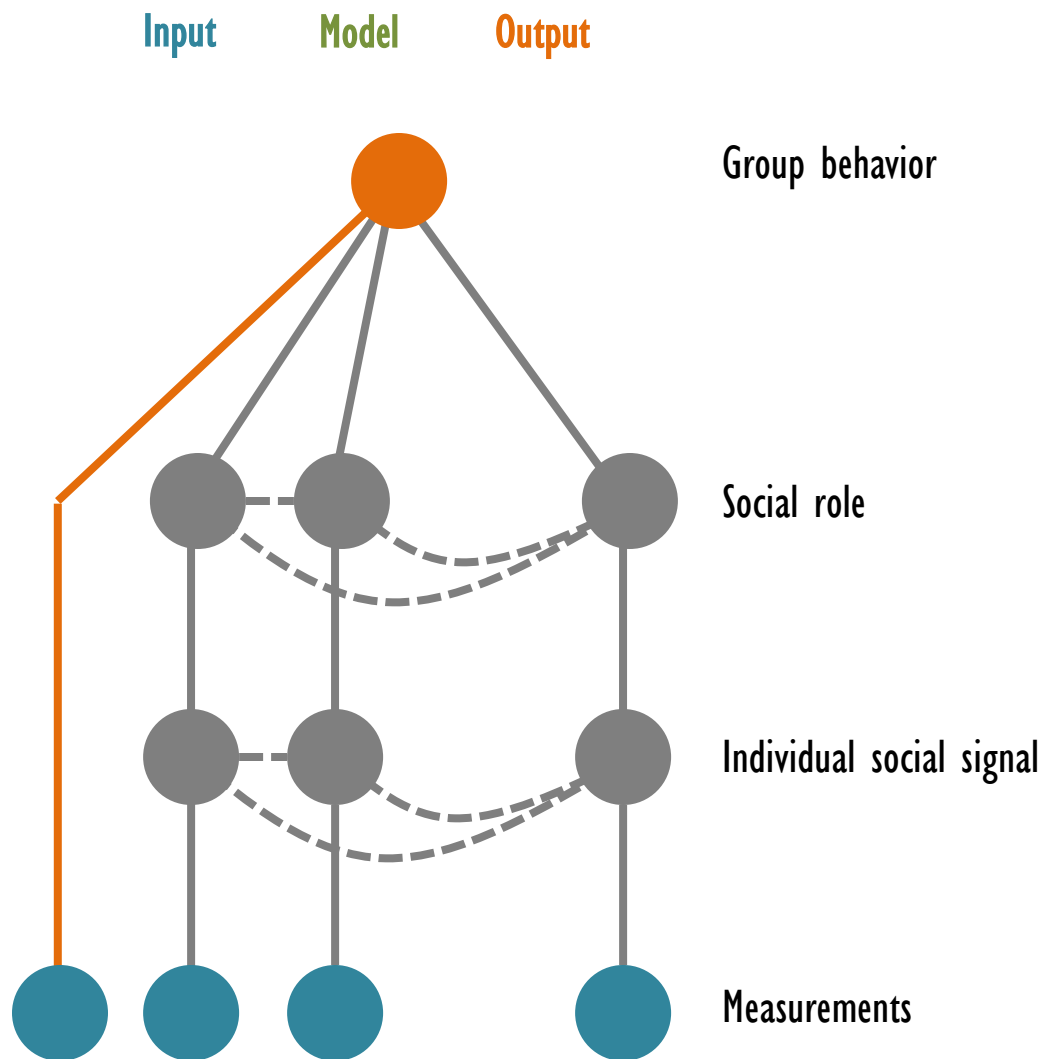


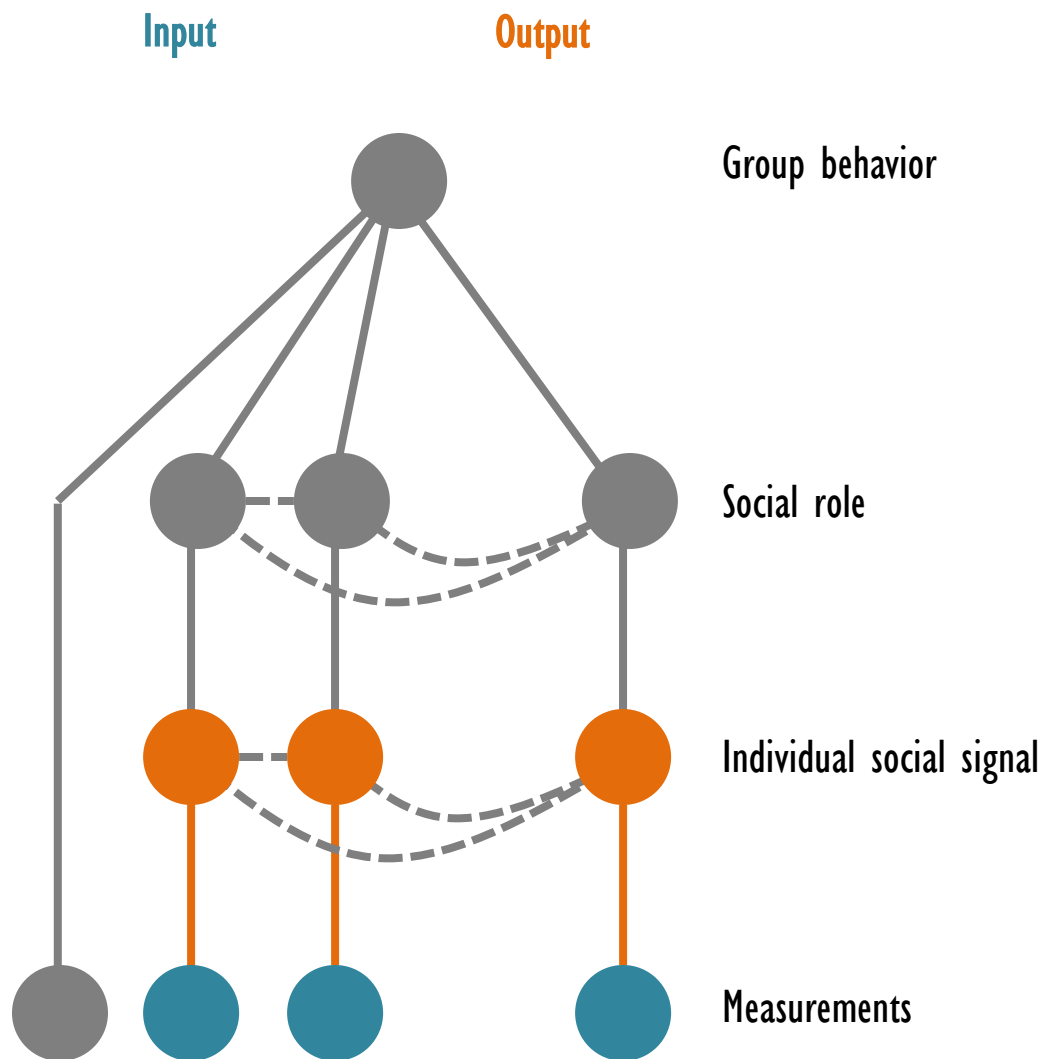


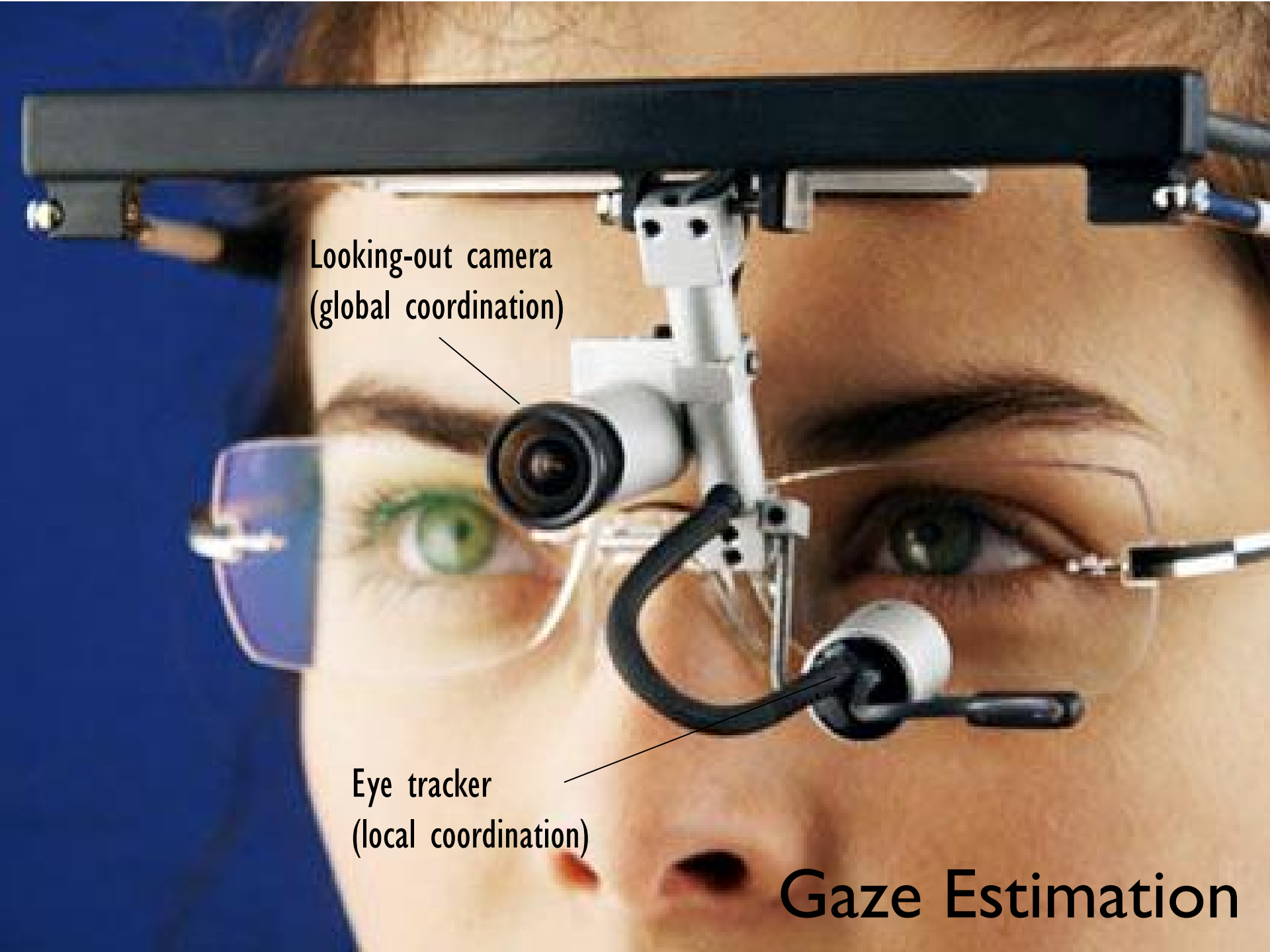










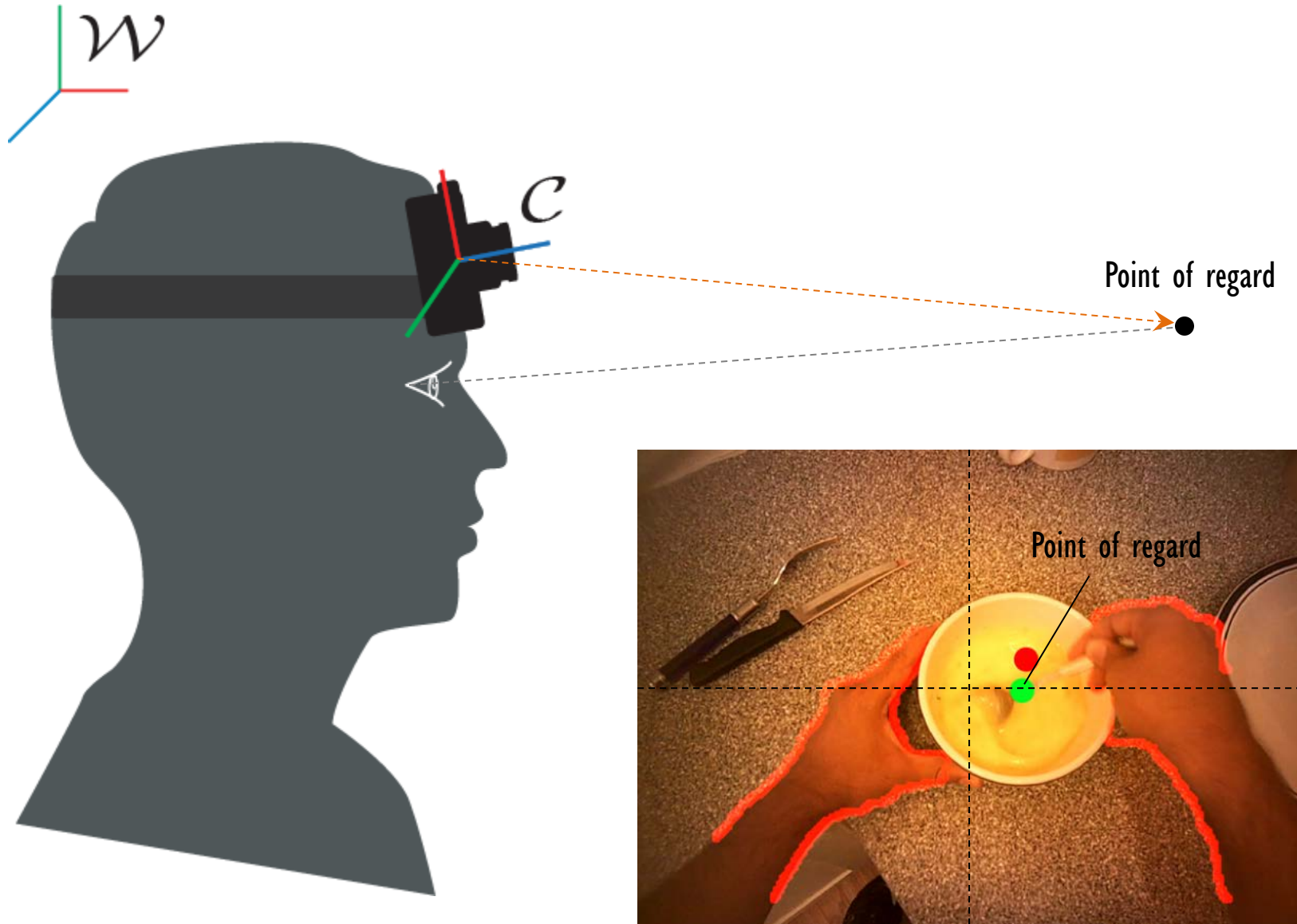
A close-up photograph of a person's face, specifically their eyes and nose, wearing a head-mounted device. The device consists of a black horizontal bar across the top of the head. A white camera with a black lens is mounted on the left side of the bar, pointing towards the person's eyes. A black cable runs from the camera, loops around the side of the head, and connects to a small white cylindrical component mounted on the right side of the bar. The person has light-colored eyes and is wearing clear safety glasses. The background is a solid blue color.

Looking-out camera
(global coordination)

Eye tracker
(local coordination)

Gaze Estimation

Gaze Estimation w/o Eye Tracker

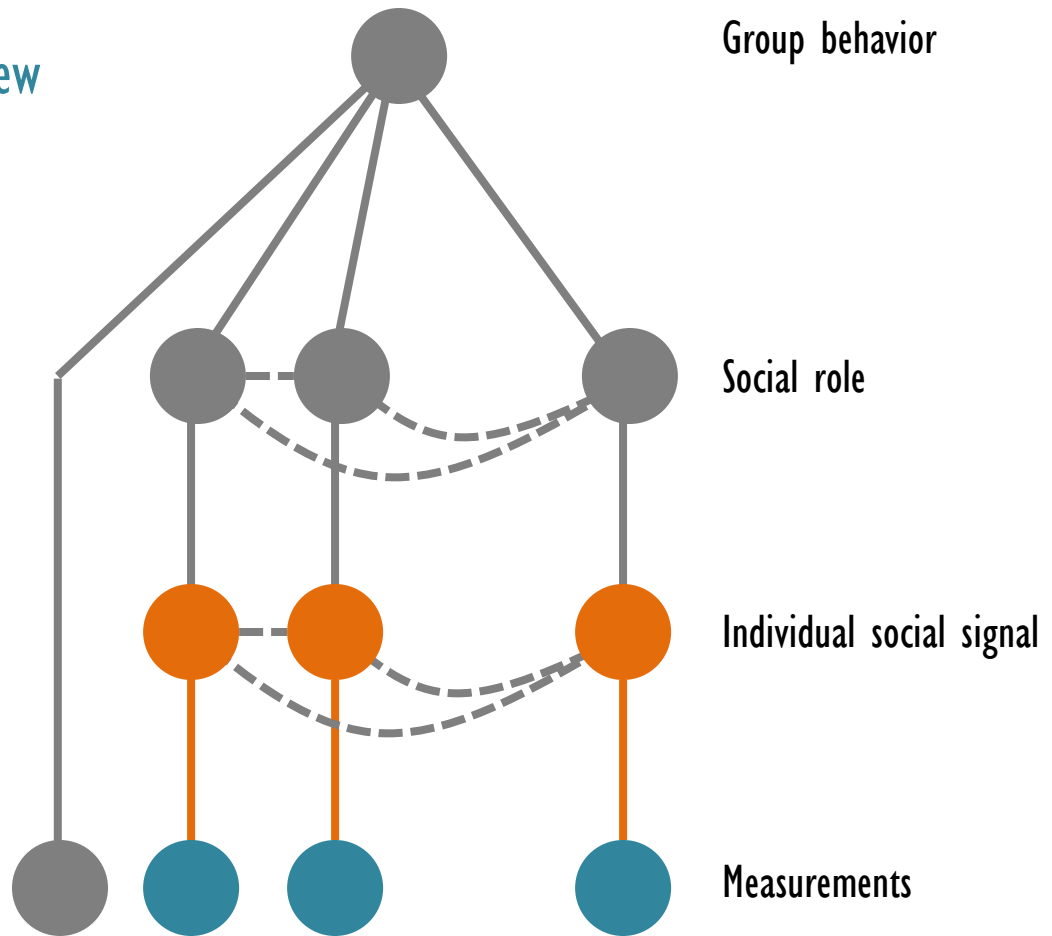
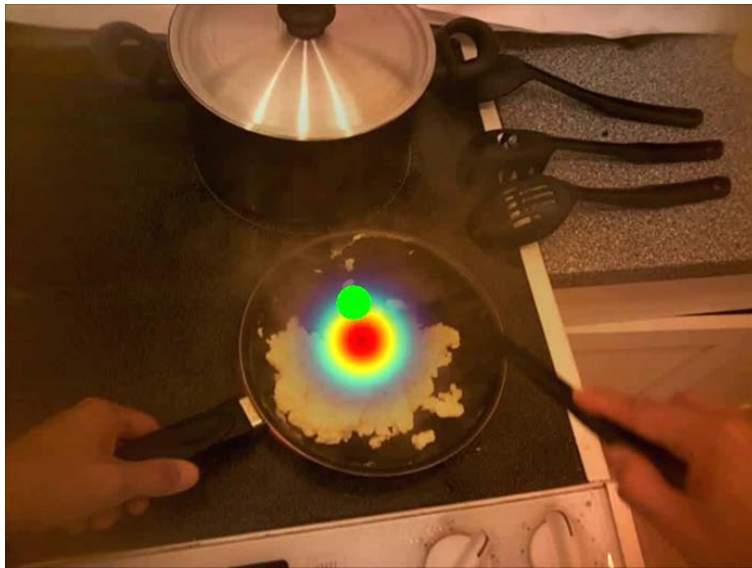


Social Signal Perception (Gaze)

[Li ICCV13]

Input: image or video of first person view

Output: localization of fixation point.

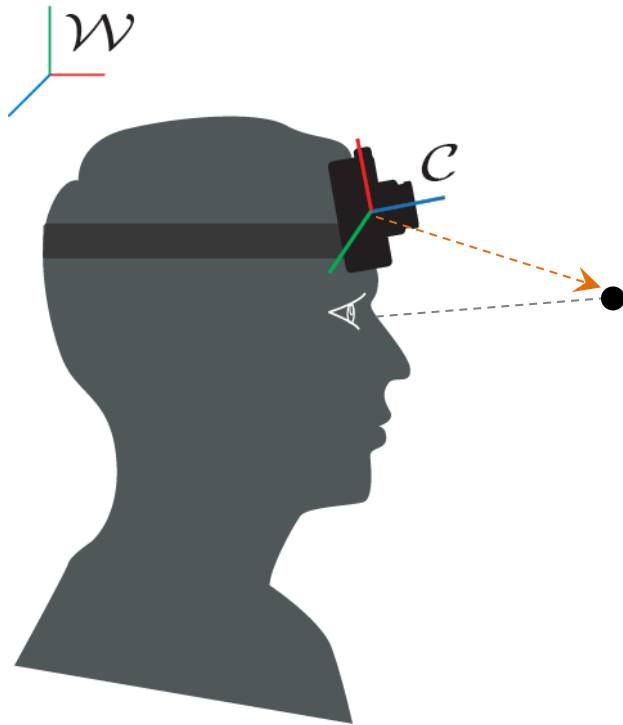


GTEA Gaze+ Dataset

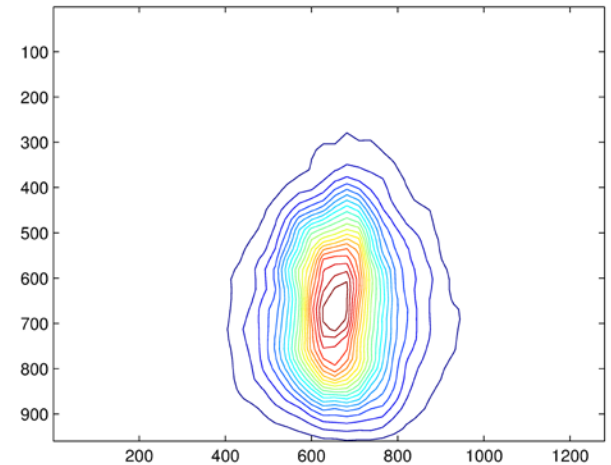
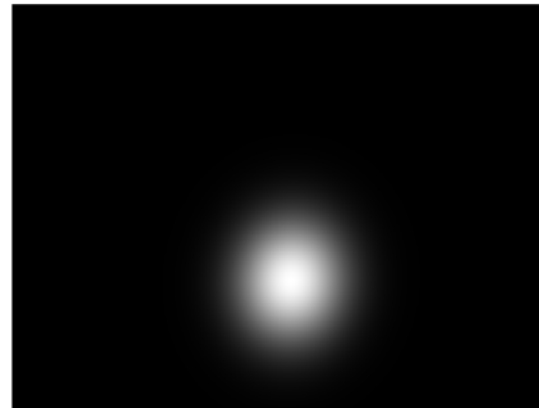


Egocentric Cue I

Eye-in-head Orientation (Center Prior)



GTEA Gaze+ Dataset



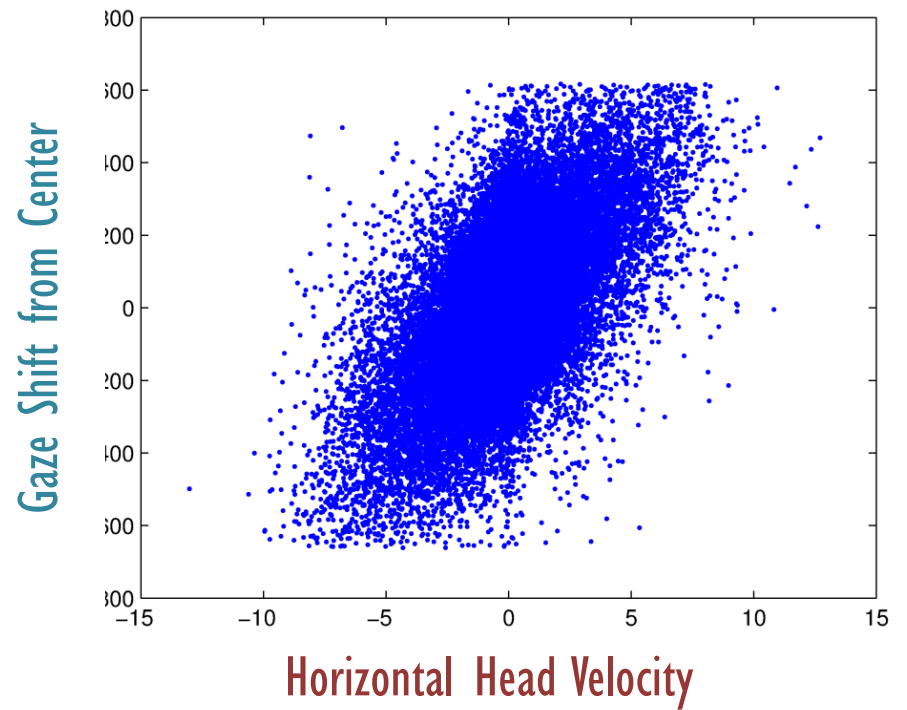
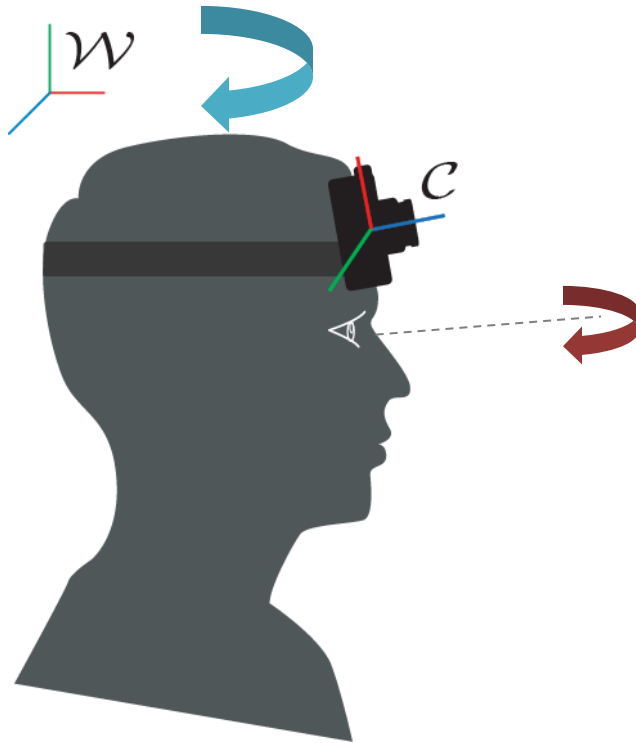
Egocentric Cue II

Head Motion



Egocentric Cue II

Head Motion



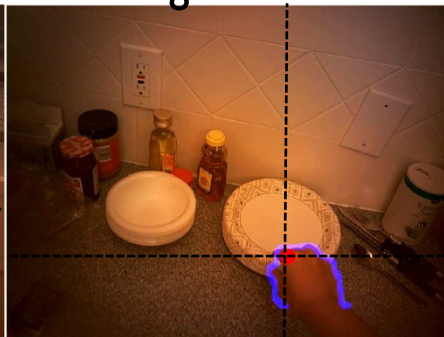
Egocentric Cue III

Hand Position

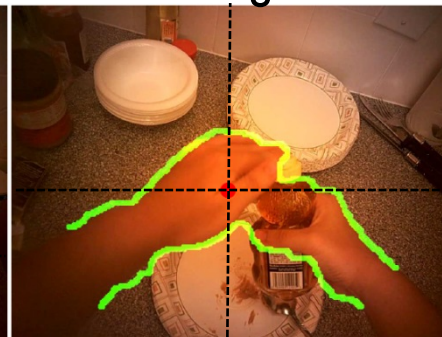
Left Hand



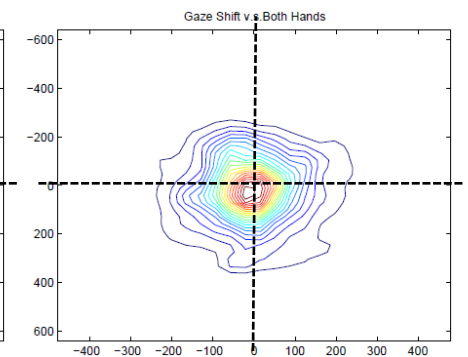
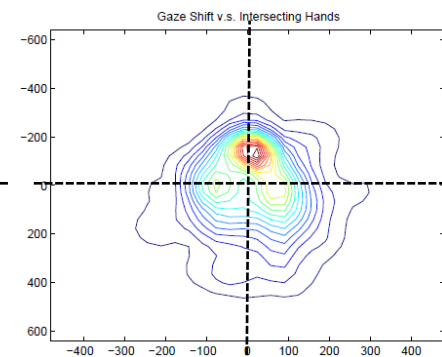
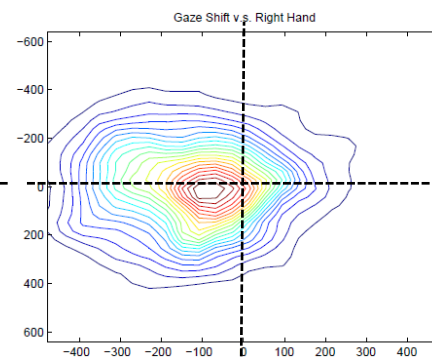
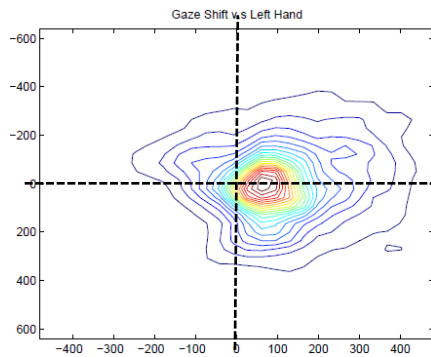
Right Hand



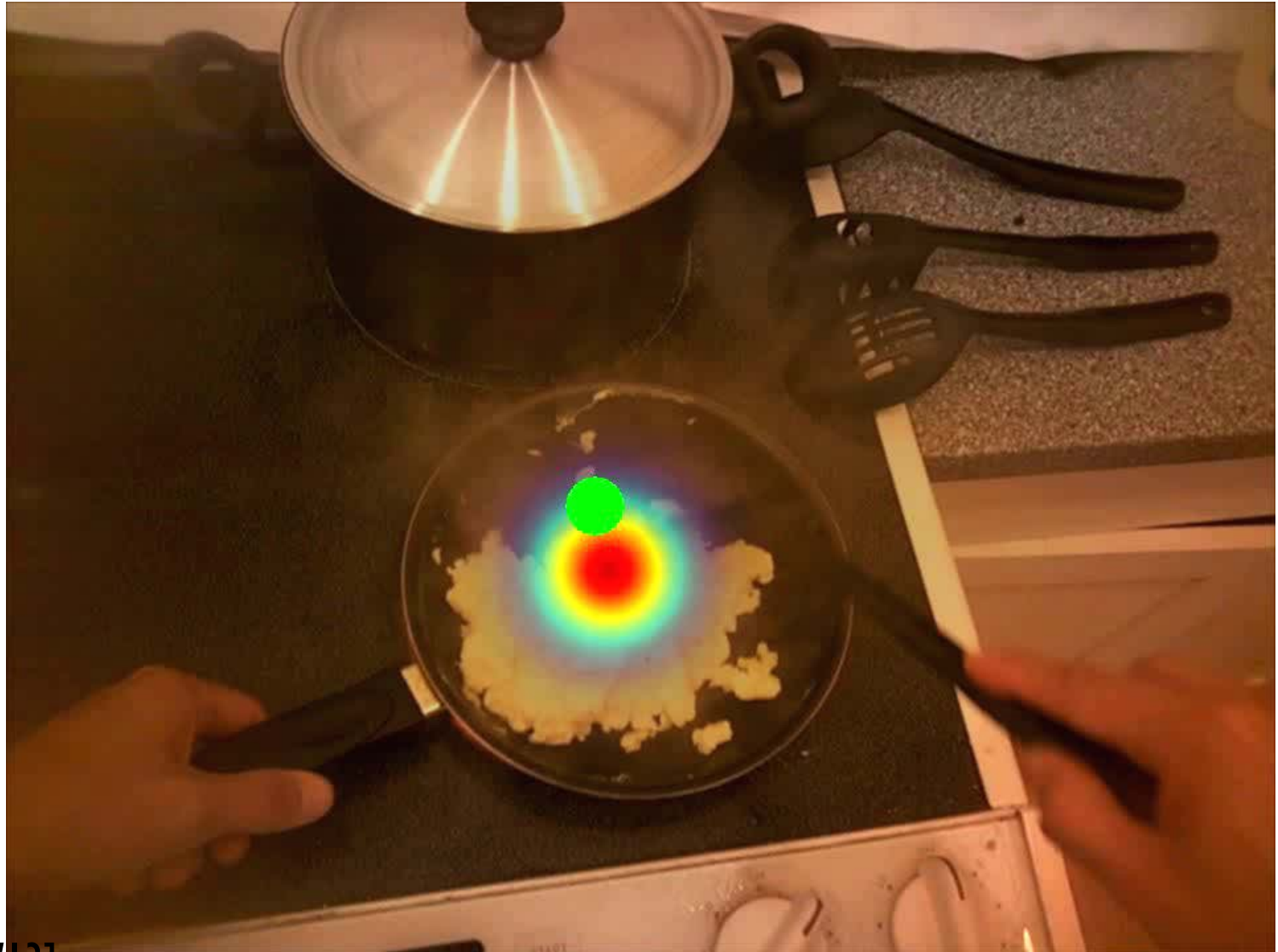
Hands Together



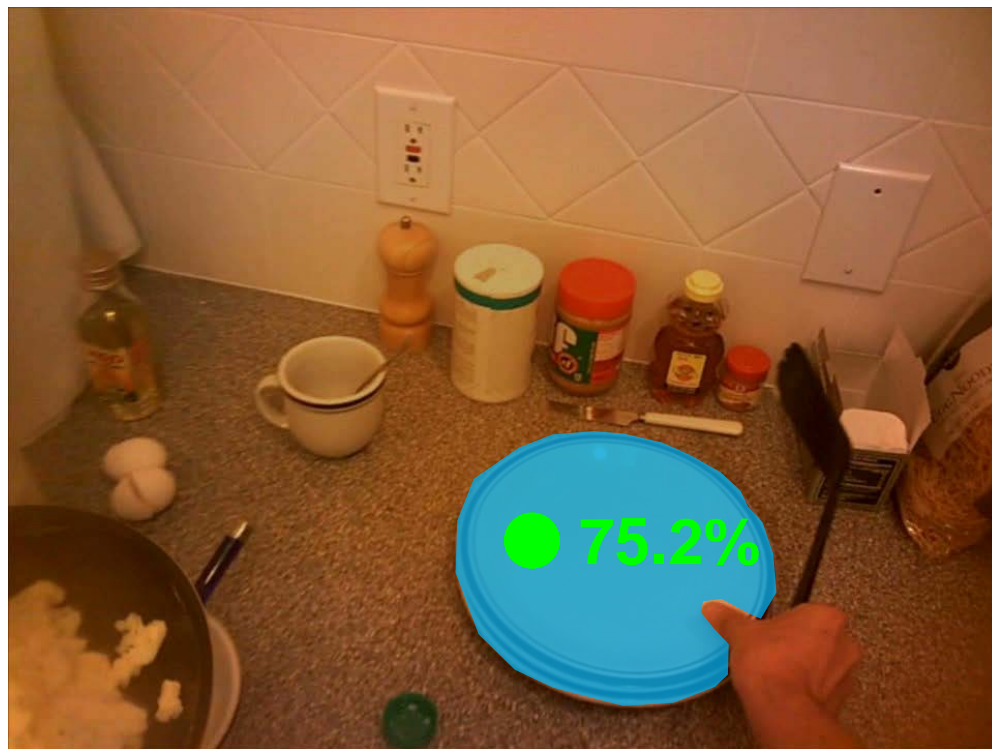
Hands Apart



Gaze Prediction



Application to Foreground Segmentation



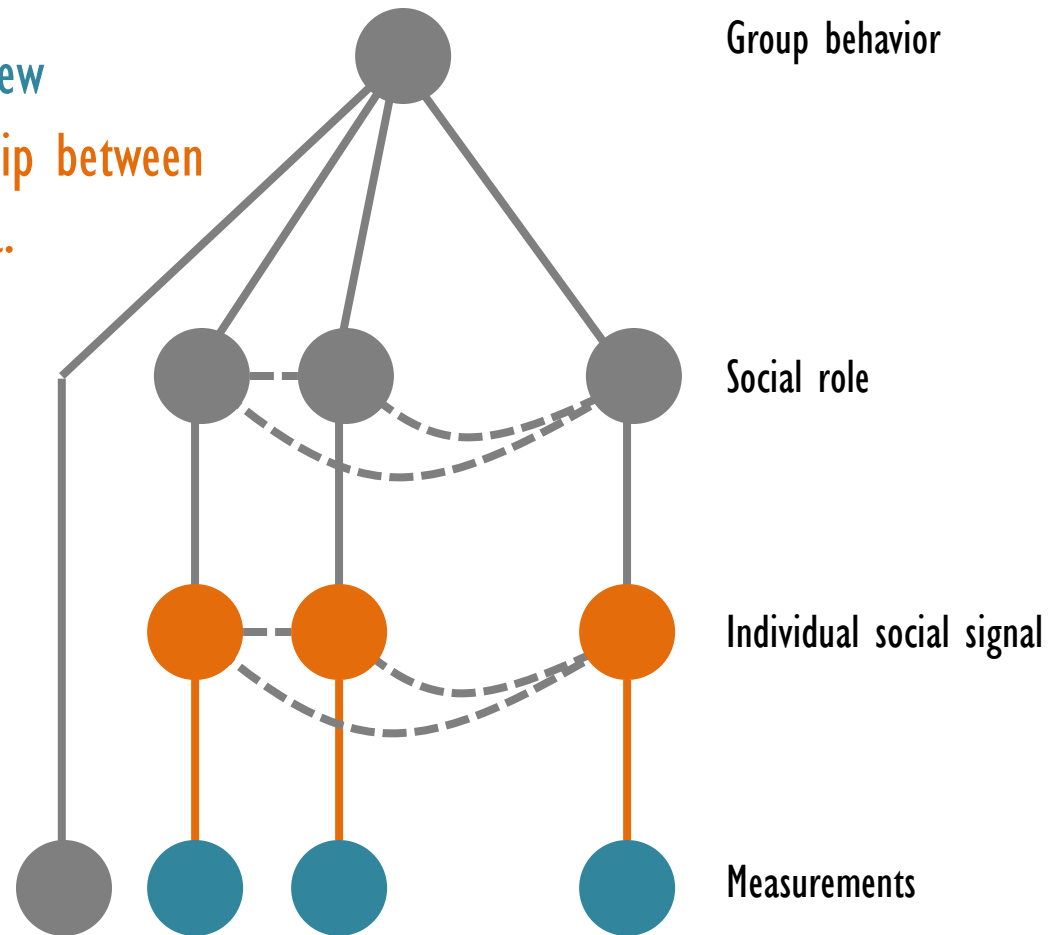
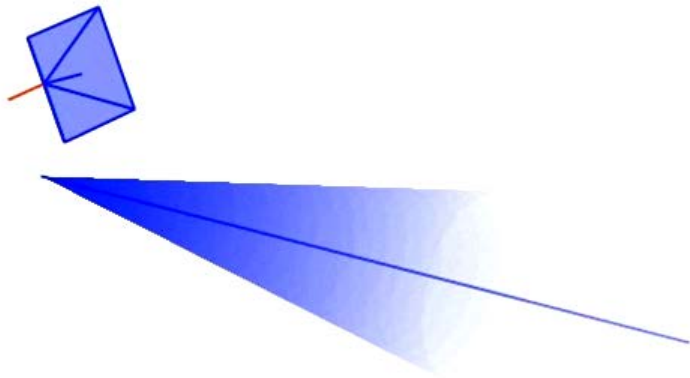
Foreground Object

Social Signal Perception (Gaze)

[Park NIPS12]

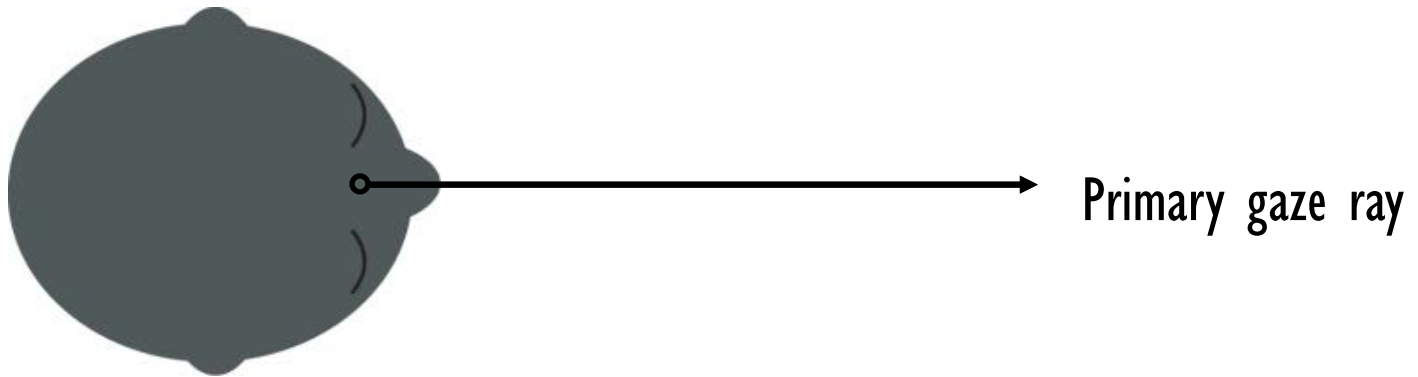
Input: image or video of first person view

Output: to calibrate the fixed relationship between
gaze and head-mounted camera.



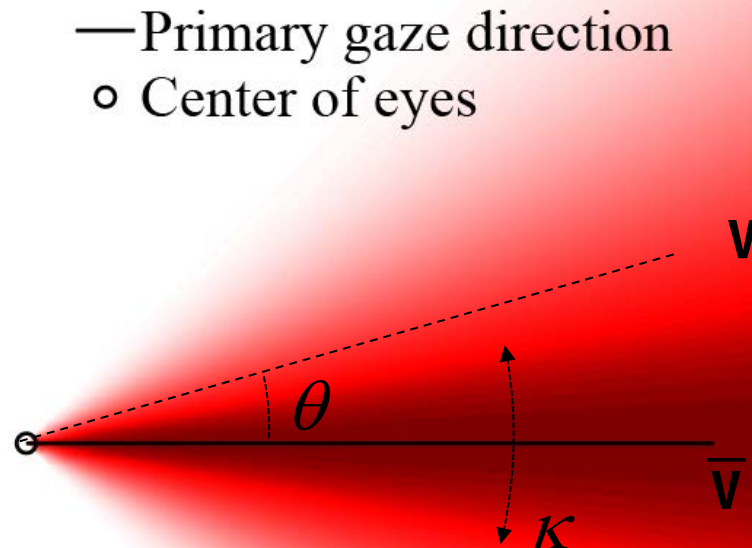
Eye-in-head Motion

[Park NIPS12]



Gaze Distribution

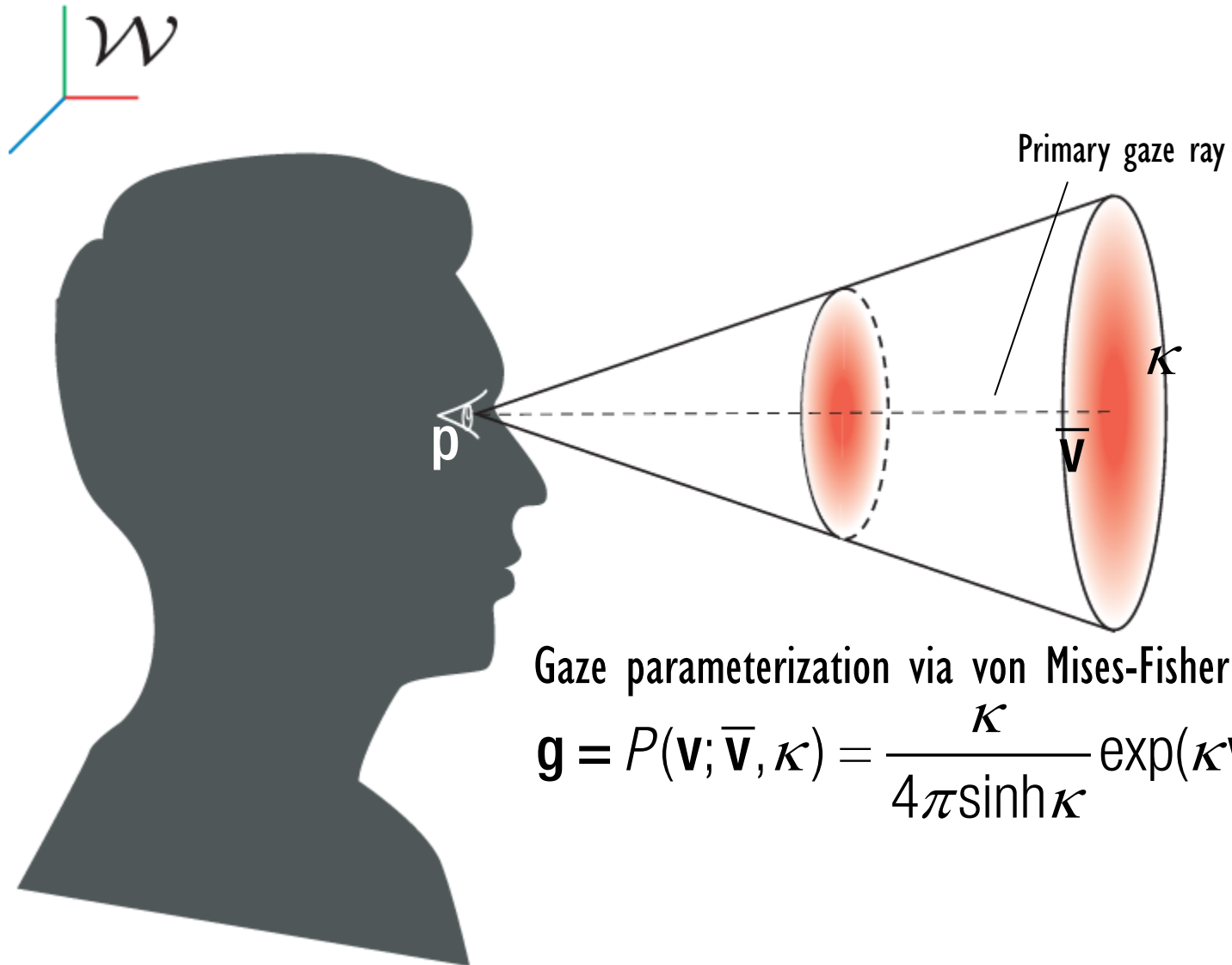
$$\cos(\theta) = \mathbf{v}^\top \bar{\mathbf{v}}$$



von Mises-Fisher distribution

$$P(\mathbf{v}; \bar{\mathbf{v}}, \kappa) = \frac{\kappa}{4\pi \sinh \kappa} \exp(\kappa \mathbf{v}^\top \bar{\mathbf{v}})$$

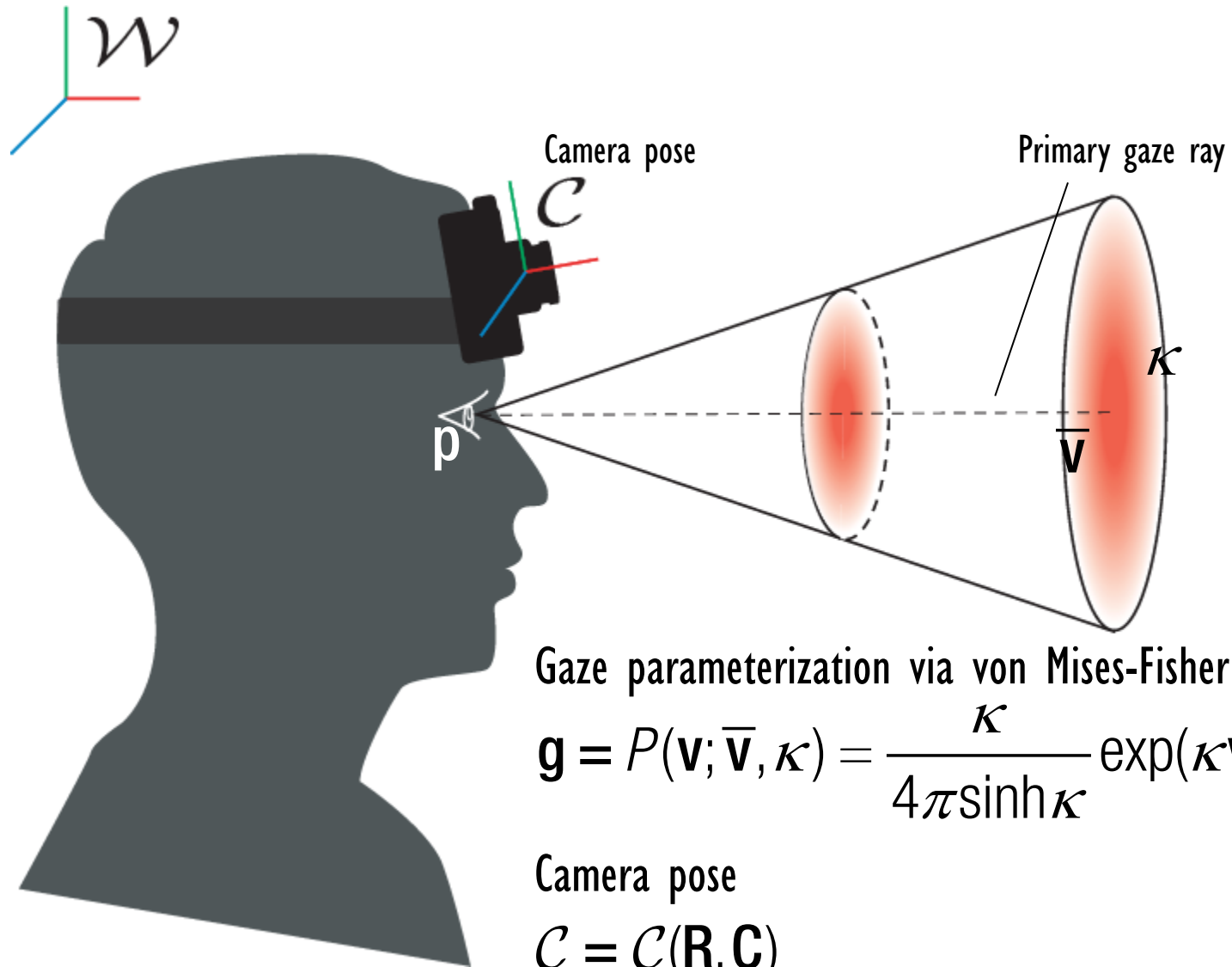
Cone-shaped Gaze Model



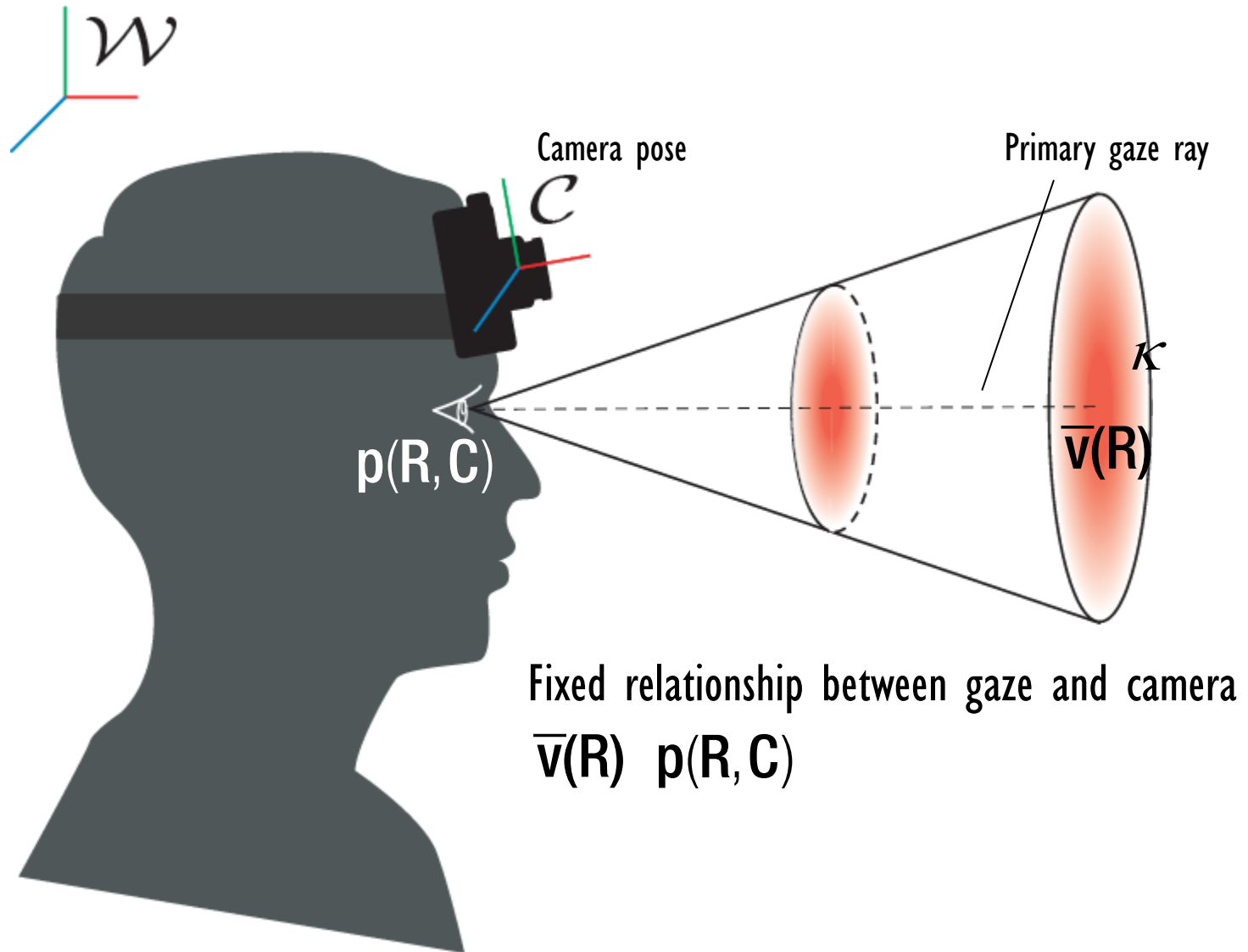
Gaze parameterization via von Mises-Fisher distribution

$$\mathbf{g} = P(\mathbf{v}; \bar{\mathbf{v}}, \kappa) = \frac{\kappa}{4\pi \sinh \kappa} \exp(\kappa \mathbf{v}^T \bar{\mathbf{v}})$$

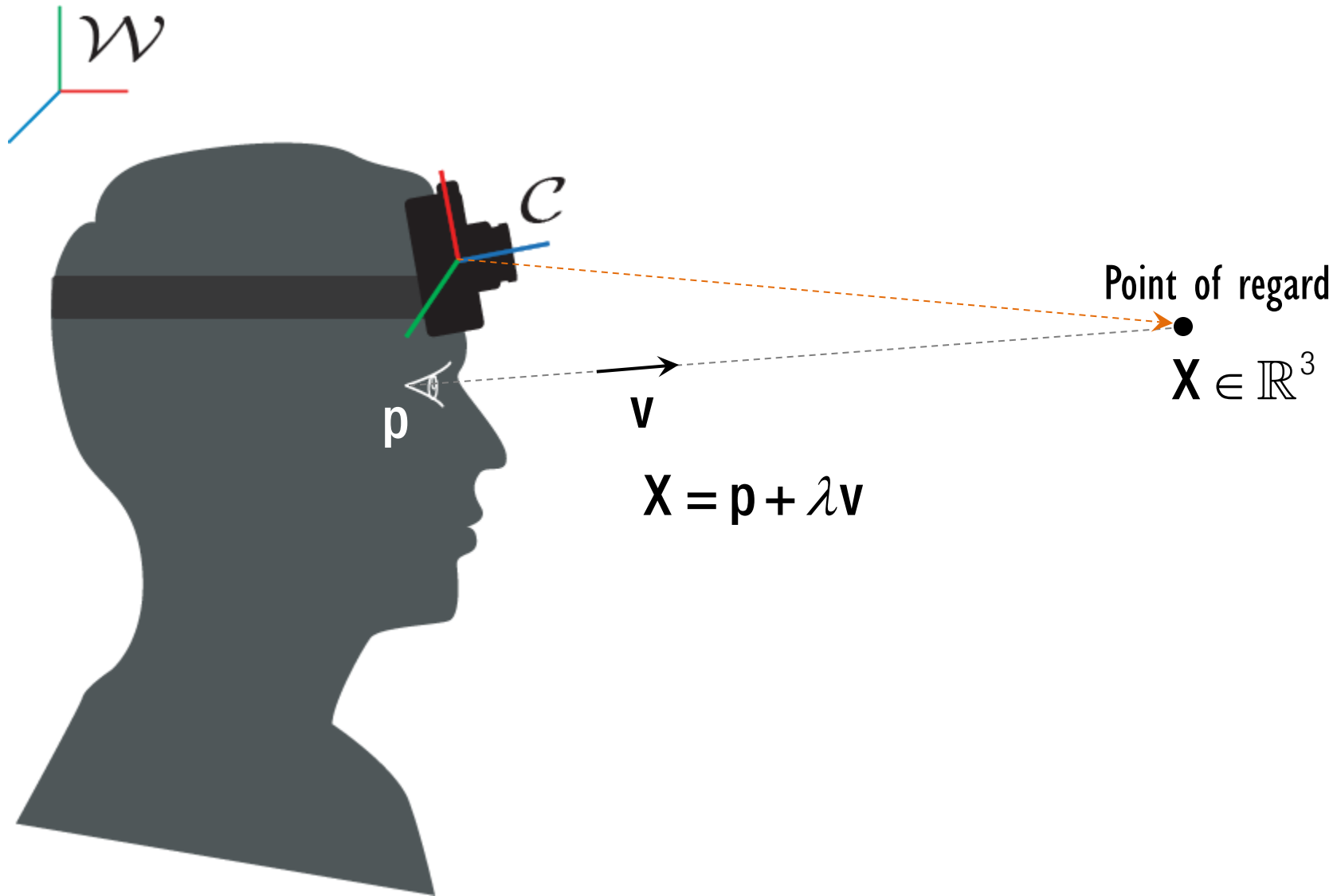
Gaze Calibration



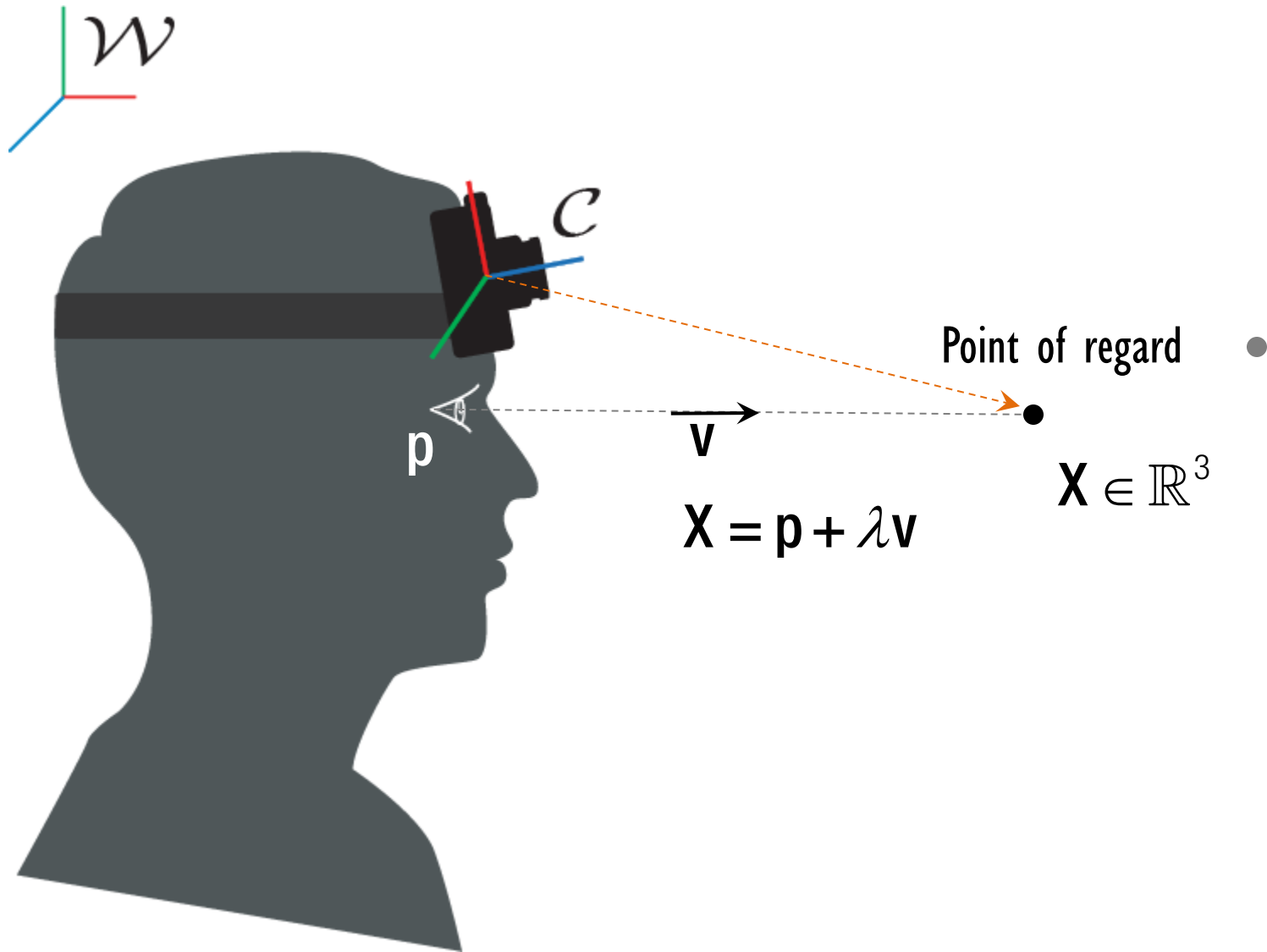
Gaze Calibration



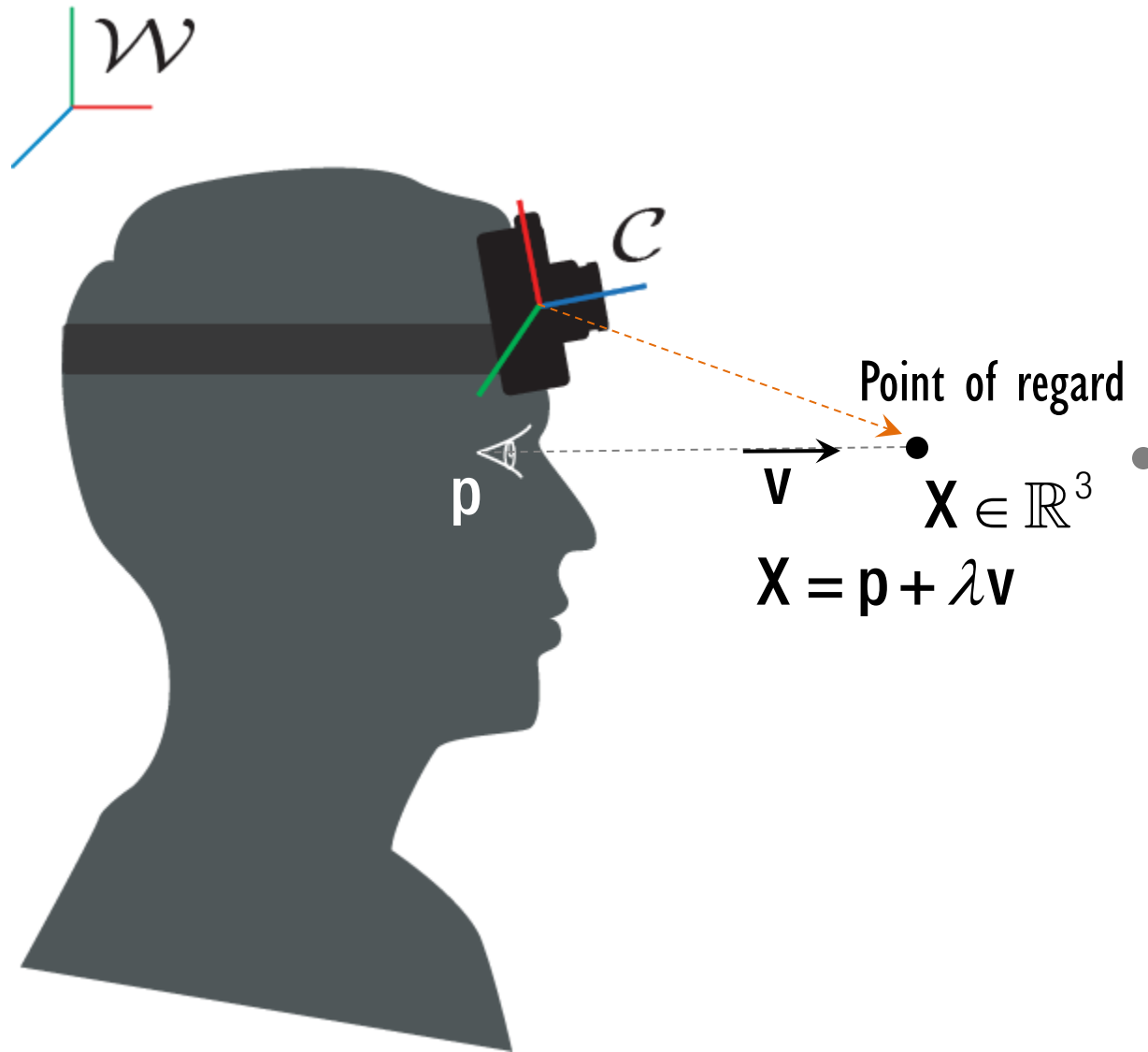
Gaze Calibration



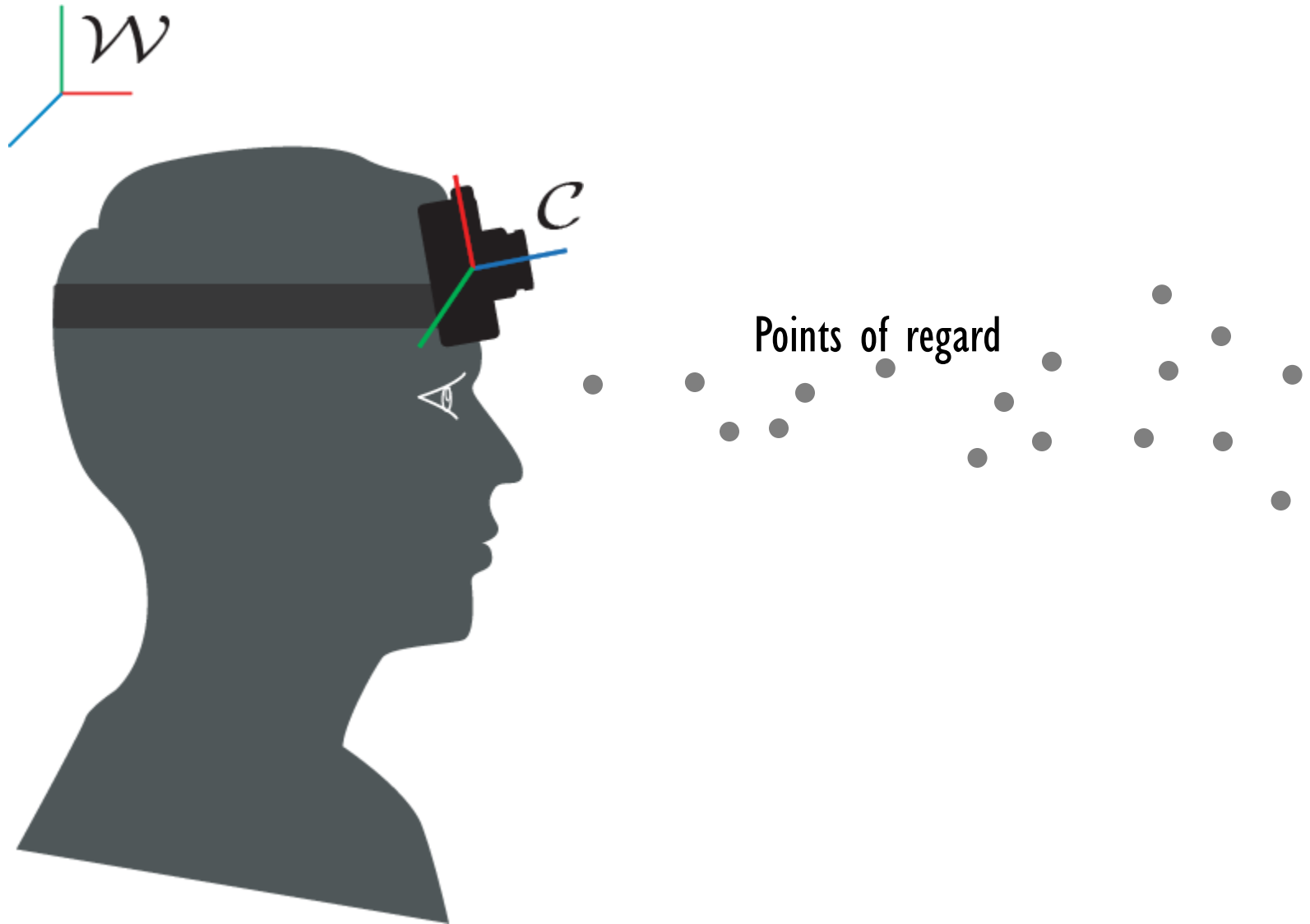
Gaze Calibration



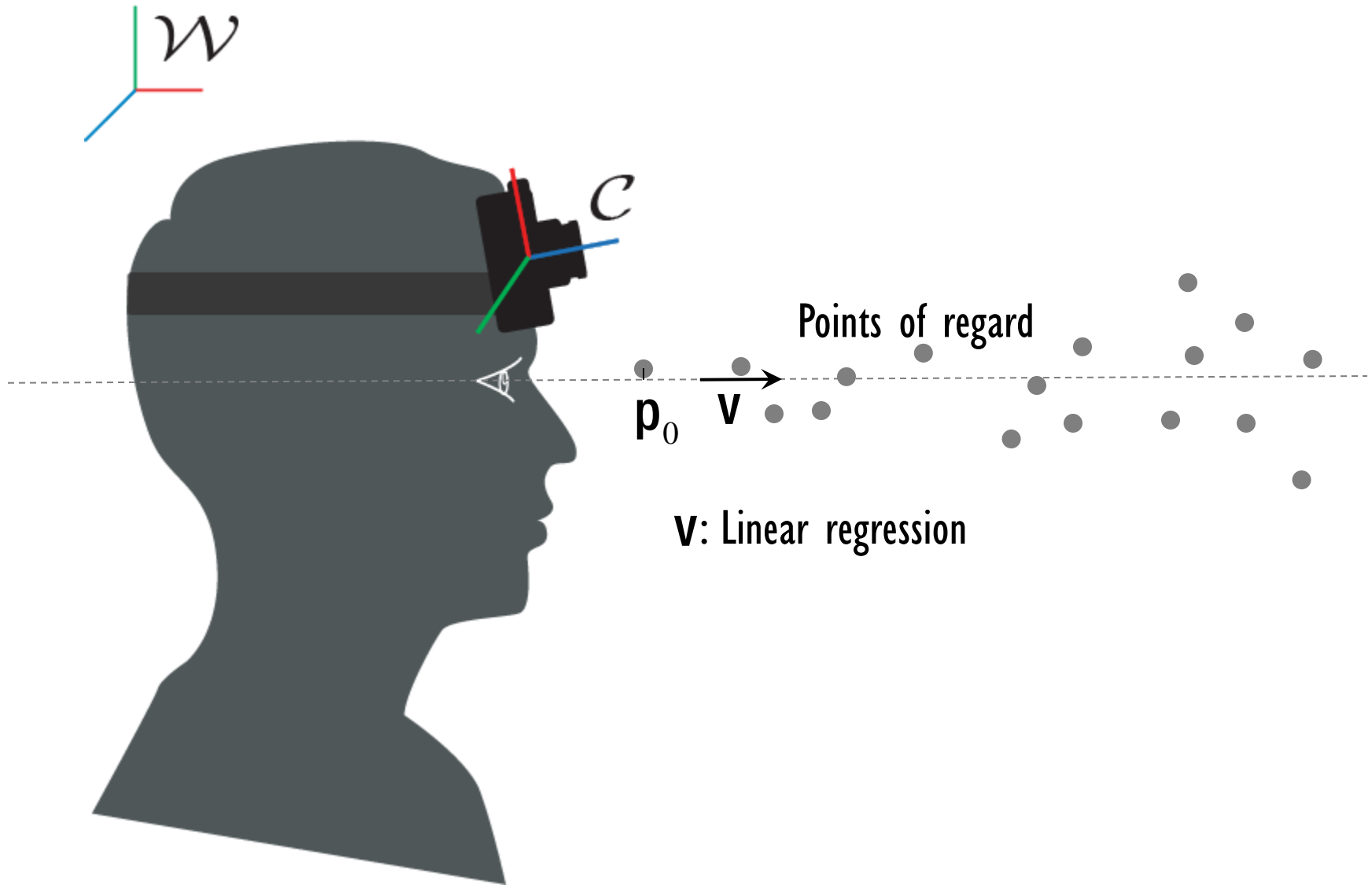
Gaze Calibration



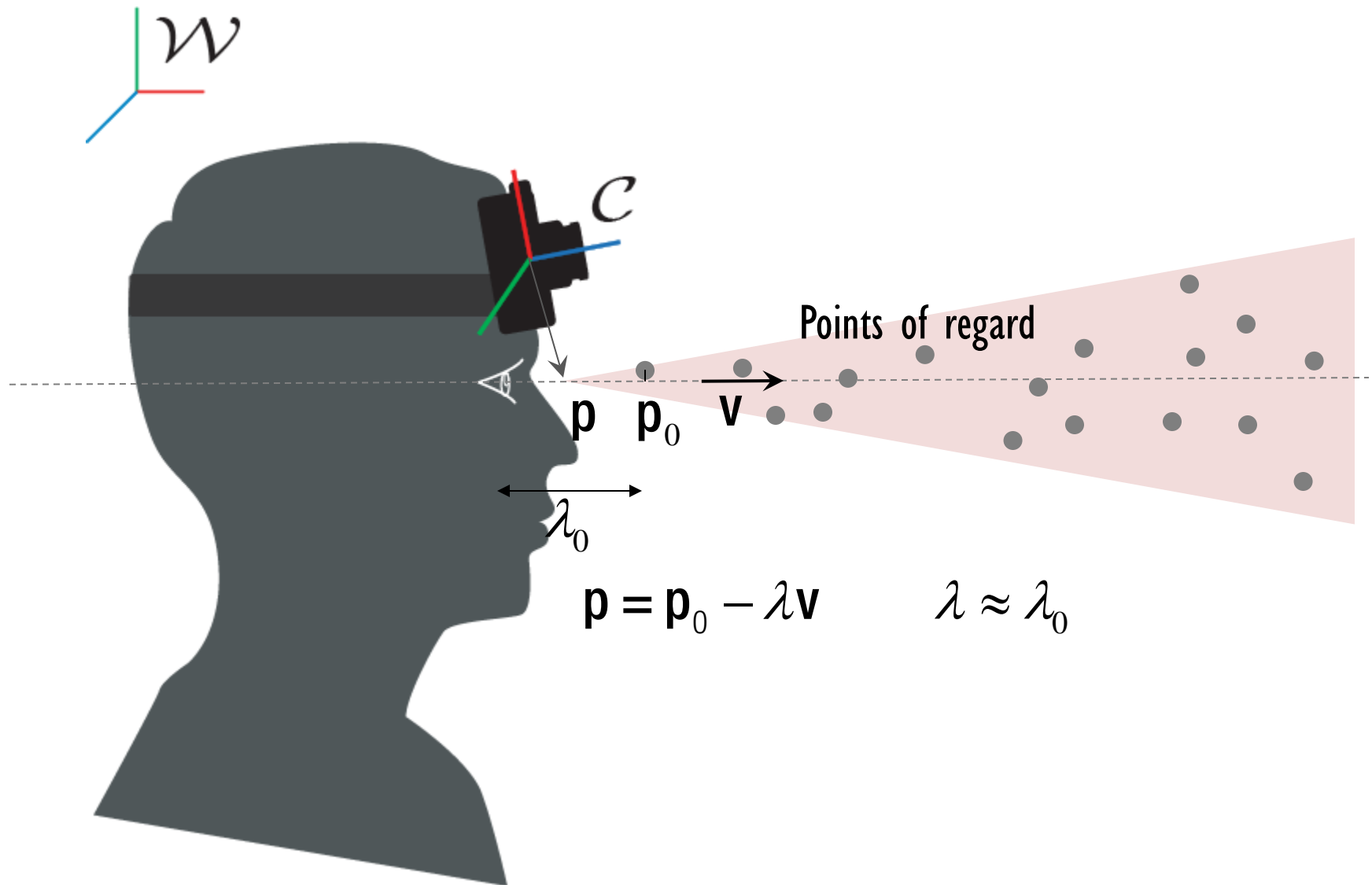
Gaze Calibration



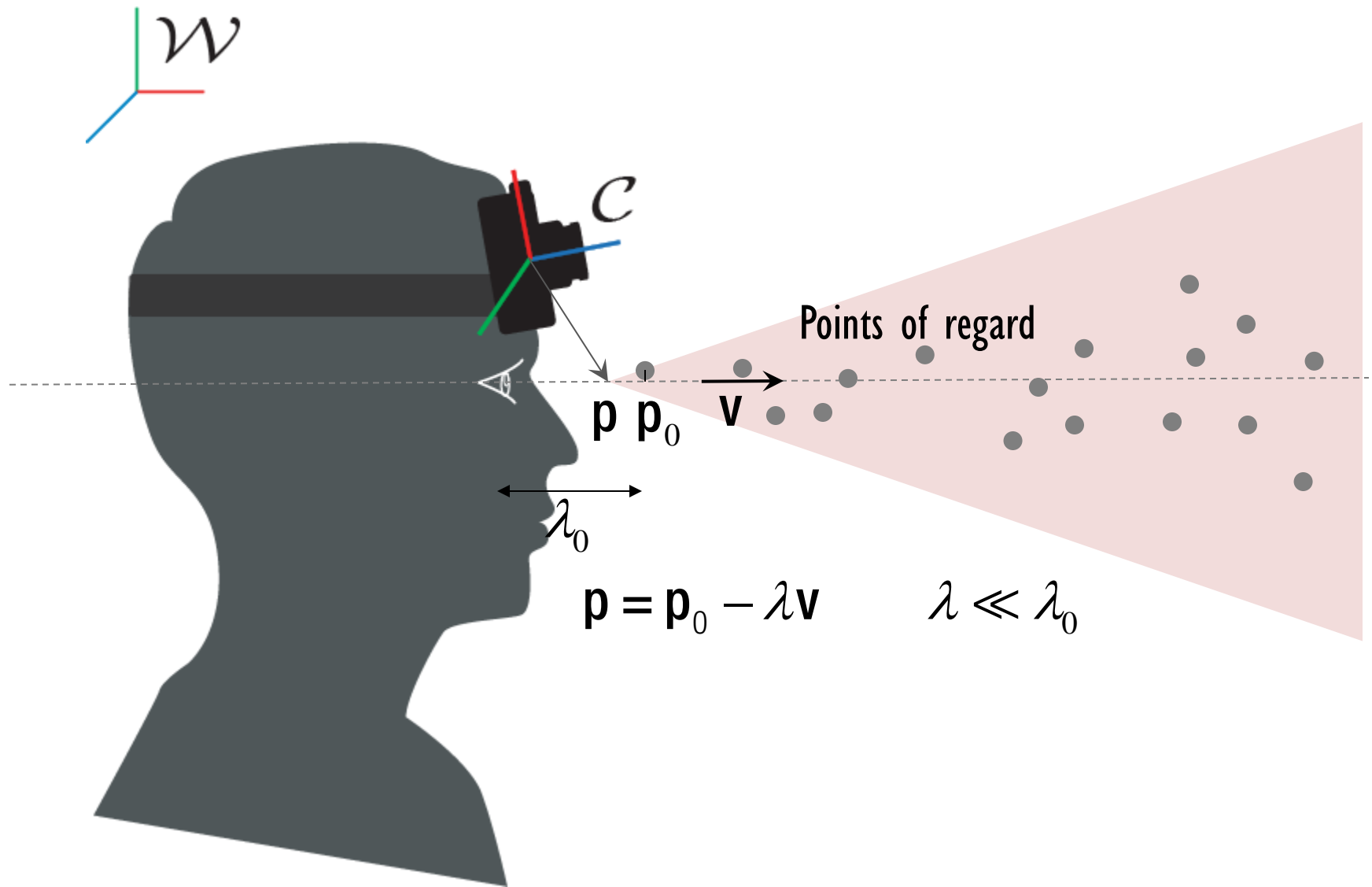
Gaze Calibration



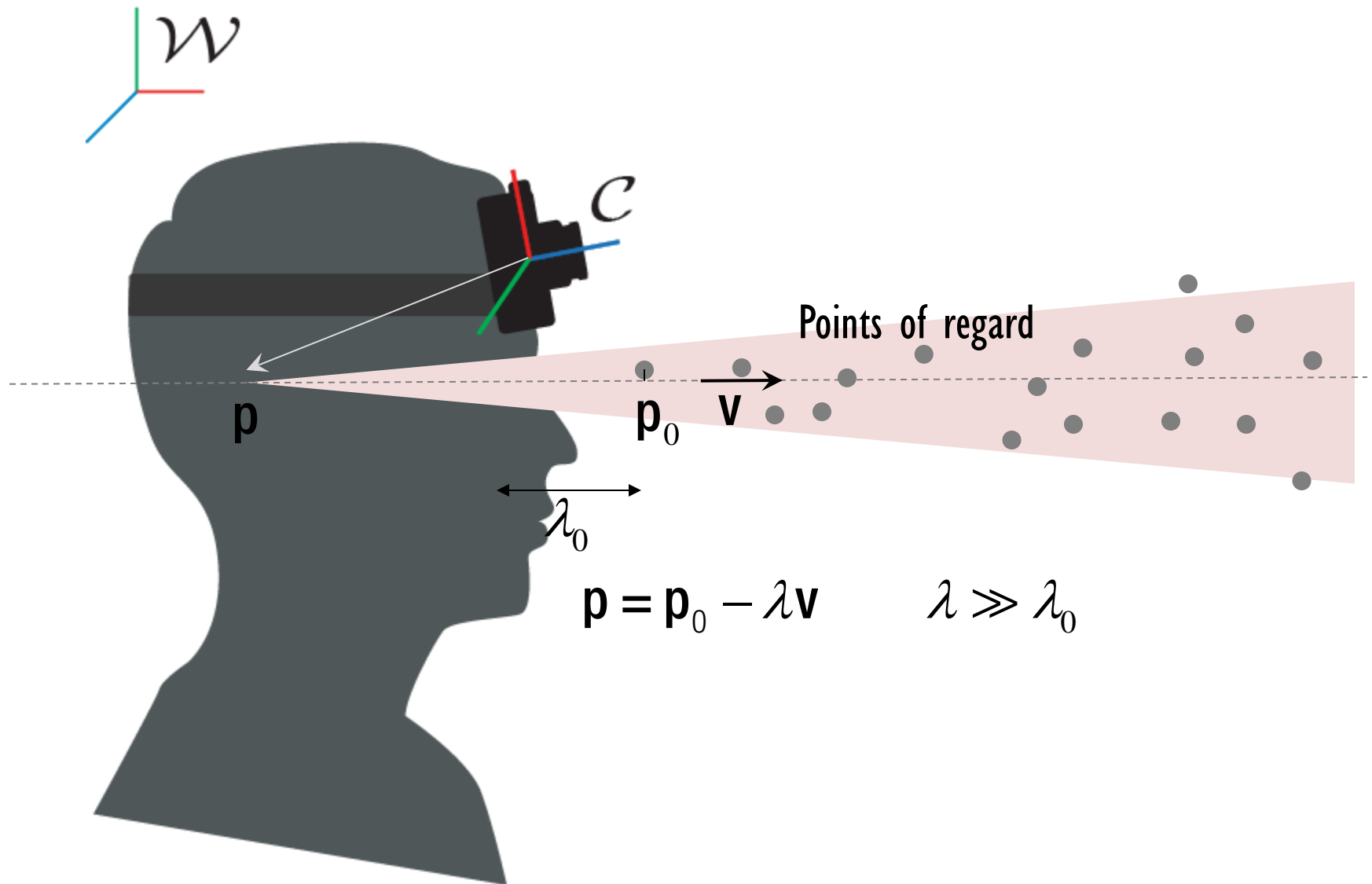
Gaze Calibration



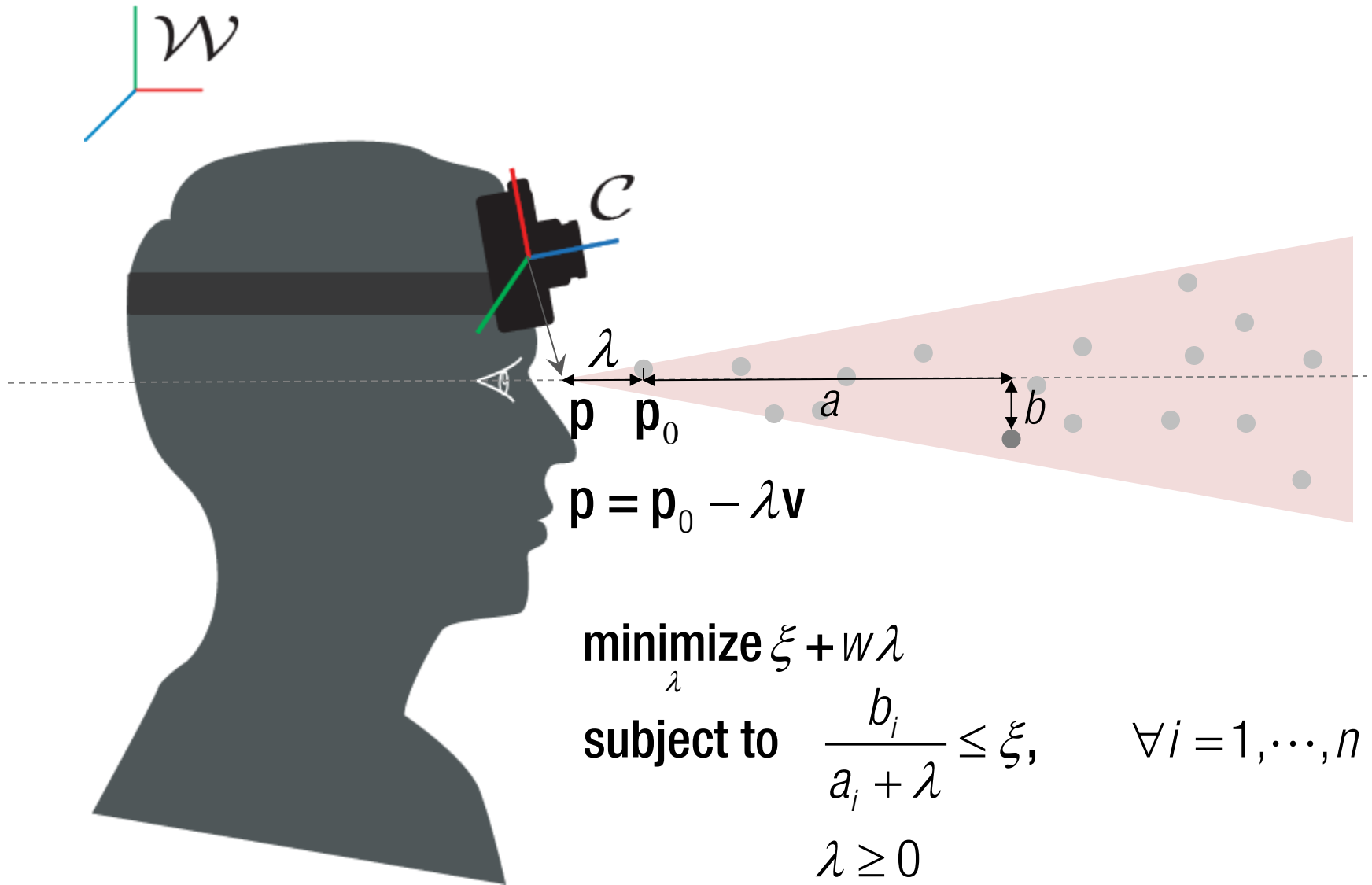
Gaze Calibration



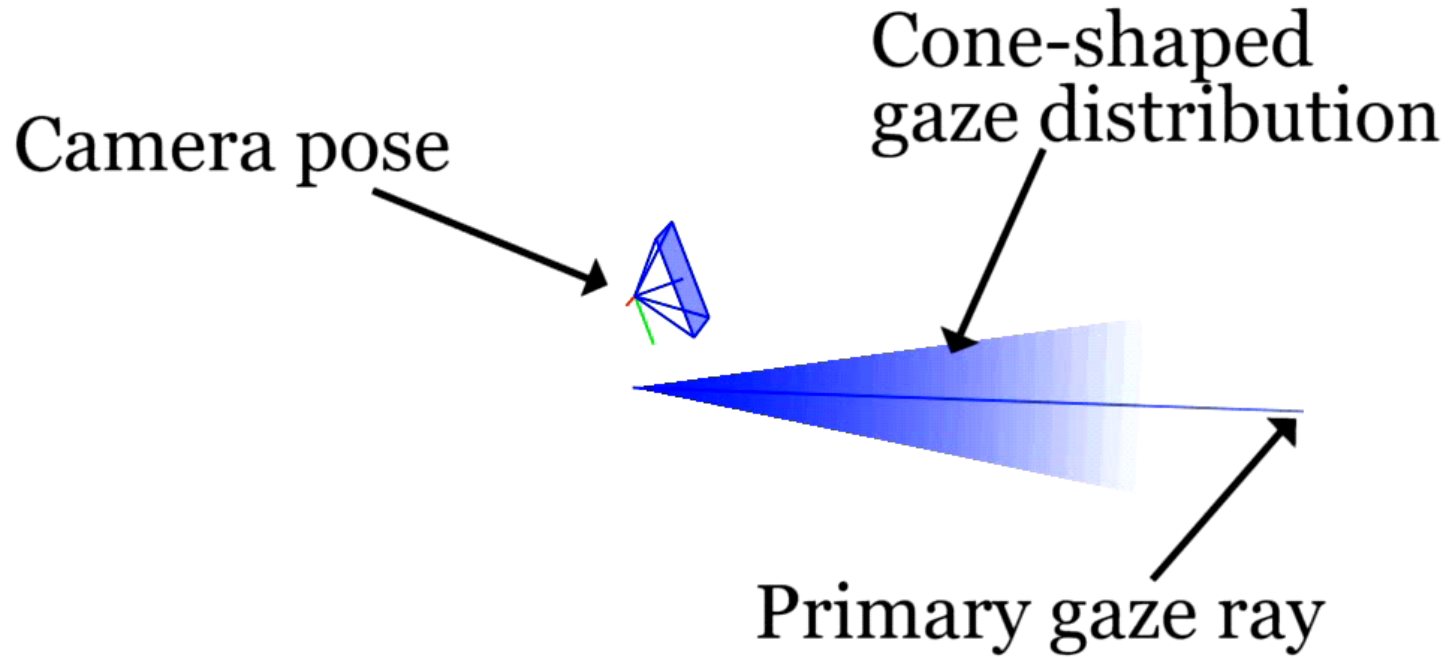
Gaze Calibration



Gaze Calibration



Gaze Ray Calibration

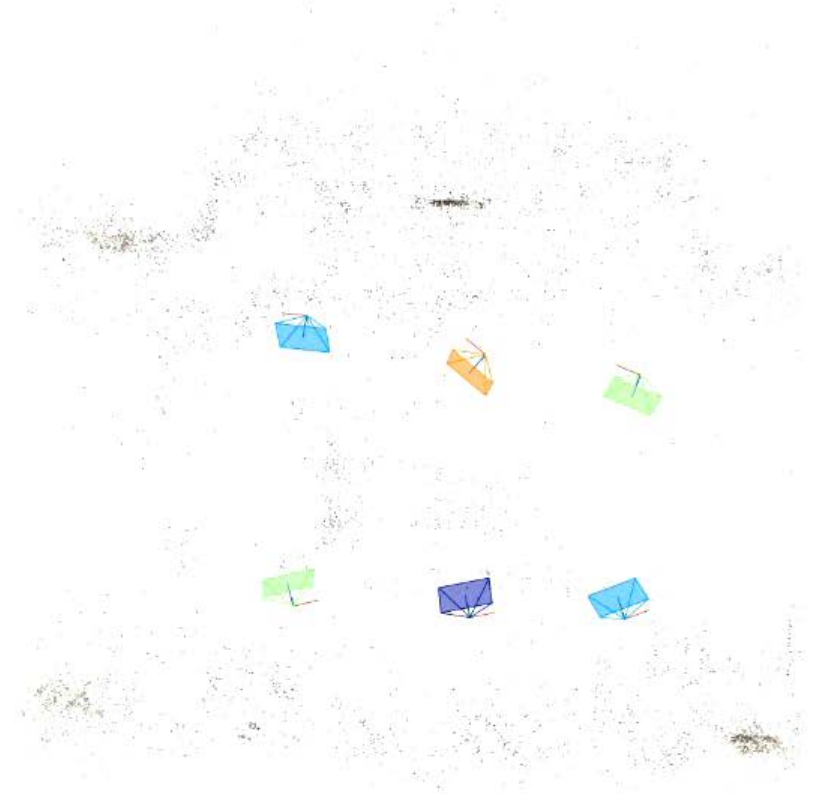


Primary gaze ray with respect to the camera pose

Gaze Ray Calibration



Back and forth motion



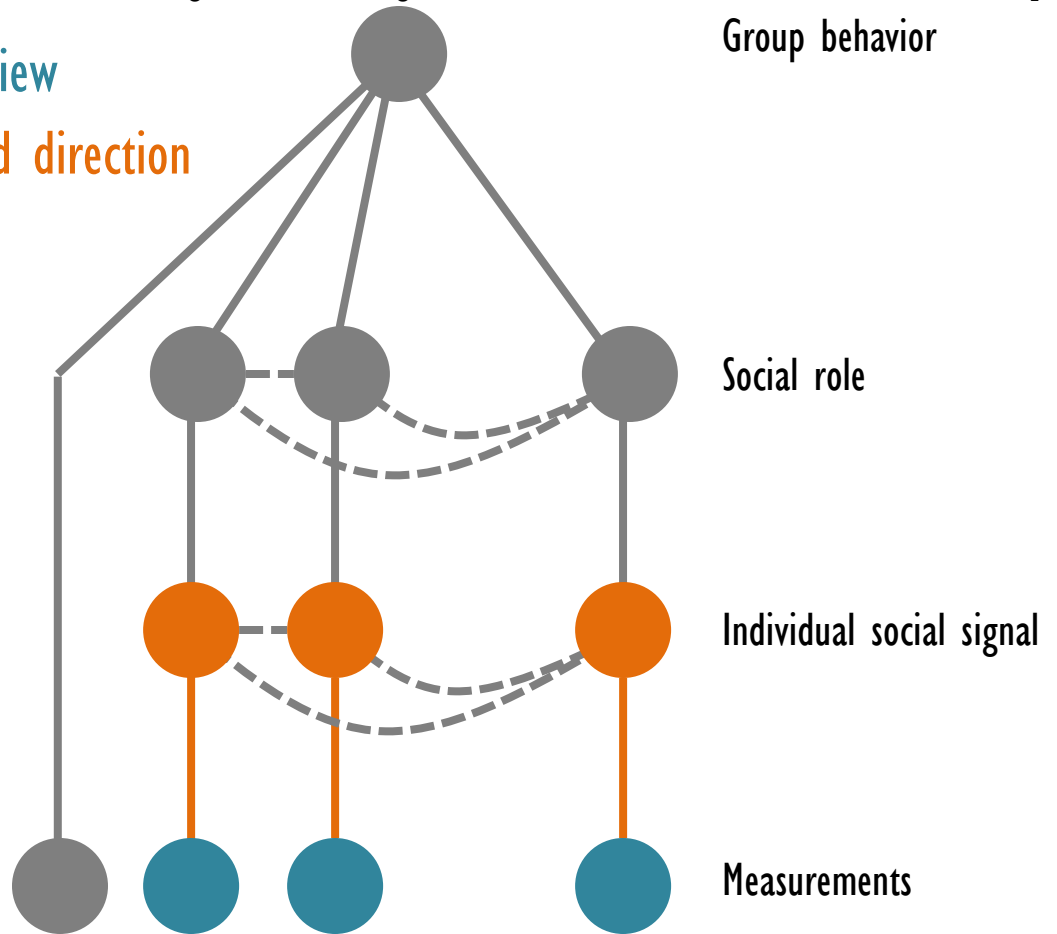
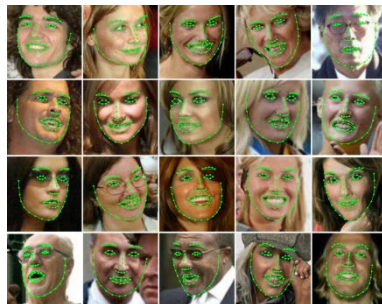
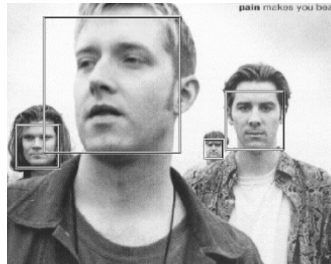
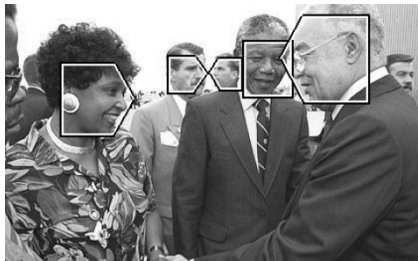
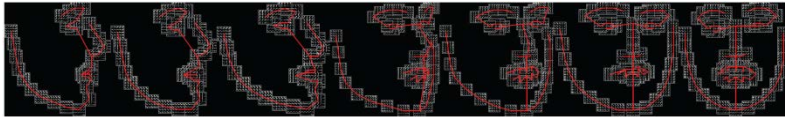
Side to side motion

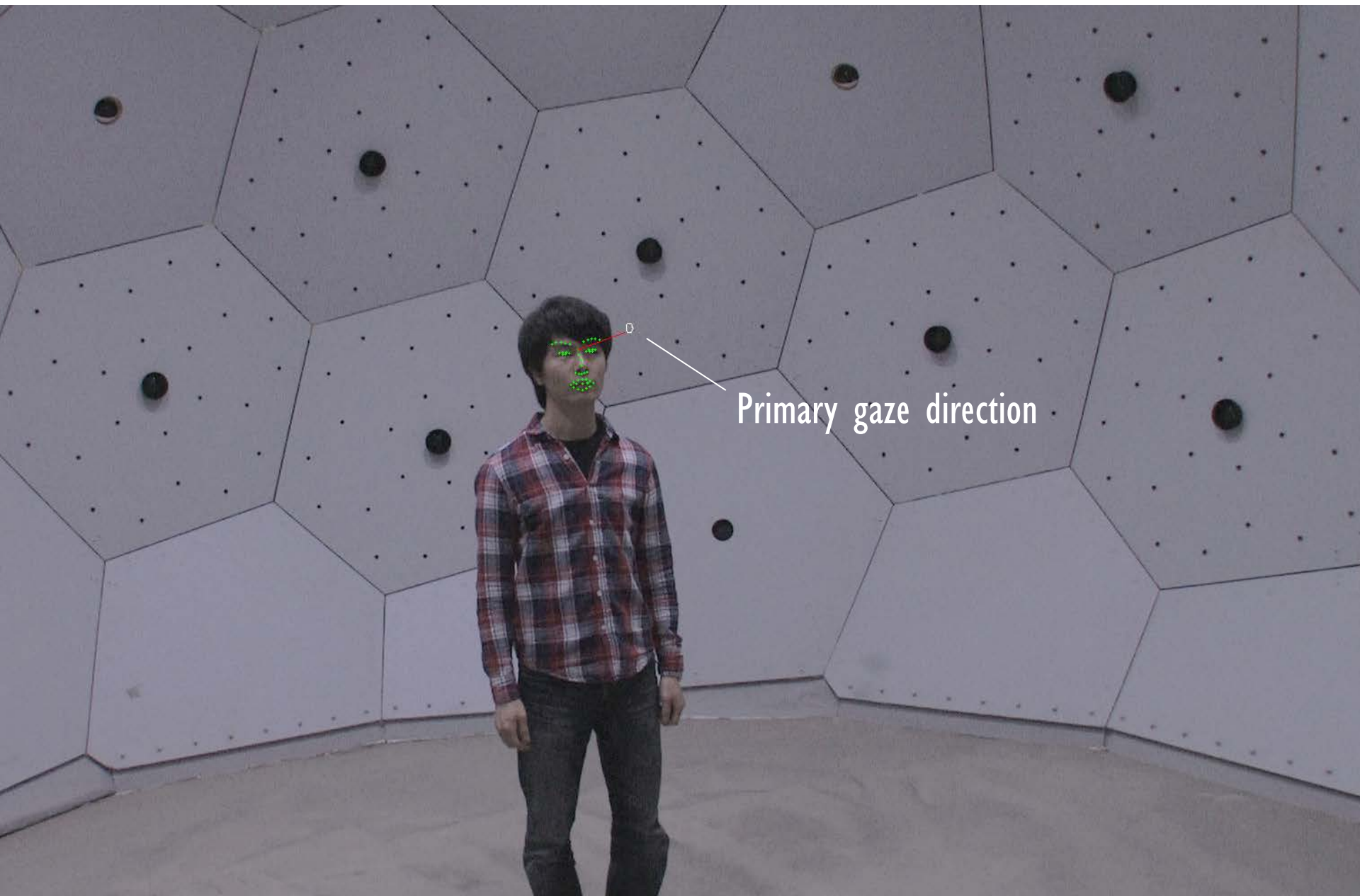
Head Detection/Alignment

[Cootes, PAMI01, Schneiderman, CVPR00, Viola, IJCV01, Matthews, IJCV04, Saragih, ICCV09, Xiong, CVPR13, Zhu, CVPR12, Marin-Jimenez, IJCV14, ...]

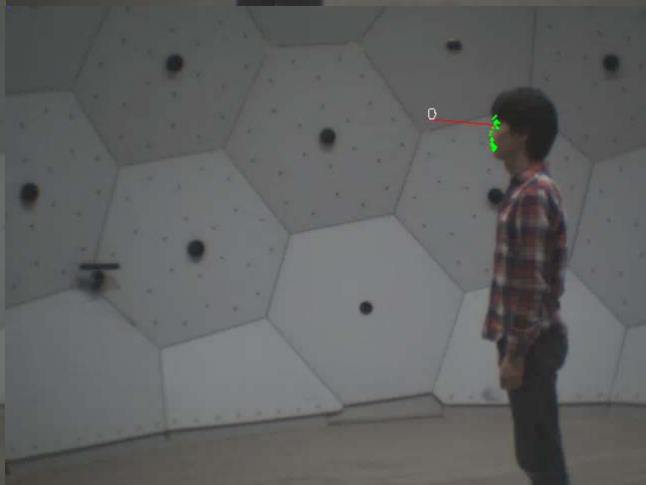
Input: image or video of third person view

Output: to find head / to estimate head direction





Primary gaze direction



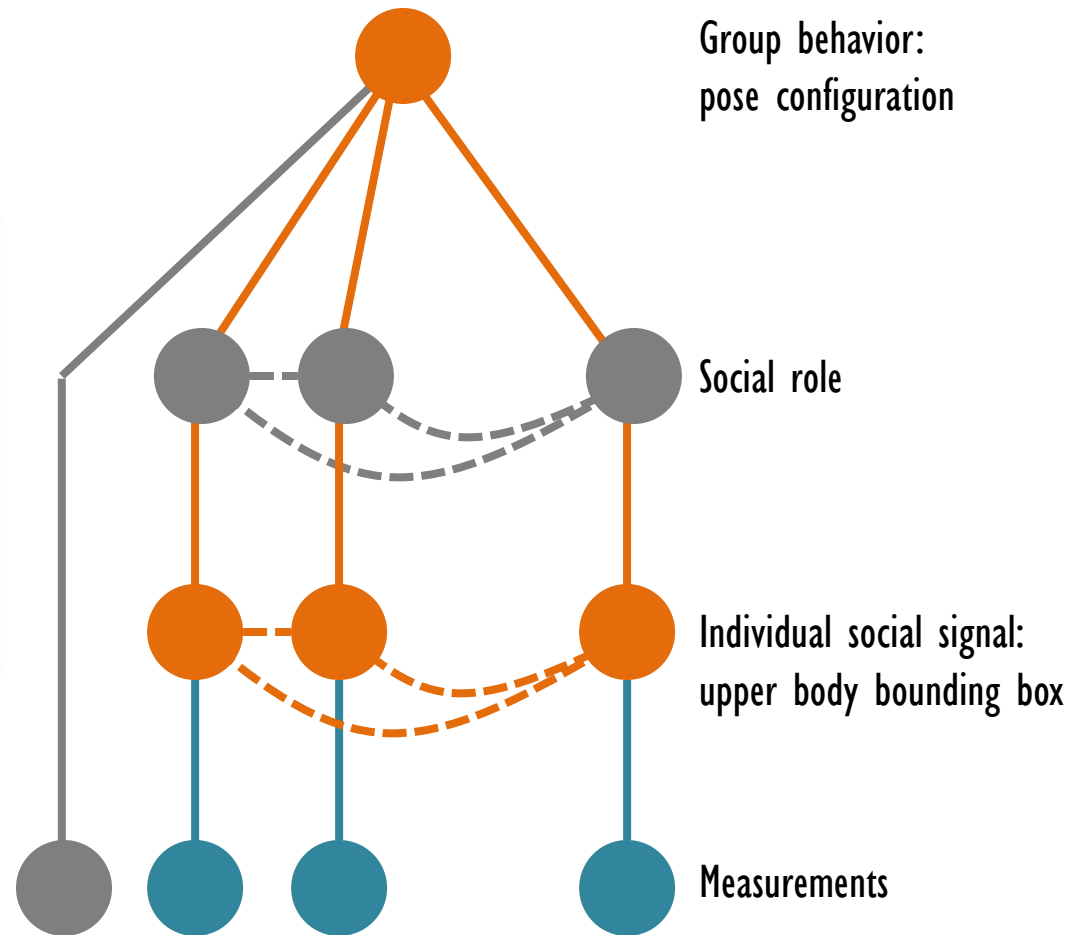
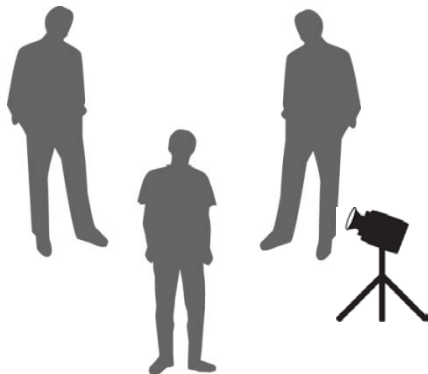


Body Configuration

[Hoai CVPR14]

Input: images of TV shows

Output: to detect social interactions





Characteristics of TV show interactions:

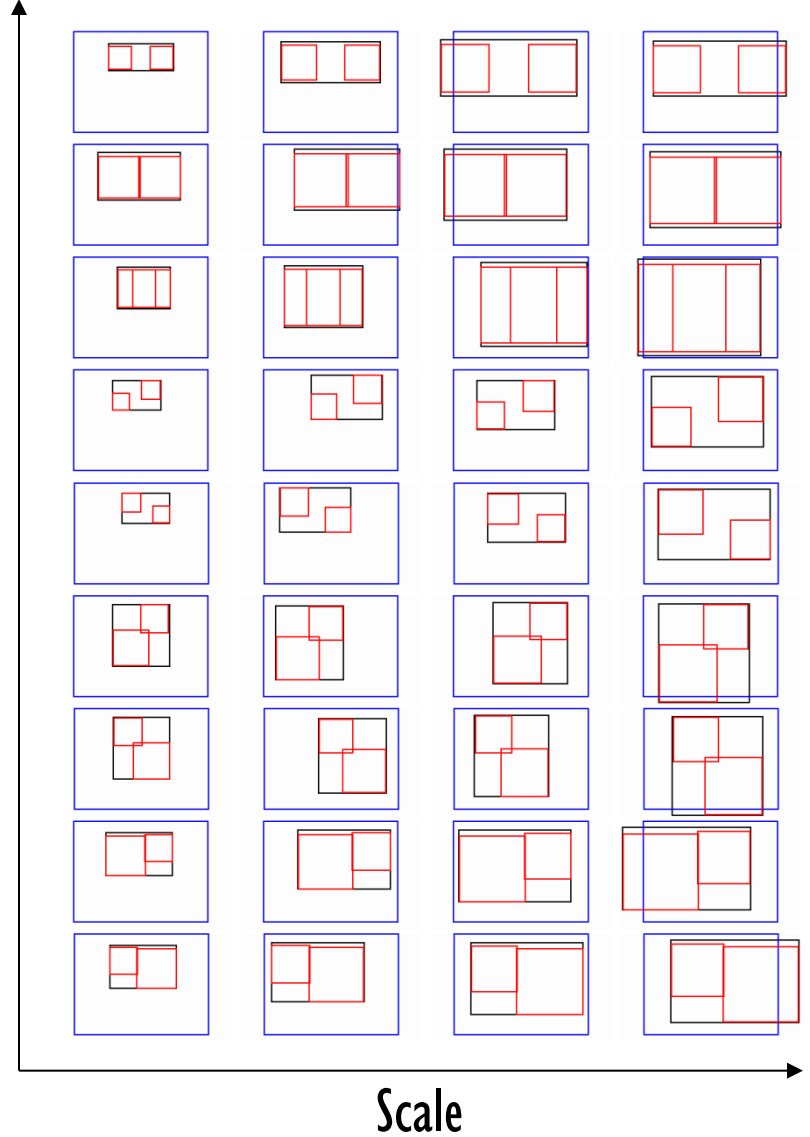
- Important contents mostly stay within frame.
- Particular camera angle is preferred.

Common Configuration of Upper Body

Two People

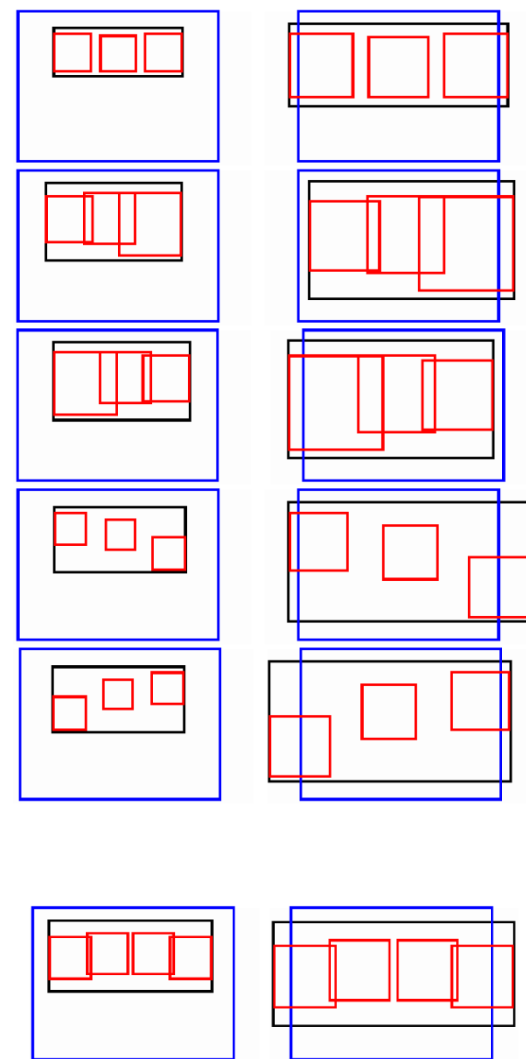


Location



Common Configuration of Upper Body

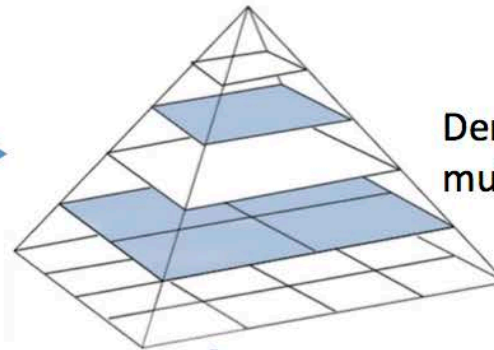
More than Three People



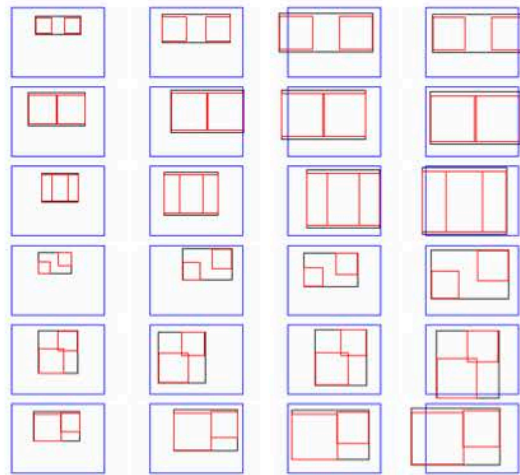
Input image



DPM



Dense detection scores at multiple location and scales



Learned configurations

Fast
inference

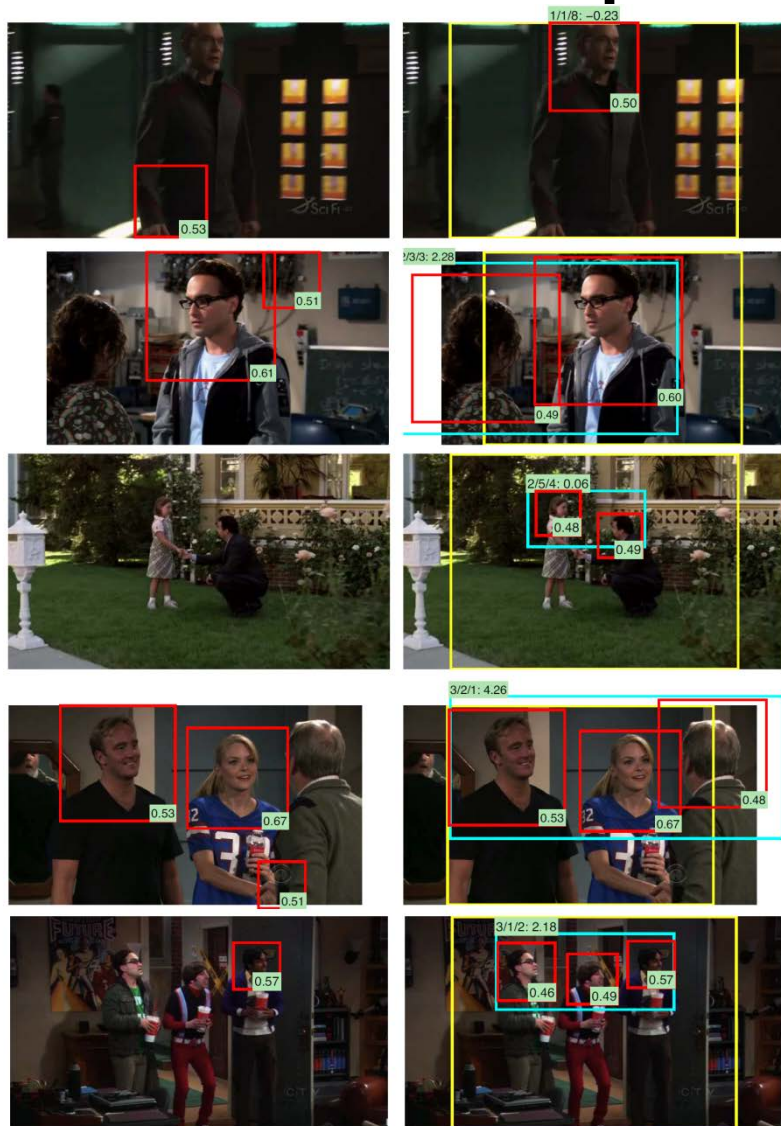
Output



Best configuration:
+ High unary scores
+ High similarity to
a common configuration

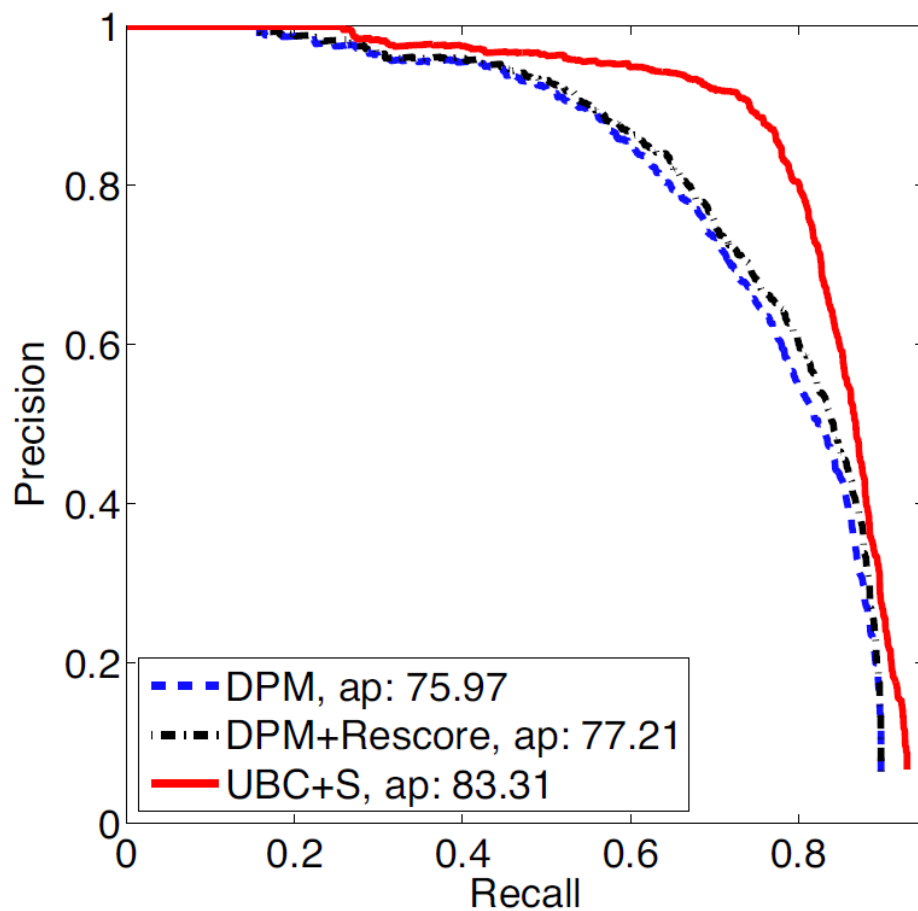
[Hoai CVPR14]

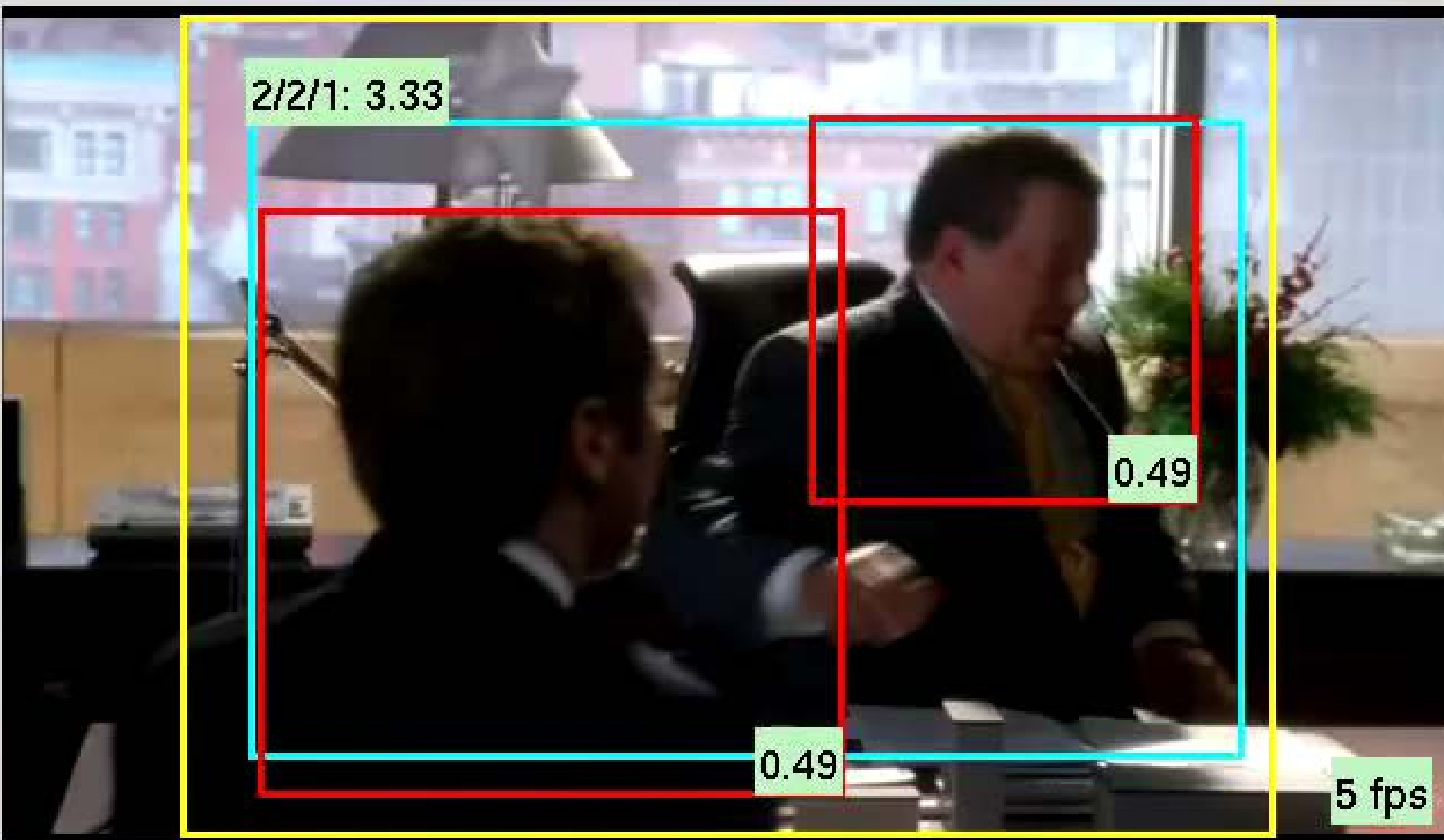
Comparison w/ DPM



DPM

Hoai, CVPR14





A video frame showing two men in an office. A yellow bounding box covers the entire scene. A cyan bounding box covers both men. Two red bounding boxes are also present: one around the man on the left and one around the man on the right. Each red box has a green label with the value '0.49'. The top-left corner has a green label with '2/2/1: 3.33'. The bottom-right corner has a green label with '5 fps'.

2/2/1: 3.33

0.49

0.49

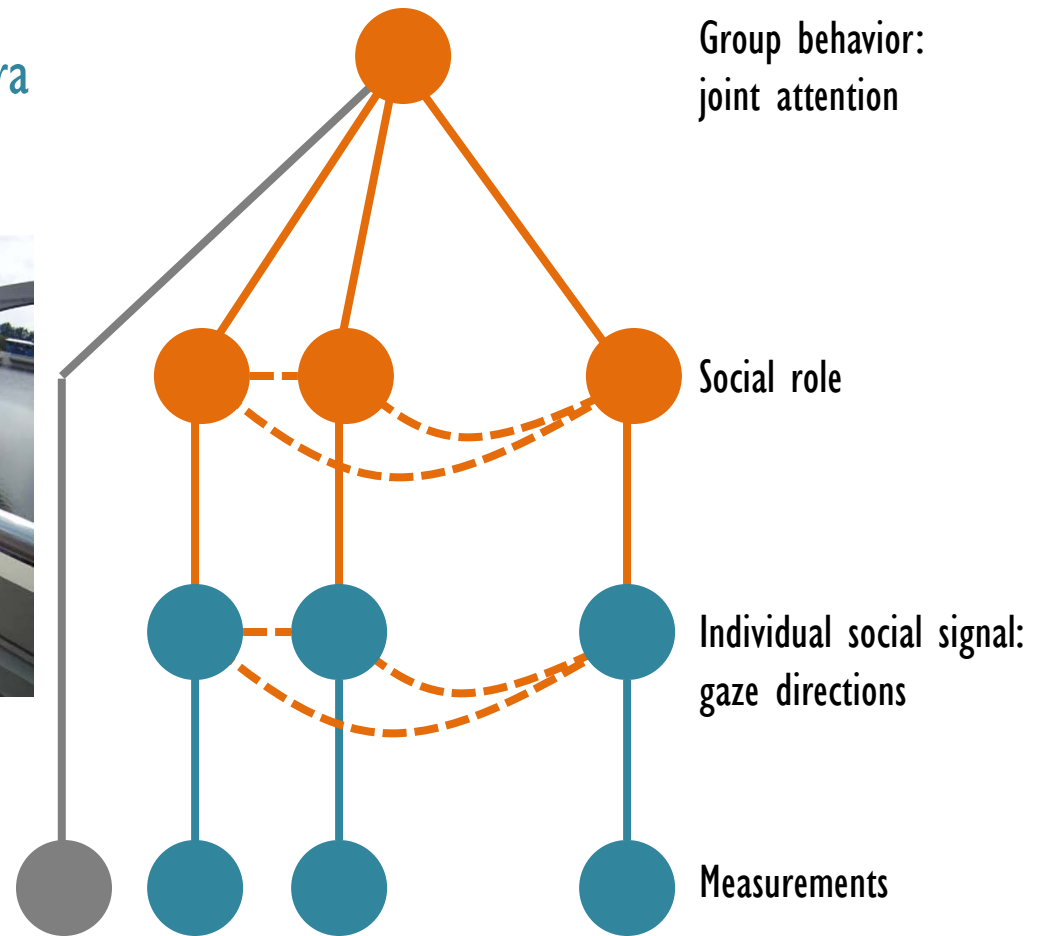
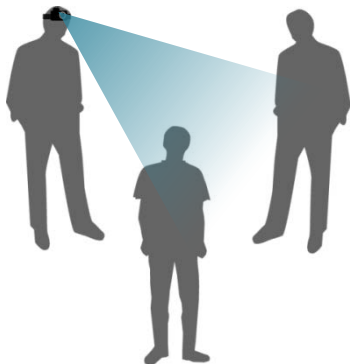
5 fps

Joint Attention

[Fathi CVPR12]

Input: images from a first person camera

Output: to detect social interactions



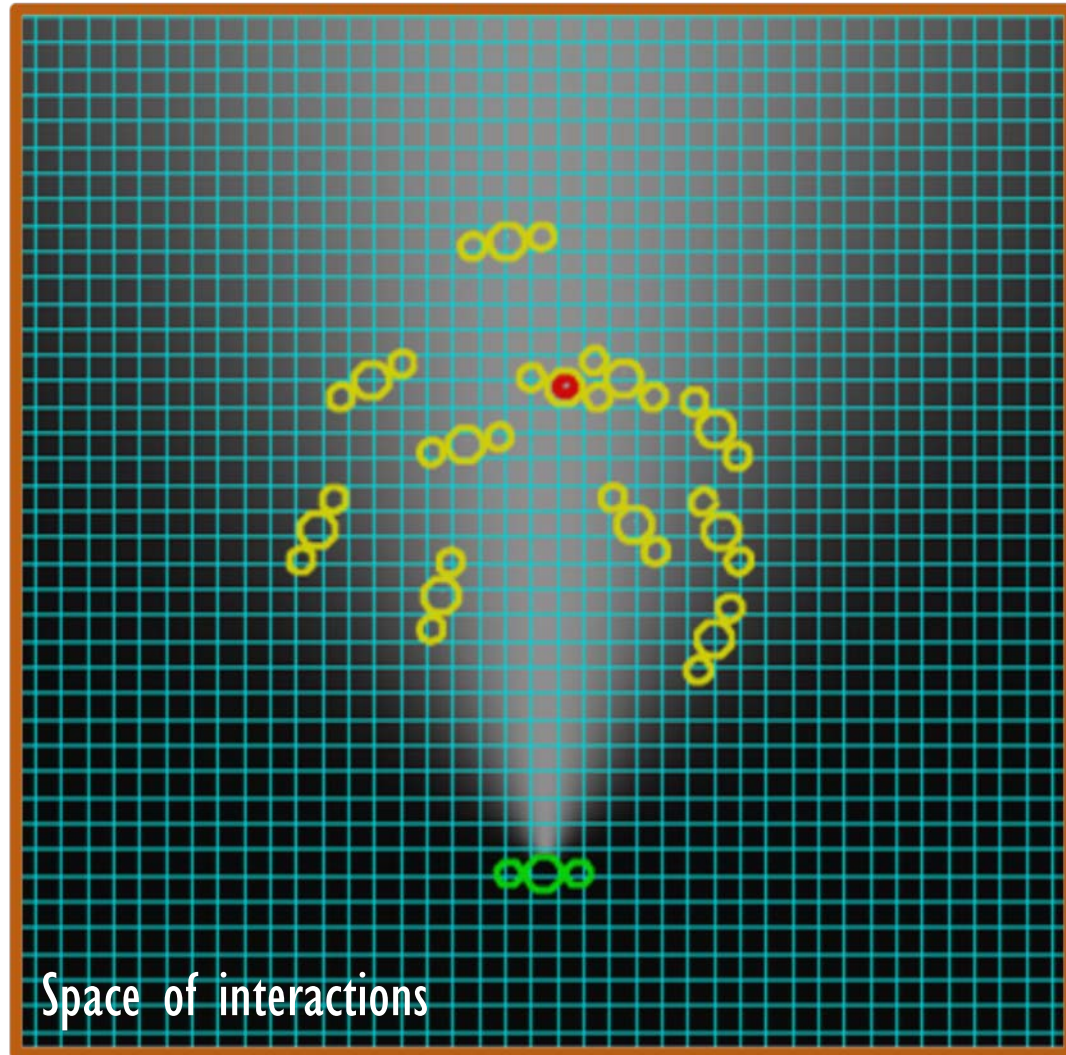
Where Do They Look?

[Fathi CVPR12]



MRF Modeling

[Fathi CVPR12]



Top view

MRF Modeling

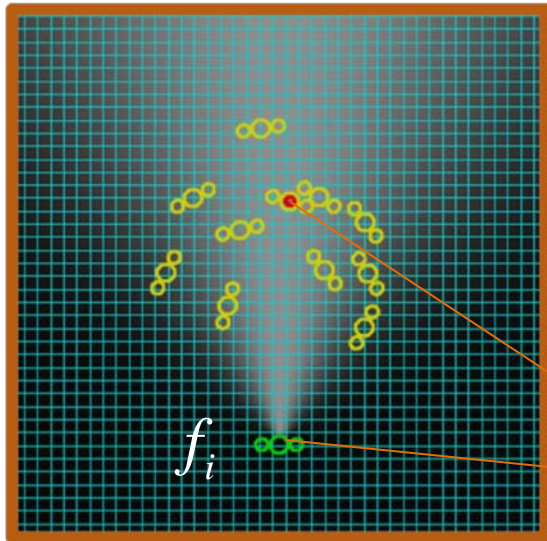
[Fathi CVPR12]

MRF modeling: unary potential

$$\phi_U(L_{f_i}, P_{f_1}, P_{f_2}, \dots, P_{f_N}) = \phi_1(L_{f_i}, P_{f_i}) \times \phi_2(L_{f_i}, P_{f_i}) \times \phi_3(L_{f_i}, P_{f_1}, \dots, P_{f_N})$$

L_{f_i} : location at which f_i is looking (label space)

P_{f_i} : position and orientation of f_i



- Head direction is aligned with the point of regard.

$$\phi_1(L_{f_i} = \ell, P_{f_i}) = \frac{1}{\sigma_1 \sqrt{2\pi}} \exp \left\{ -\frac{\|V_{f_i} - (\ell - T_{f_i})\|^2}{2\sigma_1^2} \right\}$$

- The point of regard cannot be the himself.

$$\phi_2(L_{f_i} = \ell, P_{f_i}) = \frac{1}{1 + \exp \{ -(c_2 \cdot \|\ell - P_{f_i}\|) \}}$$

- The point of regard is likely to be a face.

$$\phi_3(L_{f_i} = \ell, P_{f_1}, \dots, P_{f_N}) = \begin{cases} c_3 & \ell = P_{f_j} \forall j \neq i \\ 1 & \text{otherwise} \end{cases}$$

MRF Modeling

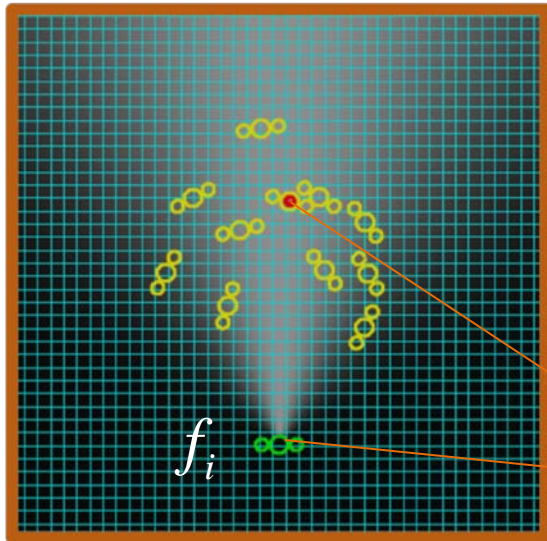
[Fathi CVPR12]

MRF modeling: unary potential

$$\phi_U(L_{f_i}, P_{f_1}, P_{f_2}, \dots, P_{f_N}) = \phi_1(L_{f_i}, P_{f_i}) \times \phi_2(L_{f_i}, P_{f_i}) \times \phi_3(L_{f_i}, P_{f_1}, \dots, P_{f_N})$$

L_{f_i} : location at which f_i is looking (label space)

P_{f_i} : position and orientation of f_i



MRF modeling: binary potential

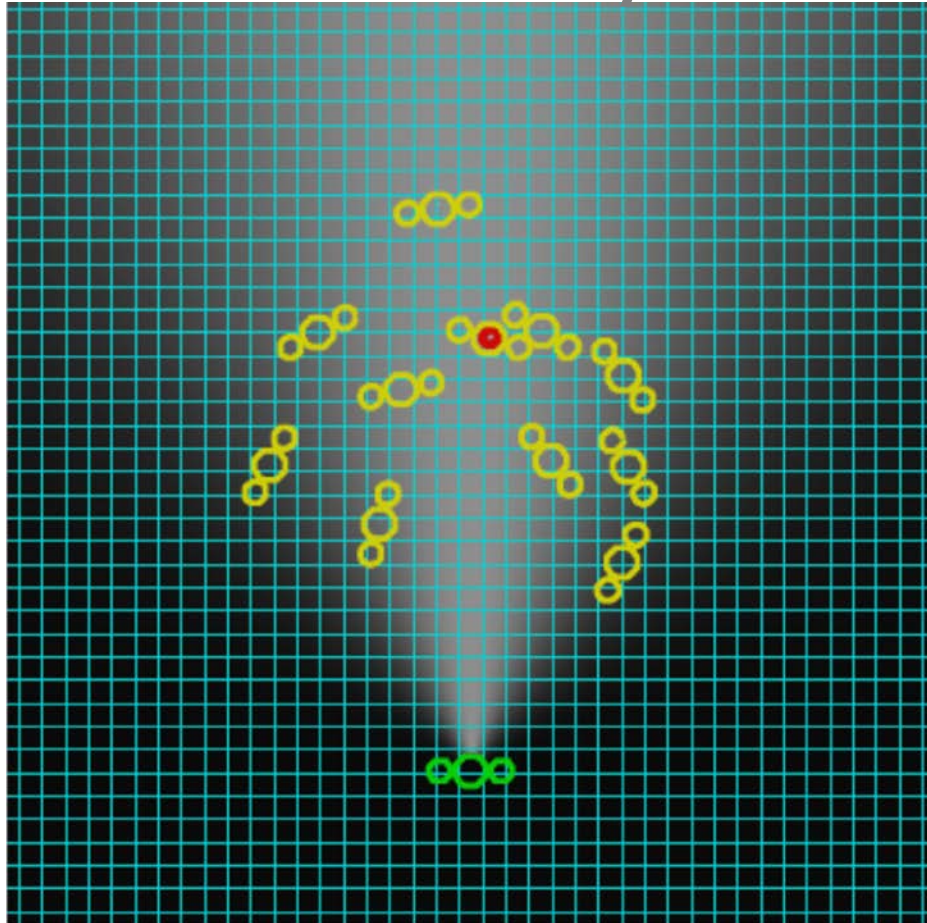
$$\phi_B(L_{f_i} = \ell_1, L_{f_j} = \ell_2) = \begin{cases} c_B & \text{if } (\ell_1 = \ell_2) \\ 1 - c_B & \text{if } (\ell_1 \neq \ell_2) \end{cases}$$

People engage joint attention.

MRF Inference

Where Do They Look?

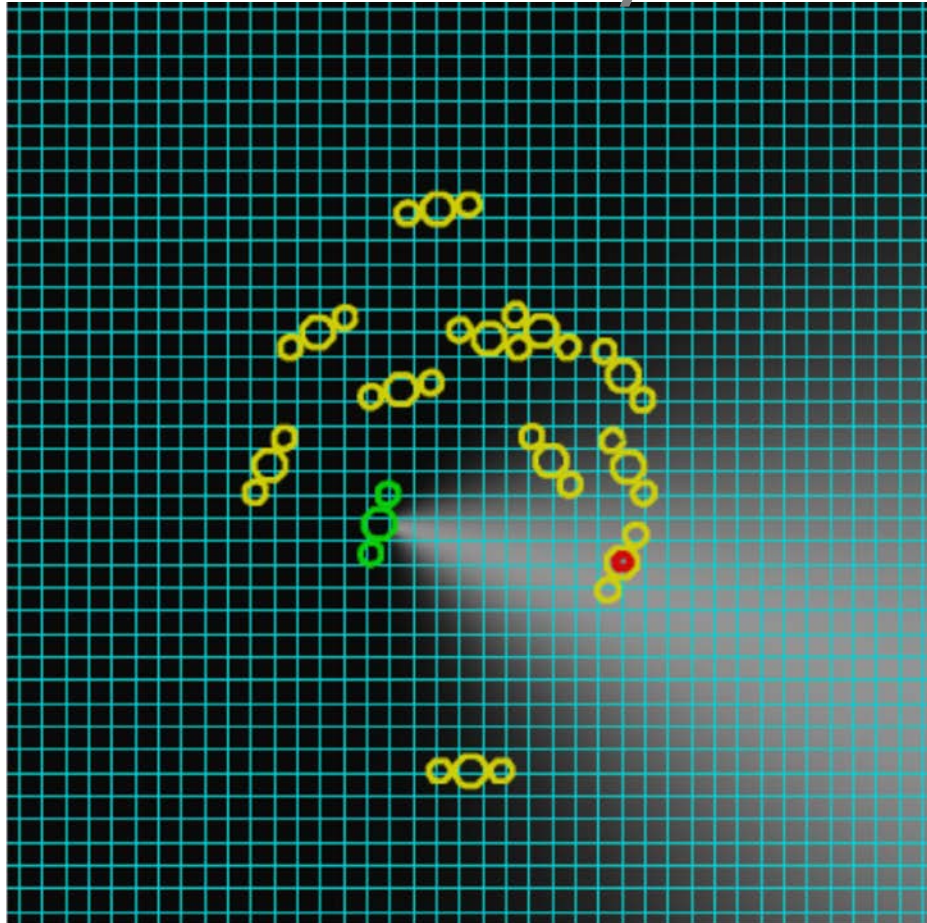
[Fathi CVPR12]



MRF Inference

Where Do They Look?

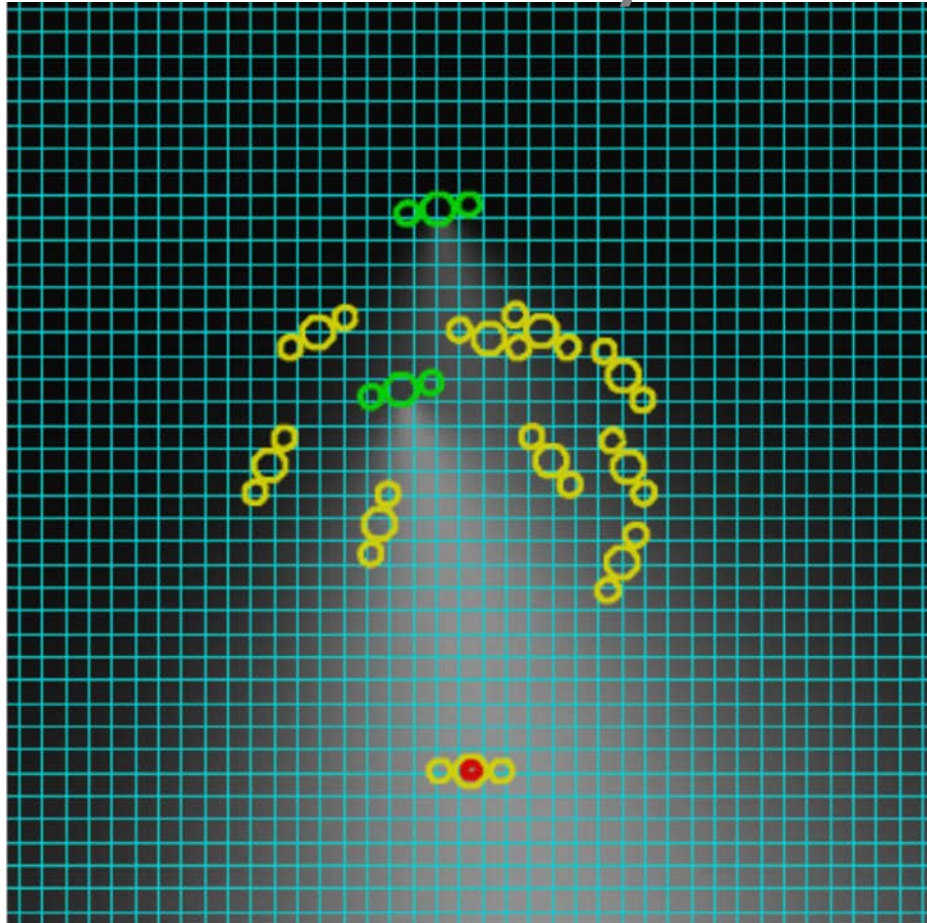
[Fathi CVPR12]



MRF Inference

Where Do They Look?

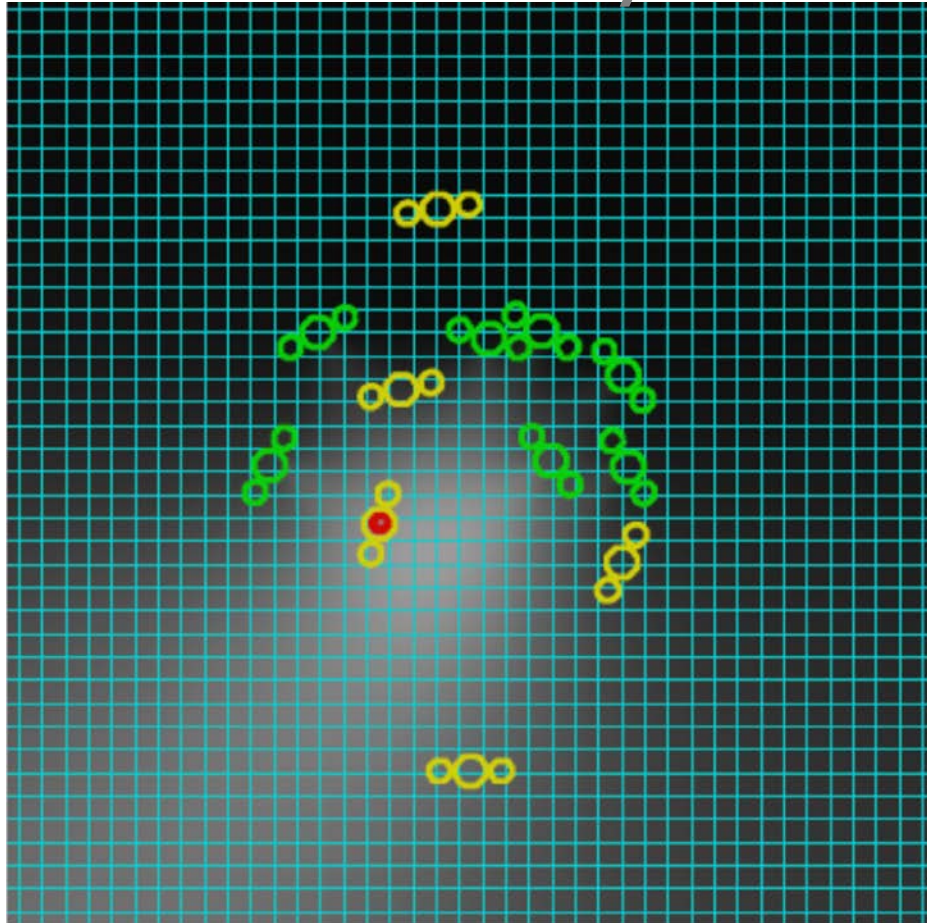
[Fathi CVPR12]



MRF Inference

Where Do They Look?

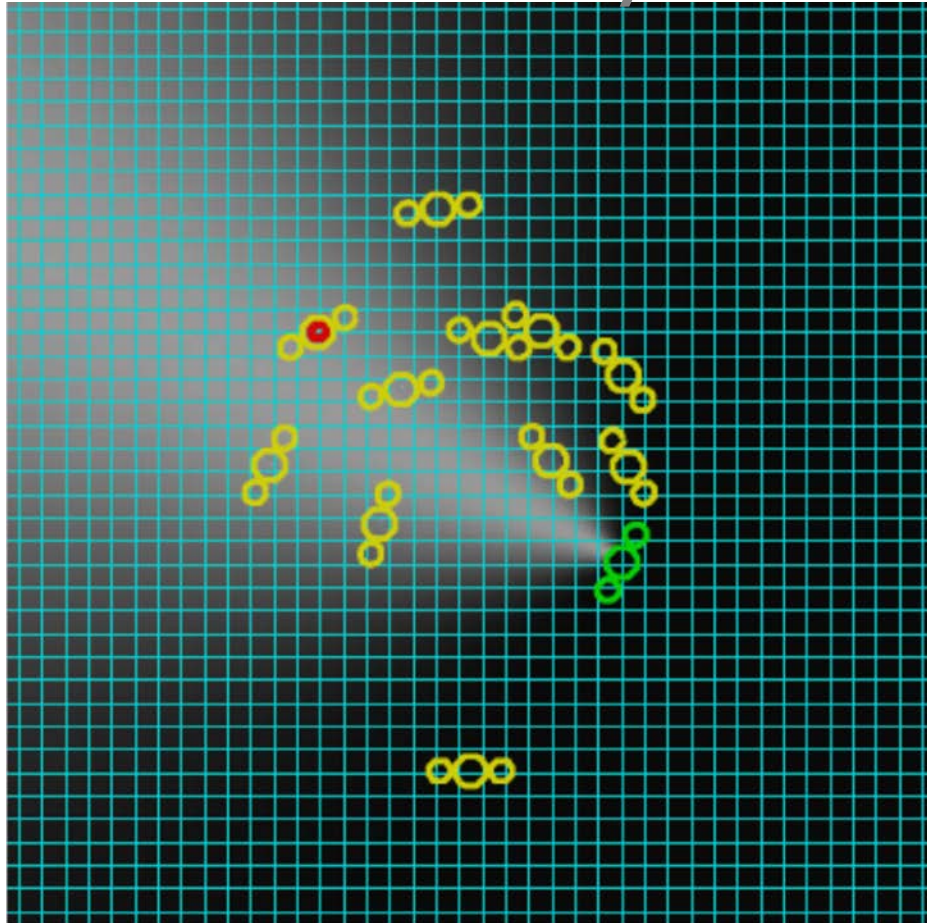
[Fathi CVPR12]



MRF Inference

Where Do They Look?

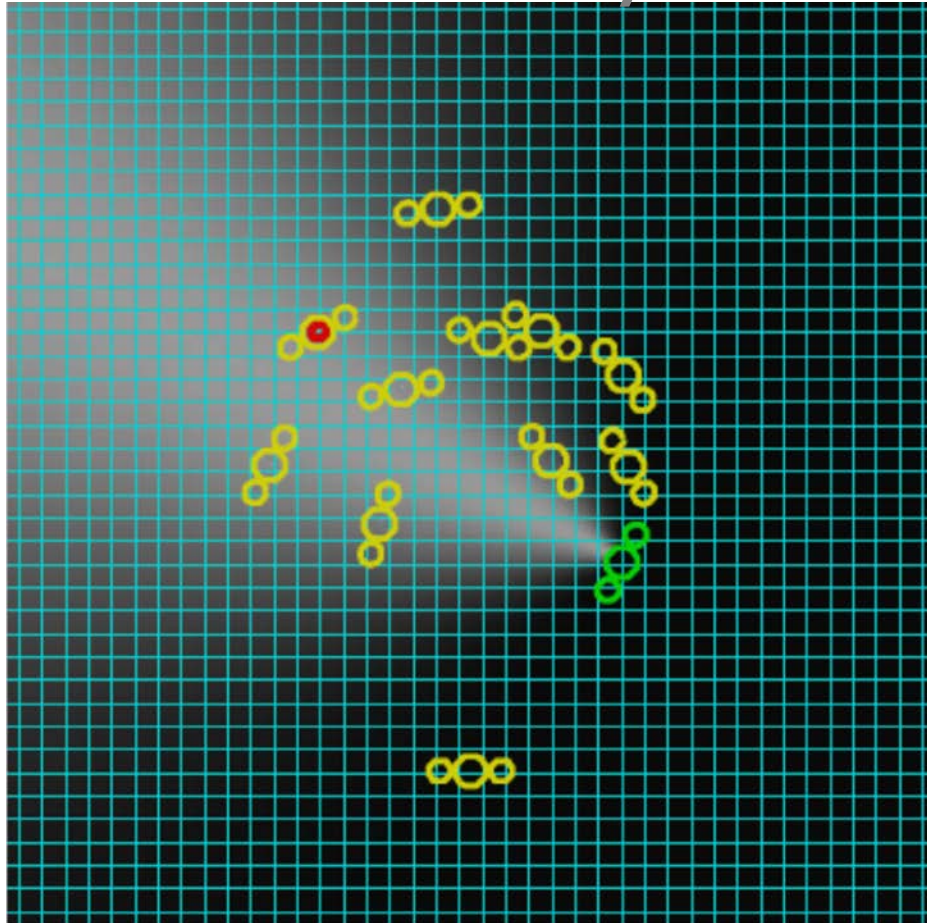
[Fathi CVPR12]



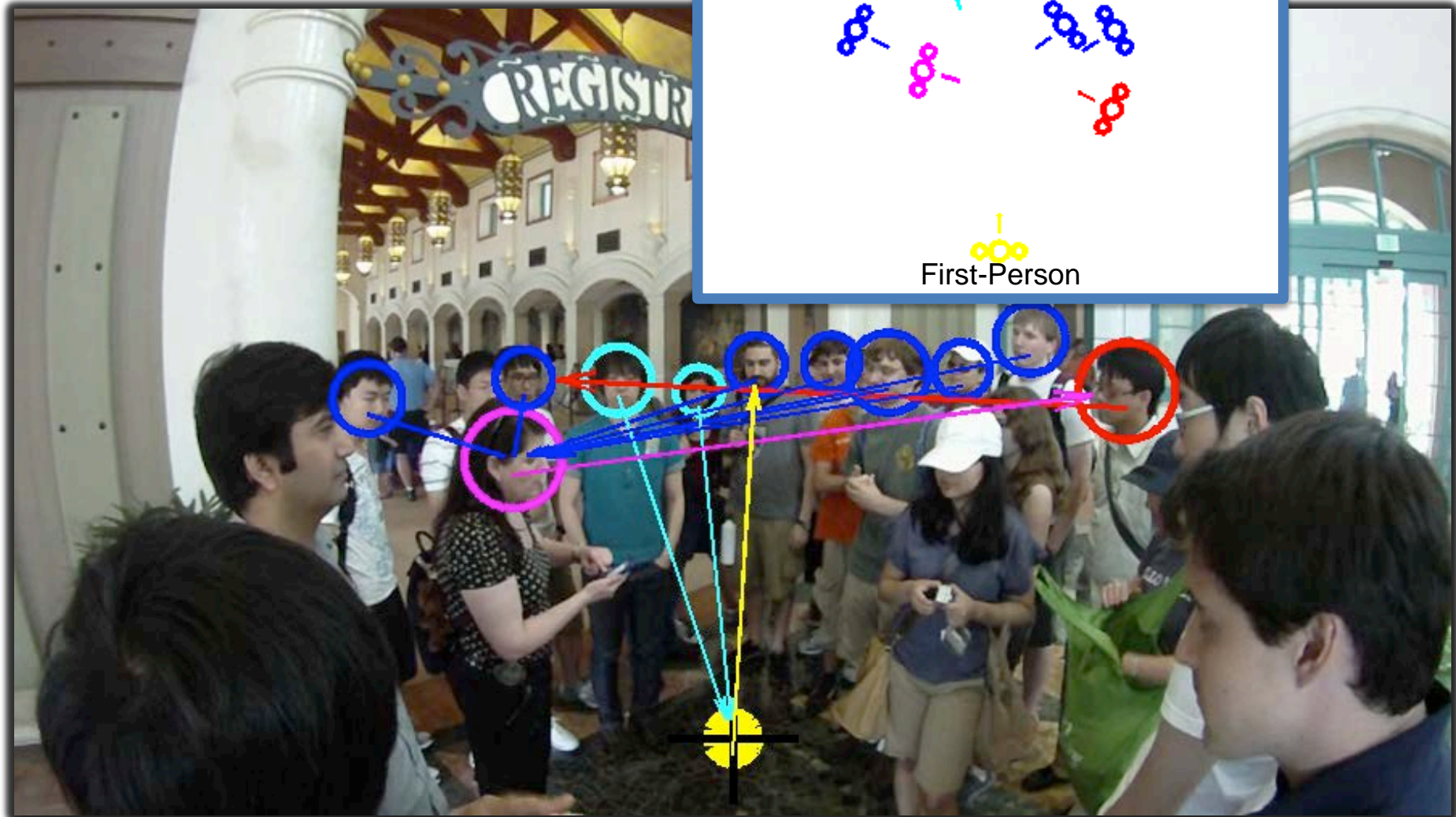
MRF Inference

Where Do They Look?

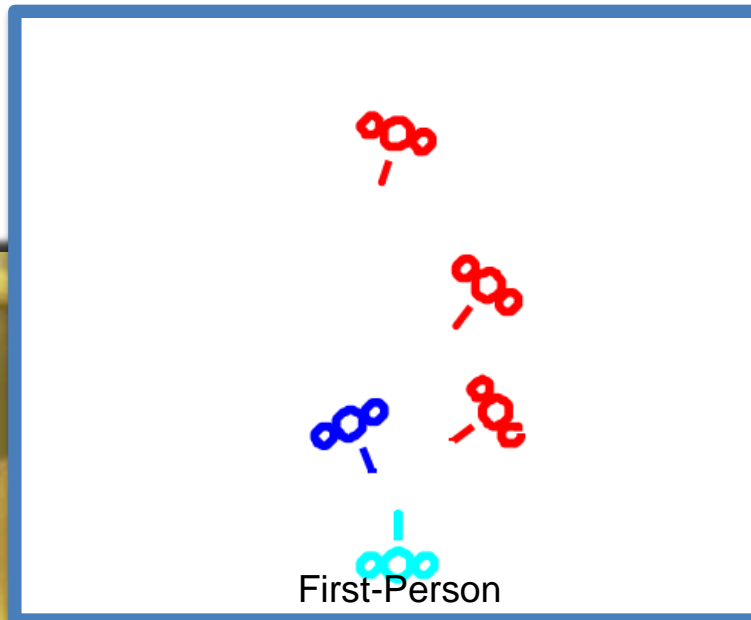
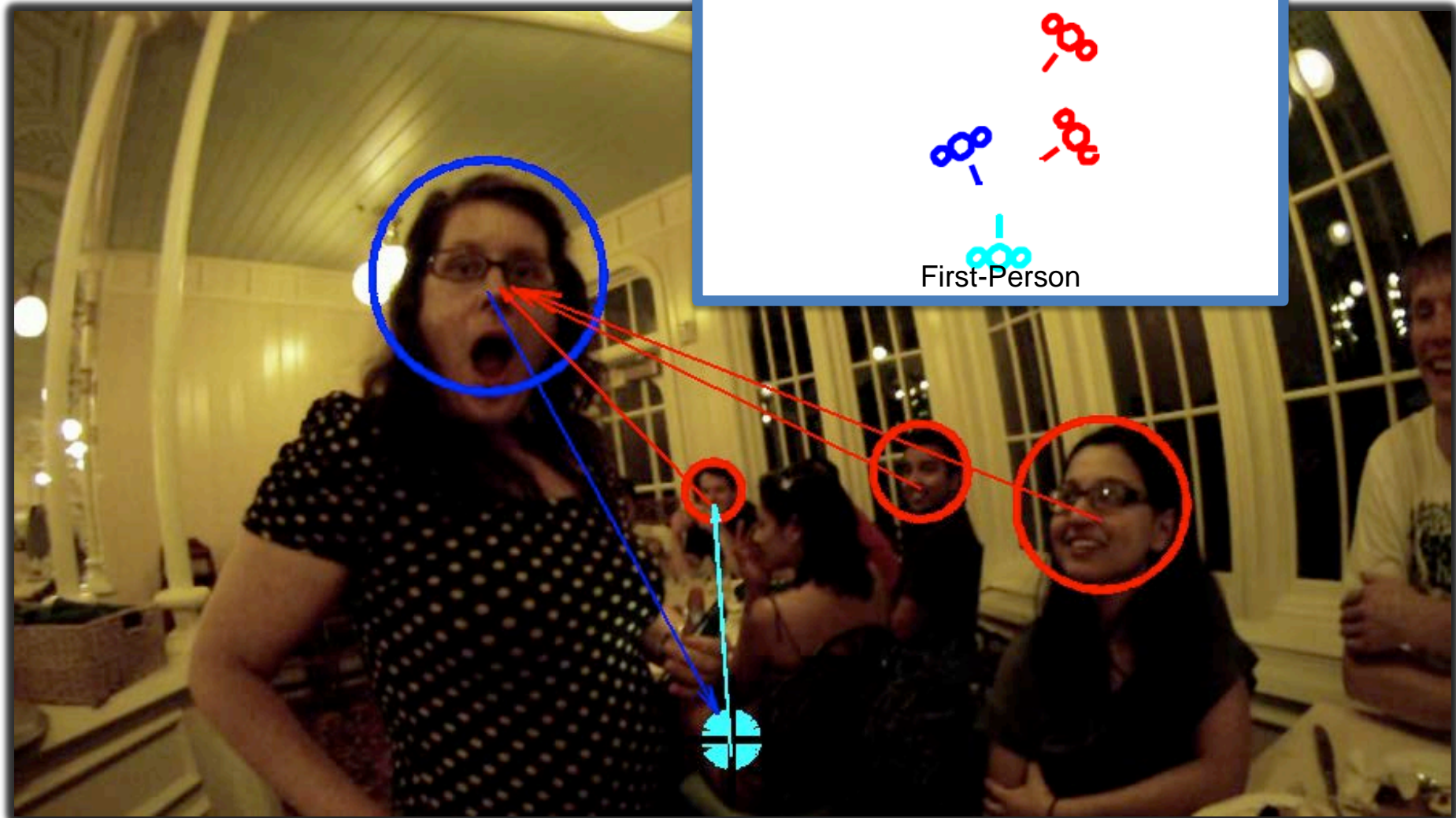
[Fathi CVPR12]



[Fathi CVPR12]



[Fathi CVPR12]



Detection of Social Interaction



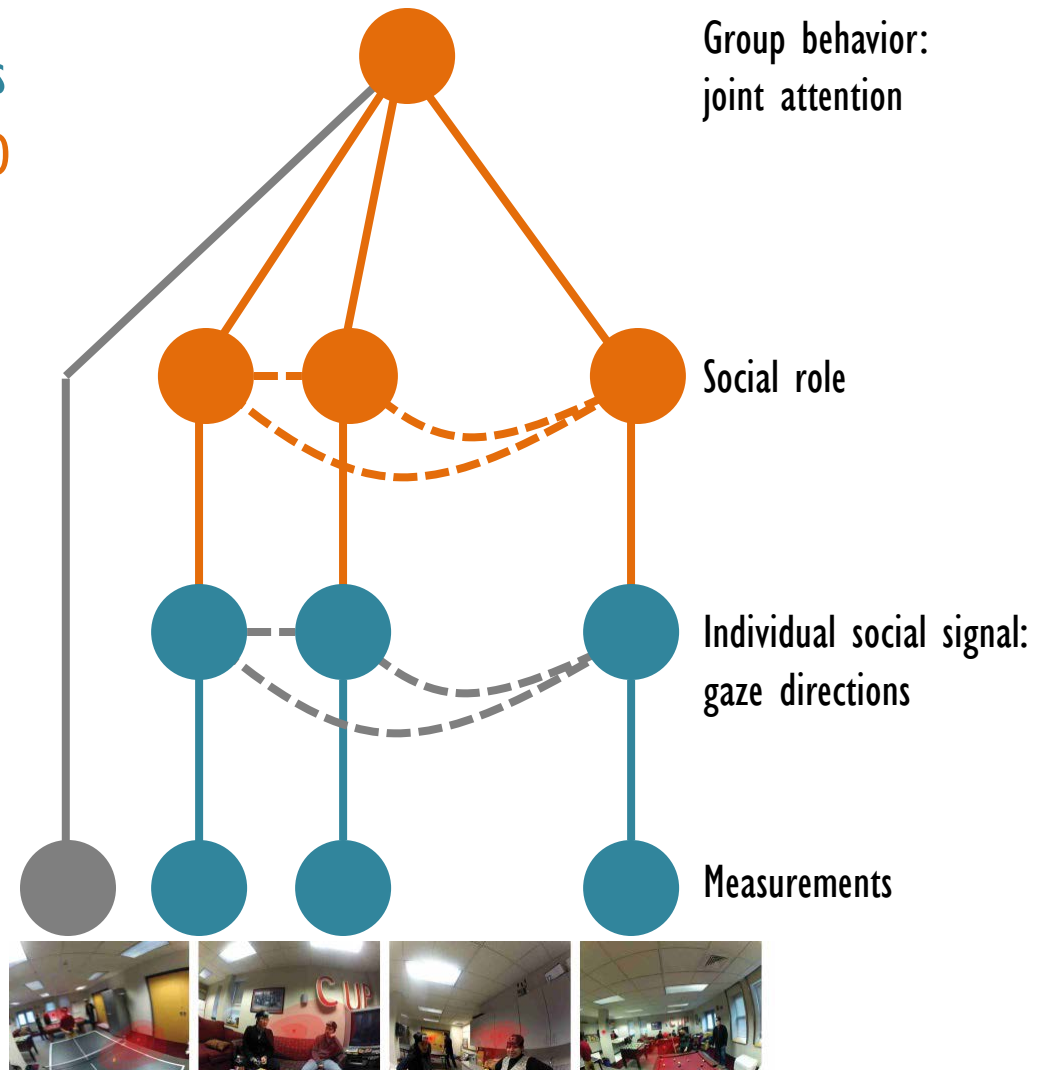
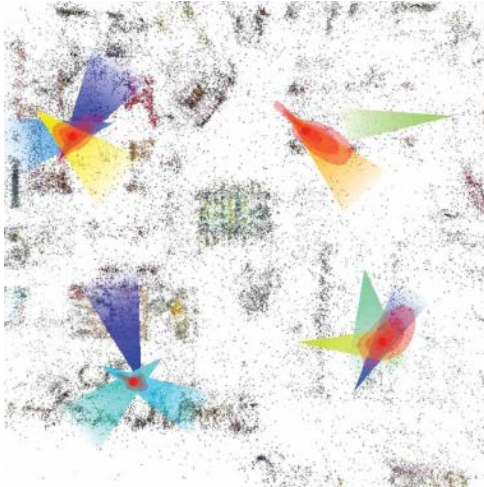
[Fathi CVPR12]

Joint Attention

[Park NIPS12]

Input: images from first person cameras

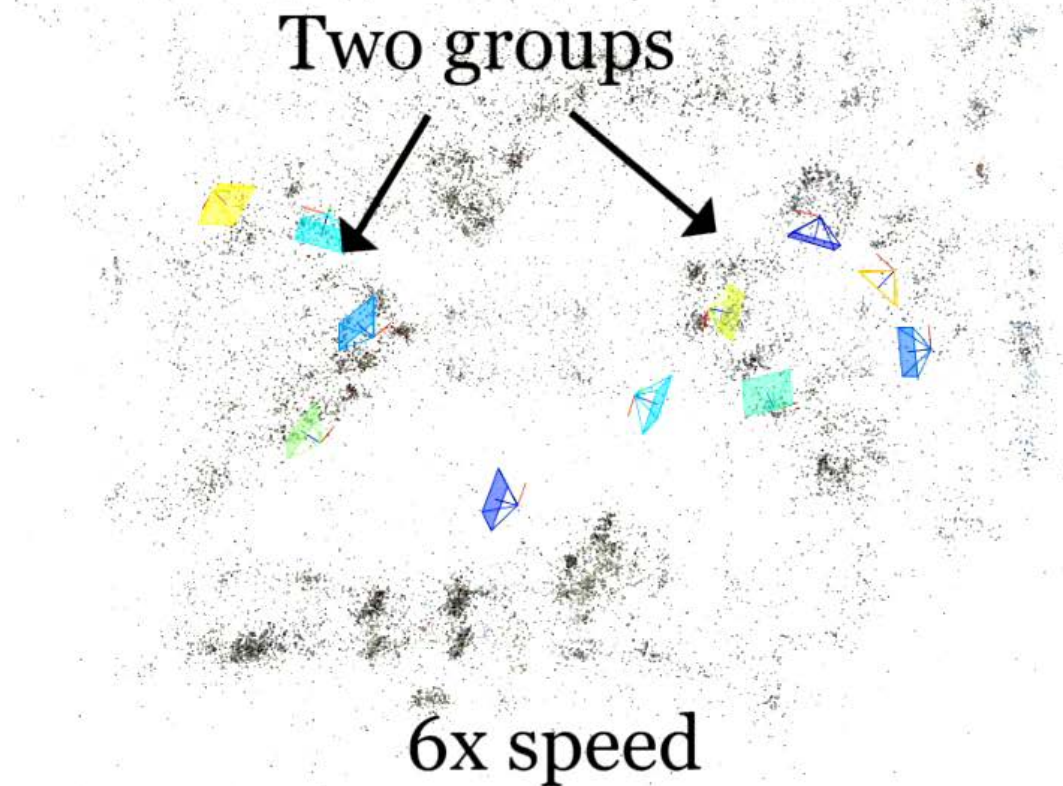
Output: to localize joint attention in 3D



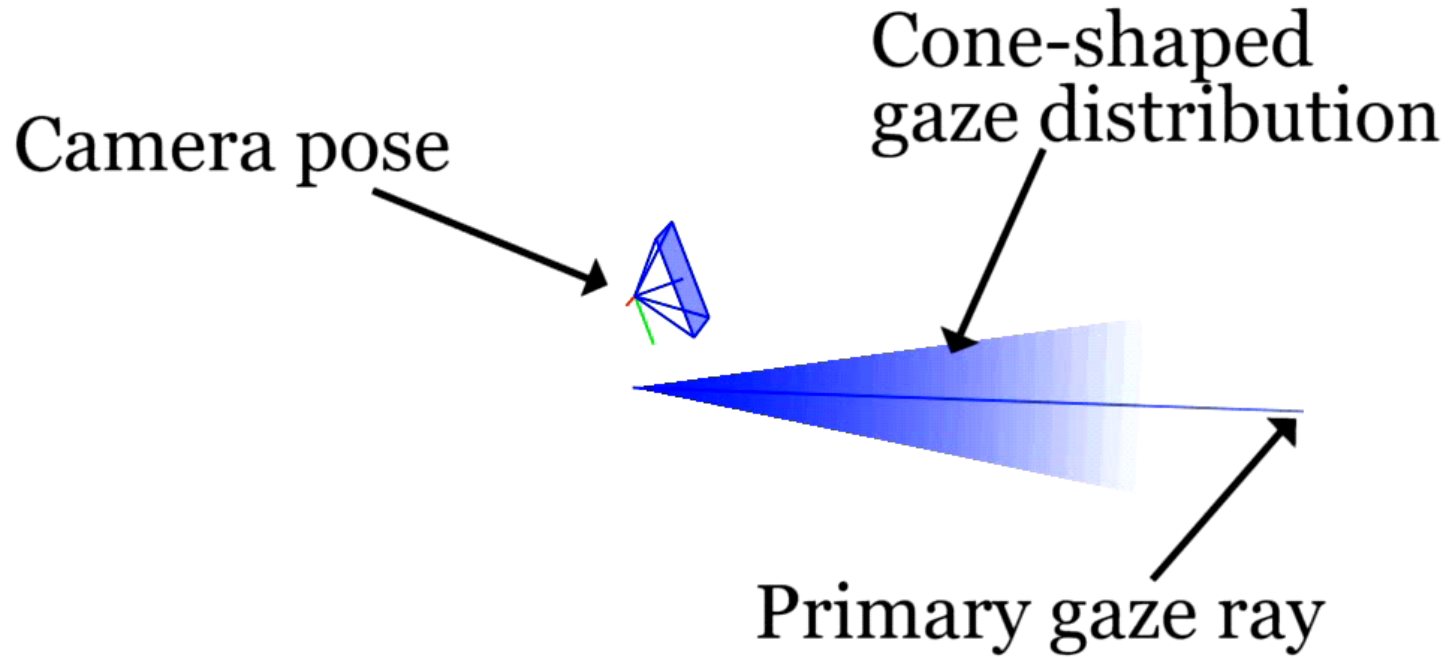
Input Video: Meeting Scene



3D Camera Pose Estimation (Structure from motion)

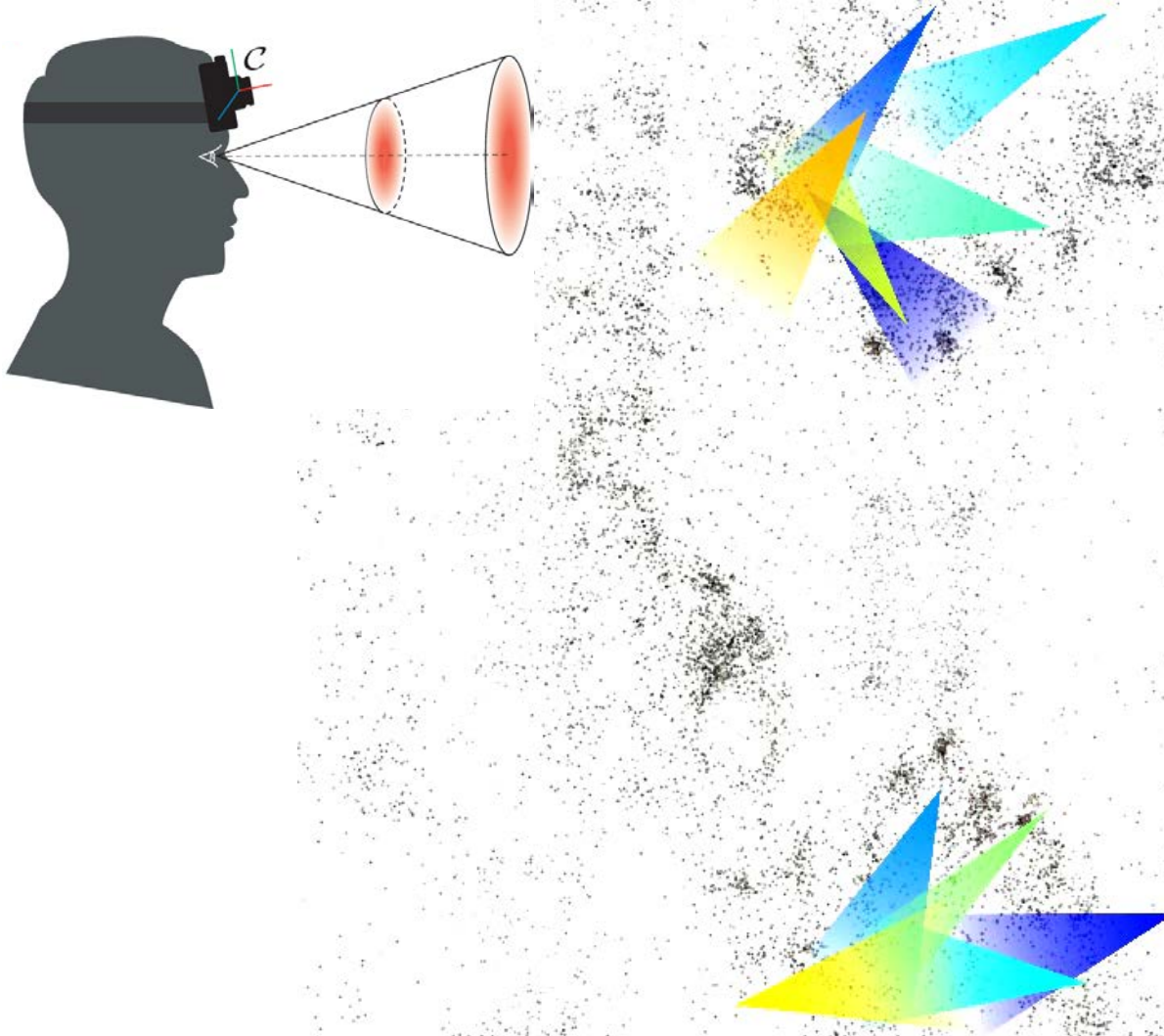


Gaze Ray Calibration

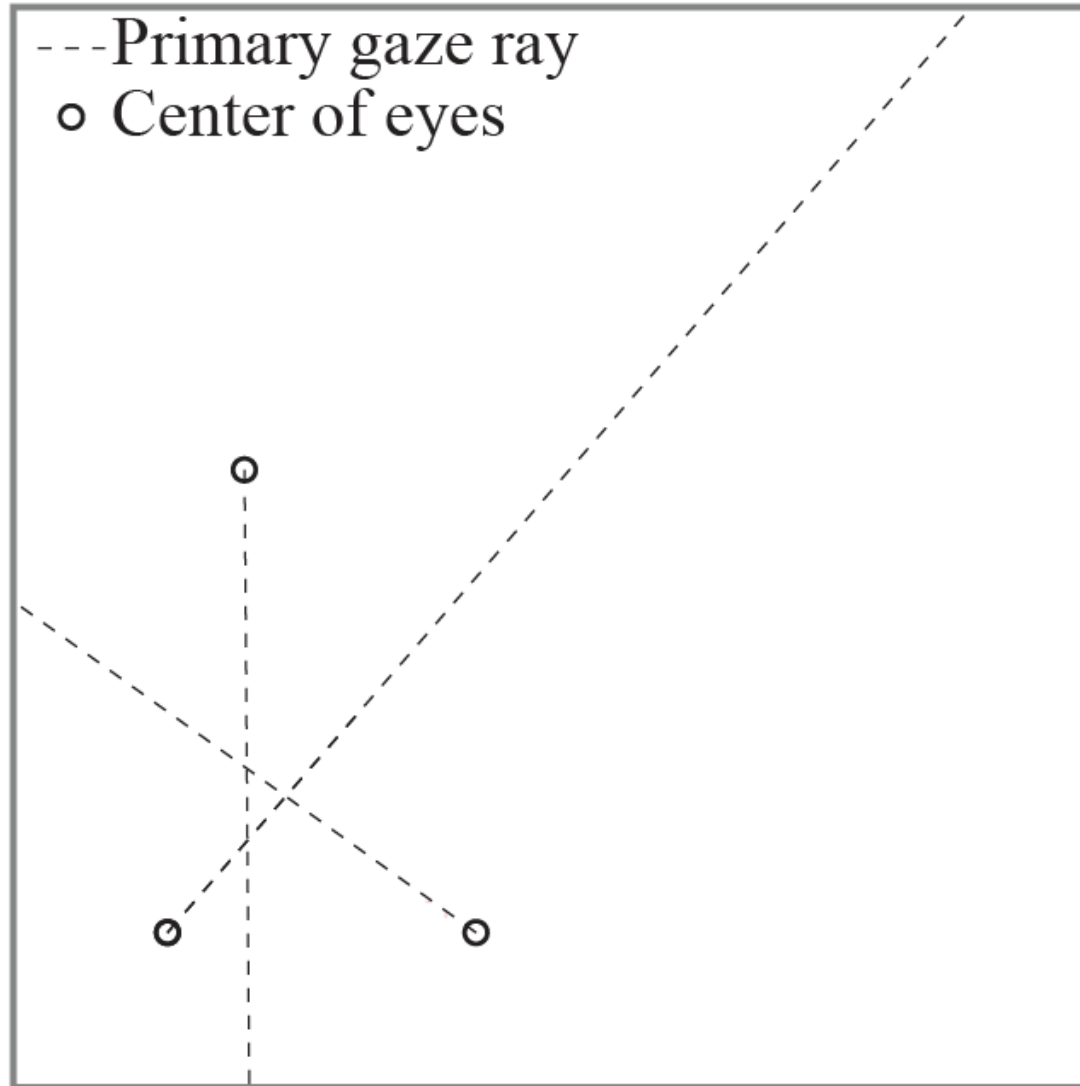


Primary gaze ray with respect to the camera pose

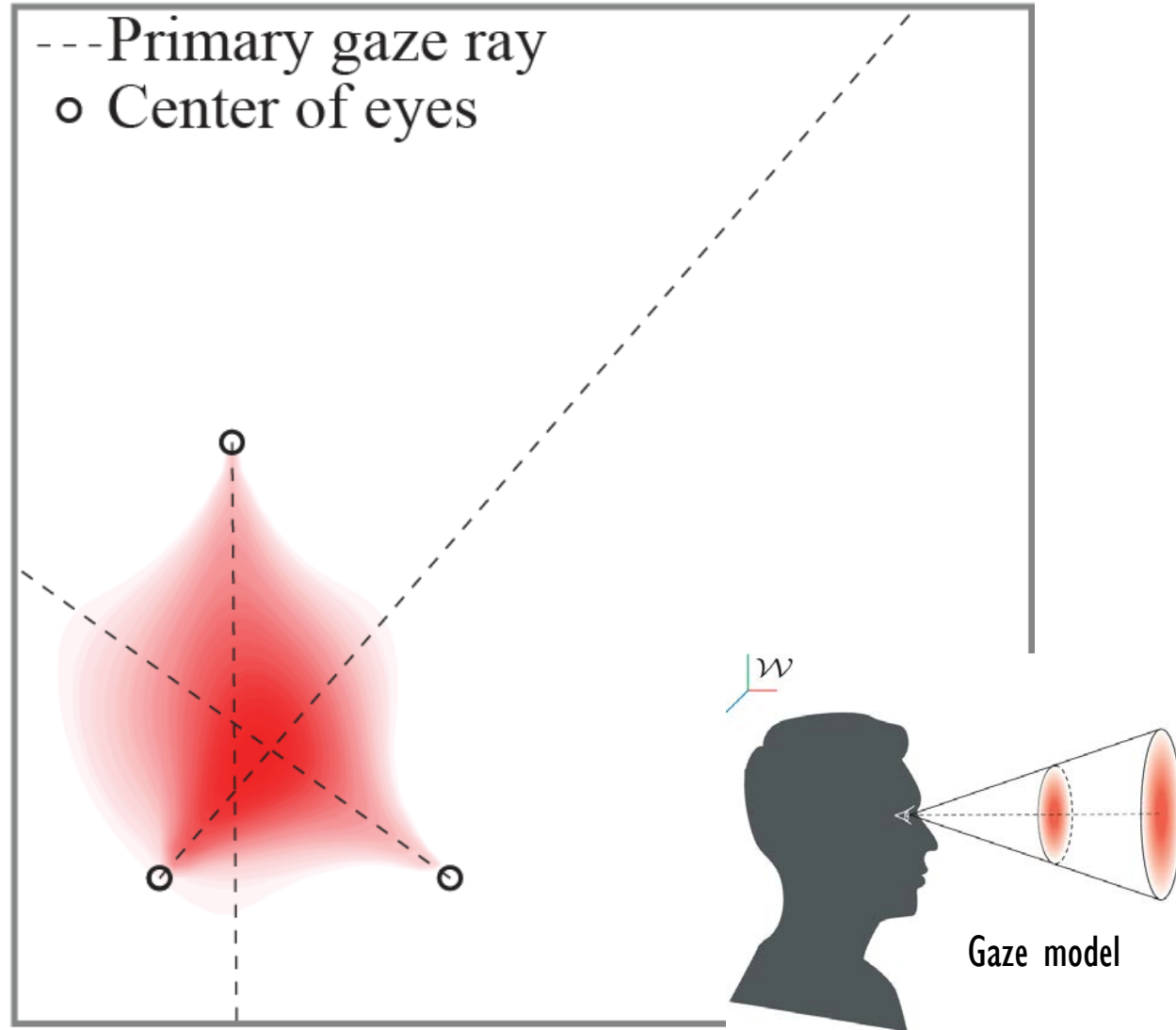
3D Gaze Registration



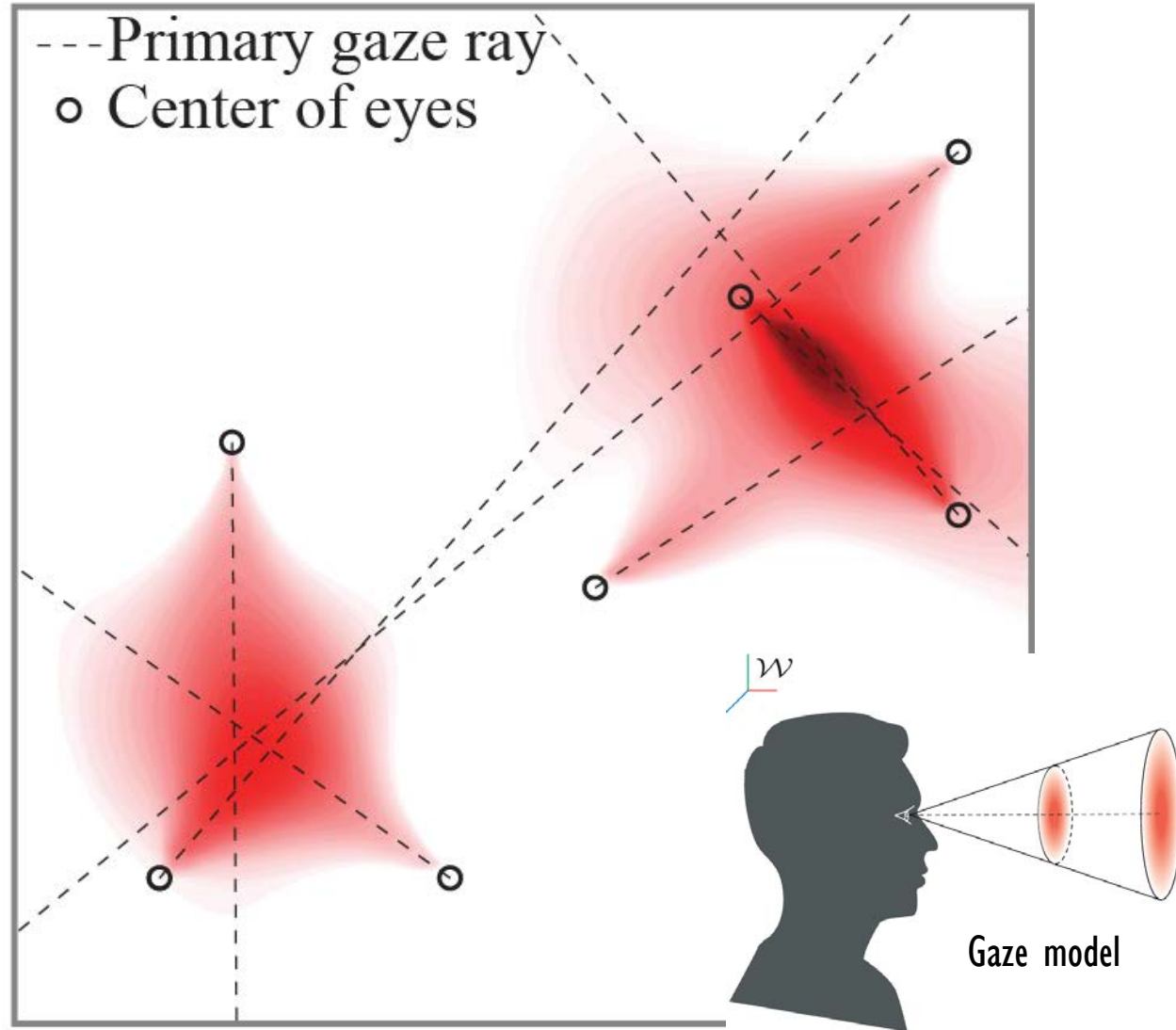
3D Gaze Registration



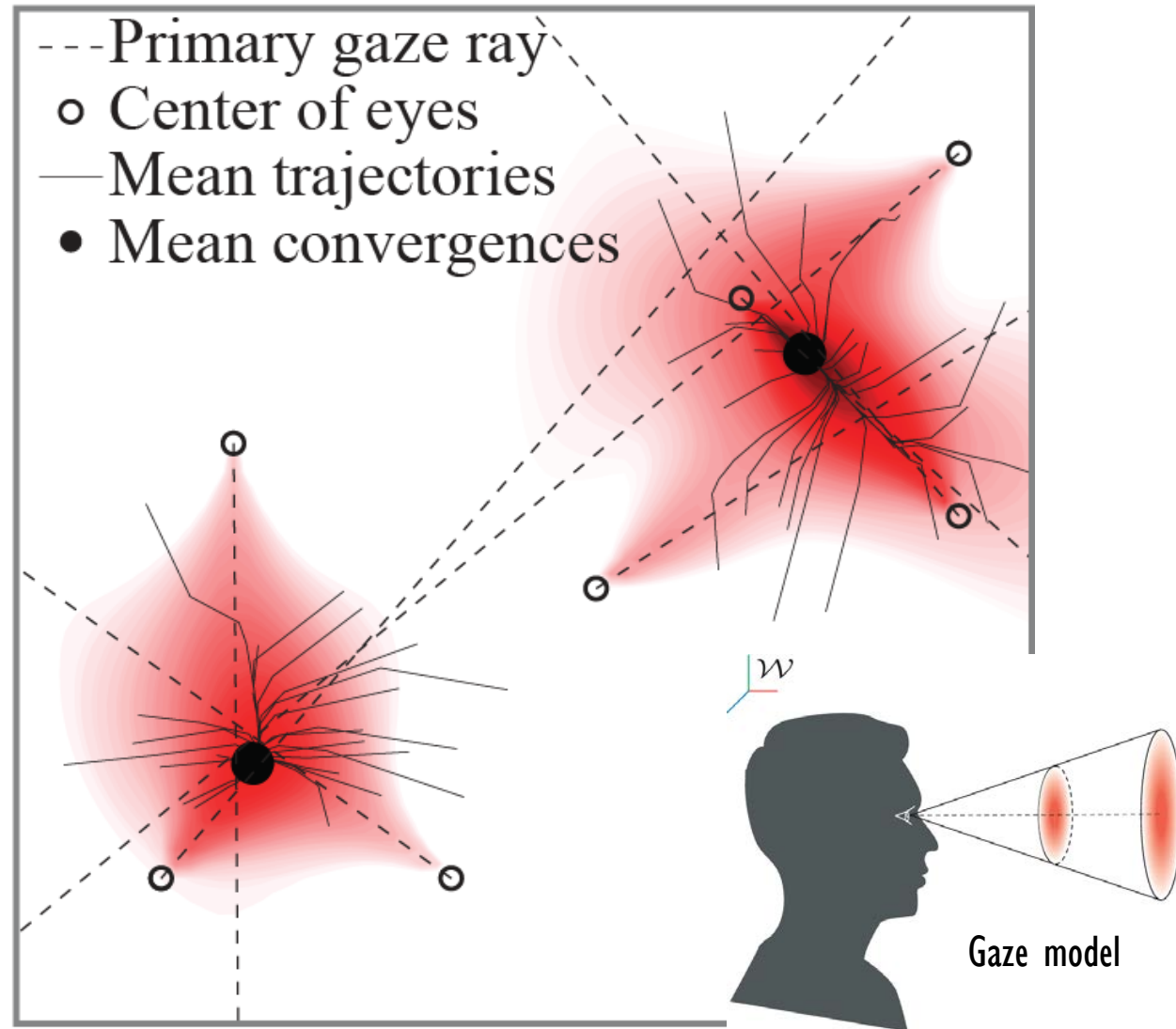
Social Saliency



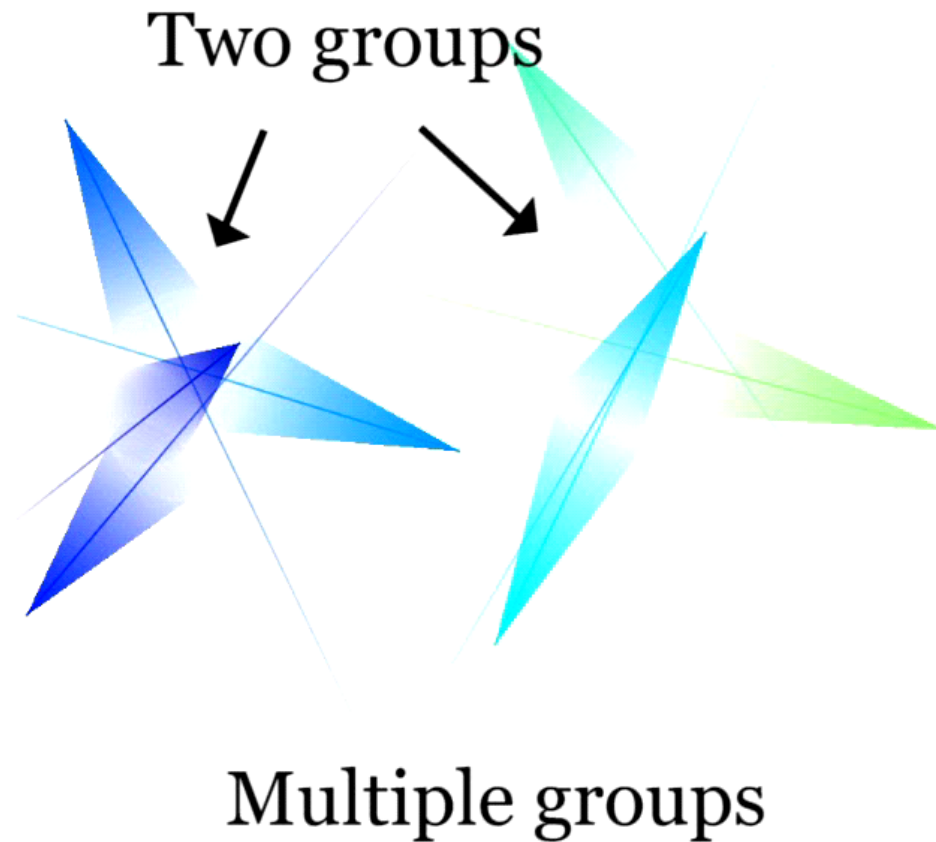
Social Saliency



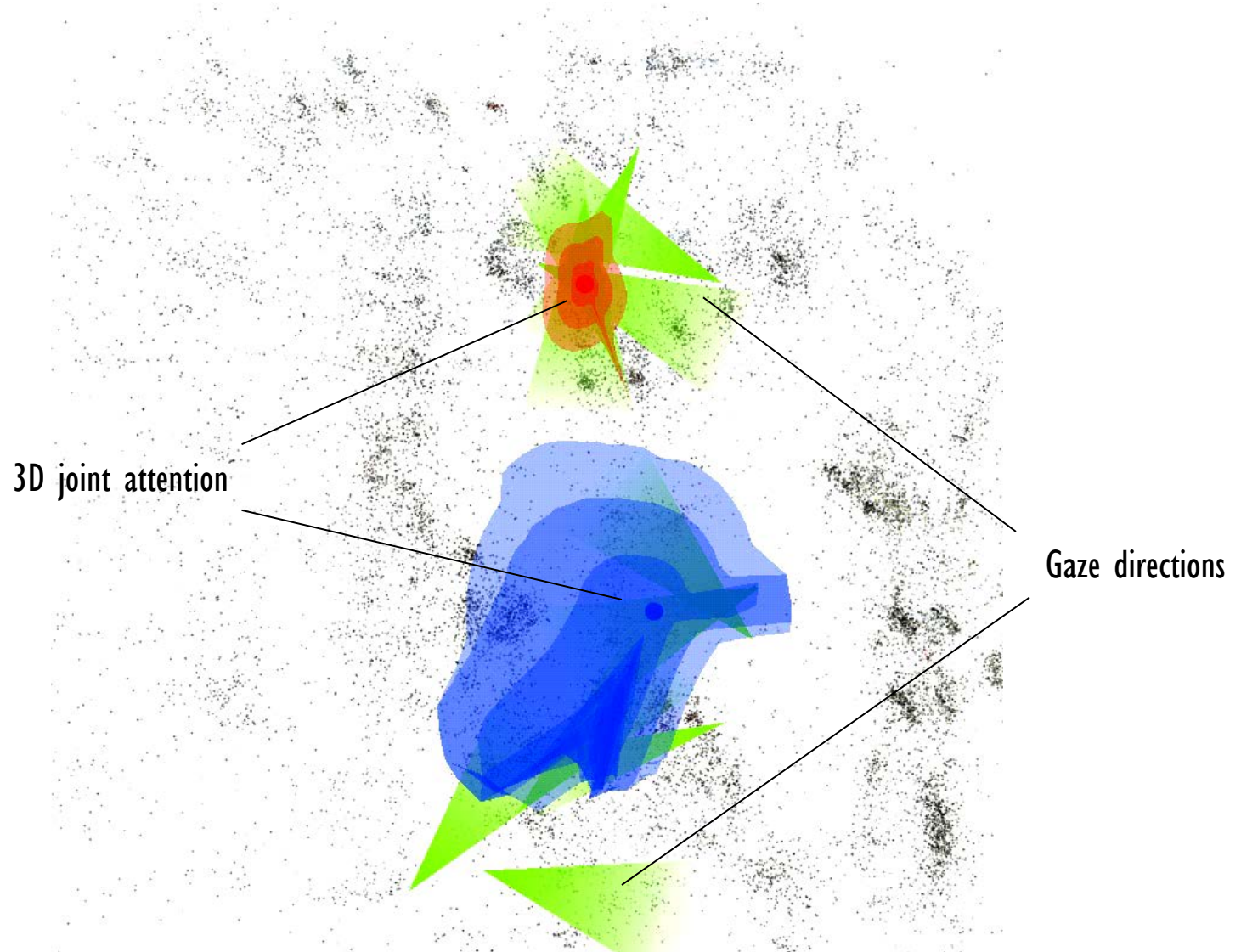
3D Joint Attention via Mode-seeking



Mode-seeking: Gaze Concurrences

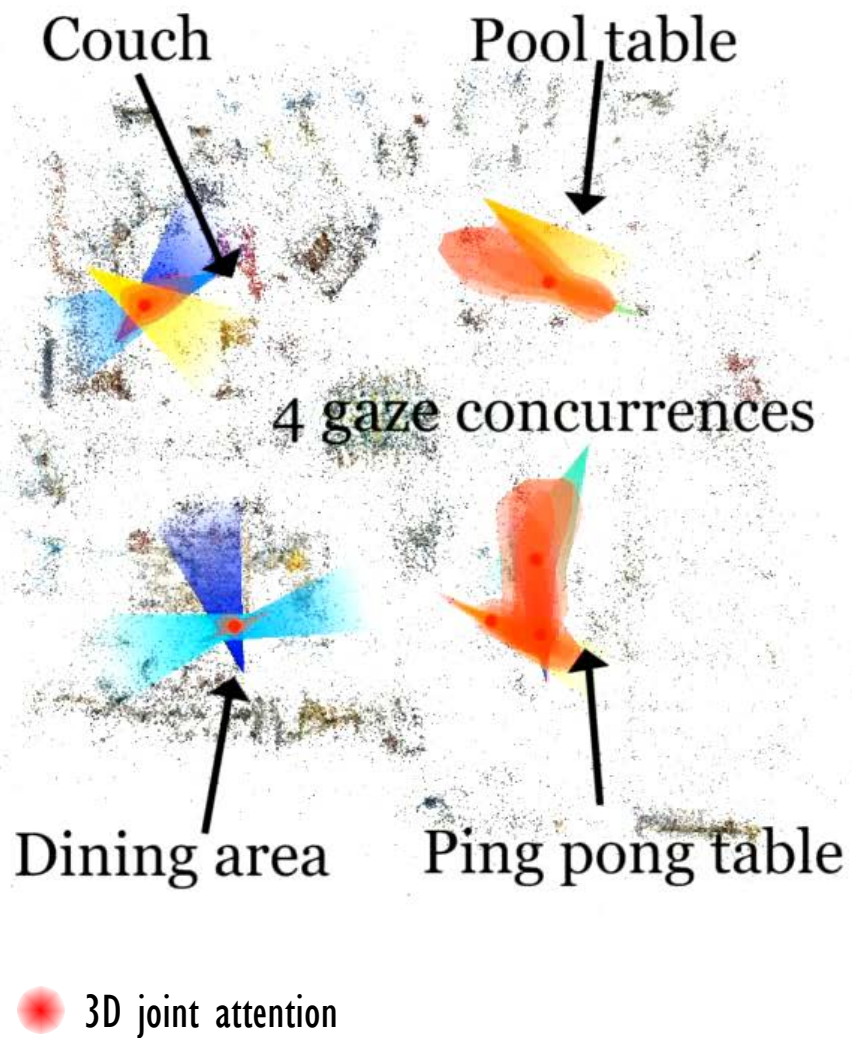


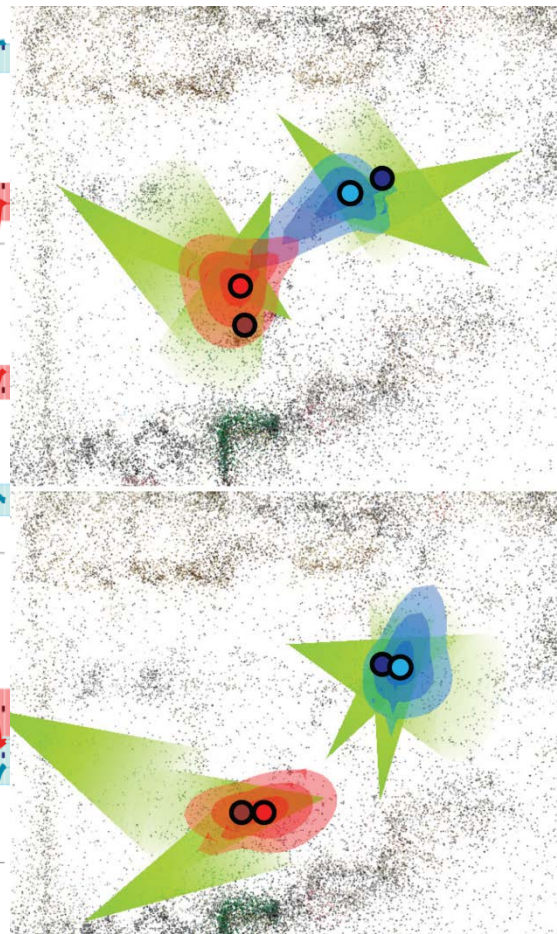
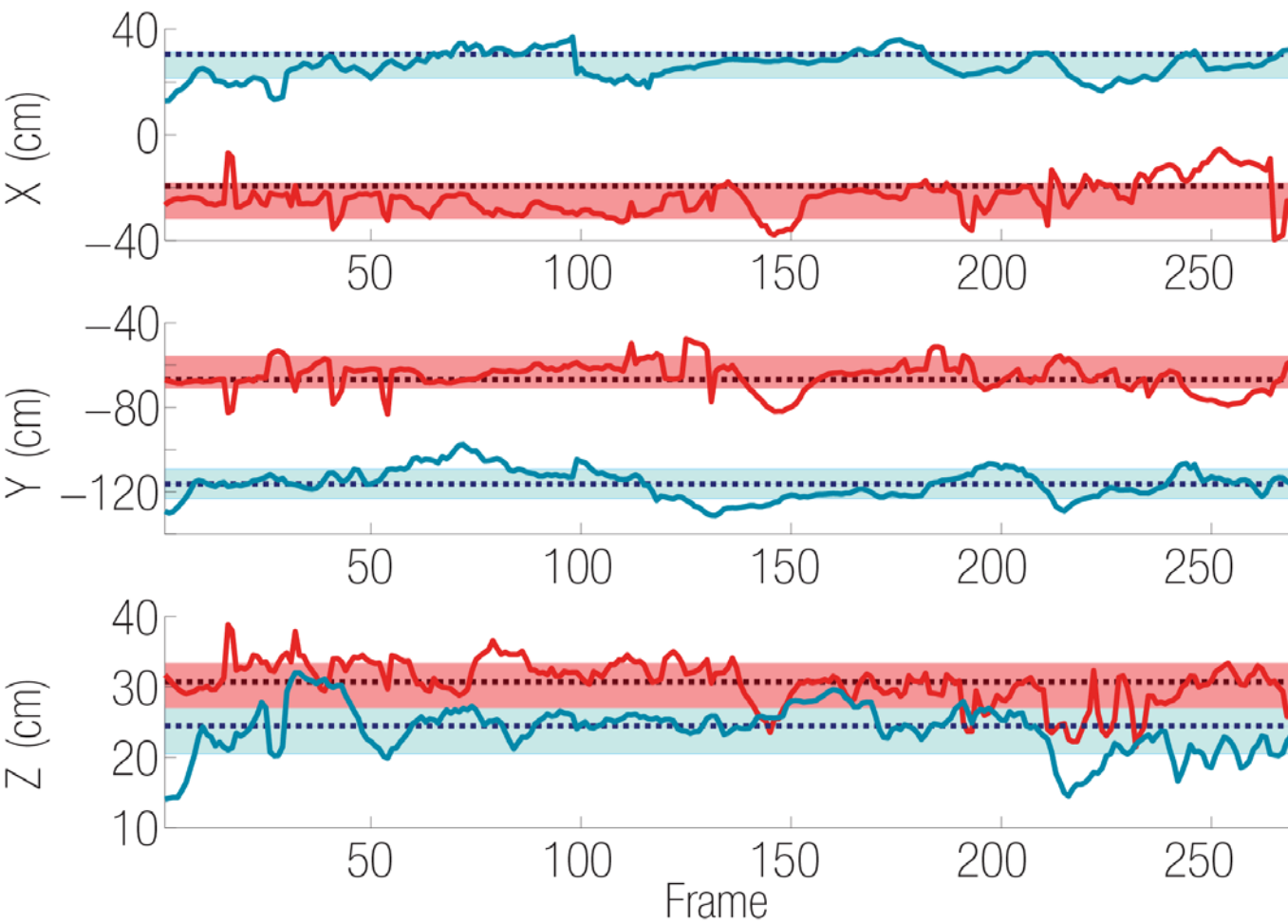
3D Joint Attention Reconstruction



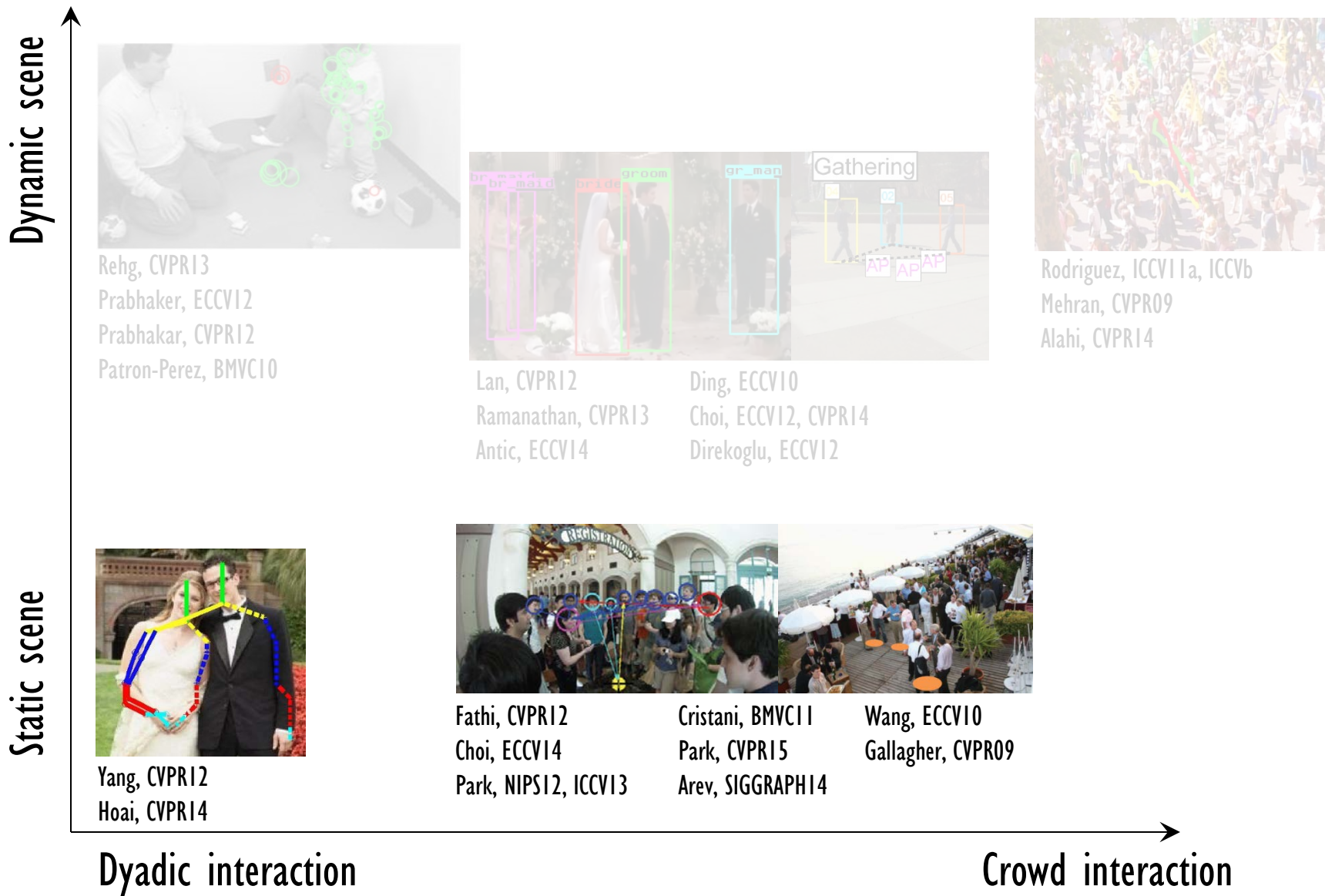


1x speed





Scene dynamism



Number of group members

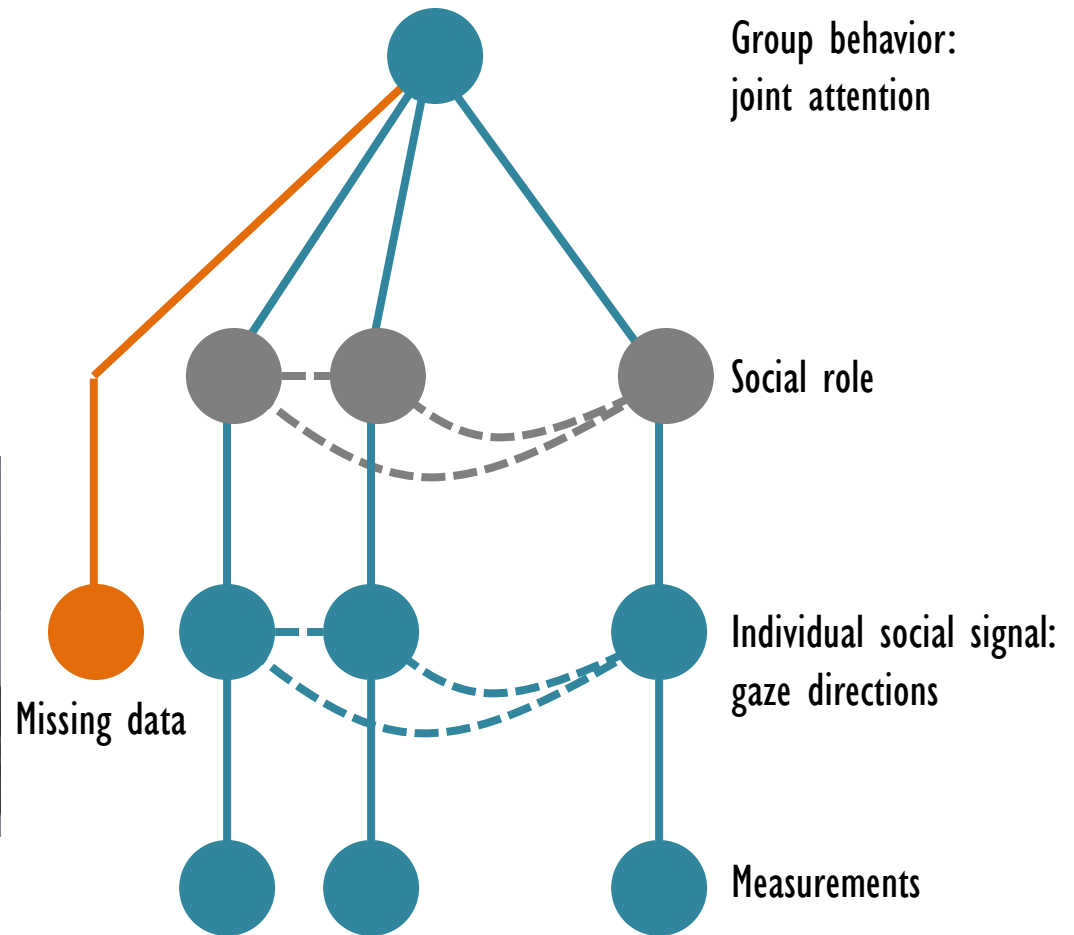
Applications of Joint Attention

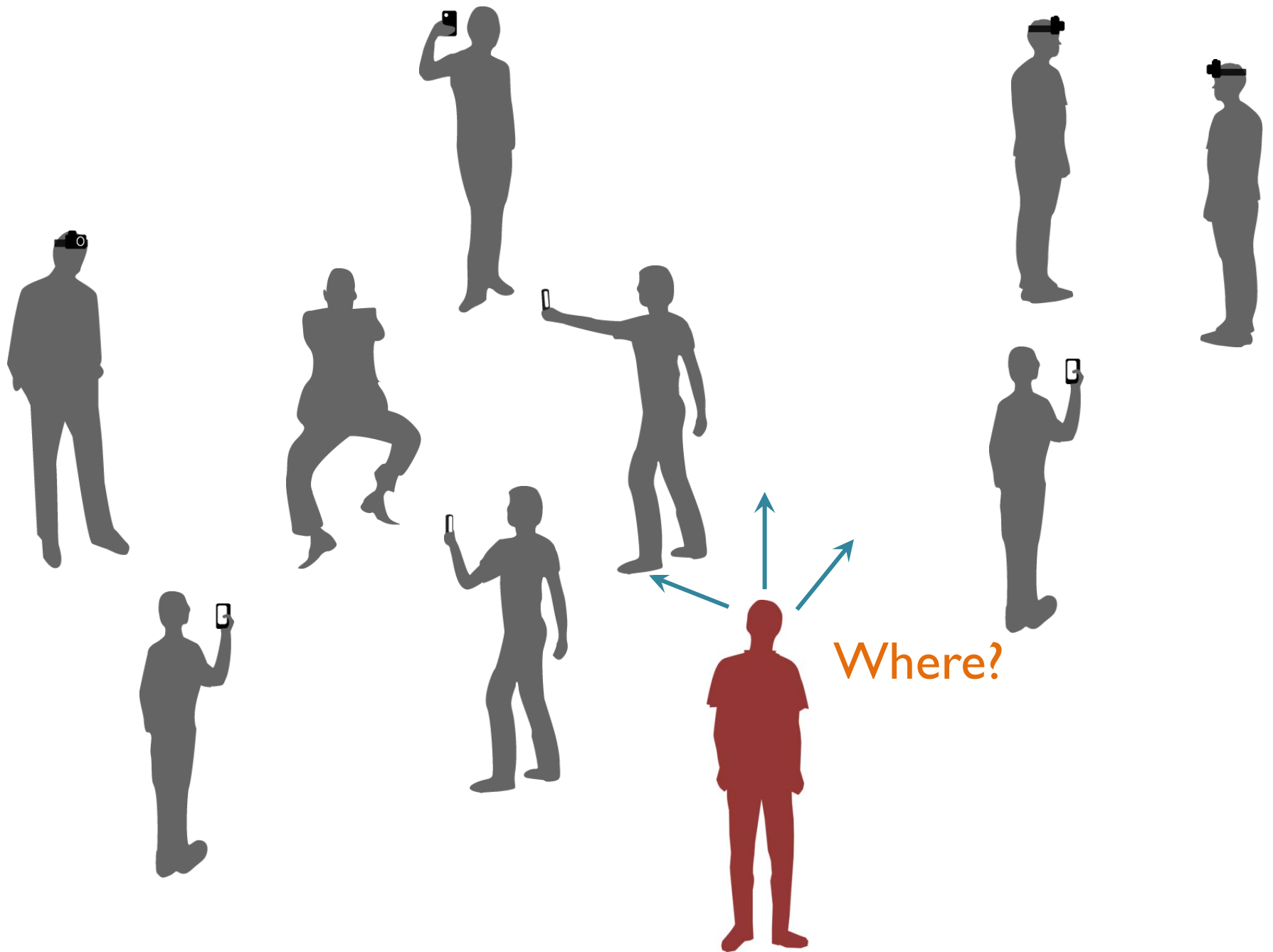
Gaze Prediction

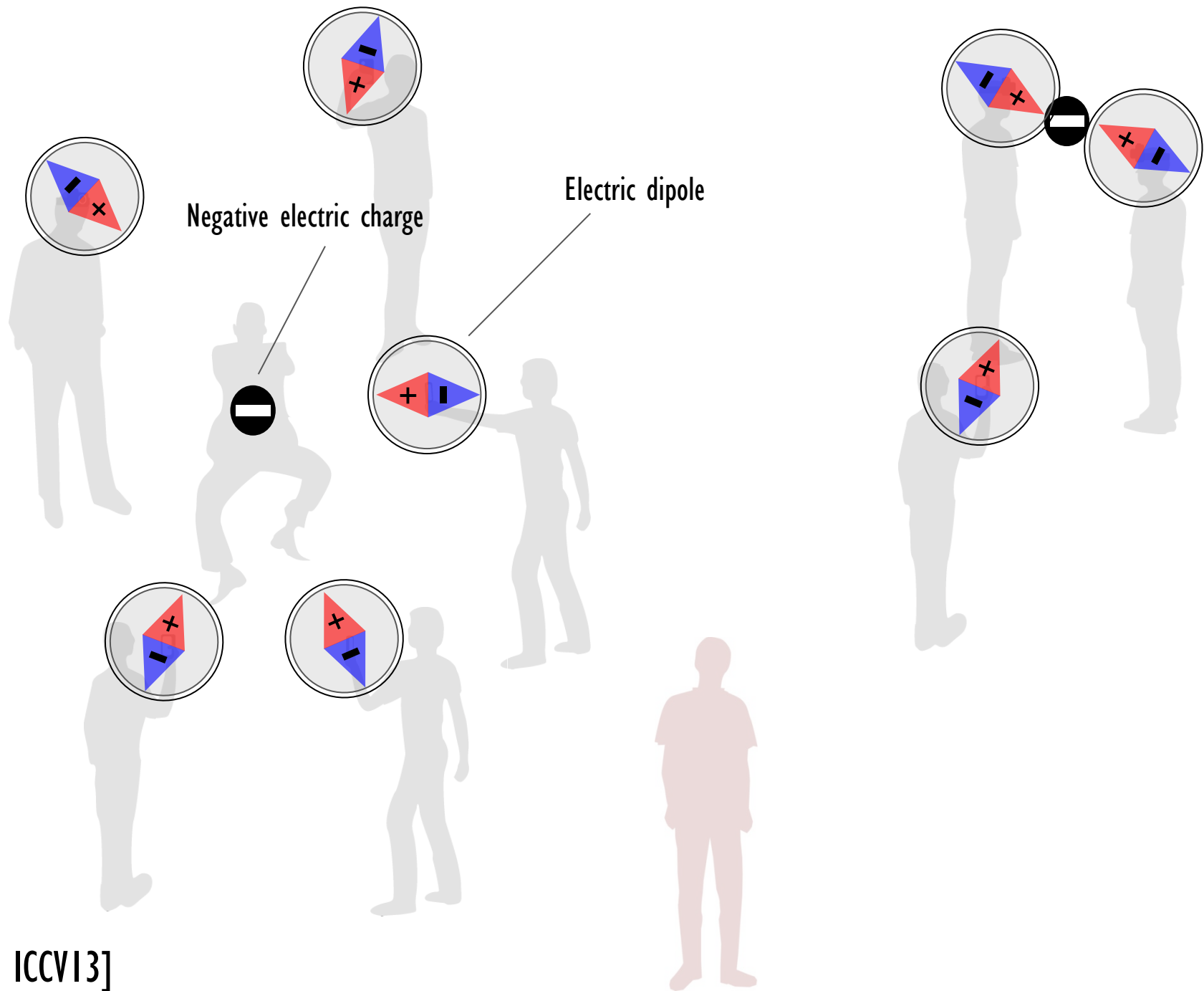
[Park ICCV13]

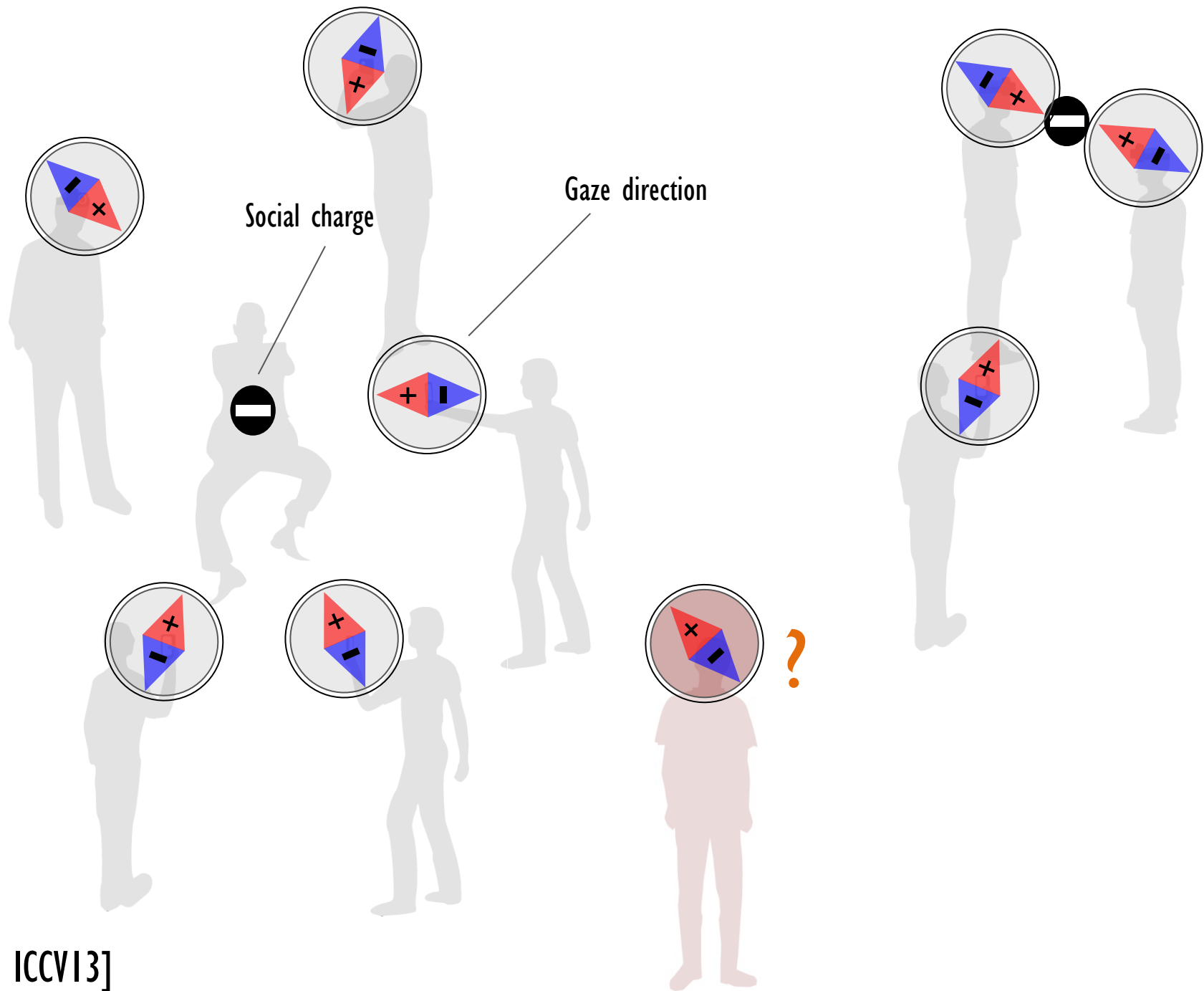
Input: images of social interactions

Output: to predict gaze direction

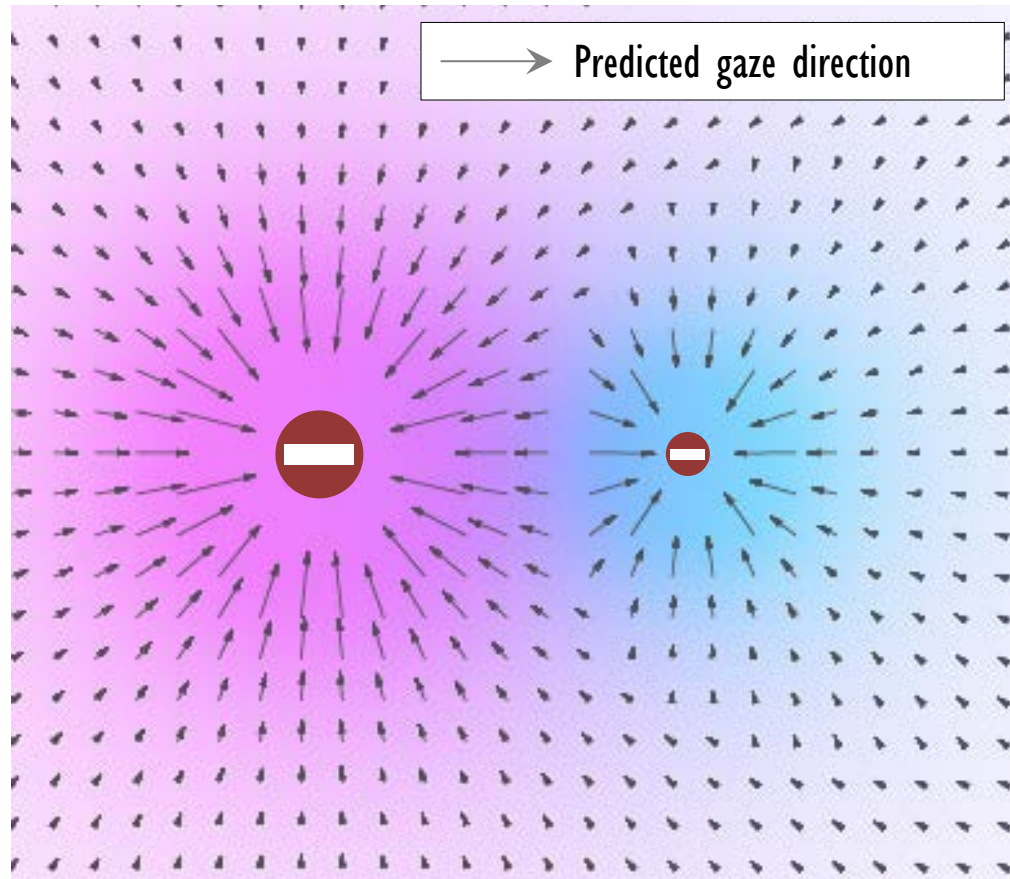








Gaze Field



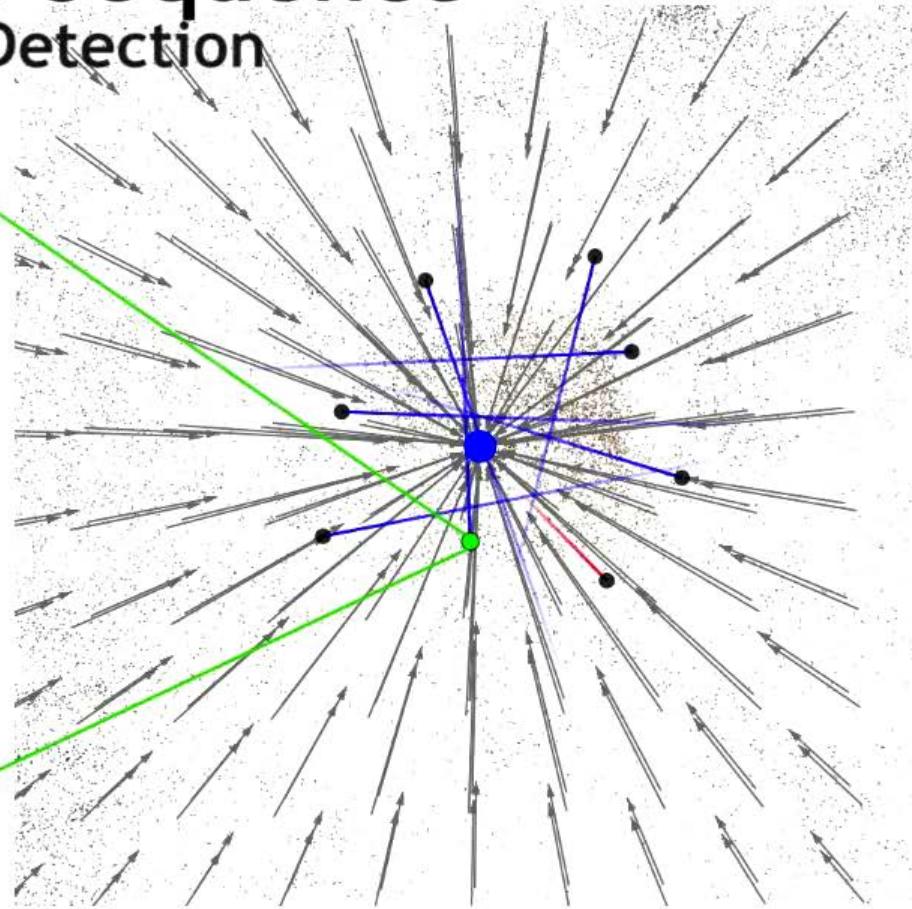
Gaze field $\mathbf{G} = \nabla \Phi \left(\text{red circle with white bar} \right)$

Social Game Sequence

Anomaly Detection



Video from the green marker (member)

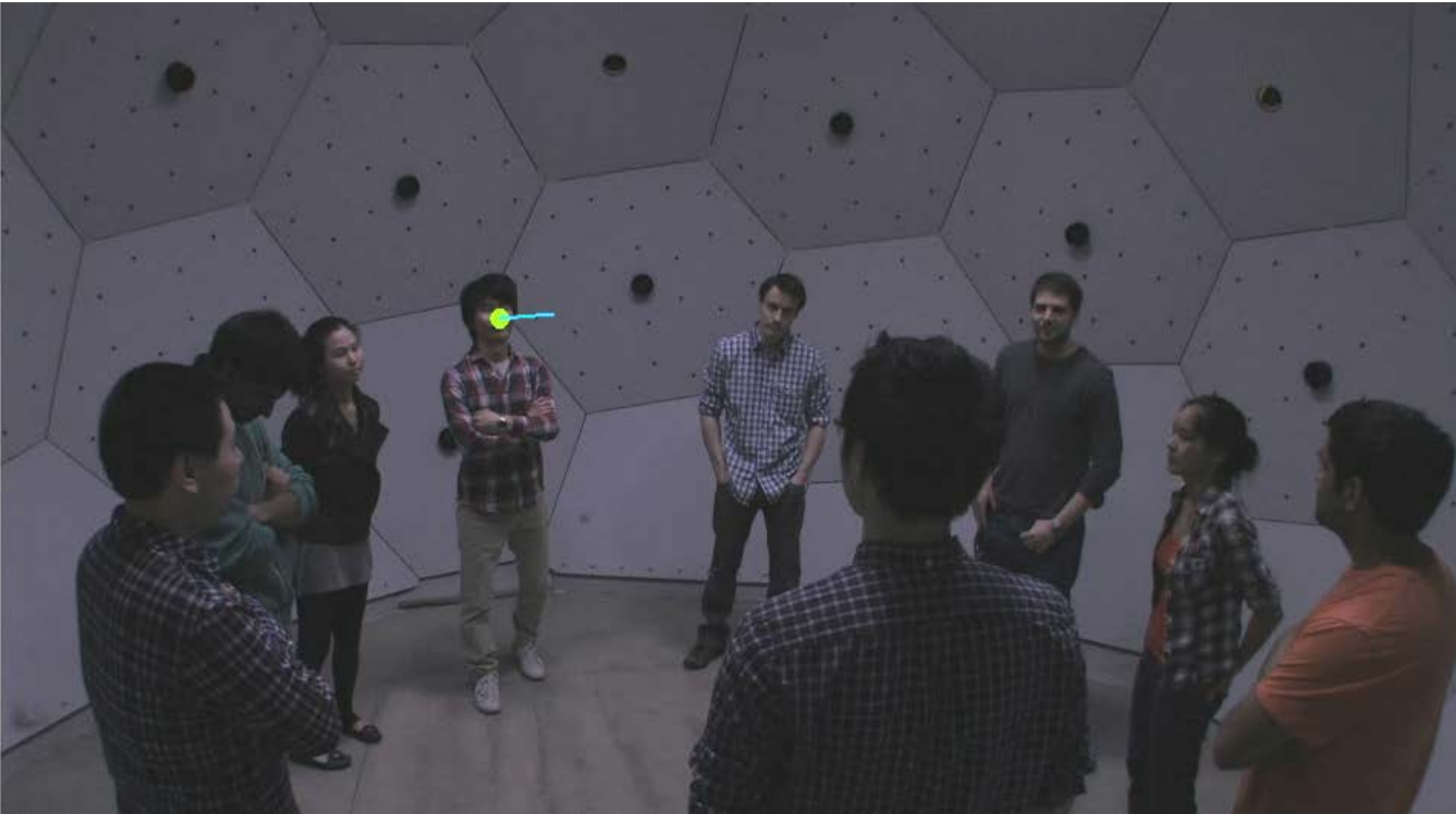


Mafia Game



Mafia Game

Prediction for Missing Data



Social Footage Editing

[Arev SIGGRAPH14]

Input: videos of social interactions

Output: to edit videos to produce a coherent story of social events.

Group behavior:
joint attention

Social role

Individual social signal:
gaze directions

Measurements

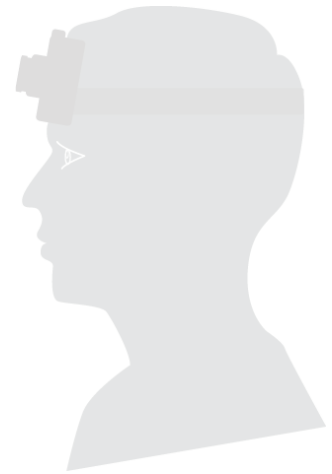
Content creation



First Person Cameras

Problems of videos taken by social cameras:

- Produce too much information to digest at once
- Are biased by an intimate and personal view



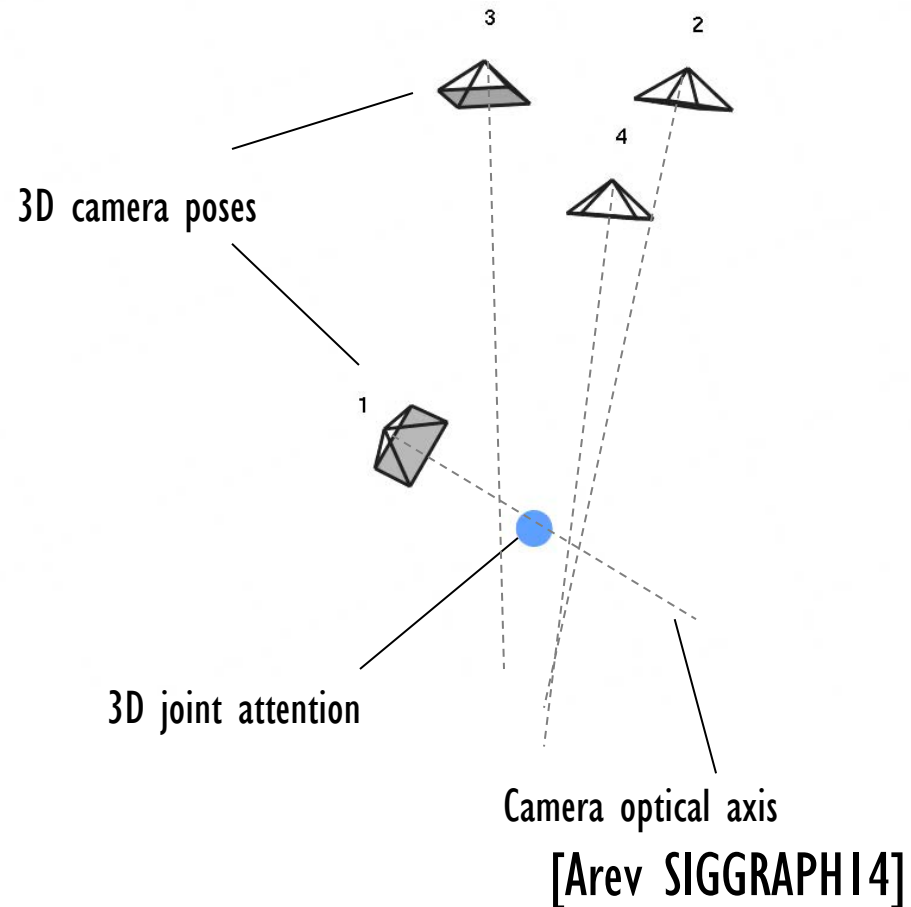
Input: Synchronized Videos

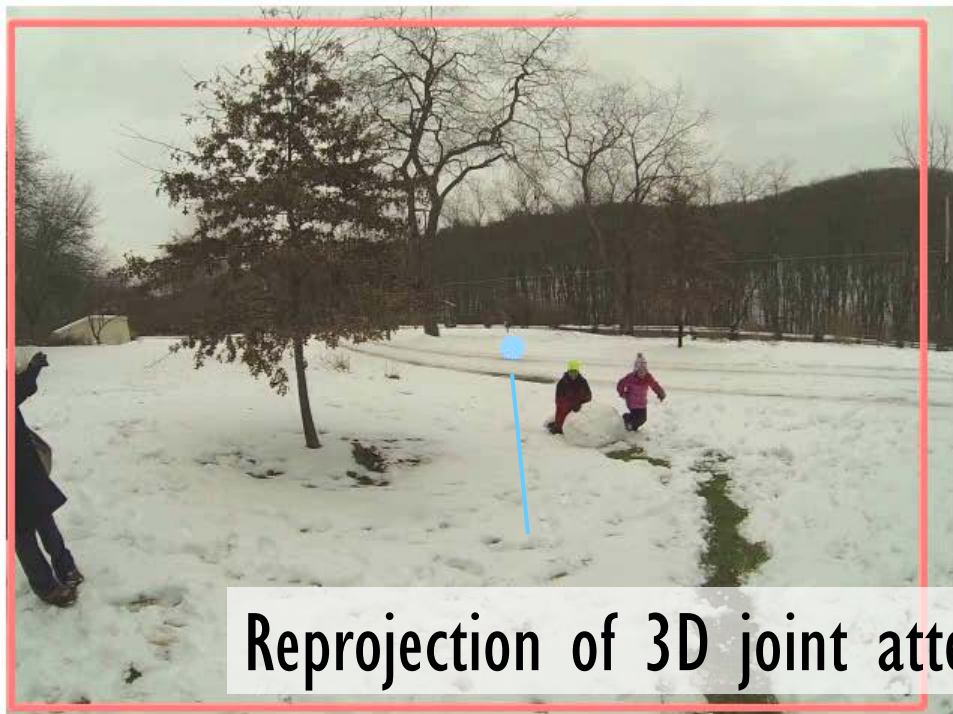


Output: Coherent Video of Event



Content: 3D Joint Attention





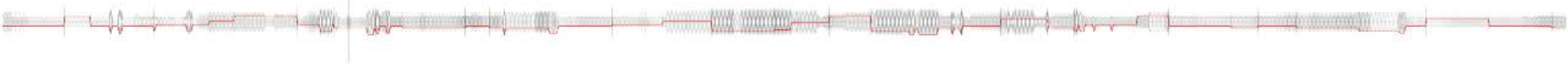
Reprojection of 3D joint attention



Automatic Video Editing



Input video feeds



Timeline

Basketball Scene



[Arev SIGGRAPH14]

Basketball Scene



Scene Summarization



Surprise Party Scene

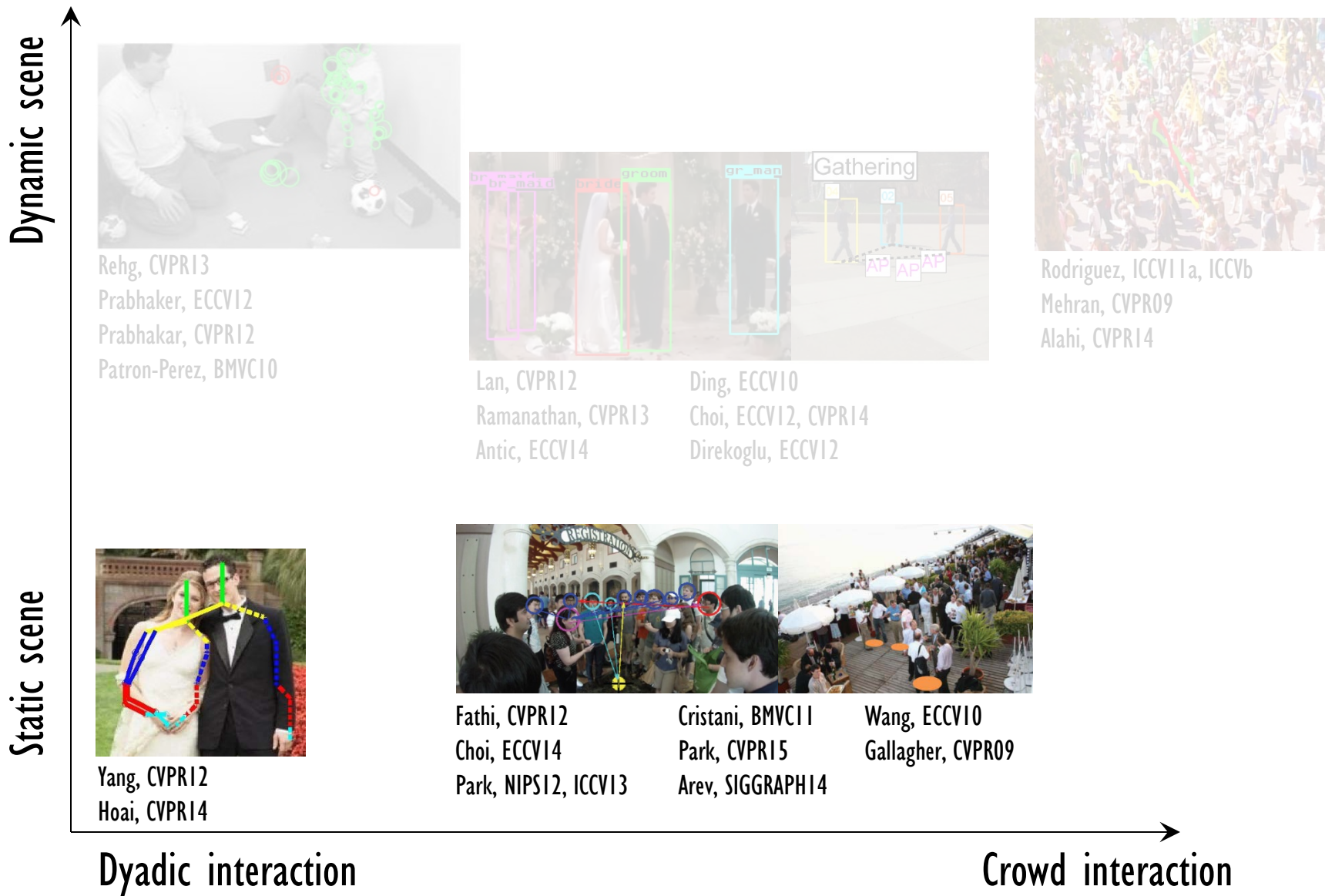


Our method



Professional Editor

Scene dynamism



Number of group members