

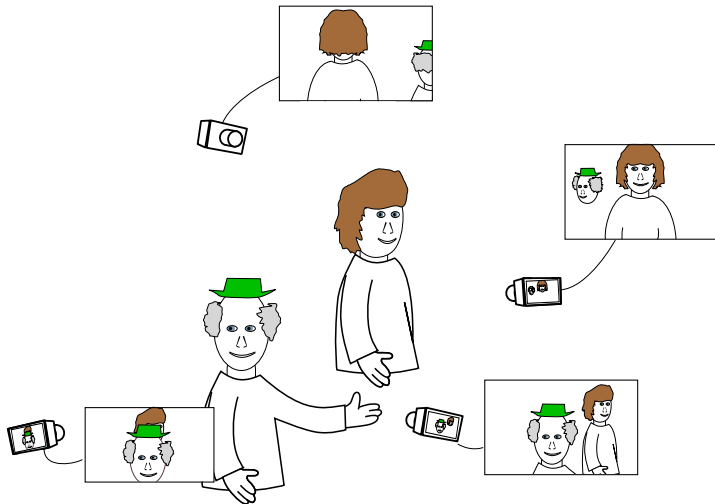
CLUSTER-BASED DISTRIBUTED FACE TRACKING IN CAMERA NETWORKS

Josiah Yoder

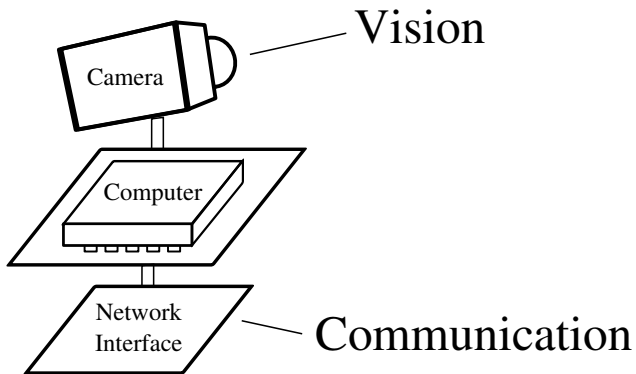
Robot Vision Lab — Purdue University

9 September 2011

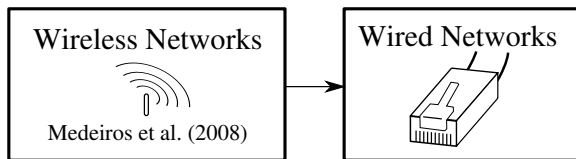
INTRODUCTION — CAMERA NETWORKS



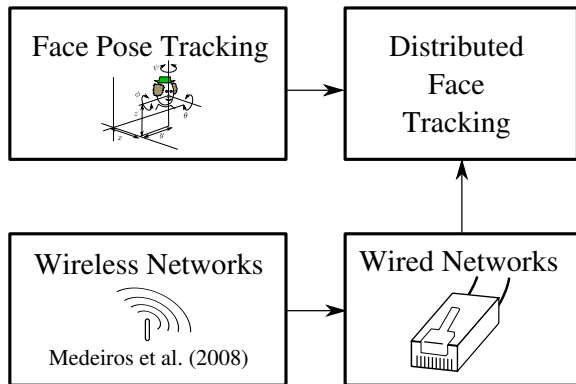
A CAMERA NODE



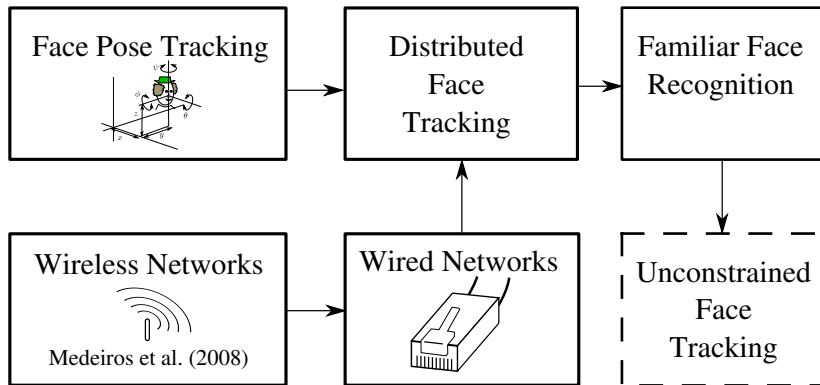
OVERVIEW



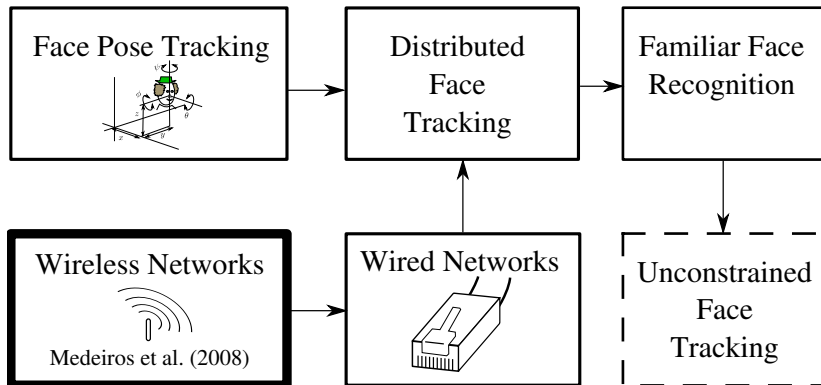
OVERVIEW



OVERVIEW



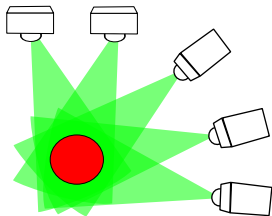
CLUSTER-BASED TRACKING IN WIRELESS NETWORKS



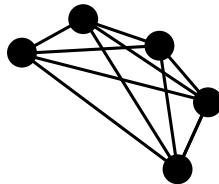
CLUSTER-BASED TRACKING IN WIRELESS NETWORKS

- Developed by Medeiros et al. (2008)
- Addresses challenges of wireless networks:
 - Limited communication range
 - Cameras tracking same target may not be able to communicate

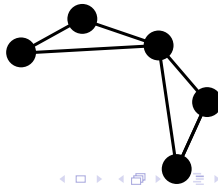
VISION GRAPHS & COMMUNICATION GRAPHS



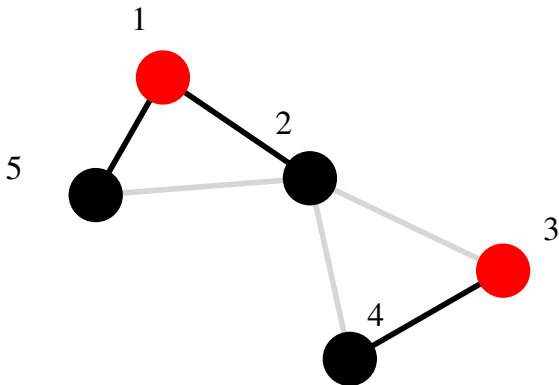
Vision Graph



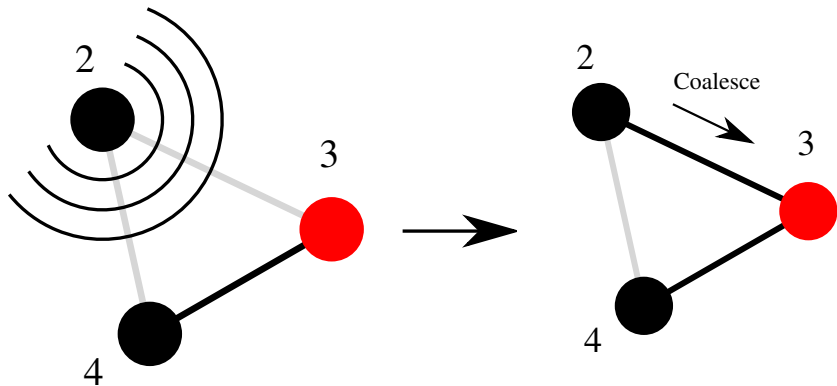
Communication Graph



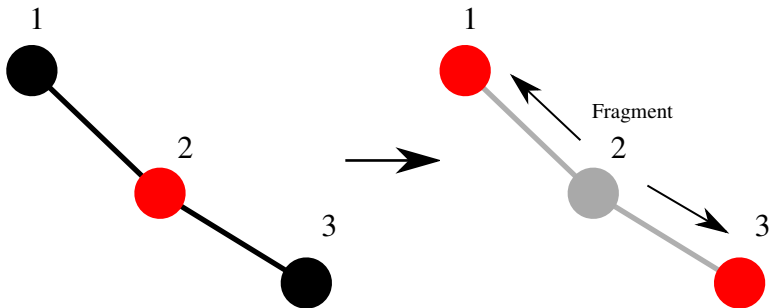
CLUSTER LEADER ELECTION



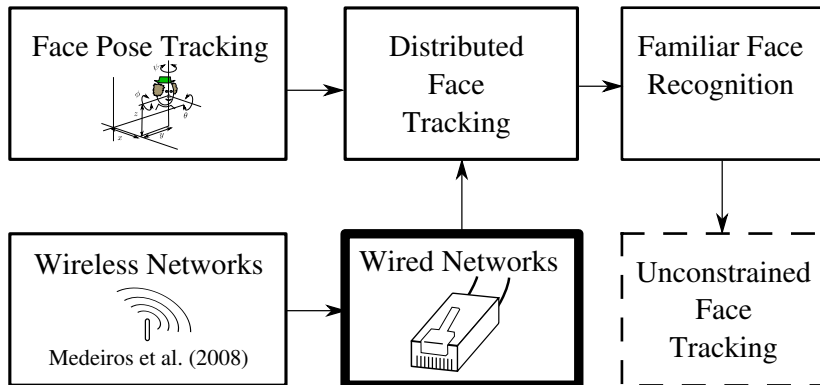
CLUSTER COALESCENCE



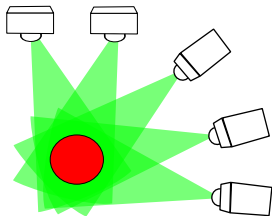
CLUSTER FRAGMENTATION



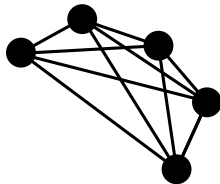
TRACKING IN WIRED NETWORKS



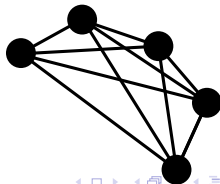
TRACKING IN WIRED NETWORKS



Vision Graph



Communication Graph

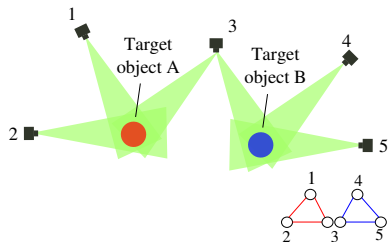


ESTABLISHING CORRESPONDING DETECTIONS

- How can cameras determine they have detected the same target?
 - Detect objects
 - Extract visual features
 - Apply similarity criterion
 - Objects “match” if criterion passes a decision threshold
- Many variations in features and matching criteria
 - Color Histograms
 - HoG features
 - SIFT features
 - Point clouds
 - Face pose estimates
 - Gabor jets
 - ...

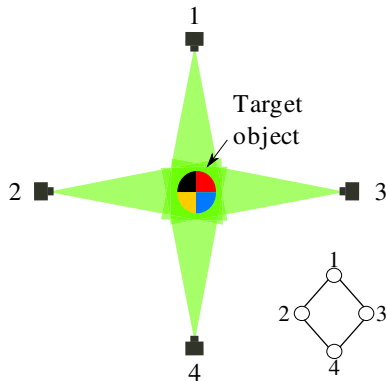
TRACKING GRAPHS

- Edge between two cameras if their detected objects pass the matching criterion
- Cameras may participate in more than one tracking graph
- Dynamic: Changes as objects move



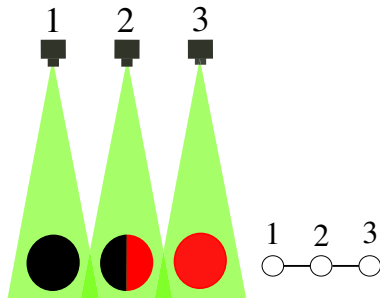
TRACKING GRAPHS: CLUSTER FORMATION

- May have missing edges in the tracking graph
- Cameras can only establish correspondence through other cameras
- For rapid cluster formation, cameras join clusters only with immediate neighbors

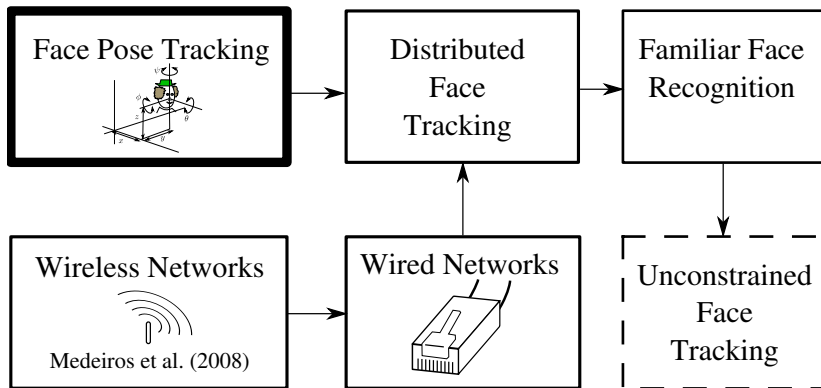


TRACKING GRAPHS: CLUSTER FORMATION

- May have false edges in tracking graph
- Leads to a single cluster tracking multiple targets
- Fixed by cluster fragmentation during propagation



A FRAMEWORK FOR MULTI-CAMERA FACE TRACKING

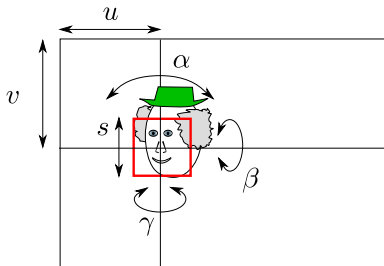
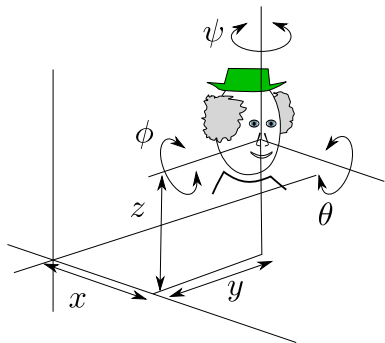


A FRAMEWORK FOR MULTI-CAMERA FACE TRACKING

A framework for face detection, pose estimation, and tracking:

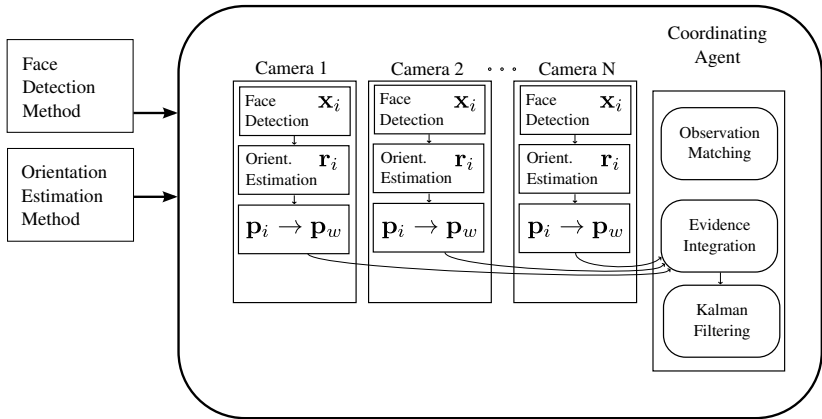
- Allows generic single-camera methods to be incorporated into a multi-camera method
- Representation of face position as a coherent 6-DOF quantity

THE 6-DOF WORLD AND IMAGE-BASED POSES



$$\begin{aligned}\mathbf{p}_w &= [x, y, z, \theta, \phi, \psi]^T \\ &= [\underbrace{x, y, z}_{\mathbf{x}_w^T}, \underbrace{\theta, \phi, \psi}_{\mathbf{r}_w^T}]^T\end{aligned}$$

$$\begin{aligned}\mathbf{p}_i &= [u, v, s, \alpha, \beta, \gamma]^T \\ &= [\underbrace{u, v, s}_{\mathbf{x}_i^T}, \underbrace{\alpha, \beta, \gamma}_{\mathbf{r}_i^T}]^T\end{aligned}$$



TRANSFORMATION OF POSITION AND ROTATION

- We can transform observations from image to world coordinates, through an invertible function \mathbf{f}

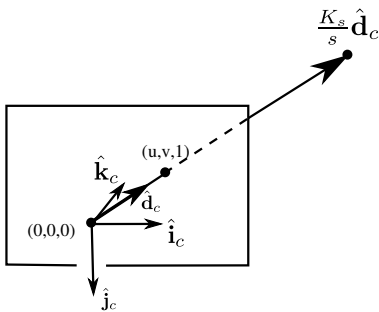
- $\mathbf{p}_w = \mathbf{f}(\mathbf{p}_i)$

- $\mathbf{p}_w = \begin{bmatrix} \mathbf{x}_w \\ \mathbf{r}_w \end{bmatrix} = \begin{bmatrix} \mathbf{f}_x(u, v, s) \\ \mathbf{f}_r(u, v, \alpha, \beta, \gamma) \end{bmatrix}$

- $\mathbf{x}_w = \begin{bmatrix} x \\ y \\ z \end{bmatrix} \quad \mathbf{r}_w = \begin{bmatrix} \theta \\ \phi \\ \psi \end{bmatrix}$

- $\mathbf{f}_x(u, v, s)$ – [Iwaki et al. 2008]
- $\mathbf{f}_r(u, v, \alpha, \beta, \gamma)$ – [Murphy-Chutorian and Trivedi 2008]

TRANSFORMATION OF POSITION



$$\mathbf{x}_w = \mathbf{f}_x(\mathbf{x}_i) = {}_w\mathbf{R}_c \left(\frac{K_s}{s} \hat{\mathbf{d}}_c \right) + {}_w\mathbf{t}_c$$

$$\hat{\mathbf{d}}_c = (u\hat{i}_c + v\hat{j}_c + \hat{k}_c) / (\sqrt{u^2 + v^2 + 1})$$

TRANSFORMATION OF ROTATION

$$\mathbf{r}_w = \mathbf{f}_r(\mathbf{x}_i, \mathbf{r}_i) = [{}_w\mathbf{R}_c {}_c\mathbf{R}_i [\mathbf{r}_i]_{3 \times 3}]_{3 \times 1}$$

- ${}_w\mathbf{R}_c$ — Camera Rotation
- ${}_c\mathbf{R}_i$ — Murphy-Chutorian & Trivedi Rotation
- $[\]_{3 \times 3}$ — Conversion to Rotation Matrix
- $[\]_{3 \times 1}$ — Conversion to Roll-Pitch-Yaw

MURPHY-CHUTORIAN & TRIVEDI ROTATION (2008)

rotation about the axis $\hat{\mathbf{k}}_c \times \hat{\mathbf{d}}_c$ by the angle $\cos^{-1}(\hat{\mathbf{k}}_c \cdot \hat{\mathbf{d}}_c)$

UNCERTAINTY MODELING

We represent observations as Gaussian distributions

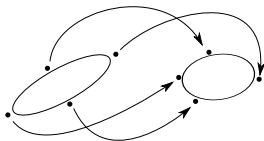
- Image-based observation:

$$\mathbf{p}_i \sim \mathcal{N}(\bar{\mathbf{p}}_i, \mathbf{C}_{\mathbf{p},i})$$

- World observation:

$$\mathbf{p}_w \sim \mathcal{N}(\bar{\mathbf{p}}_w, \mathbf{C}_{\mathbf{p},w})$$

- Transform from \mathbf{p}_i to \mathbf{p}_w using the Unscented Transform



COMPARING OBSERVATIONS

We compute the distance between j^{th} and k^{th} observation using the Mahalanobis distance

$$d(\mathbf{p}_w^j, \mathbf{p}_w^k) = (\bar{\mathbf{p}}_w^j - \bar{\mathbf{p}}_w^k)^T (\mathbf{C}_{\mathbf{p},w}^j + \mathbf{C}_{\mathbf{p},w}^k)^{-1} (\bar{\mathbf{p}}_w^j - \bar{\mathbf{p}}_w^k)$$

Distributed:

- We declare two observations consistent if $d(\mathbf{p}_w^j, \mathbf{p}_w^k) < T$

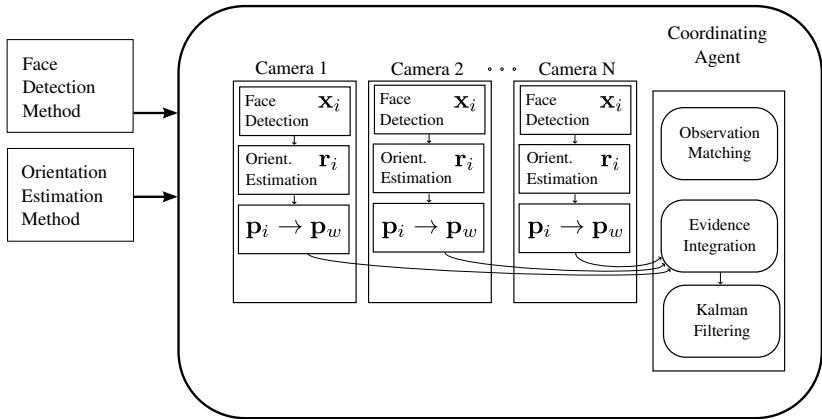
Centralized:

- We employ a feature clustering algorithm based on the distance $d(\mathbf{p}_w^j, \mathbf{p}_w^k) - T_{\text{clique}}$

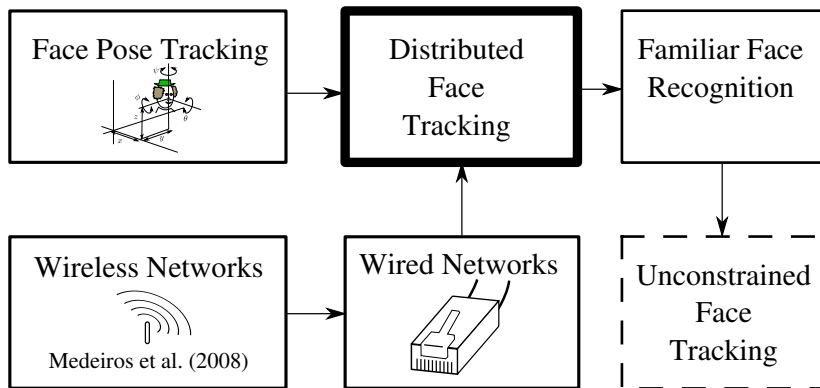
INTEGRATING OBSERVATIONS

We use a minimum-variance estimator to integrate world observations into a more accurate estimate

$$E[\hat{\mathbf{p}}_w] = (Cov[\hat{\mathbf{p}}_w]) \sum_{\mathbf{p}_w^k \in \mathcal{E}} (\mathbf{C}_{\mathbf{p},w}^k)^{-1} \bar{\mathbf{p}}_w^k$$
$$Cov[\hat{\mathbf{p}}_w] = \left(\sum_{\mathbf{p}_w^k \in \mathcal{E}} (\mathbf{C}_{\mathbf{p},w}^k)^{-1} \right)^{-1}$$



DISTRIBUTED CLUSTER-BASED FACE POSE TRACKING

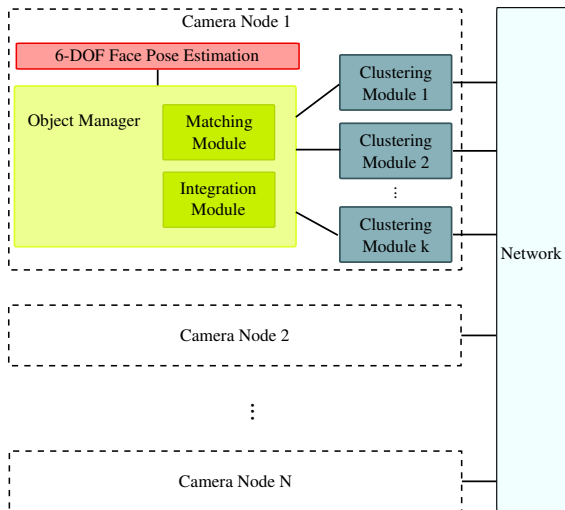


DISTRIBUTED CLUSTER-BASED FACE POSE TRACKING

Here we combine the two systems

- Multi-camera face pose tracking framework
- Cluster-based tracking protocol

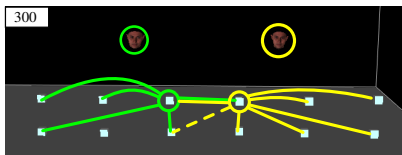
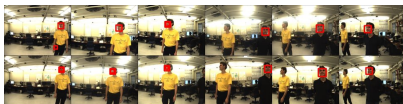
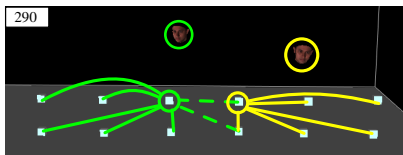
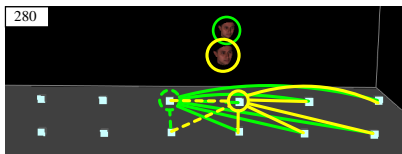
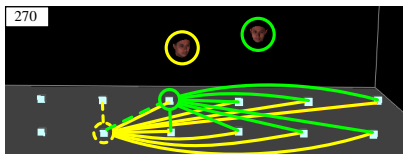
SYSTEM ARCHITECTURE



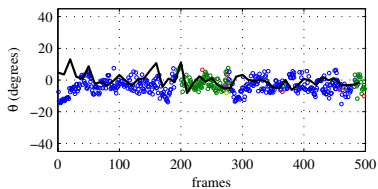
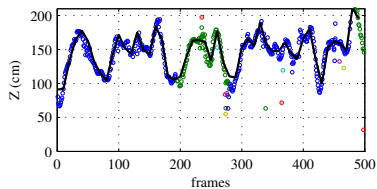
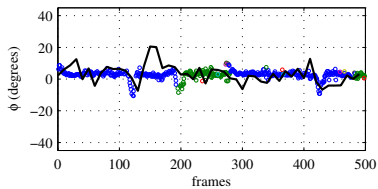
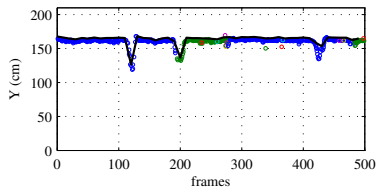
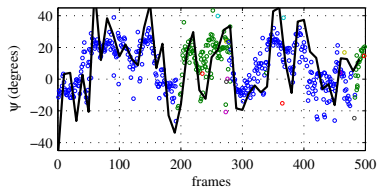
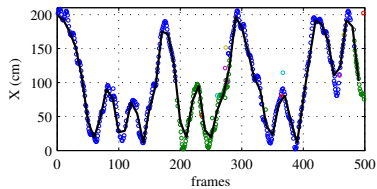
FACE POSE AS IDENTIFYING FEATURE

- Cluster-based protocol makes use of a feature to distinguish targets
- Real-time unconstrained face recognition not available
- We use current face pose for this feature

EXPERIMENTS



EXPERIMENTS



EXPERIMENTS

Comparison with a centralized system

Both

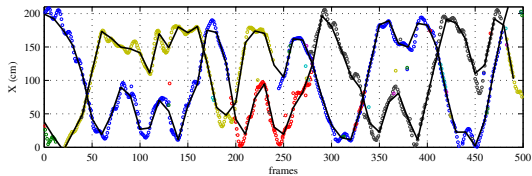
- use same synchronization of collaboration period
- use 6-DOF face pose estimation framework
- detect faces in the individual cameras

In centralized system

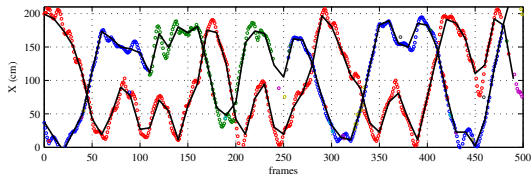
- 6-DOF observations are sent to a central server
- Correspondences are computed based on all pairwise matches

EXPERIMENTS

Distributed



Centralized



EXPERIMENTS

	<i>TP</i>	<i>FP</i>	<i>rmse_T</i> (<i>cm</i>)	<i>rmse_R</i> ($^{\circ}$)
Centralized	95 (95%)	12 (12%)	5.8	20.8
Distributed	94 (94%)	4 (4%)	6.1	18.7

FACE RECOGNITION

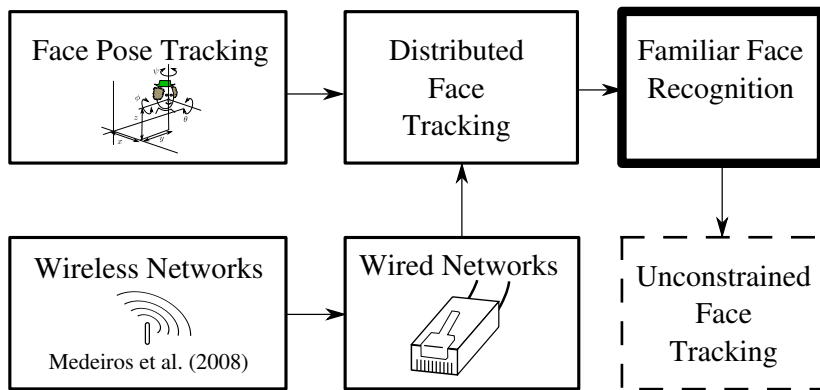
- Cluster-based protocol can also be used for other activities
- Here, we perform distributed face recognition
- Face recognition useful for multi-camera tracking
 - Associate observations from multiple cameras
 - Associate multiple tracks with the same person
 - Restore lost tracks
 - Many other applications
- Each camera performs PCA face recognition
 - Project face images into PCA space
 - Select nearest neighbor from training set (gallery) of faces
 - Send vote for that person to cluster leader
- Cluster leader counts votes and declares overall winner

EXPERIMENTS: DISTRIBUTED FACE RECOGNITION

Tracking TP / FP	Recognition TP / FP
92.4% / 7.6%	87.9% / 9.9%

- Completely distributed:
 - No central server, single point of failure
 - Scalable — Load on each link or node does not increase with network size
- Only using frontal faces ... uncommon in camera networks

HUMAN FAMILIARITY-BENCHMARKED FACE RECOGNITION DATABASE



UNCONSTRAINED FACE RECOGNITION

- Face recognition useful for multi-camera tracking
- Current algorithms poor for unconstrained face images
 - Low resolution
 - Varying pose
- Human performance is still significantly better for unconstrained images
- How to compare algorithms with the best human performance?

BENCHMARKS FOR UNCONSTRAINED FACE RECOGNITION

- The best human performance: “familiar” face recognition
- Unfair comparison?
 - People have prior knowledge unavailable to algorithms
 - Memories from previous encounters
 - Emotions, social relationships, etc.
- Can use same prior knowledge:
 - Videos of people to be pictured in the testing phase
- People can gain some sort of “familiarity” through watching the videos
- People can gain more familiarity if they chat while watching the videos (Bruce et al., 2001)

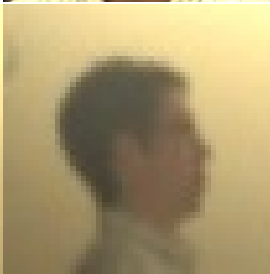
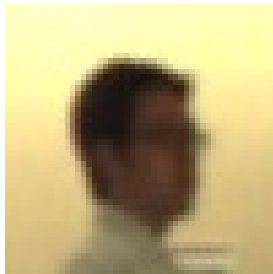
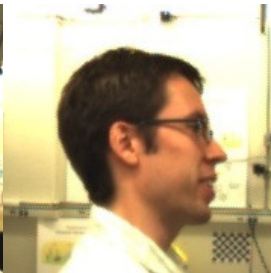
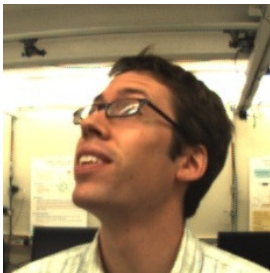
N-VIEWER FAMILIARITY BENCHMARK

- Familiarization: (training)
 - N people watch videos together
- Testing:
 - Each person performs the face matching recognition task separately
- Questions for familiarization:
 - 1 Which of your friends does he/she look like?
 - 2 What sports or hobbies might he/she like?
 - 3 What actor/actress or politician does he/she look like?
 - 4 What might his/her major be?
 - 5 Make a nickname for him/her.
 - 6 Describe his/her personality.
- Here, we use only 1-viewer familiarization.

1-VIEWER BENCHMARKED DATABASE

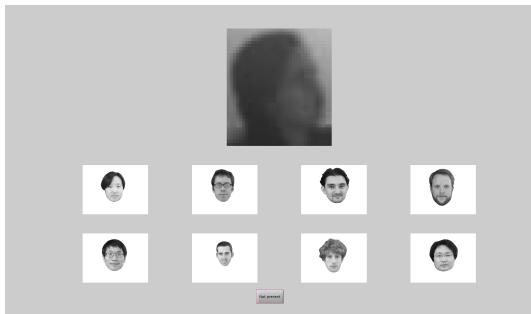


DATABASE



1-VIEWER BENCHMARKED DATABASE

- 20 subjects tested for recognition matching ability
- Unfamiliar or 1-Viewer Familiar Test
- Familiar Test

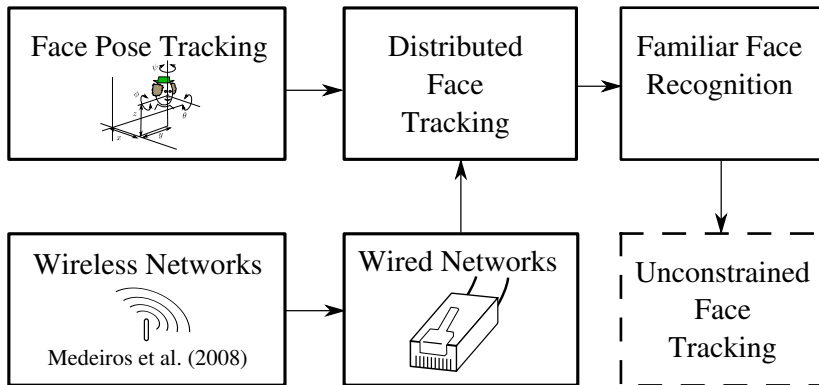


BENCHMARKED PERFORMANCE

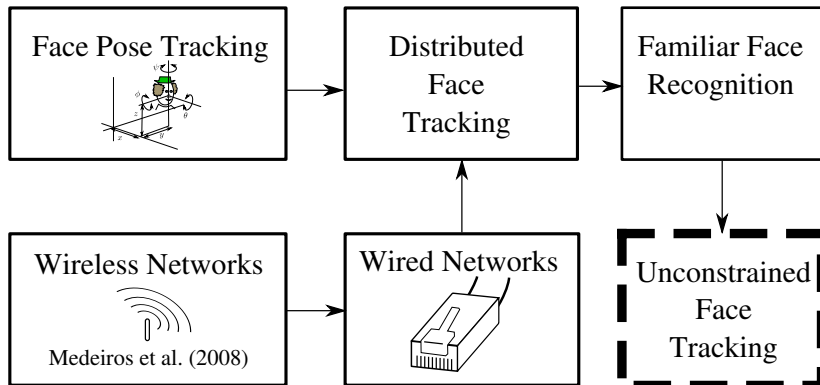
	% correct, mean (std)
Unfamiliar	54 (20)
1-Viewer Familiar	54 (17)
Previously Familiar	80 (19)

- 1-Viewer Familiarity does not improve performance
- Previous Familiarity does, even in challenging low-resolution images

CONCLUSIONS



FUTURE WORK



FUTURE WORK

Cluster-based tracking in wired camera networks

- Coalescence and fragmentation without propagation

6-DOF face pose tracking

- Use local roll-pitch-yaw axes centered around current estimate or “bootstrapped” estimate

Familiar face recognition

- Unconstrained face pose tracking in camera networks

FOR FURTHER READING I



H. Iwaki, G. Srivastava, A. Kosaka, J. Park, and A. Kak.
A novel evidence accumulation framework for robust
multi-camera person detection.

In Proceedings of the ACM/IEEE International Conference on
Distributed Smart Cameras, pages 1–10, 2008.



H. Medeiros, J. Park, and A. Kak.

Distributed object tracking using a cluster-based Kalman filter
in wireless camera networks.

In IEEE Journal of Selected Topics in Signal Processing,
volume 2, pages 448–463, 2008.

FOR FURTHER READING II



E. Murphy-Chutorian and M. Trivedi.

3d tracking and dynamic analysis of human head movements and attentional targets.

[Proceedings of the International Conference on Distributed Smart Cameras \(ICDSC'08\), 2008.](#)

PRIOR WORK

- Face Detection, Pose Estimation, and Tracking
- Face Recognition
- Cluster-Based Tracking in Camera Networks

PRIOR WORK: FACE POSE TRACKING

To track face pose, we need to:

- Detect faces
- Estimate face pose
- Track face pose

PRIOR WORK: FACE DETECTION

- Component-based detection



[Heisele et al. 2001]

- Color-based detection



- Scanning window methods



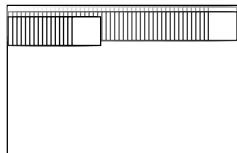
Heisele, B.; Serre, T.; Pontil, M. & Poggio, T. "Component-based face detection", Proc. CVPR, pp. 657–662. 2001.

PARAMETERS DETERMINED BY FACE DETECTION METHODS

- Every face detection method determines
 - Face position
 - Face size
 - Face rotation (approximately)
- For example, for face size can be estimated:
 - from component distance
 - from color blob size
 - from window size



[Heisele et al. 2001]



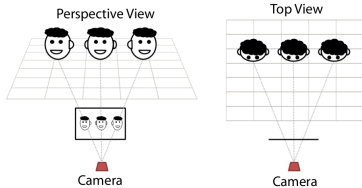
PRIOR WORK: SINGLE-CAMERA POSE ESTIMATION

- There are many pose estimation methods as well [Survey: Murphy-Chutorian and Trivedi 2009]
 - Detector arrays
 - Regression methods
 - Deformable Models
 - etc.
- Methods estimate the face rotation
 - 1-, 2-, or 3-Degrees of Freedom (DOF)
 - Often analyze cropped face image and ignore the rest of the image
 - Rotation as if cropped image were at center of camera

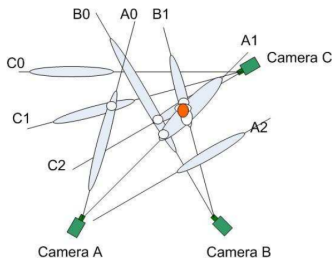
PRIOR WORK: SINGLE-CAMERA FACE TRACKING

- Pose tracking methods find a face iteratively based on the location in the previous frame
- Some of these methods analyze the cropped face image in image-based coordinates
 - Active Appearance Models
 - Appearance-template particle filters

PRIOR WORK: MULTI-CAMERA FACE POSE ESTIMATION



[Murphy-Chutorian and
Trivedi, 2008]

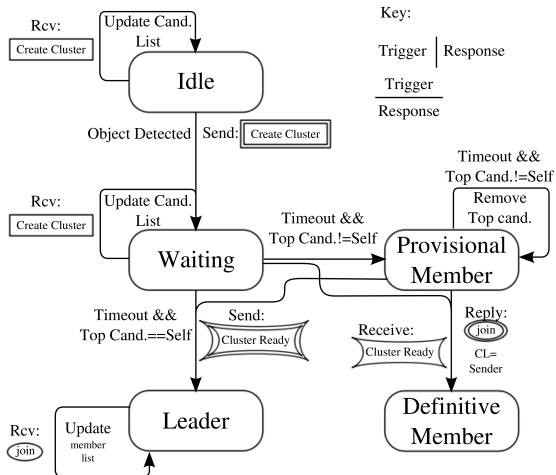


[Iwaki et al. 2008]

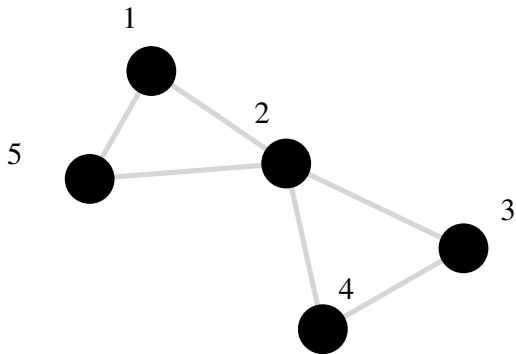
PRIOR WORK: CLUSTER-BASED TRACKING IN SENSOR NETWORKS

- General Sensor Networks:
 - Clusters used to facilitate data aggregation, e.g. for sensor monitoring
- Tracking:
 - A cluster is dedicated to tracking a single target
 - Cameras may participate in multiple clusters, track multiple targets
- Tracking protocols and systems:
 - Zhang and Cao (2004) organize clusters as Dynamic Convoy Trees
 - Blum et al. (2003) avoid creating multiple leaders tracking the same target using multi-hop communication
- No works prior to Medeiros et al. (2008) methods take into account the directional nature of camera sensors.

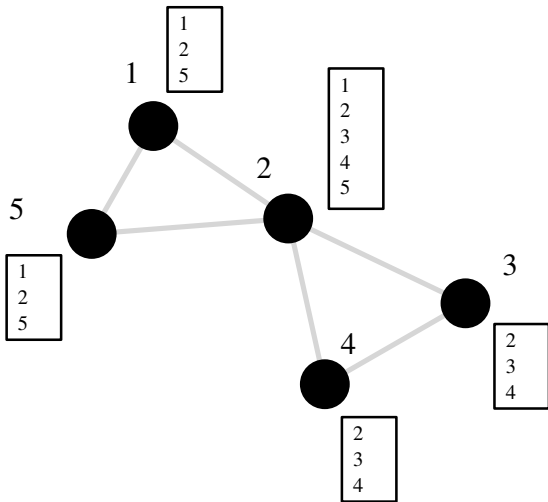
CLUSTER-BASED COMMUNICATION PROTOCOL: CLUSTER LEADER ELECTION



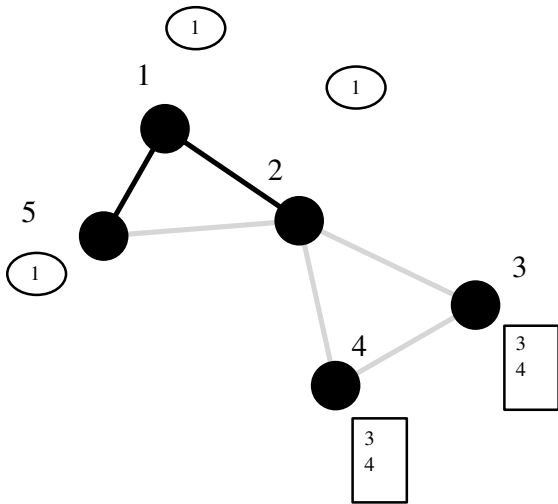
CLUSTER LEADER ELECTION



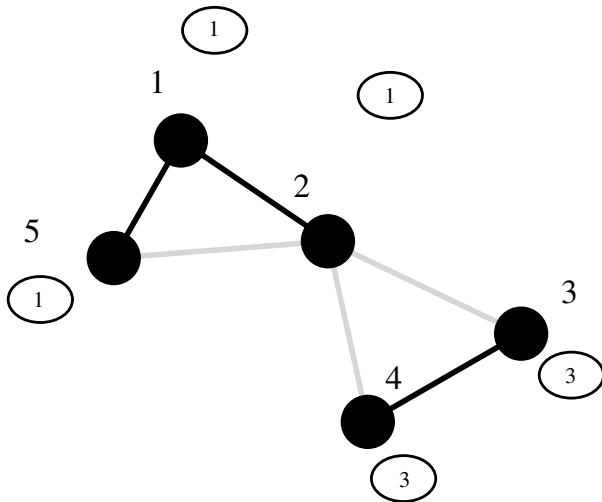
CLUSTER LEADER ELECTION



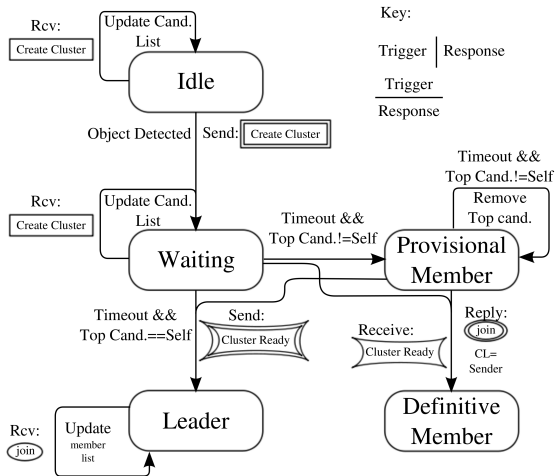
CLUSTER LEADER ELECTION



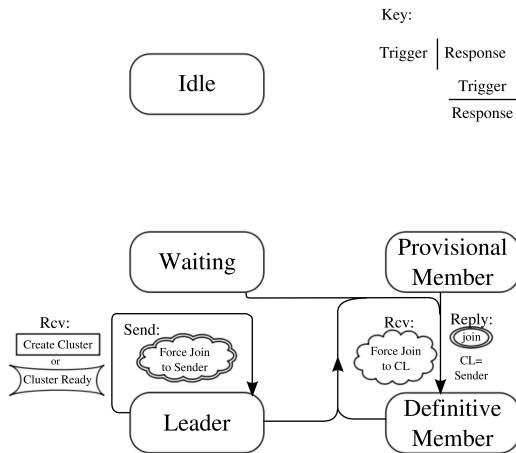
CLUSTER LEADER ELECTION



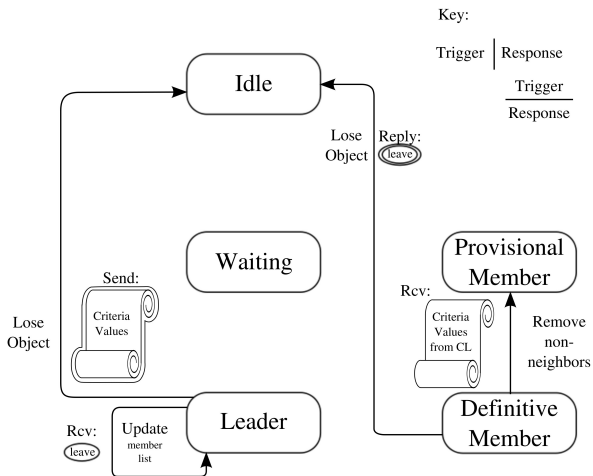
CLUSTER LEADER ELECTION



CLUSTER PROPAGATION (1)



CLUSTER PROPAGATION (2)



UNSCENTED TRANSFORMATION

- We assume that $\mathbf{C}_{\mathbf{p},i}$ is diagonal
- We take a set of $2N$ sigma-points $\dot{\mathbf{p}}_i^k$ in image-based coordinates

$$\begin{aligned}\dot{\mathbf{p}}_i^k &= \bar{\mathbf{p}}_i + \sqrt{N}\sigma_k \mathbf{e}_k & k = 1, \dots, N \\ \dot{\mathbf{p}}_i^{k+N} &= \bar{\mathbf{p}}_i + \sqrt{N}\sigma_k \mathbf{e}_k & k = 1, \dots, N\end{aligned}$$

UNSCENTED TRANSFORMATION

- Then transform each sigma point into world coordinates

$$\dot{\mathbf{p}}_w^k = \mathbf{f}(\dot{\mathbf{p}}_i^k)$$

- And compute the mean and covariance of the points in the world space

$$\bar{\mathbf{p}}_w = \frac{1}{2N} \sum_{k=1}^{2N} \dot{\mathbf{p}}_w^k$$

$$\mathbf{C}_{\mathbf{p},w} = \frac{1}{2N} \sum_{k=1}^{2N} (\dot{\mathbf{p}}_w^k - \bar{\mathbf{p}}_w) (\dot{\mathbf{p}}_w^k - \bar{\mathbf{p}}_w)^T$$