



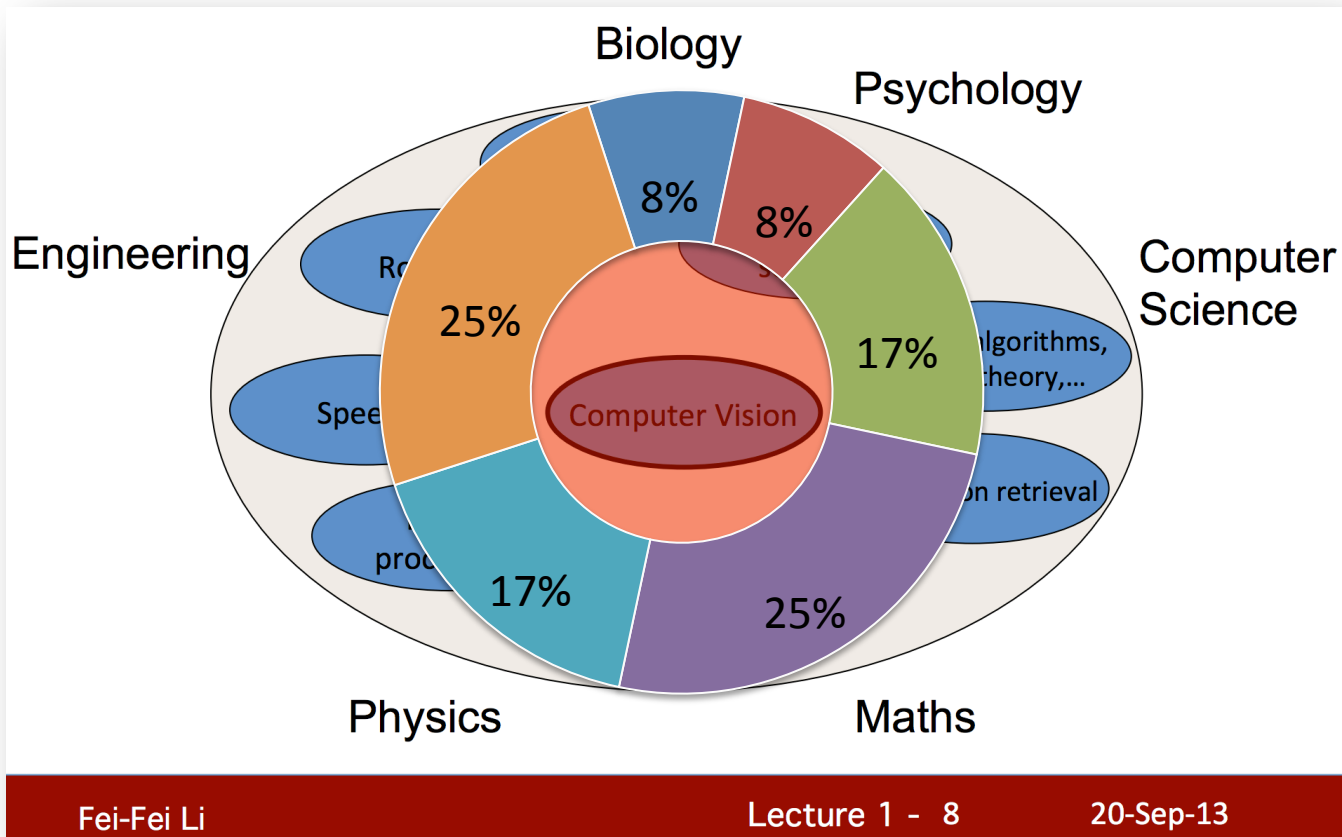
CS131

Tracking millions of people

A. Alahi
November 8th 2013

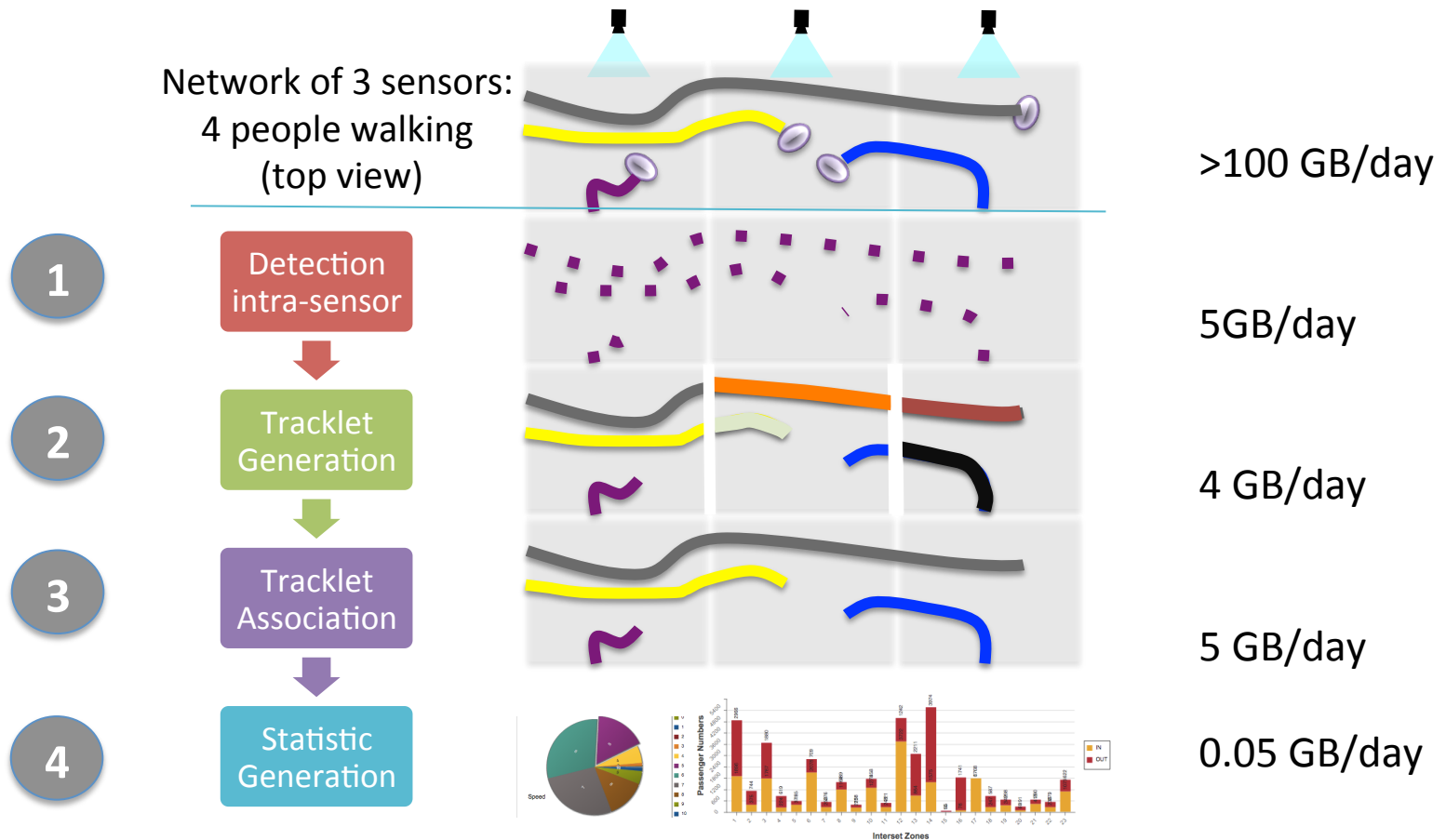


What is related to tracking people?



Outline:

From foreground extraction To tracking millions of pedestrians



Outline:

From foreground extraction To tracking millions of pedestrians

I. Detection

- I. Foreground extraction
- II. Pedestrian localization

II. Tracklet Generation

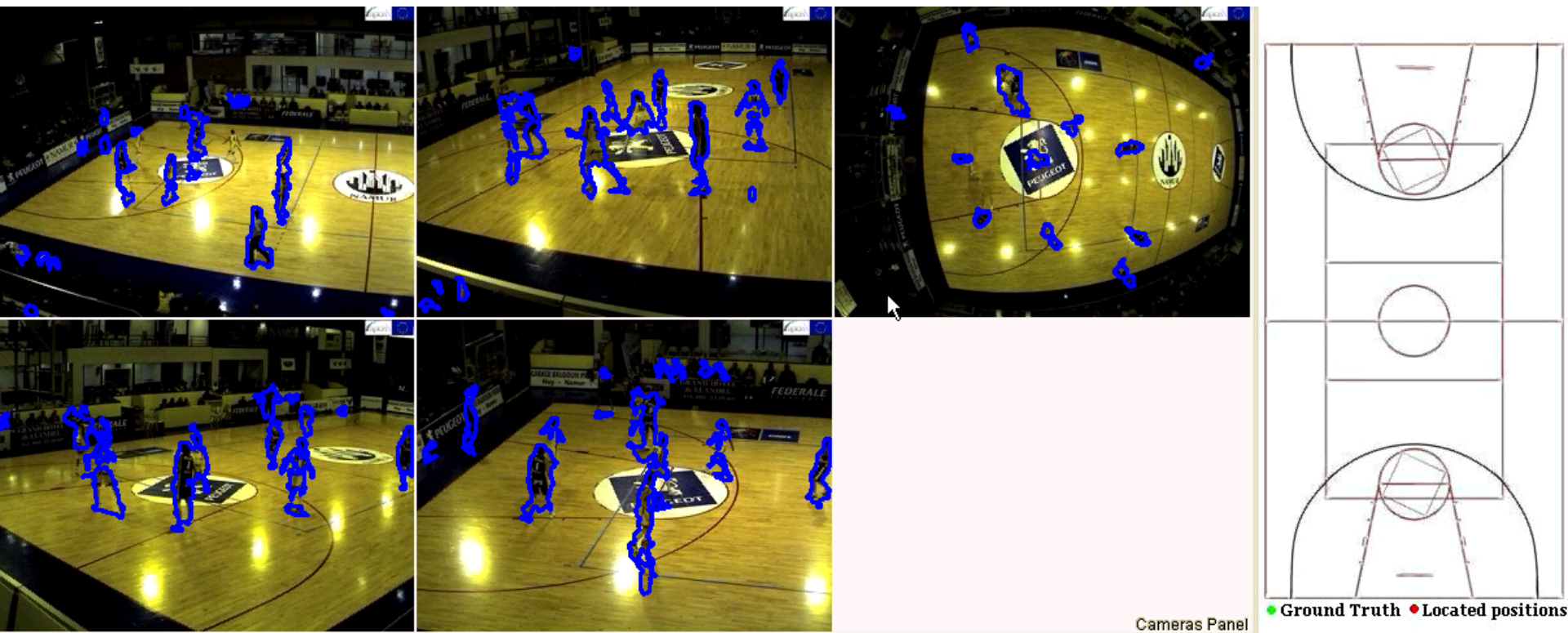
- I. Data association problem
- II. Matching appearance cues

III. Tracklet Association

- I. Modeling Social Affinities



I. Detection: Foreground extraction



- Severely degraded foreground silhouettes

- Spatially dense distribution

- Strong occlusions



I. Detection: Foreground extraction

$$F(x,y, t+1) = \begin{cases} 1 & \text{if } |I(x,y, t) - \mathbf{B}(x,y)| > T \\ 0 & \text{otherwise} \end{cases}$$

- Frame differencing

$$\mathbf{B}(x,y) = I(x,y,t-1)$$

- Mean filter

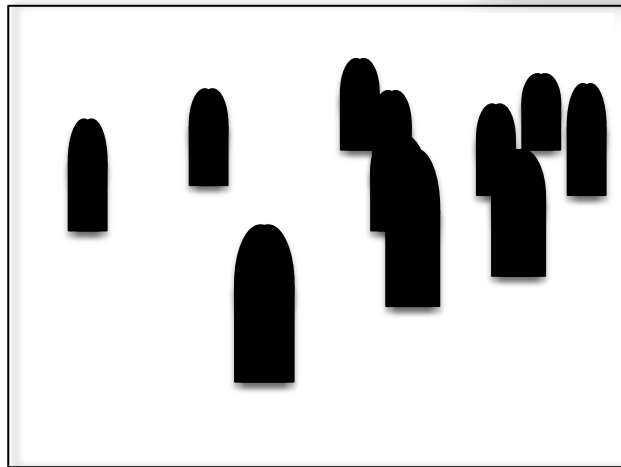
$$\mathbf{B}(x,y) = 1/N \sum_{i=1 \dots N} I(x,y,t-i)$$

- Gaussian averaging
- GMM
- ...

=> Library of 32 algorithms (*BGS library*)

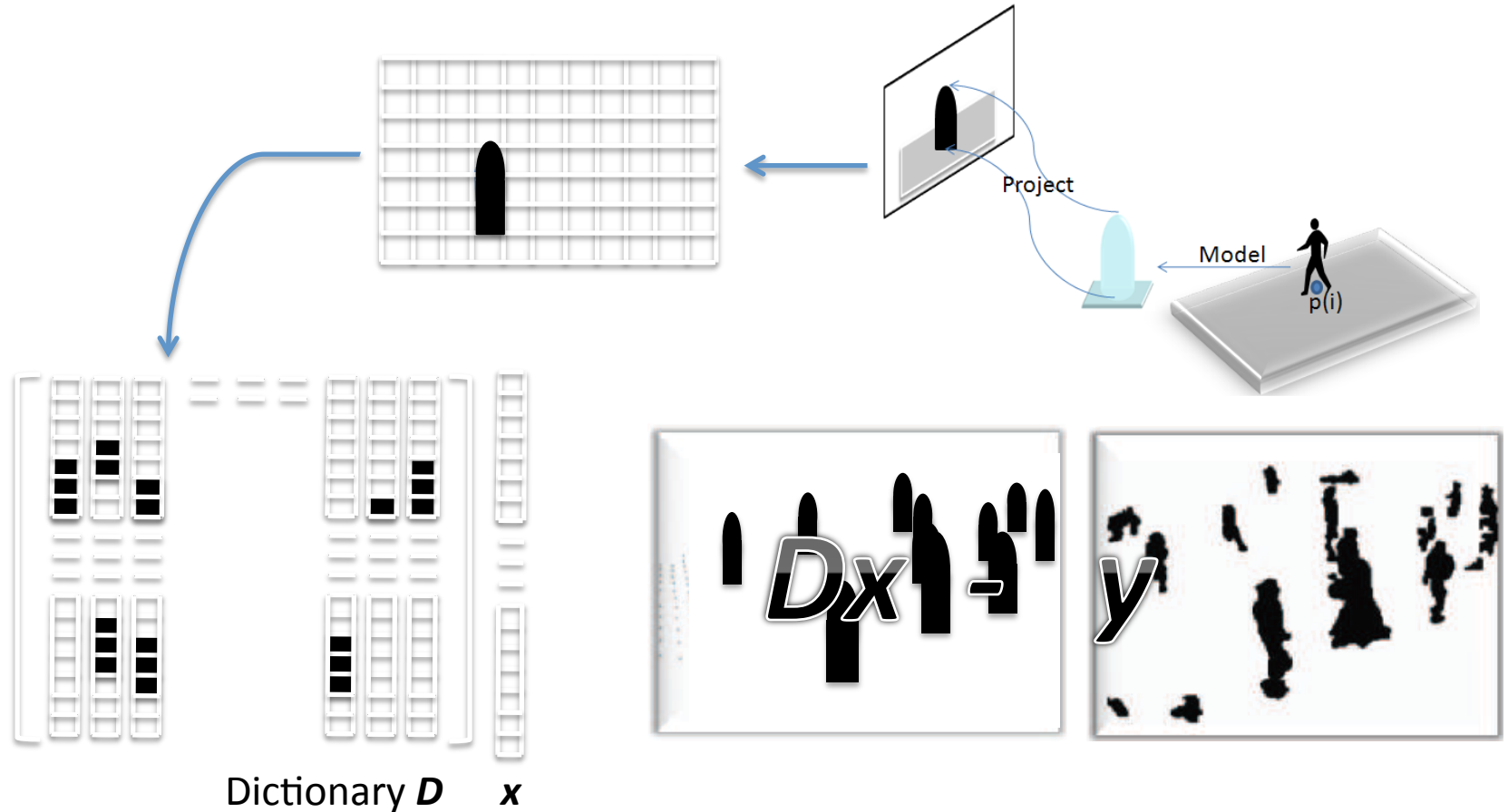


I. Detection: Pedestrian localization

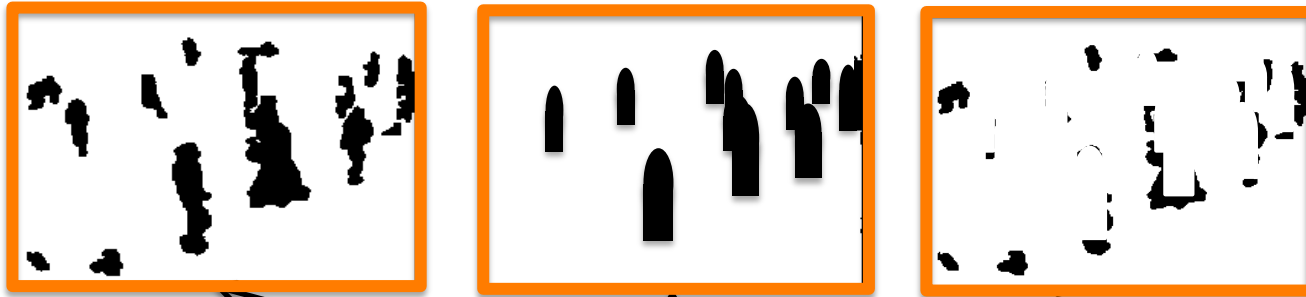


I. Detection: Calibrated Camera

- Create a dictionary D of atoms approximating the ideal foreground silhouette for every position in x



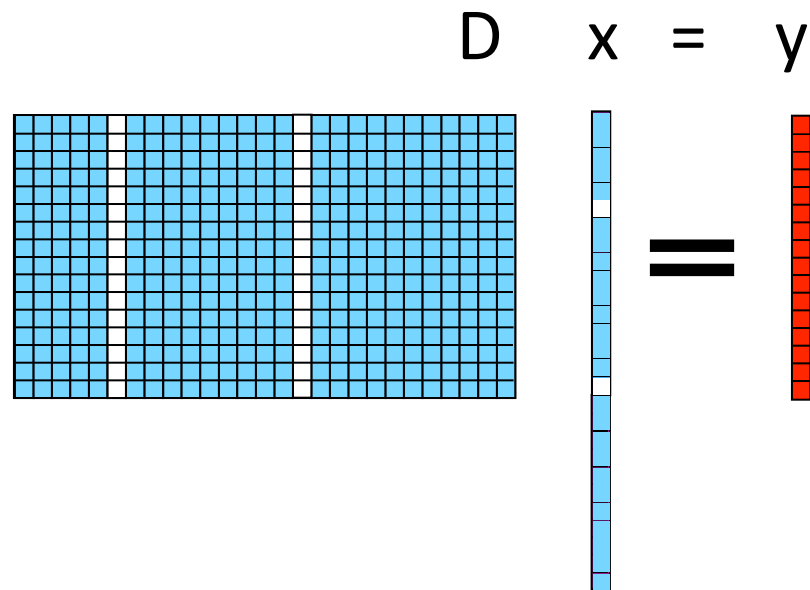
I. Detection: Sparsity driven framework



- Inverse problem:

$$y = Dx + n$$

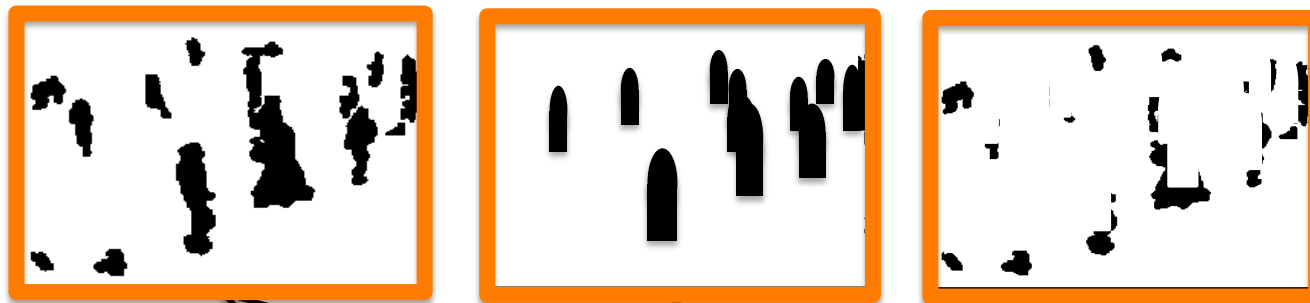
I. Detection: Greedy approach



[1] S. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," IEEE Transactions on signal processing, 1993.



I. Detection: Sparsity driven framework



- Inverse problem:

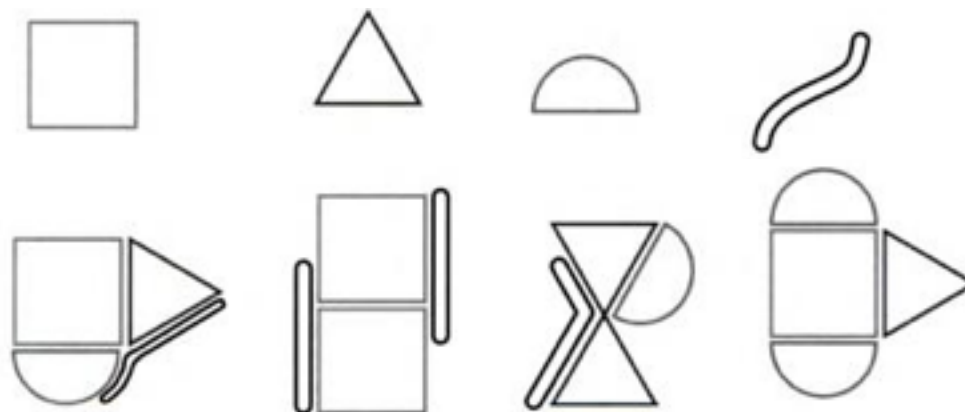
$$y = Dx + n$$

- Sparsity prior:

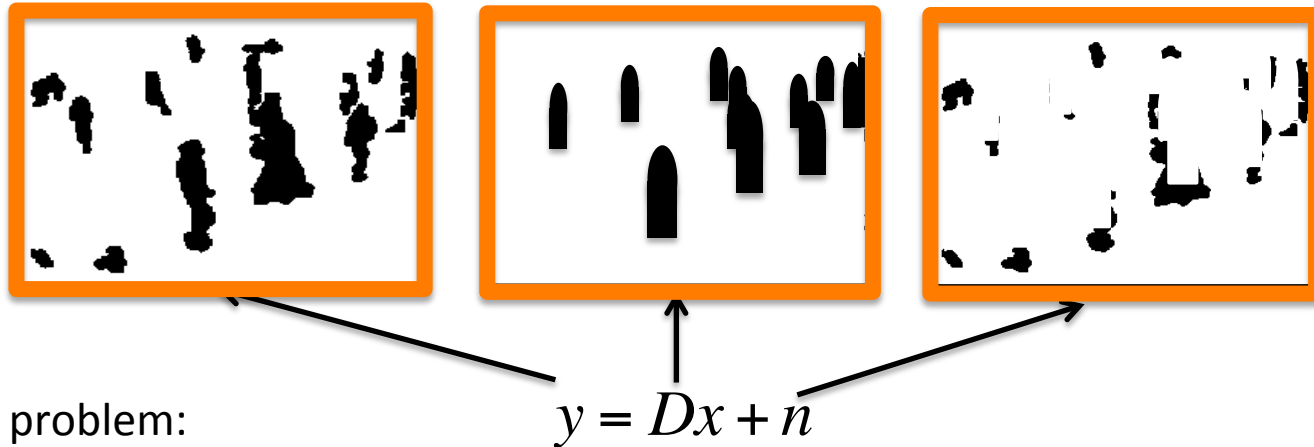
$$\min \|x\|_0 \quad \text{s. t.} \quad y = Dx + n$$

I. Detection: In praise of sparsity

“Creation is based on small number of primary, indivisible elements that combine with one another according to a few simple patterns.” [1]



I. Detection: Sparsity driven framework



- Inverse problem:

$$y = Dx + n$$

- Sparsity prior: $\min \|x\|_0$ s. t. $y = Dx + n$
- Basis Pursuit [1]: $\min \|x\|_1$ s. t. $y = Dx$
- BPDN: $\min \|x\|_1$ s. t. $\|y - Dx\| \leq n$
- Lasso: $\min \|y - Dx\|_2$ s. t. $\|x\|_1 \leq \varepsilon$



I. Detection: Pedestrian localization



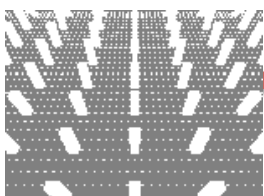
Input Sequence

	H_1	...	H_l
α_1		...	
α_2		...	
...
α_k		...	

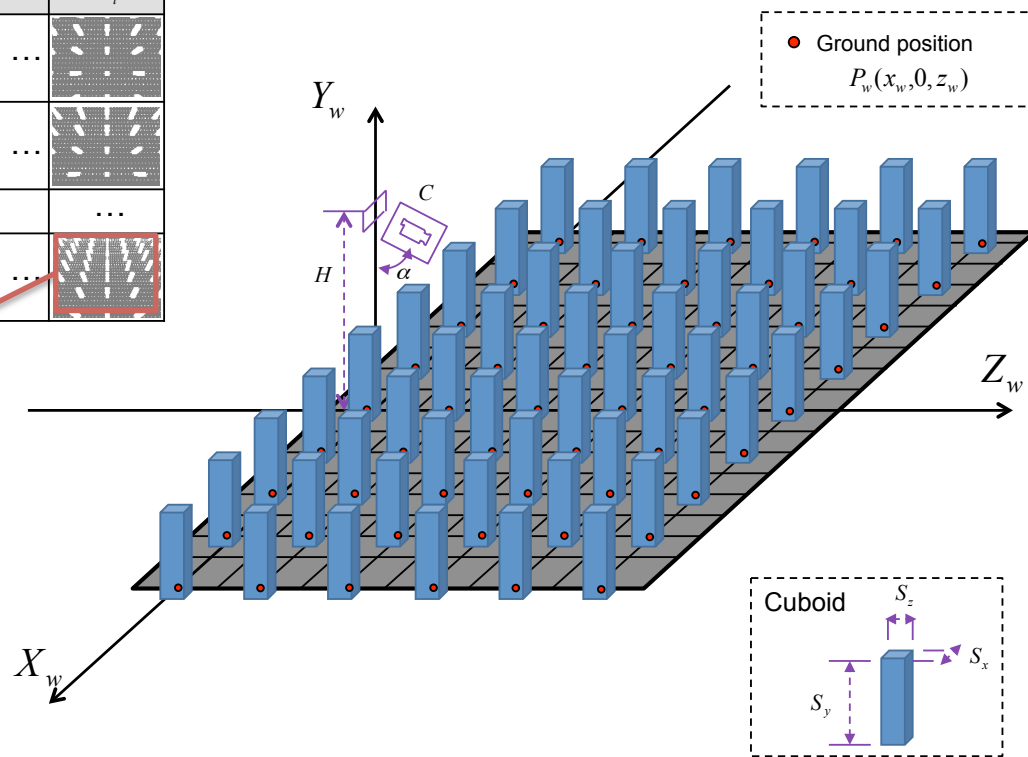
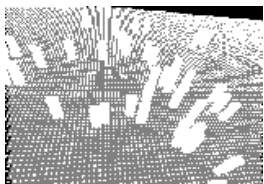
▪ Foreground extraction



▪ Used Dictionary

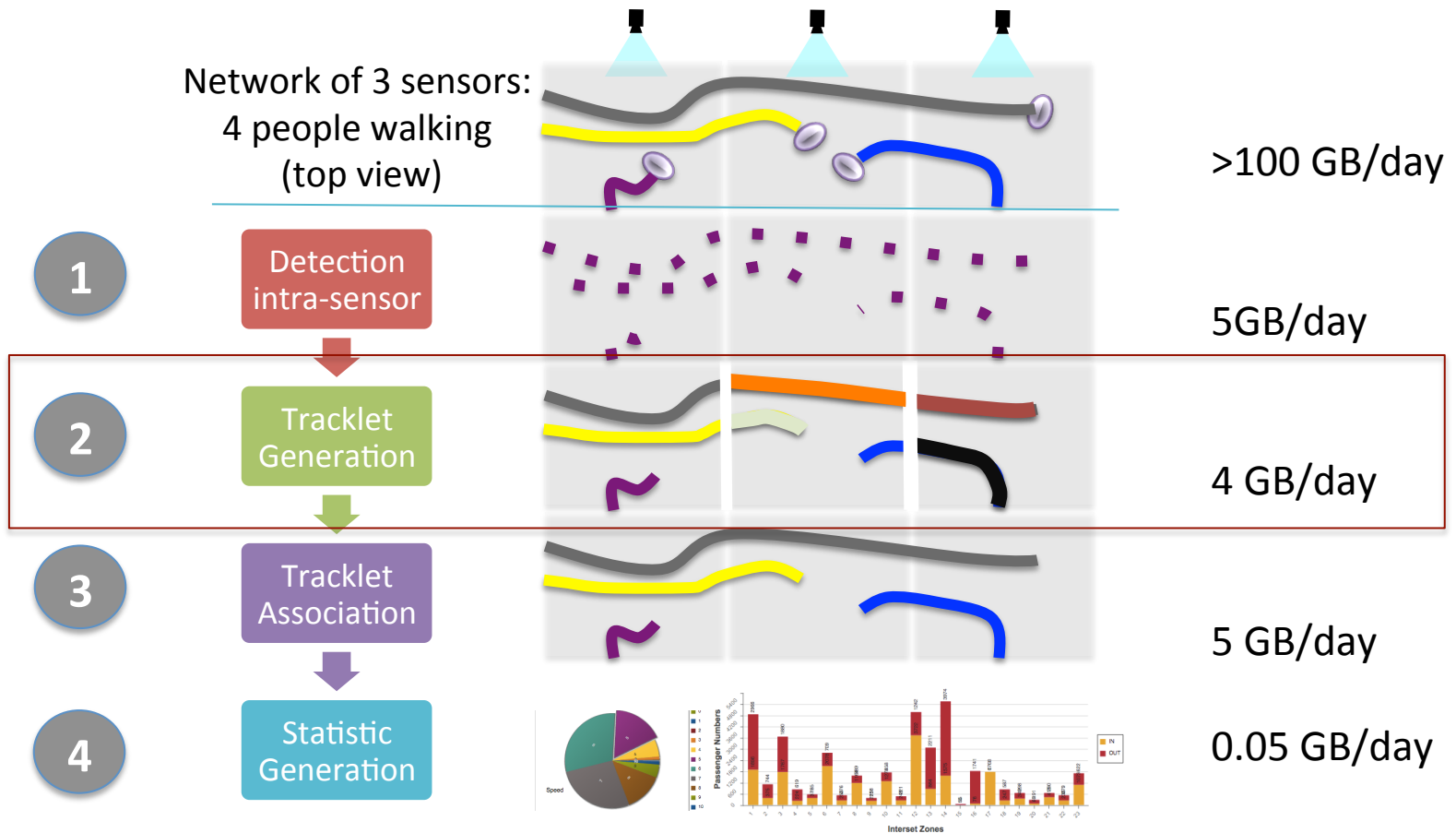


▪ Localization



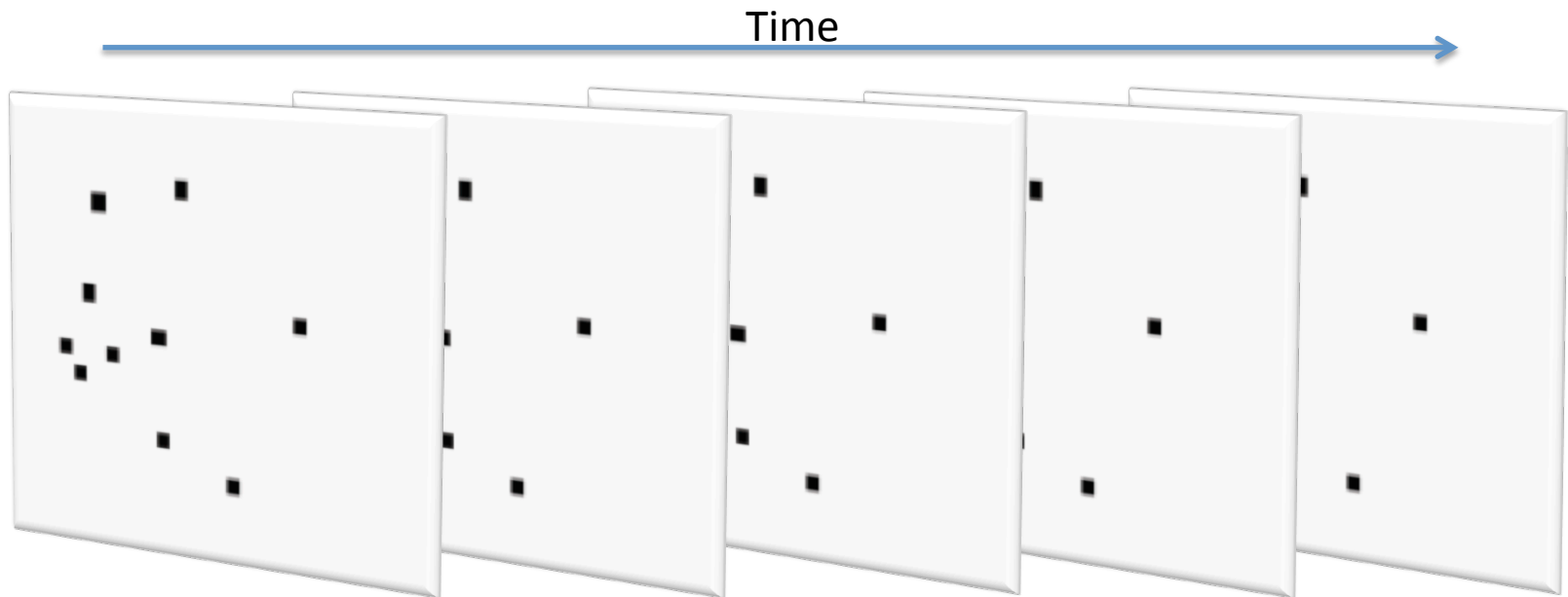
Outline:

From foreground extraction To tracking millions of pedestrians

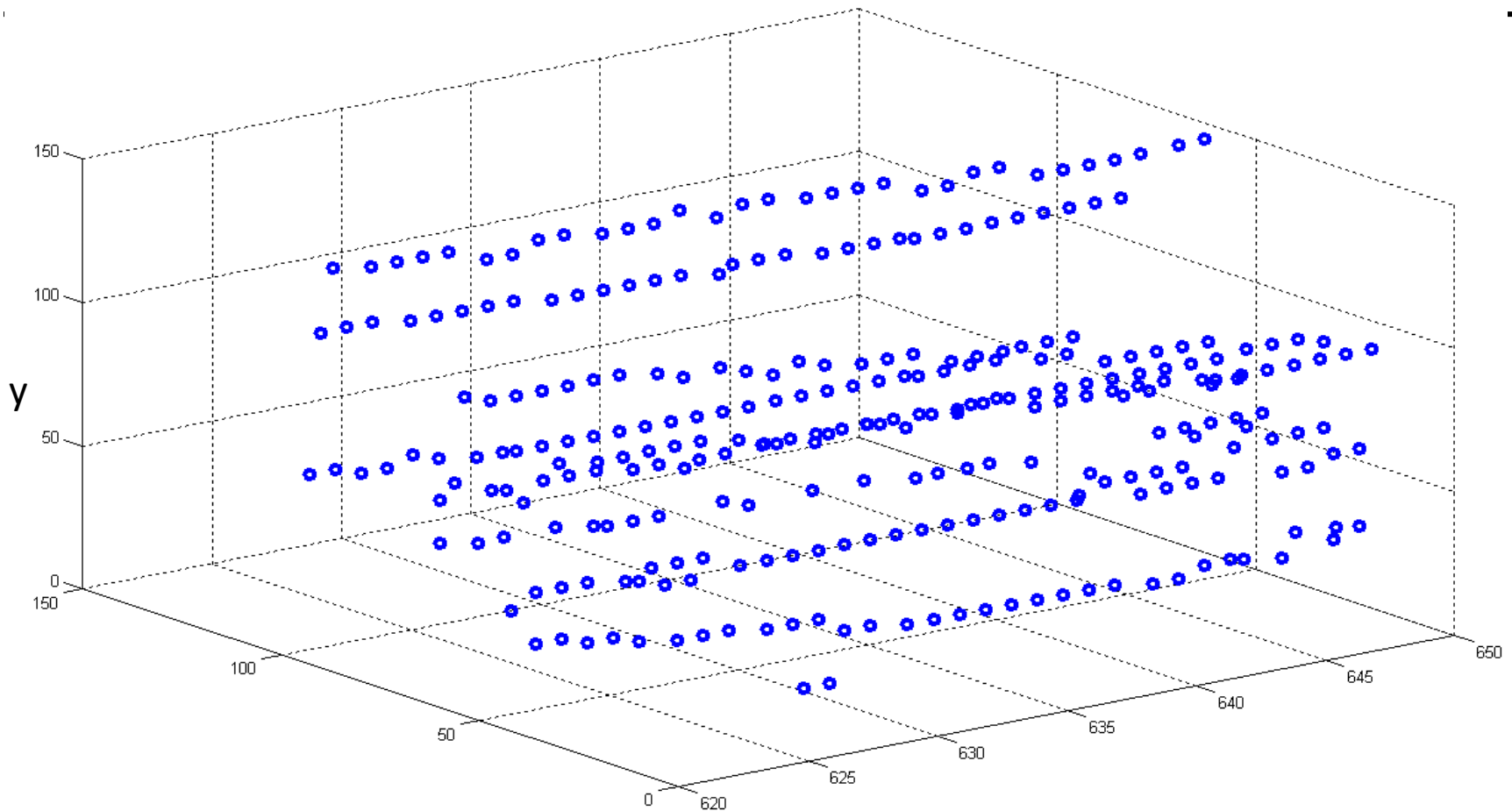


II. Tracklet generation: Data Association Problem

- Create a Directed Acyclic Graph $G = (N, E)$ where
 - N = The detected ground plane points across time
 - E = The connectivity cost between the detections (based on motion/appearance model)



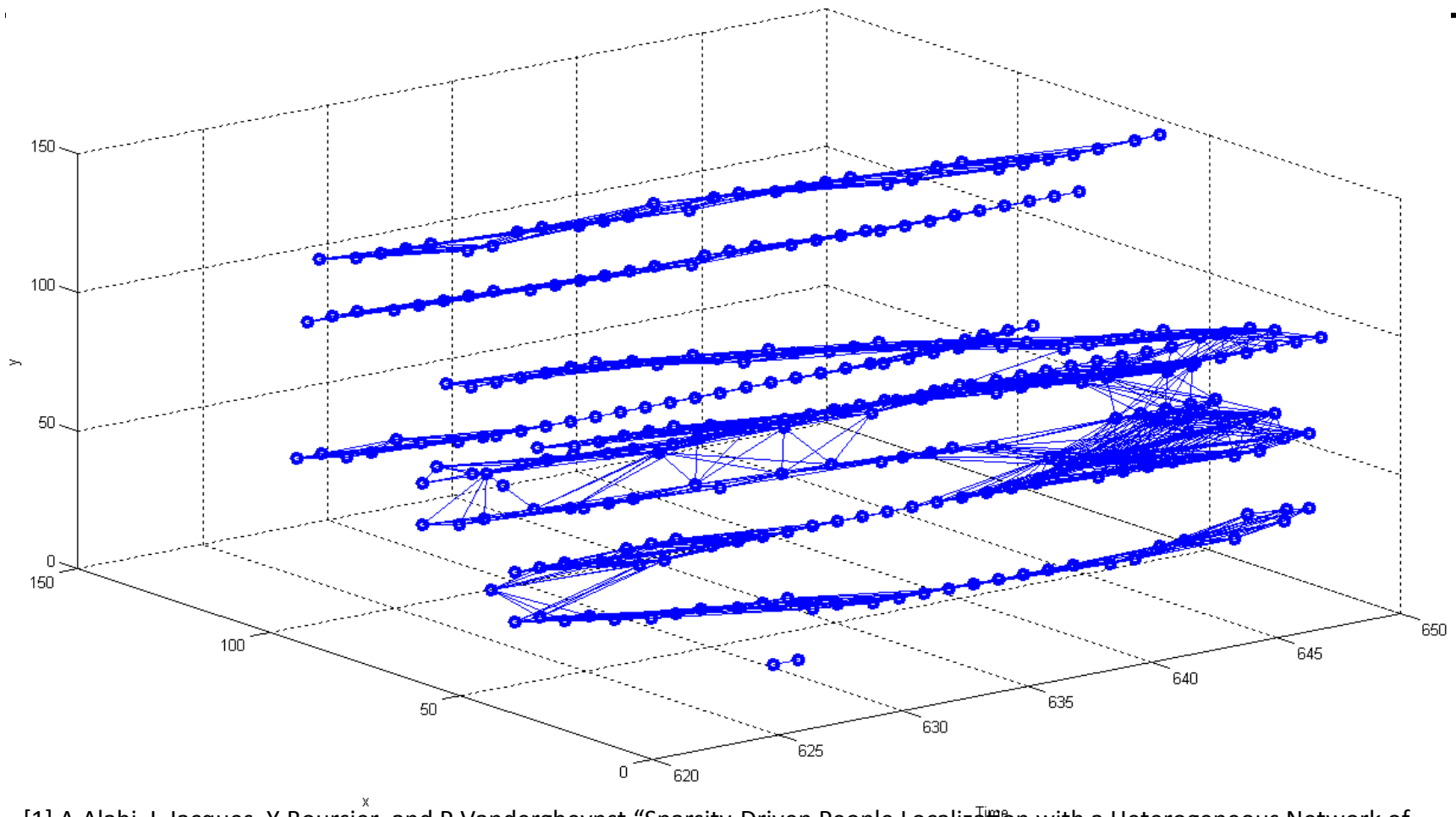
II. Tracklet generation: Data Association Problem



- [1] A. Alahi, L. Jacques, Y. Boursier, and P. Vanderghenst, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011
- [2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



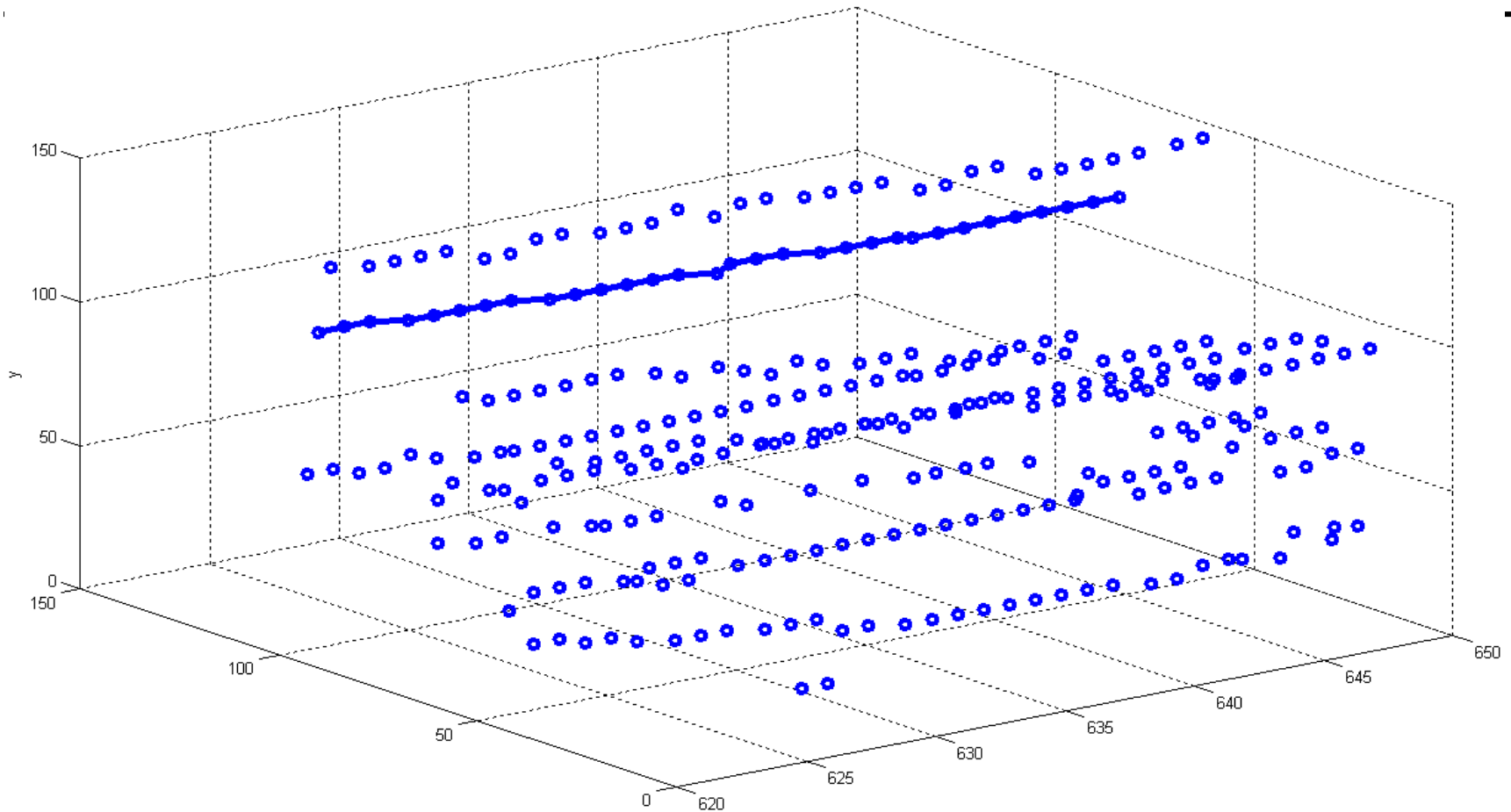
II. Tracklet generation: Create a DAG



- [1] A. Alahi, L. Jacques, Y. Boursier, and P. Vandergheynst, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011
- [2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



II. Tracklet generation: Select longest shortest path with smallest cost

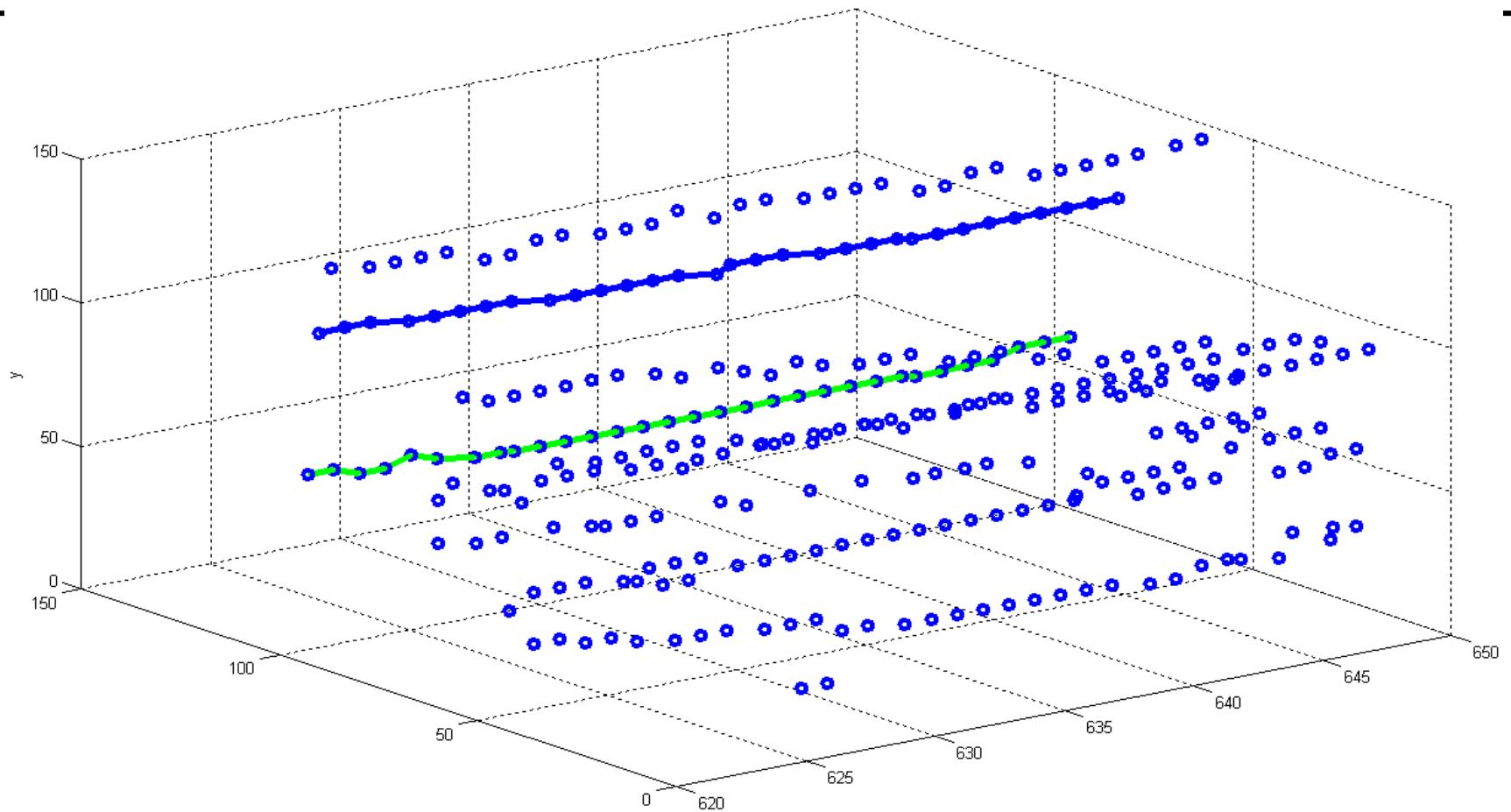


[1] A. Alahi, L. Jacques, Y. Boursier, and P. Vanderghyest, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011

[2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



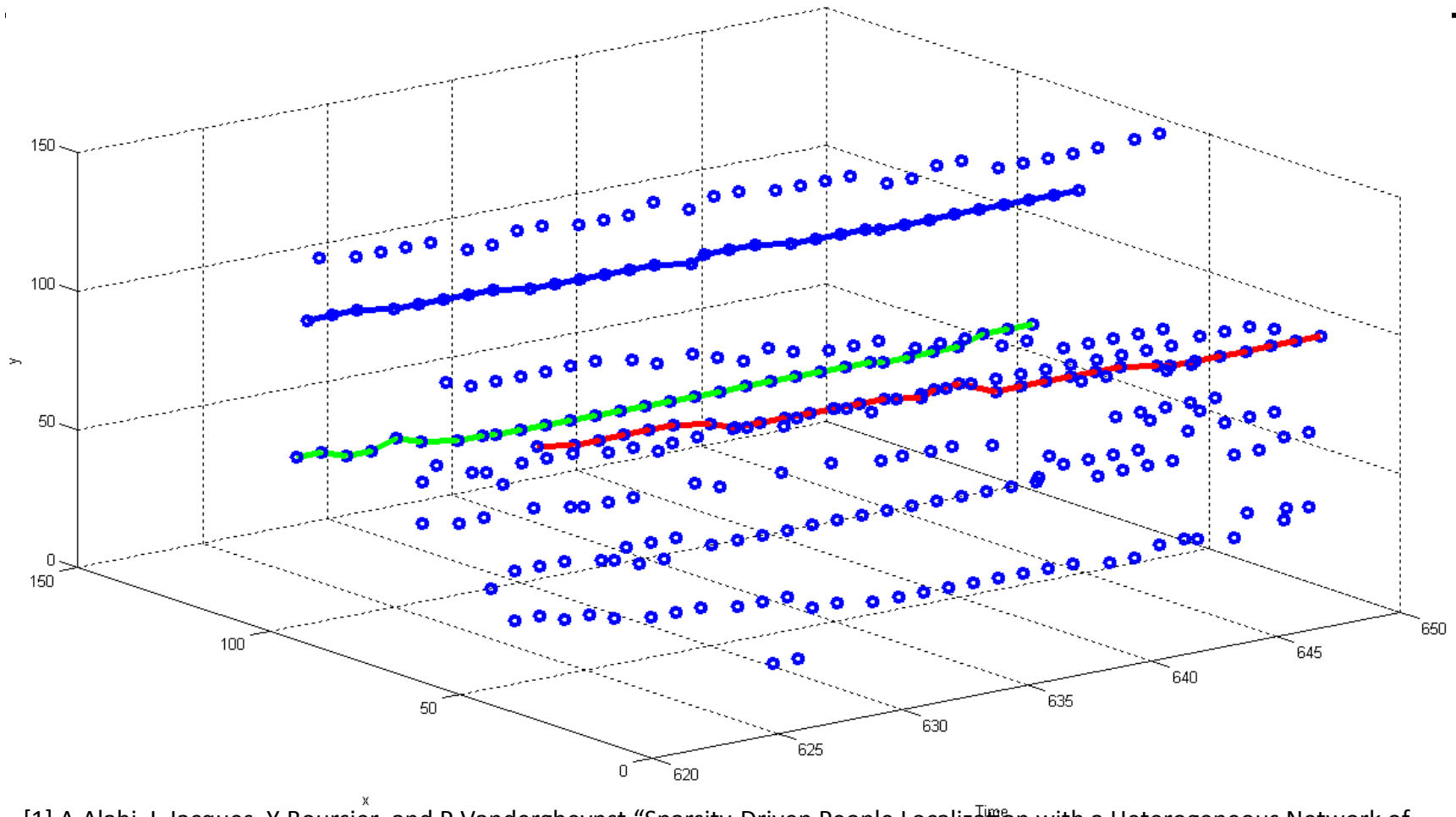
II. Tracklet generation: Iterate



- [1] A. Alahi, L. Jacques, Y. Boursier, and P. Vandergheynst, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011
- [2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



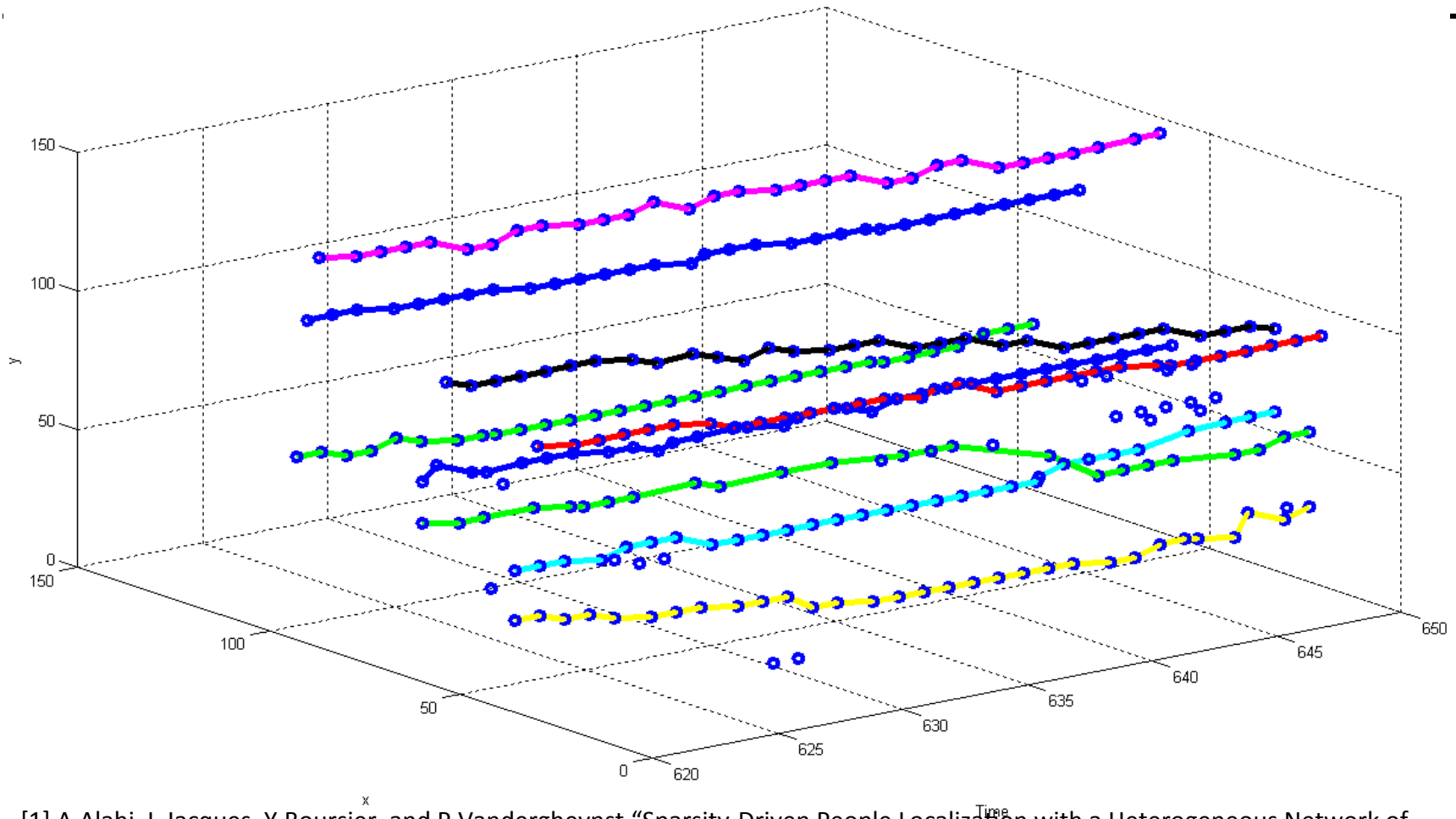
II. Tracklet generation: Iterate



- [1] A. Alahi, L. Jacques, Y. Boursier, and P. Vandergheynst, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011
- [2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



II. Tracklet generation: Till no more paths



- [1] A. Alahi, L. Jacques, Y. Boursier, and P. Vanderghelynst, "Sparsity-Driven People Localization with a Heterogeneous Network of Cameras," *Journal of Mathematical Imaging and Vision*, 2011
- [2] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua, "Multi-Camera People Tracking With a Probabilistic Occupancy Map", *PAMI* 2008



II. Tracklet Generation: Edge cost

- Motion model with social interactions
- Appearance model



II. Tracklet generation: Modeling social interactions



$$\mathbf{F}_i = \mathbf{F}_i^{Goal} + \mathbf{F}_i^{Avoidance} + \mathbf{F}_i^{Attraction} + \mathbf{F}_i^{Scene}$$

$$\mathbf{F}_i^{Avoidance} = \sum_{j \in P \setminus i} \mathbf{f}_{j \rightarrow i}^{Avoidance},$$

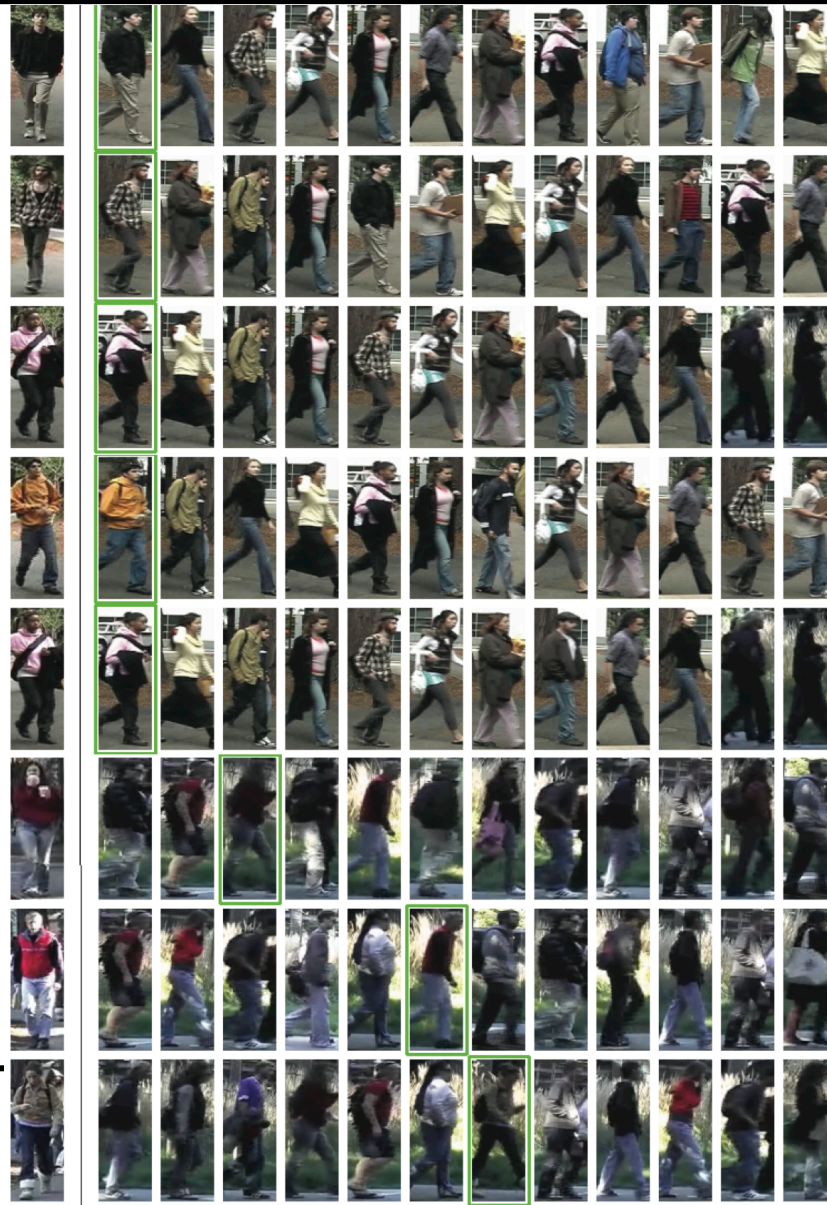
where

$$\mathbf{f}_{j \rightarrow i}^{Avoidance} = \alpha e^{\frac{d_p - d_{ij}}{\beta}} \mathbf{n}_{j \rightarrow i}$$

$$\frac{d}{dt} \mathbf{v} = \frac{\mathbf{F}_i}{m},$$

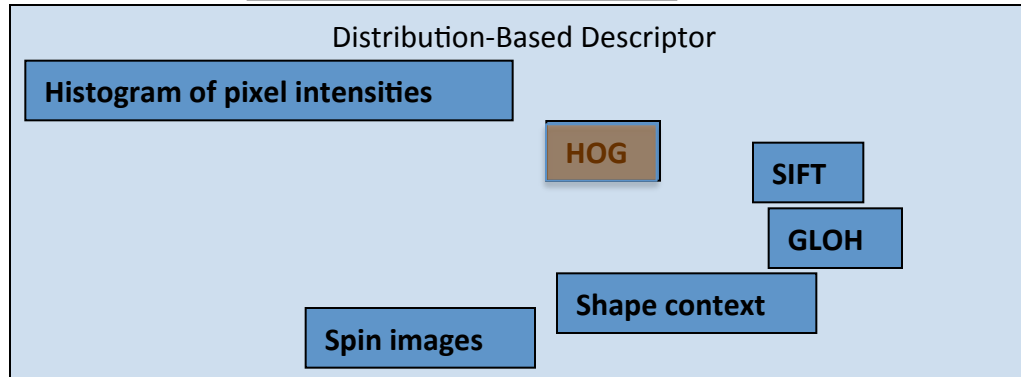


II. Tracklet Generation: Modeling appearance cues



II. Tracklet Generation: An arm-race of image descriptors

Vector of pixel intensities

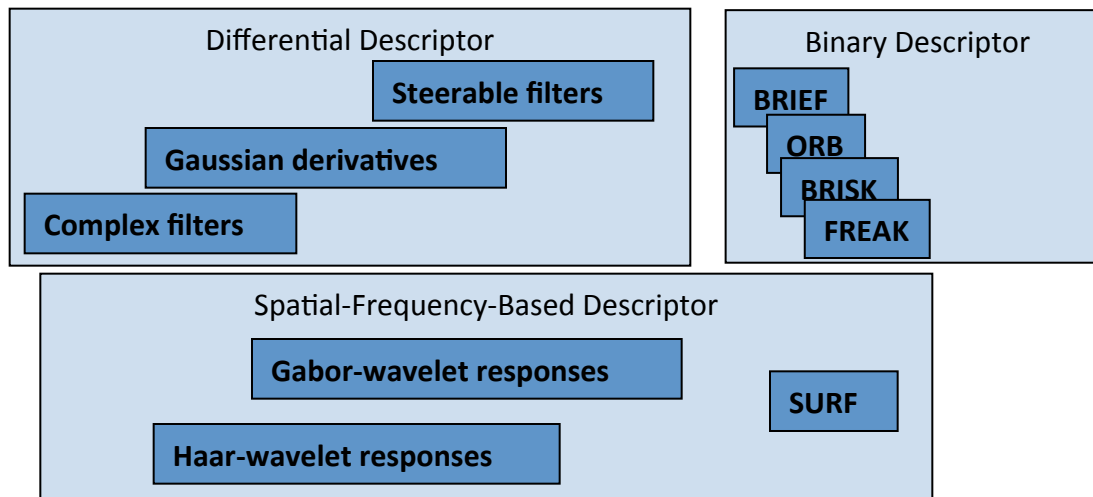


Moment invariants

Covariance of set of features

Low Performance

High Performance



[1] Gabriel, P., Hayet, J., Piater, J., Verly, J.: Object tracking using color interest points

[2] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection,"

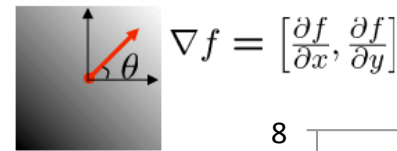
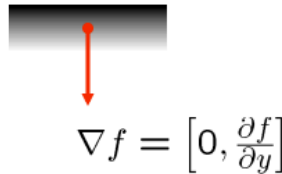
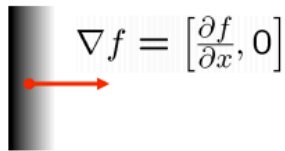
[3] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification,"



II. Tracklet Generation: HOG

Image gradient

- The gradient of an image: $\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$



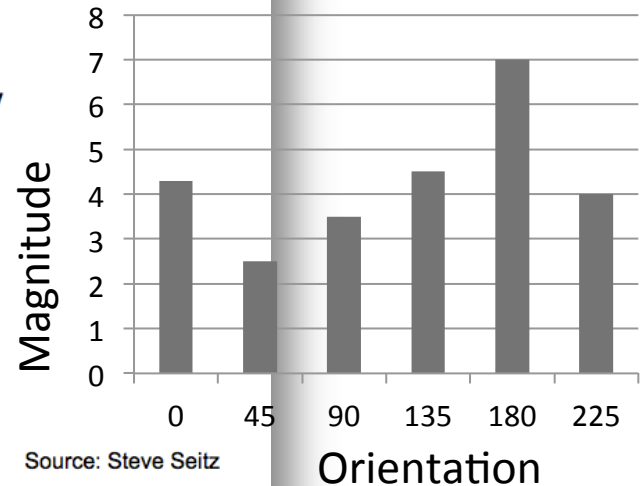
The gradient points in the direction of most rapid increase in intensity

The gradient direction is given by $\theta = \tan^{-1} \left(\frac{\partial f / \partial y}{\partial f / \partial x} \right)$

- how does this relate to the direction of the edge?

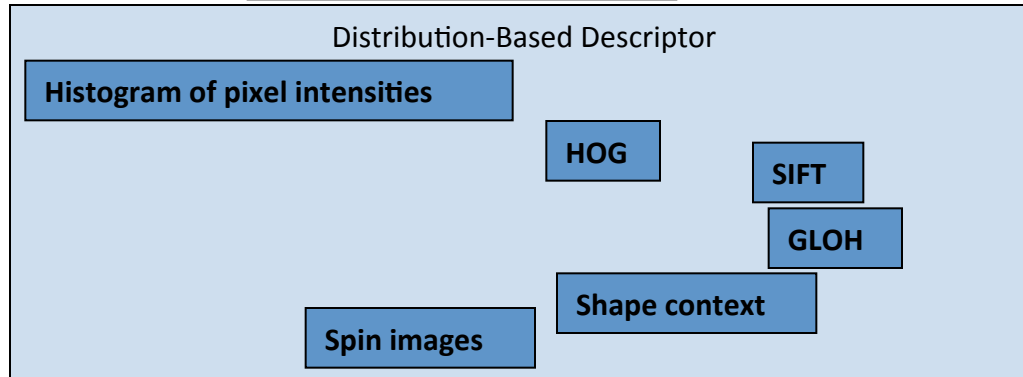
The *edge strength* is given by the gradient magnitude

$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$



II. Tracklet Generation: An arm-race of image descriptors

Vector of pixel intensities

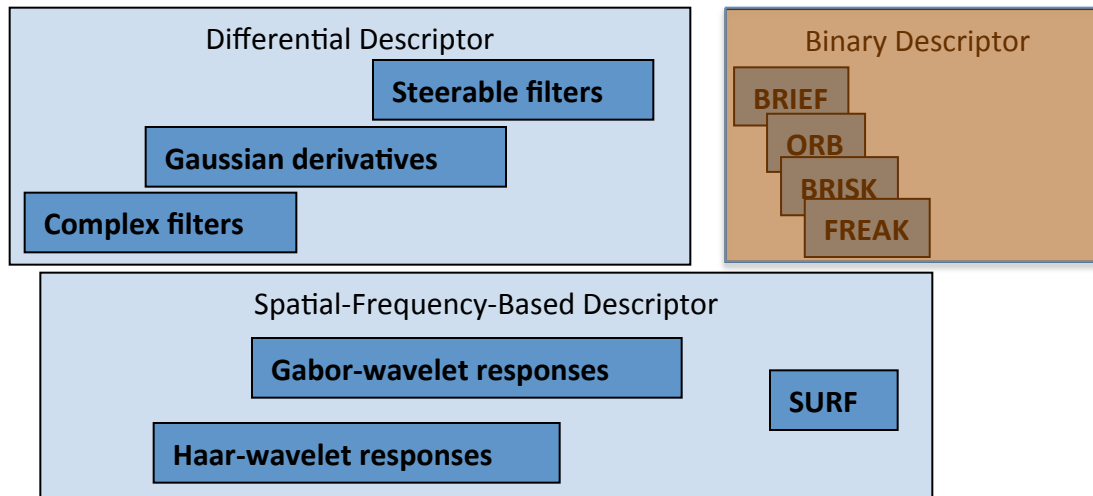


Moment invariants

Covariance of set of features

Low Performance

High Performance



[1] Gabriel, P., Hayet, J., Piater, J., Verly, J.: Object tracking using color interest points

[2] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection,"

[3] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification,"



II. Tracklet Generation: Binary descriptors



II. Tracklet Generation: BRIEF[1] / ORB[2] / BRISK[3]



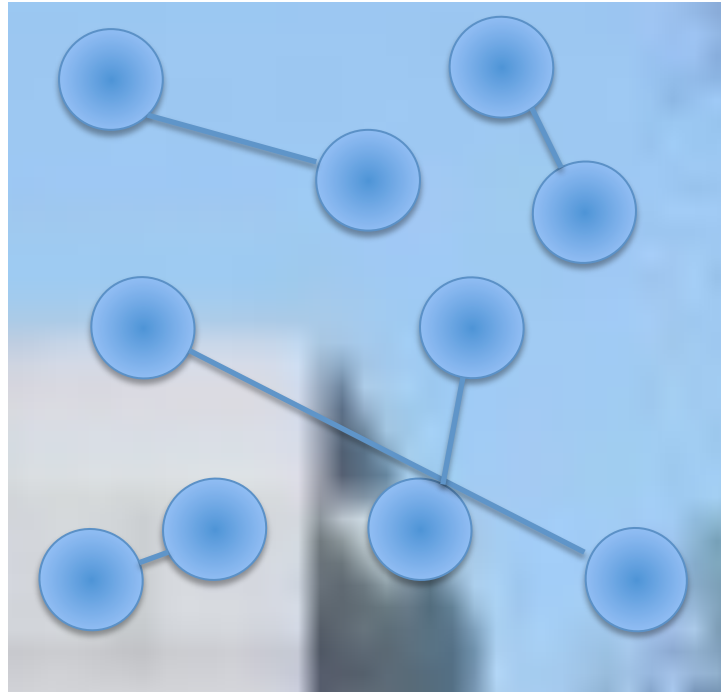
[1] Calonder, Michael, et al. "Brief: Binary robust independent elementary features." *ECCV 2010*.

[2] Rublee, Ethan, et al. "ORB: an efficient alternative to SIFT or SURF." *ICCV 2011*

[3] Leutenegger, ., et al. "BRISK: Binary robust invariant scalable keypoints." *ICCV 2011*



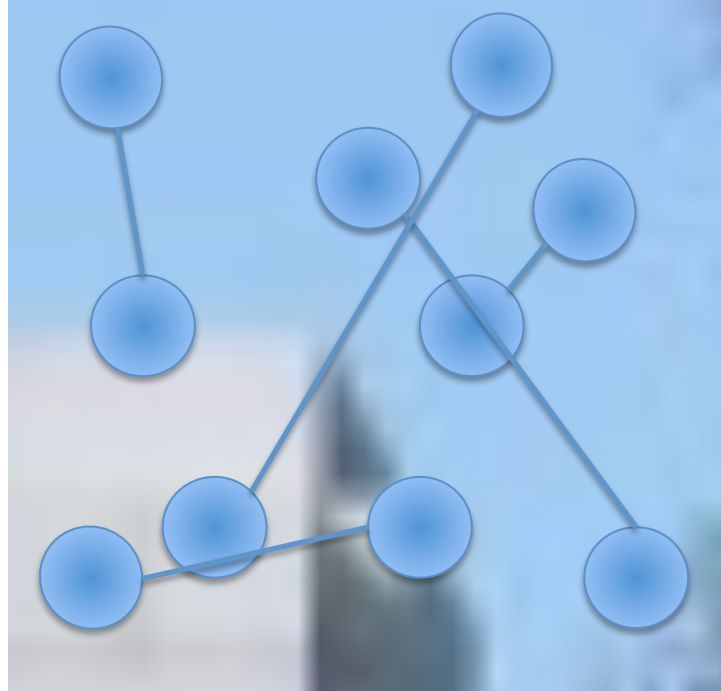
II. Tracklet Generation: **BRIEF**[1] / ORB[2] / BRISK[3]



A sequence of 1-bit DoG

II. Tracklet Generation: BRIEF[1] / **ORB**[2] / BRISK[3]

- Select
- 1) Most discriminant**
- AND
- 2) Less correlated**



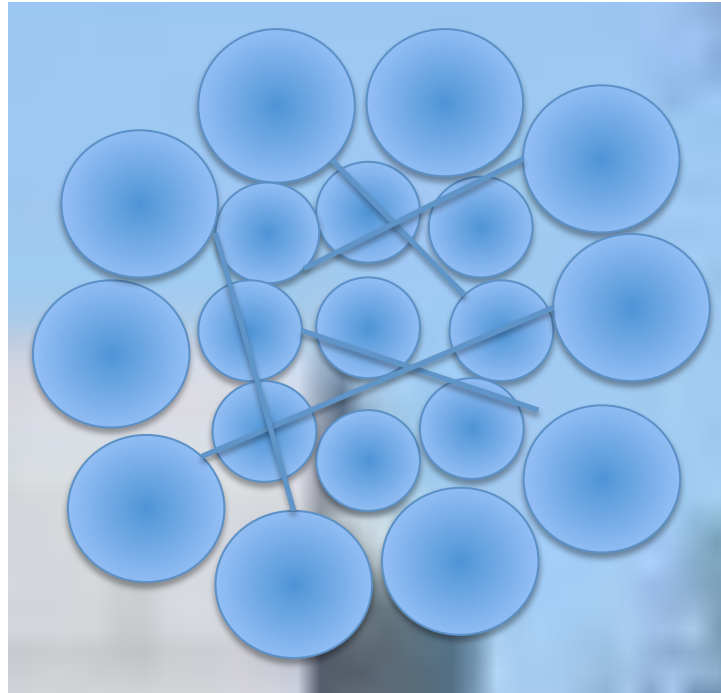
[1] Calonder, Michael, et al. "Brief: Binary robust independent elementary features." *ECCV 2010*.

[2] Rublee, Ethan, et al. "ORB: an efficient alternative to SIFT or SURF." *ICCV 2011*

[3] Leutenegger, ., et al. "BRISK: Binary robust invariant scalable keypoints." *ICCV 2011*



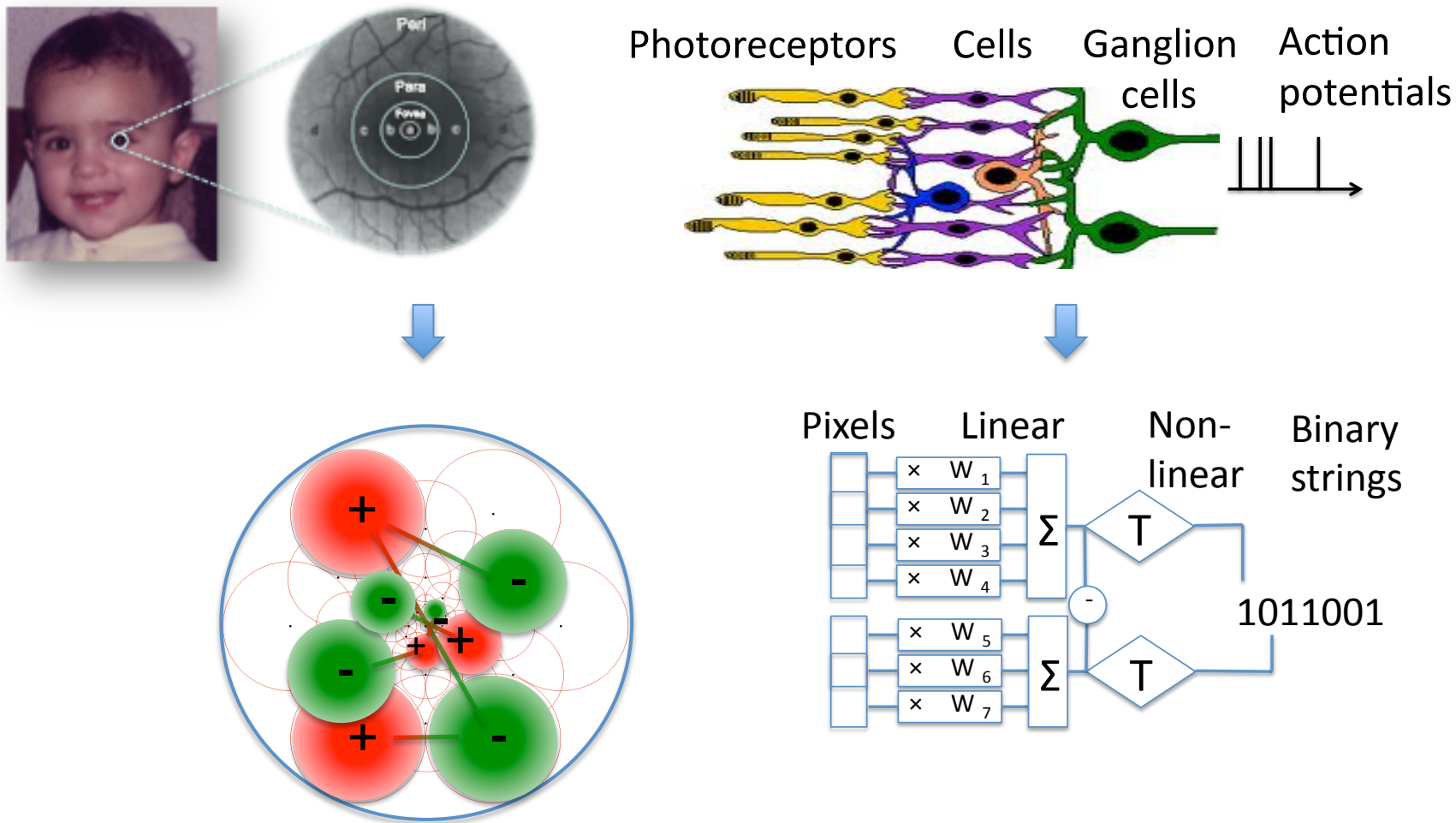
II. Tracklet Generation: BRIEF[1] / ORB[2] / BRISK[3]



-
- [1] Calonder, Michael, et al. "Brief: Binary robust independent elementary features." *ECCV 2010*.
[2] Rublee, Ethan, et al. "ORB: an efficient alternative to SIFT or SURF." *ICCV 2011*
[3] Leutenegger, . , et al. "BRISK: Binary robust invariant scalable keypoints." *ICCV 2011*



II. Tracklet Generation: Retina-inspired [1]



[1] Alahi, et al. "Freak: Fast retina keypoint." *CVPR 2012*



II. Tracklet Generation: Saccadic search

Object of interest Target image to search Matched objects

10101010...1010101

512 bits descriptor

Distance map

Distance map

10101010

010101

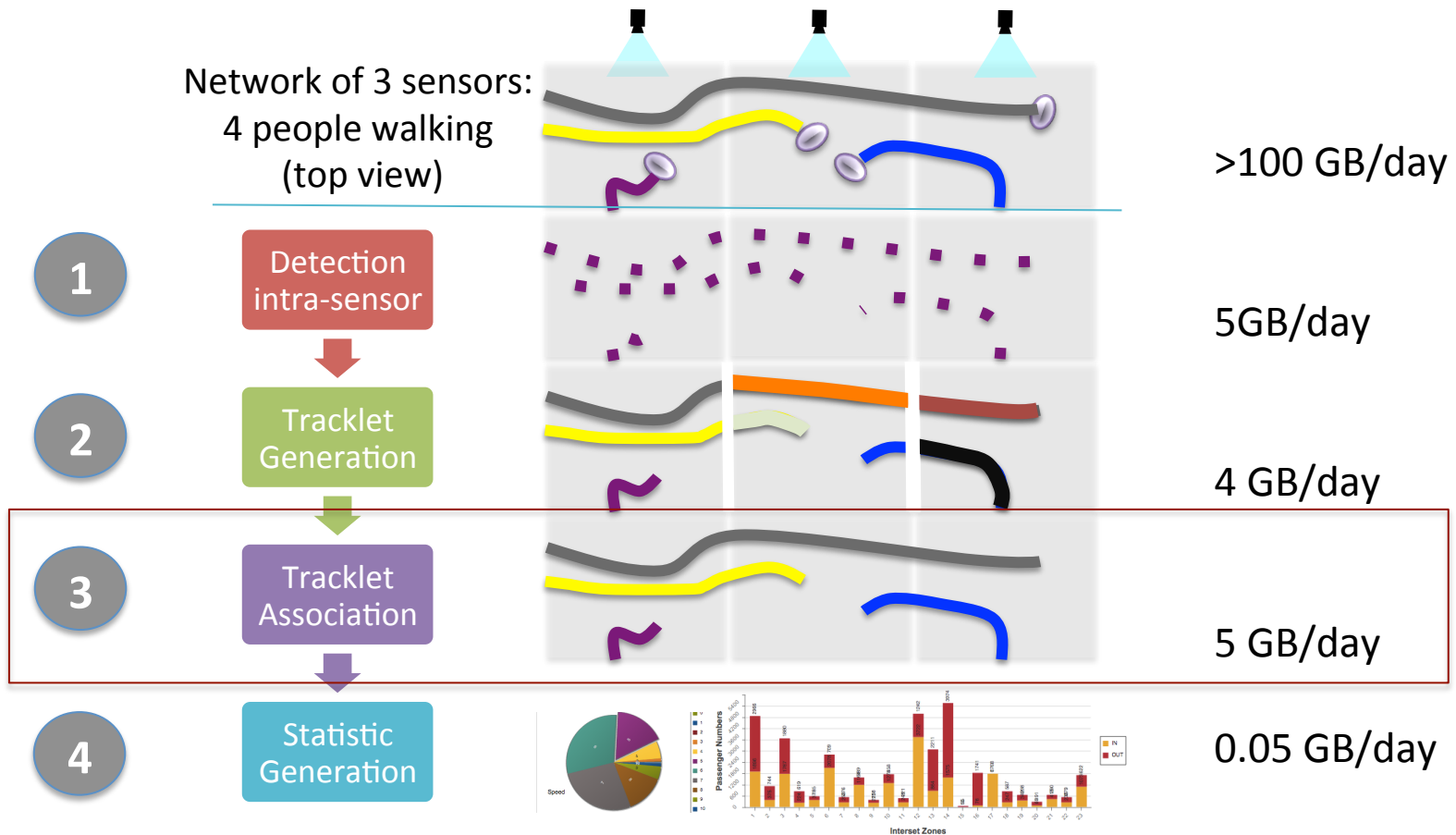
Filtering with first 128 bits Matching with last 128 bits

[1] Alahi, et al. "Freak: Fast retina keypoint." *CVPR 2012*

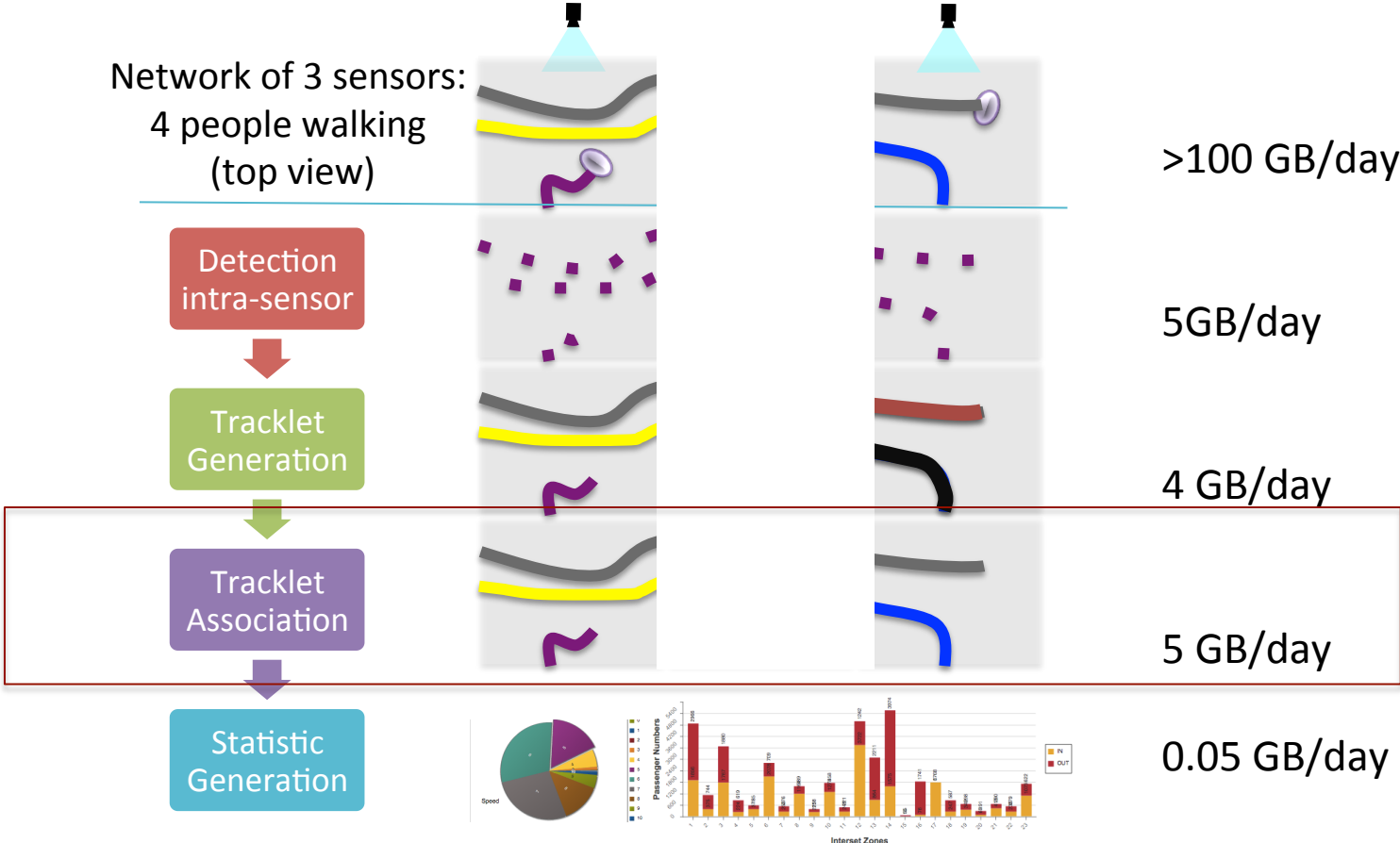


Outline:

From foreground extraction To tracking millions of pedestrians



Tracklet association in scattered network



III- Tracklet association: Problem formulation

Let:

- \mathbf{T} : all long term trajectories
- \mathbf{t} : tracklets (tracklets capture within each camera)
- Problem: Maximizing the a posteriori probability (MAP) of \mathbf{T} :

$$\mathbf{T}^* = \arg \max_{\mathbf{T}} P(\mathbf{T} | \mathbf{t}) \quad (1)$$

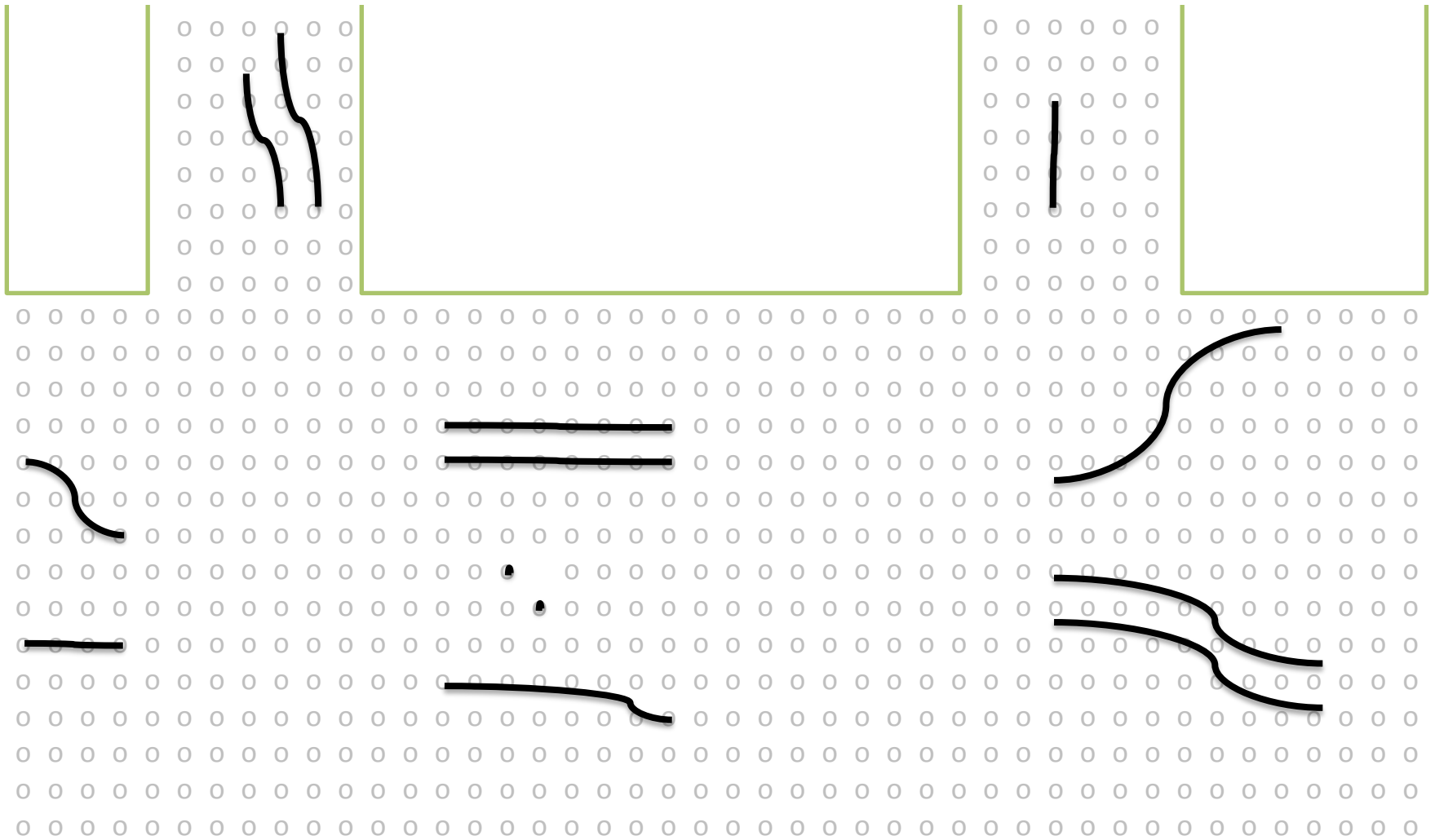
$$= \arg \max \prod_i P(\mathbf{t}_i | \mathbf{T}) P(\mathbf{T}) \quad (2),$$

where $P(\mathbf{T}) = \prod_k P(\mathbf{T}_k)$ (since trajectories do not overlap)

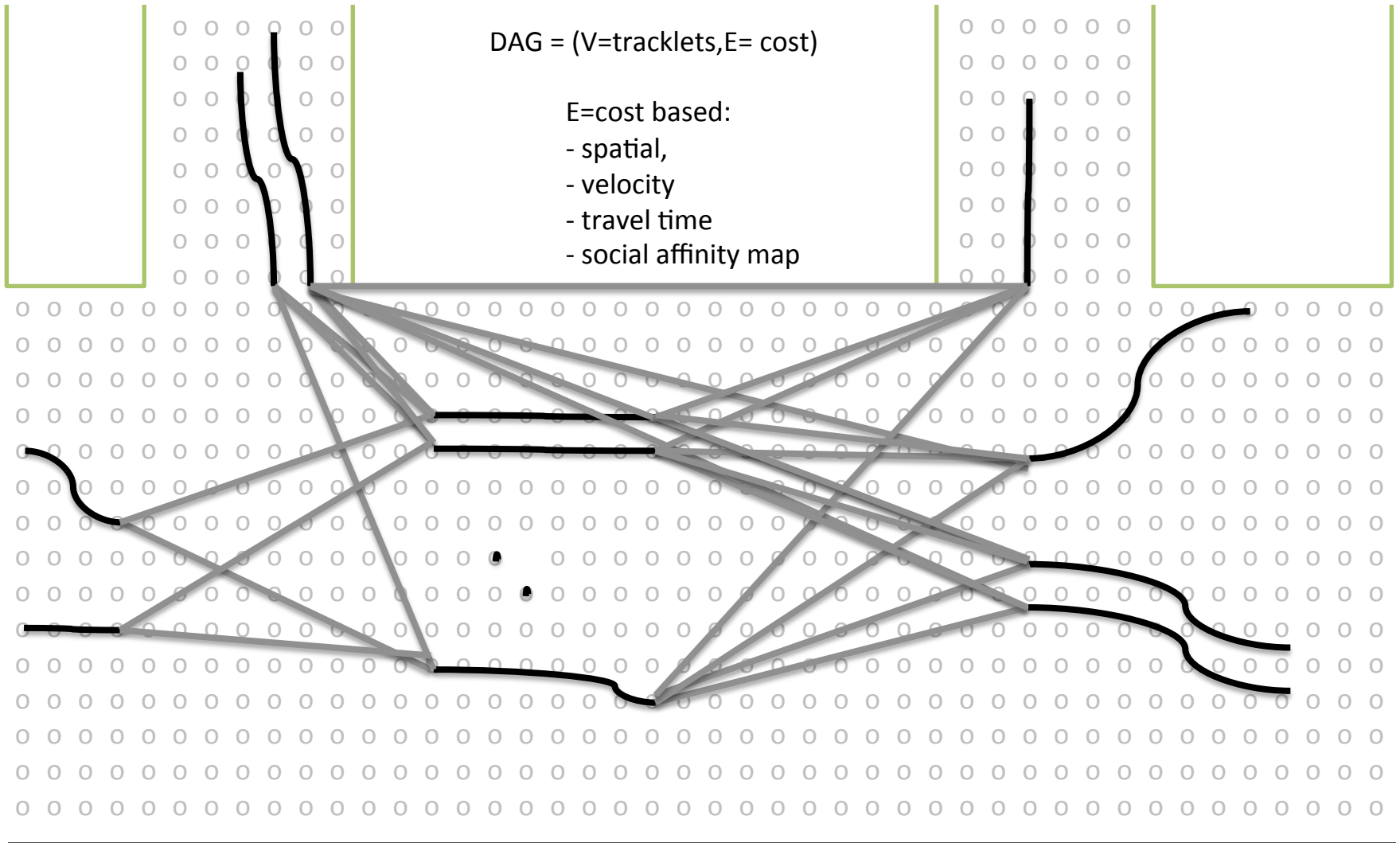
$P(\mathbf{T}_k) = P(\mathbf{t}_k^s) \dots P(\mathbf{t}_k^t | \mathbf{t}_k^{t-1}) P(\mathbf{t}_k^e)$ (markov chain)



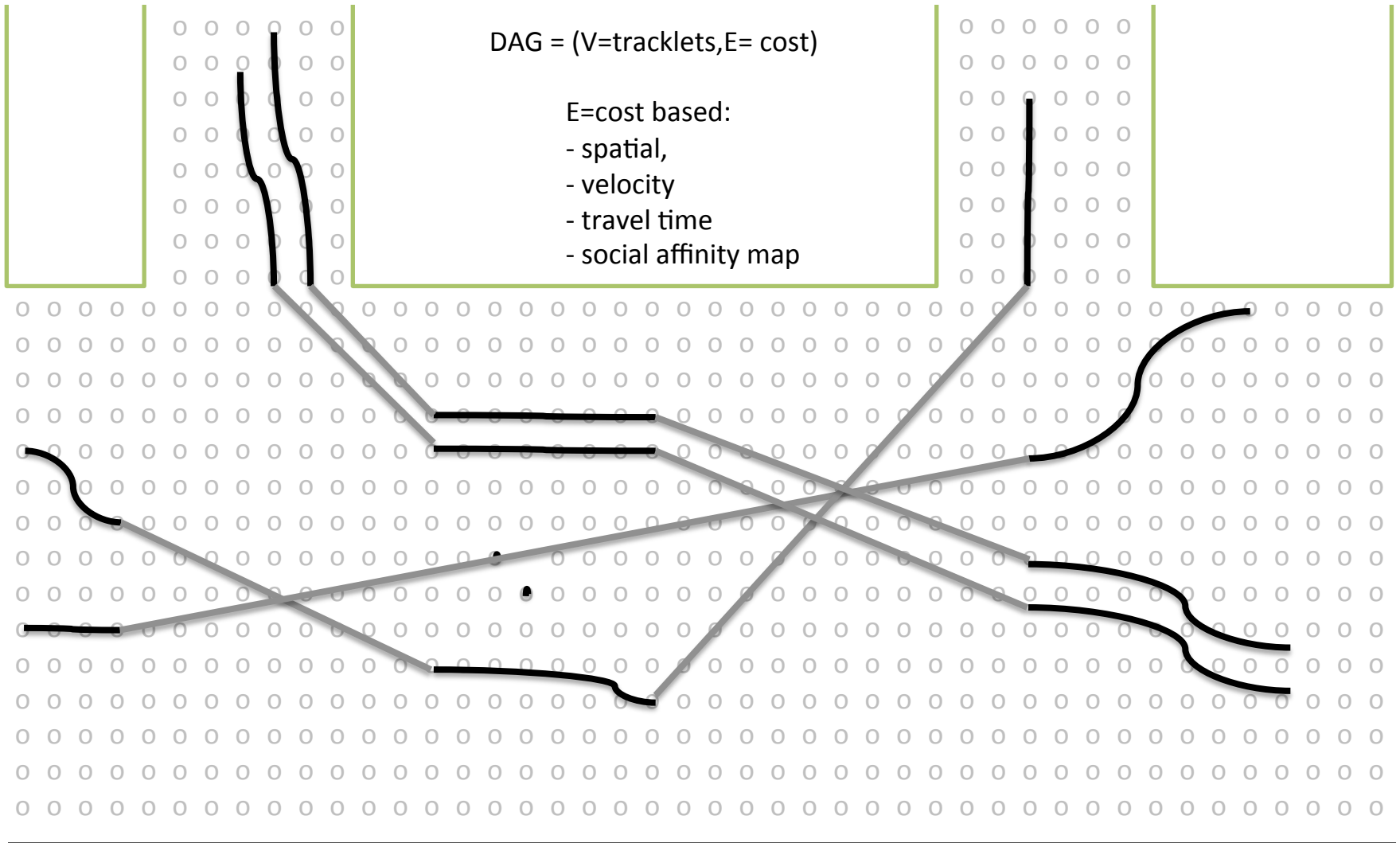
III. Tracklet association (Top view)



III. Tracklet association (Top view)



III. Tracklet association (Top view)



Conclusion

A new dimension to “Google Analytics”:
Analyzing people outside of website

