

Person Tracking and Event Recognition over A Sensor Network



Queen's University
Belfast

CSIT CENTRE
FOR SECURE
INFORMATION
TECHNOLOGIES

Persona de seguimiento y reconocimiento de eventos en una red de sensores



Queen's University
Belfast



- Motivation & Introduction
- Tracking over a sensor network
- Gender Profiling
- Multi agent surveillance architecture
- Conclusion & Summary

- "CCTV was originally seen as a preventative measure. Billions of pounds has been spent on kit, but no thought has gone into how the police are going to use the images and how they will be used in court. It's been an utter fiasco: only 3% of crimes were solved by CCTV. There's no fear of CCTV. Why don't people fear it? [They think] the cameras are not working."

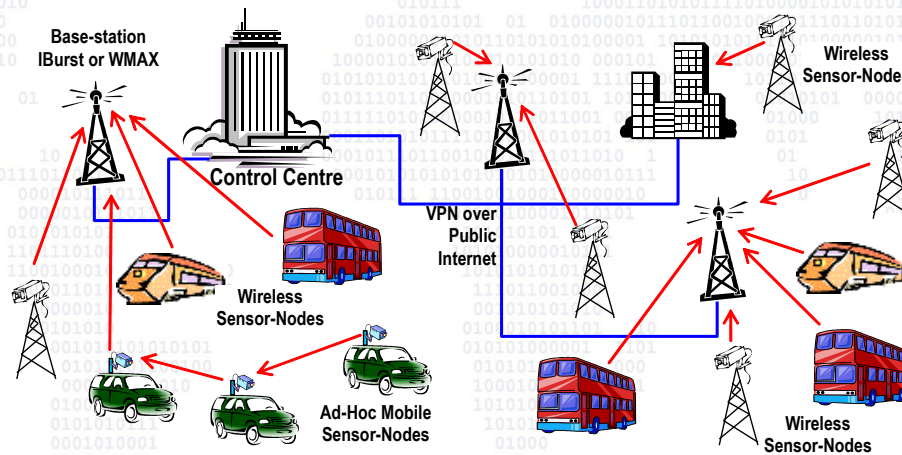
Detective Chief Inspector Mick Neville, officer in charge of the Metropolitan police unit, talking at the Security Document World Conference in London.

- “I want public transport to become world renowned for its safety and to banish the sad minority of hoodlums and trouble-makers that have blighted our buses.”
Boris Johnson, Mayor of London

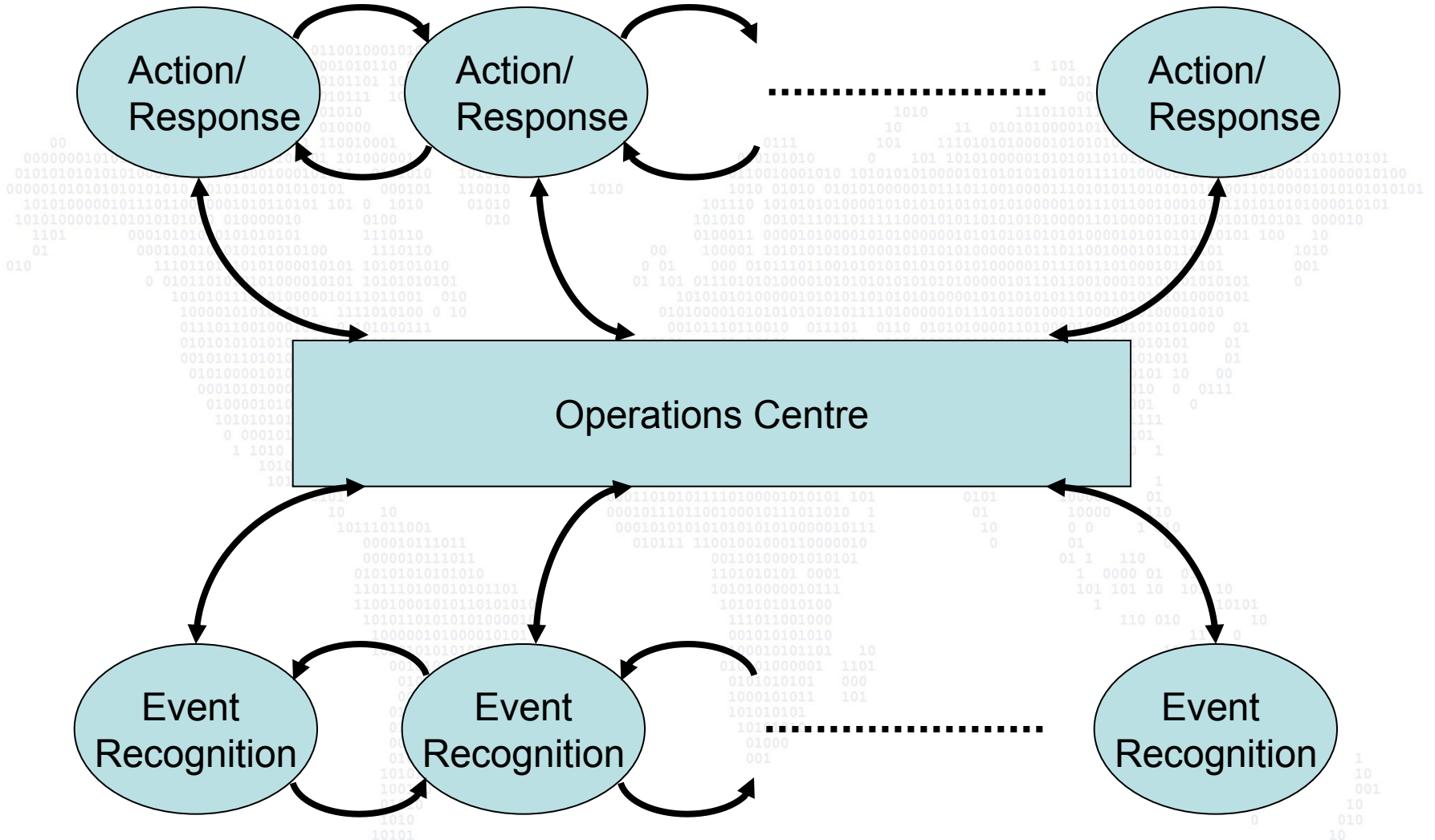
- Massive investment in CCTV in the UK
- Impact on anti-social and criminal behaviour has been minimal
 - e.g. assaults on public transport
- CCTV operates in a passive mode
- “Active” CCTV has to alert security analysts to prevent undesirable behaviour
- Greatly increase the likelihood of being caught - a major factor in crime prevention
- Persistent analysis of CCTV video footage in real-time

Secure Corridors – CCTV for buses

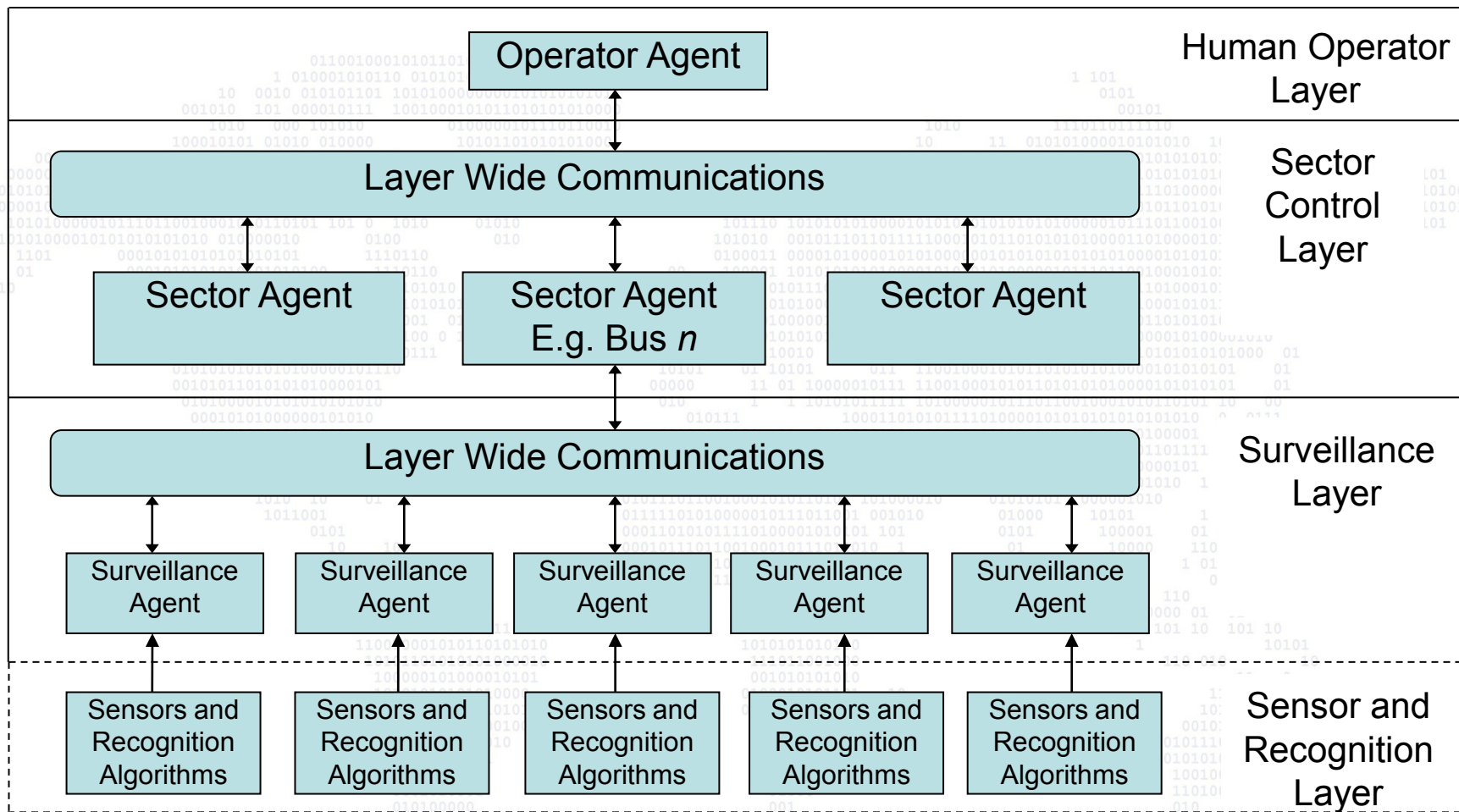
- Crime on transport platforms is a major anti-social form of crime in UK
- Aim is to reduce bus driver and passenger assaults



Event Management & Response



Architecture for Multi-Agent Surveillance



- Key requirement for active CCTV is to automatically determine the threat posed by each individual
- Focus of the computer vision community has been on behavior/action recognition
- Experienced security analysts profile individuals in the scene to determine their threat

- They identify individuals who look as though they may cause trouble
- Vast majority of offenders are young adolescent males
- Key to automatic threat assessment is:
 - to automatically measure the relative locations and motions of subjects in the scene
 - to automatically profile people in the scene based on their gender and age

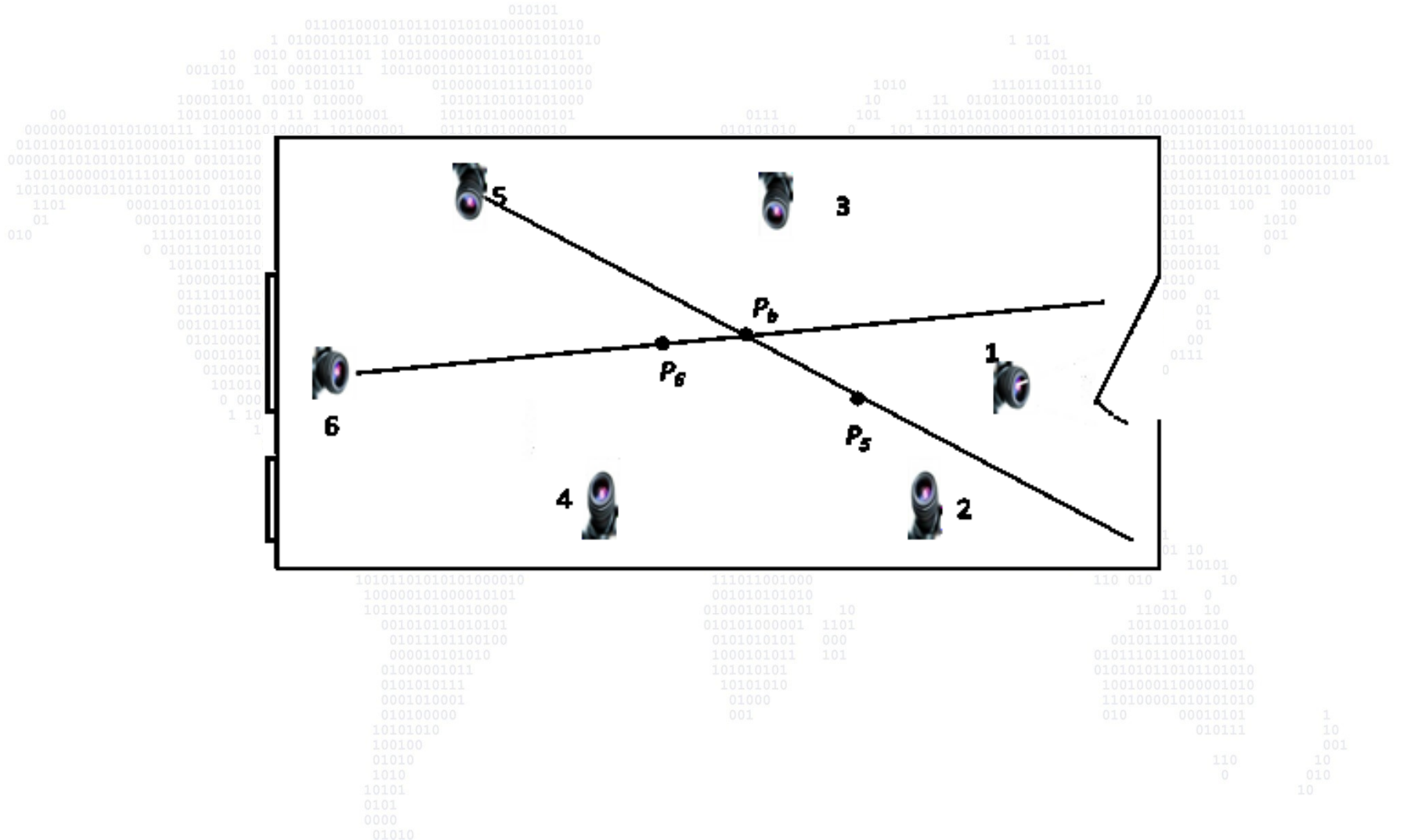
- Introduction & Motivation
- Tracking over a sensor network
- Gender Profiling
- Multi agent surveillance architecture
- Conclusion & Summary



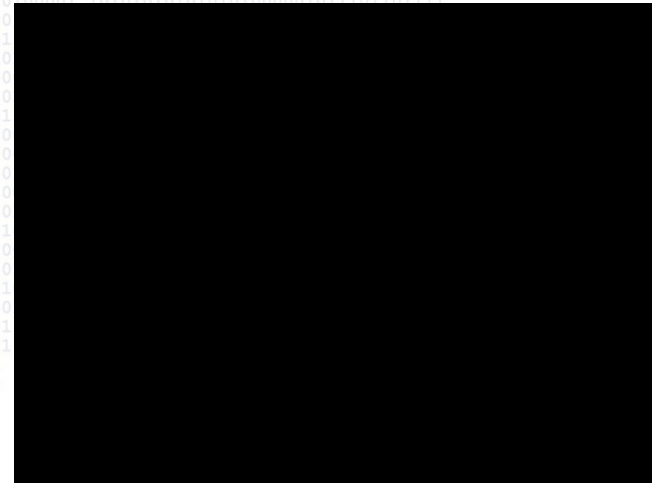
Open Spaces / Narrow Corridors



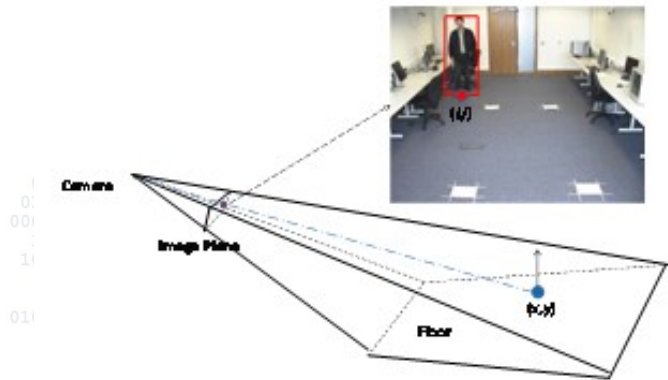
0110010001011010101010000101010
1 010001010110 0101010000101010101010
10 0010 010101101 101010000000101010101
001010 101 000010111 10010001010110101010000
1010 000 101010 0100000101110110010
100010101 01010 010000 10101101010101000
00 1010100000 0 11 110010001 1010101000010101
00000001010101010111 10101010100001 101000001 011101010000010
01101010
0 1
1010
0000
0101
1100
0100
0101
0101
1010
0101
0100
110
1100
1100
1010
0101
0101
1011
0101
011
01
010
0
00010101
010111
1
10
001
110
10
0
010
10
01010000
10101010
100100
01010
1010
1010
10101
0101
0000
01010
001
010
00010101
010111
1
10
001
110
10
0
010
10



- Moving object detection on buses likely to be poor
- Passenger obscuration by seats
- Cannot assume bottom of bounding box is coincident with floor



Single Camera

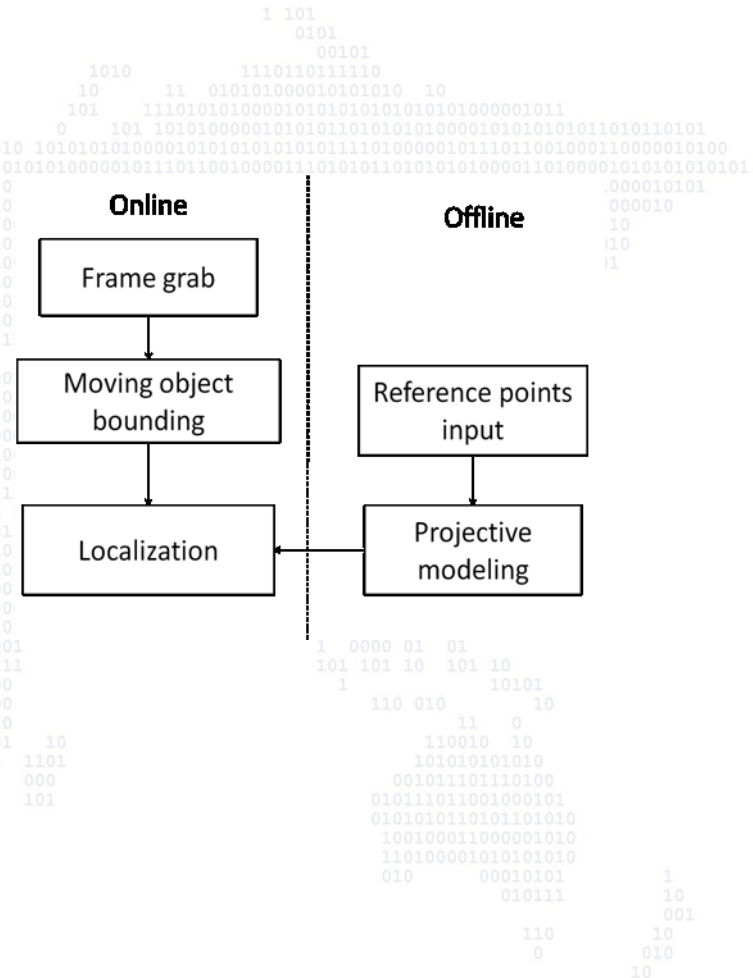


$$\begin{bmatrix} s_i \\ s_j \\ s \end{bmatrix} = \mathbf{Q} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

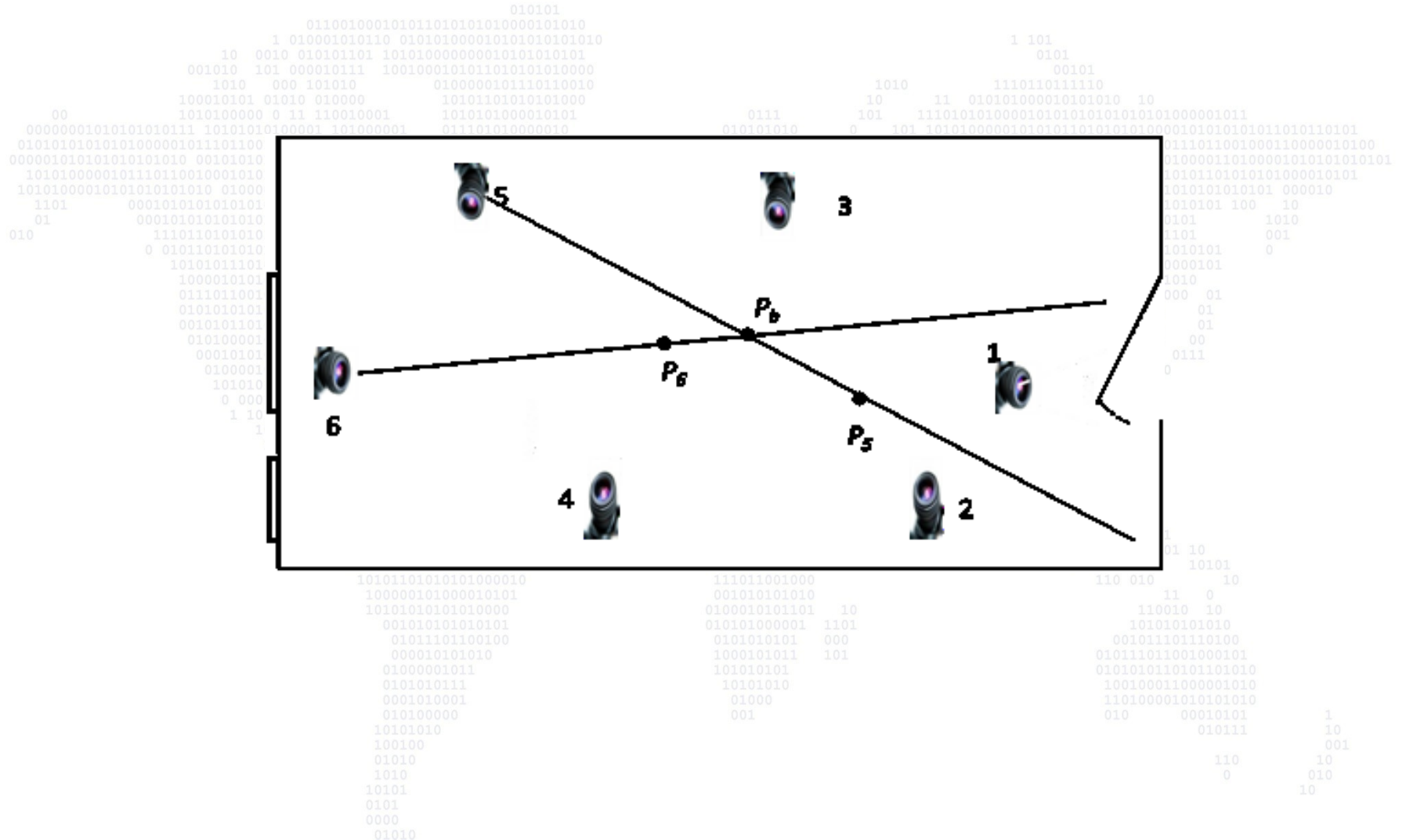
$$\mathbf{Q} = [\mathbf{R} \quad \mathbf{T}]$$

- Assumption: centre of detected object's bottom boundary is coincident with feet
- z value is zero, only need to find (x,y) coordinates
- Measure positions of four key points
- Solve for the eight unknowns in \mathbf{Q}

- Foreground extracted by background subtraction
- Bounding box placed around foreground
- Corner coordinates passed to localization module
- Solves for x and y



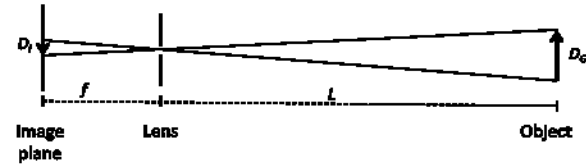
Two camera



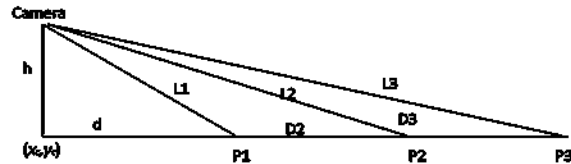
- For each camera determine its position in the room coordinates
- Determine the object position in the room coordinates for each camera
- Draw a line from the camera to object
- Object position is the intersection of lines

Two camera

- Select three points P_1, P_2 and P_3
- Determine $L_1, L_2,$ and L_3
- Know D_2 and D_3
- Solve for d , then P



(a) Pinhole model of imaging system



(b) Side view of imaging system

$$\frac{D_i}{f} = \frac{D_o}{L}$$

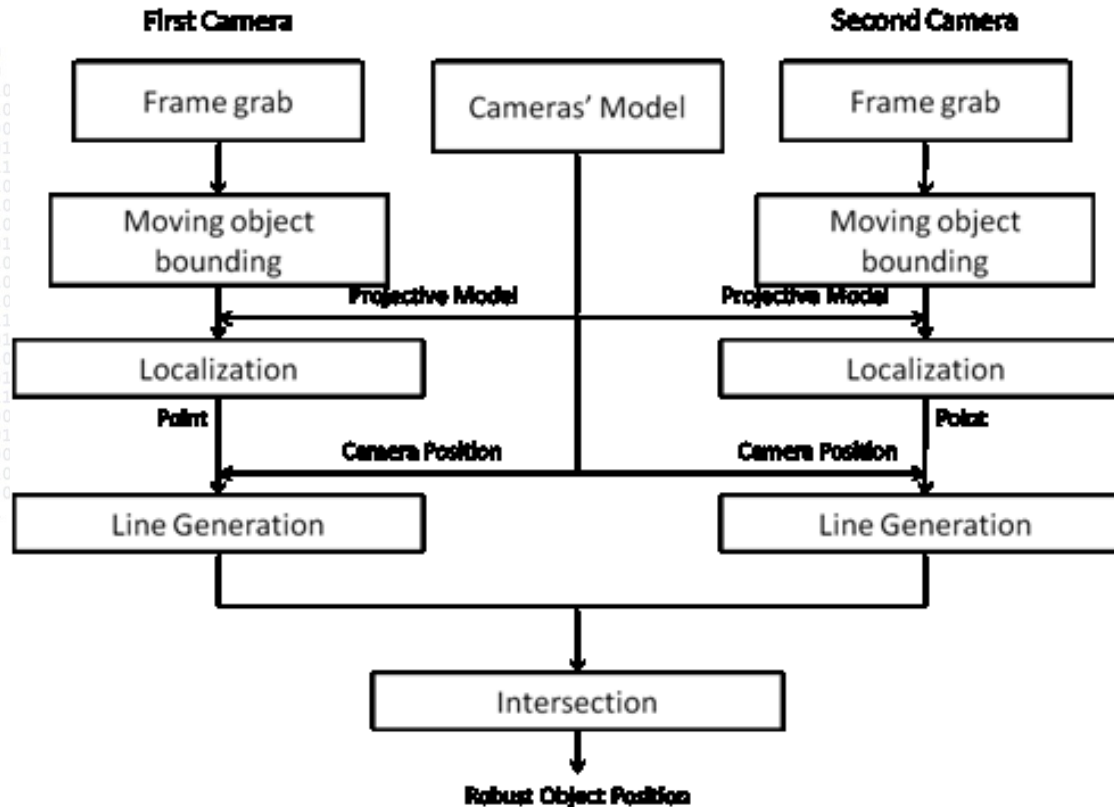
$$h^2 + d^2 = L_1^2$$

$$h^2 + (d + D_2)^2 = L_2^2$$

$$h^2 + (d + D_2 + D_3)^2 = L_3^2$$

$$\vec{x} - \vec{x}_1 + \alpha \left[\frac{\vec{x}_1 - \vec{x}_2}{|\vec{x}_1 - \vec{x}_2|} \right]$$

Two camera



```

10
001010
1010
10001010
10101000
000000010101010111 10101
010101010101010000010111011
0000010101010101010 001010
10101000001011101100100010
101010000101010101010 010
1101 00010101010101
01 0001010101010
010 11101101010
0 0101101010
101010111
10000101
01110110
01010101
0010101
0101000
000101
01000
1010
0 1
1
  
```

```

10
)10101000001011
)010000101010101011010110101
)00010111011001000110000010100
)01010100001101000010101010101
)100010101101010101000010101
)10101010101010101 000010
)101010101 100 10
)10110101 1010
)10101101 001
)101010101 0
)00001010
)0101000 01
)0101 01
)0101 01
)1 10 00
) 0 0111
) 0
)1
)1
)1
)10
)10
)10
)1 01
)0 101 10
)10101
)10 10
)11 0
)10010 10
)1010101010
)1011101110100
)11011001000101
)0101010110101101010
)00100011000001010
)110100001010101010
)10 00010101
)10111
)10
)01
)10
)010
)10
)0
)10
  
```

```

01000001011
0101010111
0001010001
010100000
10101010
100100
01010
1010
10101
0101
0000
01010
  
```

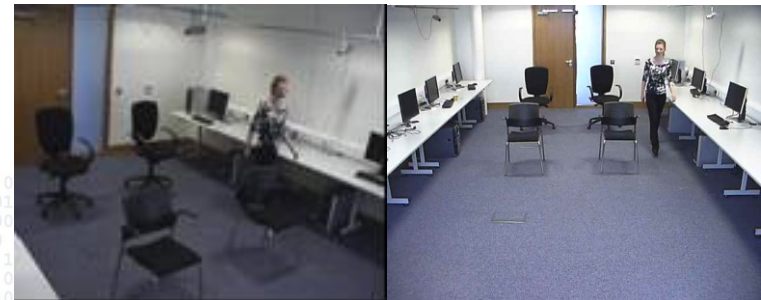
```

101010101
10101010
01000
001
  
```

```

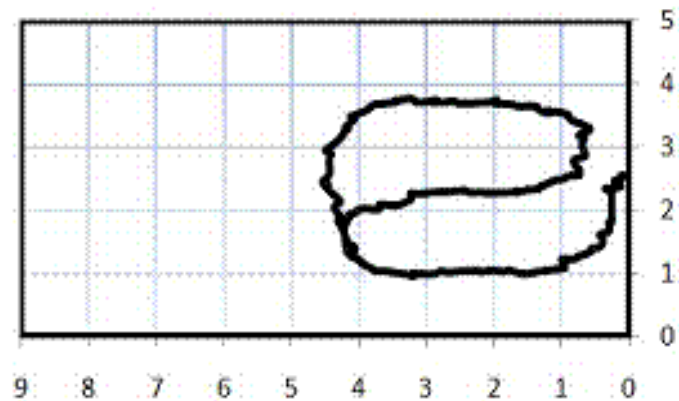
0101010110101101010
11011001000101
)0101010110101101010
)00100011000001010
)110100001010101010
)10 00010101
)10111
)10
)01
)10
)010
)10
)10
  
```

- Selected cameras 5 and 6 because they had largest coverage of any camera-pair
- 11 individuals were asked to roughly follow a route in the lab
- Chairs placed in the route
- 12400 video frames were collected
- Sequences manually analysed to give ground truth of real world position in each frame



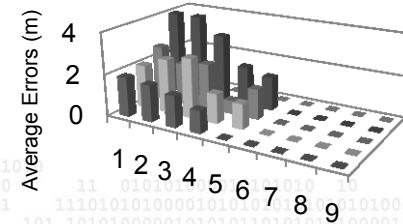
(a)

(b)

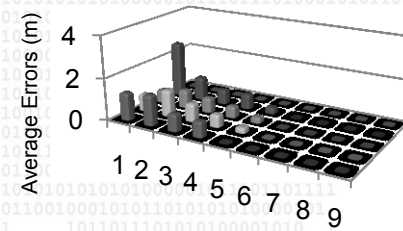


(c)

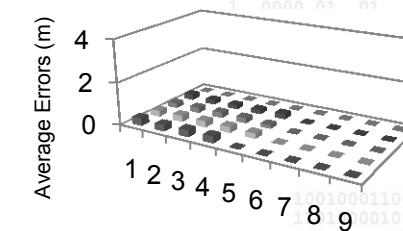
- Lab floor divided into 1m×1m blocks.
- Localization errors for each block are averaged over subjects



(a)

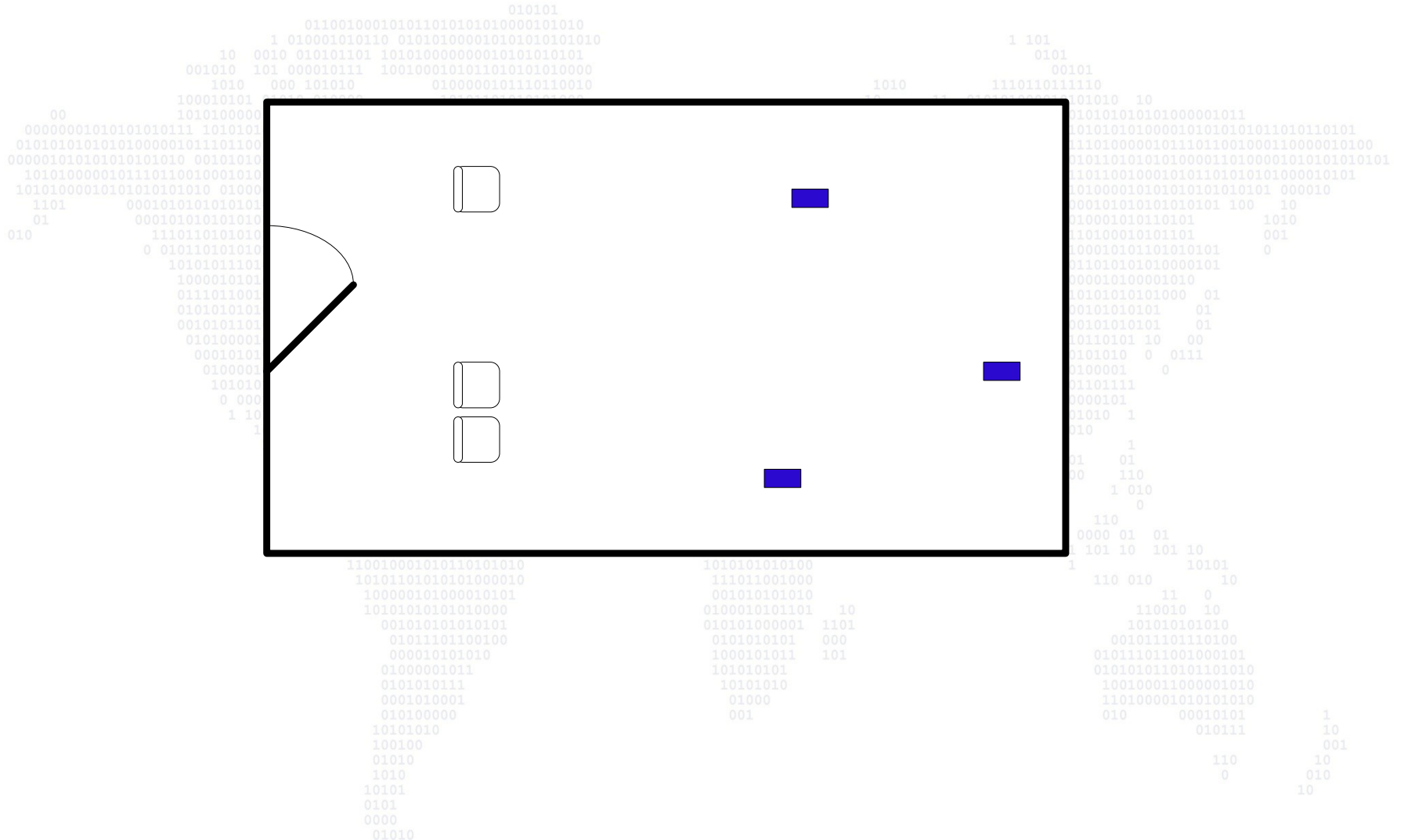


(b)



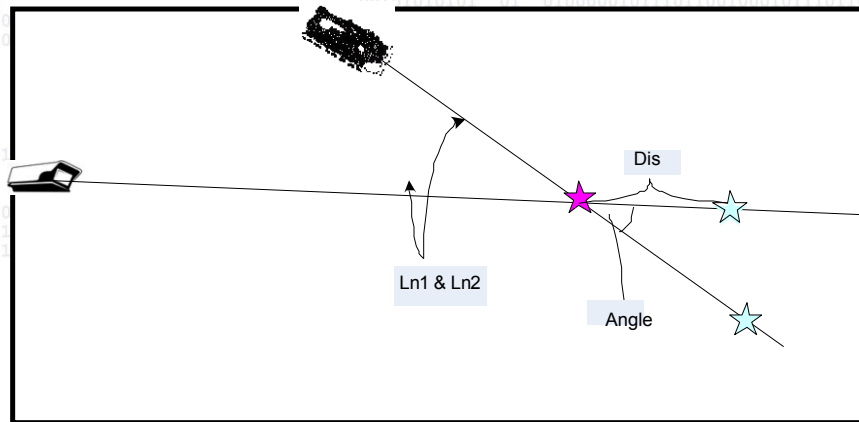
(c)

Three Camera

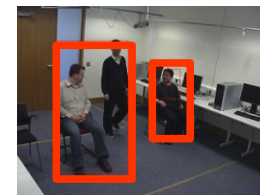
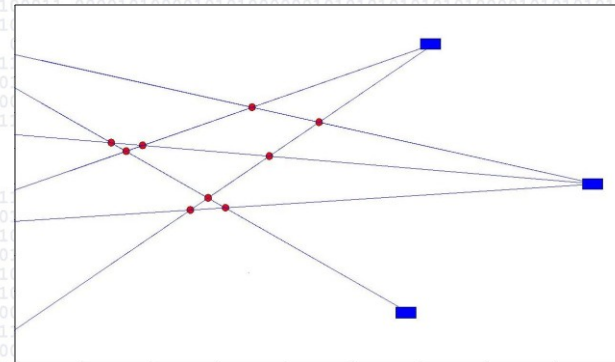
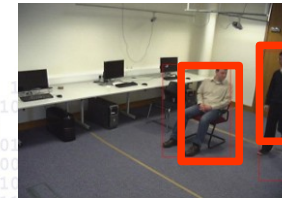


Intersection

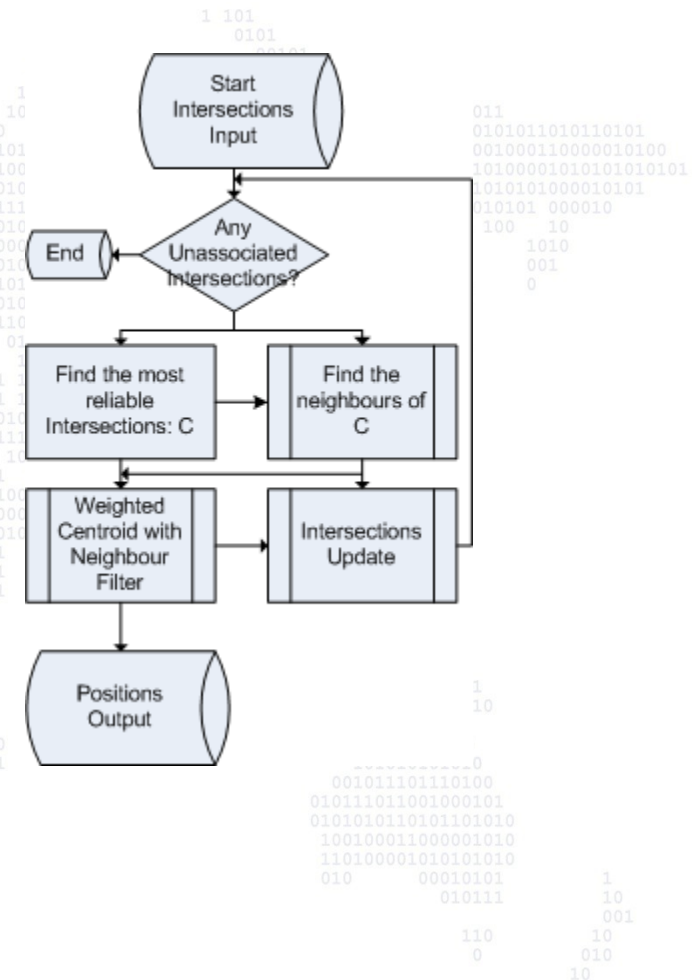
- Position: Ground floor co-ordinates, denoted by (x,y) .
- Dis: The distance to the nearest localisation point inferred by a single camera observation.
- Ln1, Ln2: The two intersecting lines generated by the two single-camera observations.
- Angle: The angle between Ln1 and Ln2.



- Multiple cameras mean multiple detections
- Some erroneous detections due to false alarms and occlusion
- Need to associate detections within subjects



- Reliability based upon
 - Cosine of angle
 - Min Dis
 - Combination of both: If $\text{min Dis} < T$ then cos
- Neighbours are detections within radius r of C
- Associated intersections removed from list



- 10 min sequence of 4 subjects entering, sitting, standing and exiting
- Corresponding to 2210 frames captured by each camera
- No. of subjects in each frame, summed over whole sequence, was 3870



Evaluation

$$T = \sum_i T_i w_i$$

$$T_i = \begin{cases} 1 & \text{if } A_i > S_i \\ A_i / S_i & \text{otherwise} \end{cases}$$

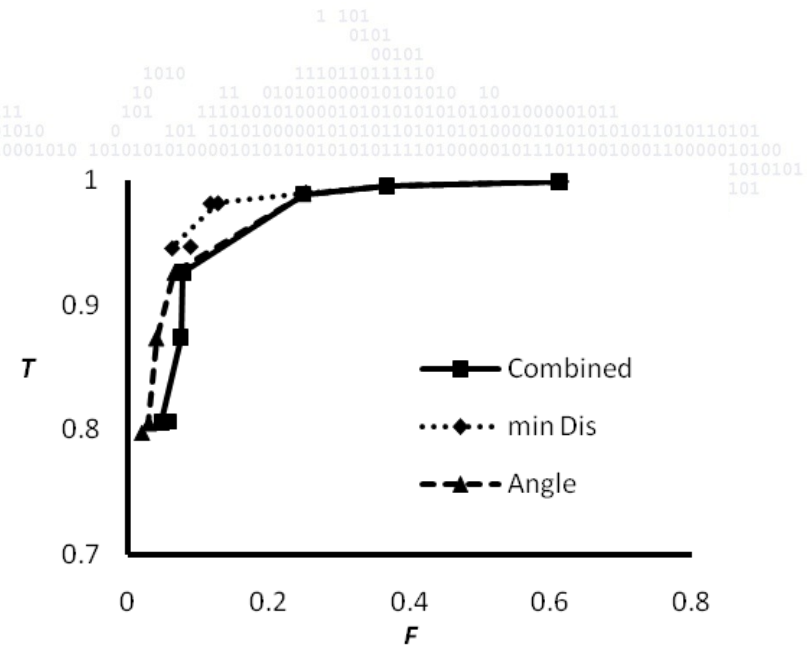
$$F = \sum_i F_i w_i$$

$$F_i = \begin{cases} \frac{A_i - S_i}{I_i - S_i} & \text{if } A_i > S_i \\ 0 & \text{otherwise} \end{cases}$$

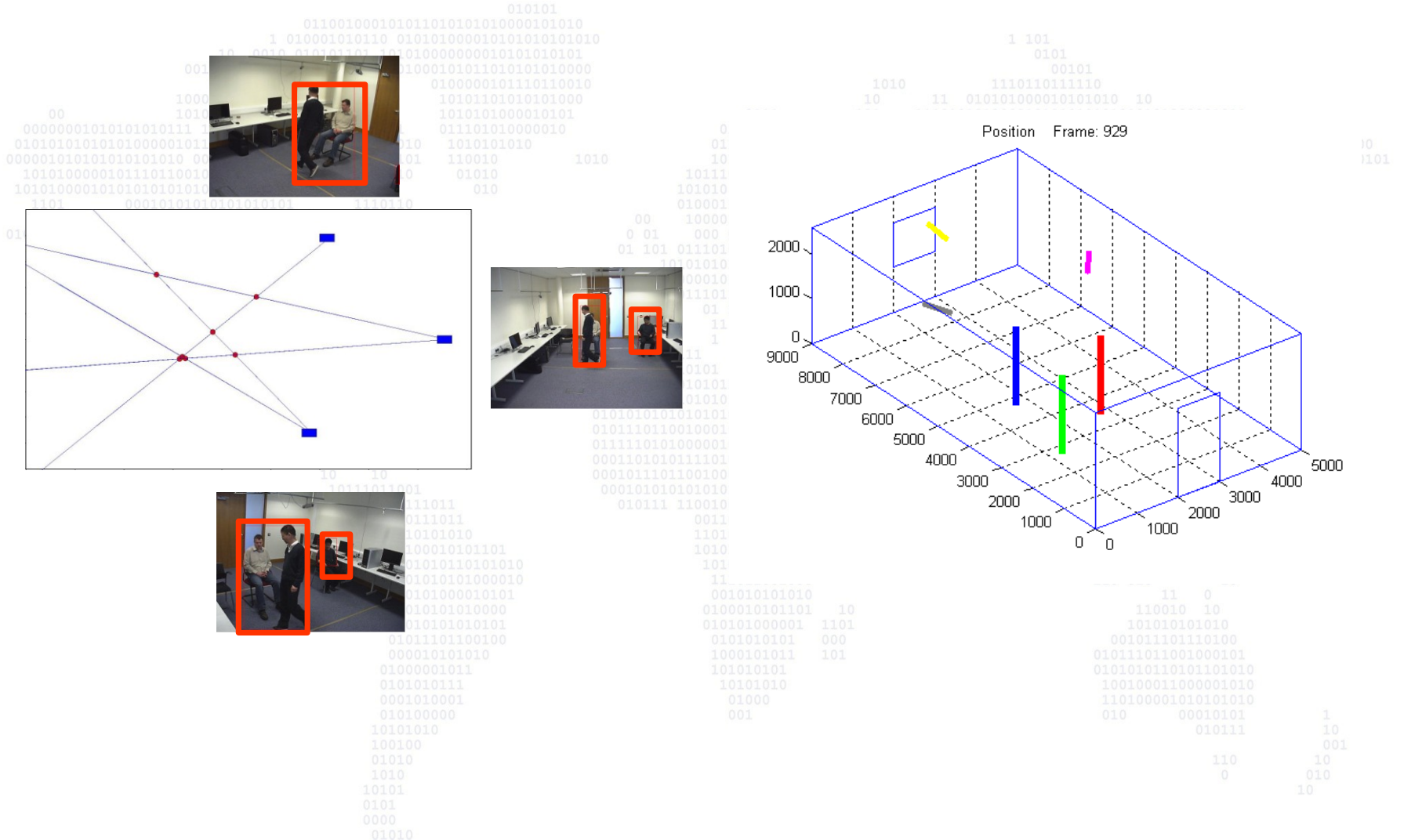
$$w_i = \frac{S_i}{\sum_i S_i}$$

- T_i is true positive rate
- S_i is the actual no. of subjects
- A_i is no. of subjects output from detection association algorithm
- F_i is false positive rate
- I_i is input no. of intersections

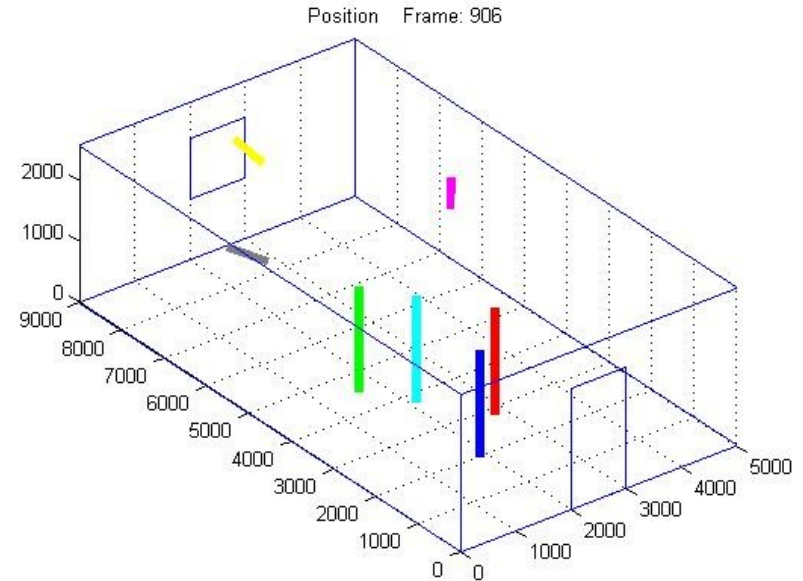
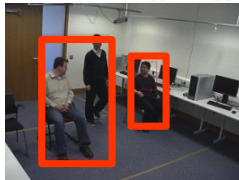
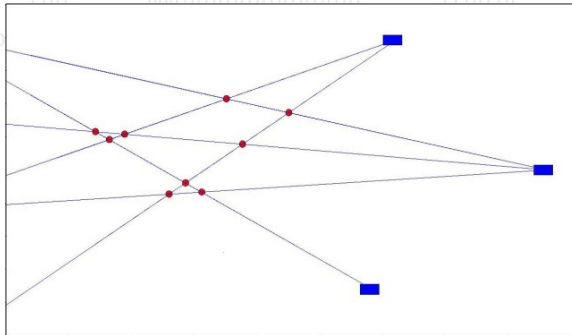
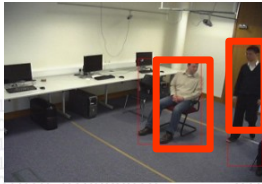
- T is a measure of underestimation
- F is a measure of overestimation
- Vary radius r to obtain ROC curves
- r varies from 10cm to 2m left to right



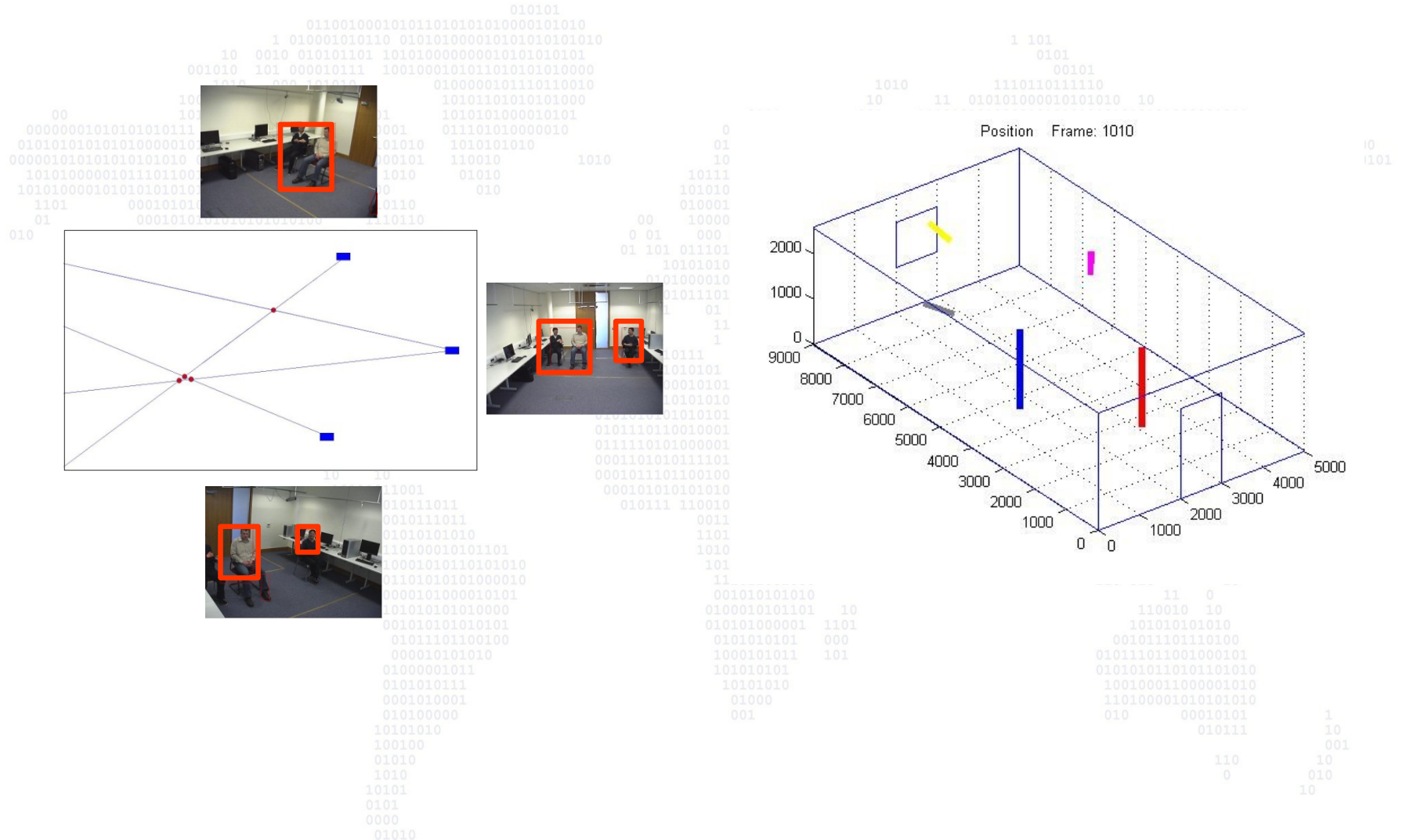
Correct Estimation



Over Estimation



Under Estimation



- Bayesian framework
 - Prediction
 - Filter
- Cannot be evaluated for most state-space models

$$p(x_t | z_{1:t-1}) = \int p(x_t | x_{t-1}) p(x_{t-1} | z_{1:t-1}) dx_{t-1}$$

$$p(x_t | z_{1:t}) = \frac{p(z_t | x_t) p(x_t | z_{1:t-1})}{\int p(z_t | x_t) p(x_t | z_{1:t-1}) dx_t}$$

$$\hat{x}_t = \arg \max_{x_t} p(x_t | z_{1:t})$$

Particle filter

- Represent posterior by set of random samples and weights
- Weights are updated according to observation likelihood
- Sequential importance sampling

$$p(x_t) \approx \frac{1}{N_p} \sum_{i=1}^{N_p} \delta(x_t - x_t^i)$$

$$p(x_t | z_{1:t}) \approx \sum_{i=1}^{N_p} w_t^i \delta(x_t - x_t^i)$$

$$w_t^i \propto w_{t-1}^i \frac{p(z_t | x_t^i) p(x_t^i | z_t)}{q(x_t^i | x_{t-1}^i, z_t)}$$

Kalman particle filter

- Kalman filter used to propagate each particle
- Steers particles towards regions with high likelihoods
- Fewer particles, less computation

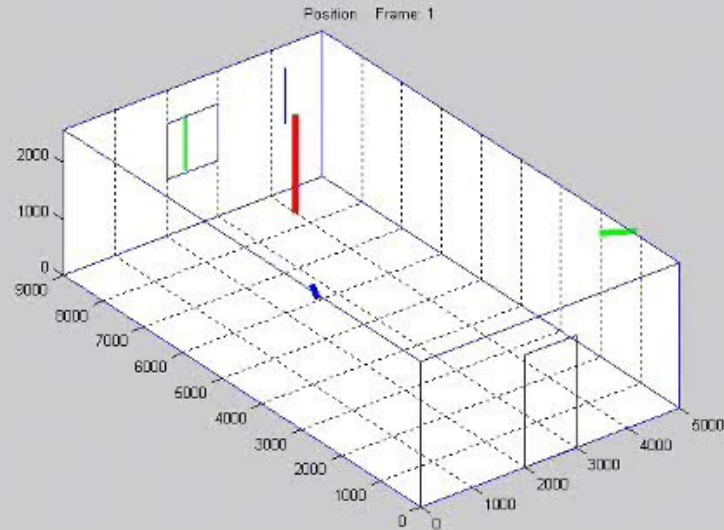
$$\bar{x}_t = Ax_{t-1} + K(y_t - H\bar{x}_t)$$

Tracking

```

10
001010
101
1000101
1010100
00
000000010101010111 1010
01010101010101000001011101
0000010101010101010 00101
1010100000101110110010001
101010000101010101010 01
1101 0001010101010
01 000101010101
010 1110110101
0 010110101
10101011
1000010
0111011
0101010
0010101
010100
00010
0100
101
0
1

```



```

011
0101011010110101
001000110000010100
10100001010101010101
1010101000010101
010101 000010
100 10
100 10
1010
001
0

```

```

1
10
0
00
101
1010
1010
1010
101
111
110
0
1
10
010
10
001
10

```

```

01010
1010
10101
0101
0000
01010

```



- Camera field of views are non-overlapping requires subject reacquisition for tracking
- Subject reacquisition is identification through applied detection, tracking and learning.
- Association of a current observed object with a previously observed object.
- Time gap could be seconds, hours, days etc.

- Online principal component analysis
- Ten components are learnt
- Temporal voting used for reacquisition

$$\mathbf{v}(0) = \mathbf{x}(1)$$

$$\mathbf{v}(t) = \frac{n-1}{n} \mathbf{v}(t-1) + \frac{1}{n} \mathbf{x}(t) \mathbf{x}^T(t) \frac{\mathbf{v}(t-1)}{\|\mathbf{v}(t-1)\|}$$

$$\mathbf{x}_2(t) = \mathbf{x}(t)_1 - \mathbf{x}_1^T(t) \frac{\mathbf{v}_1(t)}{\|\mathbf{v}_1(t)\|}$$

$$S(\mathbf{x}_t) = \sum_{i=1}^2 \dots$$

Subject Reacquisition

Time 1

Time 2

PCA Model

```

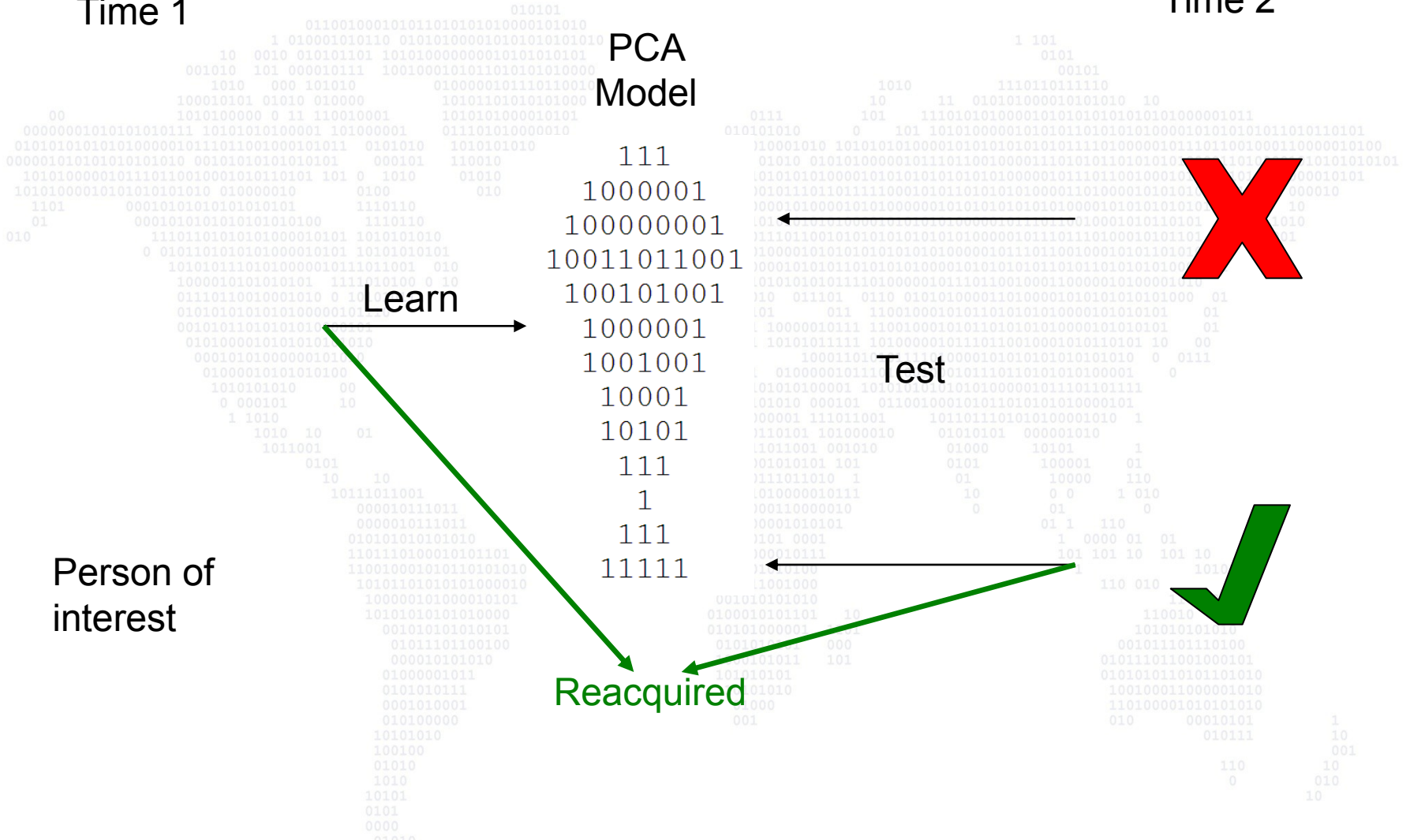
111
1000001
100000001
10011011001
100101001
1000001
1001001
10001
10101
111
1
111
11111
    
```

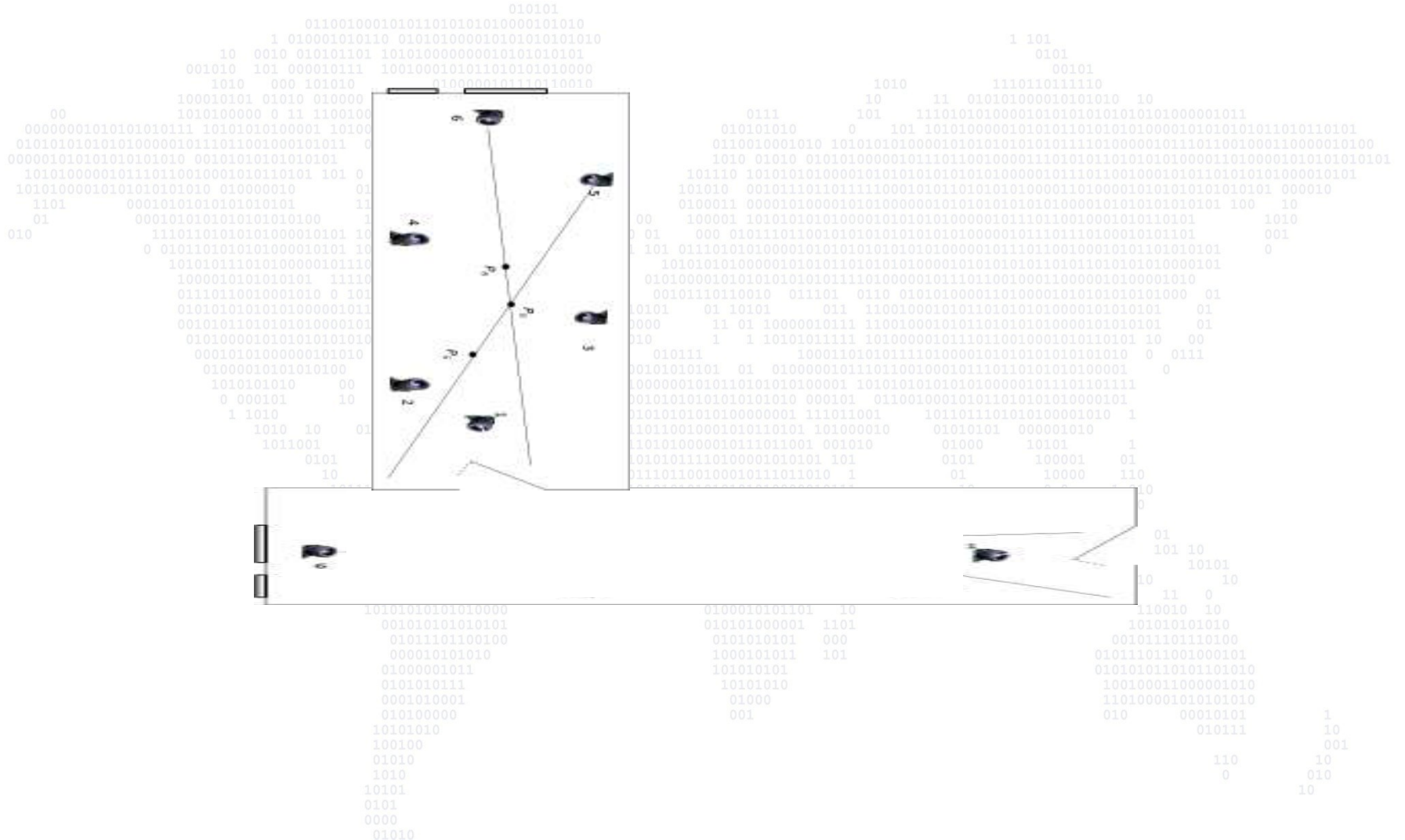
Learn

Test

Reacquired

Person of interest



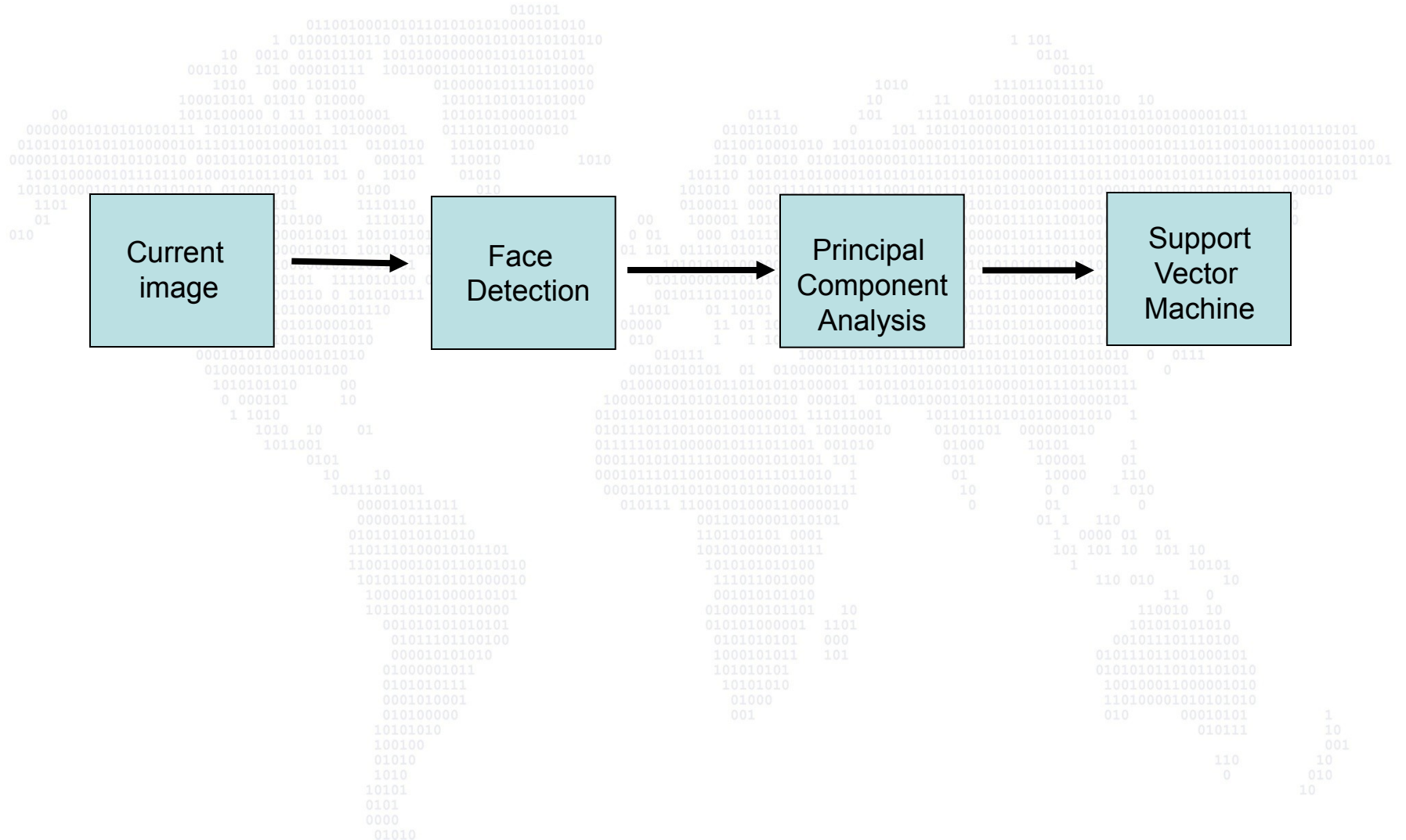


- 8 subjects at first camera
- Same 8 subjects at second camera
- Followed by 14 new subjects

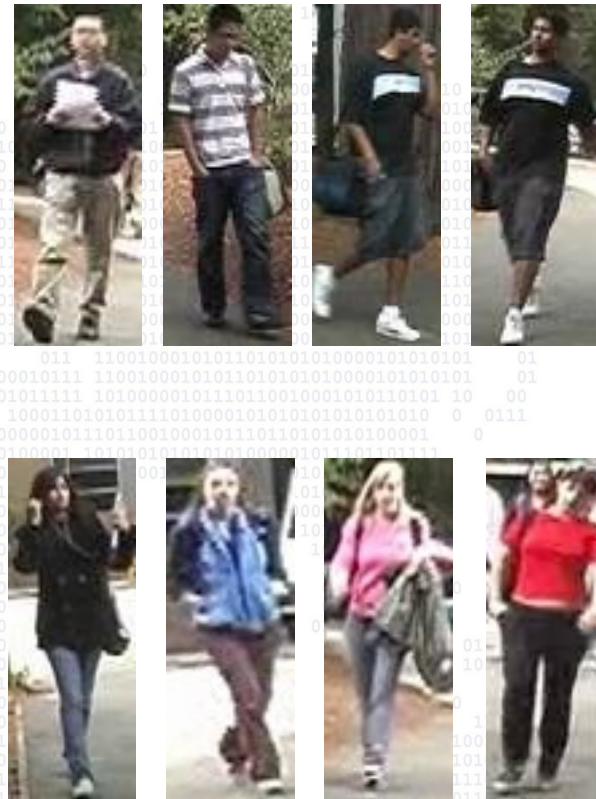
True Reacquisition (TP)	False 'Unknown' Decision (FN)	True 'Unknown' Decision (TN)	False Reacquisition (FP)
6 (75%)	2 (25%)	9 (64%)	5 (36%)

- Introduction & Motivation
- Tracking over a sensor network
- Gender Profiling
- Multi agent surveillance architecture
- Conclusion & Summary

Face-based Gender Profiling



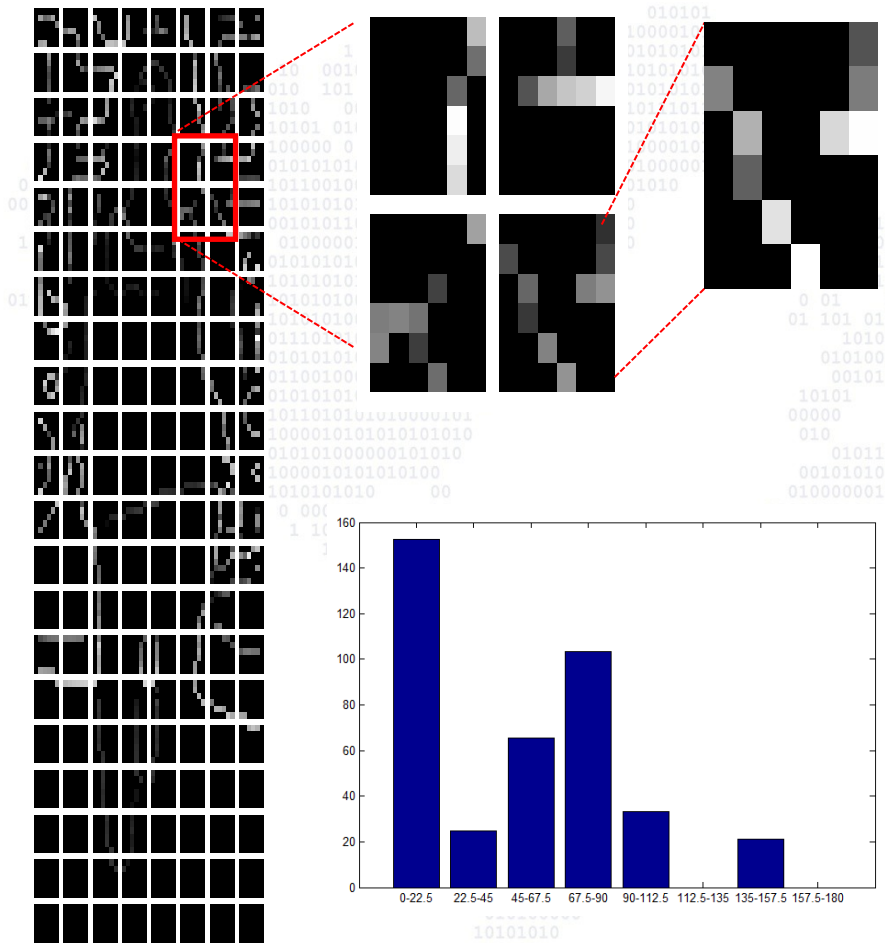
- Profiling at-a-distance
- Full body
- Based on clothes colour





- 128x48 pixel image
- Calculate Canny edge image
- Calculate gradient image then multiply with Canny

Canny Histogram of Gradients (CHoG)

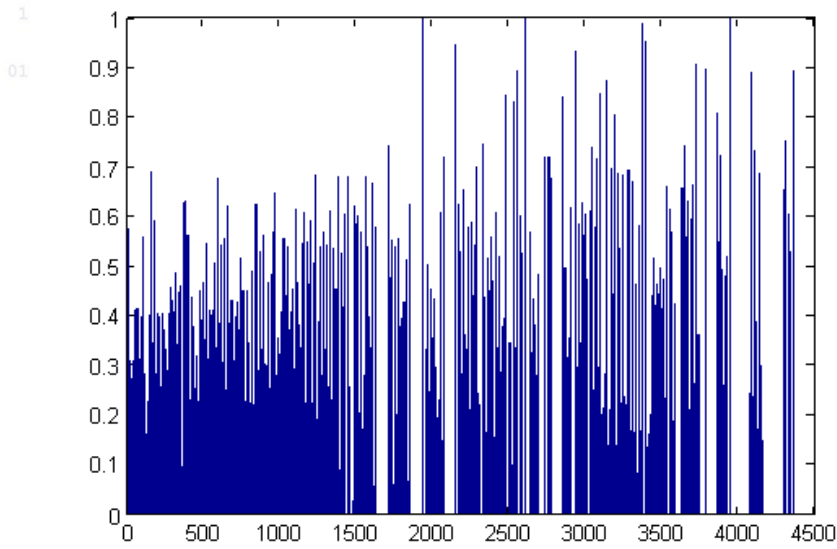


- 21x8 grid of cells of 6x6 pixels
- Histogram of cell edge orientations with 8 bins
- Bin value is sum of intensity gradients
- Cells grouped into overlapping blocks of 2x2
- Blocks overlap to give 140 (20x7) blocks.
- Concatenated HoGs from each block gives 4480 element feature vector

Average Female

```

010101
00000001010101010111 01010101000001 1010000001
001010 101 000010111 1001000101011010101010000
1010 000 101010 0100000101110110010
100010101 01010 010000 101011010101010000
00 1010100000 0 11 110010001 1010101000010101
00000001010101010111 10101010100001 1010000001 011101010000010
010101010101010000010111011001000101011 0101010 1010101010
0000010101010101010 00101010101010101 000101 110010 1010
  
```



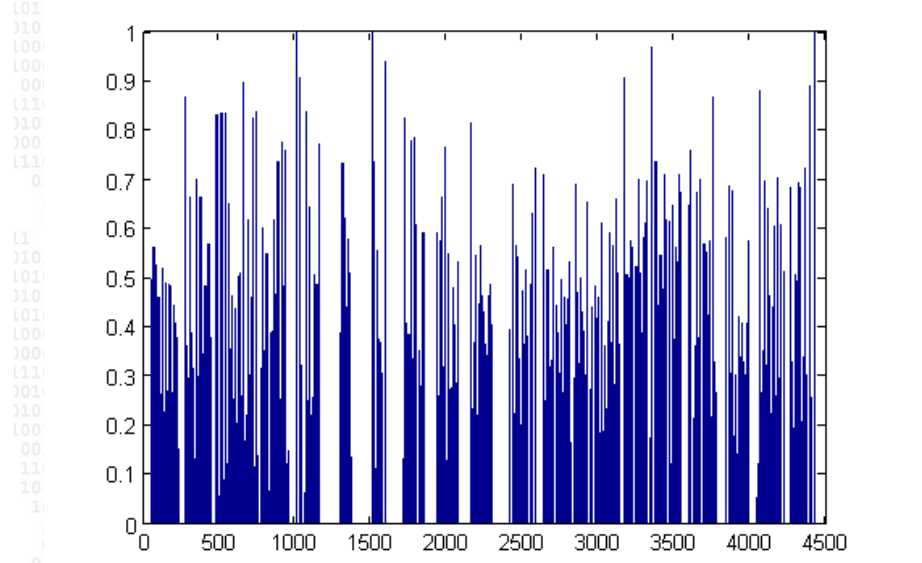
```

-----
01011101100100
000010101010
01000001011
0101010111
0001010001
010100000
10101010
100100
01010
1010
10101
0101
0000
01010
  
```

Average Male

```

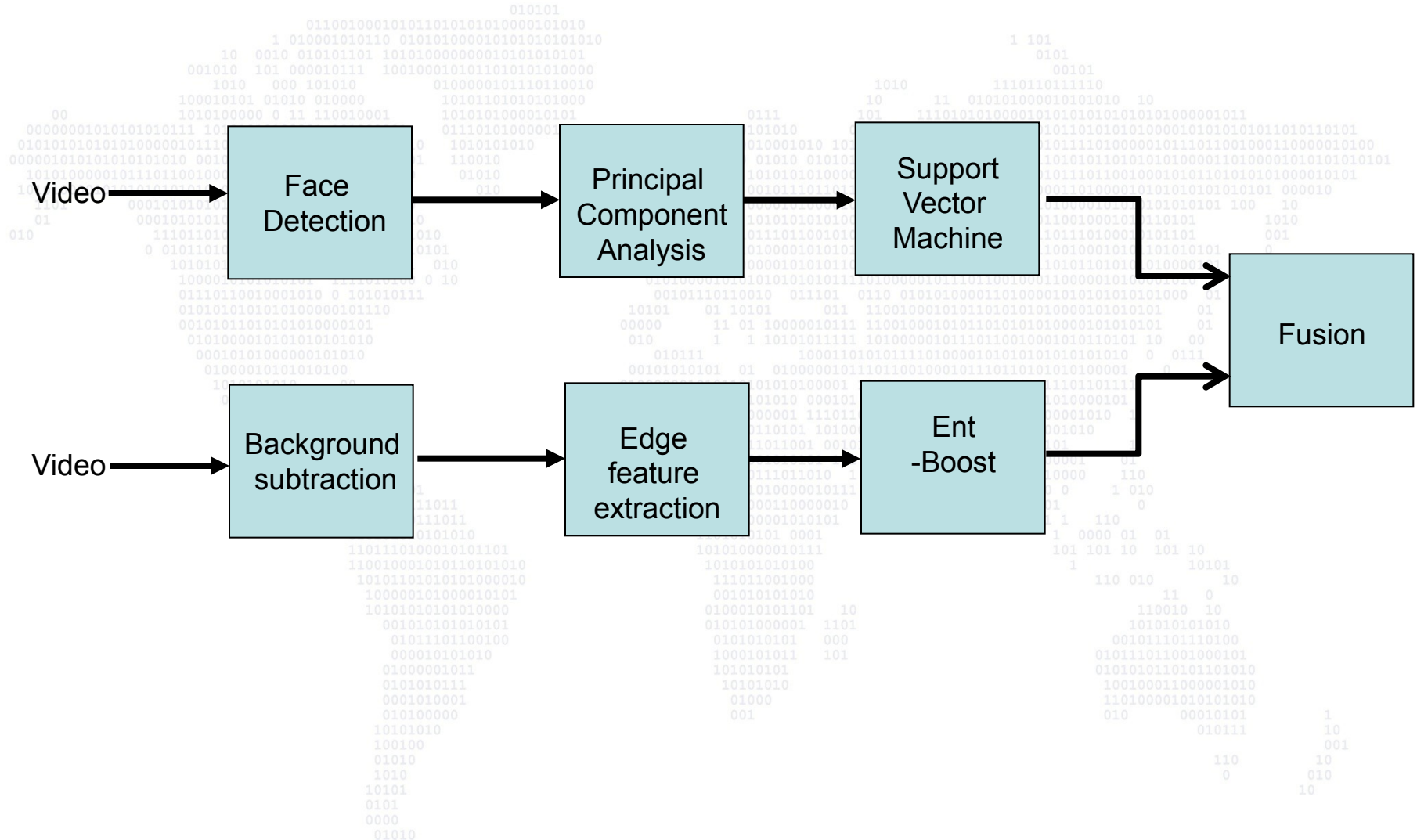
00101
1010 111011011110
10 11 010101000010101010 10
0111 101 1110101010000101010101010101000001011
010101010 0 101 1010100000101010110101010100001010101010110101
0110010001010 101010101000010101010101010111010000010111011001000110000010100
1010 01010 010101000001011101100100001110101011010101010100001101000010101010101
  
```



```

010101000001 1101
0101010101 000
1000101011 101
101010101
10101010
01000
001
1010101010
010111
010 00010101 1
010111 10
001 001
110 10
0 010
10
  
```

- Dataset of 413 images for each gender
 - MIT pedestrian recognition set
 - Viewpoint invariant pedestrian recognition set (VIPeR)
- Entropy Boost classifier
- Five-fold validation (80% training, 20% testing)
- 81% correct recognition



- Five sequences consisting of one or more subjects.
- Video resolution is 640×480 and frame rate is 10 fps (Panasonic WV-LZ62 camera).

Types	Total images	Pedestrians	Genders	Positive detections
Single-1	763	1	Female	672
Single-2	854	1	Male	807
Single-3	1013	1	Male	755
Single-4	1015	1	Female	870
Multiple-1	805	2	Female/male	730
Multiple-2	623	2	Male/male	571

- EB-Fusion approach is compared against:
 - Face based gender classification using PCA coefficients with SVM (FACE-PCA)
 - Full body based gender classification using HOG features with SVM (BODY-HOG)
 - Concatenated HOG features of face and full body components with SVM (CP-FB).

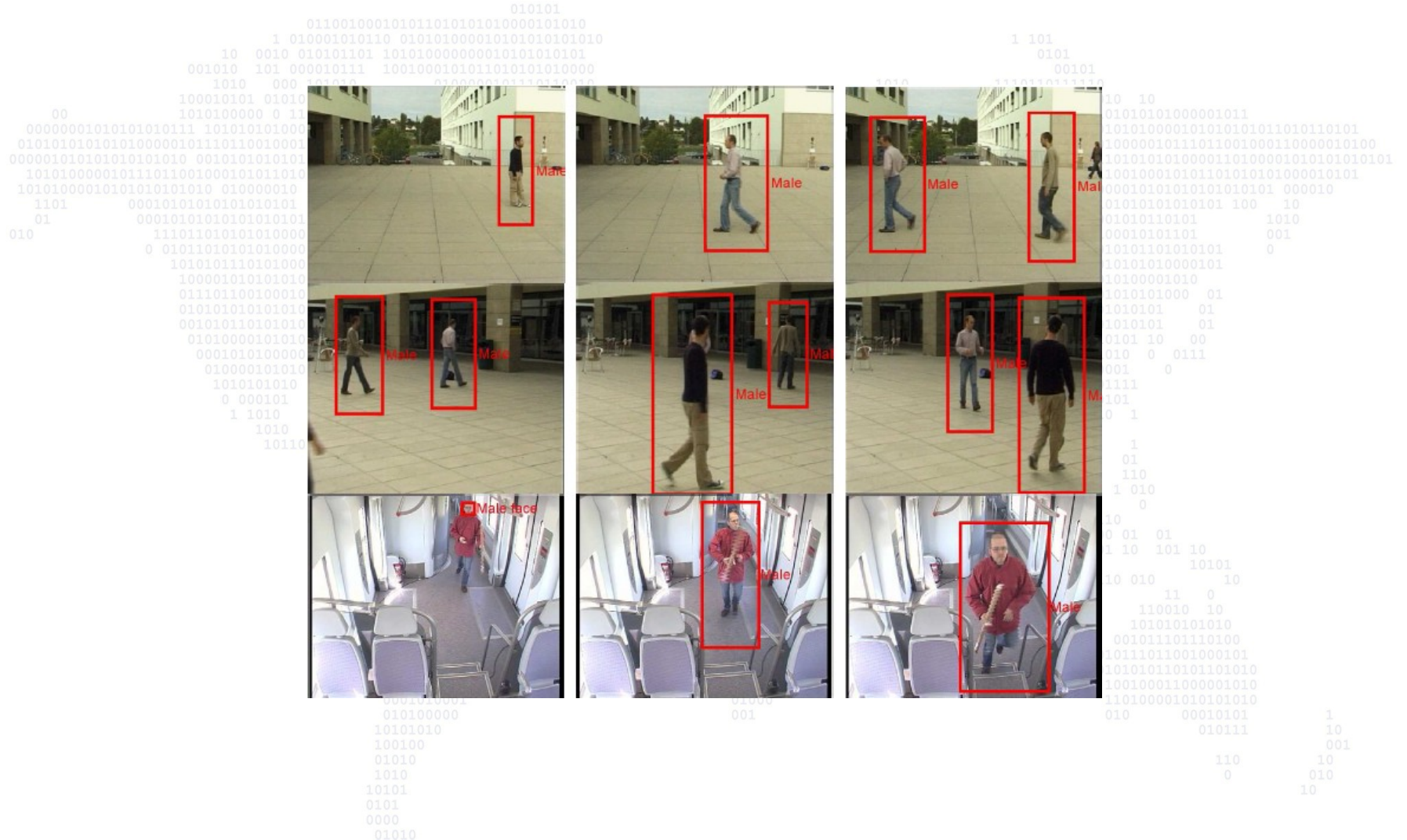
Types	EB-FUSION	CP-FB	FACE-PCA	BODY-HOG
Single-1	8.2%	16.4%	67.6%	17.8%
Single-2	9.1%	14.4%	43.2%	16.3%
Single-3	9.5%	13.2%	95.6%	14.0%
Single-4	10.1%	14.1%	94.4%	15.5%
Multiple-1	11.6%	15.4%	95.2%	16.3%
Multiple-2	9.0%	13.3%	95.9%	14.2%

Table 1 General information of the publicly accessible videos in Experiment 3.

Types	Total images	Pedestrians	Gender	Positive detections
Campus-1	2000	> 3	Male/female	1613
Campus-2	2000	> 3	Male/female	1554
Train-1	1626	> 3	Male/female	1019

Table 2 Error statistics of gender classification in Experiment 3.

Types	EB-FUSION	CP-FB	FACE-PCA	BODY-HOG
Campus-1	16.4%	18.2%	99.6%	19.8%
Campus-2	15.1%	17.4%	99.2%	19.1%
Train-1	14.7%	19.8%	76.3%	19.6%



Gender Profiling



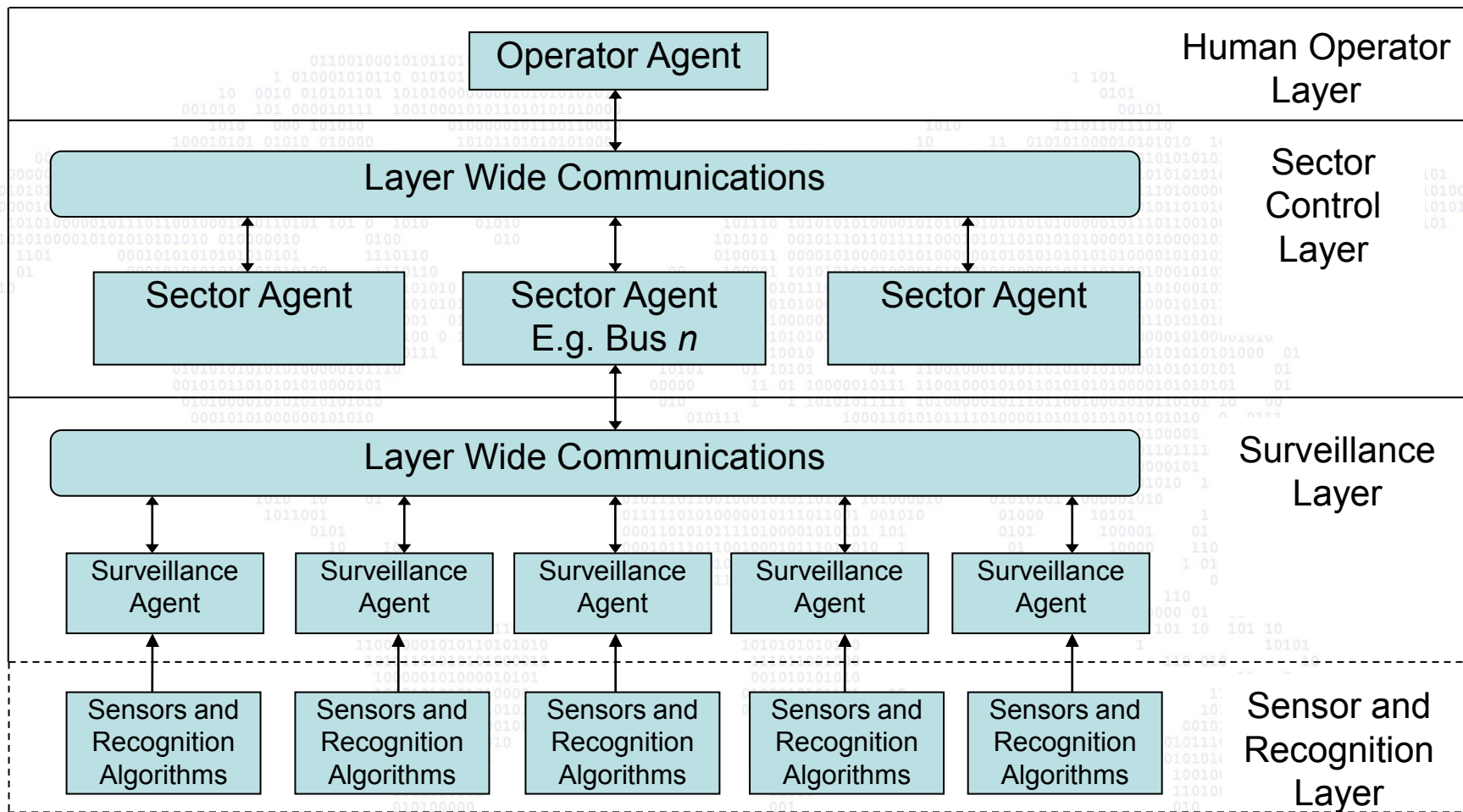
0101
0000
01010

11
101011010110101
01000110000010100
0100001010101010101
010101000010101
10101 000010
100 10
1010
001
0

0
0
01
010
010
010
01
11
10
001
10
010
10

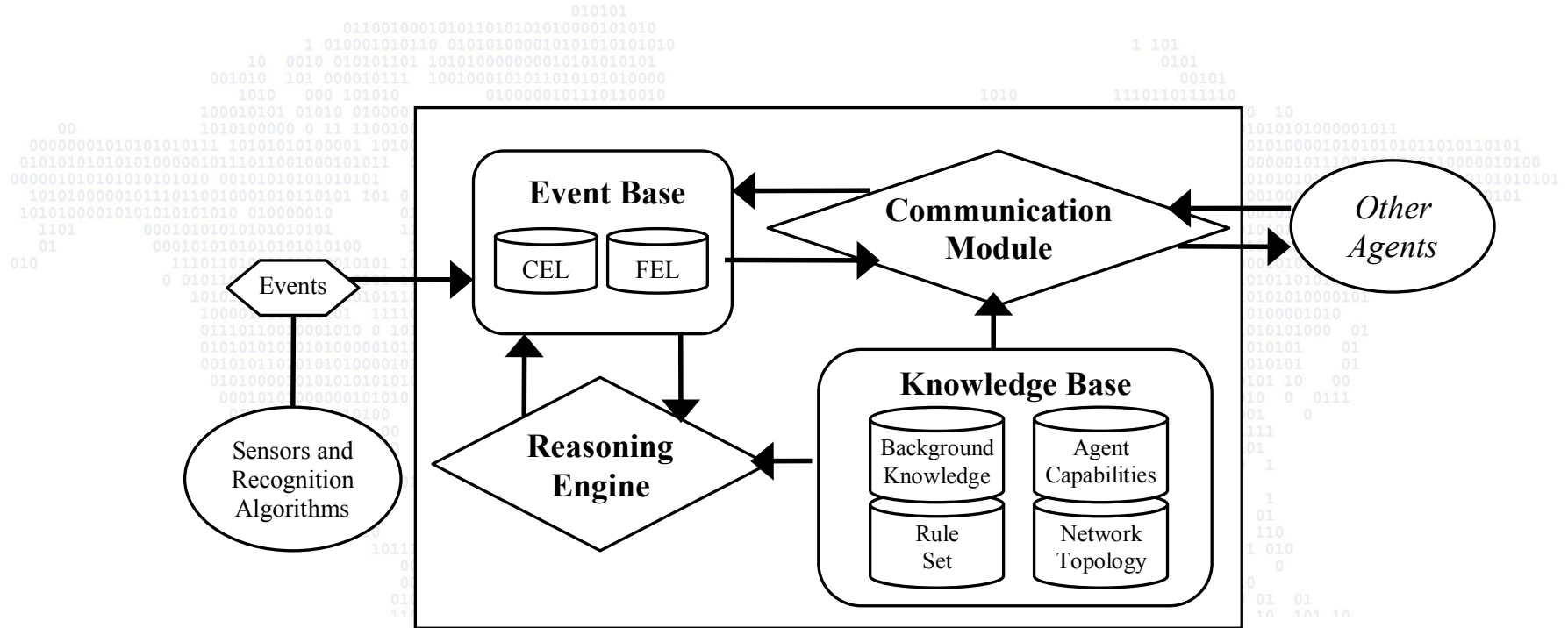
- Introduction & Motivation
- Tracking over a sensor network
- Gender Profiling
- Multi agent surveillance architecture
- Conclusion & Summary

Architecture for Multi-Agent Surveillance



- Third step in threat assessment is to combine the who's in the scene with the where in the scene
- Event management framework
- Implemented using multi-agent architecture

Surveillance Agent



- Passenger boarding
- $e_1 = (PBV, 21:05:31, 1, 0.85, 0.7, \{male\})$
- *PBV* is for event type *Person Boarding Vehicle*
- *21:05:31* is the time of occurrence *occur_T*
- *1* denotes the source, in this case a video analytics algorithm
- *0.85* is the source reliability
- *0.7* denotes the significance
- *male* is the value for the gender attribute

- Rule 1: Infers the event abusive behaviour towards driver
- $R_1 = (LS_1, EType_1, Condition_1, m_1)$
- $LS1 = (TPL, TPL + 120)$
- $EType1 = DA$ abbreviated for driver abuse
- *Condition is :*
 $ei.Etype = PL \wedge ej.Etype = PS \wedge ei.location = Drivers\ Cabin$

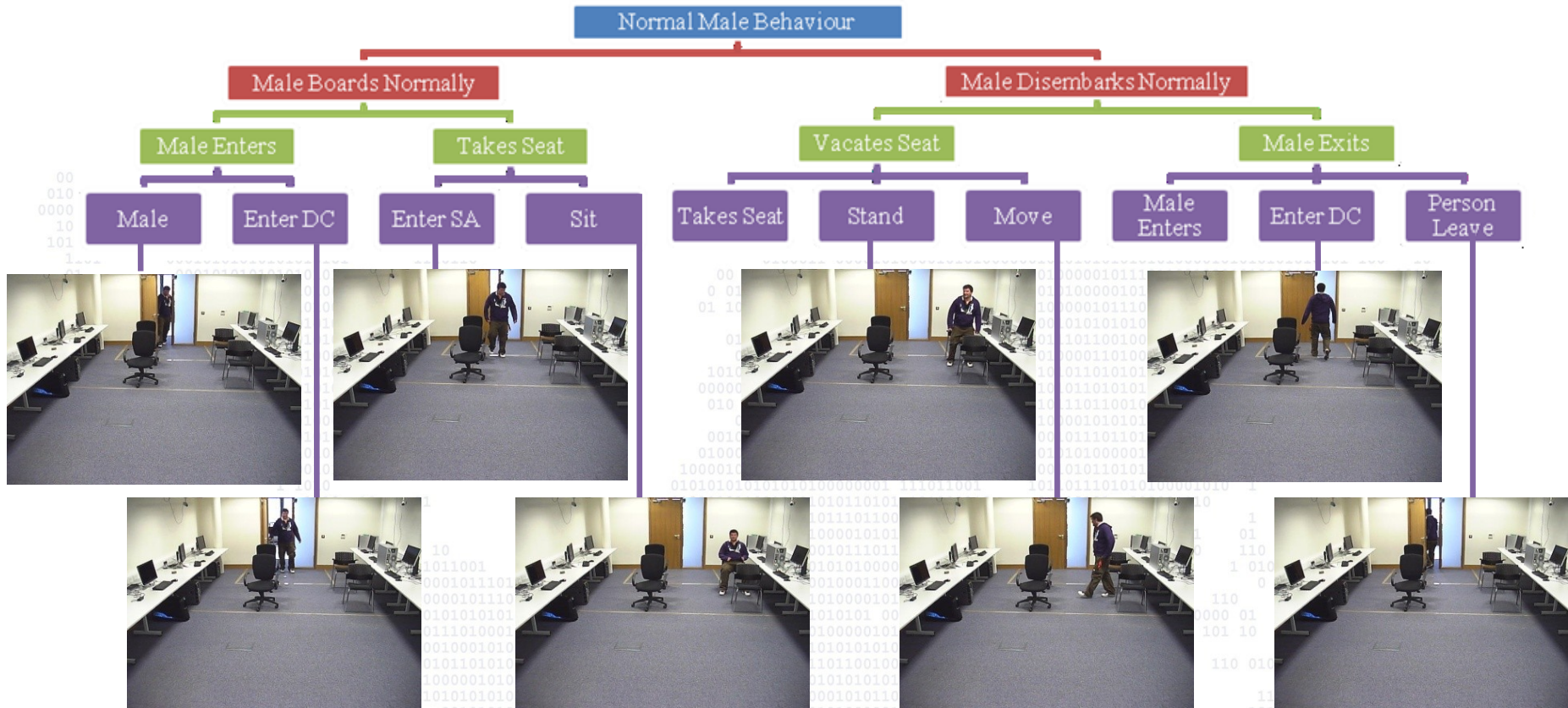
- Centralised reasoning has scalability limitations.
- Subscription-only presents a huge overhead for events that are produced frequently but are rarely needed.
- Communication protocol works on a “need to know” basis.
 - onetime event queries
 - event subscriptions

- Used where an agent only needs to know about that type of event for one individual event pattern.
- Sets other events in that pattern as a trigger for a onetime request of that event
- Reduces the communication whilst still ensuring the agent has complete knowledge for reasoning.

- 55 sequences in total
- 40 of normal passenger behaviour; person boards bus, sits and then exits the bus
- 15 of suspicious behaviour
 - 5 of people loitering in the saloon area
 - 5 of loitering in the driver's cabin
 - 5 of people obscuring face from cameras

- Lab was partitioned out into 2 sections to represent the driver's cabin and the saloon area.
- Area surrounding the door up to the first set of seats is regarded as driver's cabin
- Rest of the floor space and chairs are saloon area.
- Multi Agent System is developed using the agent middle ware JADE.
- Reasoning module for the agents and centralized reasoner developed using PROLOG.

Normal Male Behaviour

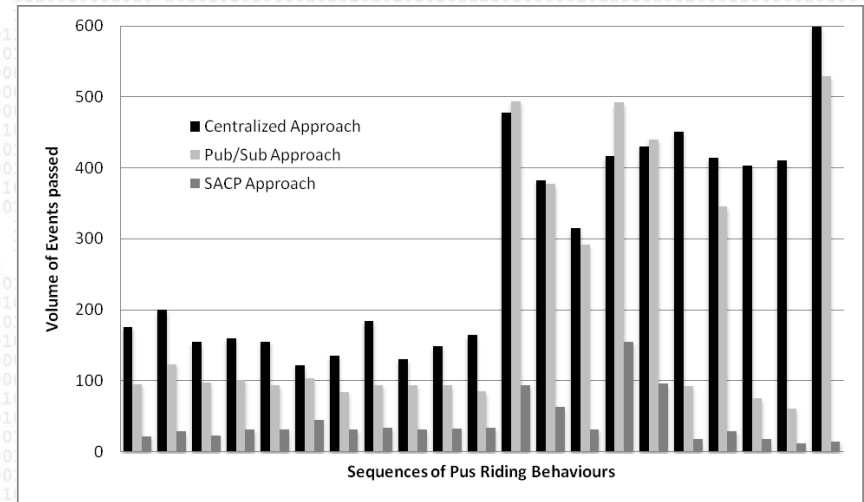


Example of a lab simulated normal male behaviour

- Ground Truth (GT)
- True Positives (TP)
- False Positives (FP)
- False Negatives (FN)
- Sensitivity ($S = TP / (TP + FN)$)
- Precision ($P = TP / (TP + FP)$)

Compound Event	GT	TP	FP	FN	S (%)	P (%)
Male Normal Behaviour	20	18	2	2	90%	90%
Female Normal Behaviour	20	16	2	4	80%	89%
Loitering Saloon Area	5	4	1	1	80%	80%
Loitering Drivers Cabin	5	4	1	1	83%	80%
Person Obscuring Face	5	4	0	1	80%	80%

- Measured events sent for reasoning
 - centralized
 - agent using publish/subscribe
 - agent using SACP



- Video analytics for tracking and gender profiling have been demonstrated
- Event recognition and composition for a sensor network demonstrated in laboratory conditions
- Preliminary evaluation of architecture for multi-agent surveillance
- Bus trial for further evaluation