



Pedestrian Tracking by Learning of Multi-view Human Color Appearance

Huynh Loc Huu

Content

- ▶ 1. Introduction
 - 1.1 Applicable Service
 - 1.2 Challenges
 - 1.3 Motivation
 - 1.3 Proposed method
- ▶ 2. Overview of the system-Frameworks
 - 2.1 Overview of the system
 - 2.2 Formulation a tracking system as MAP
 - 2.3 Reversible jump Markov chain Monte Carlo (RJMCMC)
 - 2.4 Posterior probability
- ▶ 3. Implementation of Tracking System
 - 3.1 The Background Likelihood
 - 3.3 The Color Likelihood
 - 3.4 The Prior
- ▶ 4. Evaluation
- ▶ 5. Conclusion and Future works

1. Introduction

- Locating the accurate position of all persons in the observed area.
- Different people are enclosed by different color rectangles.
- Each person is enclosed by the same color in all camera-views.



Figure 1.1: Result of Pedestrian Tracking by Using Multiple Cameras

1.1 Applicable Service

- Pedestrian surveillance for security purpose:
 - ◆ Detecting and identify tragic events, potential accidents.
 - ◆ Automatically reporting to the security center.
 - ◆ Reduce the human workload of security authority.
- The research of analysis customer behavior in supermarket by navigating the interesting location of the customers.
- Analyzing team sports' tactical information.



Figure 1.2: Pedestrian Surveillance in Supermarket

1.2 Challenges



Figure 1.3: Obstacles



Figure 1.4: People is occluded



Figure 1.5: Identities switching

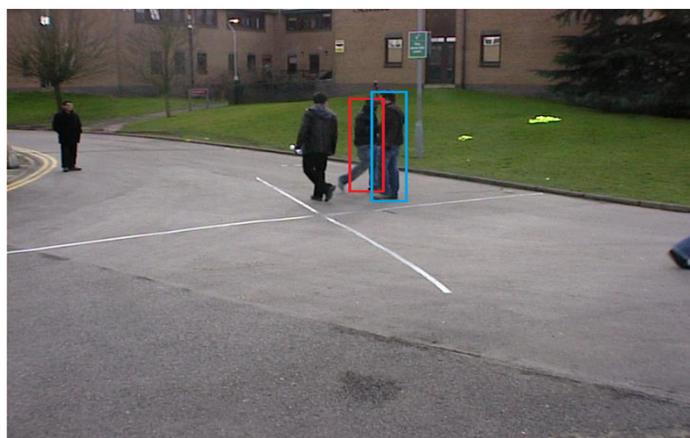


Figure 1.6: Interacting between people

1.3 Motivation



Figure 1.7: Noisy made it very hard to detect people



Figure 1.8: Accurate detection by using the information from color

- We are motivated from the problems of previous methods which used only the information from foreground mask. By the using of color appearance, we could detect people more accurately and help to overcome the identities switching problem.
- The tendency of moving such as spatial location and velocity also can be a very important cue to smooth the tracking trajectories.

1.4 Proposed method

- ▶ Handling occlusion by using multiple cameras.
- ▶ We applied color appearance learning into the tracking system.
 - The advantages is fusing color appearance information from all views without directly reconstructing 3D shape.
 - The appearance of people is learned and updated automatically.
- ▶ Human model is the mixing of ellipse and rectangle.
- ▶ The tracking model is comprehensively formulated as the modified Bayesian tracking model.
- ▶ The tendency of moving such as spatial location and velocity are used to smooth the tracking trajectories.



2.1 Overview of the System

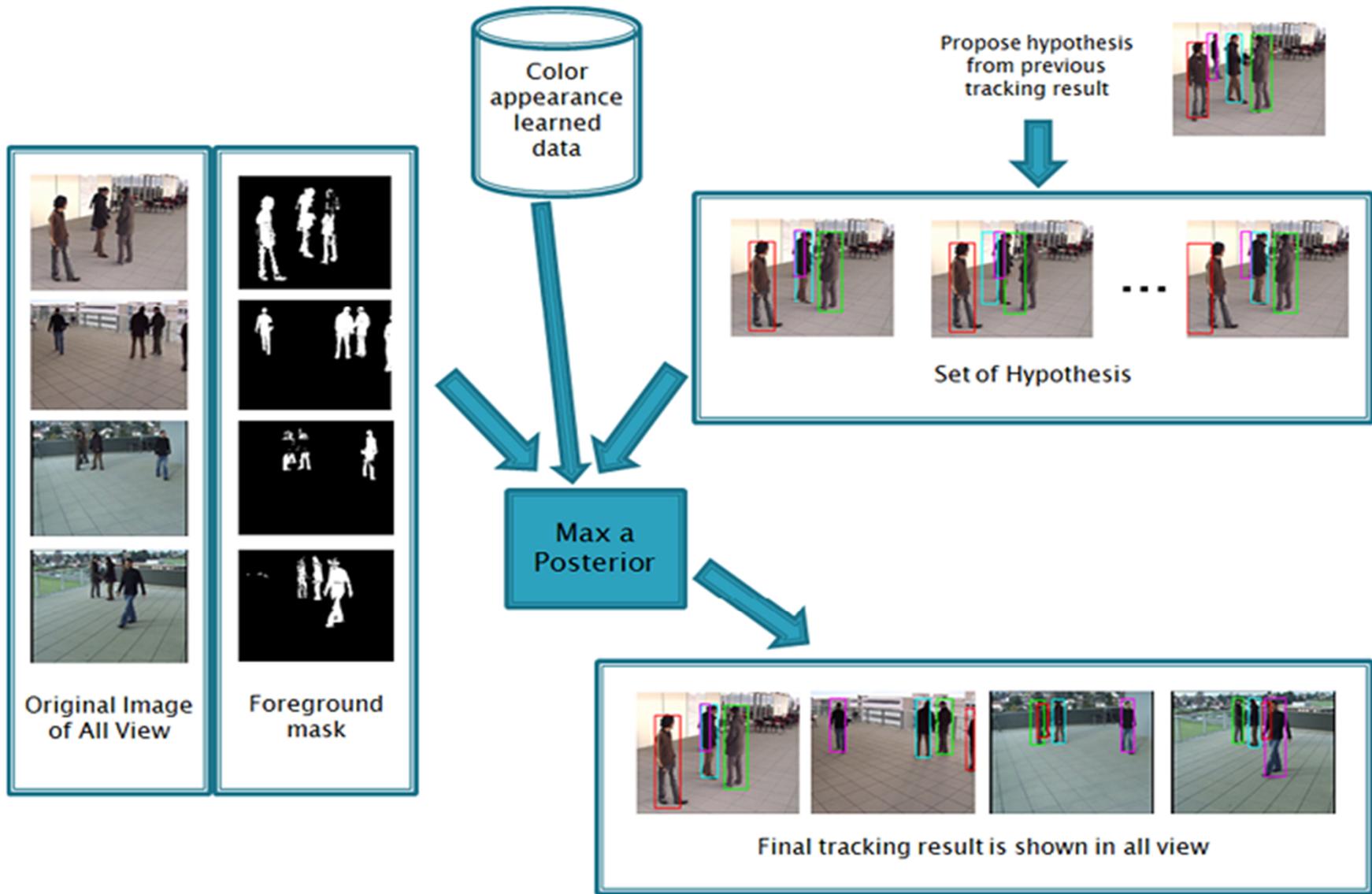


Figure 2.1: Overview of the system

2.2 Formulation a tracking system as MAP

- ◆ At time t , the state of multi targets is defined as:

$$X^{*(t)} = \{ x_1^{(t)}, x_2^{(t)}, \dots x_n^{(t)} \}, n > 0$$

- ◆ At target i , the state includes the position c and the height of person h :

$$x_i^{(t)} = \{ c_i, h_i \}$$

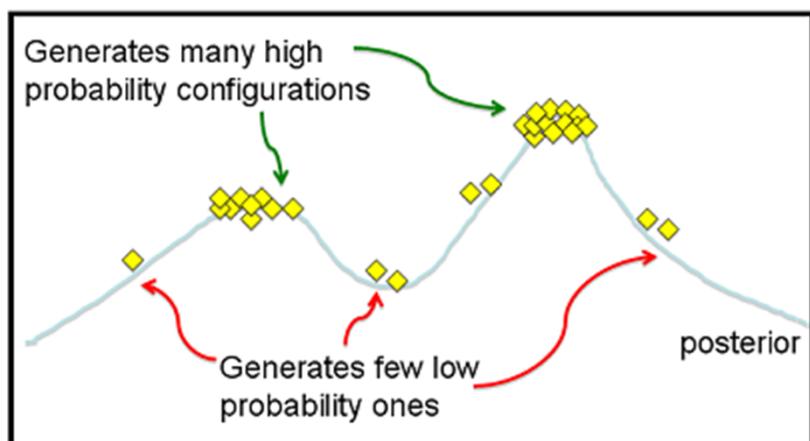
- ◆ Formulating the tracking problem as computing the maximum a posterior (MAP) estimation

$$X^{*(t)} = \operatorname{argmax}_{X^{(t)} \in \omega} P(X^{(t)} | I^{(t)}, X^{*(t-1)}) \quad (2.1)$$



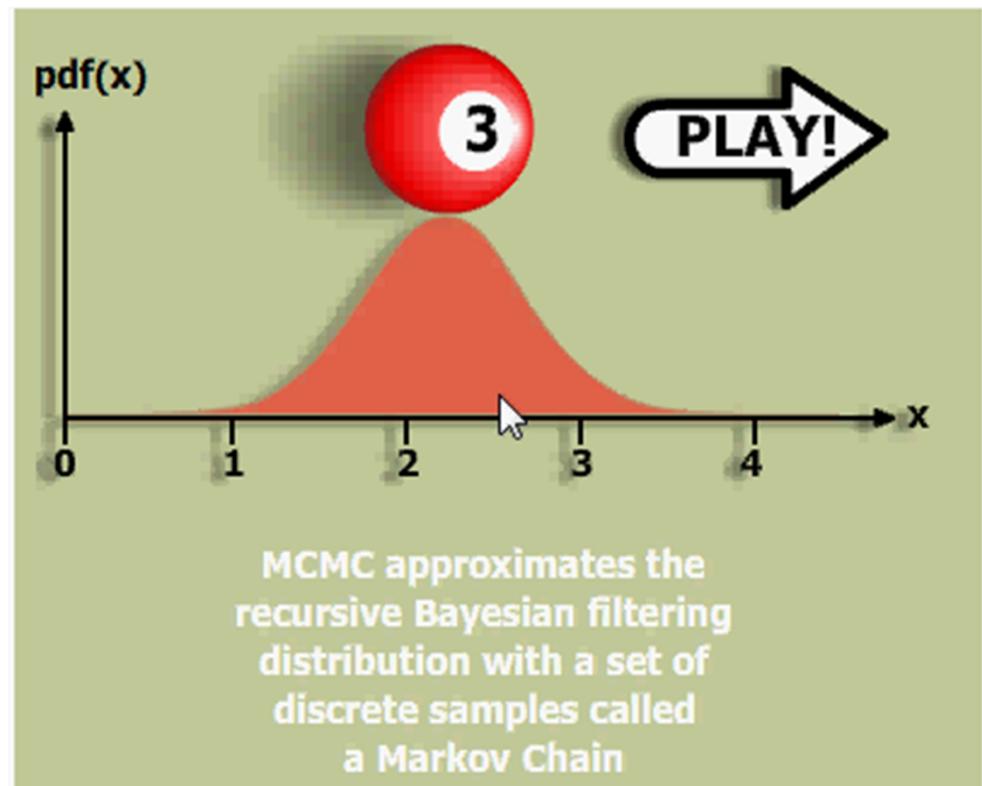
2.3 Reversible jump Markov chain Monte Carlo (RJMCMC)

- ▶ The growing of number of people also increases extremely number of hypothesis. Therefore, direct computation of the posterior for all hypothesis is intractable.
- To solve this problem, we used RJMCMC approximate the posterior distribution.



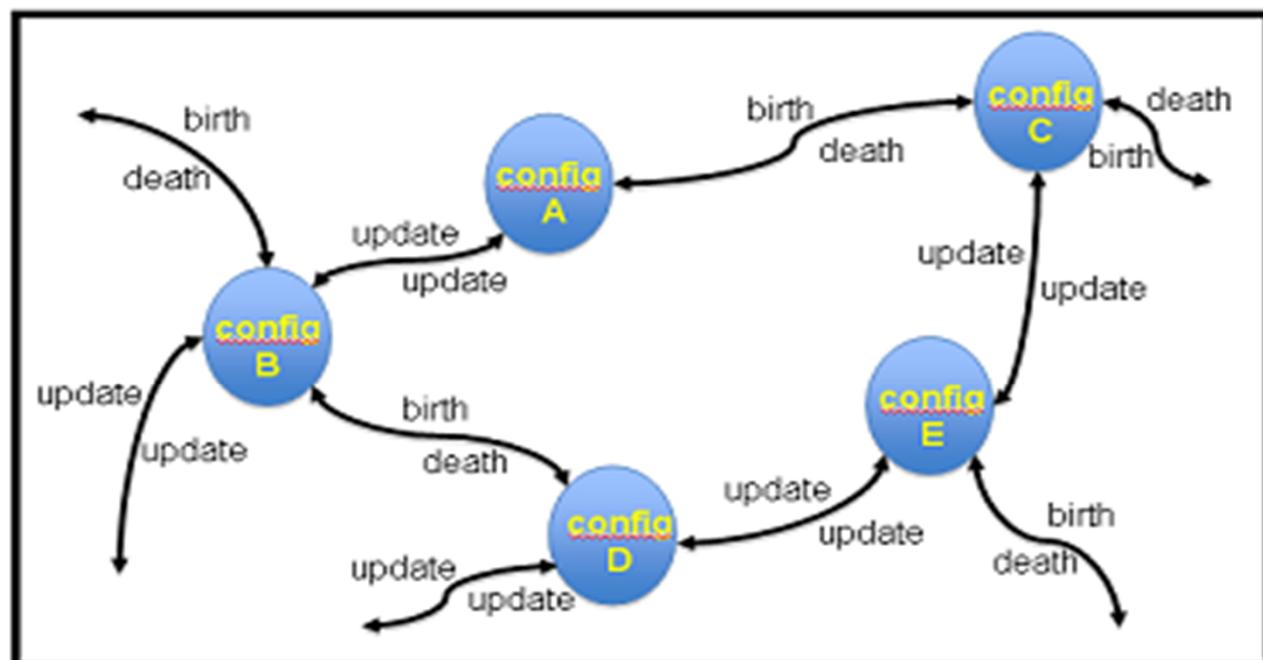
Markov Chain Monte Carlo

Sampling the unknown posterior function
with the finite number of sample



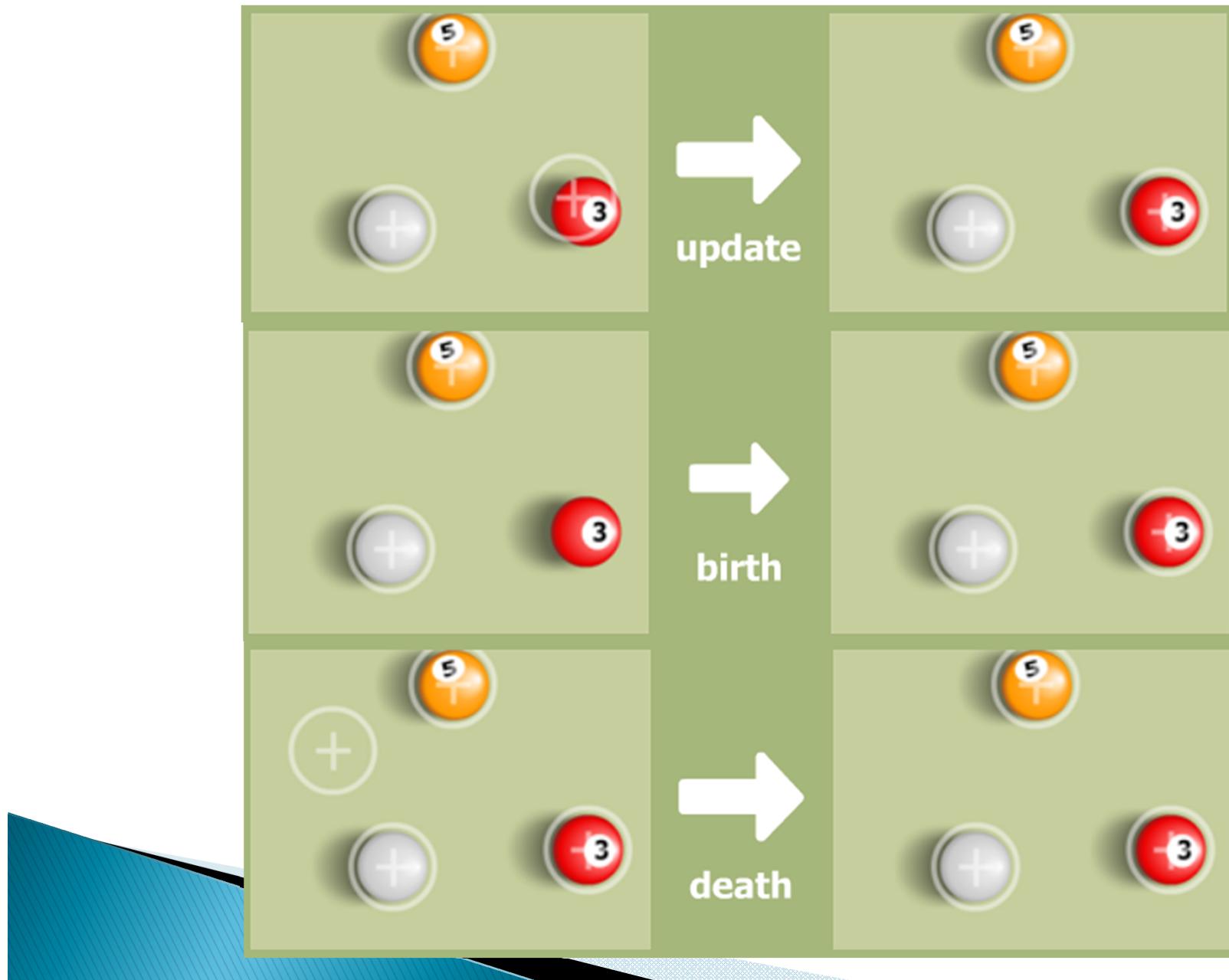
2.3 Reversible jump Markov chain Monte Carlo (RJMCMC)

RJMCMC extends the classical MCMC to handle the unknown number of objects



Exploring the configuration space

2.3 Reversible jump Markov chain Monte Carlo (RJMCMC)



2.4 Posterior probability as modified Bayesian model

- ◆ The posterior probability of each hypothesis is decomposed into a likelihood term and a prior term

$$P(X^{(t)} | I^{(t)}, X^{*(t-1)}) = P(X^{(t)} | X^{*(t-1)}) * P(I^{(t)} | X^{(t)}) \quad (2.2)$$

The prior
term

The likelihood
term

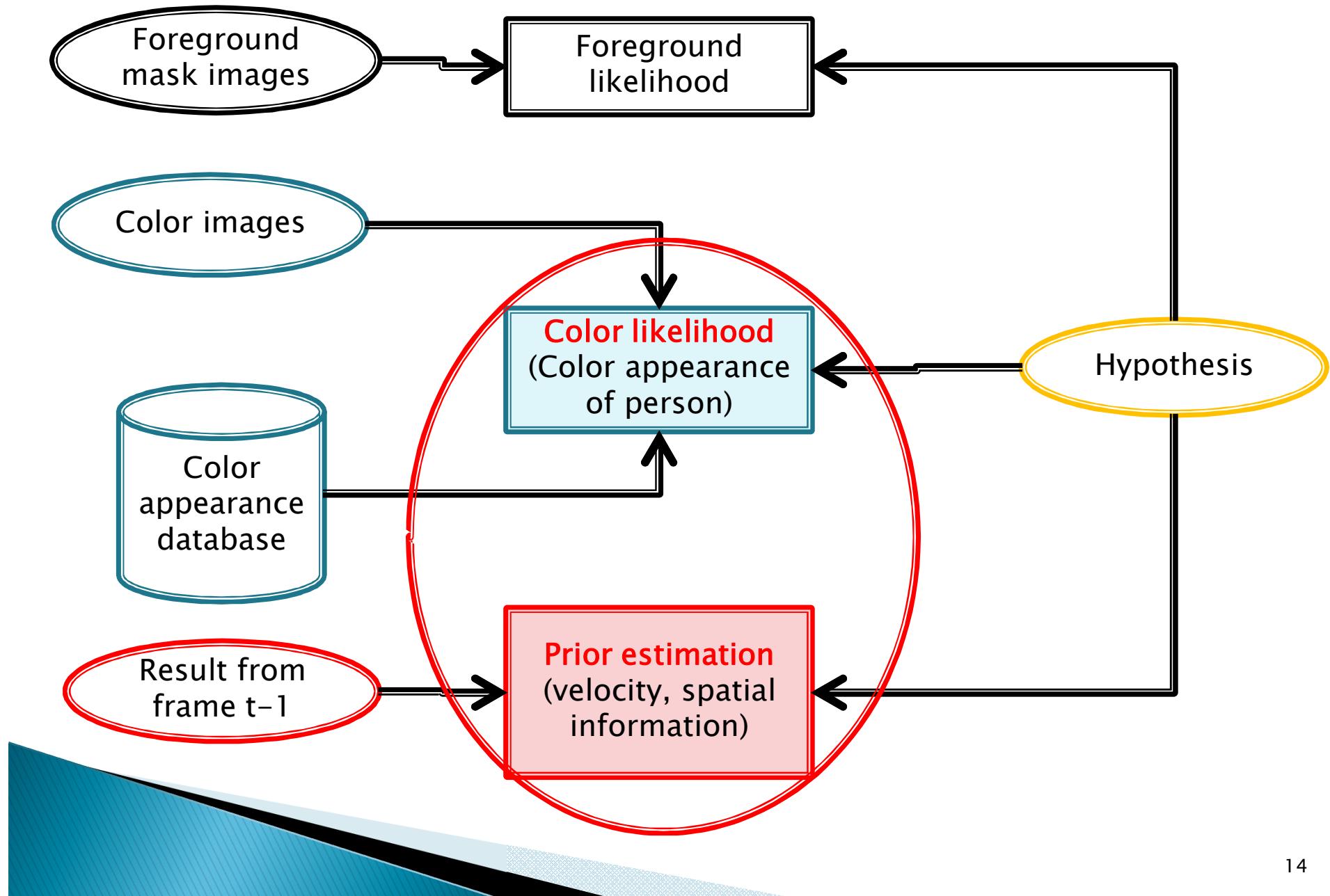
- ◆ Our likelihood probability is decomposed into the summarization of the background likelihood and the color likelihood:

$$P(I^{(t)} | X^{(t)}) = \alpha * P(BI^{(t)} | X^{(t)}) + \beta * P(CI^{(t)} | X^{(t)}, \text{color_database}) \quad (2.3)$$

The background
likelihood term

The color likelihood
term

2.4 The components of posterior computation



3.1 The Background Likelihood

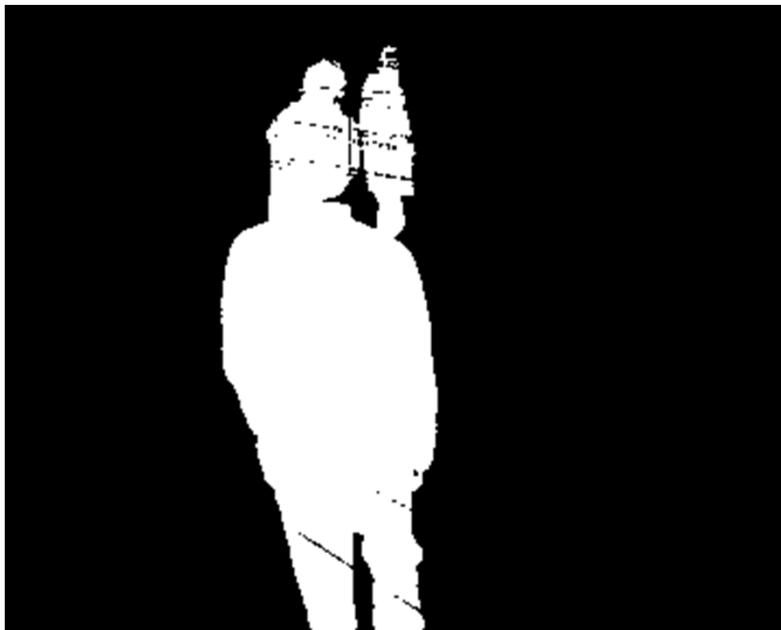


Figure 3.1:
Foreground mask



Figure 3.2:
Human model

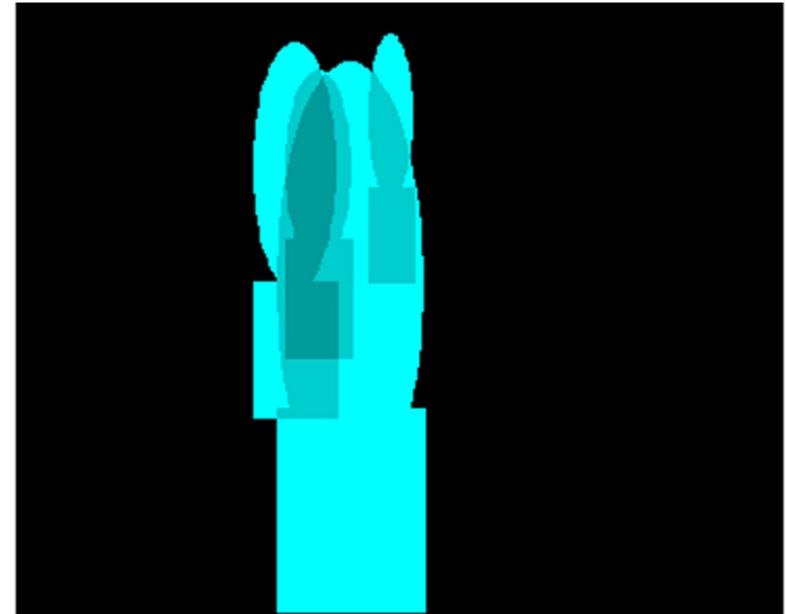
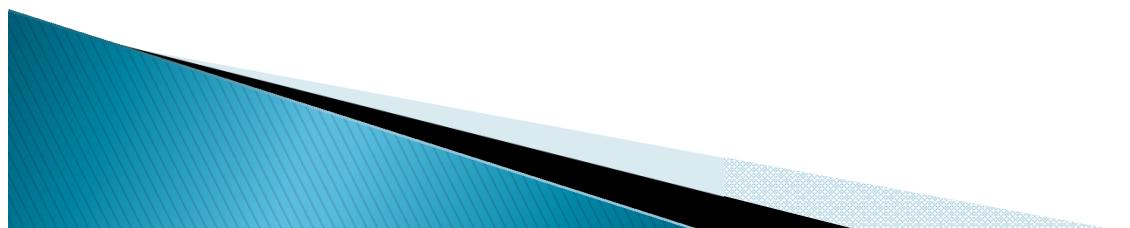
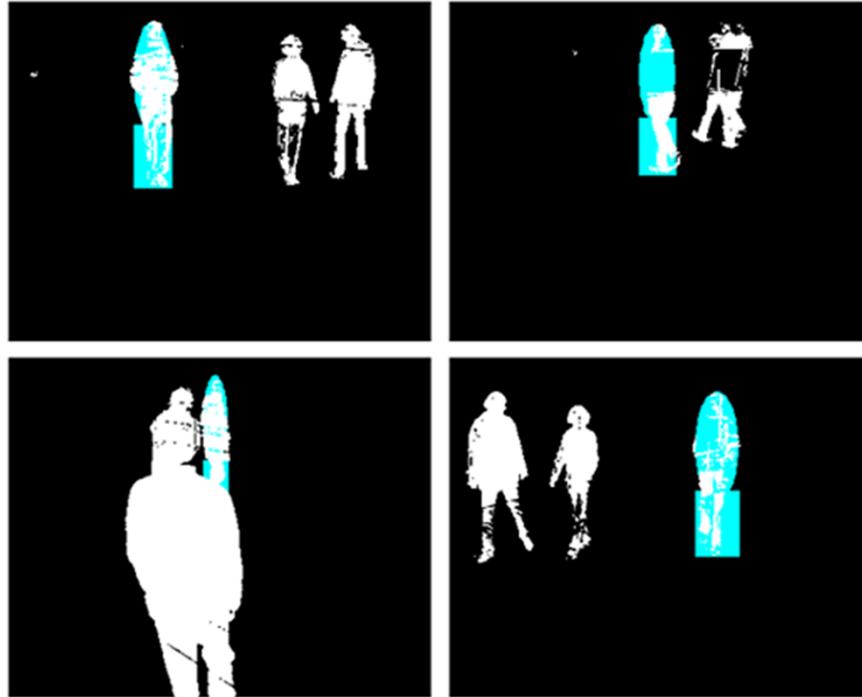


Figure 3.2: Synthesis image
of the hypothesis by
projecting human model into
the image



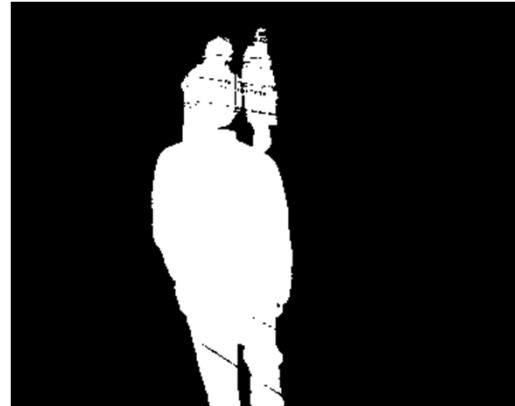
3.1 The Background Likelihood



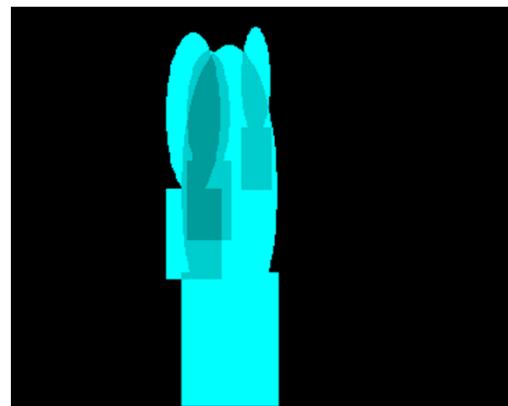
- ◆ For fusing all camera views, the general background likelihood probability is computed by summarizing the background likelihood probability of all camera views v :

$$P(BI^{(t)} | X^{(t)}) = \sum_{v=1}^V e^{-\rho(BI_v^{(t)}, SBI_v^{(t)})} \quad (3.1)$$

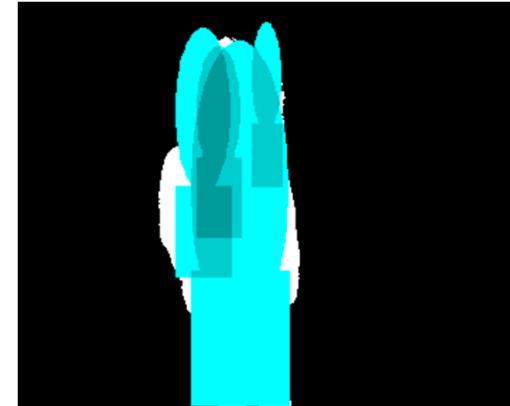
3.1 The Background Likelihood



Foreground image



Synthesis image of hypothesis



Distance between
Foreground and Synthesis

- ◆ For each view, the normalized pseudo-distance between the foreground mask BI and the synthesis foreground image SBI is computed by counting the number of foreground pixel in the area of synthesis image

$$\rho(BI, SBI) = \frac{\sum_{i=1, j=1}^{Width, Height} p_{ij}^{BI} * (1 - p_{ij}^{SBI}) + p_{ij}^{SBI} * (1 - p_{ij}^{BI})}{\sum_{i=1, j=1}^{Width, Height} p_{ij}^{SBI}} \quad (3.2)$$

3.2 The Color Likelihood



Human model

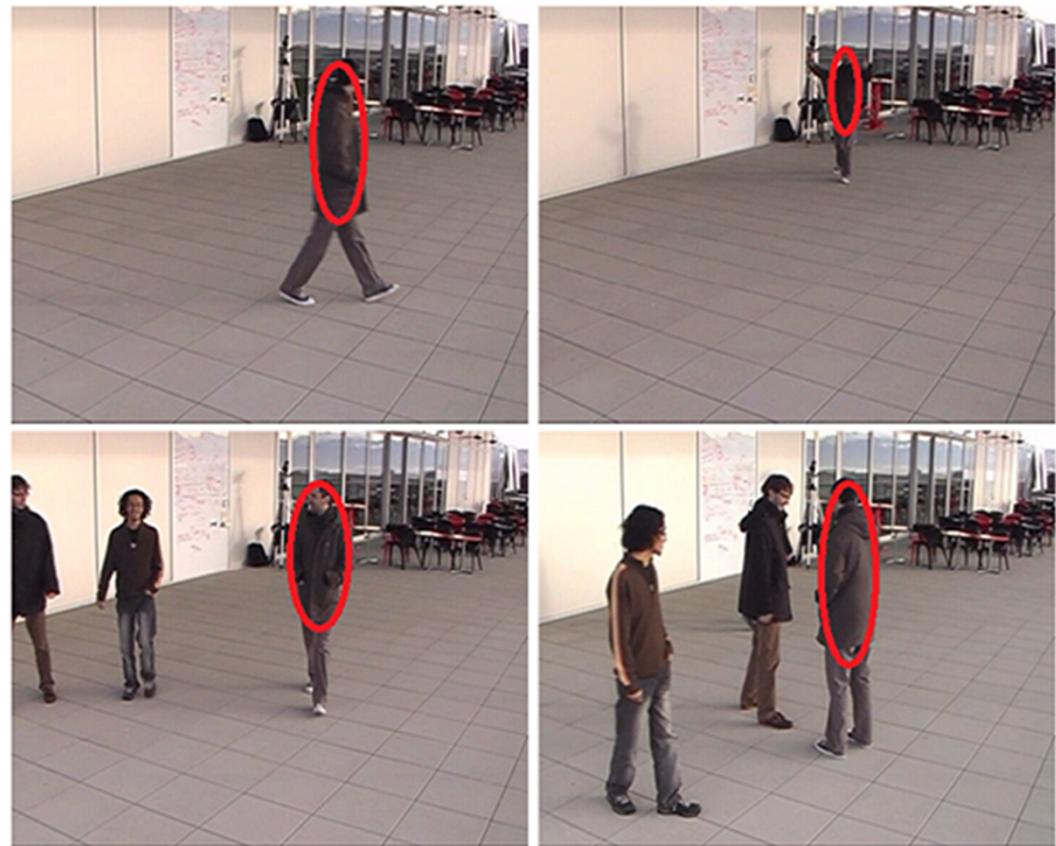


Figure 3.4: Ellipse shape for color matching



3.2 The Color Likelihood

- ◆ The color matching step has been done separately each view.
- ◆ Result of all views are summarized to contribute the final color likelihood probability following the below formulation:

$$P(CI^{(t)} | X^{(t)}, \text{color_database}) = \sum_{v=1}^V e^{-\sigma(CI_v^{(t)} | X^{(t)}, \text{color_database}_v^{(t)})} \quad (3.3)$$



Figure 3.5: The difference in color distribution between different views

3.2 The Color Likelihood

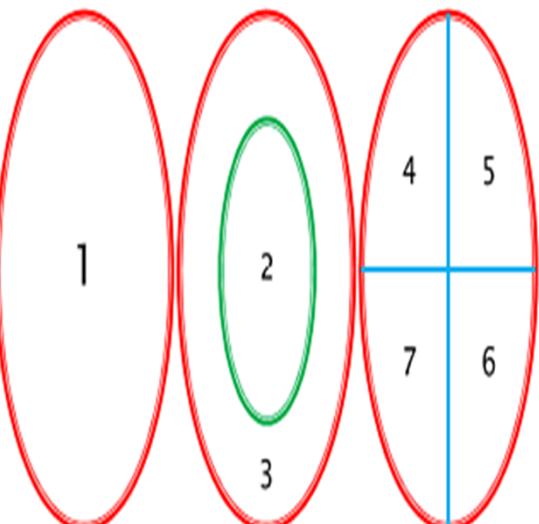
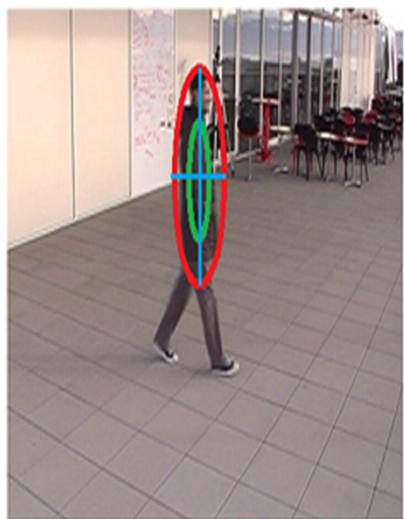
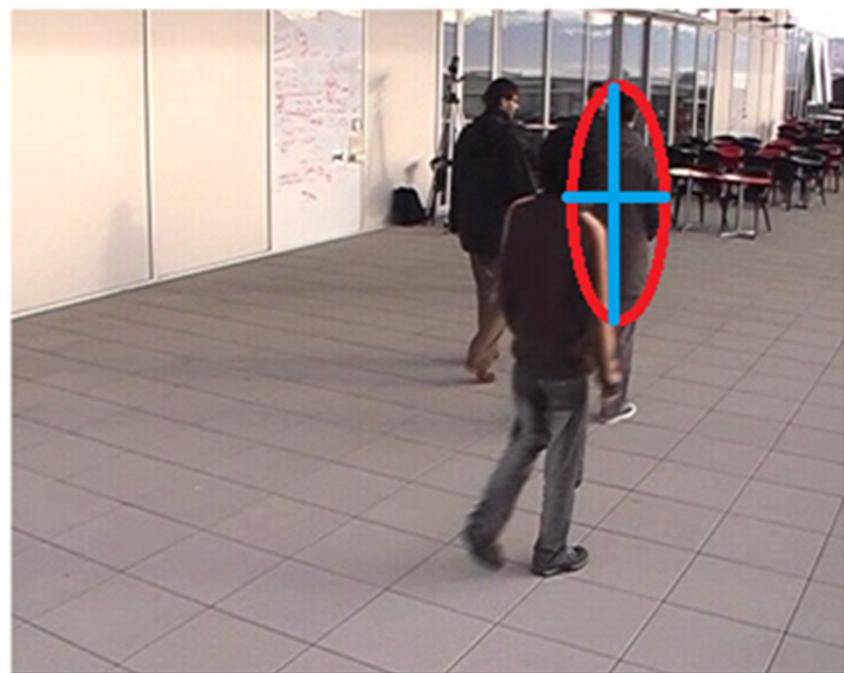


Figure 3.6: Seven parts of ellipse model

Figure 3.7: Visible and Invisible parts of ellipse



3.2 The Color Likelihood

- ◆ The color matching step for each person available in the camera view is computed by the summarization of all 7 parts using Bhattachayya distance.

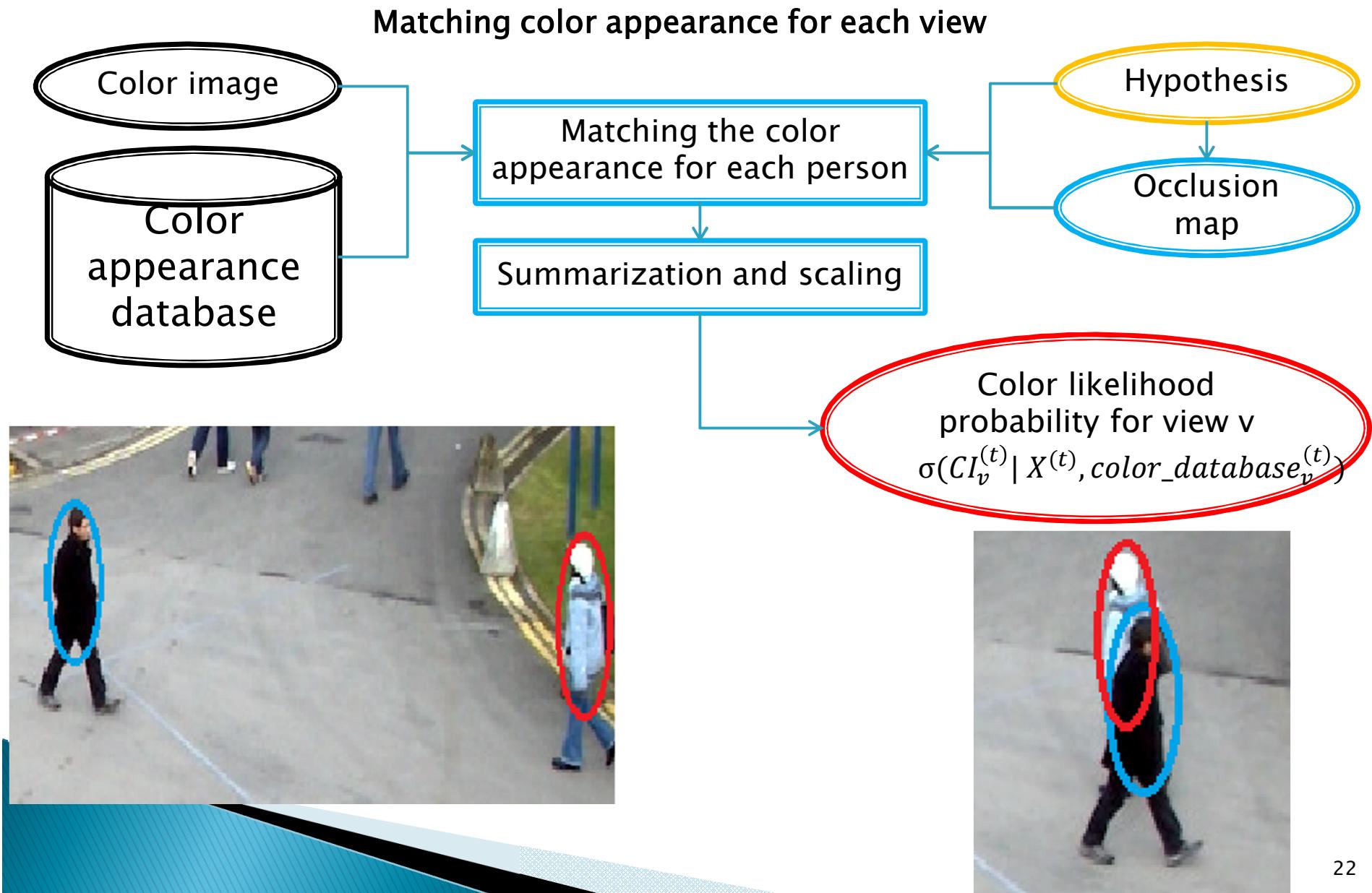
$$\tau(\text{CI} | x^{(t)}) = \sum_{p=1}^7 \delta_p \text{BH}(image_his_p^{(t)}, database_his_p^{(t)}) \quad (3.4)$$

- ◆ The color matching for each view v is the summarization and scaling of the matching of all targets.

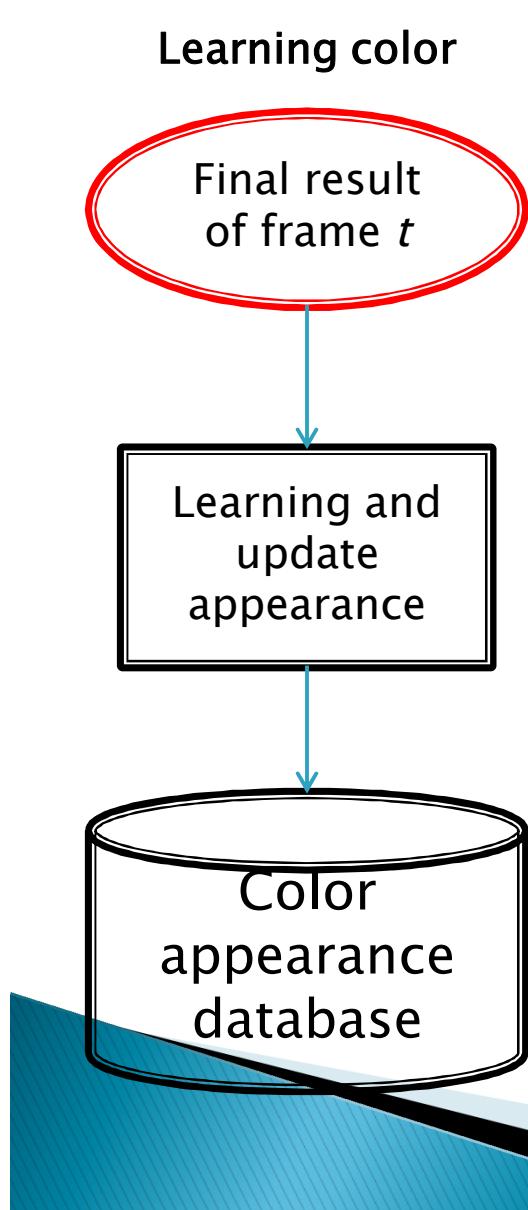
$$\sigma(CI_v^{(t)} | X^{(t)}, \text{color_database}_v^{(t)}) = \frac{\sum_{i=1}^n \tau(CI_v^{(t)} | x_i^{(t)}) * visible_i}{\sum_{i=1}^n visible_i} \quad (3.5)$$



3.2 The Color Likelihood



3.2 The Color Likelihood



Color appearance from previous old frames



Result of tracking at frame t

- Store the black coat person and update him with the new color at frame t .
- Don't update the blue one because she is invisible in the new frame.

3.3 The Prior

- ◆ We realize that tendency of moving is also the important information in pedestrian tracking.
- ◆ Our assumptions are related to three kinds of penalty for this prior information.
 - ◆ Dynamic model
 - ◆ Mutual exclusion model
 - ◆ Ghost detection

$$P(X^{(t)} | X^{*(t-1)}) = e^{-\text{dyn}(X^{(t)} | X^{*(t-1)}, X^{*(t-2)}, \dots)} * e^{-\text{exc}(X^{(t)})} * e^{-\text{gho}(BI^{(t)} | X^{(t)})} \quad (3.6)$$



3.3 The Prior

- ◆ People do not change their velocity in the short time. So the velocity v of people in the batch of 15 frames should be equal. **The dynamic model** is:

$$\text{dyn}(X^{(t)} | X^{*(t-1)}, X^{*(t-2)}, \dots) = \sum_{i=1}^n \sum_{f=t-15}^t |v_i^f - v_i^{f+1}|^2$$

$$v_i^f = c_i^{f+1} - c_i^f. \quad (3.7)$$

- ◆ Two objects cannot occupy the same space simultaneously. **The mutual exclusion model** is:

$$\text{exc}(X^{(t)}) = \sum_{i \neq j} \frac{\text{close_dist}^2}{|c_i^t - c_j^t|^2} \quad (3.8)$$

- ◆ If the number of foreground pixels covered by the human model of each human state is too small, **the ghost detection** will determine a human hypothesis of human is not valid.

$$\text{gho}(BI^{(t)} | x_i^{(t)}) = \begin{cases} 1 & \exists v, s.t. \frac{|SBI^v(x_i^{(t)}) \cap BI^v|}{|SBI^v(x_i^{(t)})|} \leq \text{ghost_const} \\ 0 & \text{otherwise} \end{cases} \quad (3.9)$$

4.1 Evaluation–Multiple Object Detection Precision

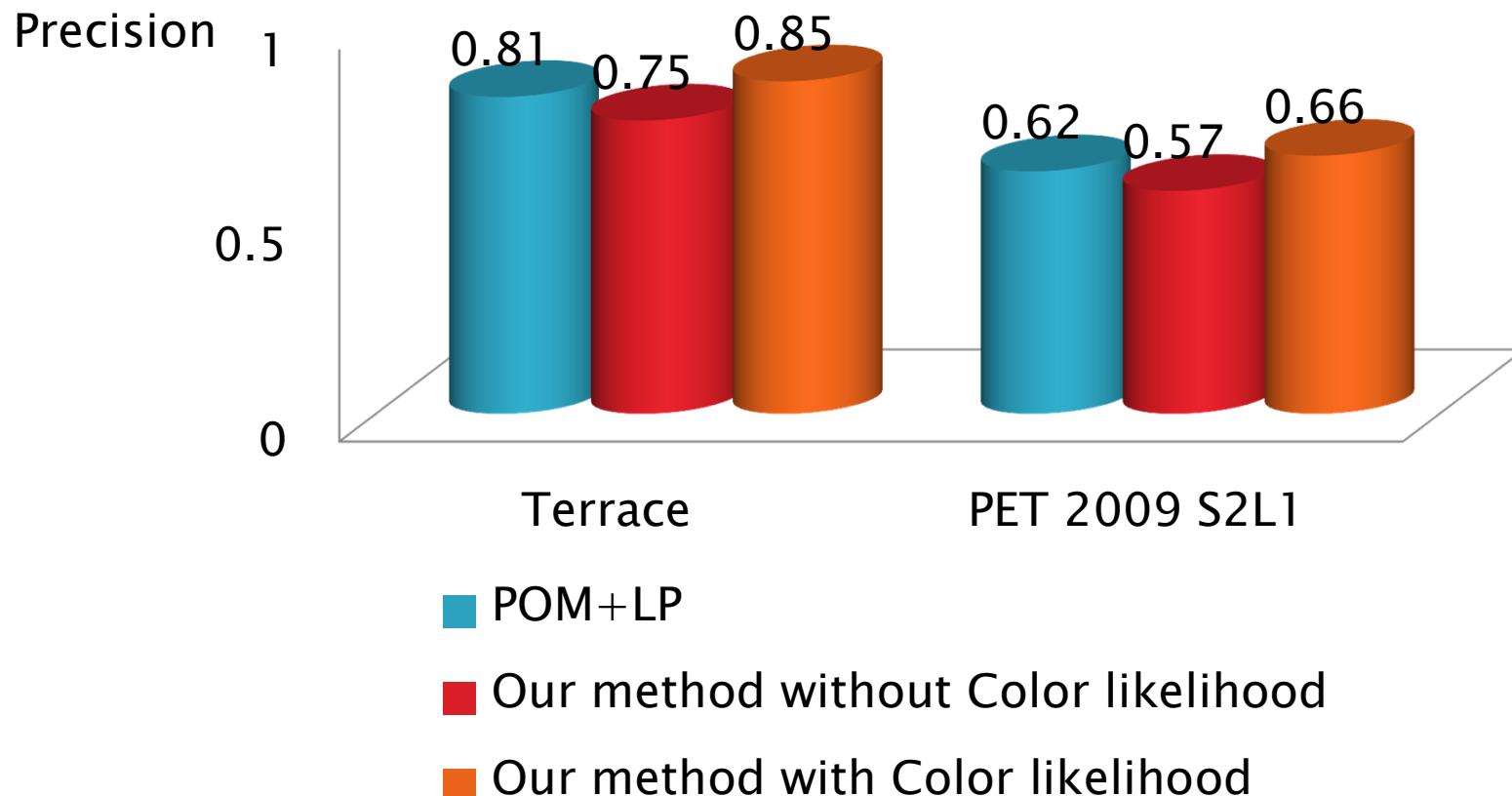


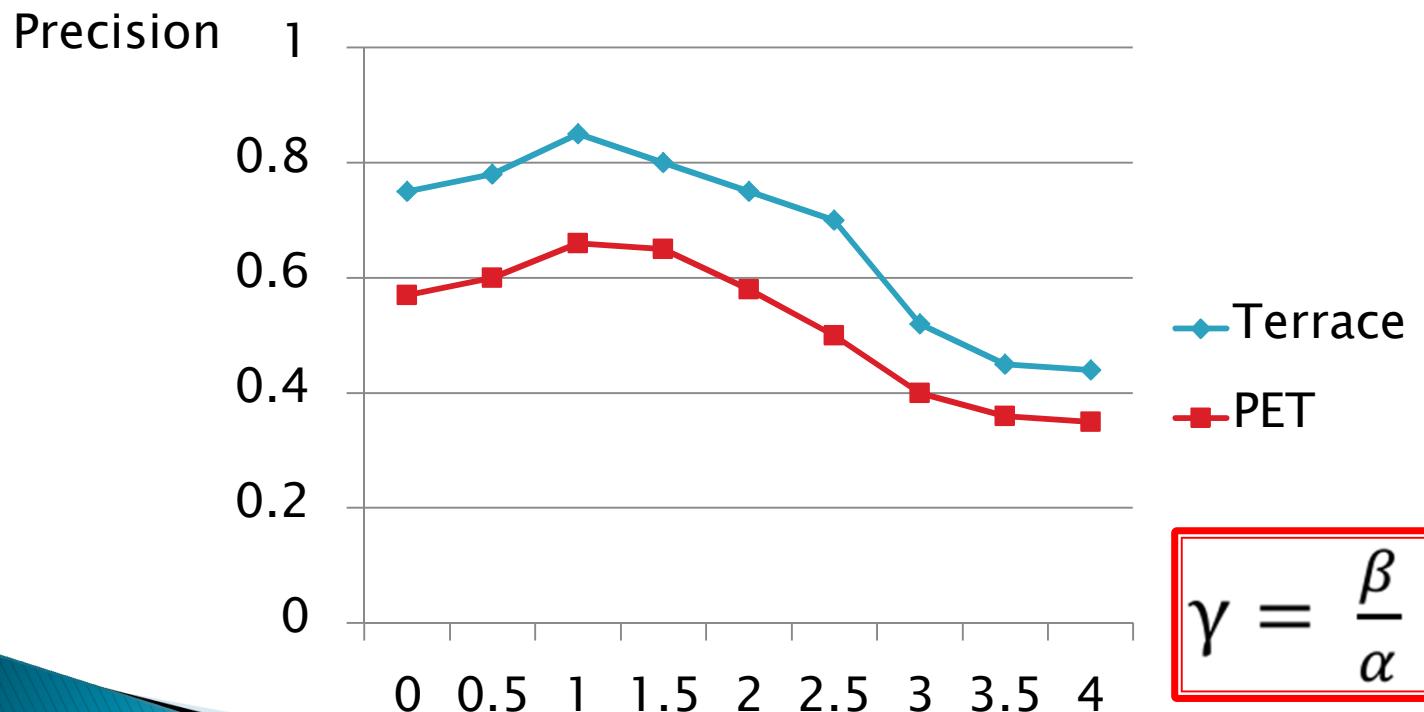
Figure 4.1: The Multiple Object Detection Precision (MODP) evaluation between POM+LP and our proposed methods

4.2 Evaluation–Determining α and β

The background likelihood

The color likelihood

$$P(I^{(t)} | X^{(t)}) = \alpha * P(BI^{(t)} | X^{(t)}) + \beta * P(CI^{(t)} | X^{(t)}, \text{color_database})$$



$$\gamma = \frac{\beta}{\alpha}$$

Figure 4.2: Illustration for changing the ratio

5.1 Conclusion

- A vision system capable of detecting multiple people, tracking them and automatically learning their color appearance.
- The color appearance was very powerful not only in improving the precision of locating the people, but also overcoming the identities switching problem.
- The tendency of moving such as spatial location and velocity was smoothing the tracking trajectories.
- The system also has some weakness:
 - The system could not recovery when the system lost its targets because the detection step is not separate from the tracking step.
 - The system is depended on the accuracy of calibration parameters. The error of calibration decreased the precision of the system



5.2 Future works

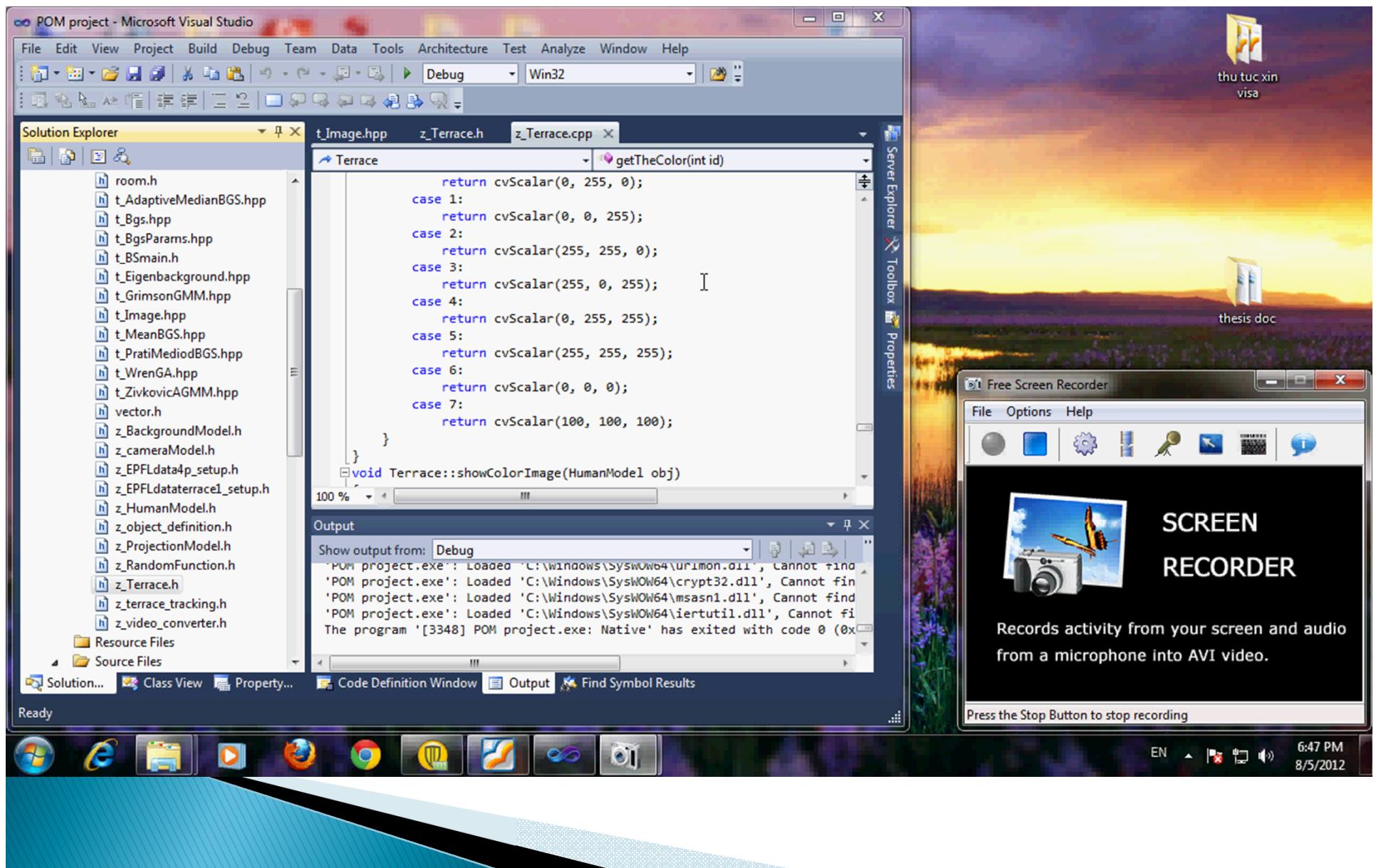
- Using more sophisticated appearance, detecting the people falling down, or jumping up.
- Using more powerful features.
- Using faster processor with more efficient implementation to reach real-time performance.
- Combining MCMC with faster deterministic local search.



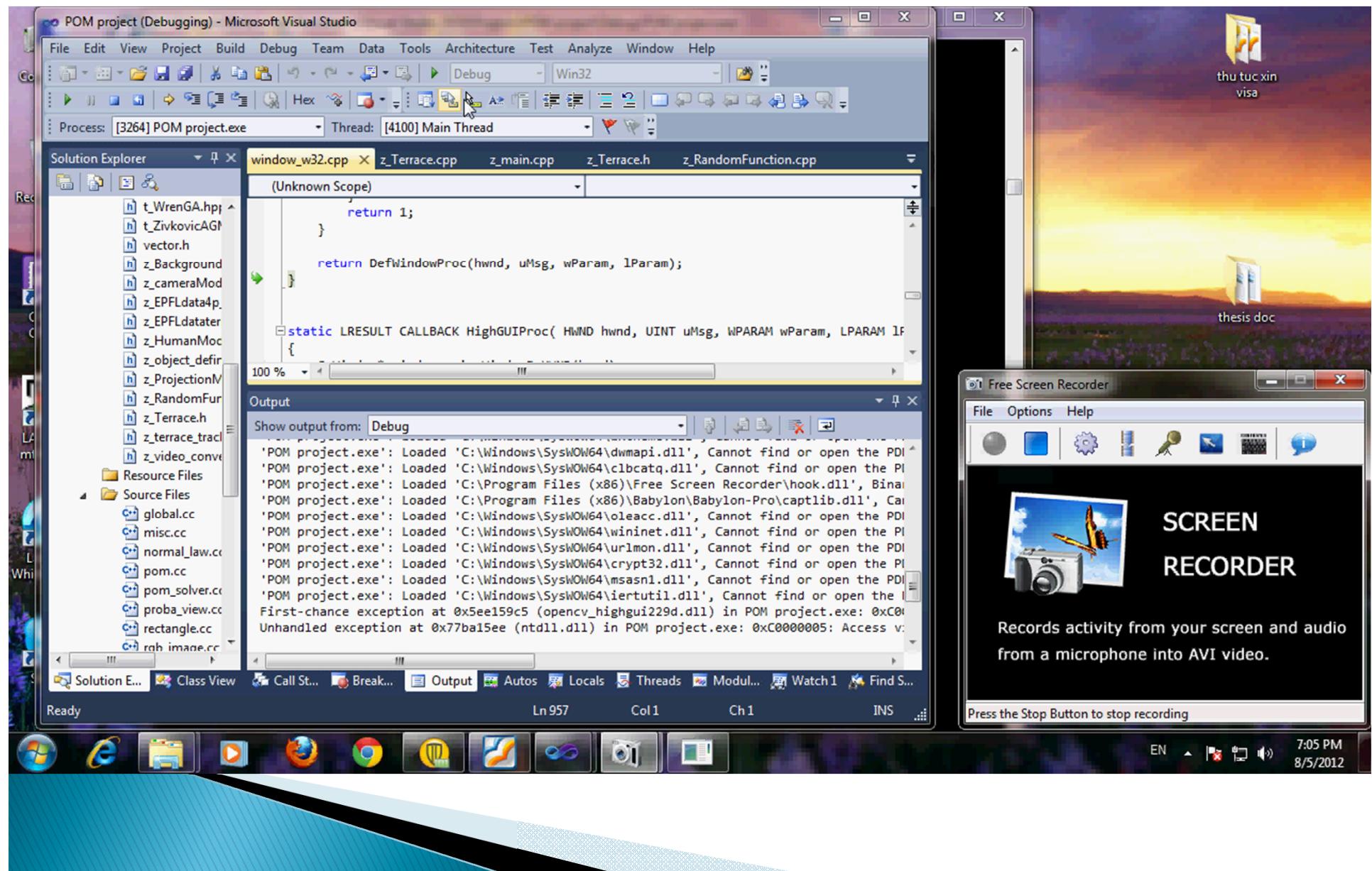
Thank you
ご静聴ありがとうございます



Demo



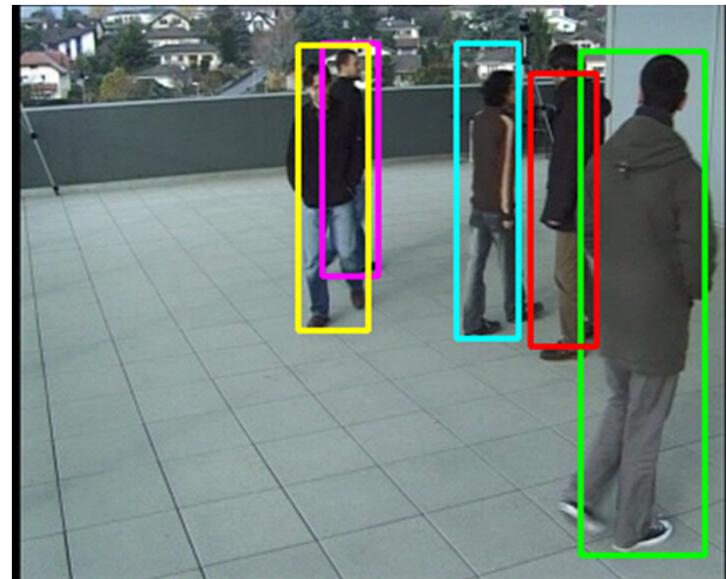
Demo



Identities switching without using color likelihood

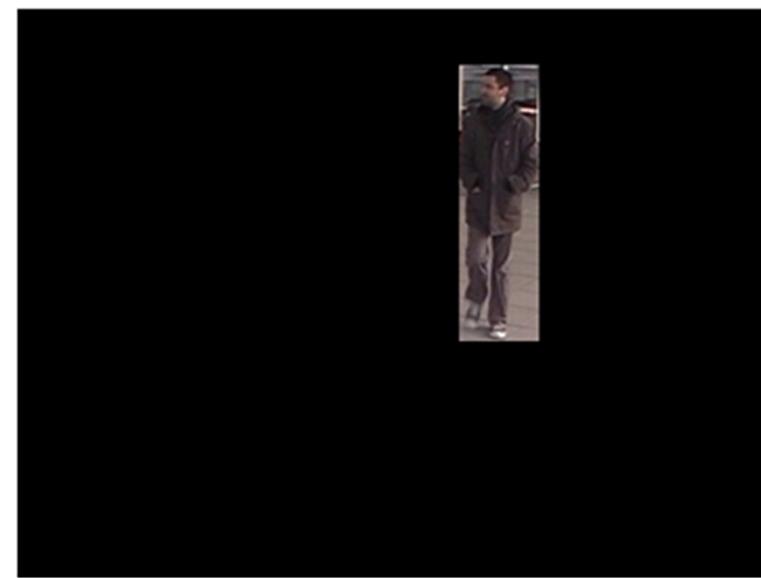
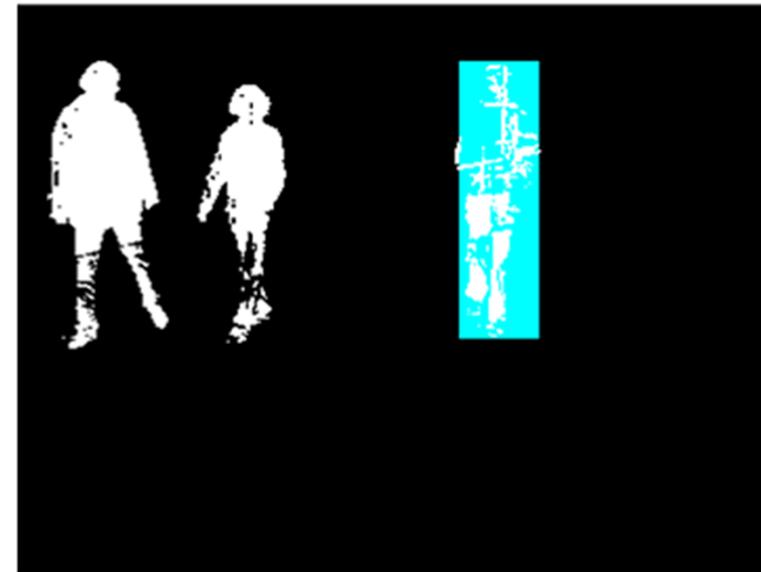
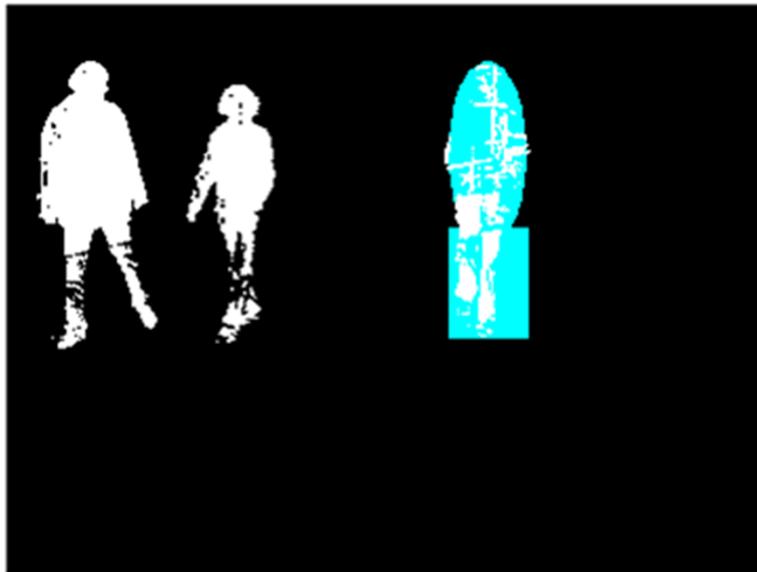


Tracking result at frame $t-1$

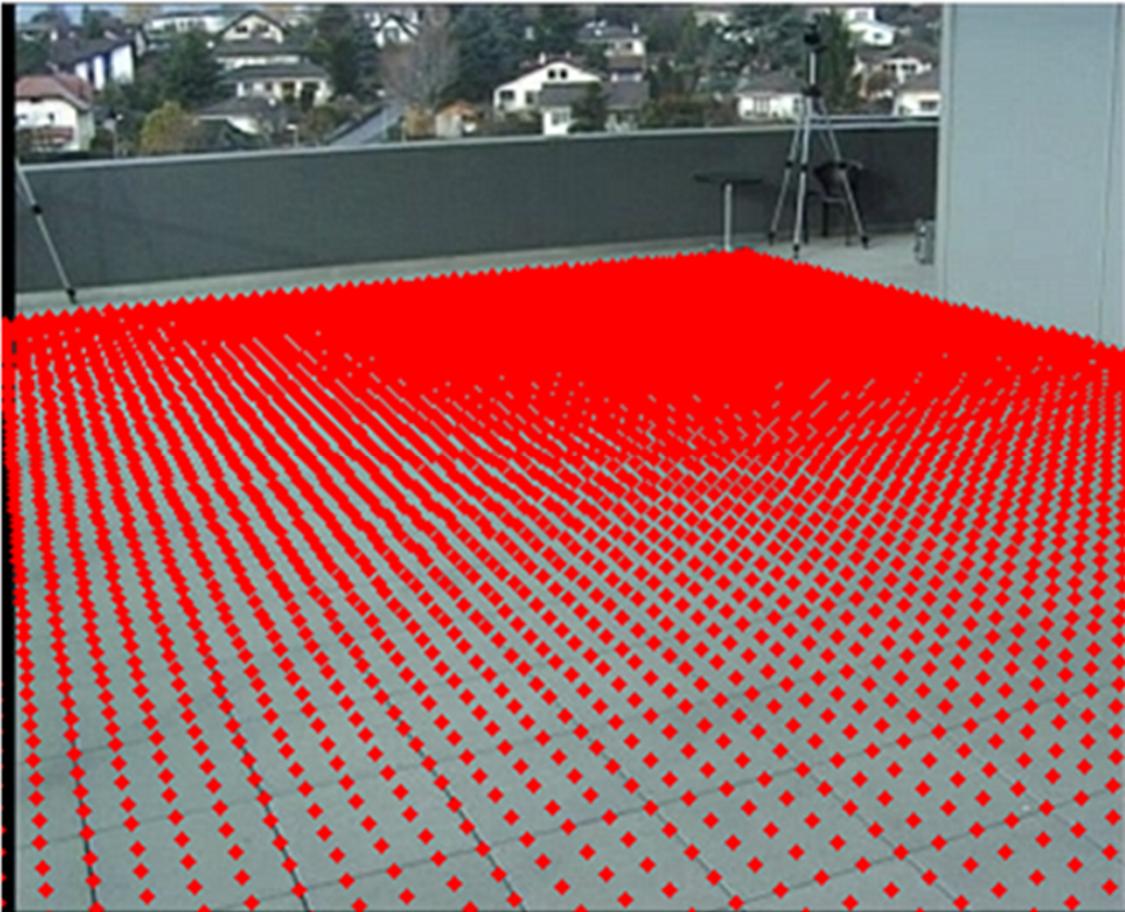


Tracking result at frame t

Human model approximate better



Grid Discretization and Camera Calibration

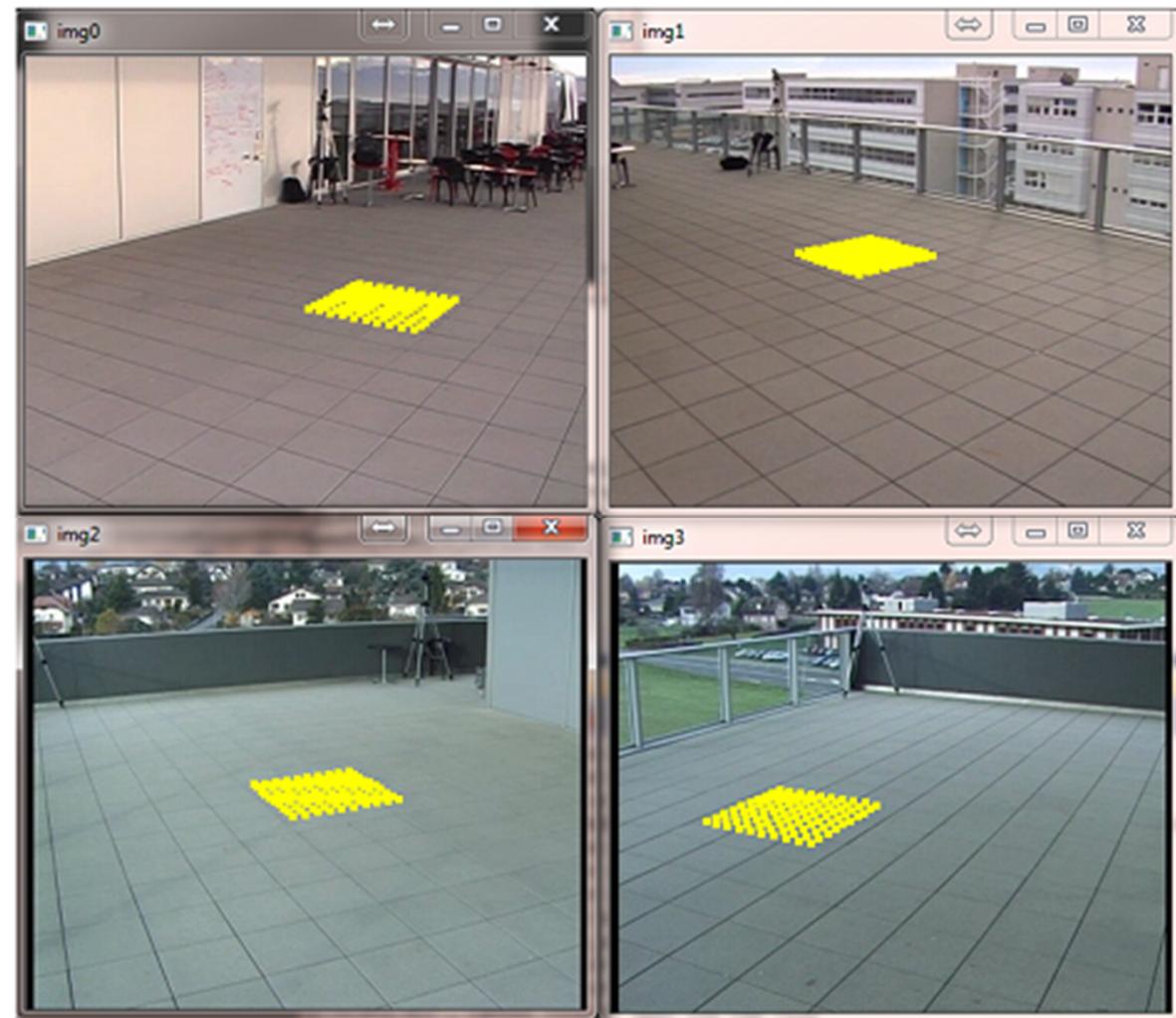


Ground plane is discretized into a finite number of locations. The red points are the locations

- The ground is discretized into grid of location.
- The advantage of this discretization is that we could implement our method as a discrete problem.
- The speed of the system is also increased by preventing the complexity of continuous computation.

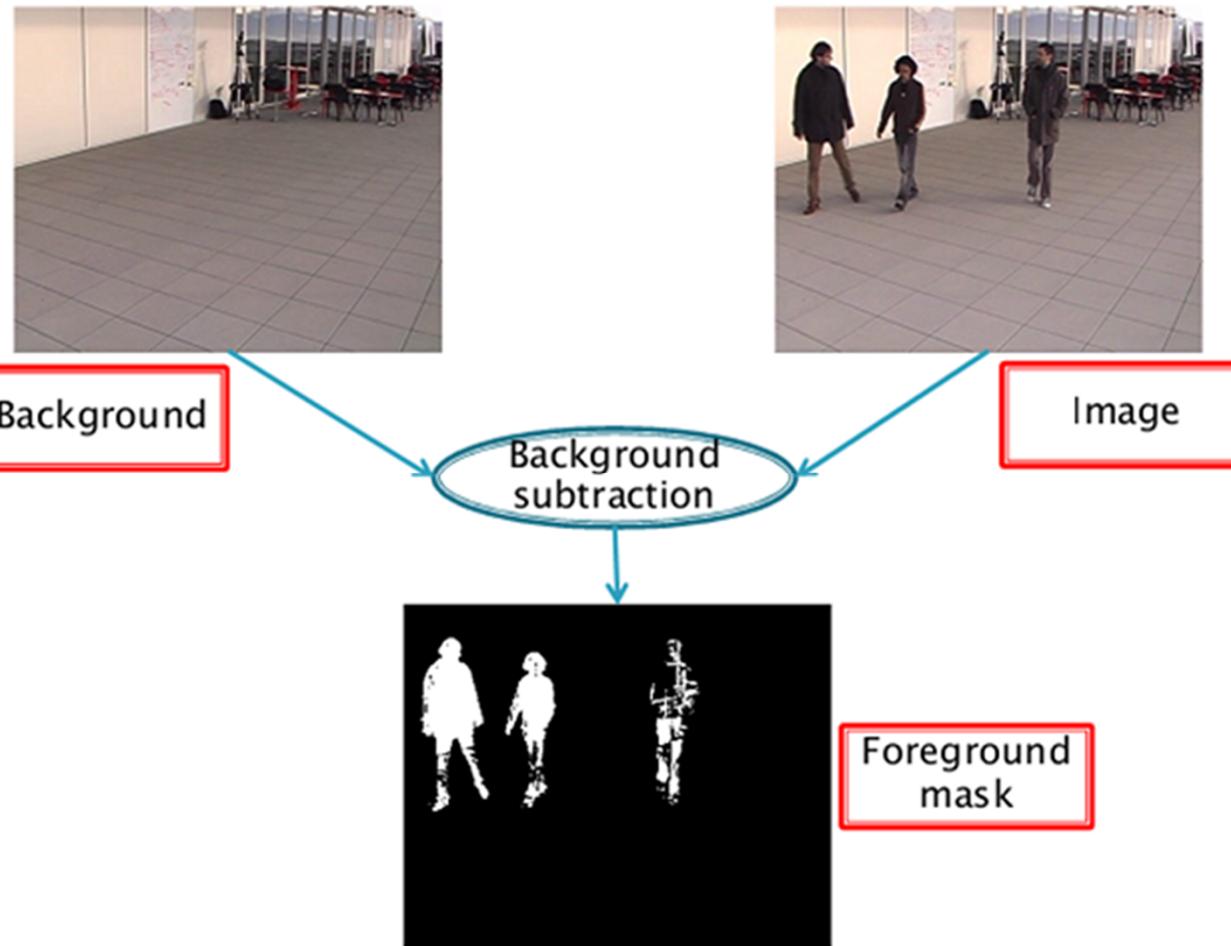
Grid Discretization and Camera Calibration

- Calibration is the method which finding the corresponding points between different views.
- We used the calibration parameters provided by the dataset.
- There are many methods to calibrate multiple cameras such as Tsai or auto calibration.



Corresponding points in all views show the result of calibration

Background Subtraction



The background subtraction