

## Visual Tracking

### Computer Vision I CSE 252A Lecture 17

CSE 252A, Fall 2014

Computer Vision I

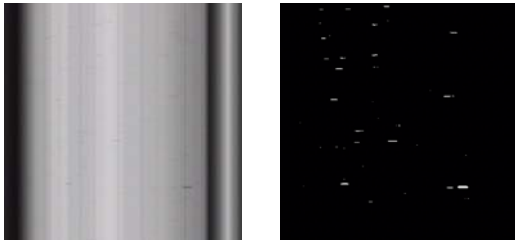
## Announcements

- Read Chapter 11 of Forsyth & Ponce
- Homework 3 is due today by 11:59 PM
- Homework 4 is due Dec 18, 11:59 PM
- Please complete evaluations
  - Course
  - TA

CSE 252A, Fall 2014

Computer Vision I

## Foreground/Background Segmentation



CSE 252A, Fall 2014

Computer Vision I

Input Image



CSE 252A, Fall 2014

Computer Vision I

## Background Image

Initially,

$$B_0(x, y) = \text{median}\{I_n(x, y) | 0 \leq n < N\}$$

Thereafter, updated in a sequential estimation of the mean manner

$$B_n(x, y) = \begin{cases} (1-\alpha)B_{n-1}(x, y) + \alpha I_n(x, y) & \text{if } I_n(x, y) \in \text{background} \\ B_{n-1}(x, y) & \text{otherwise} \end{cases}$$

CSE 252A, Fall 2014

Computer Vision I

## Background Image



CSE 252A, Fall 2014

Computer Vision I

### Difference Image

$$D_n^2(x, y) = (I_n(x, y) - B_{n-1}(x, y))^2$$



CSE 252A, Fall 2014

Computer Vision I

### Segmentation Image

Statistic-based threshold operation: assumes independent Gaussian distributions  
 Sequential estimation of the mean of squares

$$J_n^2(x, y) = \begin{cases} (1-\alpha)J_{n-1}^2(x, y) + \alpha I_n^2(x, y) & \text{if } I_n(x, y) \in \text{background} \\ J_{n-1}^2(x, y) & \text{otherwise} \end{cases}$$

Variance estimation (approximation)

$$\sigma_n^2(x, y) = J_n^2(x, y) - B_n^2(x, y)$$

Threshold

$$F_n(x, y) = \begin{cases} 1 & \text{if } D_n^2(x, y) \geq T \sigma_n^2(x, y) \\ 0 & \text{otherwise} \end{cases}$$

CSE 252A, Fall 2014

Computer Vision I

### Segmentation Image



CSE 252A, Fall 2014

Computer Vision I

## Visual Tracking



### Main Challenges

1. 3-D Pose Variation
2. Occlusion of the target
3. Illumination variation
4. Camera jitter
5. Expression variation etc.

[ Ho, Lee, Kriegman ]

CSE 252A, Fall 2014

Computer Vision I

## Main tracking notions

- State: usually a finite number of parameters (a vector) that characterizes the "state" (e.g., location, size, pose, deformation of thing being tracked).
- Dynamics: How does the state change over time? How is that change constrained?
- Representation: How do you represent the thing being tracked
- Prediction: Given the state at time  $t-1$ , what is an estimate of the state at time  $t$ ?
- Correction: Given the predicted state at time  $t$ , and a measurement at time  $t$ , update the state.
- Initialization – what is the state at time  $t=0$ ?

CSE 252A, Fall 2014

Computer Vision I

## What is state?

- 2-D image location,  $\Phi=(u, v)$
- Image location + scale  $\Phi=(u, v, s)$
- Image location + scale + orientation  $\Phi=(u, v, s, \theta)$
- Affine transformation
- 3-D pose
- 3-D pose plus internal shape parameters (some may be discrete).
  - e.g., for a face, 3-D pose +facial expression using FACS + eye state (open/closed).
- Collections of control points specifying a spline
- Above, but for multiple objects (e.g. tracking a formation of airplanes).
- Augment above with temporal derivatives  $(\phi, \dot{\phi})$

CSE 252A, Fall 2014

Computer Vision I

## State Examples:

- object is ball, state is 3D position+velocity, measurements are stereo pairs
- object is person, state is body configuration, measurements are frames
- What is state here?



CSE 252A, Fall 2014

Computer Vision I

## Example: Blob Tracker

- From input image  $I(u,v)$  (color?) at time  $t$ , create a binary image by applying a function  $f(I(u,v))$ .
- Clean up binary image using morphological operators
- Perform connected component exploration to find “blobs.” – connected regions.
- Compute their moments (mean and covariance of coordinates of region), and use as state
- Using state estimate from  $t-1$  and perform “data association” to identify state in from  $t$ .

CSE 252A, Fall 2014

Computer Vision I

## Blob Tracking in IR Images



- Threshold about body temperature
- Connected component analysis
- Position, scale, orientation of regions
- Temporal coherence

CSE 252A, Fall 2014

Computer Vision I

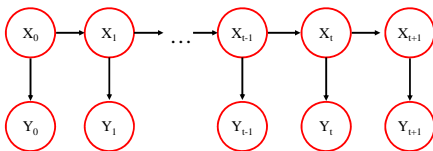
## Tracking: Probabilistic framework

- Very general model:
  - We assume there are moving objects, which have an underlying state  $X$
  - There are measurements  $Y$ , some of which are functions of this state
  - There is a clock
    - at each tick, the state changes:  $X_{t-1}, X_t, X_{t+1}$
    - at each tick, we get a new observation:  $Y_{t-1}, Y_t, Y_{t+1}$

CSE 252A, Fall 2014

Computer Vision I

## Tracking State



- Instead of “knowing state” at each instant, we treat the state as random variables  $X_i$  characterized by a pdf  $P(X_i)$  or perhaps conditioned on other Random Variables e.g.,  $P(X_i / X_{i-1})$ , etc.
- The observation (measurement)  $Y_i$  is a random variable conditioned on the state  $P(Y_i / X_i)$
- Generally, we don’t observe the state – it’s hidden.

CSE 252A, Fall 2014

Computer Vision I

## Three main steps



- **Prediction:** we have seen  $y_0, \dots, y_{i-1}$  — what state does this set of measurements predict for the  $i$ 'th frame? to solve this problem, we need to obtain a representation of  $P(X_i | Y_0 = y_0, \dots, Y_{i-1} = y_{i-1})$ .
- **Data association:** Some of the measurements obtained from the  $i$ -th frame may tell us about the object’s state. Typically, we use  $P(X_i | Y_0 = y_0, \dots, Y_{i-1} = y_{i-1})$  to identify these measurements.
- **Correction:** now that we have  $y_i$  — the relevant measurements — we need to compute a representation of  $P(X_i | Y_0 = y_0, \dots, Y_i = y_i)$ .

We can try to express these conditional distributions parametrically, sample the distribution, or estimate the mode.

CSE 252A, Fall 2014

Computer Vision I

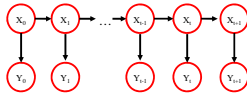
## Simplifying Assumptions

- Only the immediate past matters: formally, we require

$$P(\mathbf{X}_i | \mathbf{X}_1, \dots, \mathbf{X}_{i-1}) = P(\mathbf{X}_i | \mathbf{X}_{i-1})$$

- Measurements depend only on the current state: we assume that  $\mathbf{Y}_i$  is conditionally independent of all other measurements given  $\mathbf{X}_i$ . This means that

$$P(\mathbf{Y}_1 | \mathbf{Y}_1, \dots, \mathbf{Y}_{i-1}, \mathbf{X}_i) = P(\mathbf{Y}_i | \mathbf{X}_i)$$



CSE 252A, Fall 2014

Computer Vision I

## Tracking as induction

- Assume data association is done
  - Sometimes challenging in cluttered scenes. See work by Christopher Rasmussen on Joint Probabilistic Data Association Filters (JPDAF).
- Do correction for the 0<sup>th</sup> frame
- Assume we have corrected estimate for i<sup>th</sup> frame
  - show we can do prediction for i+1 frame, correction for i+1 frame

CSE 252A, Fall 2014

Computer Vision I

## Base case

$P(y | x)$  is our observation model – The probability of  $y$  given  $x$ . For example,  $P(y | x)$  might be a Gaussian with mean  $x$ .

Firstly, we assume that we have  $P(\mathbf{X}_0)$  ← Prior distribution of initial state

And, we make a measurement  $\mathbf{y}_0$

$$\begin{aligned} P(\mathbf{X}_0 | \mathbf{Y}_0 = \mathbf{y}_0) &= \frac{P(\mathbf{y}_0 | \mathbf{X}_0) P(\mathbf{X}_0)}{P(\mathbf{y}_0)} \\ &= \frac{P(\mathbf{y}_0 | \mathbf{X}_0) P(\mathbf{X}_0)}{\int P(\mathbf{y}_0 | \mathbf{X}_0) P(\mathbf{X}_0) d\mathbf{X}_0} \\ &\propto P(\mathbf{y}_0 | \mathbf{X}_0) P(\mathbf{X}_0) \end{aligned}$$

CSE 252A, Fall 2014

Computer Vision I

## Induction step: State Prediction

Given  $P(\mathbf{X}_{i-1} | \mathbf{y}_0, \dots, \mathbf{y}_{i-1})$ .

### Prediction

Prediction involves representing

$$P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1})$$

Our independence assumptions make it possible to write

$$\begin{aligned} P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) &= \int P(\mathbf{X}_i, \mathbf{X}_{i-1} | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) d\mathbf{X}_{i-1} \\ &= \int P(\mathbf{X}_i | \mathbf{X}_{i-1}, \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) P(\mathbf{X}_{i-1} | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) d\mathbf{X}_{i-1} \\ &= \int P(\mathbf{X}_i | \mathbf{X}_{i-1}) P(\mathbf{X}_{i-1} | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) d\mathbf{X}_{i-1} \end{aligned}$$

CSE 252A, Fall 2014

Computer Vision I

## Induction step: State Correction

In prediction, we estimated the state  $\mathbf{X}_i$  given the measurements up to  $i-1$ . Now we get the measure at time  $i$  called  $\mathbf{y}_i$ .

### Correction

Correction involves obtaining a representation of

$$P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_i)$$

Our independence assumptions make it possible to write

$$\begin{aligned} P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_i) &= \frac{P(\mathbf{X}_i, \mathbf{y}_0, \dots, \mathbf{y}_i)}{P(\mathbf{y}_0, \dots, \mathbf{y}_i)} \\ &= \frac{P(\mathbf{y}_i | \mathbf{X}_i, \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) P(\mathbf{y}_0, \dots, \mathbf{y}_{i-1})}{P(\mathbf{y}_0, \dots, \mathbf{y}_i)} \\ &= \frac{P(\mathbf{y}_i | \mathbf{X}_i) P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) P(\mathbf{y}_0, \dots, \mathbf{y}_{i-1})}{P(\mathbf{y}_0, \dots, \mathbf{y}_i)} \\ &= \frac{P(\mathbf{y}_i | \mathbf{X}_i) P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1})}{\int P(\mathbf{y}_i | \mathbf{X}_i) P(\mathbf{X}_i | \mathbf{y}_0, \dots, \mathbf{y}_{i-1}) d\mathbf{X}_i} \end{aligned}$$

CSE 252A, Fall 2014

Computer Vision I

## How is this formulation used

1. It's ignored. At each time instant, the state is estimated (perhaps a maximum likelihood estimate or something non-probabilistic)
2. The conditional distributions are represented by some convenient parametric form (e.g., Gaussian).
3. The PDF's are represented non-parametrically, and sampling techniques are used.

CSE 252A, Fall 2014

Computer Vision I

## Linear dynamic models

- Use notation  $\sim$  to mean “has the pdf of”,  $N(a, b)$  is a normal distribution with mean  $a$  and covariance  $b$ .
- A linear dynamic model has the form

$$\mathbf{x}_i = N(\mathbf{D}_{i-1} \mathbf{x}_{i-1}; \Sigma_{d_i})$$

$$\mathbf{y}_i = N(\mathbf{M}_i \mathbf{x}_i; \Sigma_{m_i})$$

CSE 252A, Fall 2014

Computer Vision I

## Examples

- Points moving with constant velocity
- Periodic motion
- Etc.
- Points moving with constant acceleration

CSE 252A, Fall 2014

Computer Vision I

## Points moving with constant velocity

- We have

$$u_i = u_{i-1} + \Delta t v_{i-1} + \varepsilon_i \quad \text{Position}$$

$$v_i = v_{i-1} + \zeta_i \quad \text{Velocity}$$

– (the Greek letters denote noise terms)

- Stack  $(u, v)$  into a single state vector

$$\begin{pmatrix} u \\ v \end{pmatrix}_i = \begin{pmatrix} 1 & \Delta t \\ 0 & 1 \end{pmatrix} \begin{pmatrix} u \\ v \end{pmatrix}_{i-1} + \text{noise}$$

$\mathbf{D}_i$

which is the form we had above

CSE 252A, Fall 2014

Computer Vision I

## Points moving with constant acceleration

- We have

$$u_i = u_{i-1} + \Delta t v_{i-1} + \varepsilon_i$$

$$v_i = v_{i-1} + \Delta t a_{i-1} + \zeta_i$$

$$a_i = a_{i-1} + \xi_i$$

– (the Greek letters denote noise terms)

- Stack  $(u, v)$  into a single state vector

$$\begin{pmatrix} u \\ v \\ a \end{pmatrix}_i = \begin{pmatrix} 1 & \Delta t & 0 \\ 0 & 1 & \Delta t \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} u \\ v \\ a \end{pmatrix}_{i-1} + \text{noise}$$

– which is the form we had above

CSE 252A, Fall 2014

Computer Vision I

## The Kalman Filter

- Key ideas:
  - Linear models interact uniquely well with Gaussian noise - make the prior Gaussian, everything else Gaussian and the calculations are easy
  - Gaussians are really easy to represent – Just a mean vector mean and covariance matrix.

CSE 252A, Fall 2014

Computer Vision I

## The Kalman Filter in 1D

- Dynamic Model

$$x_i \sim N(d_i x_{i-1}, \sigma_{d_i}^2)$$

$$y_i \sim N(m_i x_i, \sigma_{m_i}^2)$$

- Notation

mean of  $P(X_i | y_0, \dots, y_{i-1})$  as  $\bar{X}_i^-$  — Predicted mean

Corrected mean — mean of  $P(X_i | y_0, \dots, y_i)$  as  $\bar{X}_i^+$

the standard deviation of  $P(X_i | y_0, \dots, y_{i-1})$  as  $\sigma_i^-$   
of  $P(X_i | y_0, \dots, y_i)$  as  $\sigma_i^+$

CSE 252A, Fall 2014

Computer Vision I

## Prediction for 1-D Kalman filter

- The new state is obtained by
  - multiplying old state by known constant
  - adding zero-mean noise
- Therefore, predicted mean for new state is
  - constant times mean for old state
- Predicted variance is
  - sum of constant<sup>2</sup> times old state variance and noise variance

### Because:

- Old state is normal random variable,
- Multiplying normal rv by constant implies
  - mean is multiplied by a constant
  - variance is multiplied by square of constant
- Adding zero mean noise adds zero to the mean,
- Adding rv's adds variance

CSE 252A, Fall 2014

Computer Vision I

### Dynamic Model:

$$x_i \sim N(d_i x_{i-1}, \sigma_d)$$

$$y_i \sim N(m_i x_i, \sigma_{m_i})$$

Start Assumptions:  $\bar{x}_0$  and  $\sigma_0^-$  are known  
Update Equations: Prediction

$$\bar{x}_i = d_i \bar{x}_{i-1}$$

$$\sigma_i^- = \sqrt{\sigma_d^2 + (d_i \sigma_{i-1}^-)^2}$$

Update Equations: Correction

$$x_i^+ = \left( \frac{\bar{x}_i \sigma_{m_i}^2 + m_i y_i (\sigma_i^-)^2}{\sigma_{m_i}^2 + m_i^2 (\sigma_i^-)^2} \right)$$

$$\sigma_i^+ = \sqrt{\left( \frac{\sigma_{m_i}^2 (\sigma_i^-)^2}{(\sigma_{m_i}^2 + m_i^2 (\sigma_i^-)^2)} \right)}$$

CSE 252A, F

puter Vision I

## Correction for 1D Kalman filter

- Pattern match to identities given in book
  - basically, guess the integrals, get:

- Notice:

- if measurement noise is small, we rely mainly on the measurement,
- if it's large, mainly on the prediction

$$x_i^+ = \left( \frac{\bar{x}_i \sigma_{m_i}^2 + m_i y_i (\sigma_i^-)^2}{\sigma_{m_i}^2 + m_i^2 (\sigma_i^-)^2} \right)$$

$$\sigma_i^+ = \sqrt{\left( \frac{\sigma_{m_i}^2 (\sigma_i^-)^2}{(\sigma_{m_i}^2 + m_i^2 (\sigma_i^-)^2)} \right)}$$

CSE 252A, Fall 2014

Computer Vision I

## Multi-variate Kalman Filter

### Dynamic Model:

$$x_i \sim N(D_i x_{i-1}, \Sigma_d)$$

$$y_i \sim N(M_i x_i, \Sigma_{m_i})$$

Start Assumptions:  $\bar{x}_0$  and  $\Sigma_0^-$  are known  
Update Equations: Prediction

$$\bar{x}_i = D_i \bar{x}_{i-1}$$

$$\Sigma_i^- = \Sigma_d + D_i \Sigma_{i-1}^- D_i$$

Update Equations: Correction

$$K_i = \Sigma_i^- M_i^T [M_i \Sigma_i^- M_i^T + \Sigma_{m_i}]^{-1}$$

$$\bar{x}_i^+ = \bar{x}_i + K_i [y_i - M_i \bar{x}_i]$$

$$\Sigma_i^+ = I d - K_i M_i \Sigma_i^-$$

CSE 252A, Fall 2014

ter Vision I

## Tracking Modalities

(Define the features  $Y_i$ )

- Color
  - Histogram [Birchfield 1998; Bradski 1998]
  - Volume [Wren *et al.*, 1995; Bregler, 1997; Darrell, 1998]
- Shape
  - Deformable curve [Kass *et al.* 1988]
  - Template [Blake *et al.* 1993; Birchfield 1998]
  - Example-based [Cootes *et al.*, 1993; Baumberg & Hogg, 1994]
- Appearance
  - Correlation [Lucas & Kanade, 1981; Shi & Tomasi, 1994]
  - Photometric variation [Hager & Belhumeur, 1998]
  - Outliers [Black *et al.*, 1998; Hager & Belhumeur, 1998]
  - Nonrigidity [Black *et al.*, 1998; Sclaroff & Isidoro, 1998]
- Motion
  - Background model [Wren *et al.*, 1995; Rosales & Sclaroff, 1999; Stauffer & Grimson, 1999]
  - Optical flow [Cutler & Turk]
  - Egomotion [Sawhney & Ayer, 1996; Irani & Anandan, 1998]
- Stereo
  - Blob correlation [Azarbayejani & Pentland, 1996]
  - Disparity map [Kanade *et al.*, 1996; Konolige, 1997; Darrell *et al.*, 1998]

CSE 252A, Fall 2014

Computer Vision I

## Color Blob tracking



- Color-based tracker gets lost on white knight: Same Color

CSE 252A, Fall 2014

Computer Vision I

## Snakes: Active Contours

- Contour  $C$ : continuous curve on smooth surface in  $\mathcal{R}^3$
- Snake  $S$ : projection of  $C$  to image
- Curve types
  - Edge between regions on surface with contrasting properties
  - Line that contrasts with surface properties on both side
  - Silhouette of surface against contrasting background
- General Algorithm:
  - Perform edge detection
  - Fit parametric or non-parametric curve to data

CSE 252A, Fall 2014

Computer Vision I

## Snakes: Basic Approach

- Parameterize a closed contour
- $\mathbf{r}(s) = \mathbf{q}'\mathbf{B}(s)$  or  $\mathbf{r}(s) = \mathbf{U}(s)\mathbf{Q}$
- Given a predicted state  $\mathbf{q}$ , search radially for edges
- Solve a least squares problem for new state

$$\mathbf{Q} = (q_0^x \quad q_n^x \quad q_0^y \quad q_n^y)$$

$$\mathbf{U}(s) = \begin{pmatrix} \mathbf{B}(s)^T & 0 \\ 0 & \mathbf{B}(s)^T \end{pmatrix}$$



CSE 252A, Fall 2014

Computer Vision I

## Tracker Composition: Only Shape (Snakes)

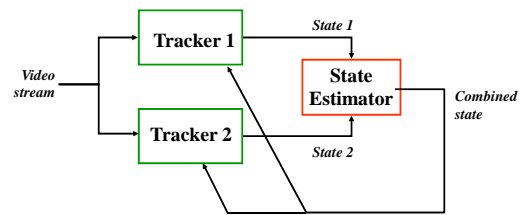


- Geometry-based tracker gets lost on black pawn: Same shape

CSE 252A, Fall 2014

Computer Vision I

## Tracker Composition



CSE 252A, Fall 2014

Computer Vision I

## Tracker Composition: Color and Shape

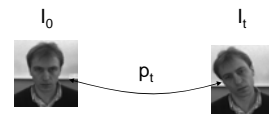


- Combining Trackers => Robustness
- Trackers in video, IR and range

CSE 252A, Fall 2014

Computer Vision I

## Visual Tracking using regions



Variability model:  $I_t = g(I_0, p_t)$

Incremental Estimation: From  $I_0$ ,  $I_{t+1}$  and  $p_t$  compute  $\Delta p_{t+1}$

$$\| I_0 - g(I_{t+1}, p_{t+1}) \|^2 \implies \min$$

CSE 252A, Fall 2014

Computer Vision I

## Tracking using Textured Regions

- Mean intensity difference between  $I$  and affine warp of template image [Shi & Tomasi, 1994]



$$\psi_{region}(x, y) = \sum_{(x,y) \in W} (I_R(x, y) - I_C(x, y))^2$$

CSE 252A, Fall 2014

Computer Vision I

## Image Warping

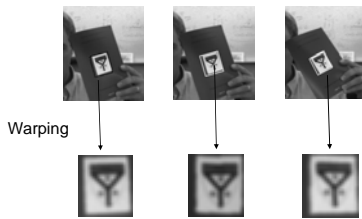
- Warping is a change of coordinates:
  - $J(u, v) = I(f(u, v, p), g(u, v, p))$
- Always prefer to warp to destination to avoid gaps
- Two interpolation schemes
  - nearest neighbor
  - bilinear
- $J(\mathbf{u}) = I(A \mathbf{u})$
- Note that we can “unroll” the loop to avoid the matrix multiply
- For much of tracking, nearest neighbor works well

CSE 252A, Fall 2014

Computer Vision I

## Template tracking: Planar Case

Planar Object => Affine motion model:  $u'_i = A u_i + d$



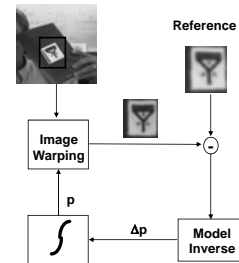
$$I_t = g(\rho_t, I_0)$$

CSE 252A, Fall 2014

Computer Vision I

## Hager/Toyama: Tracking Cycle

- Prediction
  - Prior states predict new appearance
- Image warping
  - Generate a “normalized view”
- Model inverse
  - Compute error from nominal
- State integration
  - Apply correction to state



CSE 252A, Fall 2014

Computer Vision I

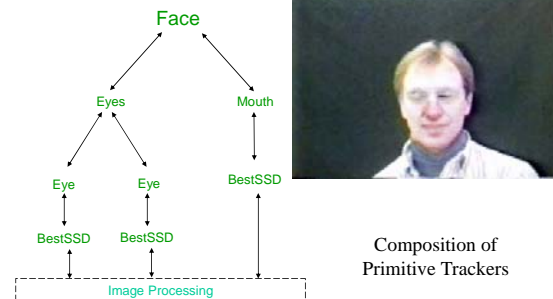
## SSD Tracking



CSE 252A, Fall 2014

Computer Vision I

## XVision: A tracking System



Composition of Primitive Trackers

CSE 252A, Fall 2014

Computer Vision I