



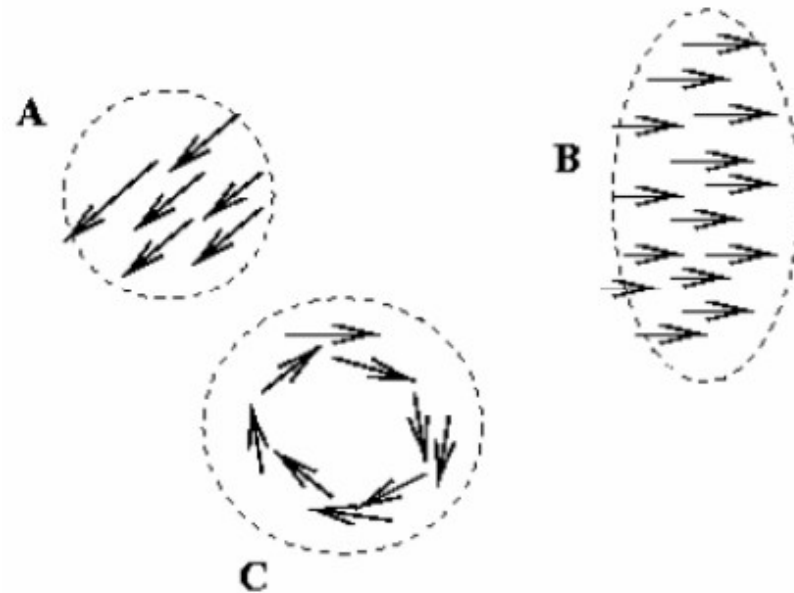
UNIVERSITÀ
CA' FOSCARI
VENEZIA

Motion and Tracking

Andrea Torsello
DAIS
Università Ca' Foscari
via Torino 155,
30172 Mestre (VE)

Motion Segmentation

- Segment the video into multiple coherently moving objects

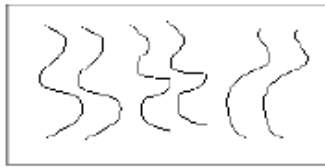


Motion and Perceptual Organization

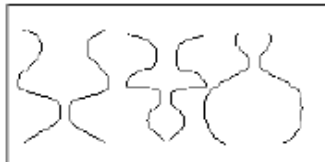
- Sometimes, motion is the only cue



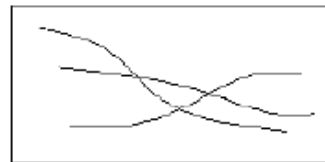
UNIVERSITÀ
CA' FOSCARI
VENEZIA



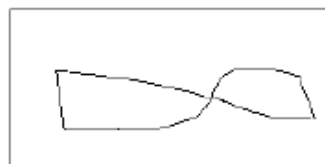
Parallelism



Symmetry



Continuity



Closure



Not grouped



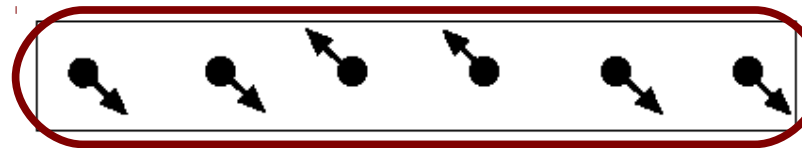
Proximity



Similarity



Similarity



Common Fate



Common Region



Motion and Perceptual Organization

- Sometimes, motion is the foremost cue



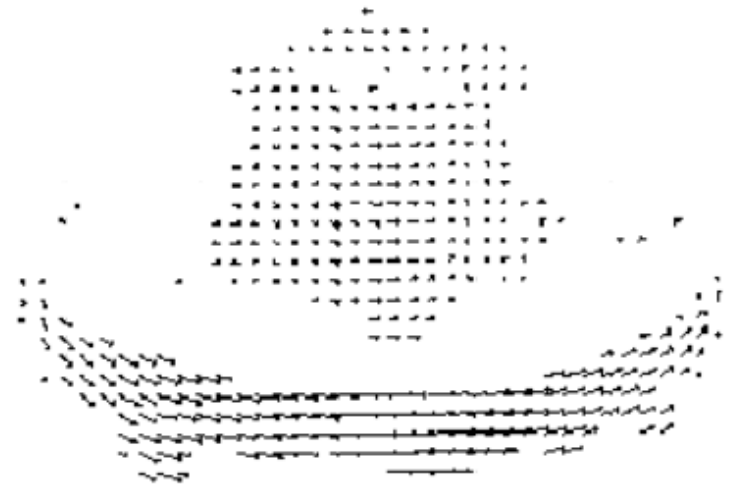
UNIVERSITÀ
CA' FOSCARI
VENEZIA

Motion Field

- The motion field is the projection of the 3D scene motion into the image



UNIVERSITÀ
CA' FOSCARI
VENEZIA



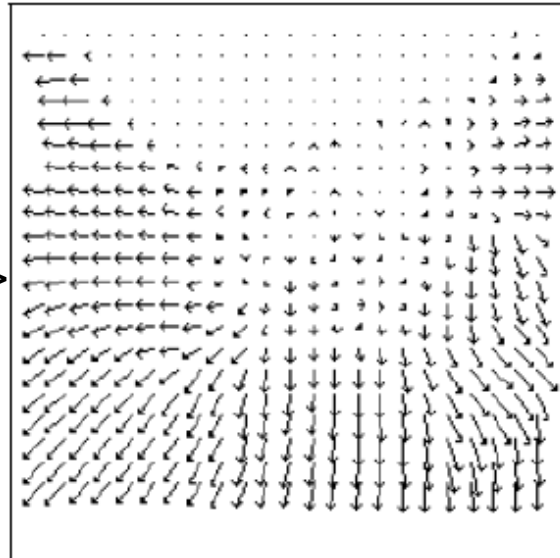
Motion field + camera motion



UNIVERSITÀ
CA' FOSCARI
VENEZIA



⇒



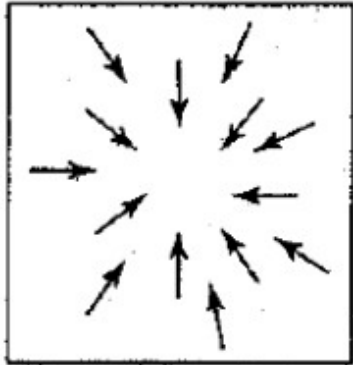
Length of flow
vectors inversely
proportional to
depth Z of 3d
point

points closer to the camera move more
quickly across the image plane

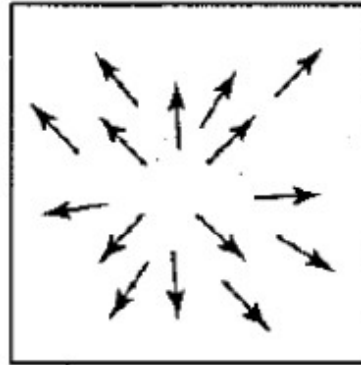
Motion field + camera motion



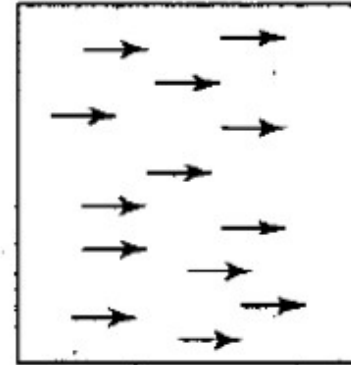
UNIVERSITÀ
CA' FOSCARI
VENEZIA



Zoom out



Zoom in



Pan right to left

Motion Estimation Techniques



UNIVERSITÀ
CA' FOSCARI
VENEZIA

- Direct methods
 - Directly recover image motion at each pixel from spatio-temporal image brightness variations
 - Dense motion fields, fields but sensitive to appearance variations
 - Suitable for video and when image motion is small
- Feature-based methods
 - Extract visual features (corners, textured areas) and track them over multiple frames
 - Sparse motion fields, but more robust tracking
 - Suitable when image motion is large (10s of pixels)

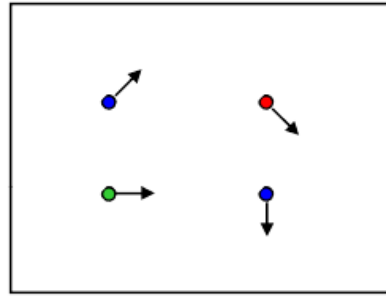
Optical Flow

- **Optical flow** is the apparent motion of brightness patterns in the image
- Ideally, optical flow would be the same as the motion field
- Have to be careful: apparent motion can be caused by lighting changes without any actual motion

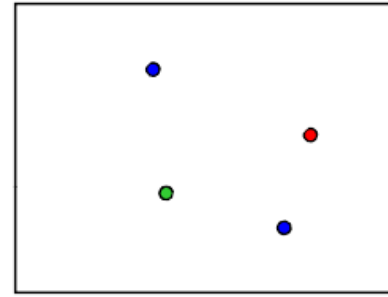


UNIVERSITÀ
CA' FOSCARI
VENEZIA

Optical Flow



$H(x, y)$



$I(x, y)$

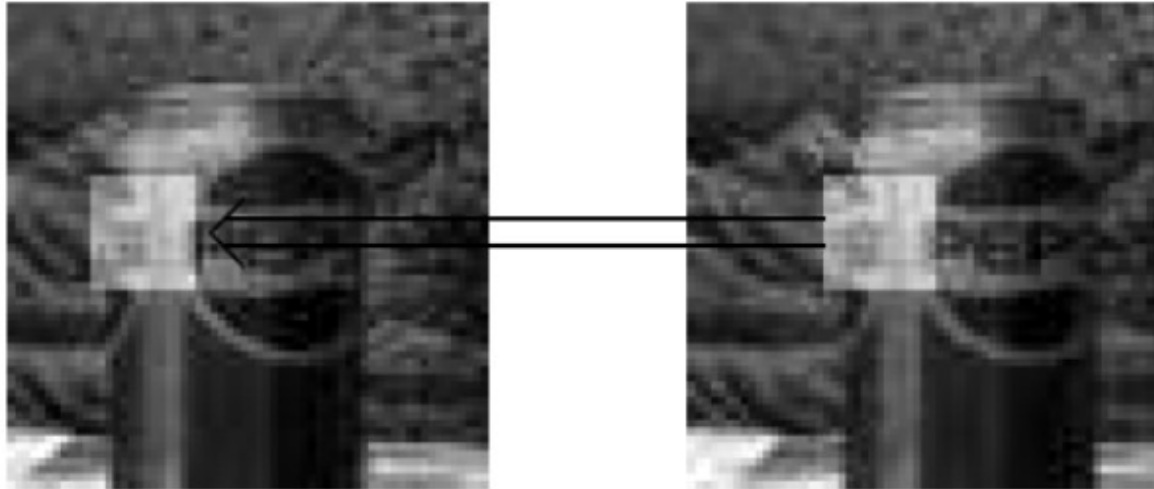
- How to estimate pixel motion from image H to image I ?
 - Solve pixel correspondence problem: given a pixel in H , H look for nearby pixels of the same color in I
- Key assumptions
 - color constancy: a point in H looks the same in I
 - For grayscale images, this is brightness constancy
 - small motion: points do not move very far
- This is called the optical flow problem



Brightness Constancy

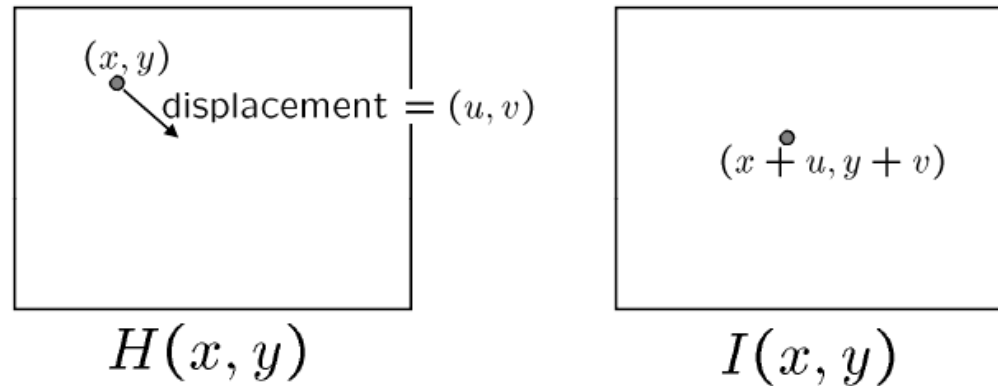


UNIVERSITÀ
CA' FOSCARI
VENEZIA



The highlighted region in the right image looks roughly the same as the region in the left image

Optical Flow Constraints



- Brightness constancy:

$$H(x, y) = I(x + u, y + v)$$

- Small motion:

$$I(x + u, y + v) \approx I(x, y) + \frac{\partial I}{\partial x} u + \frac{\partial I}{\partial y} v$$

- Combining these equations

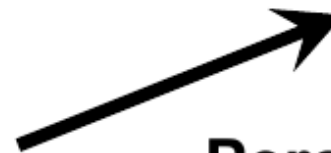
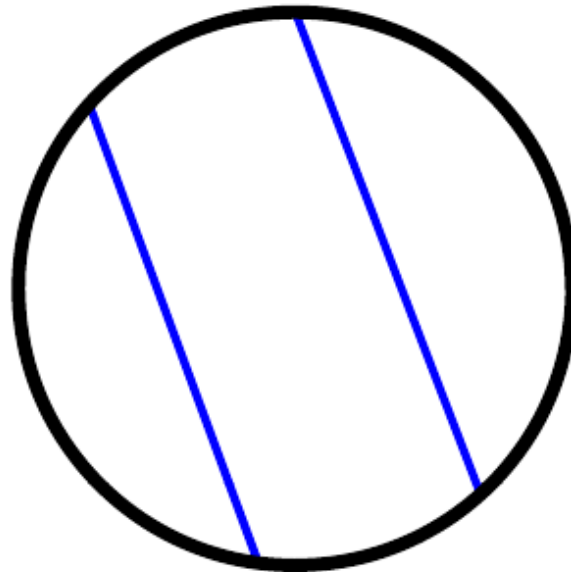
$$\begin{aligned} 0 &= I(x + u, y + v) - H(x, y) \\ &\approx (I(x, y) - H(x, y)) + I_x u + I_y v \\ &\approx I_t + \nabla I(u, v)^T \end{aligned}$$



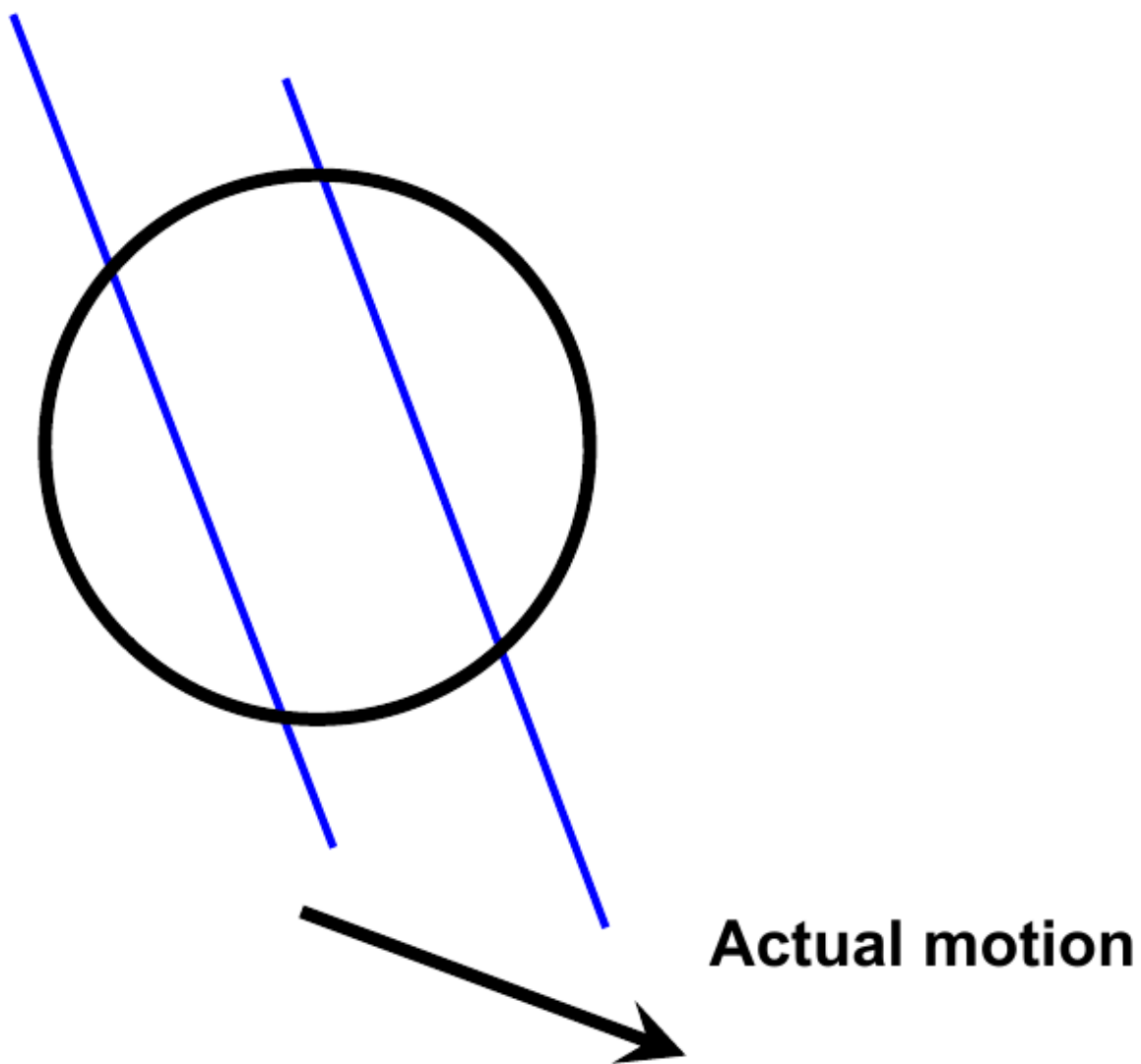
The Aperture Problem



UNIVERSITÀ
CA' FOSCARI
VENEZIA



Perceived motion



The barber pole illusion



UNIVERSITÀ
CA' FOSCARI
VENEZIA



The barber pole illusion



UNIVERSITÀ
CA' FOSCARI
VENEZIA



Solving the Aperture Problem

- How to get more equations for a pixel?
- Spatial coherence constraint: pretend the pixel's neighbors have the same (u,v)
 - If we use a 5x5 window, that gives us 25 equations per pixel

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$



Solving the Aperture Problem

- Prob: we have more equations than unknowns
- Solution: solve least squares problem

$$\begin{matrix} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{matrix} \longrightarrow \text{minimize } \|Ad - b\|^2$$

- Solved by $\begin{matrix} (A^T A) & d = A^T b \\ 2 \times 2 & 2 \times 1 & 2 \times 1 \end{matrix}$

$$\begin{matrix} \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} & \begin{bmatrix} u \\ v \end{bmatrix} = - & \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix} \\ A^T A & & A^T b \end{matrix}$$





$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$ $A^T b$

- When is this solvable?
 - $A^T A$ should be invertible
 - $A^T A$ should not be too small
 - eigenvalues λ_1 and λ_2 of $A^T A$ should not be too small
 - $A^T A$ should be well-conditioned
 - λ_1 / λ_2 should not be too large ($\lambda_1 =$ larger eigenvalue)

Edge



- gradients very large or very small
- large λ_1 , small λ_2



UNIVERSITÀ
CA' FOSCARI
VENEZIA

Low-texture region



- gradients have small magnitude
- small λ_1 , small λ_2



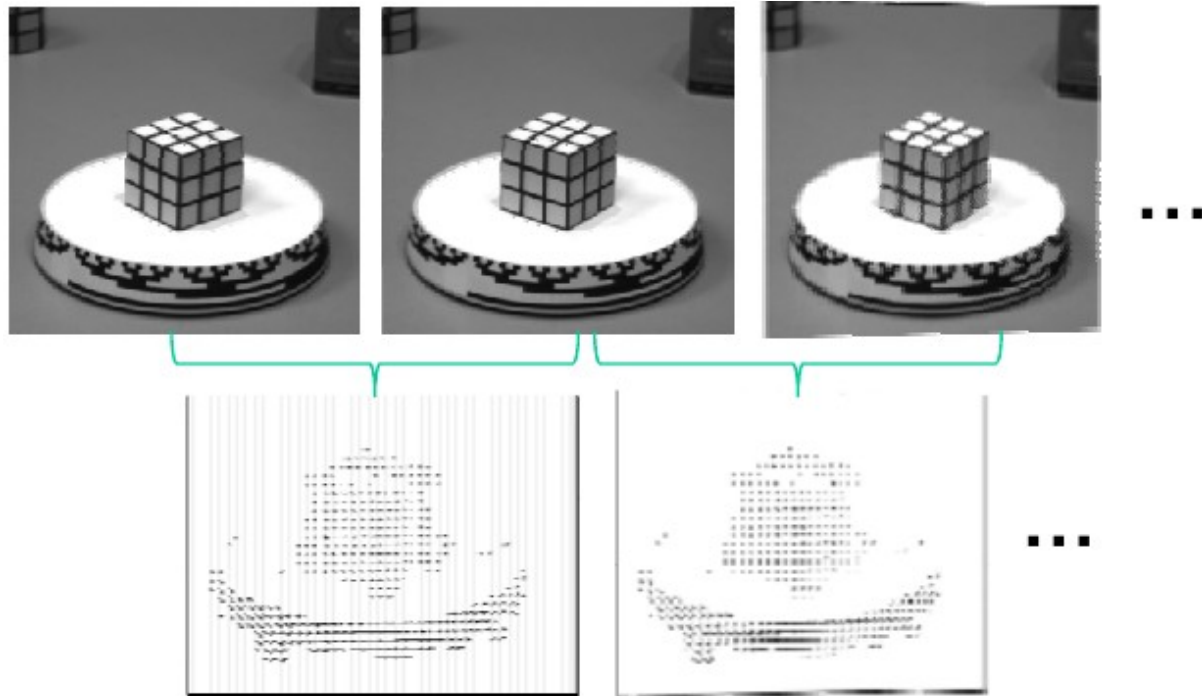
UNIVERSITÀ
CA' FOSCARI
VENEZIA



- gradients are different, large magnitudes
- large λ_1 , large λ_2

Optical flow for tracking

- If we have more than just a pair of frames, we could compute flow from one to the next:



- But flow only reliable for small motions, and we may have occlusions, textureless regions that yield bad estimates anyway...



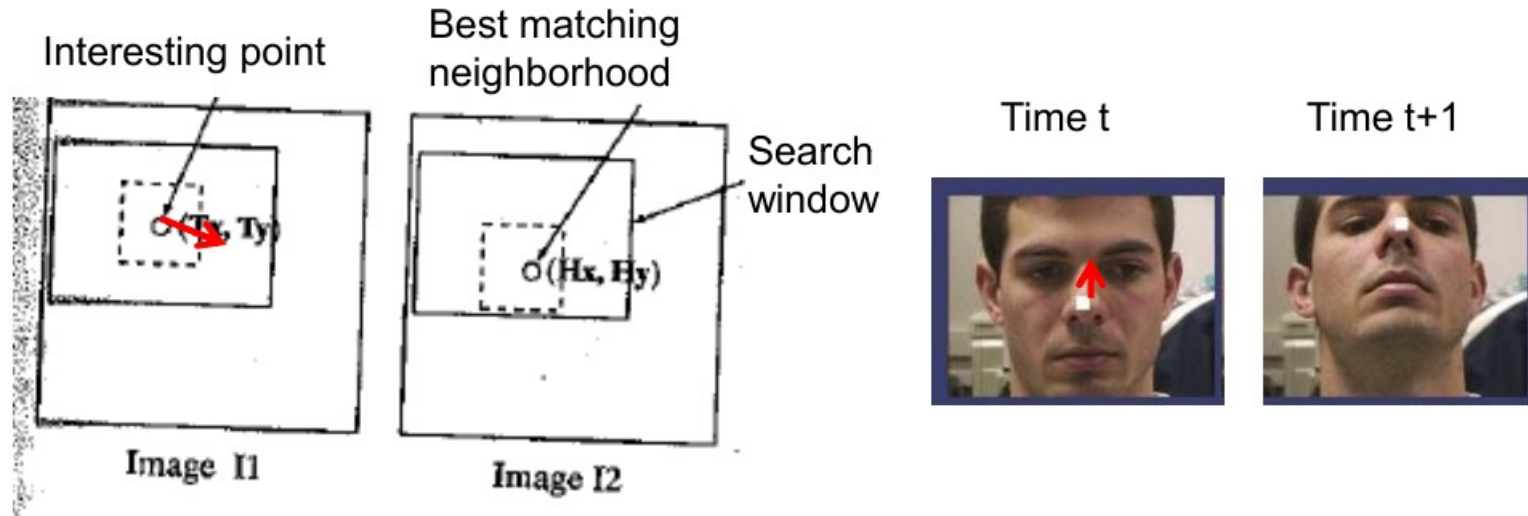
Motion Estimation Techniques



UNIVERSITÀ
CA' FOSCARI
VENEZIA

- Direct methods
 - Directly recover image motion at each pixel from spatio-temporal image brightness variations
 - Dense motion fields, fields but sensitive to appearance variations
 - Suitable for video and when image motion is small
- Feature-based methods
 - Extract visual features (corners, textured areas) and track them over multiple frames
 - Sparse motion fields, but more robust tracking
 - Suitable when image motion is large (10s of pixels)

Feature-based matching for motion



- Search window is centered at the point where we last saw the feature, in image I1
- Best match = position where we have the highest normalized cross-correlation value
- Where should the search window be placed?
 - Near match at previous frame
 - More generally, taking into account the expected **dynamics** of the object



Detection vs. tracking



UNIVERSITÀ
CA' FOSCARI
VENEZIA



$t=1$



$t=2$

...



$t=20$

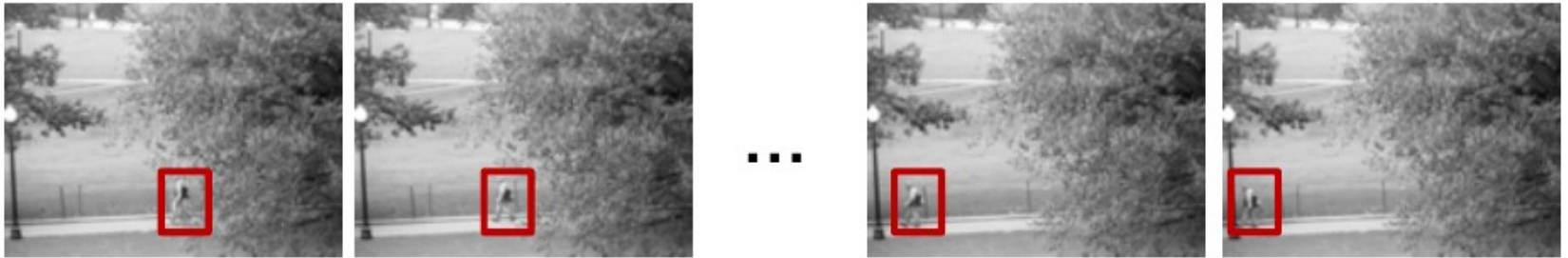


$t=21$

Detection vs. tracking



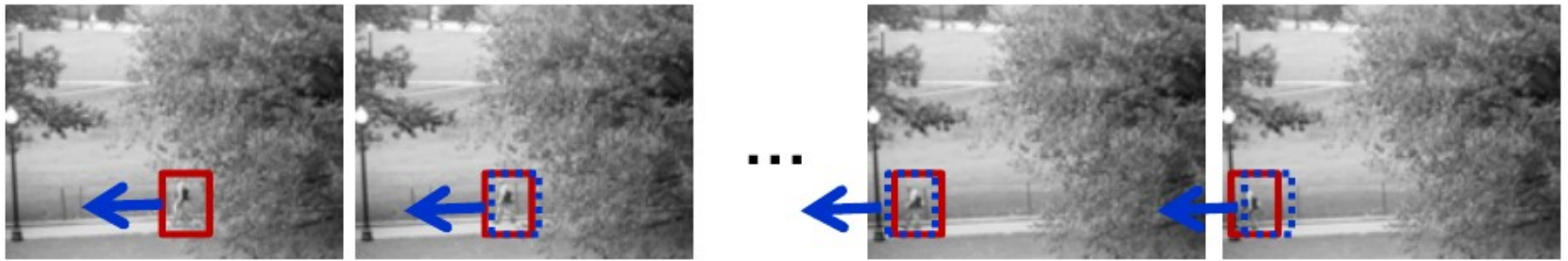
UNIVERSITÀ
CA' FOSCARI
VENEZIA



Detection: We detect the object independently in each frame and can record its position over time, e.g., based on blob's centroid or detection window coordinates



UNIVERSITÀ
CA' FOSCARI
VENEZIA



Tracking with dynamics: We use image measurements to estimate position of object, but also incorporate position predicted by dynamics, i.e., our expectation of object's motion pattern.

Tracking with dynamics

- Use model of expected motion to predict where objects will occur in next frame, even before seeing the image.
- **Intent:**
 - Do less work looking for the object, restrict the search.
 - Get improved estimates since measurement noise is tempered by smoothness, dynamics priors.
- **Assumption: continuous motion patterns:**
 - Camera is not moving instantly to new viewpoint
 - Objects do not disappear and reappear in different places in the scene
 - Gradual change in pose between camera and scene



Tracking as inference

- The hidden state consists of the true parameters we care about, denoted X .
- The measurement is our noisy observation that results from the underlying state, denoted Y .
- At each time step, state changes (from X_{t-1} to X_t) and we get a new observation Y_t
- Our goal: recover most likely state X_t given
 - All observations seen so far.
 - Knowledge about dynamics of state transitions.



Independence Assumptions

- Only immediate past state influences current state

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

- Measurements at time i only depend on the current state

$$P(\mathbf{Y}_i, \mathbf{Y}_j, \dots, \mathbf{Y}_k | \mathbf{X}_i) = P(\mathbf{Y}_i | \mathbf{X}_i) P(\mathbf{Y}_j, \dots, \mathbf{Y}_k | \mathbf{X}_i)$$



Tracking via deformable contours



UNIVERSITÀ
CA' FOSCARI
VENEZIA

- Use final contour/model extracted at frame t as an initial solution for frame $t+1$
- Evolve initial contour to fit exact object boundary at frame $t+1$
- Repeat, initializing with most recent frame.

