



3D Reconstruction with Computer Vision  
Meeting 21: Video and Tracking Redux  
CS 378 Fall 2014, UT Austin, Bryan Klingner, 6 November

Grow 0

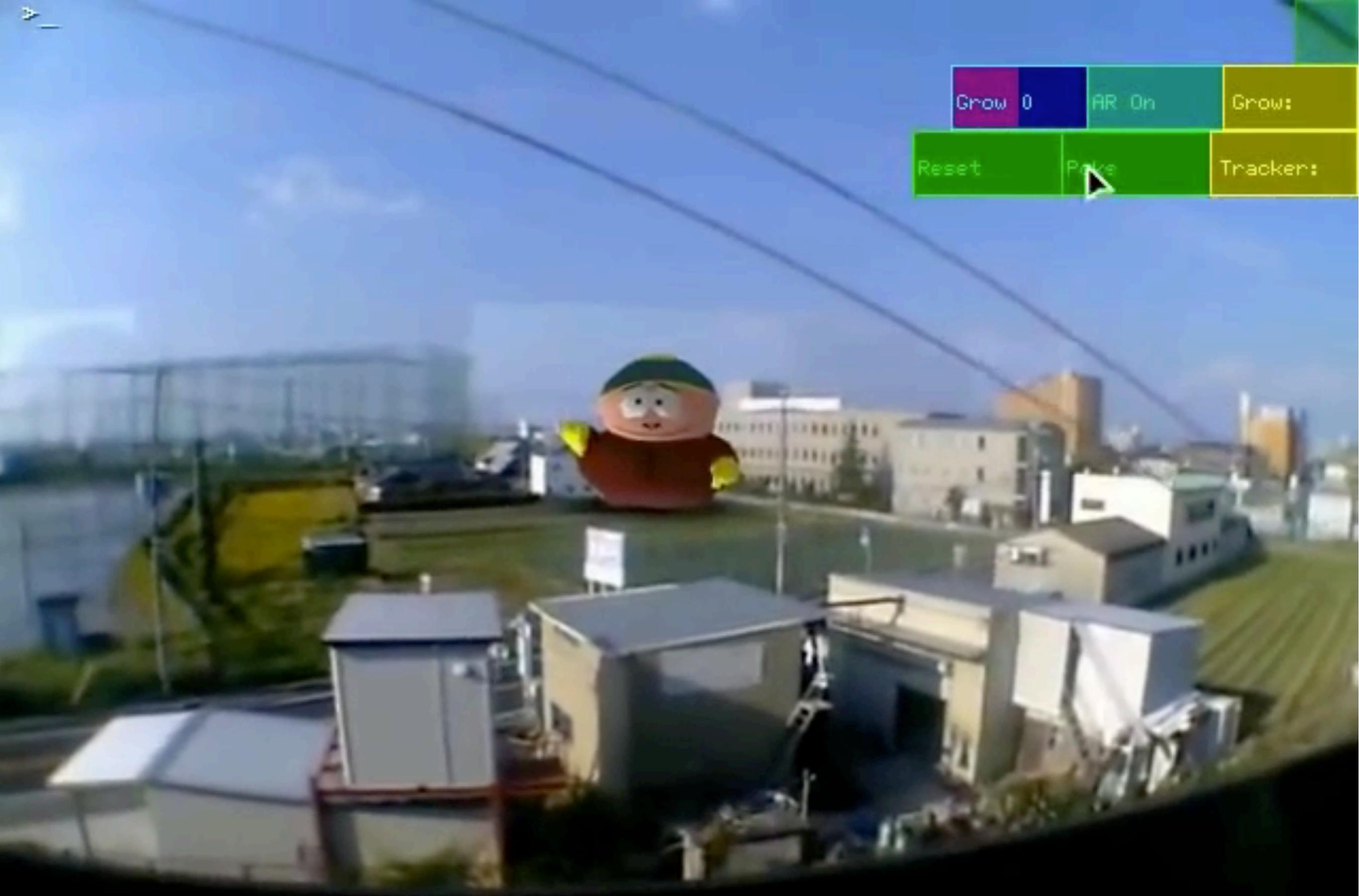
AR On

Grow:

Reset

Poke

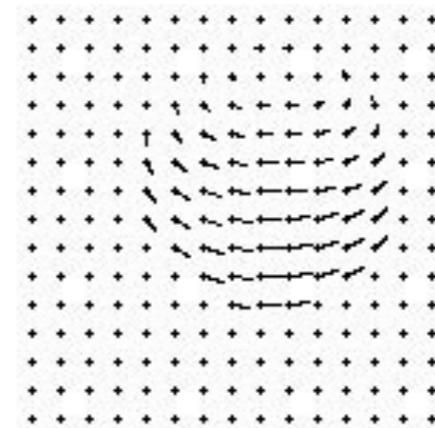
Tracker:



On the train to Kyoto

# Motion and tracking

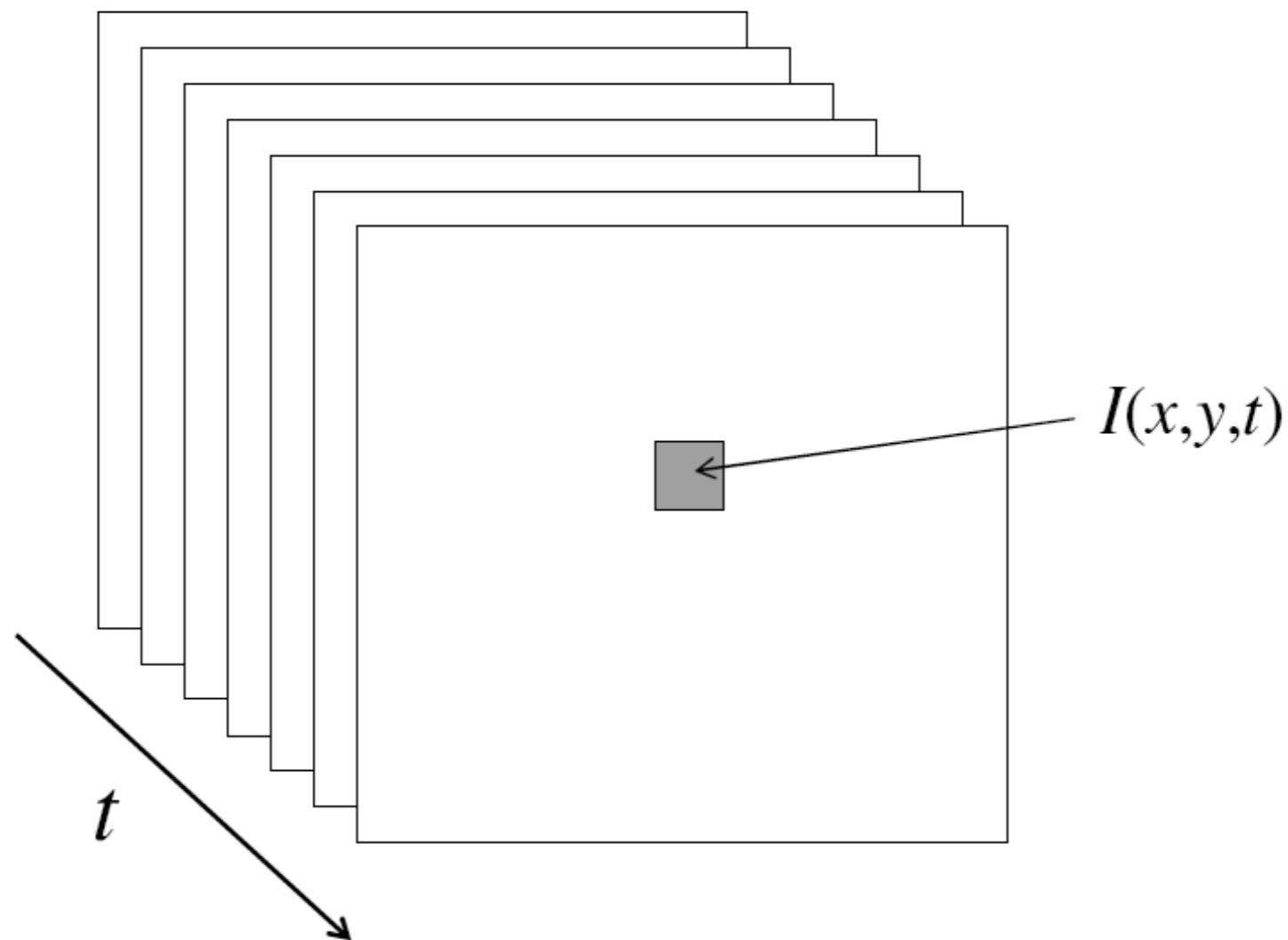
Tracking objects, video analysis, low level motion



Tomas Izo

# Video

- A video is a sequence of frames captured over time
- Now our image data is a function of space ( $x, y$ ) and time ( $t$ )

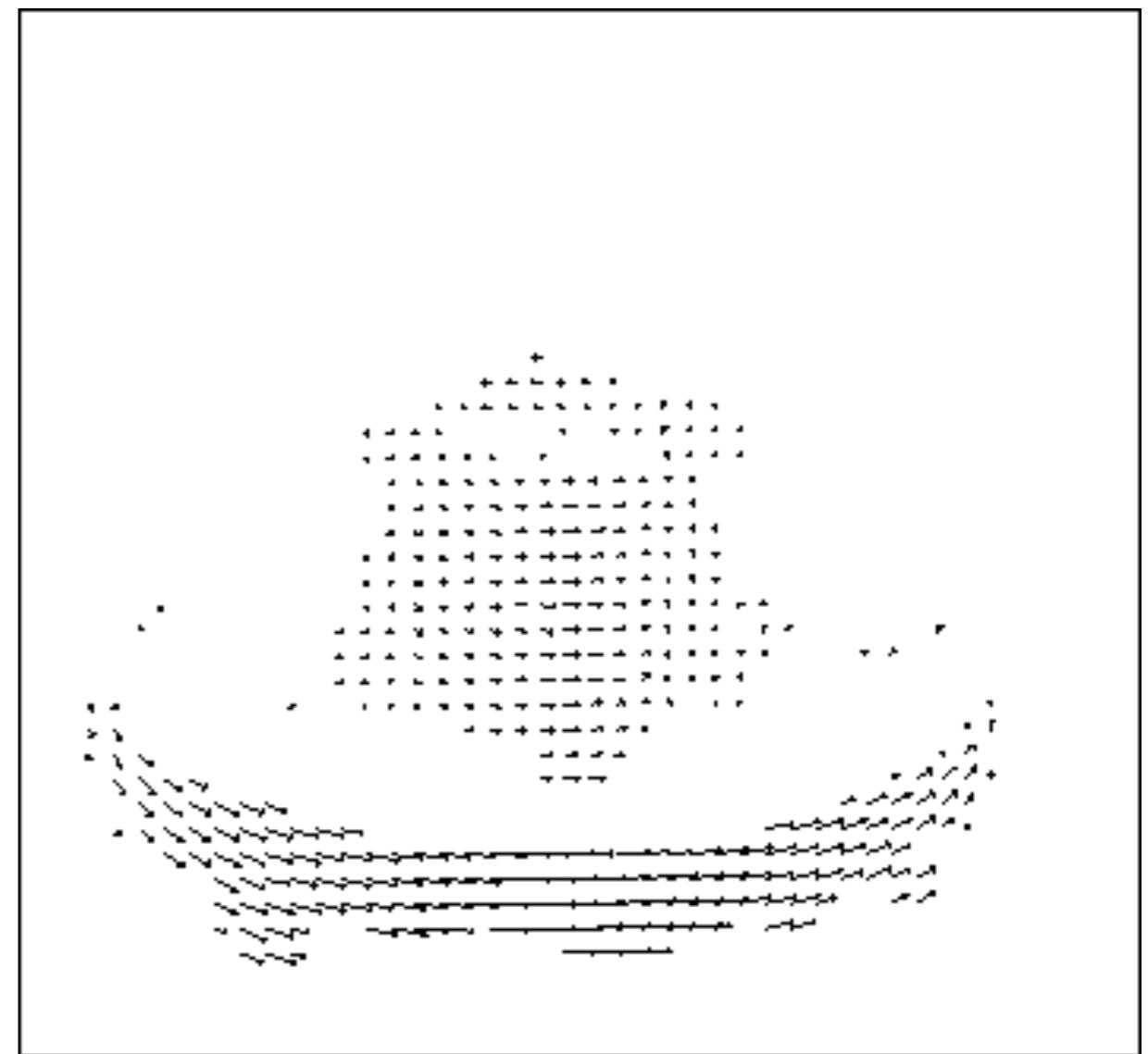
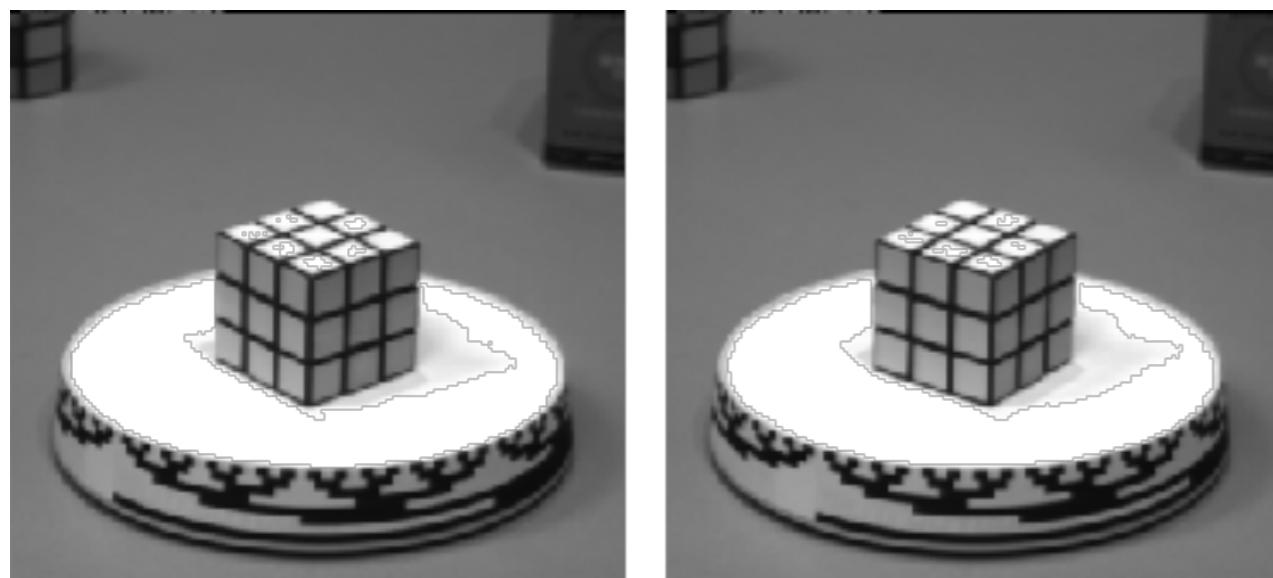


# Uses of motion

- Estimating 3D structure
- Segmenting objects based on motion cues
- Learning dynamical models
- Recognizing events and activities
- Improving video quality (motion stabilization)

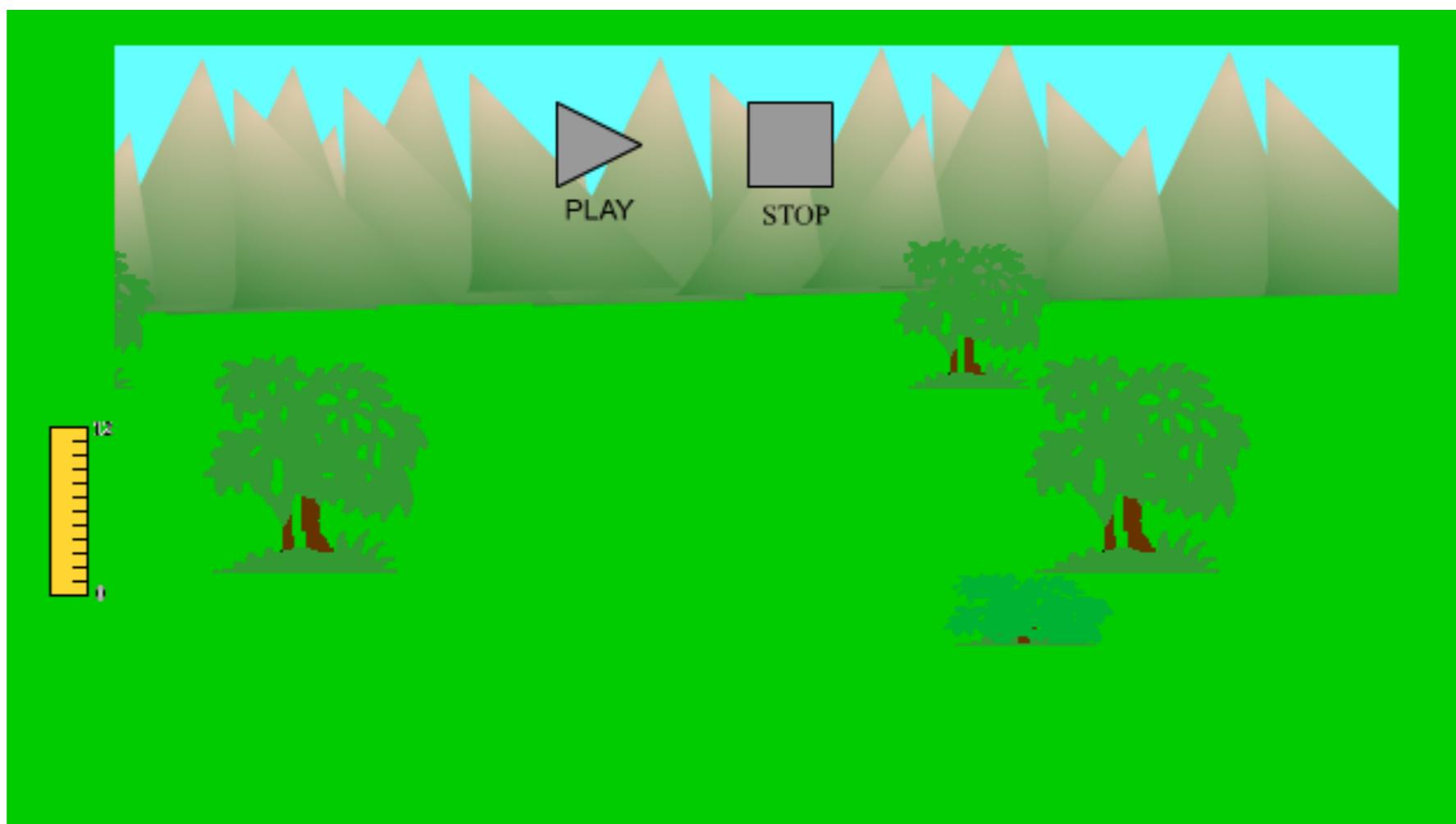
# Motion field

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

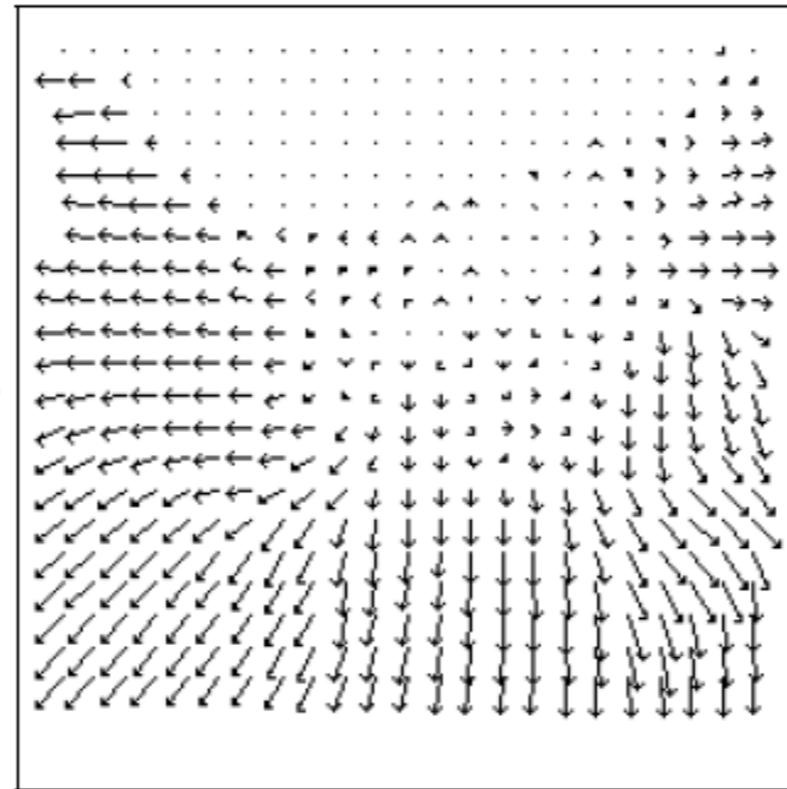


# Motion parallax

[http://psych.hanover.edu/KRANTZ/MotionParallax/  
MotionParallax.html](http://psych.hanover.edu/KRANTZ/MotionParallax/MotionParallax.html)



# Motion field + camera motion

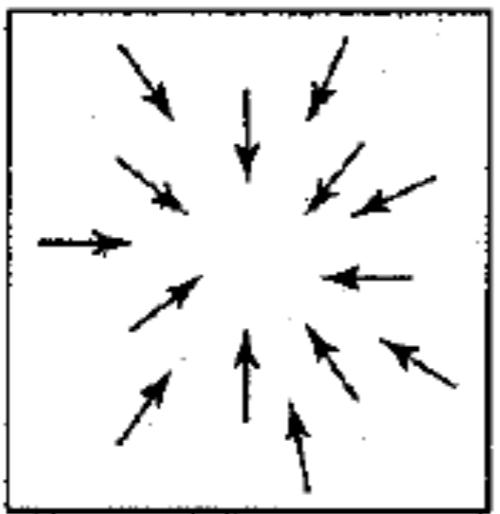


Length of flow vectors inversely proportional to depth  $Z$  of 3d point

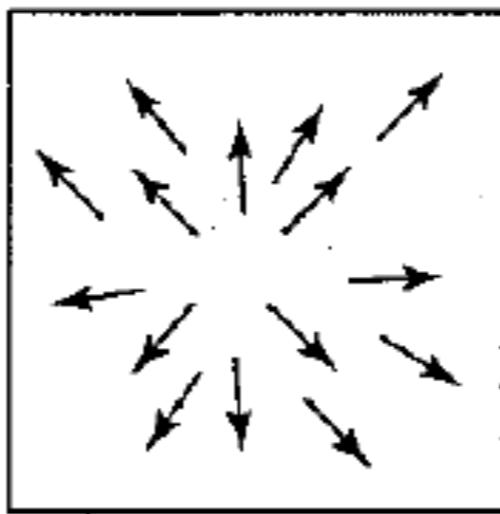
Figure 1.2: Two images taken from a helicopter flying through a canyon and the computed optical flow field.

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

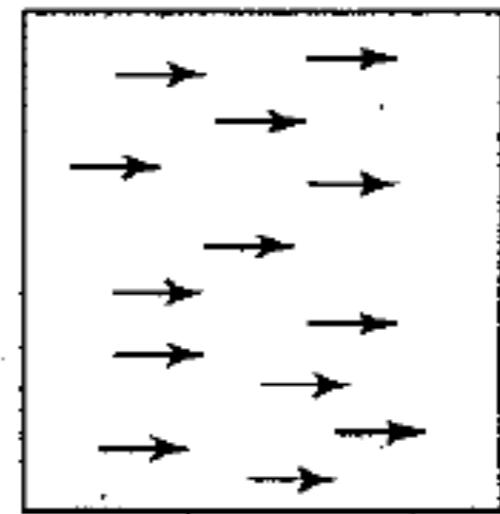
# Motion field + camera motion



**Zoom out**



**Zoom in**



**Pan right to left**

Main App... Move... Colors Video

- Camera Input Options

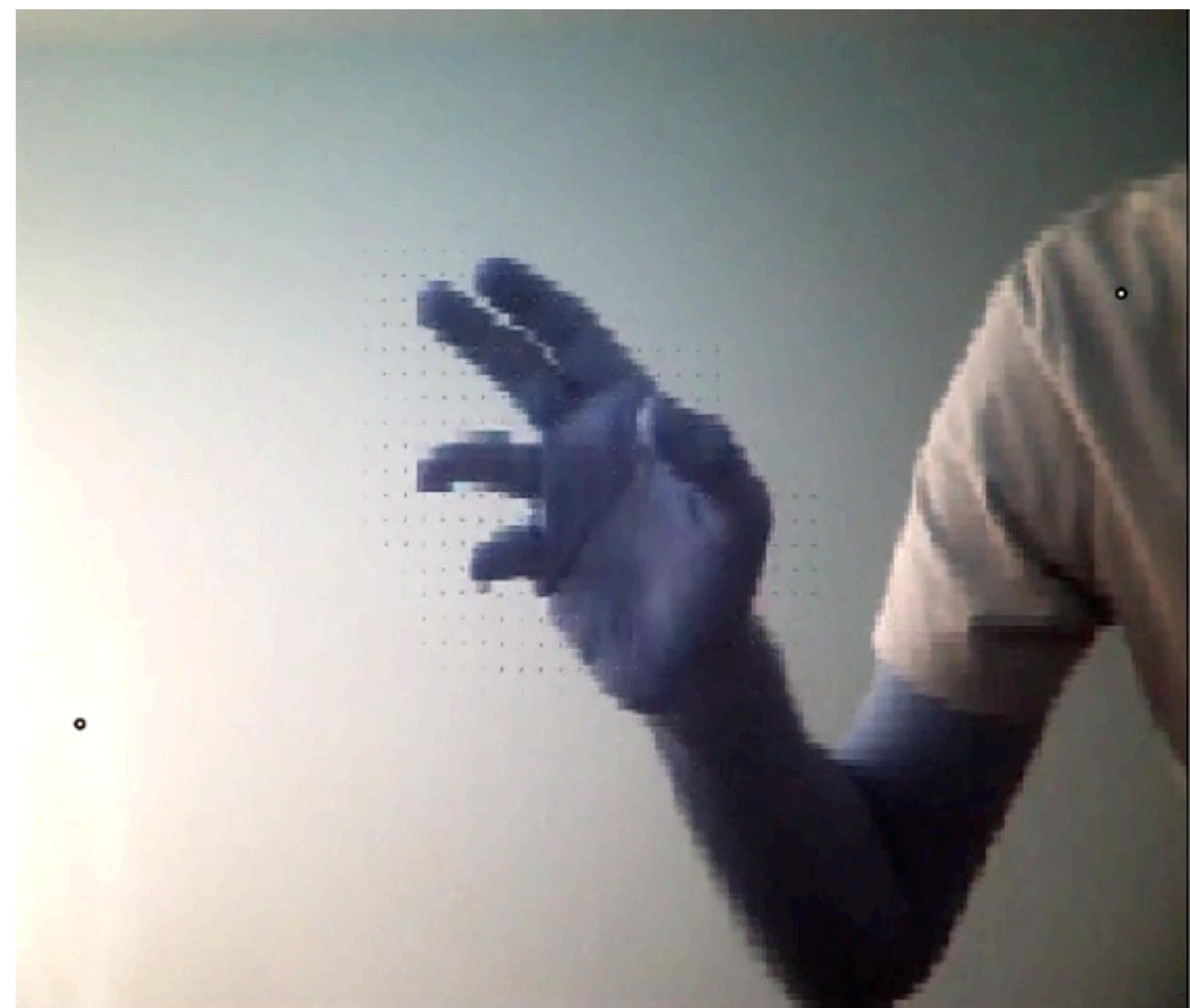
Switch camera input  
render motion fx  
render video  
flipH

fx color



State

grid  
particle speed  
threads  
particle force  
particle contour  
optical flow  
contour field



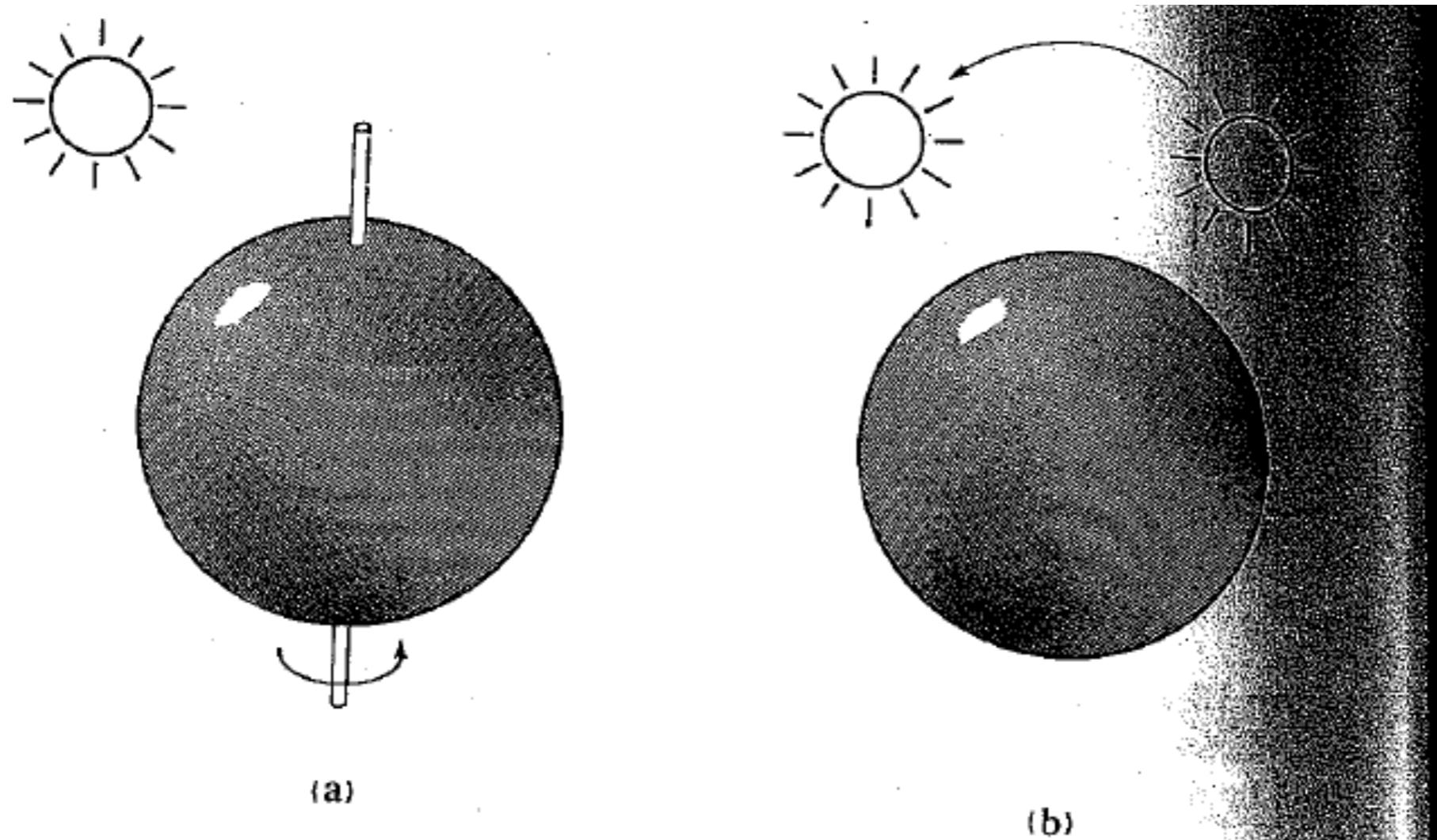
# Motion estimation techniques

- Direct methods
  - Directly recover image motion at each pixel from spatio-temporal image brightness variations
  - Dense motion 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

- Definition: 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

# Apparent motion != motion field

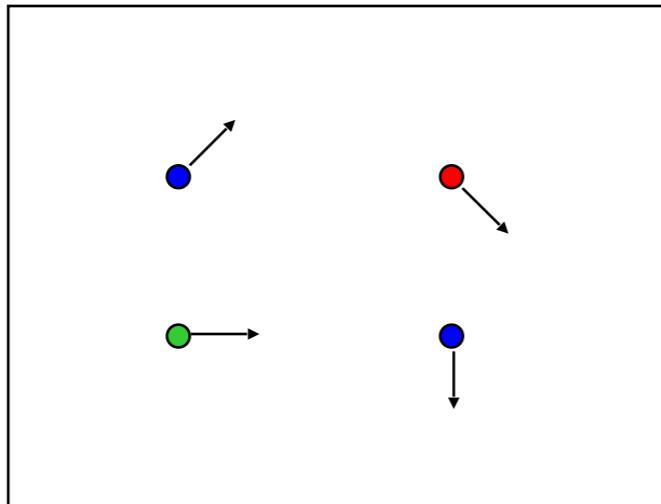


**Figure 12-2.** The optical flow is not always equal to the motion field. In (a) a smooth sphere is rotating under constant illumination—the image does not change, yet the motion field is nonzero. In (b) a fixed sphere is illuminated by a moving source—the shading in the image changes, yet the motion field is zero.

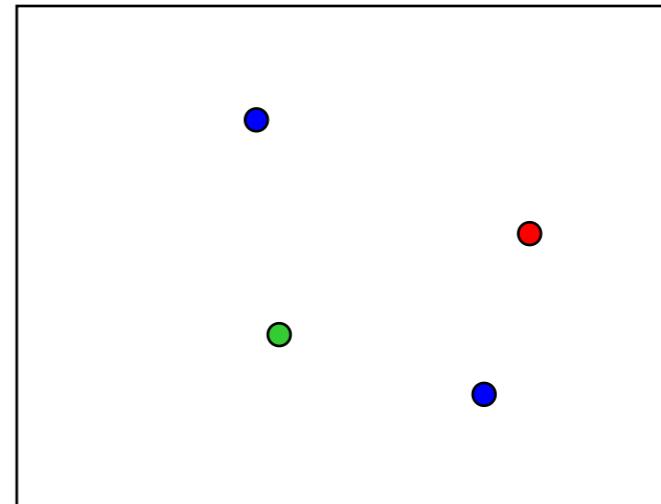
Figure from Horn book

# Problem definition: 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  $H$ , 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

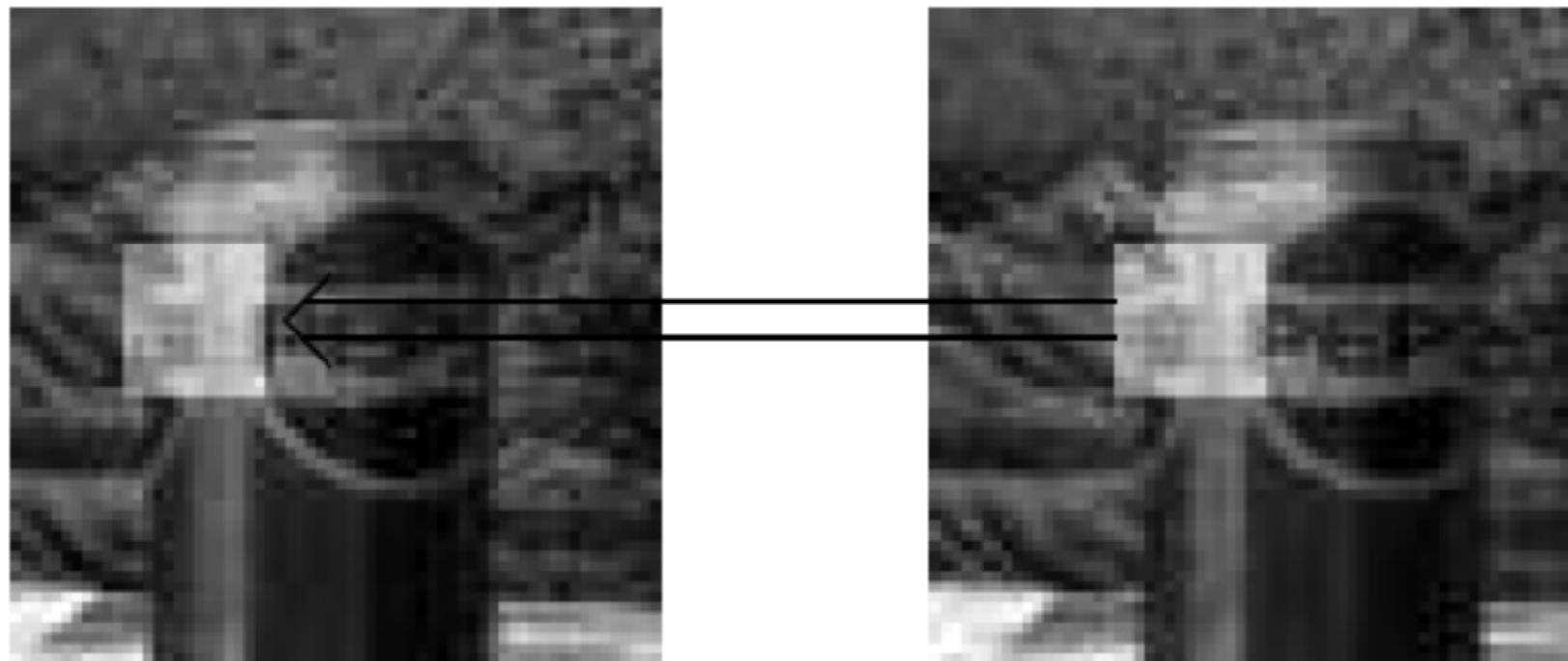
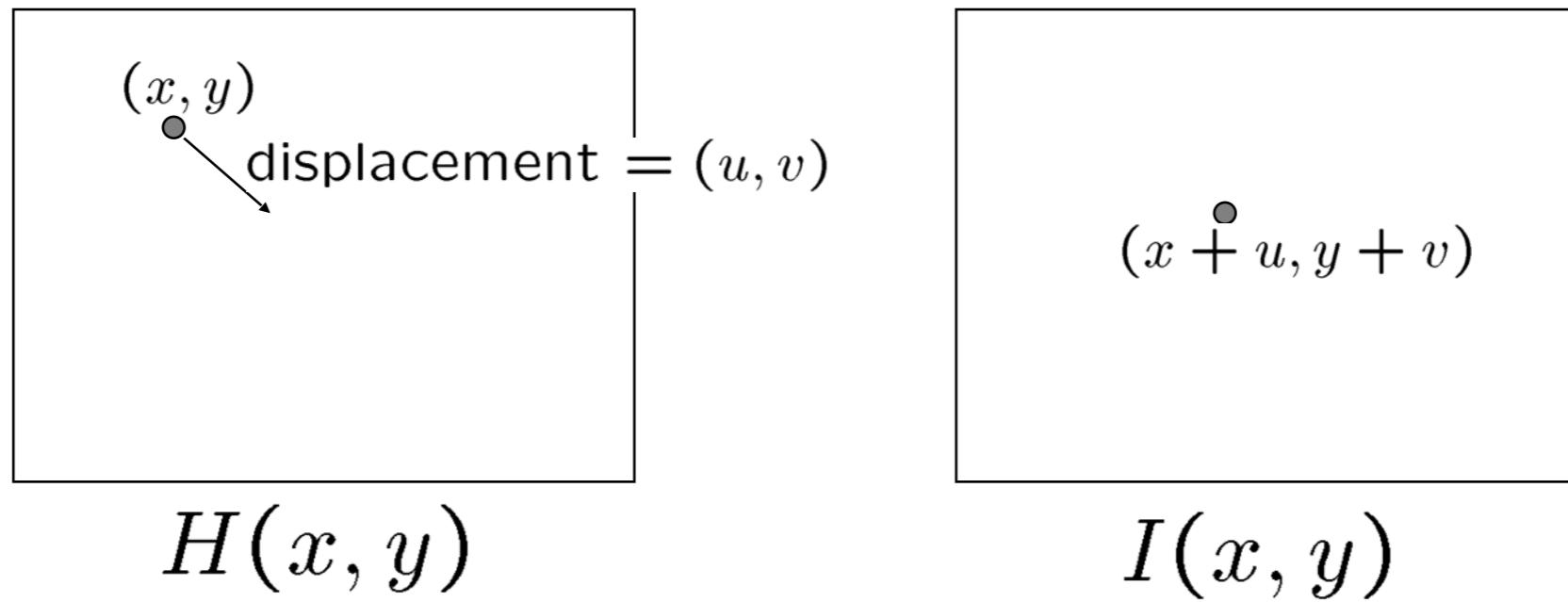


Figure 1.5: Data conservation assumption. The highlighted region in the right image looks roughly the same as the region in the left image, despite the fact that it has moved.

# Optical flow constraints (grayscale images)

---



Let's look at these constraints more closely

- brightness constancy: Q: what's the equation?

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

- small motion:

$$\begin{aligned} I(x+u, y+v) &= I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \text{higher order terms} \\ &\approx I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v \end{aligned}$$

# Optical flow equation

---

Combining these two equations

$$0 = I(x + u, y + v) - H(x, y)$$

shorthand:  $I_x = \frac{\partial I}{\partial x}$

$$\approx I(x, y) + I_x u + I_y v - H(x, y)$$

$$\approx (I(x, y) - H(x, y)) + I_x u + I_y v$$

$$\approx I_t + I_x u + I_y v$$

$$\approx I_t + \nabla I \cdot [u \ v]$$

# Optical flow equation

---

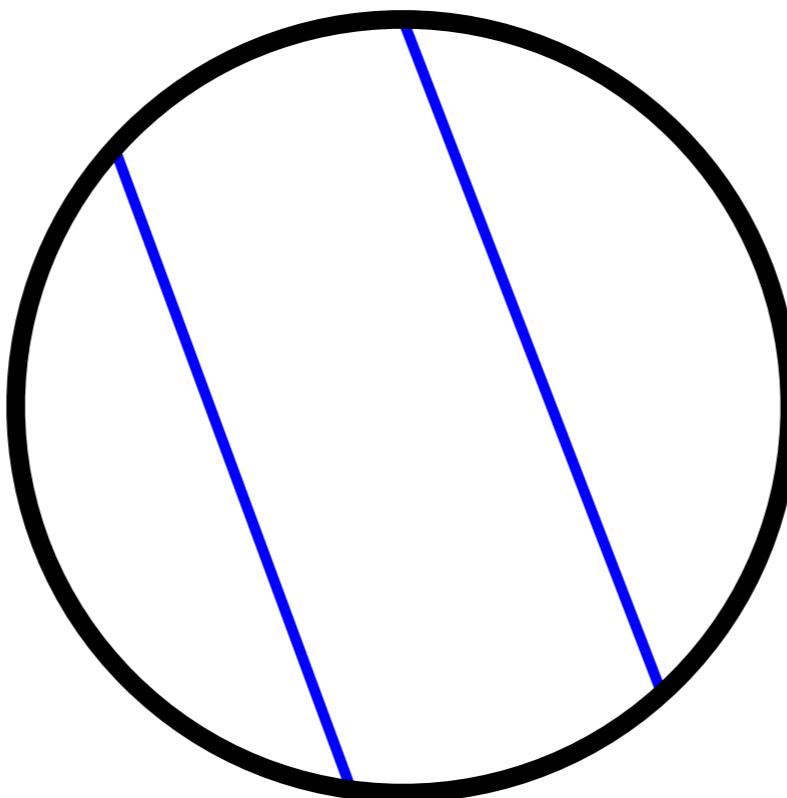
$$0 = I_t + \nabla I \cdot [u \ v]$$

Q: how many unknowns and equations per pixel?

Intuitively, what does this ambiguity mean?

# The aperture problem

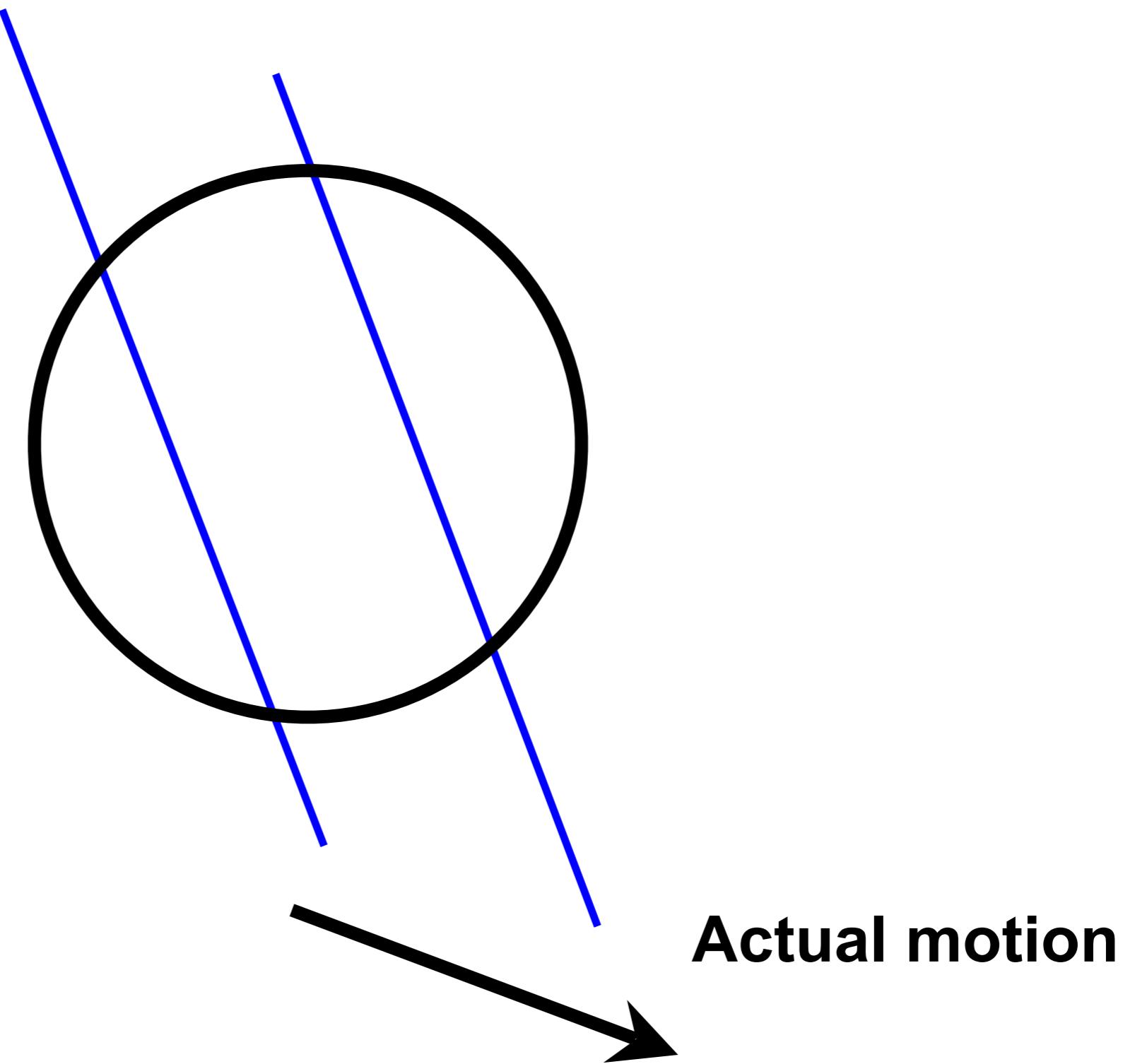
---



**Perceived motion**

# The aperture problem

---



# The barber pole illusion



[http://en.wikipedia.org/wiki/Barberpole\\_illusion](http://en.wikipedia.org/wiki/Barberpole_illusion)

# The barber pole illusion



[http://en.wikipedia.org/wiki/Barberpole\\_illusion](http://en.wikipedia.org/wiki/Barberpole_illusion)

# Solving the aperture problem (grayscale image)

- How to get more equations for a pixel?
- **Spatial coherence constraint:** pretend the pixel's neighbors have the same  $(u,v)$

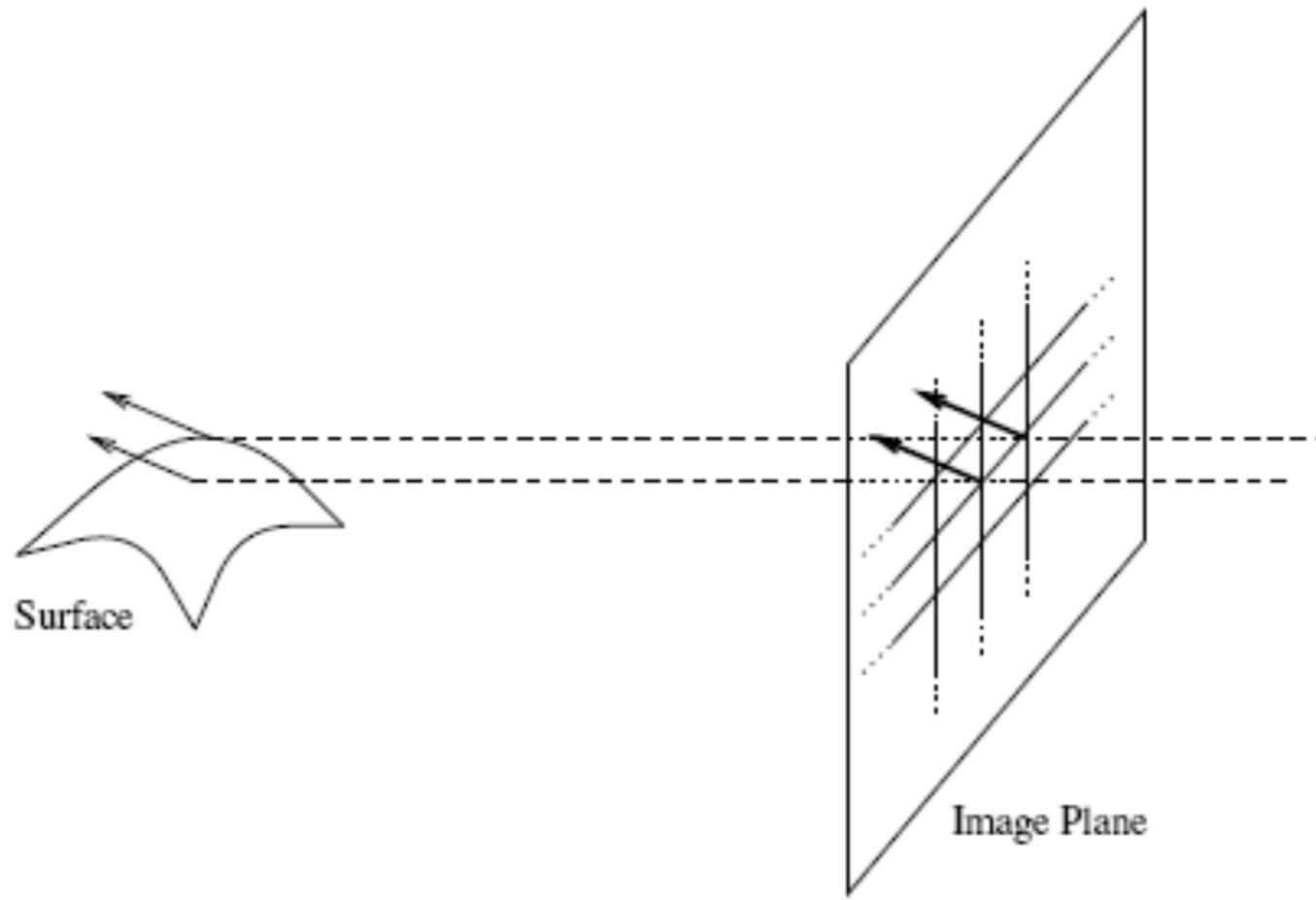


Figure 1.7: Spatial coherence assumption. Neighboring points in the image are assumed to belong to the same surface in the scene.

# Solving the aperture problem (grayscale image)

- 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  $5 \times 5$  window, that gives us 25 equations per pixel

$$0 = I_t(\mathbf{p}_i) + \nabla I(\mathbf{p}_i) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix}$$

$$\begin{array}{ccc} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{array}$$

# Solving the aperture problem

---

Prob: we have more equations than unknowns

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

Solution: solve least squares problem

- minimum least squares solution given by solution (in  $d$ ) of:

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

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

- The summations are over all pixels in the  $K \times K$  window
- This technique was first proposed by Lucas & Kanade (1981)

# Conditions for solvability

---

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

When is this solvable?

- $A^T A$  should be invertible
- $A^T A$  should not be too small
  - eigenvalues  $\lambda_1$  and  $\lambda_2$  of  $A^T A$  should not be too small
- $A^T A$  should be well-conditioned
  - $\lambda_1 / \lambda_2$  should not be too large ( $\lambda_1$  = larger eigenvalue)

# Edge



- gradients very large or very small
- large  $\lambda_1$ , small  $\lambda_2$

# Low-texture region



- gradients have small magnitude
- small  $\lambda_1$ , small  $\lambda_2$

# High-texture region



- gradients are different, large magnitudes
- large  $\lambda_1$ , large  $\lambda_2$

Main App... Move... Colors Video

- Camera Input Options

Switch camera input  
render motion fx  
render video  
flipH

fx color

H -21474  
S 0  
V 247  
A 112

State

grid  
particle speed  
threads  
particle force  
particle contour  
optical flow  
contour field



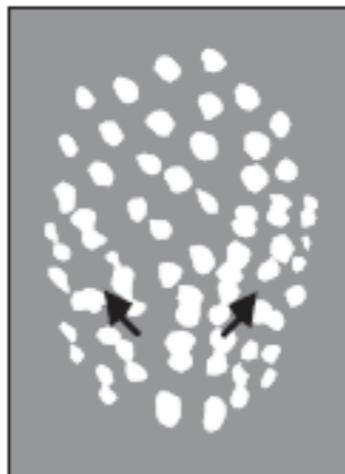
# Using optical flow: recognizing facial expressions



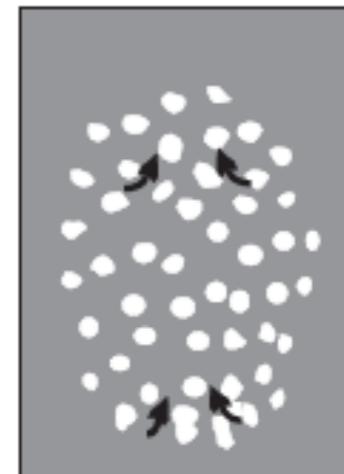
Disgust



Sadness



Happiness



Sadness



happiness



fear



Surprise



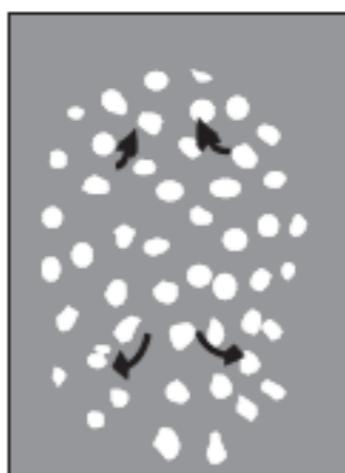
Anger



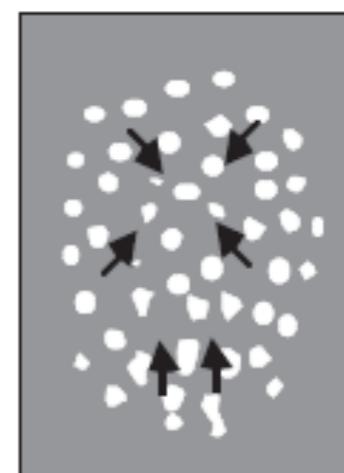
Anger



Surprise



Fear



Disgust

**Recognizing Human Facial Expression  
(1994)**

by Yaser Yacoob, Larry S. Davis

# Applying optical flow: video stabilization



# Applying optical flow: video stabilization



# Applying optical flow: video stabilization



# Motion vs. Stereo: Similarities

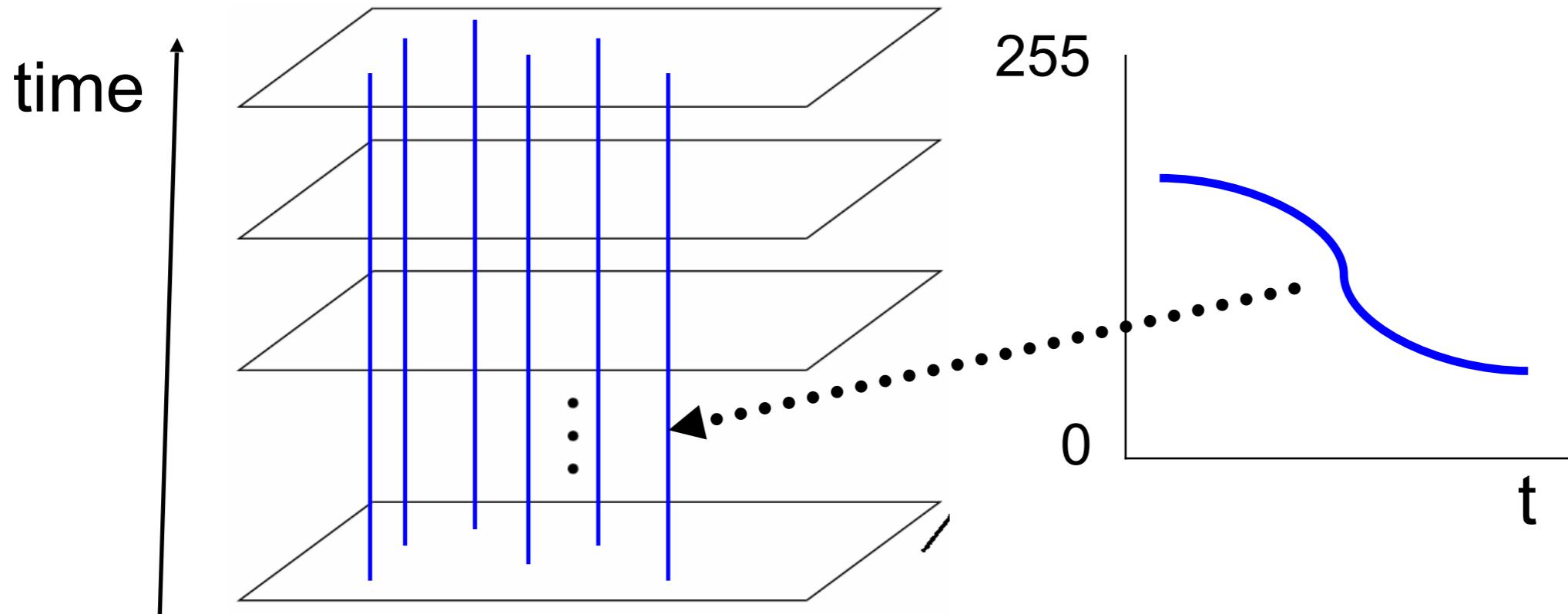
- Both involve solving
  - Correspondence: disparities, motion vectors
  - Reconstruction

# Motion vs. Stereo: Differences

- **Motion:**
  - Uses velocity: consecutive frames must be close to get good approximate time derivative
  - 3d movement between camera and scene not necessarily single 3d rigid transformation
- **Whereas with stereo:**
  - Could have any disparity value
  - View pair separated by a single 3d transformation

# Video as an “Image Stack”

---



Can look at video data as a spatio-temporal volume

- If camera is stationary, each line through time corresponds to a single ray in space

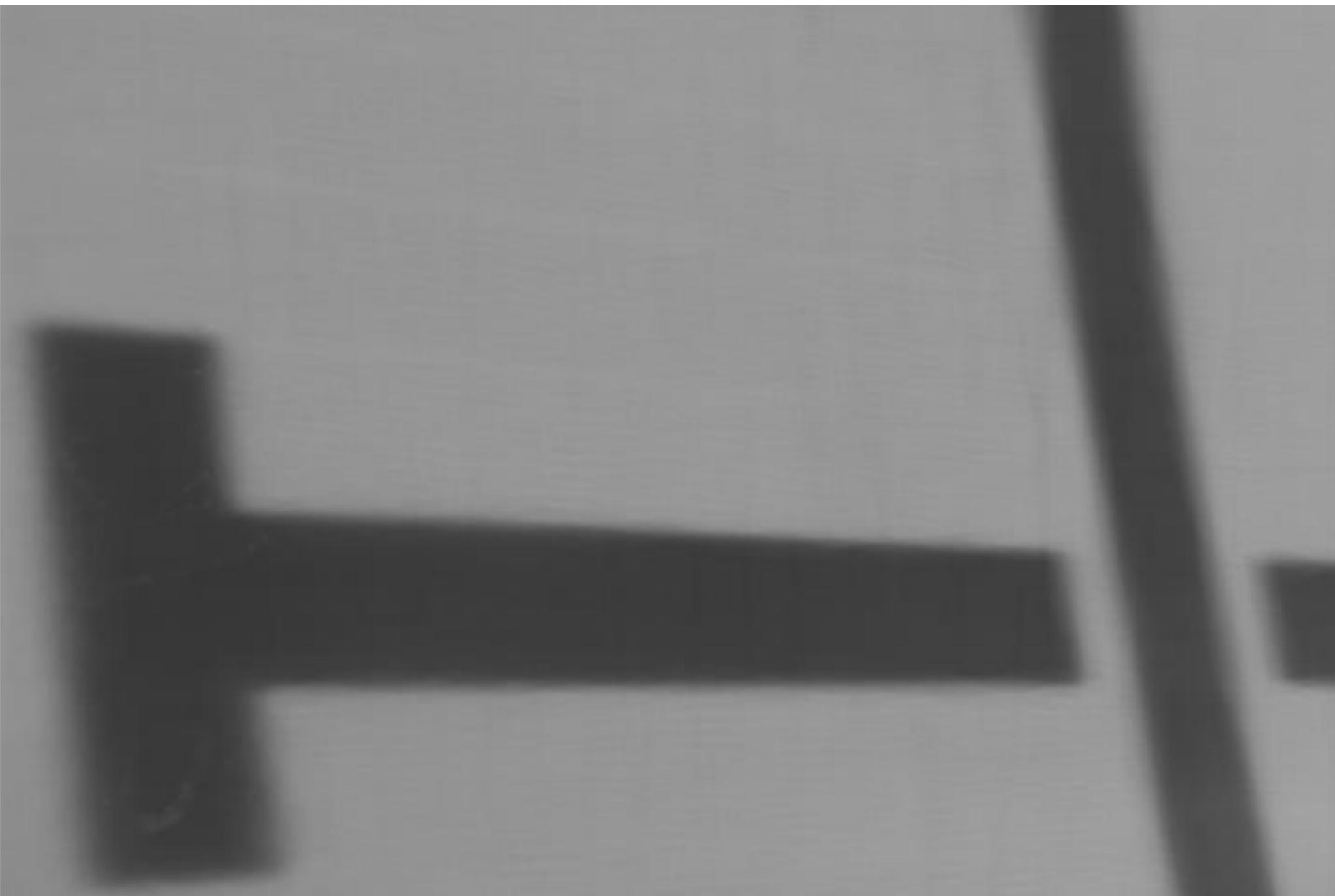
# Input Video

---



# Average Image

---



# Background Subtraction

- ▶ Given an image (mostly likely to be a video frame), we want to identify the **foreground objects** in that image!



## Motivation

- ▶ In most cases, objects are of interest, not the scene.
- ▶ Makes our life easier: less processing costs, and less room for error.

# Background subtraction

- Simple techniques can do ok with static camera
- ...But hard to do perfectly
- Widely used:
  - Traffic monitoring (counting vehicles, detecting & tracking vehicles, pedestrians),
  - Human action recognition (run, walk, jump, squat),
  - Human-computer interaction
  - Object tracking

# Simple Approach

Image at time  $t$ :

$$I(x, y, t)$$



Background at time  $t$ :

$$B(x, y, t)$$



-

$$| > Th$$

1. Estimate the background for time  $t$ .
2. Subtract the estimated background from the input frame.
3. Apply a threshold,  $Th$ , to the absolute difference to get the **foreground mask**.

# Frame Differencing

- Background is estimated to be the previous frame.  
Background subtraction equation then becomes:

$$B(x, y, t) = I(x, y, t - 1)$$



$$|I(x, y, t) - I(x, y, t - 1)| > Th$$

- Depending on the object structure, speed, frame rate and global threshold, this approach may or may **not** be useful (usually **not**).



-



$$| > Th$$

# Frame Differencing

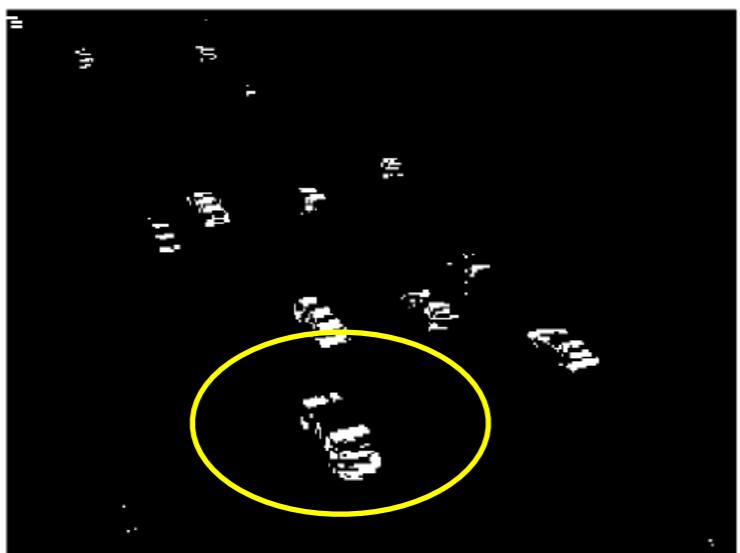
$Th = 25$



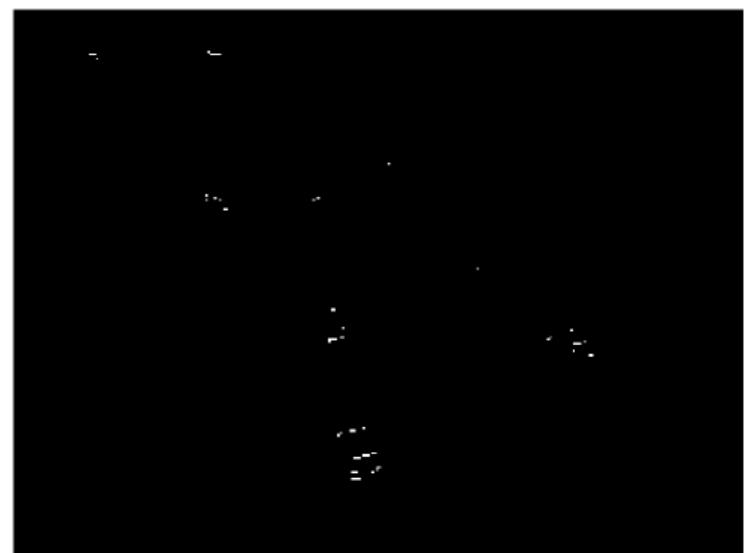
$Th = 50$



$Th = 100$



$Th = 200$



# Mean Filter

- ▶ In this case the background is the mean of the previous  $n$  frames:

$$B(x, y, t) = \frac{1}{n} \sum_{i=0}^{n-1} I(x, y, t - i)$$
$$\downarrow$$
$$|I(x, y, t) - \frac{1}{n} \sum_{i=0}^{n-1} I(x, y, t - i)| > Th$$

- ▶ For  $n = 10$ :

Estimated Background



Foreground Mask



# Median Filter

- ▶ Assuming that the background is more likely to appear in a scene, we can use the median of the previous  $n$  frames as the background model:

$$B(x, y, t) = \text{median}\{I(x, y, t - i)\}$$



$$|I(x, y, t) - \text{median}\{I(x, y, t - i)\}| > Th \text{ where } i \in \{0, \dots, n - 1\}.$$

- ▶ For  $n = 10$ :

Estimated Background



Foreground Mask



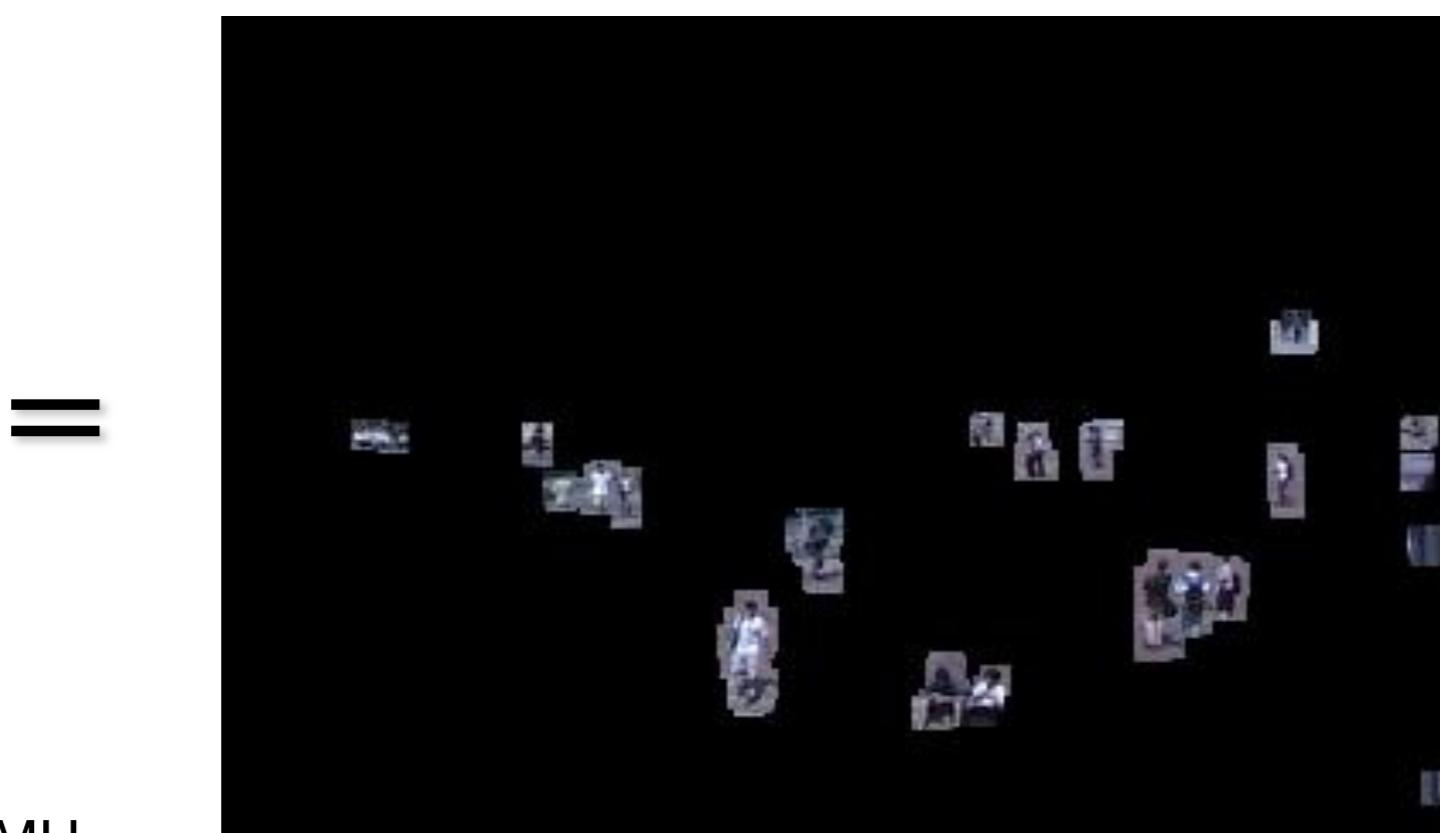
# Average/Median Image

---



# Background Subtraction

---



# Pros and cons

## Advantages:

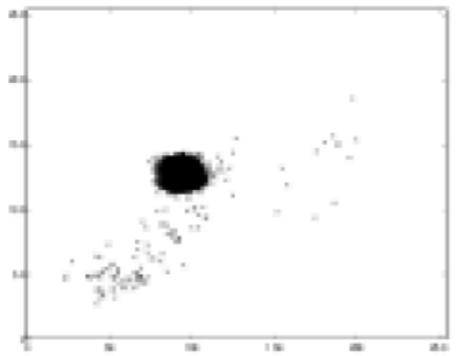
- Extremely easy to implement and use!
- All pretty fast.
- Corresponding background models need not be constant, they change over time.

## Disadvantages:

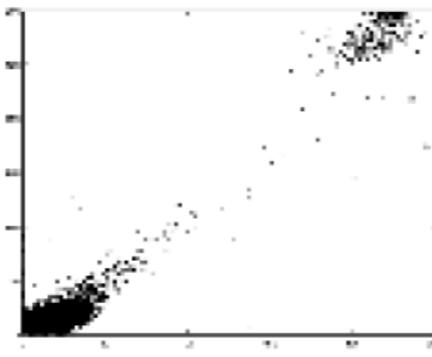
- Accuracy of frame differencing depends on object speed and frame rate
- Median background model: relatively high memory requirements.
- Setting global threshold Th...

*When will this basic approach fail?*

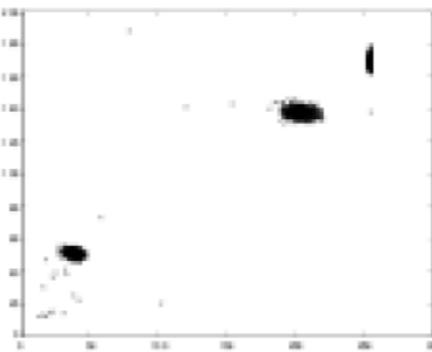
# Background mixture models



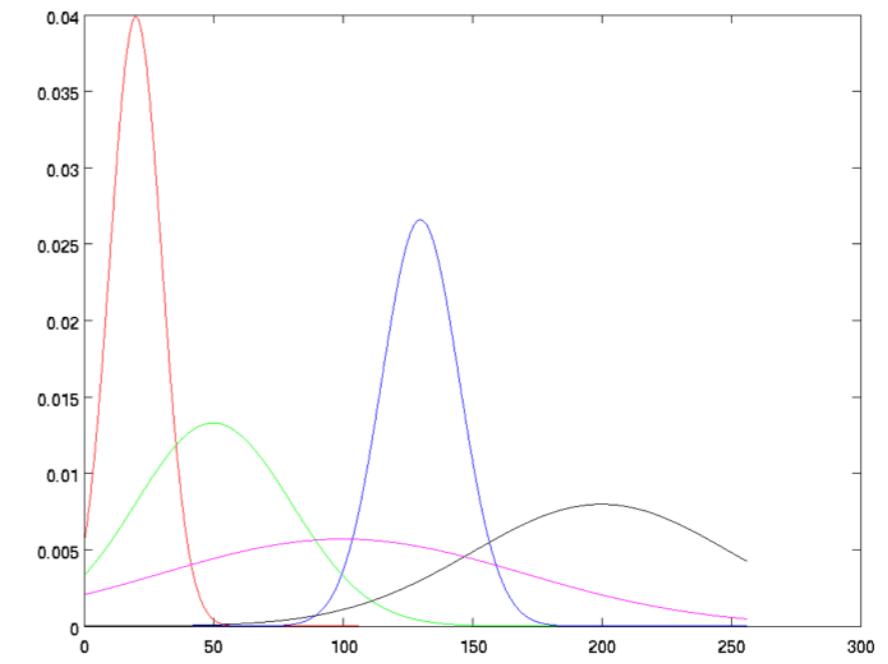
(a)



(b)



(c)



Idea: model each background pixel with a *mixture* of Gaussians; update its parameters over time.

# Background subtraction with depth



*How can we select foreground pixels based on depth information?*

# Human activity in video

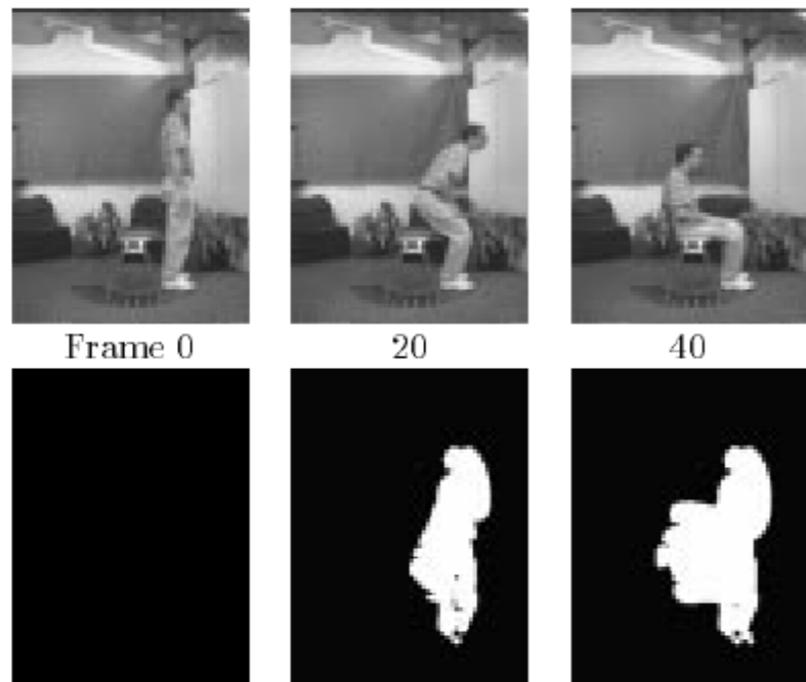
No universal terminology, but approximately:

- “**Actions**”: atomic motion patterns -- often gesture-like, single clear-cut trajectory, single nameable behavior (e.g., sit, wave arms)
- “**Activity**”: series or composition of actions (e.g., interactions between people)
- “**Event**”: combination of activities or actions (e.g., a football game, a traffic accident)

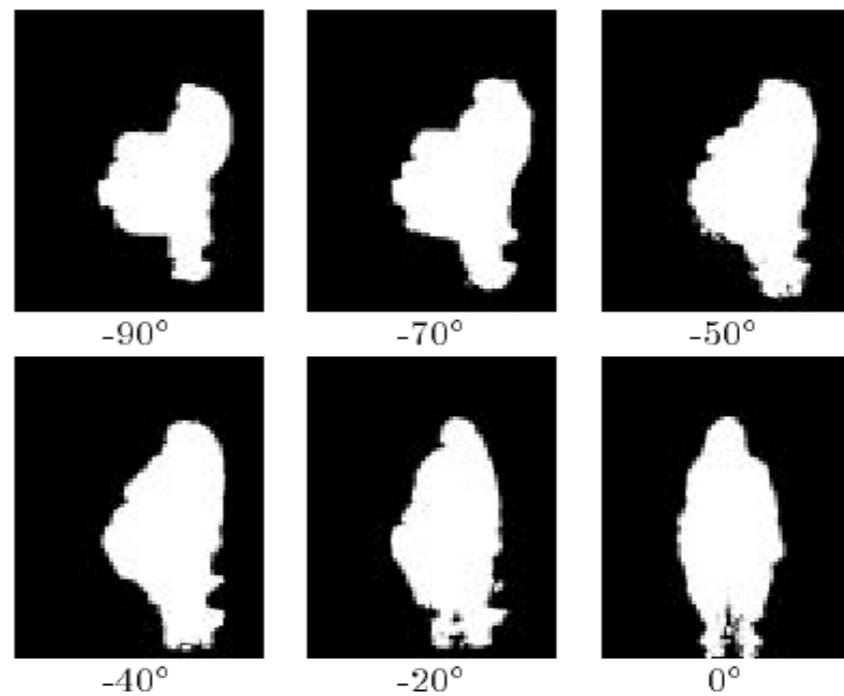
# Motion Energy Images

$$E_\tau(x, y, t) = \bigcup_{i=0}^{\tau-1} D(x, y, t - i)$$

$D(x, y, t)$ : Binary image sequence indicating motion locations

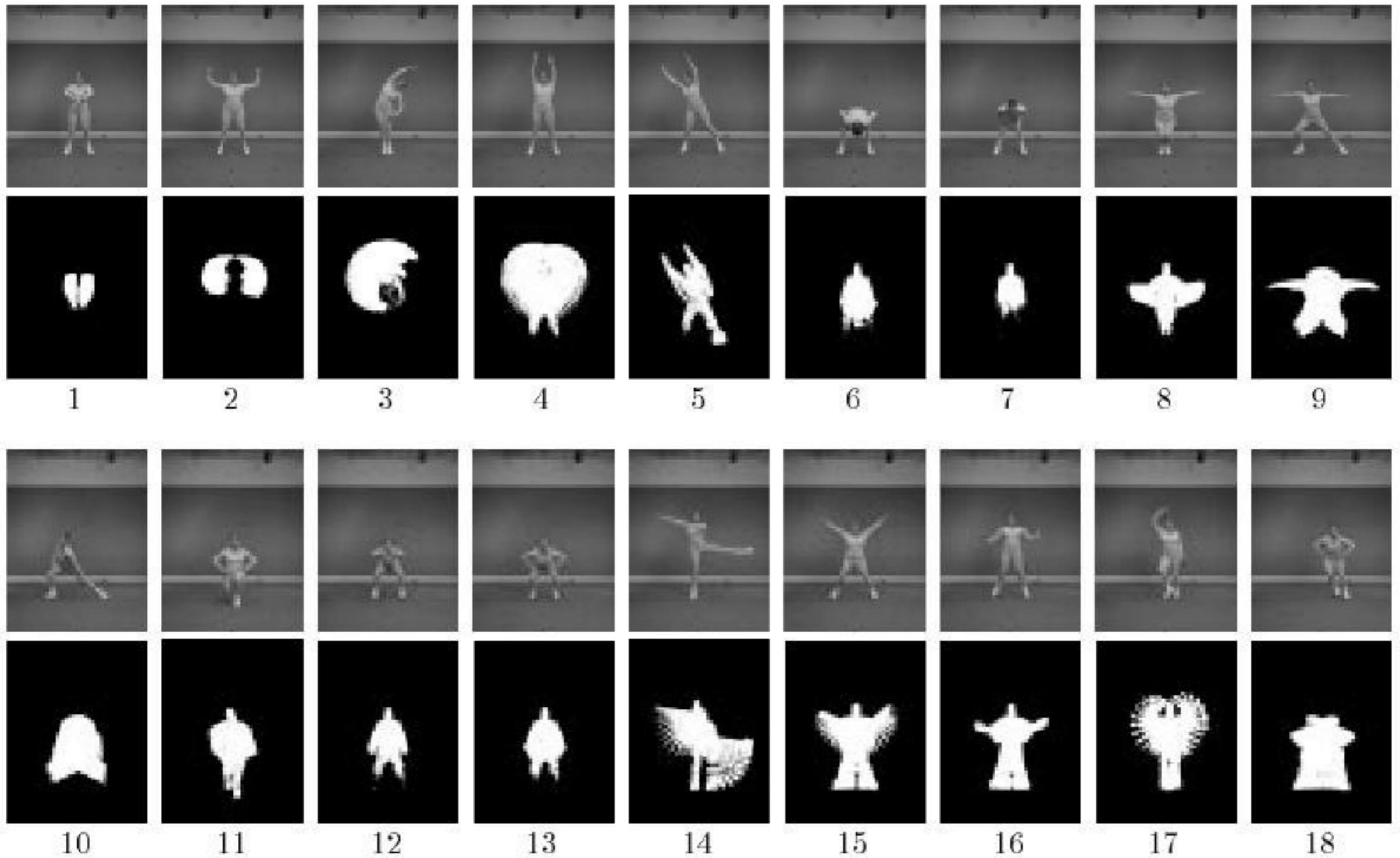


**Figure 2:** Example of someone sitting. Top row contains key frames; bottom row is cumulative motion images starting from Frame 0.



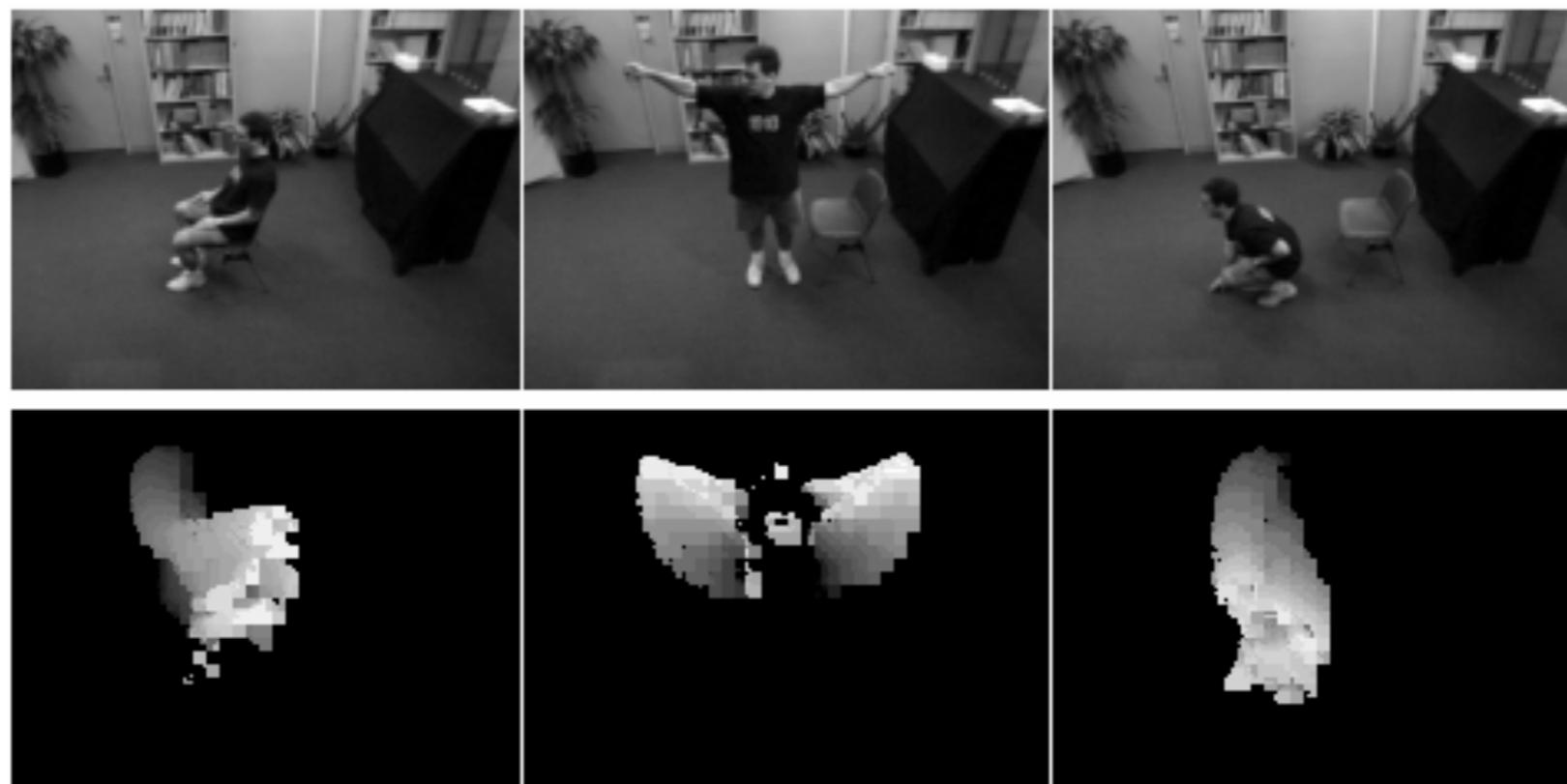
**Figure 3:** MEIs of sitting action over 90° viewing angle. The smooth change implies only a coarse sampling of viewing direction is necessary to recognize the action from all angles.

# Motion Energy Images



# Motion History Images

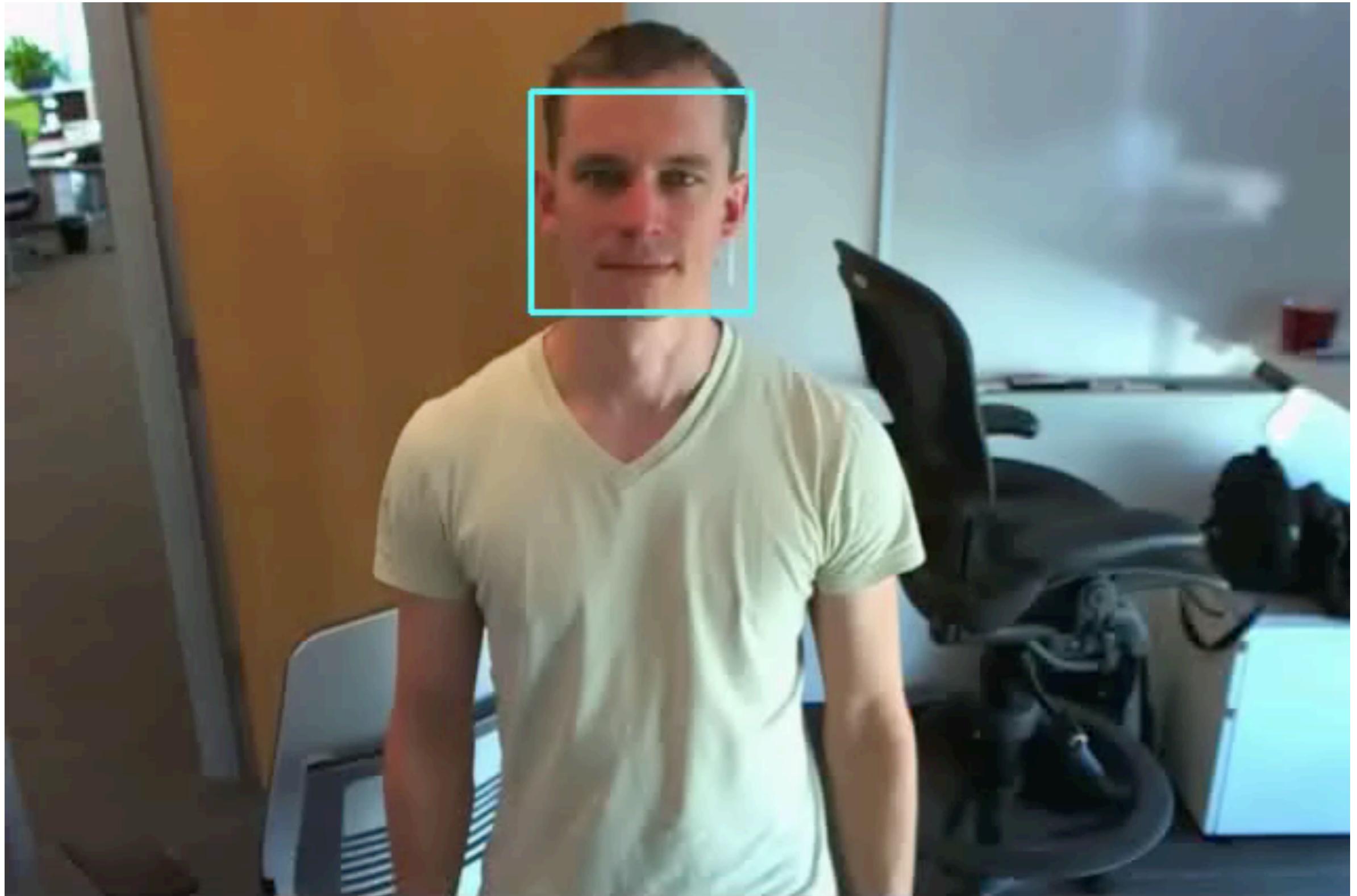
$$H_{\tau}(x, y, t) = \begin{cases} \tau & \text{if } D(x, y, t) = 1 \\ \max(0, H_{\tau}(x, y, t - 1) - 1) & \text{otherwise} \end{cases}$$



# Summary

- **Background subtraction:**
  - Essential low-level processing tool to segment moving objects from static camera's video
- **Action recognition:**
  - Increasing attention to actions as motion and appearance patterns
  - For instrumented/constrained environments, relatively simple techniques allow effective gesture or action recognition

# Project 3: Tracking



- **Code reviews due Friday, 7 November, 11:59pm**
- **Revisions due Tuesday, 11 November, 11:59pm**