

# Lecture 13: Tracking motion features – optical flow

Professor Fei-Fei Li  
Stanford Vision Lab

# What we will learn today?

- Introduction
- Optical flow
- Feature tracking
- Applications
- (Supplementary) Technical note

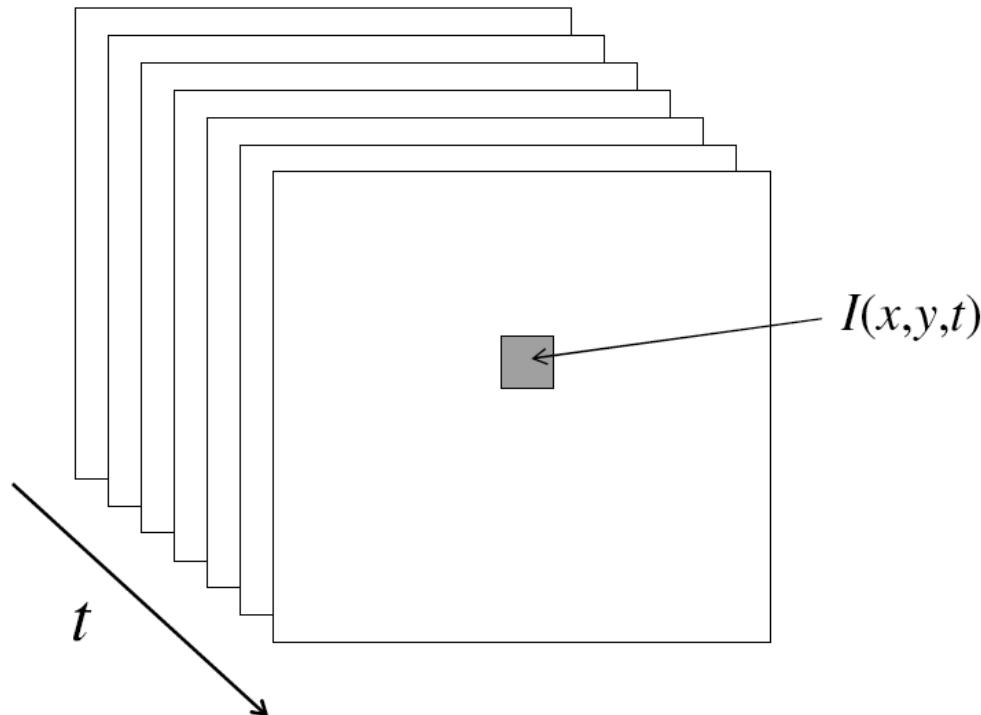
**Reading:** [Szeliski] Chapters: 8.4, 8.5

[Fleet & Weiss, 2005]

<http://www.cs.toronto.edu/pub/jepson/teaching/vision/2503/opticalFlow.pdf>

# From images to videos

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

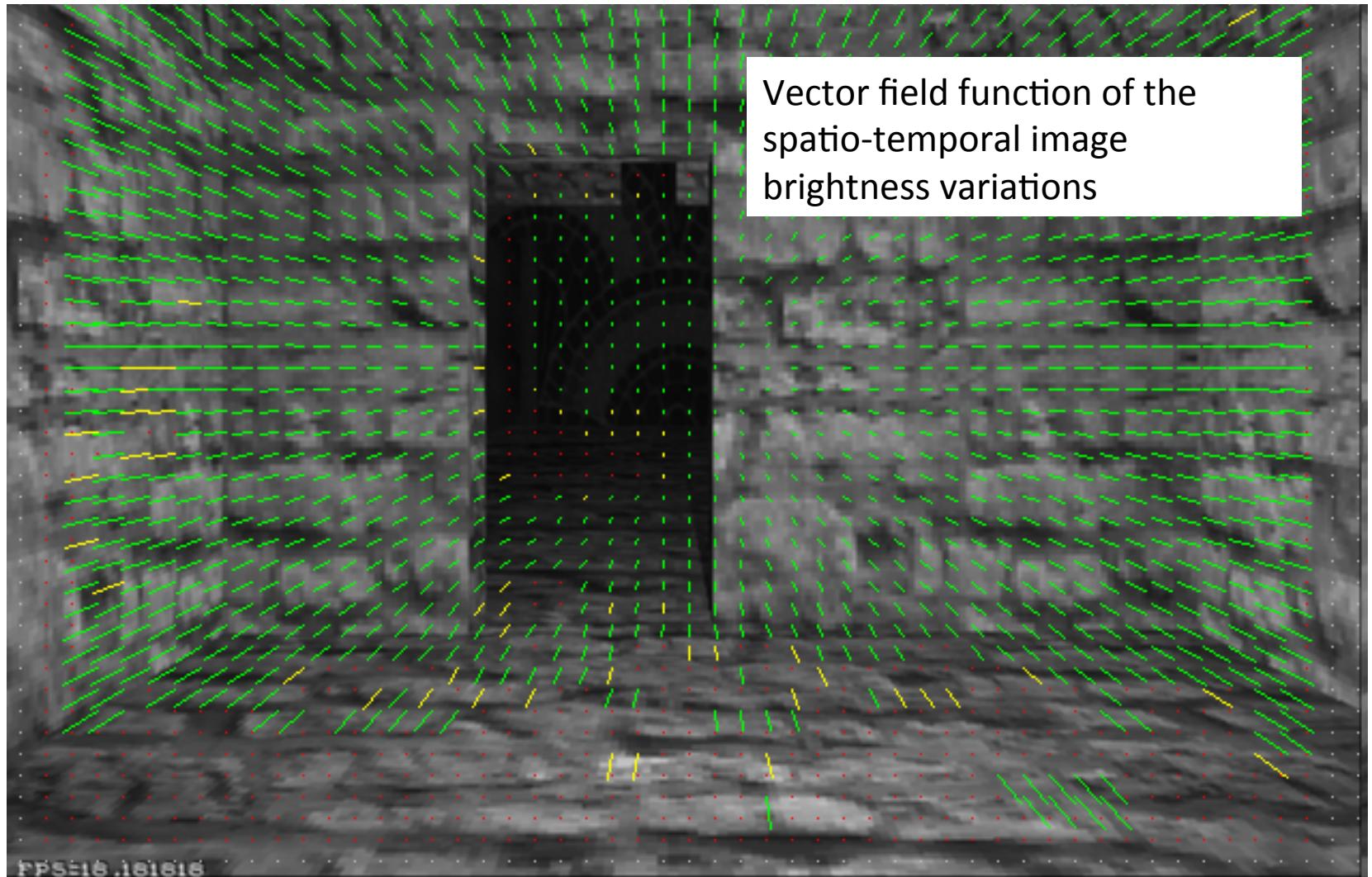


# Motion estimation techniques

- Optical flow
  - Recover image motion at each pixel from spatio-temporal image brightness variations (optical flow)
- Feature-tracking
  - Extract visual features (corners, textured areas) and “track” them over multiple frames



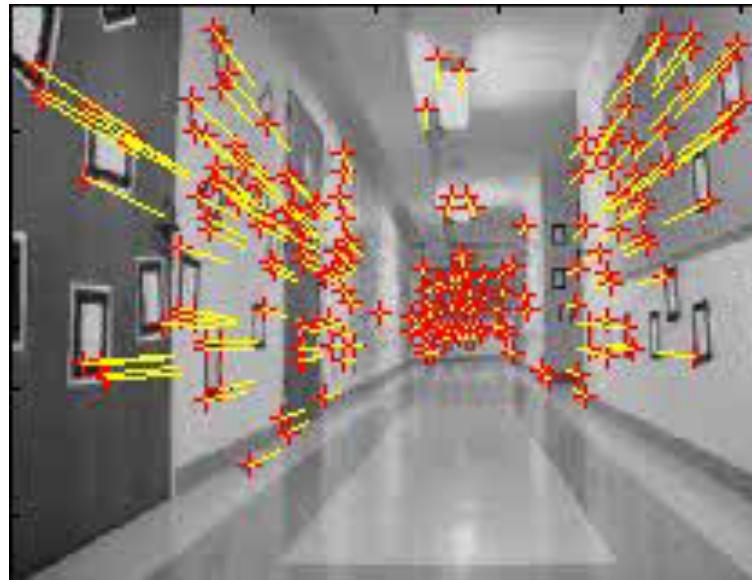
# Optical flow



FPS: 10.161616

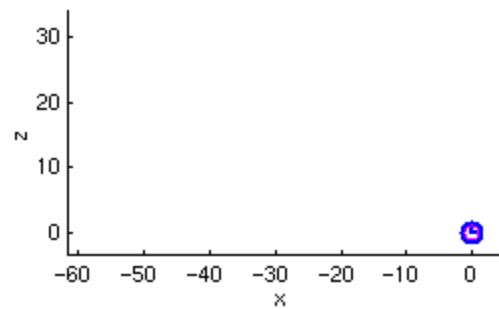
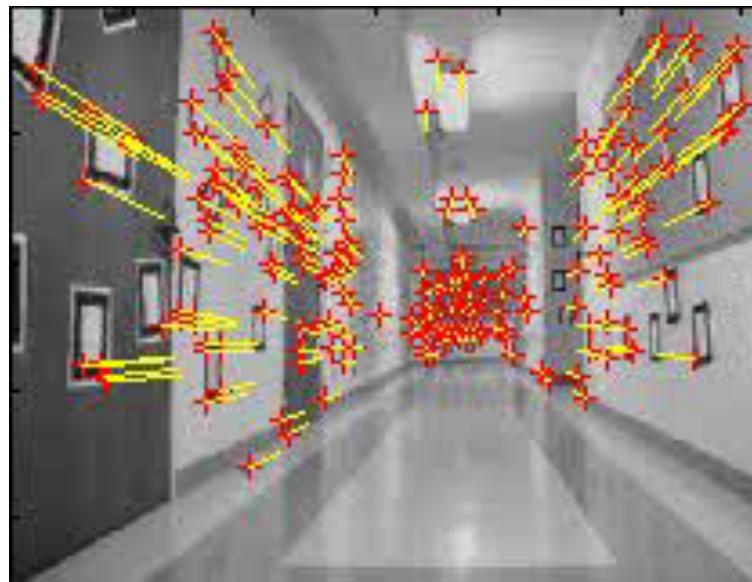
Picture courtesy of Selim Temizer - Learning and Intelligent Systems (LIS) Group, MIT

# Feature-tracking



Courtesy of Jean-Yves Bouguet – Vision Lab, California Institute of Technology

# Feature-tracking



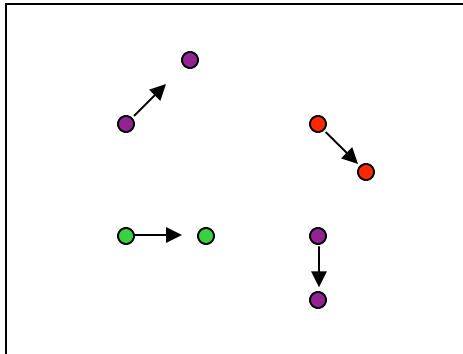
Courtesy of Jean-Yves Bouguet – Vision Lab, California Institute of Technology

# Optical flow

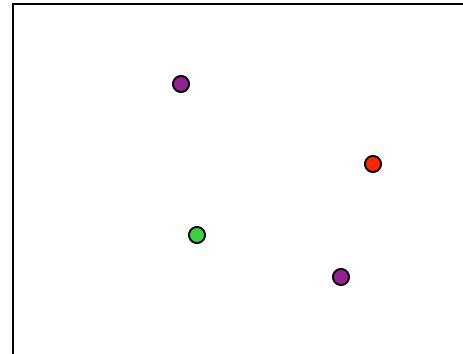
- Definition: optical flow is the *apparent* motion of brightness patterns in the image
- Note: apparent motion can be caused by lighting changes without any actual motion
  - Think of a uniform rotating sphere under fixed lighting vs. a stationary sphere under moving illumination

**GOAL:** Recover image motion at each pixel from optical flow

# Estimating optical flow



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

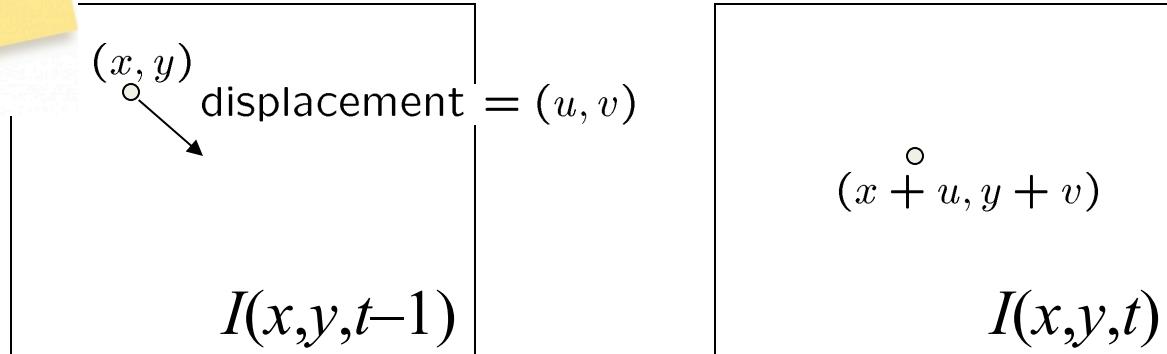


$I(x,y,t)$

- Given two subsequent frames, estimate the apparent motion field  $u(x,y), v(x,y)$  between them
- Key assumptions
  - **Brightness constancy:** projection of the same point looks the same in every frame
  - **Small motion:** points do not move very far
  - **Spatial coherence:** points move like their neighbors

technical  
note

# The brightness constancy constraint



- Brightness Constancy Equation:

$$I(x, y, t - 1) = I(x + u(x, y), y + v(x, y), t)$$

Linearizing the right side using Taylor expansion:

$$I(x + u, y + v, t) \approx I(x, y, t - 1) + \boxed{I_x} \cdot u(x, y) + I_y \cdot v(x, y) + \boxed{I_t}$$

Image derivative along x

$$I(x + u, y + v, t) - I(x, y, t - 1) = I_x \cdot u(x, y) + I_y \cdot v(x, y) + I_t$$

$$\text{Hence, } I_x \cdot u + I_y \cdot v + I_t \approx 0 \rightarrow \nabla I \cdot [u \ v]^T + I_t = 0$$

technical  
note

# The brightness constancy constraint

Can we use this equation to recover image motion ( $u, v$ ) at each pixel?

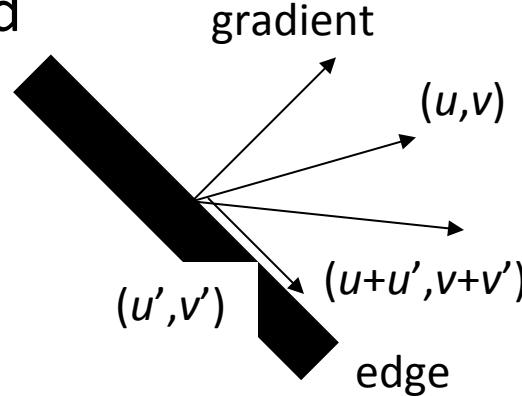
$$\nabla I \cdot [u \ v]^T + I_t = 0$$

- How many equations and unknowns per pixel?
  - One equation (this is a scalar equation!), two unknowns ( $u, v$ )

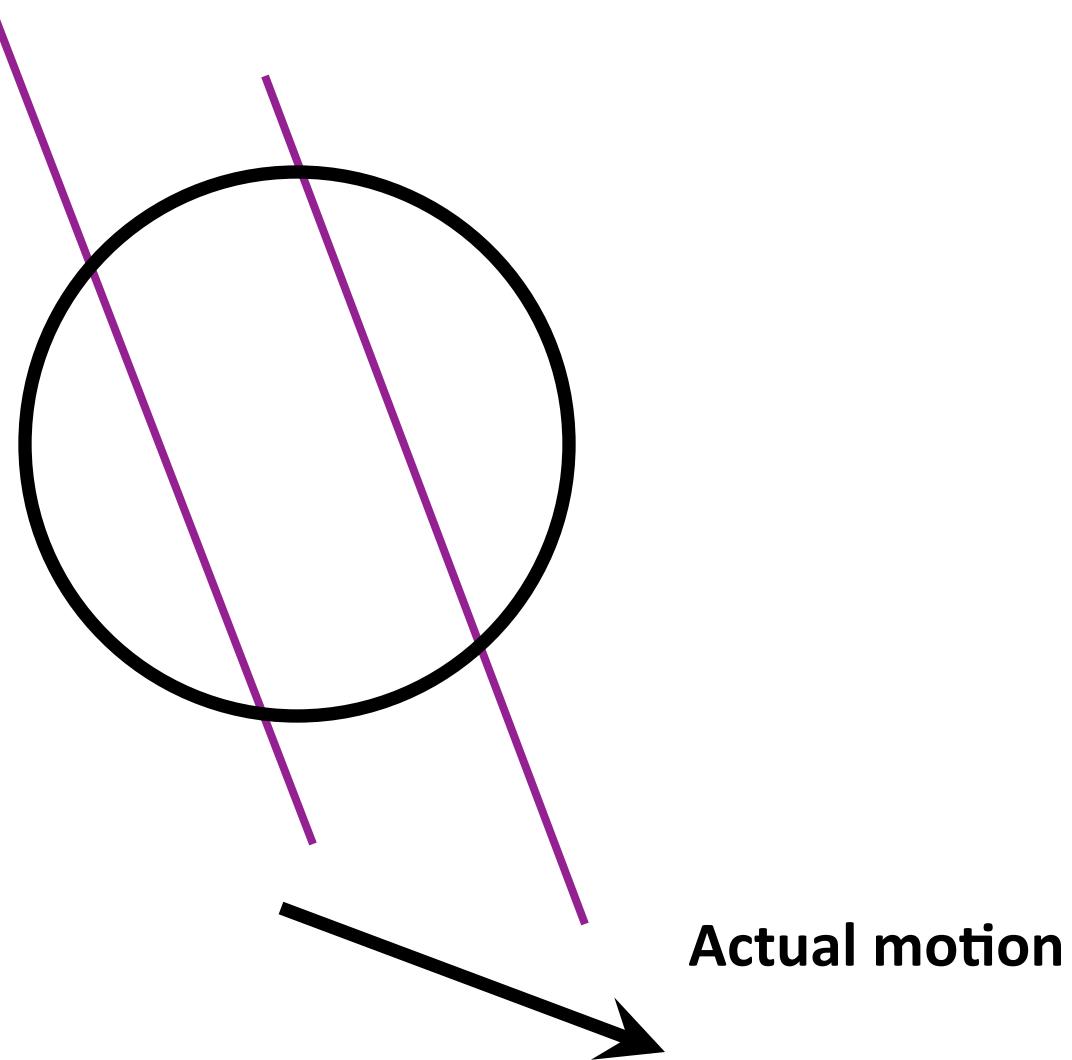
The component of the flow perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

If  $(u, v)$  satisfies the equation,  
so does  $(u+u', v+v')$  if

$$\nabla I \cdot [u' \ v']^T = 0$$

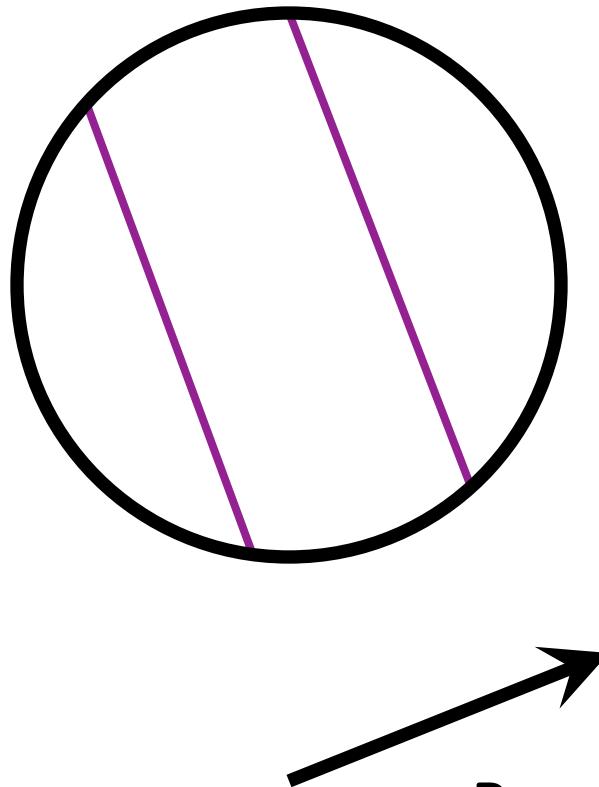


# The aperture problem



Source: Silvio Savarese

# The aperture problem



Source: Silvio Savarese

# The barber pole illusion



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

Source: Silvio Savarese

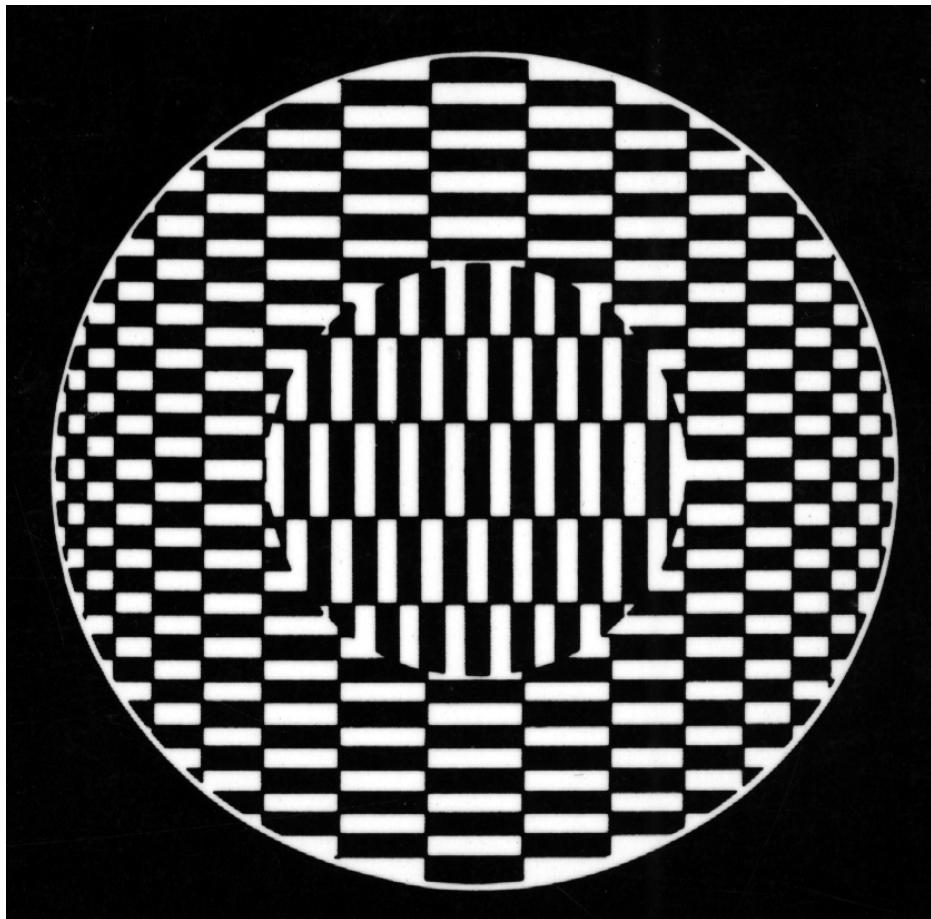
# The barber pole illusion



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

Source: Silvio Savarese

# Aperture problem cont'd



\* From Marc Pollefeys COMP 256 2003

technical  
note

# Solving the ambiguity...

B. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

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

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

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

Source: Silvio Savarese

# Lucas-Kanade flow

- Overconstrained linear system:

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

$A \quad d = b$   
 $25 \times 2 \quad 2 \times 1 \quad 25 \times 1$

# Conditions for solvability

- When is this system solvable?
  - What if the window contains just a single straight edge?

# Lucas-Kanade flow

- Overconstrained linear system

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix} \quad \begin{matrix} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{matrix}$$

Least squares solution for  $d$  given by  $(A^T A)^{-1} d = A^T b$

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

$$A^T A \qquad \qquad \qquad A^T b$$

The summations are over all pixels in the  $K \times K$  window

# Conditions for solvability

- Optimal  $(u, v)$  satisfies Lucas-Kanade equation

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

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

## When is This Solvable?

- $A^T A$  should be invertible
- $A^T A$  should not be too small due to noise
  - 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)

Does this remind anything to you?

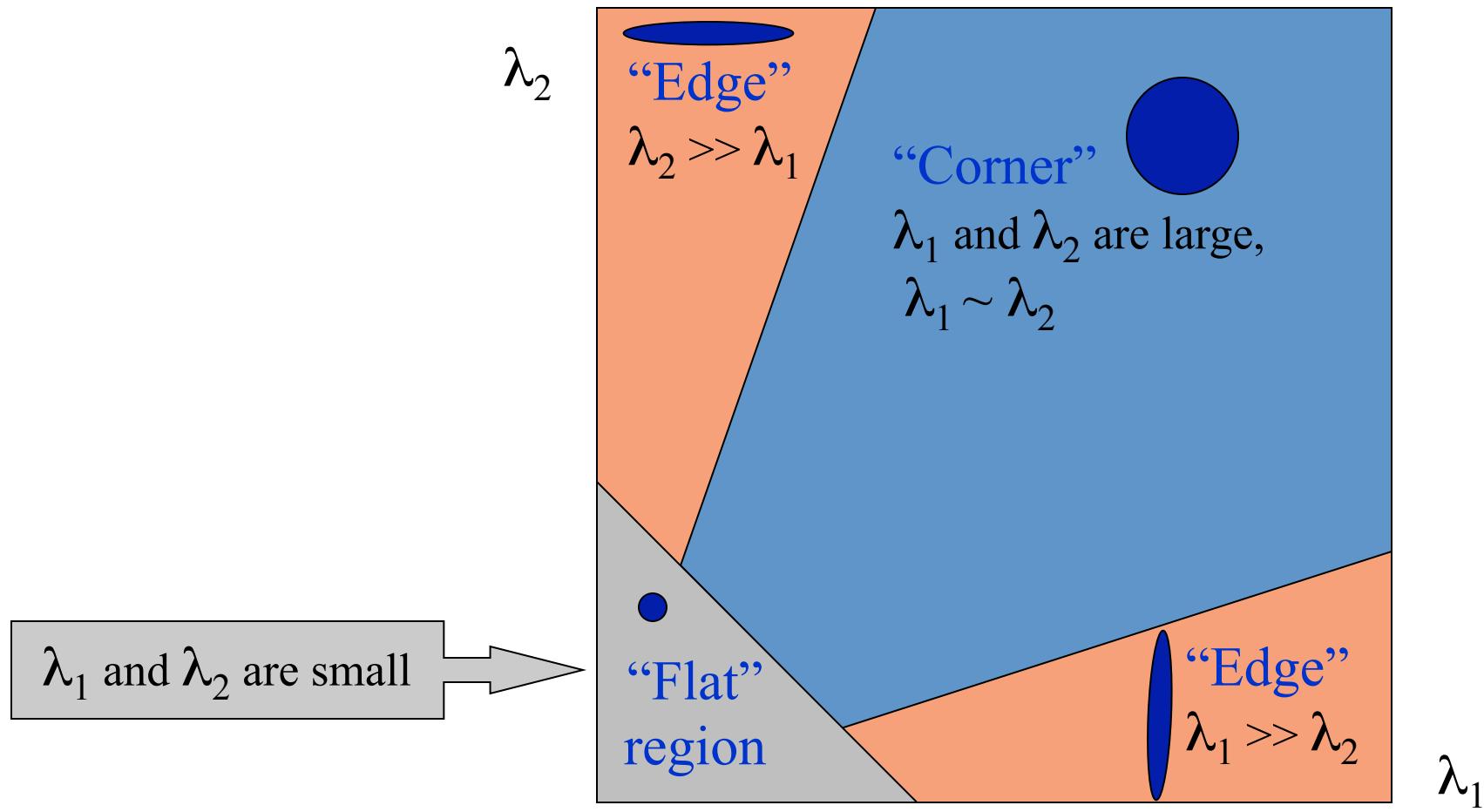
# $M = A^T A$ is the *second moment matrix* ! (Harris corner detector...)

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = \sum \begin{bmatrix} I_x \\ I_y \end{bmatrix} [I_x \ I_y] = \sum \nabla I (\nabla I)^T$$

- Eigenvectors and eigenvalues of  $A^T A$  relate to edge direction and magnitude
  - The eigenvector associated with the larger eigenvalue points in the direction of fastest intensity change
  - The other eigenvector is orthogonal to it

# Interpreting the eigenvalues

Classification of image points using eigenvalues of the second moment matrix:



technical  
note

# Edge



$$\sum \nabla I (\nabla I)^T$$

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

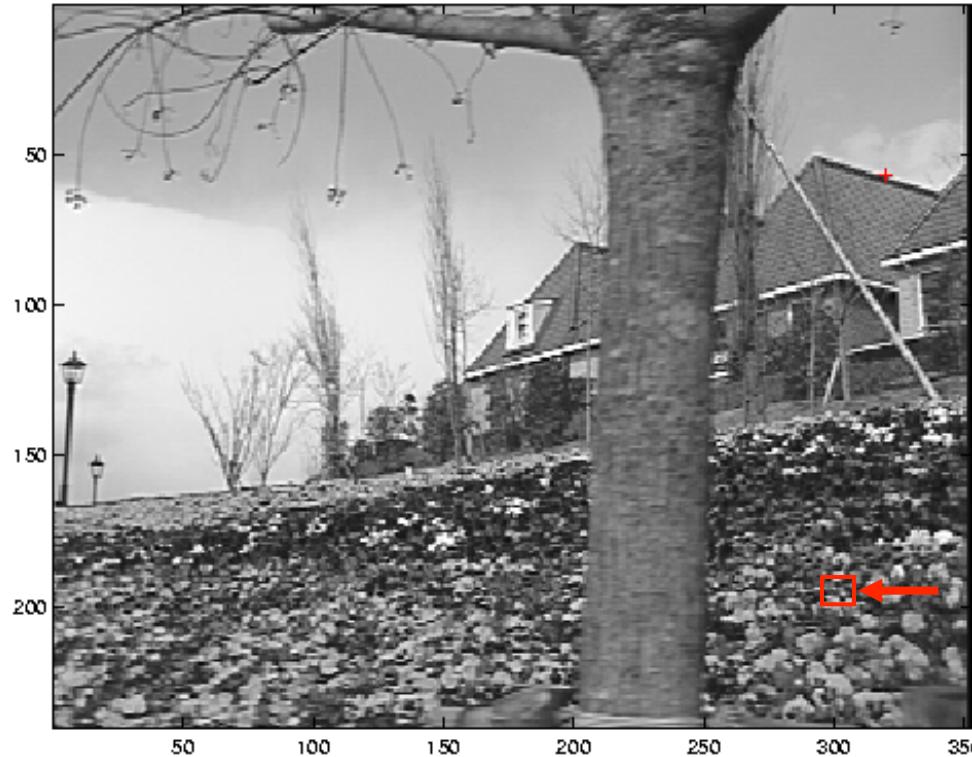
# Low-texture region



$$\sum \nabla I (\nabla I)^T$$

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

# High-texture region



$$\sum \nabla I (\nabla I)^T$$

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

# What we will learn today?

- Introduction
- Optical flow
- Feature tracking
- Applications
- (Supplementary) Technical note

# What are good features to track?

- Can measure “quality” of features from just a single image
- Hence: tracking Harris corners (or equivalent) guarantees small error sensitivity!

→ Implemented in Open CV

# Recap

- Key assumptions (Errors in Lucas-Kanade)

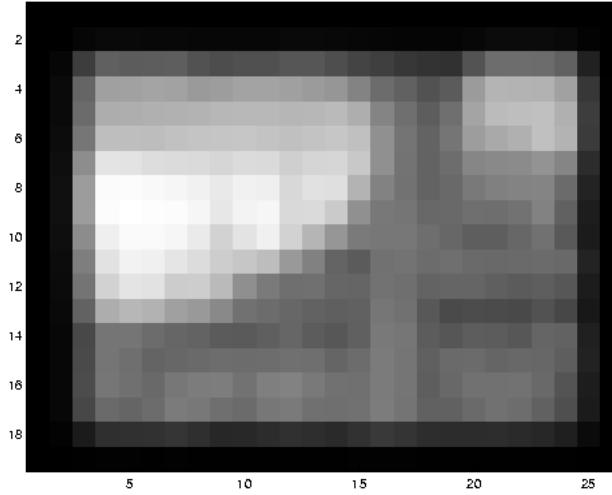
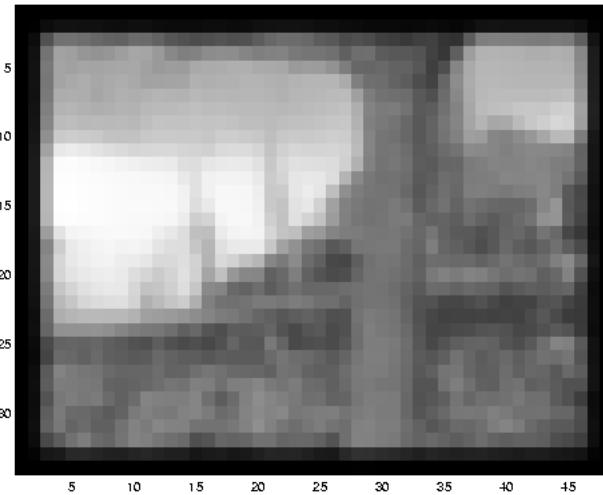
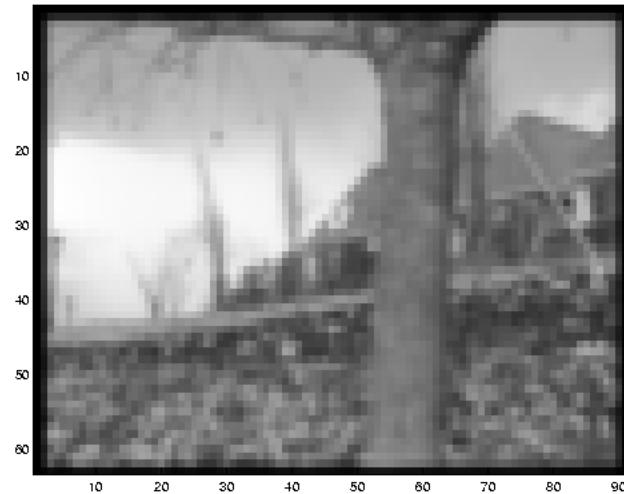
- **Small motion:** points do not move very far
- **Brightness constancy:** projection of the same point looks the same in every frame
- **Spatial coherence:** points move like their neighbors

# Revisiting the small motion assumption



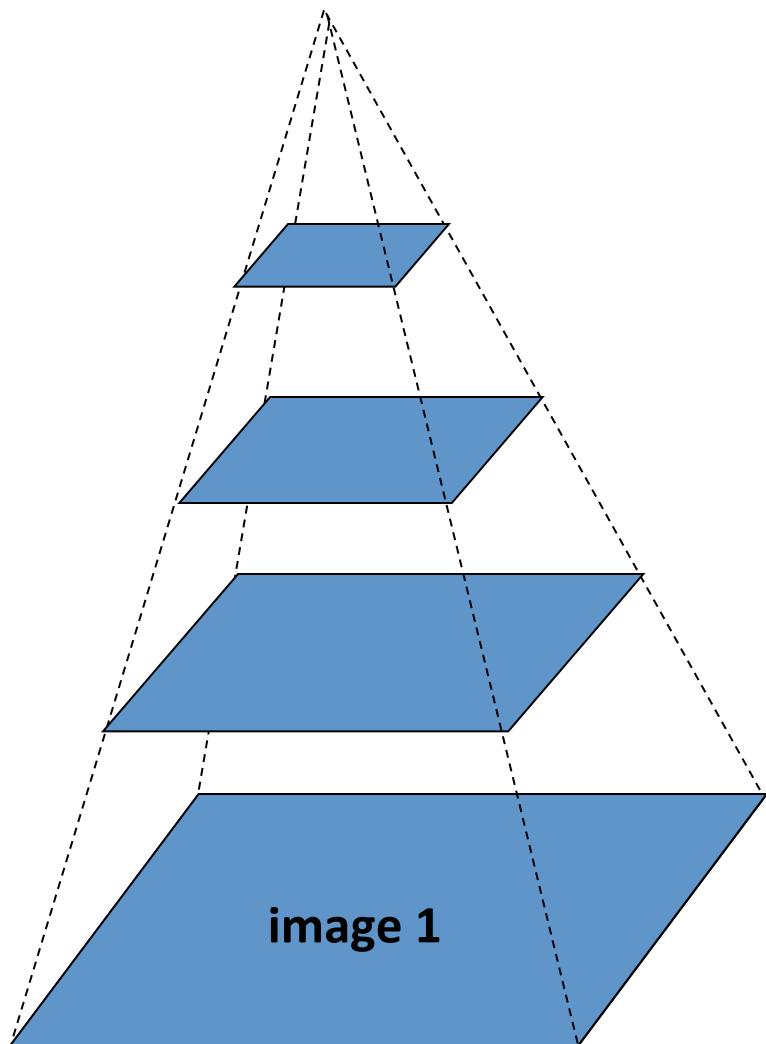
- Is this motion small enough?
  - Probably not—it's much larger than one pixel (2<sup>nd</sup> order terms dominate)
  - How might we solve this problem?

# Reduce the resolution!



\* From Khurram Hassan-Shafique CAP5415 Computer Vision 2003

# Coarse-to-fine optical flow estimation



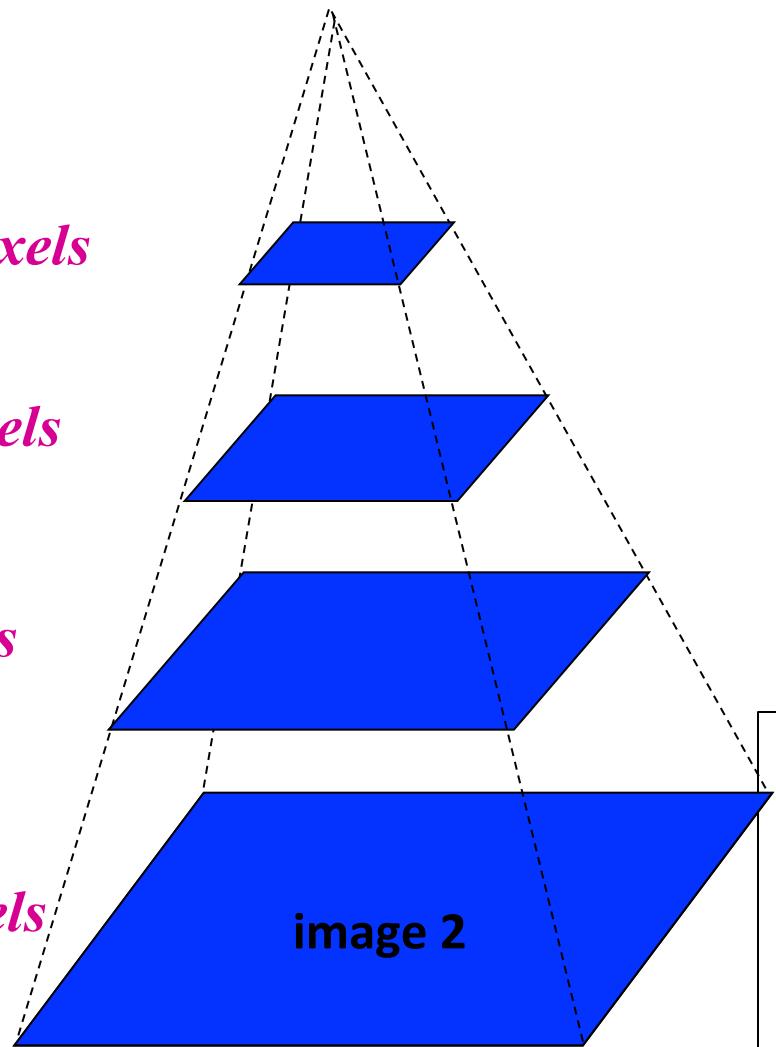
Gaussian pyramid of image 1

$u=1.25 \text{ pixels}$

$u=2.5 \text{ pixels}$

$u=5 \text{ pixels}$

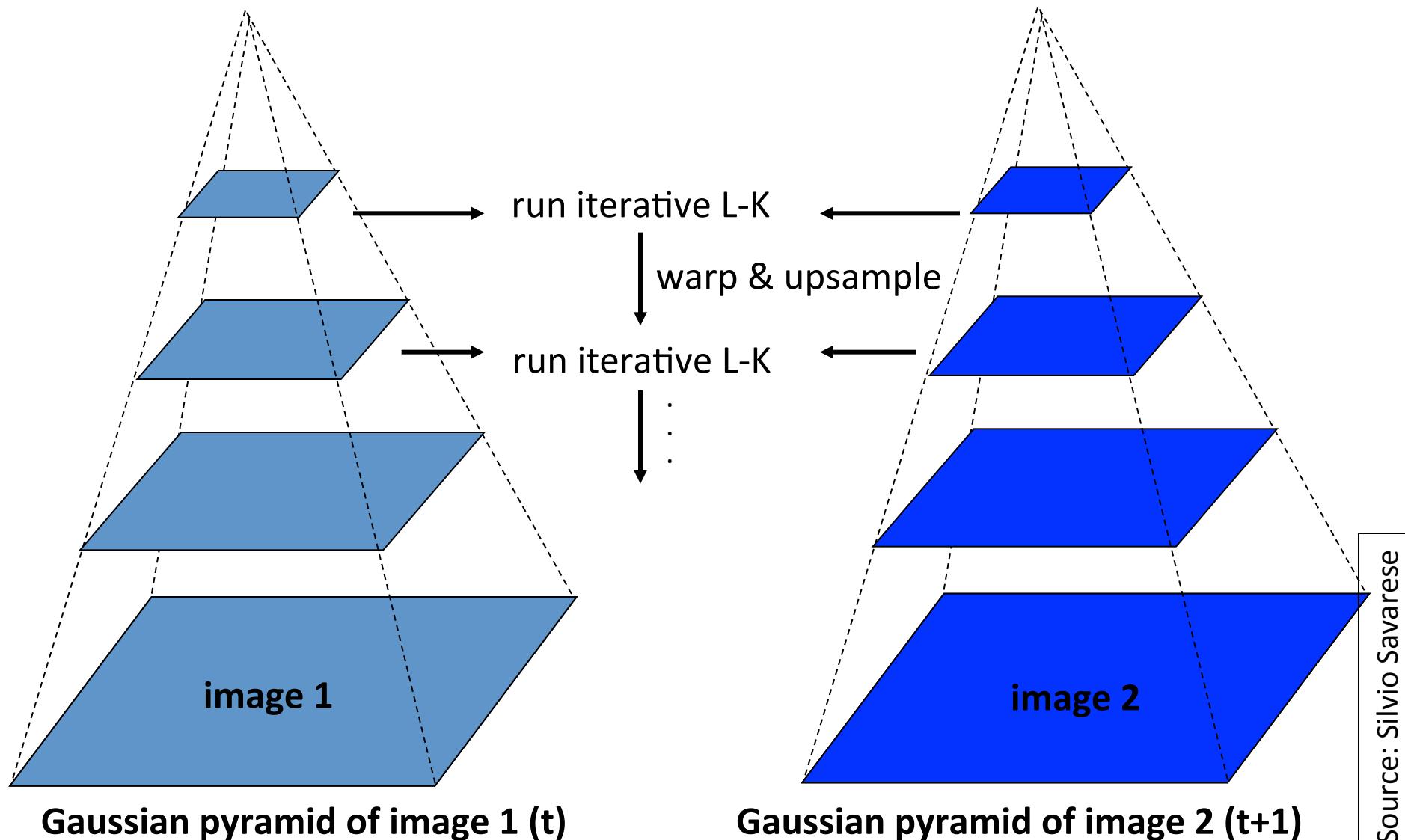
$u=10 \text{ pixels}$



Gaussian pyramid of image 2

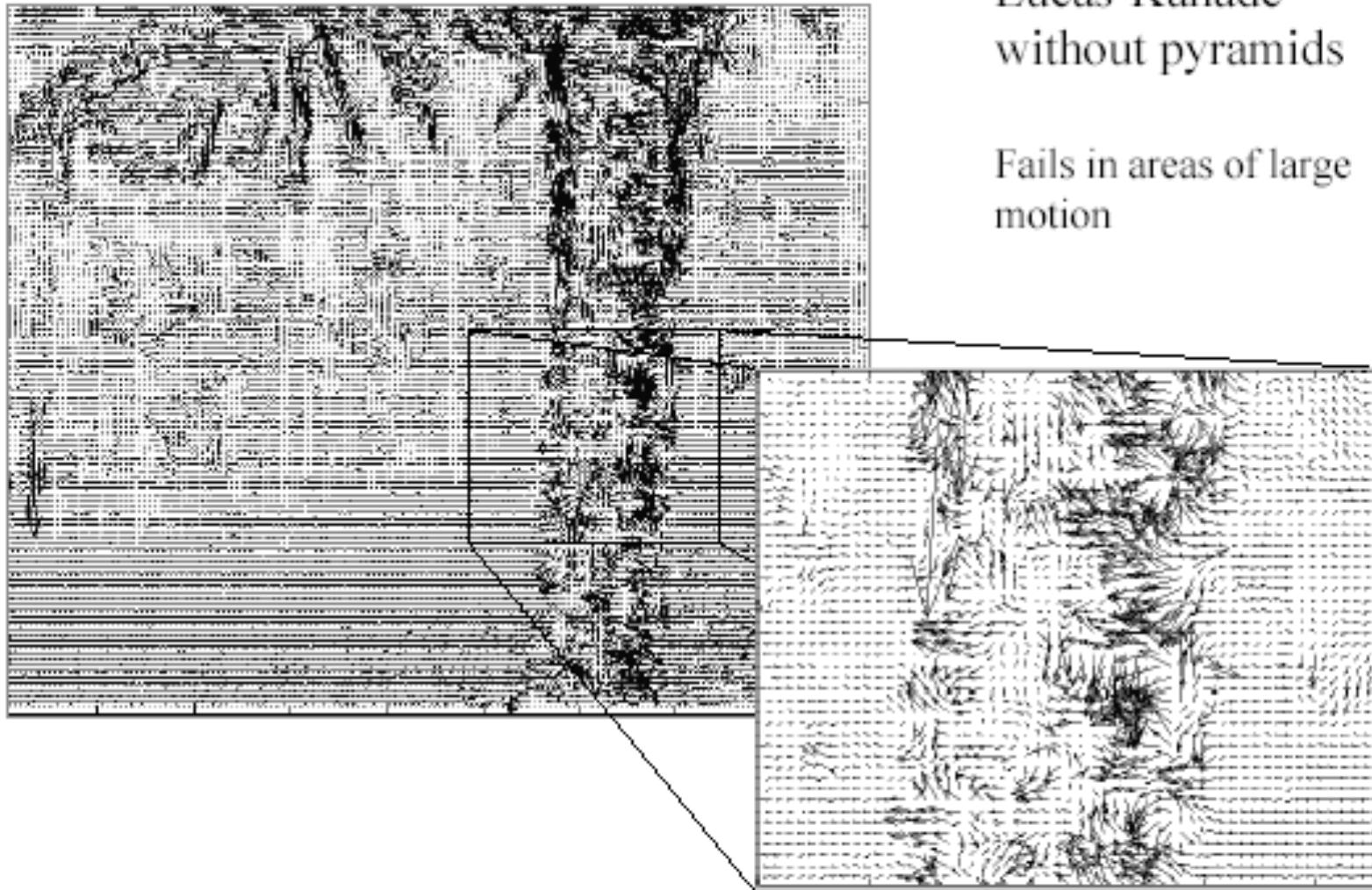
Source: Silvio Savarese

# Coarse-to-fine optical flow estimation



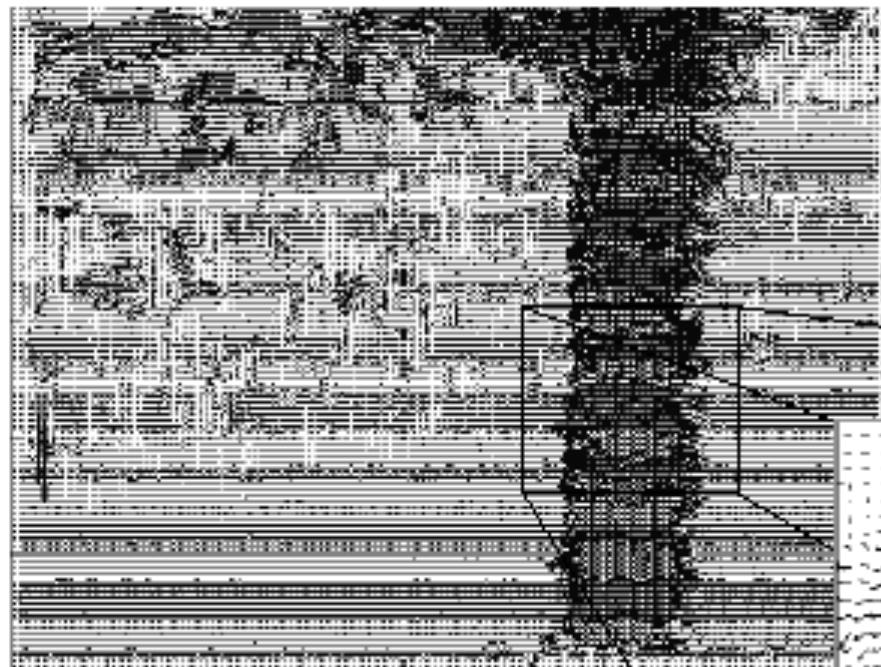
Source: Silvio Savarese

# Optical Flow Results

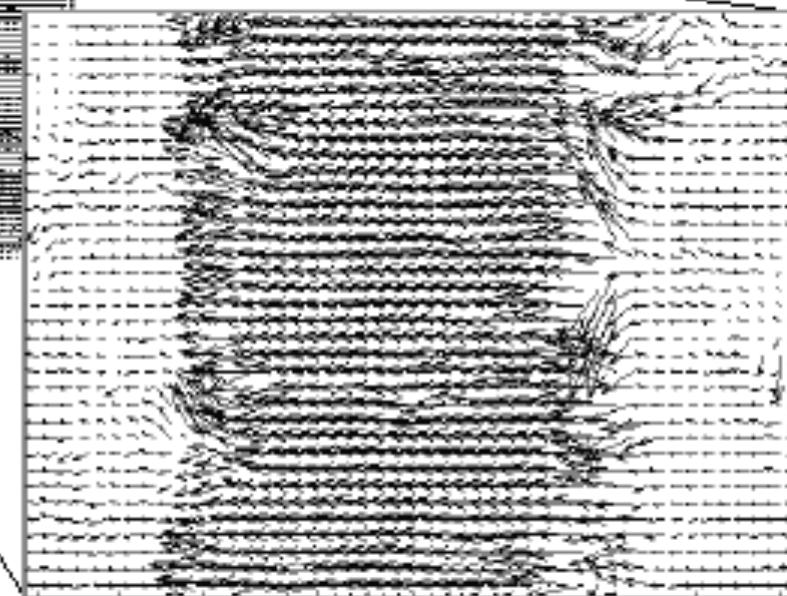


\* From Khurram Hassan-Shafique CAP5415 Computer Vision 2003

# Optical Flow Results



Lucas-Kanade with Pyramids



- <http://www.ces.clemson.edu/~stb/klt/>
- OpenCV

# Recap

- Key assumptions (Errors in Lucas-Kanade)
  - **Small motion:** points do not move very far
  - **Brightness constancy:** projection of the same point looks the same in every frame
  - **Spatial coherence:** points move like their neighbors

# Reminder: Gestalt – common fate



Common Fate

# Motion segmentation

- How do we represent the motion in this scene?



Source: Silvio Savarese

# Motion segmentation

J. Wang and E. Adelson. Layered Representation for Motion Analysis. *CVPR* 1993.

- Break image sequence into “layers” each of which has a coherent (affine) motion



Source: Silvio Savarese

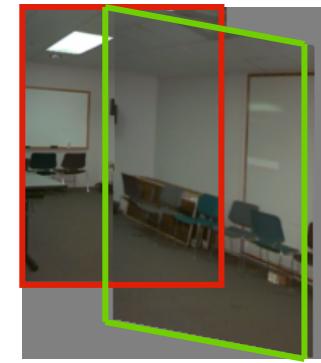
# Affine motion

$$u(x, y) = a_1 + a_2 x + a_3 y$$

$$v(x, y) = a_4 + a_5 x + a_6 y$$

- Substituting into the brightness constancy equation:

$$I_x \cdot u + I_y \cdot v + I_t \approx 0$$

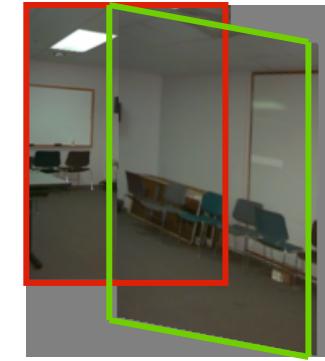


# Affine motion

$$u(x, y) = a_1 + a_2x + a_3y$$

$$v(x, y) = a_4 + a_5x + a_6y$$

- Substituting into the brightness constancy equation:



$$I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t \approx 0$$

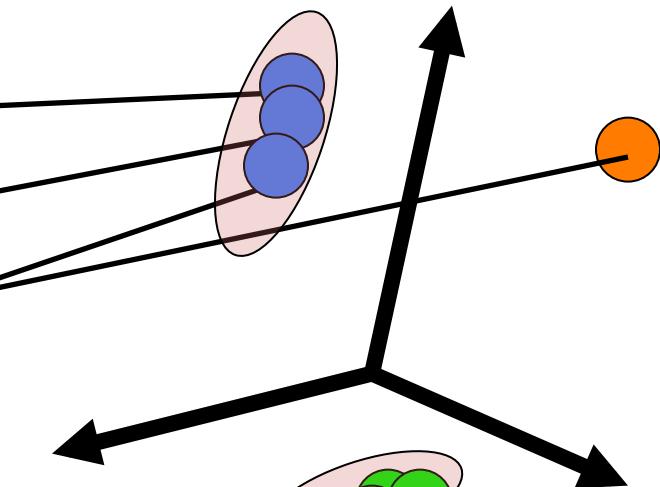
- Each pixel provides 1 linear constraint in 6 unknowns
- Least squares minimization:

$$Err(\vec{a}) = \sum [I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t]^2$$

Source: Silvio Savarese

# How do we estimate the layers?

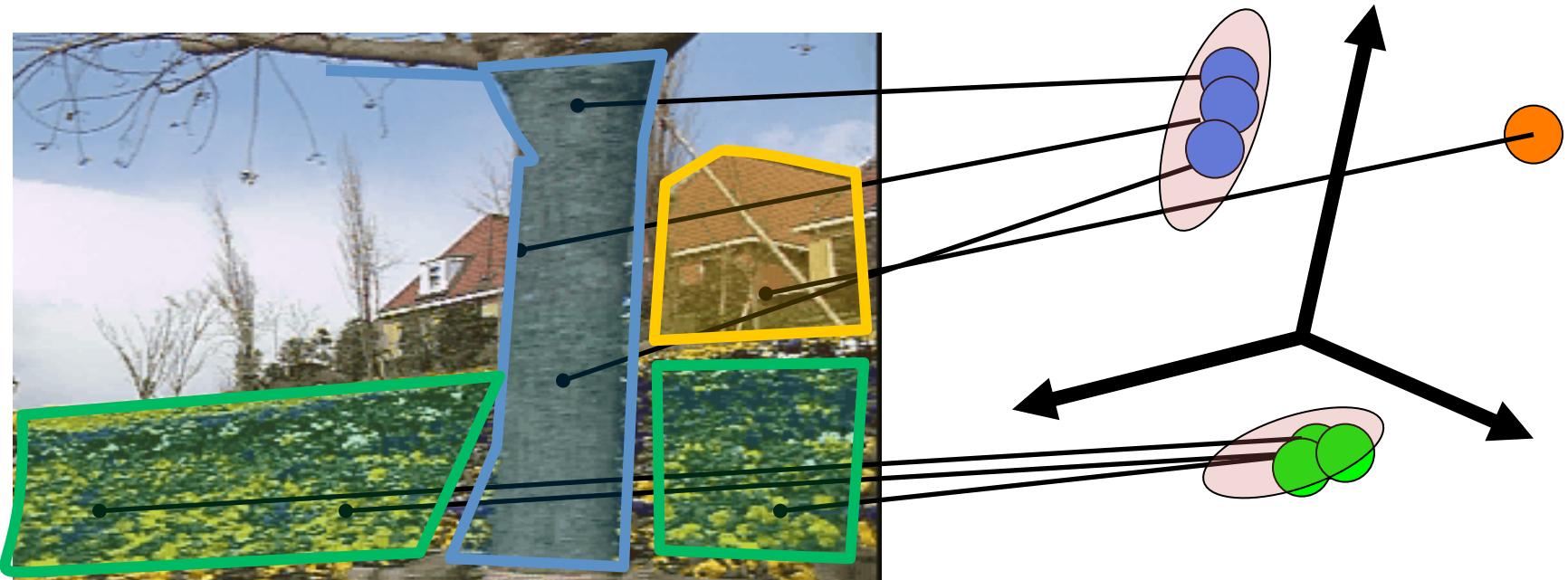
- 1. Obtain a set of initial affine motion hypotheses
  - Divide the image into blocks and estimate affine motion parameters in each block by least squares
    - Eliminate hypotheses with high residual error
  - Map into motion parameter space
  - Perform k-means clustering on affine motion parameters
    - Merge clusters that are close and retain the largest clusters to obtain a smaller set of hypotheses to describe all the motions in the scene



Source: Silvio Savarese

# How do we estimate the layers?

- 1. Obtain a set of initial affine motion hypotheses
  - Divide the image into blocks and estimate affine motion parameters in each block by least squares
    - Eliminate hypotheses with high residual error
  - Map into motion parameter space
  - Perform k-means clustering on affine motion parameters
    - Merge clusters that are close and retain the largest clusters to obtain a smaller set of hypotheses to describe all the motions in the scene



Source: Silvio Savarese

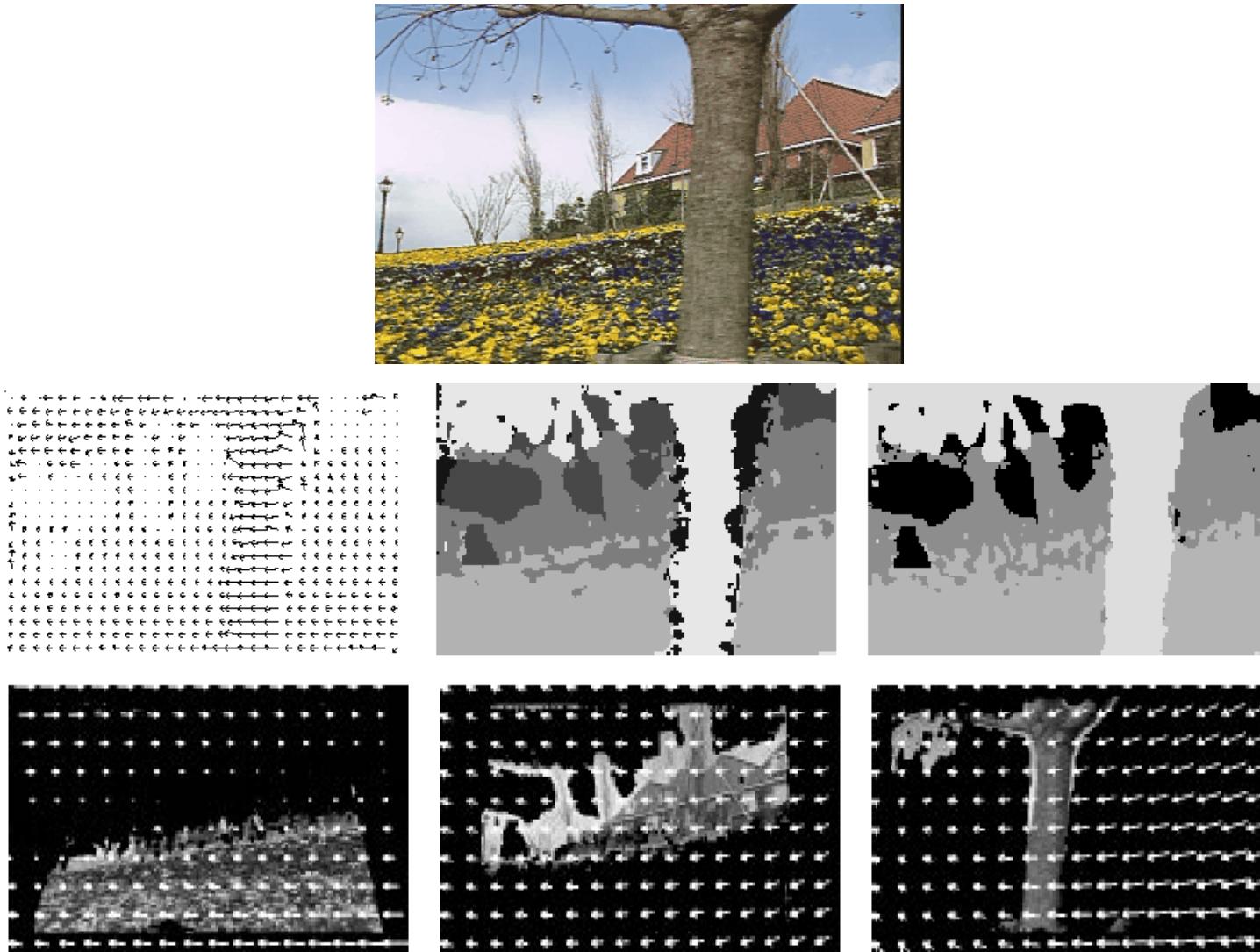
# How do we estimate the layers?

- 1. Obtain a set of initial affine motion hypotheses
  - Divide the image into blocks and estimate affine motion parameters in each block by least squares
    - Eliminate hypotheses with high residual error
  - Map into motion parameter space
  - Perform k-means clustering on affine motion parameters
    - Merge clusters that are close and retain the largest clusters to obtain a smaller set of hypotheses to describe all the motions in the scene

## 2. Iterate until convergence:

- Assign each pixel to best hypothesis
  - Pixels with high residual error remain unassigned
- Perform region filtering to enforce spatial constraints
- Re-estimate affine motions in each region

# Example result



J. Wang and E. Adelson. [Layered Representation for Motion Analysis](#). CVPR 1993.

Source: Silvio Savarese

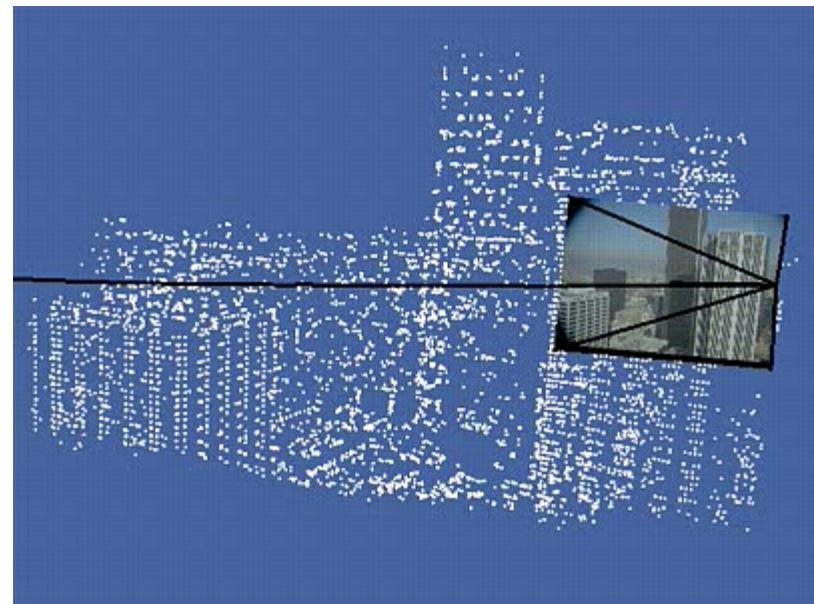
# What we will learn today?

- Introduction
- Optical flow
- Feature tracking
- Applications
- (Supplementary) Technical note

# Uses of motion

- Tracking features
- Segmenting objects based on motion cues
- Learning dynamical models
- Improving video quality
  - Motion stabilization
  - Super resolution
- Tracking objects
- Recognizing events and activities

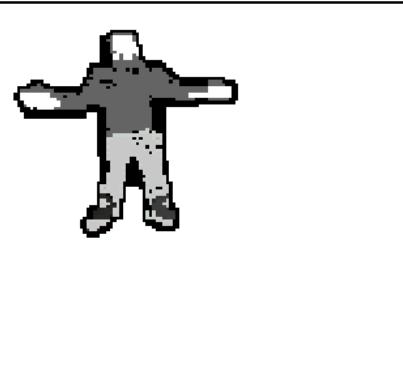
# Estimating 3D structure



Source: Silvio Savarese

# Segmenting objects based on motion cues

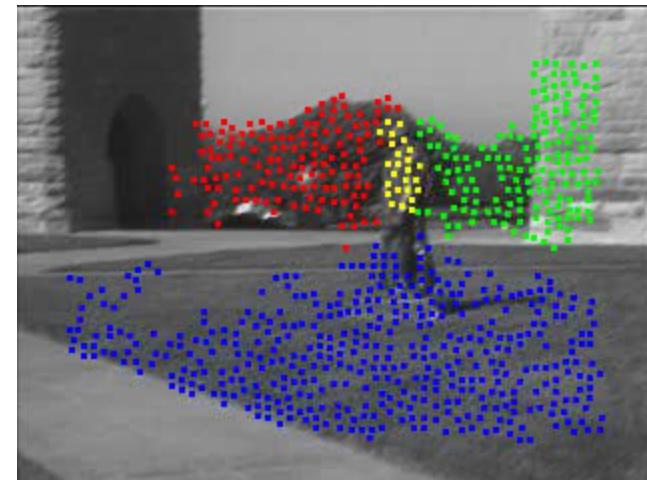
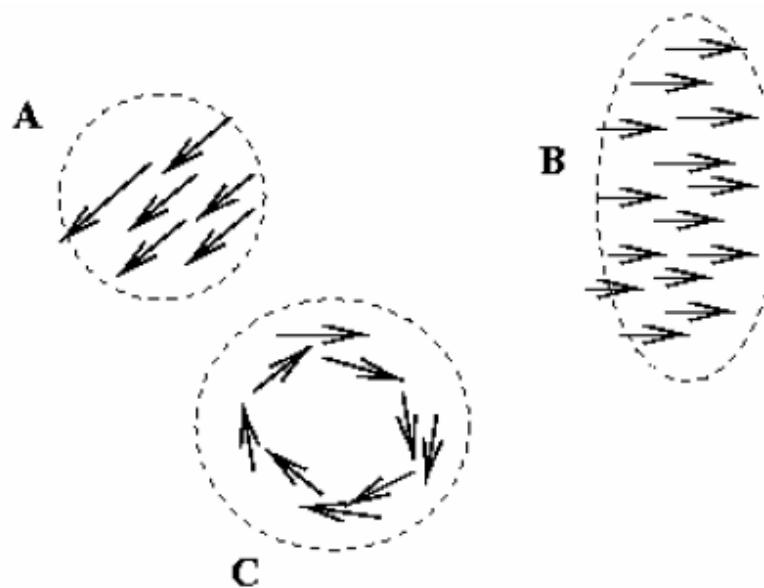
- Background subtraction
  - A static camera is observing a scene
  - Goal: separate the static *background* from the moving *foreground*



Source: Silvio Savarese

# Segmenting objects based on motion cues

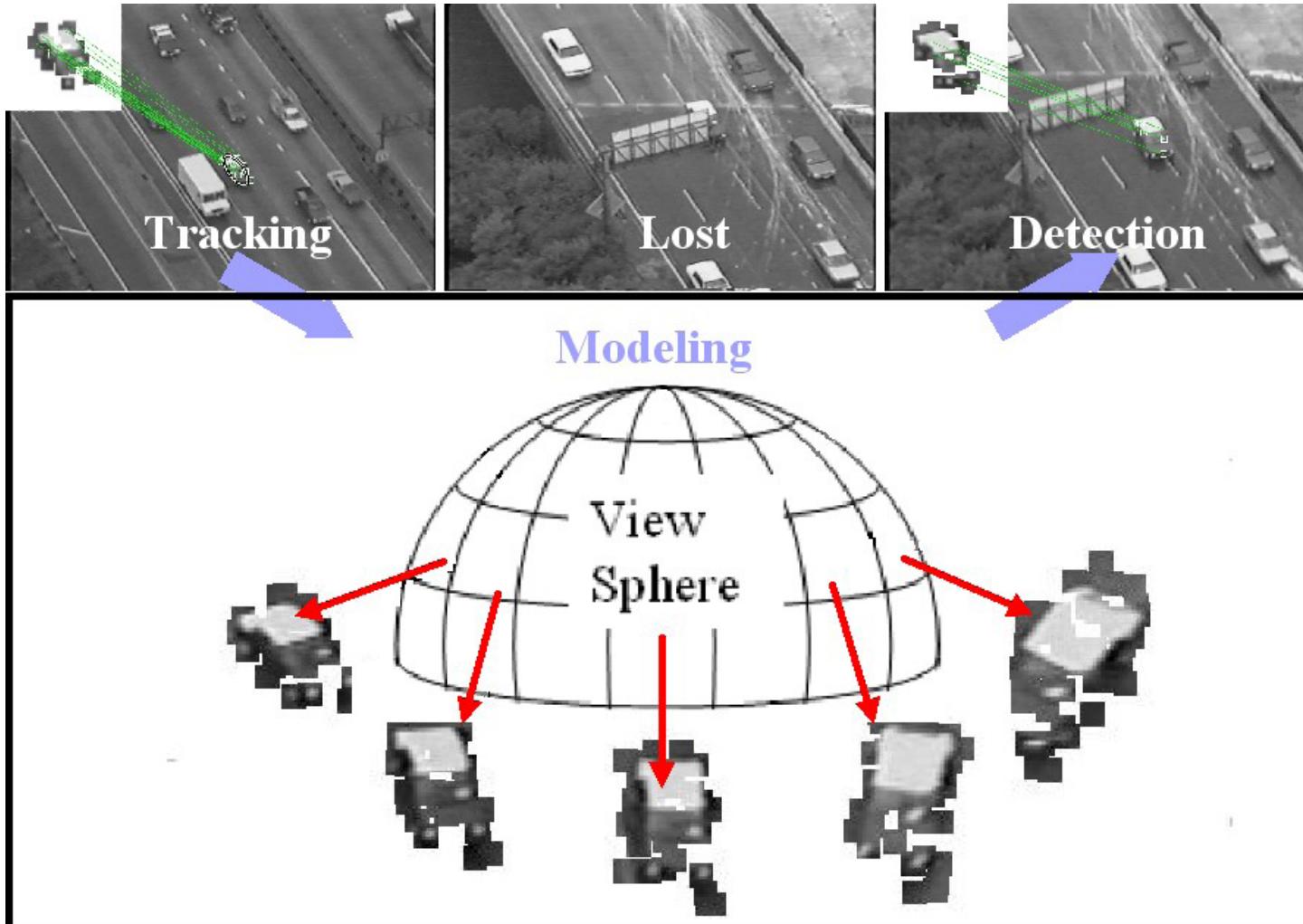
- Motion segmentation
  - Segment the video into multiple *coherently* moving objects



S. J. Pundlik and S. T. Birchfield, Motion Segmentation at Any Speed, Proceedings of the British Machine Vision Conference (BMVC) 2006

Source: Silvio Savarese

# Tracking objects



Z.Yin and R.Collins, "On-the-fly Object Modeling while Tracking," *IEEE Computer Vision and Pattern Recognition (CVPR '07)*, Minneapolis, MN, June 2007.

Source: Silvio Savarese

# Synthesizing dynamic textures



# Super-resolution

Example: A set of low quality images

Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S
Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S
Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S	Most of the test data o couple of exceptions. 1 low-temperature solder investigated (or some manufacturing technol nonwetting of 40In40Sb microstructural coarse mal cycling of 58Bi42S

Source: Silvio Savarese

# Super-resolution

Each of these images looks like this:

Most of the test data are couple of exceptions. low-temperature solder investigated (or some manufacturing technology) re-wetting of 400-1050 microstructural coarse and cycling of FSP at 425

Source: Silvio Savarese

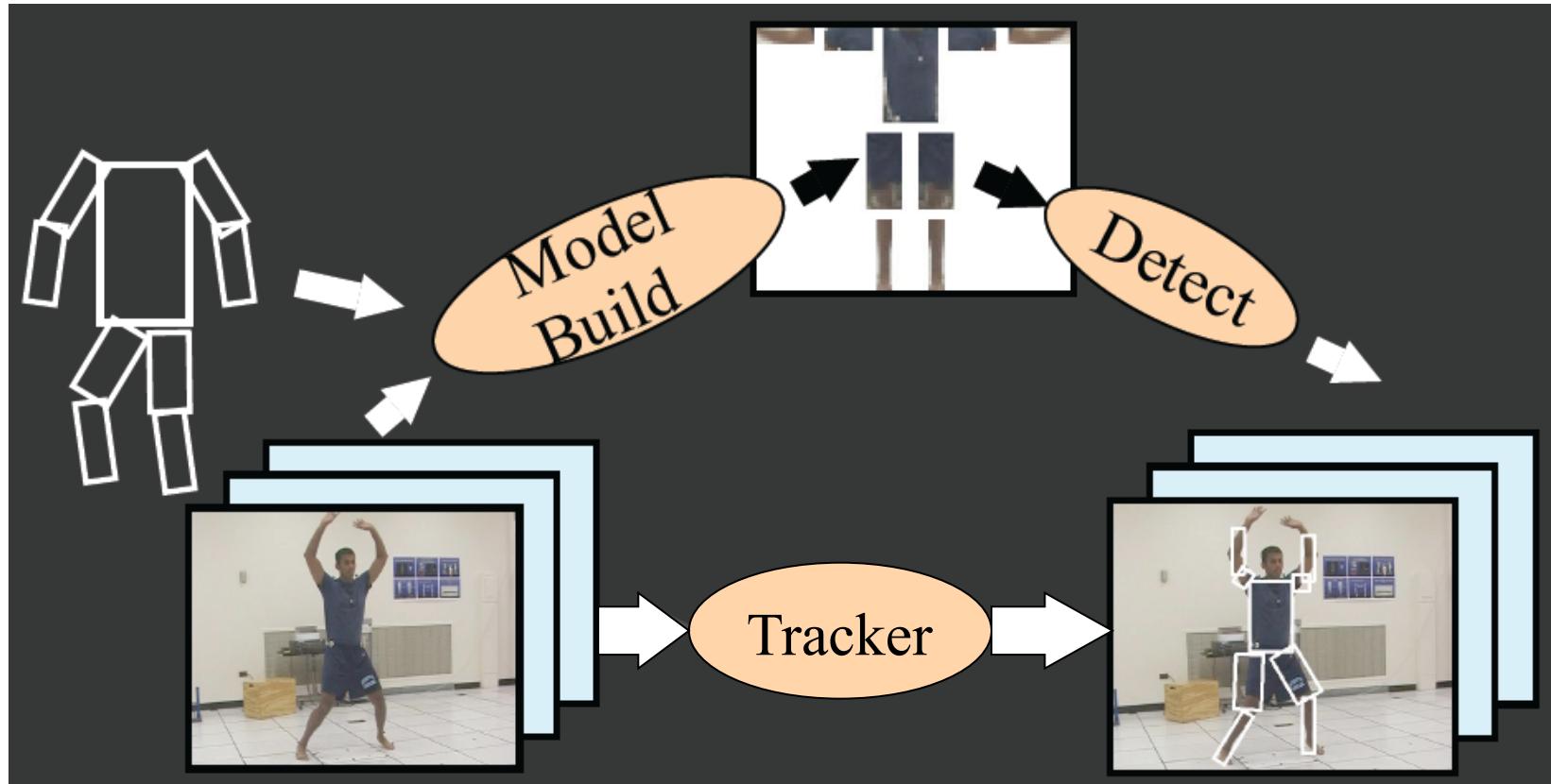
# Super-resolution

The recovery result:

Most of the test data obtained in the literature are in agreement with the model, with a couple of exceptions. These exceptions are related to low-temperature solder joints, which have been either investigated (or some cases not) by different manufacturing technologies. In the case of nonwetting of 40In40Sn solder joints, the model predicts a microstructural coarsening due to thermal cycling of 58Bi42Sb solder joints.

Source: Silvio Savarese

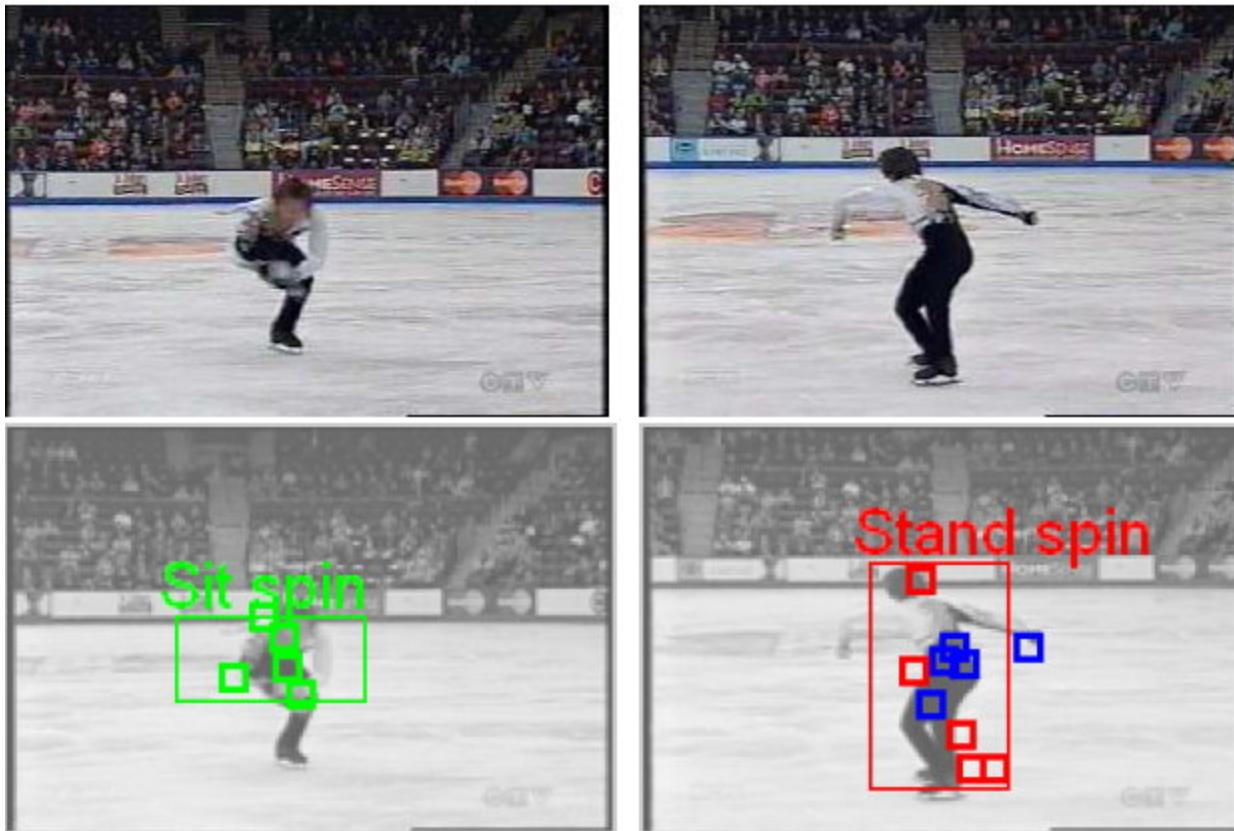
# Recognizing events and activities



D. Ramanan, D. Forsyth, and A. Zisserman. [Tracking People by Learning their Appearance](#). PAMI 2007.

Source: Silvio Savarese

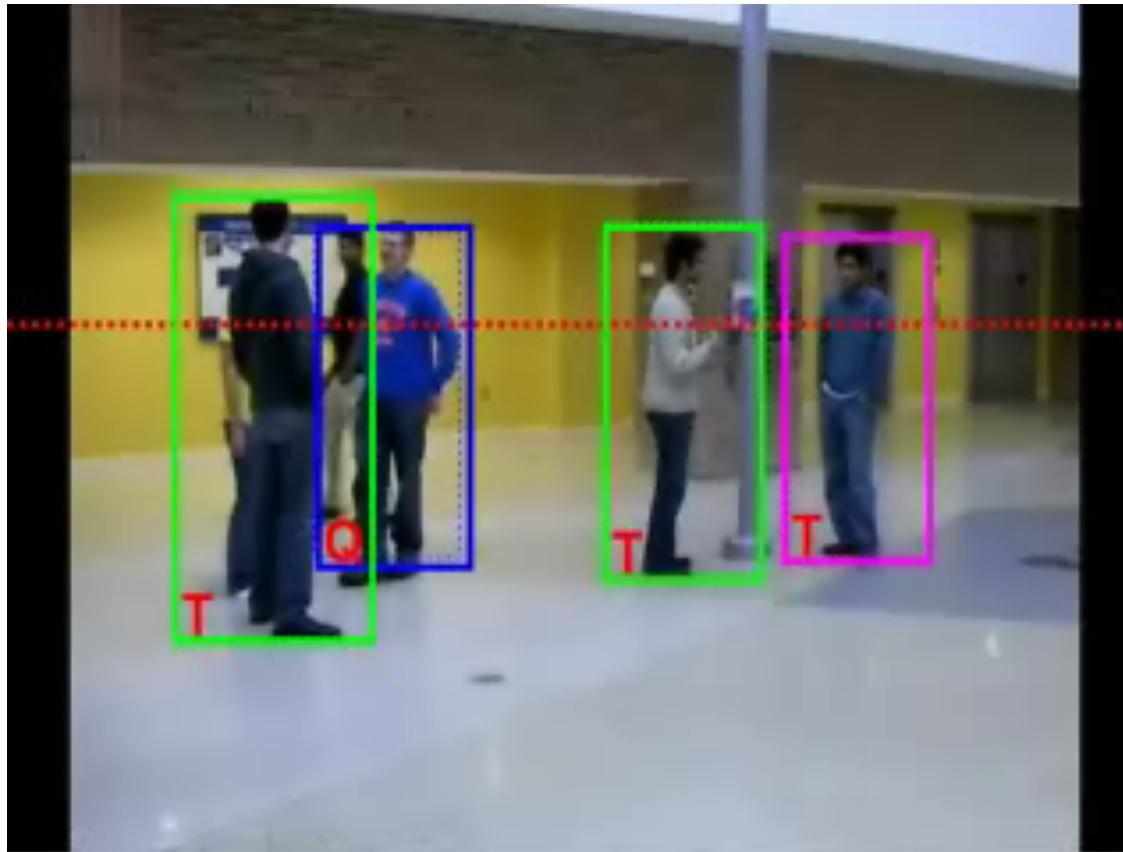
# Recognizing events and activities



Juan Carlos Niebles, Hongcheng Wang and Li Fei-Fei, **Unsupervised Learning of Human Action Categories Using Spatial-Temporal Words**, ([BMVC](#)), Edinburgh, 2006.

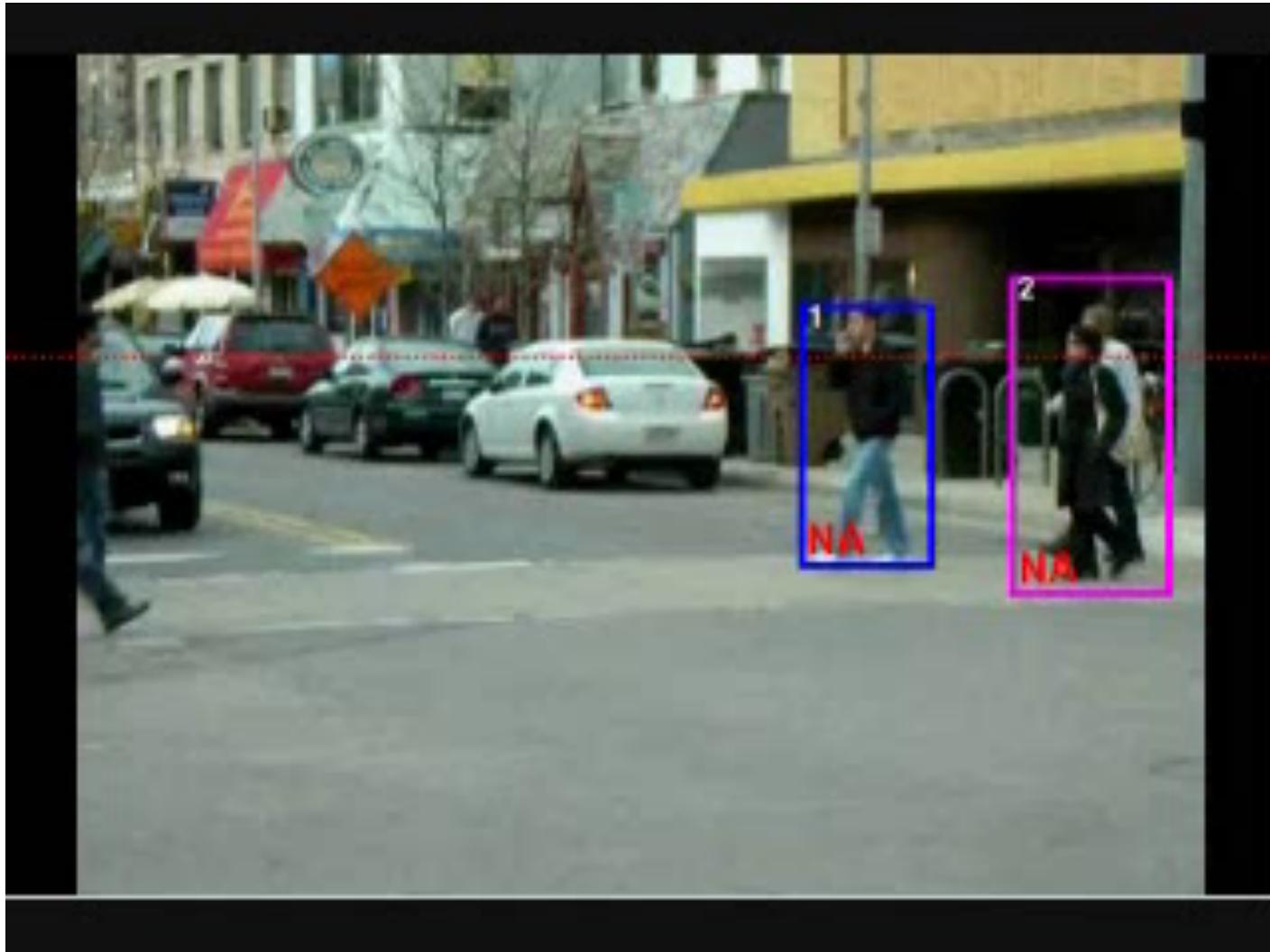
# Recognizing events and activities

Crossing – Talking – Queuing – Dancing – jogging

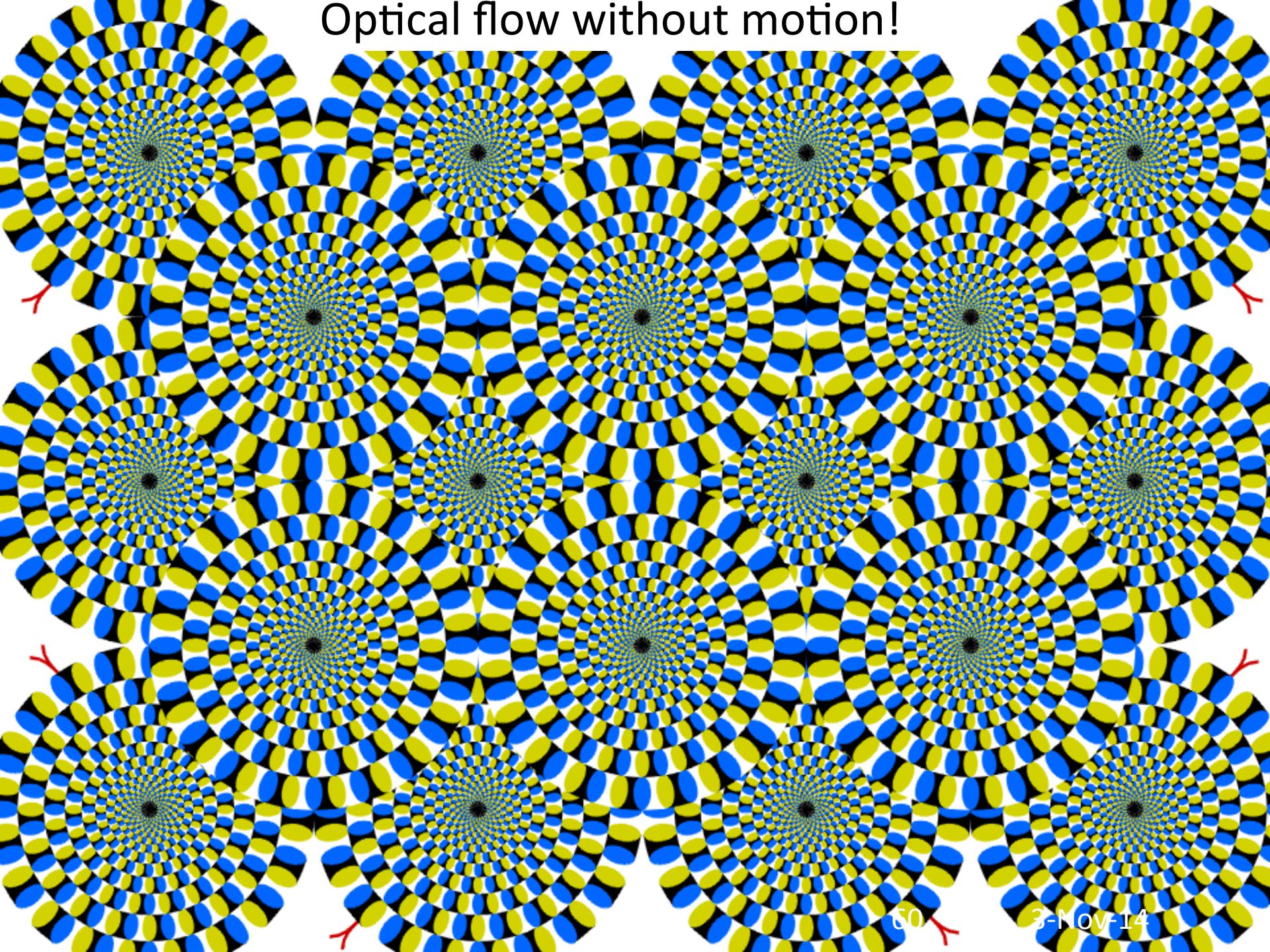


W. Choi & K. Shahid & S. Savarese WMC 2010

Source: Silvio Savarese



W. Choi, K. Shahid, S. Savarese, "What are they doing? : Collective Activity Classification Using Spatio-Temporal Relationship Among People", 9th International Workshop on Visual Surveillance (VWS09) in conjunction with ICCV 09



Optical flow without motion!

# What we have learned today?

- Introduction
- Optical flow
- Feature tracking
- Applications
- (Supplementary) Technical note

**Reading:** [Szeliski] Chapters: 8.4, 8.5

[Fleet & Weiss, 2005]

<http://www.cs.toronto.edu/pub/jepson/teaching/vision/2503/opticalFlow.pdf>

# Optical Flow Estimation

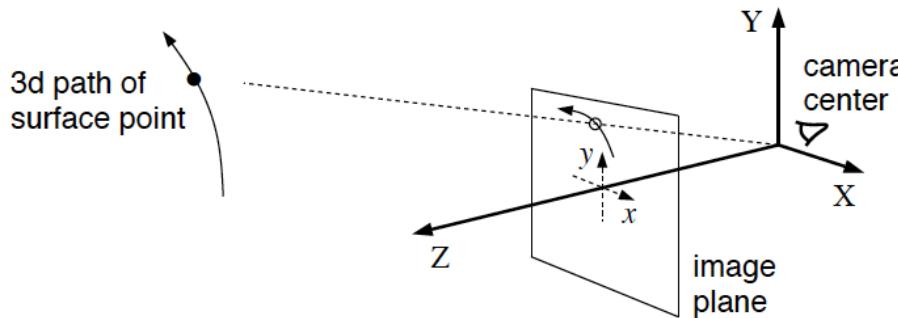
technical  
note

**Goal:** Introduction to image motion and 2D optical flow estimation.

## Motivation:

- Motion is a rich source of information about the world:
  - segmentation
  - surface structure from parallax
  - self-motion
  - recognition
  - understanding behavior
  - understanding scene dynamics
- Other correspondence / registration problems:
  - stereo disparity (short and wide baseline)
  - computer-assisted surgery
  - multiview alignment for mosaicing or stop-frame animation

# Introduction to Image Motion



A 3D point  $\vec{X}$  follows a space-time path  $\vec{X}(t)$ . Its velocity is  $\vec{V}$  is

$$\vec{V} = \frac{d\vec{X}(t)}{dt} = \left( \frac{dX(t)}{dt}, \frac{dY(t)}{dt}, \frac{dZ(t)}{dt} \right)^T.$$

Perspective projection (for nodal length  $f$ ) of the 3D path onto the image plane produces a 2D path,

$$\vec{x}(t) = (x(t), y(t)) = \left( \frac{fX(t)}{Z(t)}, \frac{fY(t)}{Z(t)} \right),$$

the instantaneous 2D velocity of which is

$$\begin{aligned} \vec{u} &= \left( \frac{dx(t)}{dt}, \frac{dy(t)}{dt} \right)^T \\ &= \frac{f}{Z(t)} \left( \frac{dX(t)}{dt}, \frac{dY(t)}{dt} \right)^T - \frac{f}{Z^2(t)} \frac{dZ(t)}{dt} (X(t), Y(t))^T. \end{aligned}$$

# Optical Flow

technical note

## Two Key Problems:

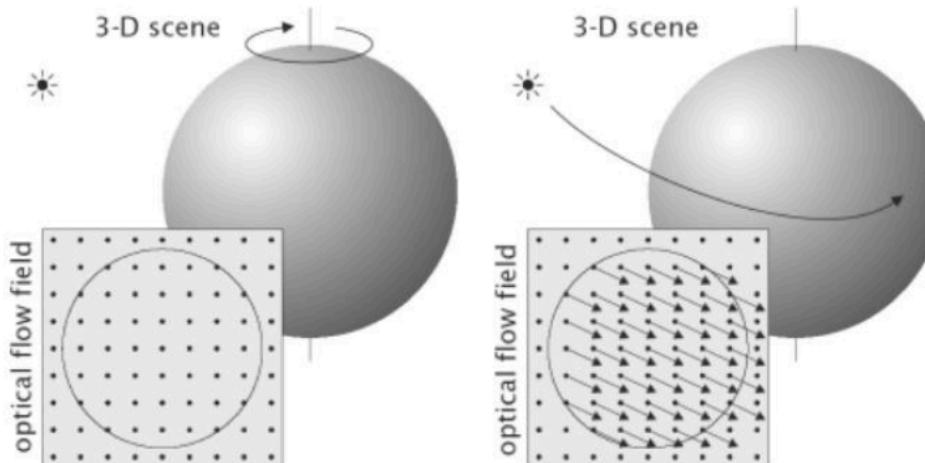
1. Determine what image property to track.
2. Determine how to track it

**Brightness constancy:** More precisely, let's track points of constant brightness, assuming that surface radiance is constant over time:

$$I(x, y, t + 1) = I(x - u_1, y - u_2, t).$$

Brightness constancy is often assumed by researchers, and often violated by Mother Nature; so the resulting optical flow field is sometimes a very poor approximation to the 2D motion field.

For example, a rotating Lambertian sphere with a static light source produces a static image. But a stationary sphere with a moving light source produces drifting intensities (figure from Jahne et al, 1999).



# Gradient-Based Motion Estimation

Let  $f(x)$  denote a 1D greylevel function of spatial position, and assume  $f(x)$  is just translated by  $d$  between time 1 and time 2:

$$f_2(x) = f_1(x - d).$$

We can express the shifted signal as a Taylor expansion of  $f_1$  about  $x$ :

$$f_1(x - d) = f_1(x) - d f'_1(x) + O(d^2 f''_1),$$

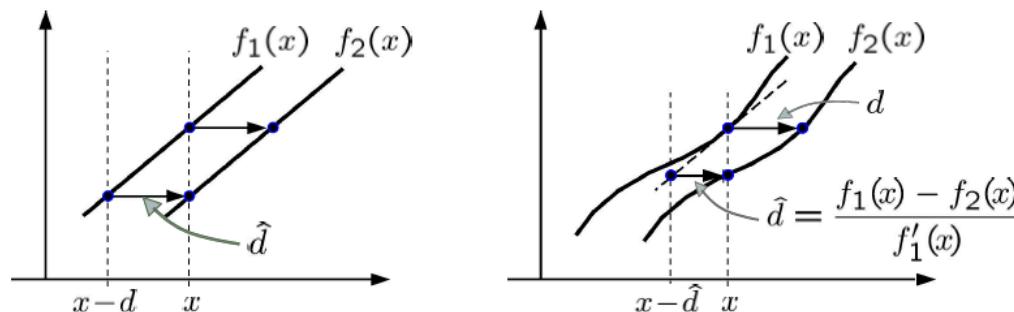
in which case the difference between the two signals is given by

$$f_2(x) - f_1(x) = -d f'_1(x) + O(d^2 f''_1).$$

Then, a first-order approximation to the displacement is

$$\hat{d} = \frac{f_1(x) - f_2(x)}{f'_1(x)}.$$

For linear signals the first-order estimate is exact.



## Gradient Constraint Equation

In two spatial dimensions, the first-order approximation is:

$$\begin{aligned} f(x + u_1, y + u_2, t + 1) \approx \\ f(x, y, t) + u_1 f_x(x, y, t) + u_2 f_y(x, y, t) + f_t(x, y, t) \quad (1) \end{aligned}$$

Taking the image difference with times  $t$  and  $t + 1$ , then yields:

$$u_1 f_x(x, y, t) + u_2 f_y(x, y, t) + f_t(x, y, t) = 0.$$

Written in vector form, with  $\vec{\nabla} f \equiv (f_x, f_y)^T$ :

$$\vec{u}^T \vec{\nabla} f(x, y, t) + f_t(x, y, t) = 0.$$

When the duration between frames is large, it is sometimes more appropriate to use only spatial derivatives in the Taylor series approximation in (1). Then one obtains a different approximation

$$\vec{u}^T \vec{\nabla} f(x, y, t) + \Delta f(x, y, t) = 0$$

where  $\Delta f(x, y, t) = f(x, y, t + 1) - f(x, y, t)$ .

# Brightness Conservation

One can also derive the gradient constraint equation directly from brightness conservation.

Let  $(x(t), y(t))$  denote a space-time path along which the image intensity remains constant; i.e., the time-varying image  $f$  satisfies

$$f(x(t), y(t), t) = c$$

Taking the total derivative of both sides gives

$$\frac{d}{dt} f(x(t), y(t), t) = 0$$

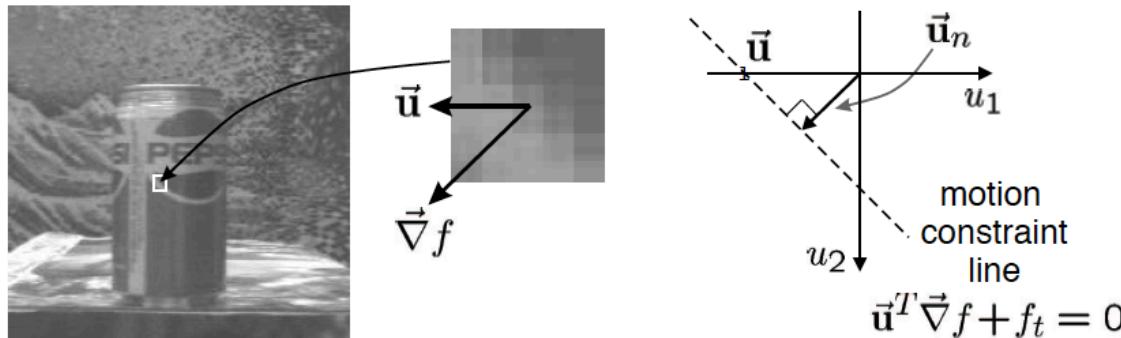
The total derivative is given in terms of its partial derivatives

$$\begin{aligned}\frac{d}{dt} f(x(t), y(t), t) &= \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} + \frac{\partial f}{\partial t} \frac{dt}{dt} \\ &= f_x u_1 + f_y u_2 + f_t \\ &= \bar{\mathbf{u}}^T \vec{\nabla} f + f_t \\ &= 0\end{aligned}$$

This derivation assumes that there is no aliasing, so that one can in principle reconstruct the continuous underlying signal and its derivatives in space and time. There are many situations in which this is a reasonable assumption. But with common video cameras temporal aliasing is often a problem with many video sequences, as we discuss later in these notes.

## Normal Velocity

The gradient constraint provides one constraint in two unknowns. It defines a line in velocity space:



The gradient constrains the velocity in the direction normal to the local image orientation, but does not constrain the tangential velocity. That is, it uniquely determines only the normal velocity:

$$\vec{u}_n = \frac{-f_t}{||\vec{\nabla}f||} \frac{\vec{\nabla}f}{||\vec{\nabla}f||}$$

When the gradient magnitude is zero, we get no constraint!

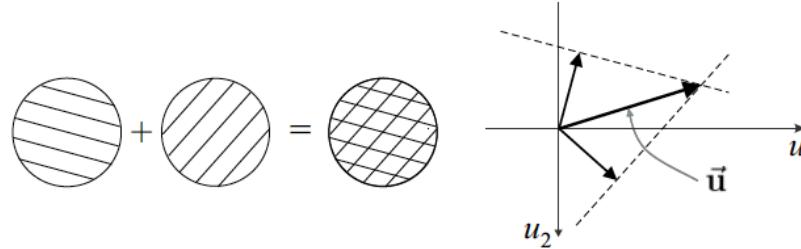
In any case, further constraints are required to estimate both elements of the 2D velocity  $\vec{u} = (u_1, u_2)^T$ .

# Area-Based Regression

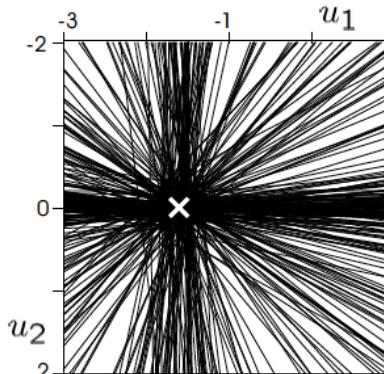
**Smoothness Assumption:** 2D motion is smooth in the local image neighborhood of the point at which we are estimating optical flow.

E.g., Let's say we have gradient constraints at two adjacent pixels,  $\vec{x}_1$  and  $\vec{x}_2$ , which have the same velocity,  $\vec{u} = (u_1, u_2)^T$ :

$$\begin{bmatrix} f_x(x_1, y_1, t) & f_y(x_1, y_1, t) \\ f_x(x_2, y_2, t) & f_y(x_2, y_2, t) \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \end{pmatrix} + \begin{bmatrix} f_t(x_1, y_1, t) \\ f_t(x_2, y_2, t) \end{bmatrix} = \vec{0},$$



More generally, we can use many constraints from within a local region about the point at which we wish to estimate the optical flow. Here is an example from the Pepsi sequence:



## Area-Based Regression (Linear LS)

technical note

We will not satisfy all constraints exactly. Rather we seek the velocity that minimizes the squared error in each constraint (called the least-squares (LS) velocity estimate):

$$E(u_1, u_2) = \sum_{x,y} g(x, y) [u_1 f_x(x, y, t) + u_2 f_y(x, y, t) + f_t(x, y, t)]^2$$

where  $g(x, y)$  is a low-pass window that helps give more weight to constraints at the center of the region.

**Solution:** Differentiate  $E$  with respect to  $(u_1, u_2)$  and set to zero:

$$\frac{\partial E(u_1, u_2)}{\partial u_1} = \sum_{xy} g(x, y) [u_1 f_x^2 + u_2 f_x f_y + f_x f_t] = 0$$

$$\frac{\partial E(u_1, u_2)}{\partial u_2} = \sum_{xy} g(x, y) [u_2 f_y^2 + u_1 f_x f_y + f_y f_t] = 0$$

The constraints thus yield two linear equations for  $u_1$  and  $u_2$ .

In matrix notation, these *normal equations* and their solution are

$$\mathbf{M} \vec{\mathbf{u}} + \vec{\mathbf{b}} = \vec{0} , \quad \hat{\mathbf{u}} = -\mathbf{M}^{-1} \vec{\mathbf{b}}$$

(assuming  $\mathbf{M}^{-1}$  exists), where

$$\mathbf{M} = \sum g \begin{pmatrix} f_x \\ f_y \end{pmatrix} (f_x, f_y) = \begin{bmatrix} \sum g f_x^2 & \sum g f_x f_y \\ \sum g f_x f_y & \sum g f_y^2 \end{bmatrix}$$

$$\vec{\mathbf{b}} = \sum g f_t \begin{pmatrix} f_x \\ f_y \end{pmatrix} = \begin{pmatrix} \sum g f_x f_t \\ \sum g f_y f_t \end{pmatrix}$$

## Area-Based Regression (Implementation)

Since we want the optical flow at each pixel, we can compute the components of the normal equations with a set of image operators. This is much preferred to looping over image pixels.

1. First, compute 3 image gradient images at time t, corresponding to the gradient measurements:

$$f_x(\vec{x}), \quad f_y(\vec{x}), \quad f_t(\vec{x})$$

2. Point-wise, we compute the quadratic functions of the derivative images. This produces five images equal to:

$$f_x^2(\vec{x}), \quad f_y^2(\vec{x}), \quad f_x(\vec{x})f_y(\vec{x}), \quad f_t(\vec{x})f_y(\vec{x}), \quad f_t(\vec{x})f_x(\vec{x})$$

3. Blurring the quadratic images corresponds to the accumulating of local constraints under the spatial (or spatiotemporal) support window  $g$  above. This produces five images, each of which contains an element of the normal equations at a specific image locations:

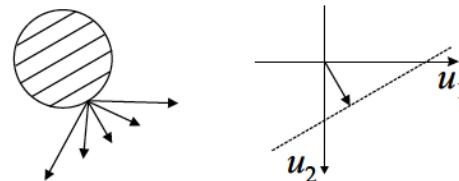
$$\begin{aligned} & g(\vec{x}) * f_x^2(\vec{x}), \quad g(\vec{x}) * f_y^2(\vec{x}), \\ & g(\vec{x}) * [f_x(\vec{x}) f_y(\vec{x})], \quad g(\vec{x}) * [f_t(\vec{x}) f_y(\vec{x})], \quad g(\vec{x}) * [f_t(\vec{x}) f_x(\vec{x})] \end{aligned}$$

4. Compute the two images containing the components of optical flow at each pixel. This is given in closed form since the inverse of the normal matrix (i.e.,  $M$  above) is easily expressed in closed form.

# Aperture Problem

technical note

Even with the collection of constraints from within a region, the estimation will be undetermined when the matrix  $M$  is singular:



When all image gradients are parallel, then the normal matrix for the least-squares solution becomes singular (rank 1). E.g., for gradients  $m(x, y)\vec{n}$ , where  $m(x, y)$  is the gradient magnitude at pixel  $(x, y)$

$$M = \left( \sum_{x,y} g(x, y) m^2(x, y) \right) \vec{n} \vec{n}^T$$

**Aperture Problem:** For sufficiently small spatial apertures the normal matrix,  $M$ , will be singular (rank deficient).

**Generalized Aperture Problem:** For small apertures  $M$  becomes singular, but for sufficiently large apertures the 2D motion field may deviate significantly from the assumed motion model.

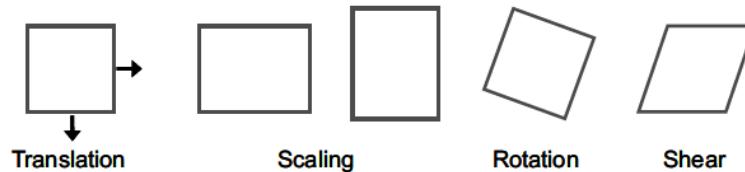
**Heuristic Diagnostics:** Use singular values,  $s_1 \geq s_2 \geq s_3$ , of  $\sum_{x,y} g \vec{h} \vec{h}^T$ , for  $\vec{h} = (f_x, f_y, f_t)^T$ , to assess the rank of the local image structure:

- if  $s_2$  is "too small" then only normal velocity is available
- if  $s_3$  is "large" then the gradient constraints do not lie in a plane so a single velocity will not fit the data well.

Such "bounds" will depend on the LS region size and the noise.

# Higher-Order Motion Models

Constant flow model within an image neighborhood is often a poor model, especially when the regions become larger. Affine models often provide a more suitable model of local image deformation.



Affine flow for an image region centered at location  $\vec{x}_0$  is given by

$$\vec{u}(\vec{x}) = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} (\vec{x} - \vec{x}_0) + \begin{pmatrix} a_5 \\ a_6 \end{pmatrix} = \mathbf{A}(\vec{x}; \vec{x}_0) \vec{a}$$

where  $\vec{a} = (a_1, a_2, \dots, a_6)^T$  and

$$\mathbf{A}(\vec{x}; \vec{x}_0) = \begin{bmatrix} x - x_0 & y - y_0 & 0 & 0 & 1 & 0 \\ 0 & 0 & x - x_0 & y - y_0 & 0 & 1 \end{bmatrix}$$

Then the gradient constraint becomes:

$$\begin{aligned} 0 &= \vec{u}(x, y)^T \vec{\nabla} f(x, y, t) + f_t(x, y, t) \\ &= \vec{a}^T \mathbf{A}(x, y)^T \vec{\nabla} f(x, y, t) + f_t(x, y, t), \end{aligned}$$

so, as above, the weighted least-squares solution for  $\vec{a}$  is given by:

$$\hat{\vec{a}} = \mathbf{M}^{-1} \vec{b}$$

where, for weighting with the spatial window  $g$ ,

$$\mathbf{M} = \sum_{x,y} g(\vec{x}) \mathbf{A}^T \vec{\nabla} f \vec{\nabla} f^T \mathbf{A}, \quad \vec{b} = - \sum_{x,y} g(\vec{x}) \mathbf{A}^T \vec{\nabla} f f_t$$