

Feature Tracking and Optical Flow

Computer Vision

CS 143, Brown

James Hays

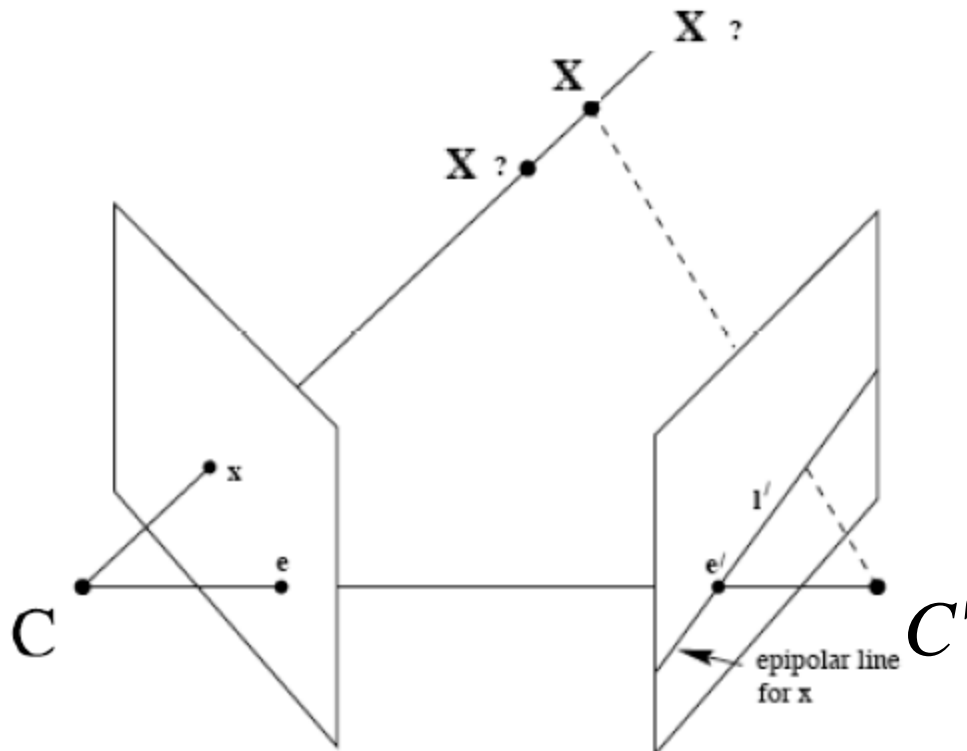
Many slides adapted from Derek Hoesim, Lana Lazebnik, Silvio Savese, who in turn adapted slides from Steve Seitz, Rick Szeliski, Martial Hebert, Mark Pollefeys, and others

Recap: Stereo

- Input: two views
- Output: dense depth measurements
- Outline:
 - If views are not rectified, estimate the epipolar geometry
 - Use a stereo correspondence algorithm to match patches across images while respecting the epipolar constraints
 - Depth is a function of disparity
 - $z = (\text{baseline} * f) / (X - X')$

Recap: Epipoles

- Point x in left image corresponds to **epipolar line l'** in right image
- Epipolar line passes through the epipole (the intersection of the cameras' baseline with the image plane)



Recap: Fundamental Matrix

- Fundamental matrix maps from a point in one image to a line in the other

$$\mathbf{l}' = \mathbf{F}\mathbf{x} \quad \mathbf{l} = \mathbf{F}^\top \mathbf{x}'$$

- If \mathbf{x} and \mathbf{x}' correspond to the same 3d point \mathbf{X} :

$$\mathbf{x}'^\top \mathbf{F}\mathbf{x} = 0$$

- \mathbf{F} can be estimated from 8 $(\mathbf{x}, \mathbf{x}')$ matches by simply setting up a system of linear equations.

8-point algorithm

1. Solve a system of homogeneous linear equations

a. Write down the system of equations

$$\mathbf{x}^T F \mathbf{x}' = 0$$

$$x'x f_{11} + x'y f_{12} + x' f_{13} + y'x f_{21} + y'y f_{22} + y' f_{23} + x f_{31} + y f_{32} + f_{33} = 0$$

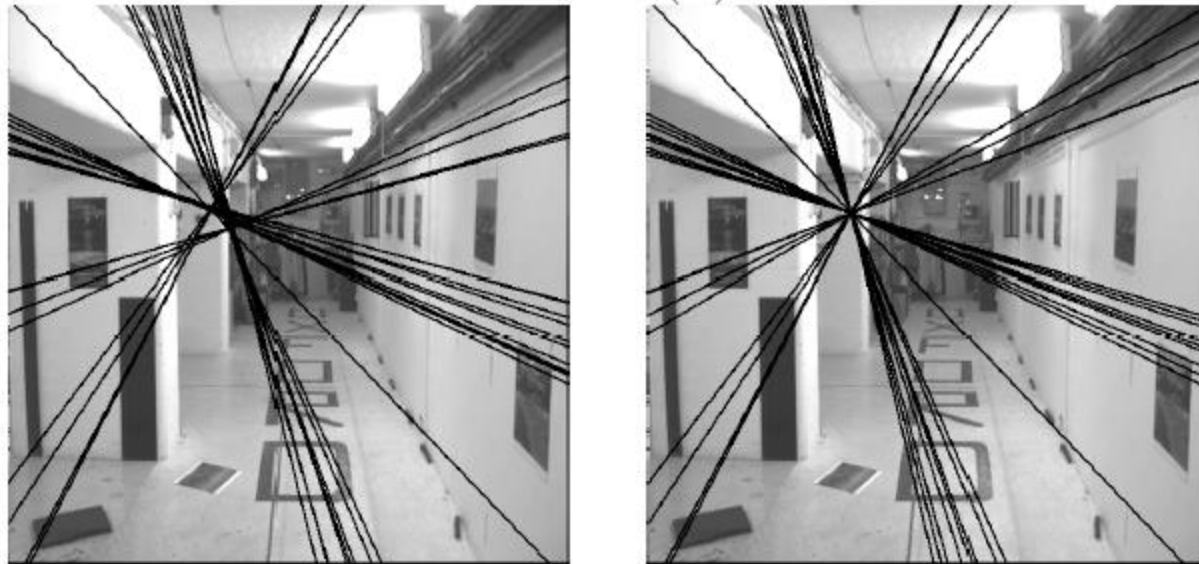
$$A\mathbf{f} = \begin{bmatrix} x'_1 x_1 & x'_1 y_1 & x'_1 & y'_1 x_1 & y'_1 y_1 & y'_1 & x_1 & y_1 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x'_n x_n & x'_n y_n & x'_n & y'_n x_n & y'_n y_n & y'_n & x_n & y_n & 1 \end{bmatrix} \begin{pmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \\ f_{33} \end{pmatrix} = \mathbf{0}$$

8-point algorithm

1. Solve a system of homogeneous linear equations
 - a. Write down the system of equations
 - b. Solve \mathbf{f} from $A\mathbf{f}=\mathbf{0}$

Need to enforce singularity constraint

Fundamental matrix has rank 2 : $\det(\mathbf{F}) = 0$.



Left : Uncorrected \mathbf{F} – epipolar lines are not coincident.

Right : Epipolar lines from corrected \mathbf{F} .

8-point algorithm

1. Solve a system of homogeneous linear equations
 - a. Write down the system of equations
 - b. Solve \mathbf{f} from $\mathbf{A}\mathbf{f}=\mathbf{0}$
2. Resolve $\det(\mathbf{F}) = 0$ constraint using SVD

Matlab:

```
[U, S, V] = svd(F);  
S(3,3) = 0;  
F = U*S*V';
```


Recap: Structure from Motion

- Input: Arbitrary number of uncalibrated views
- Output: Camera parameters and sparse 3d points

Photo synth

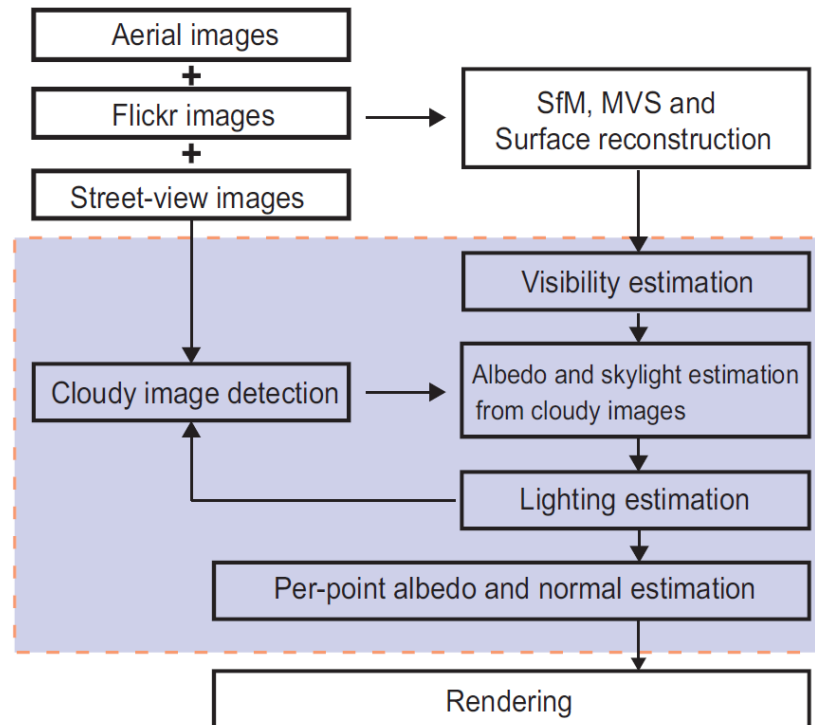
Noah Snavely, Steven M. Seitz, Richard Szeliski, "[Photo tourism: Exploring photo collections in 3D](#)," SIGGRAPH 2006



<http://photosynth.net/>

Can we do SfM and Stereo?

- Yes, numerous systems do.
- One example of the state of the art:



<http://www.youtube.com/watch?v=NdeD4cjLI0c>

This class: recovering motion

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

Two problems, one registration method

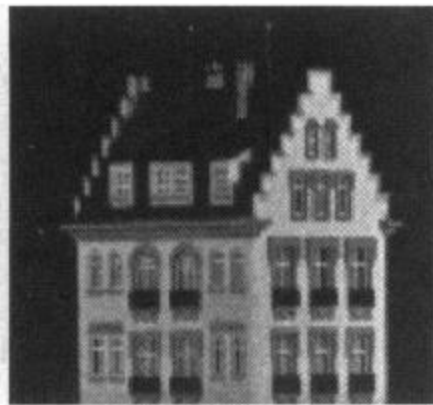
B. Lucas and T. Kanade. [An iterative image registration technique with an application to stereo vision.](#) In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

Feature tracking

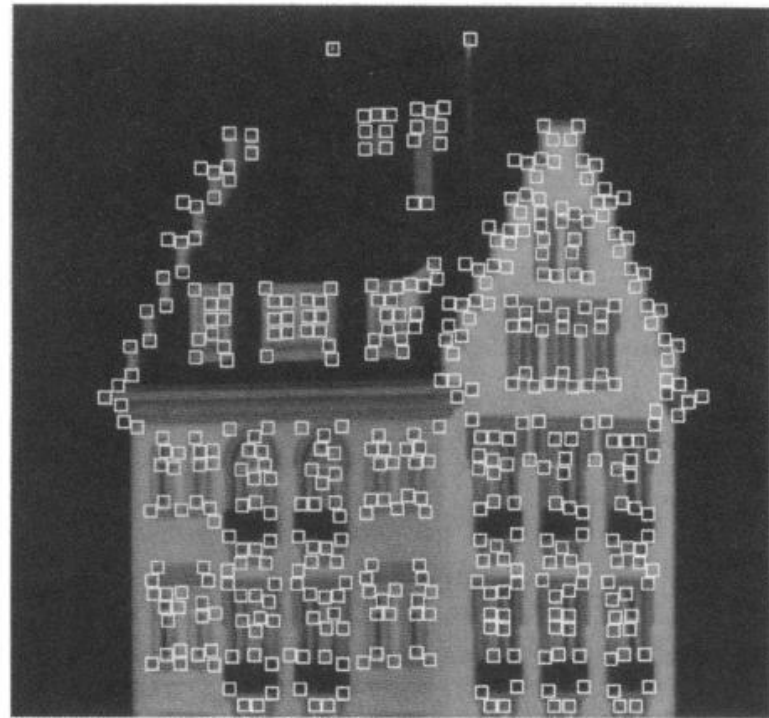
- Many problems, such as structure from motion require matching points
- If motion is small, tracking is an easy way to get them



60



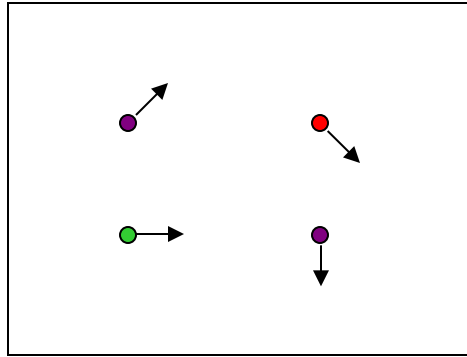
150



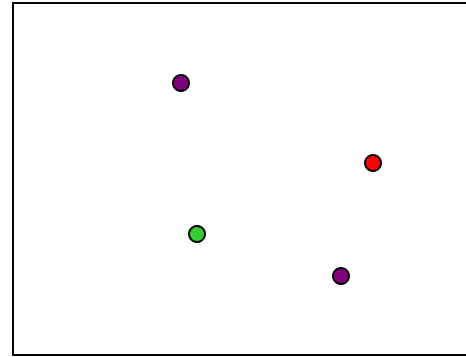
Feature tracking

- Challenges
 - Figure out which features can be tracked
 - Efficiently track across frames
 - Some points may change appearance over time (e.g., due to rotation, moving into shadows, etc.)
 - Drift: small errors can accumulate as appearance model is updated
 - Points may appear or disappear: need to be able to add/delete tracked points

Feature tracking



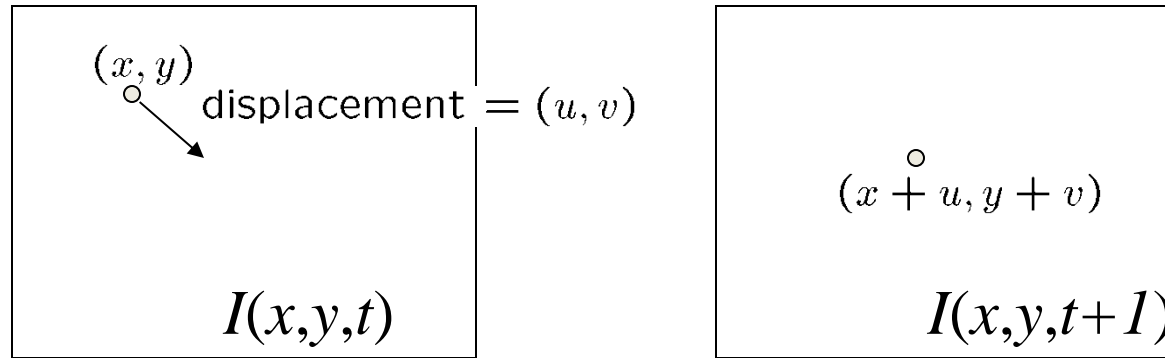
$I(x,y,t)$



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

- Given two subsequent frames, estimate the point translation
- Key assumptions of Lucas-Kanade Tracker
 - **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

The brightness constancy constraint



- Brightness Constancy Equation:

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

Take Taylor expansion of $I(x+u, y+v, t+1)$ at (x, y, t) to linearize the right side:

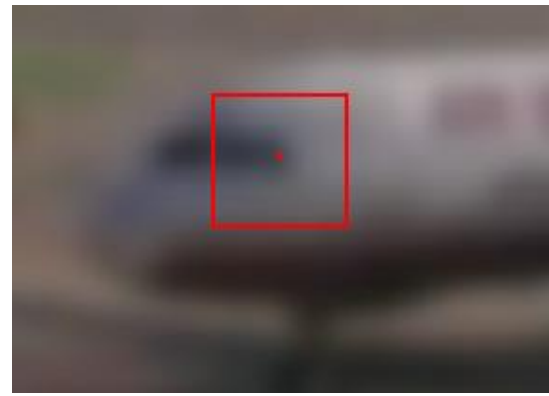
$$I(x + u, y + v, t + 1) \approx I(x, y, t) + \overset{\text{Image derivative along x}}{I_x} \cdot u + I_y \cdot v + \overset{\text{Difference over frames}}{I_t}$$
$$I(x + u, y + v, t + 1) - I(x, y, t) = +I_x \cdot u + I_y \cdot v + I_t$$

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

How does this make sense?

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

- What do the static image gradients have to do with motion estimation?



The brightness constancy constraint

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

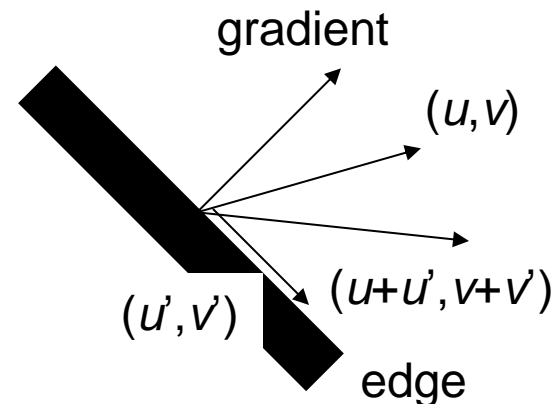
$$\nabla I \cdot [u \ v]^T + I_t = 0$$

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

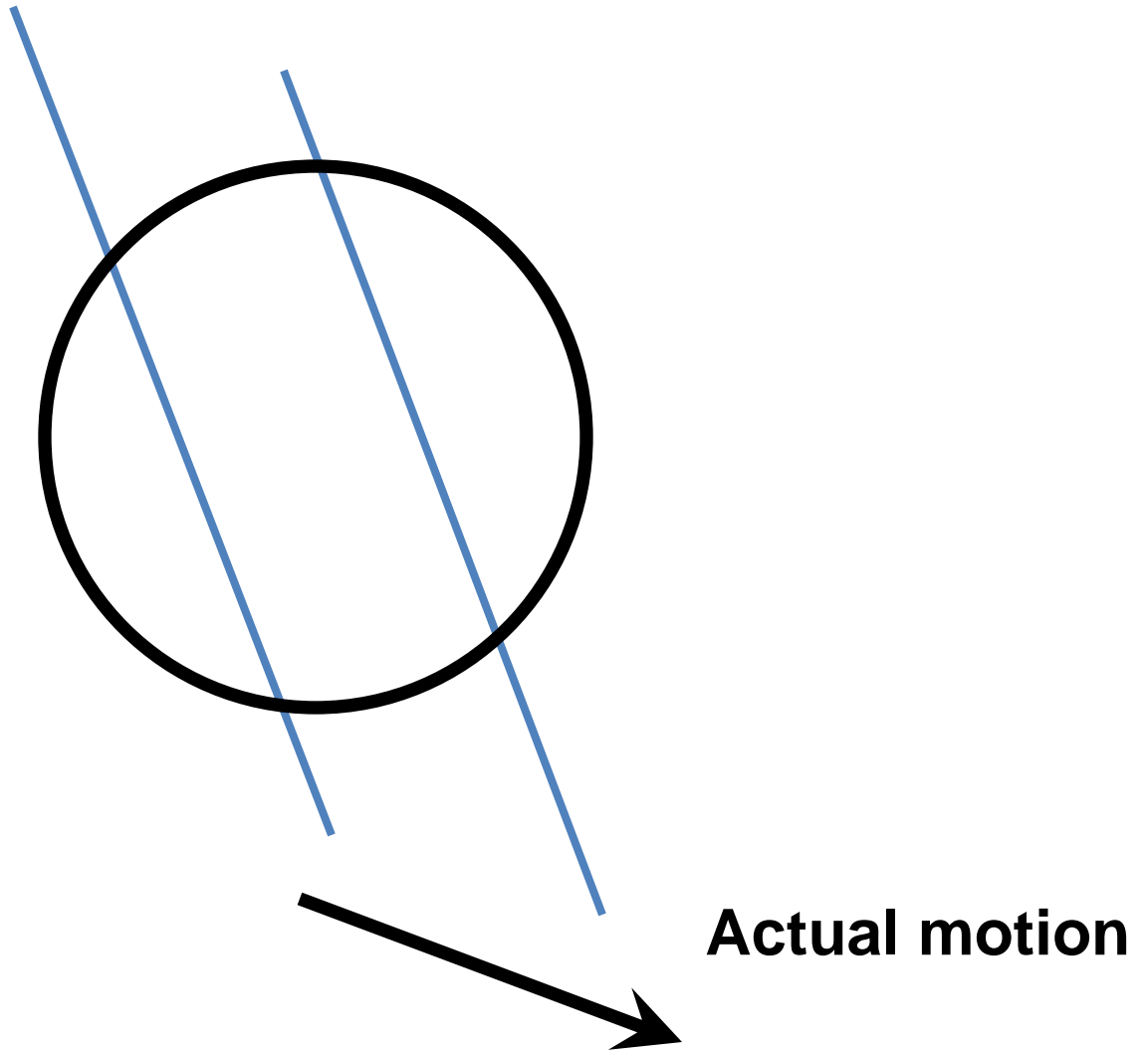
The component of the motion perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

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

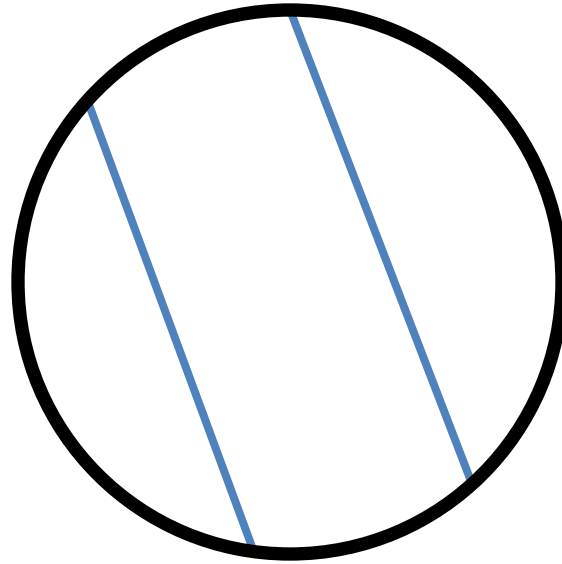
$$\nabla I \cdot [u' \ v']^T = 0$$



The aperture problem

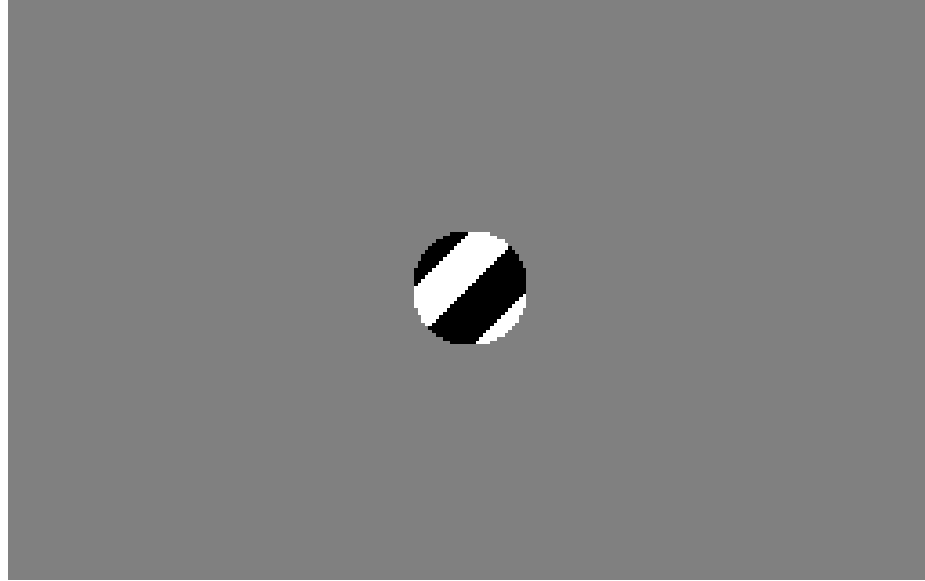


The aperture problem



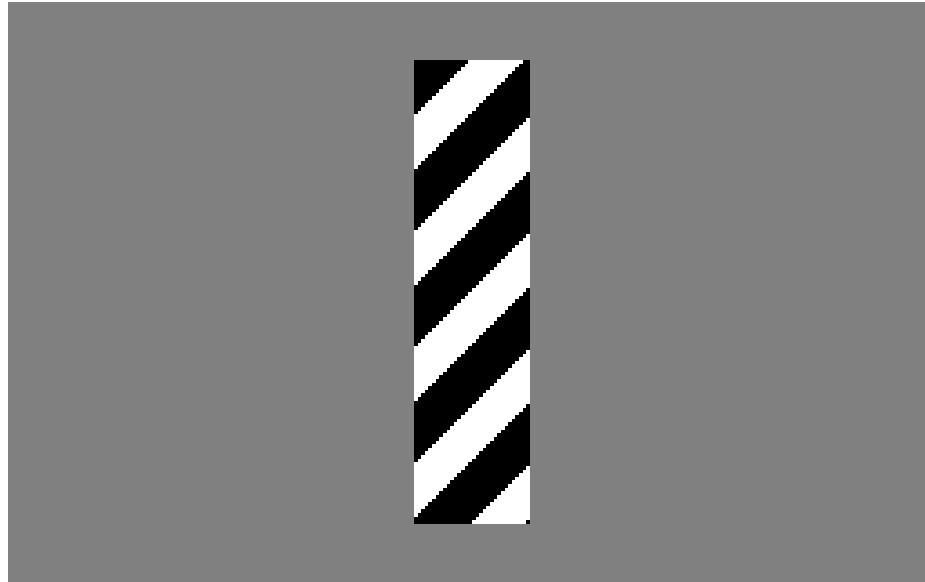
Perceived motion

The barber pole illusion



http://en.wikipedia.org/wiki/Barberpole_illusion

The barber pole illusion



http://en.wikipedia.org/wiki/Barberpole_illusion

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

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

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

Solving the ambiguity...

- Least squares problem:

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix} \quad \begin{matrix} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{matrix}$$

Matching patches across images

- Overconstrained linear system

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix} \quad \begin{matrix} A & d = b \\ 25 \times 2 & 2 \times 1 & 25 \times 1 \end{matrix}$$

Least squares solution for d given by $(A^T A) d = A^T b$

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

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

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? I.e., what are good points to track?

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

Does this remind you of anything?

Criteria for Harris corner detector

Low-texture region



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

- gradients have small magnitude
- small λ_1 , small λ_2

Edge



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

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

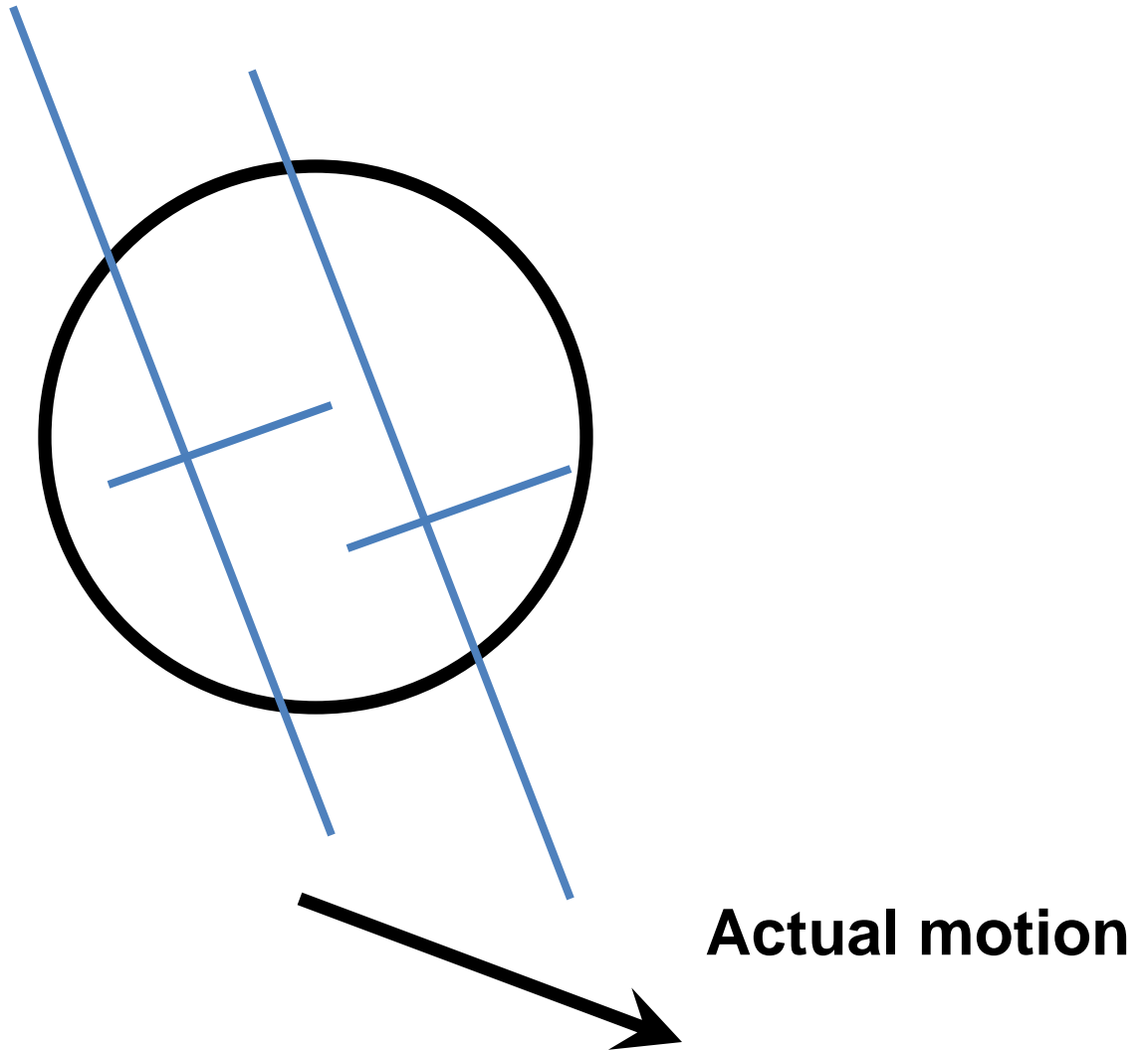
High-texture region



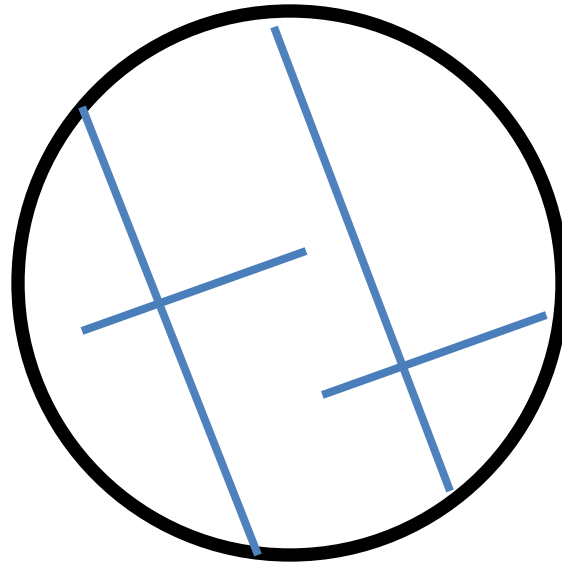
$$\sum \nabla I (\nabla I)^T$$

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

The aperture problem resolved



The aperture problem resolved



Perceived motion

Dealing with larger movements: Iterative refinement

1. Initialize $(x', y') = (x, y)$
2. Compute (u, v) by

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

Original (x, y) position

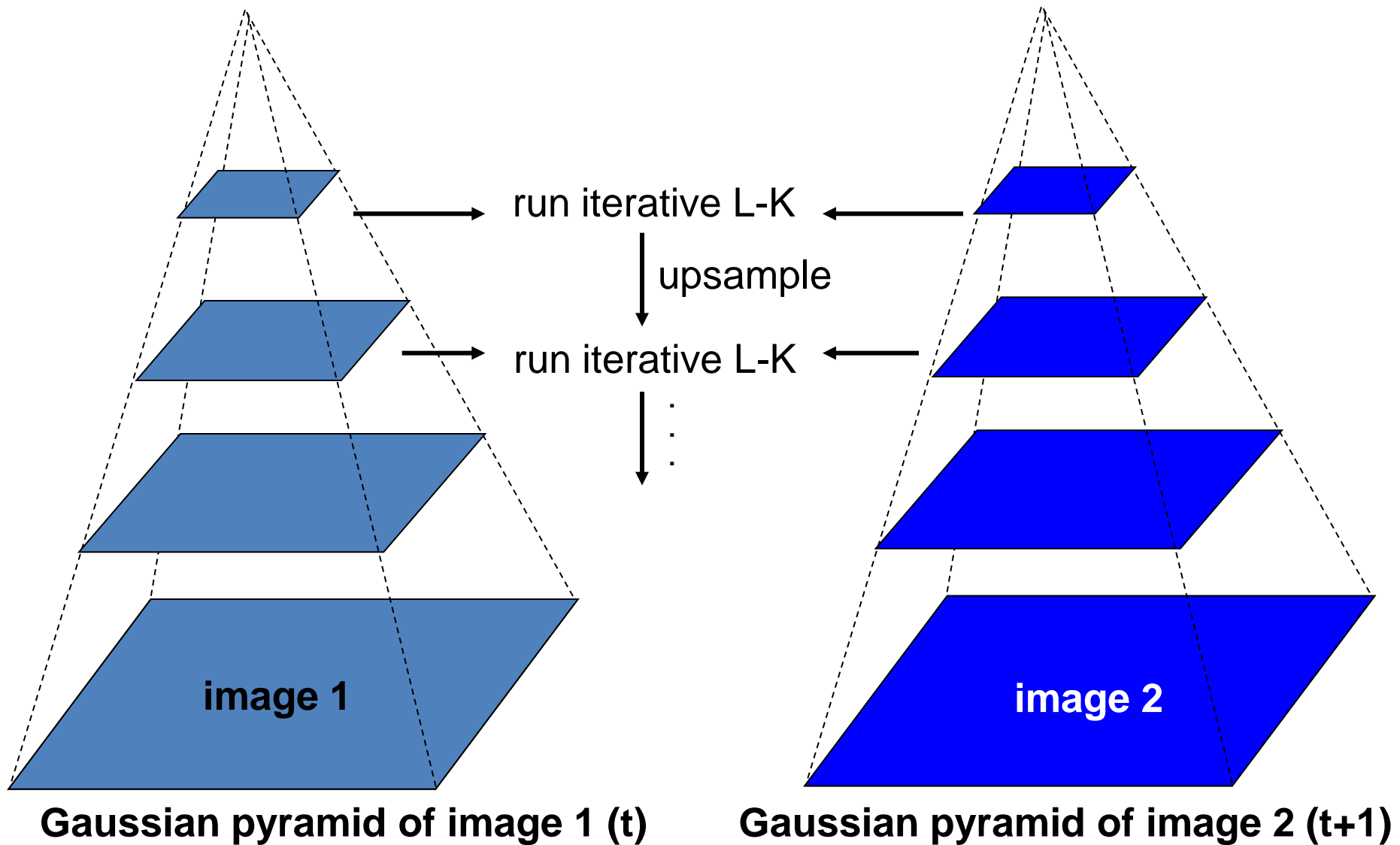
$I_t = I(x', y', t+1) - I(x, y, t)$

2nd moment matrix for feature patch in first image

displacement

3. Shift window by (u, v) : $x' = x' + u$; $y' = y' + v$;
4. Recalculate I_t
5. Repeat steps 2-4 until small change
 - Use interpolation for subpixel values

Dealing with larger movements: coarse-to-fine registration



Shi-Tomasi feature tracker

- Find good features using eigenvalues of second-moment matrix (e.g., Harris detector or threshold on the smallest eigenvalue)
 - Key idea: “good” features to track are the ones whose motion can be estimated reliably
- Track from frame to frame with Lucas-Kanade
 - This amounts to assuming a translation model for frame-to-frame feature movement
- Check consistency of tracks by *affine* registration to the first observed instance of the feature
 - Affine model is more accurate for larger displacements
 - Comparing to the first frame helps to minimize drift

Tracking example

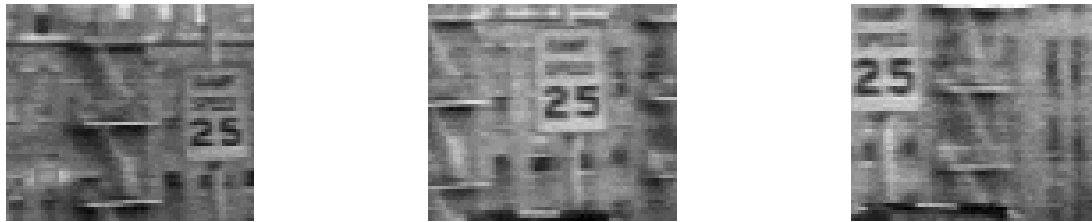


Figure 1: Three frame details from Woody Allen's *Manhattan*. The details are from the 1st, 11th, and 21st frames of a subsequence from the movie.

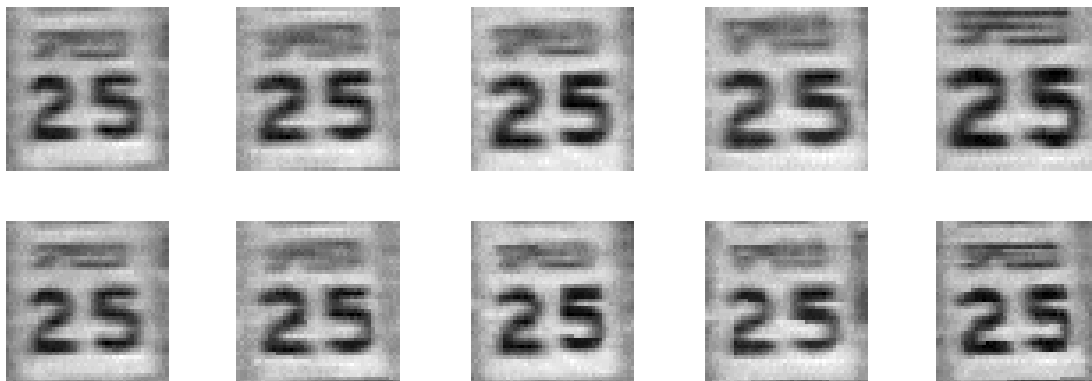


Figure 2: The traffic sign windows from frames 1,6,11,16,21 as tracked (top), and warped by the computed deformation matrices (bottom).

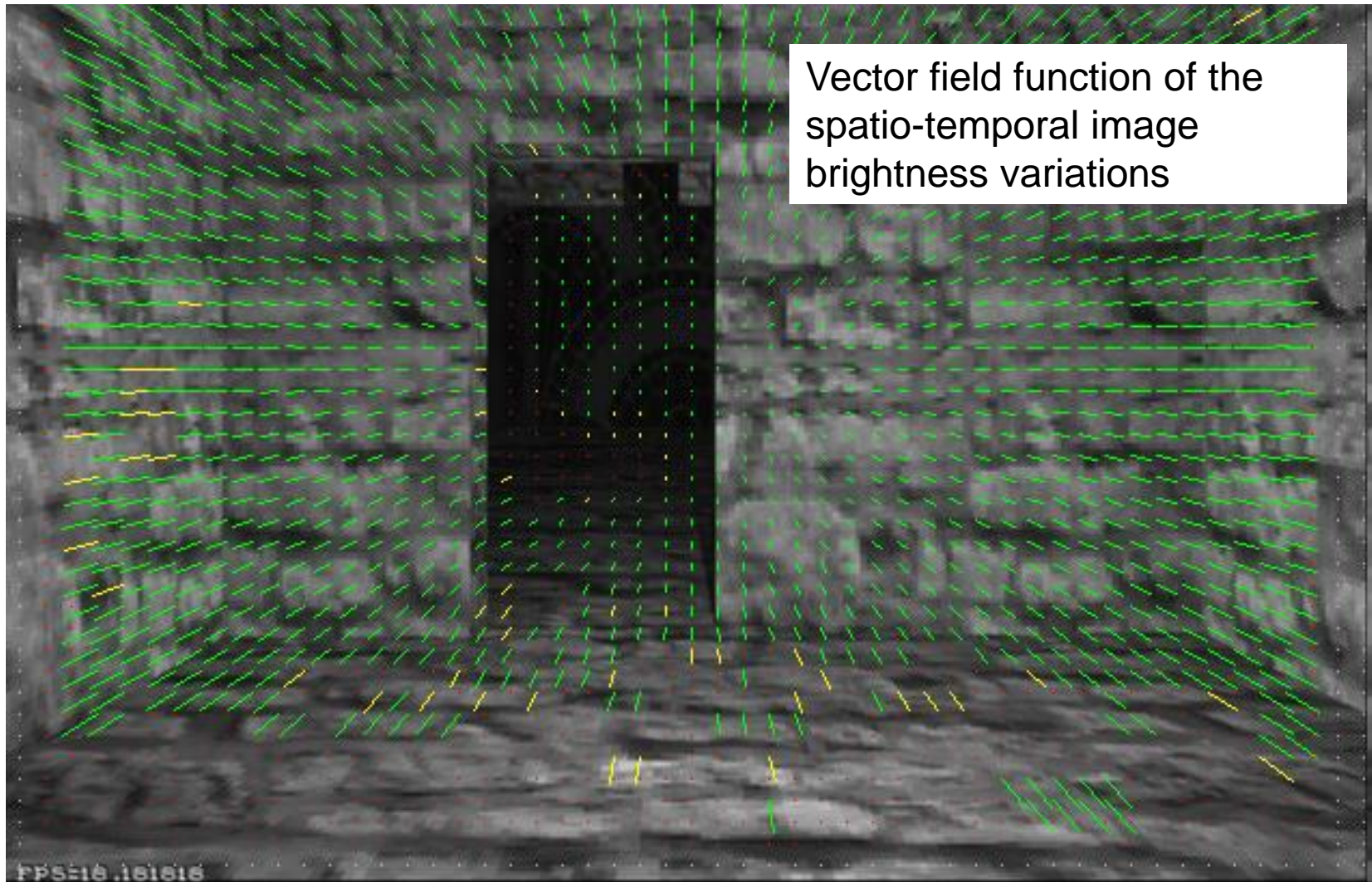
Summary of KLT tracking

- Find a good point to track (harris corner)
- Use intensity second moment matrix and difference across frames to find displacement
- Iterate and use coarse-to-fine search to deal with larger movements
- When creating long tracks, check appearance of registered patch against appearance of initial patch to find points that have drifted

Implementation issues

- Window size
 - Small window more sensitive to noise and may miss larger motions (without pyramid)
 - Large window more likely to cross an occlusion boundary (and it's slower)
 - 15x15 to 31x31 seems typical
- Weighting the window
 - Common to apply weights so that center matters more (e.g., with Gaussian)

Optical flow



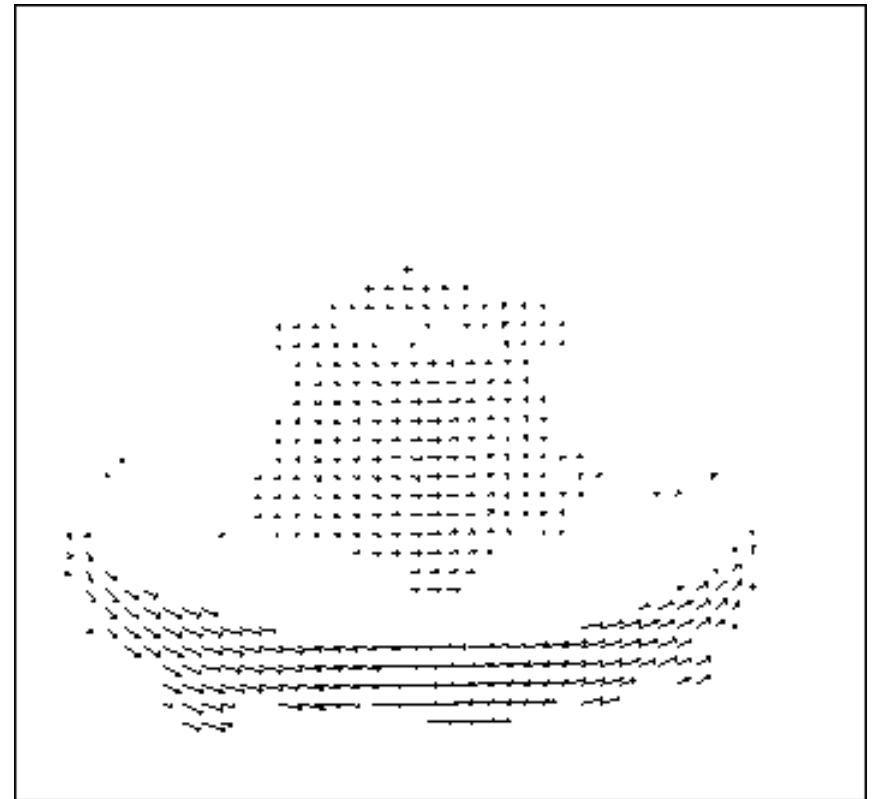
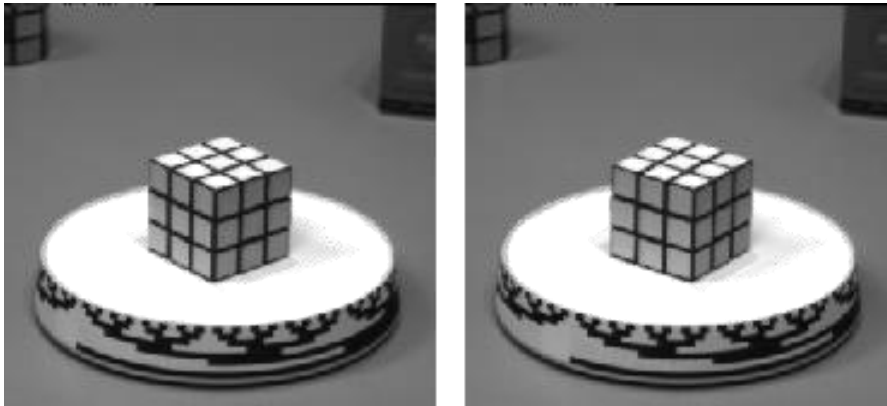
Picture courtesy of Selim Temizer - Learning and Intelligent Systems (LIS) Group, MIT

Uses of motion

- Estimating 3D structure
- Segmenting objects based on motion cues
- Learning and tracking 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



What would the motion field of a non-rotating ball moving towards the camera look like?

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
 - Think of a uniform rotating sphere under fixed lighting vs. a stationary sphere under moving illumination

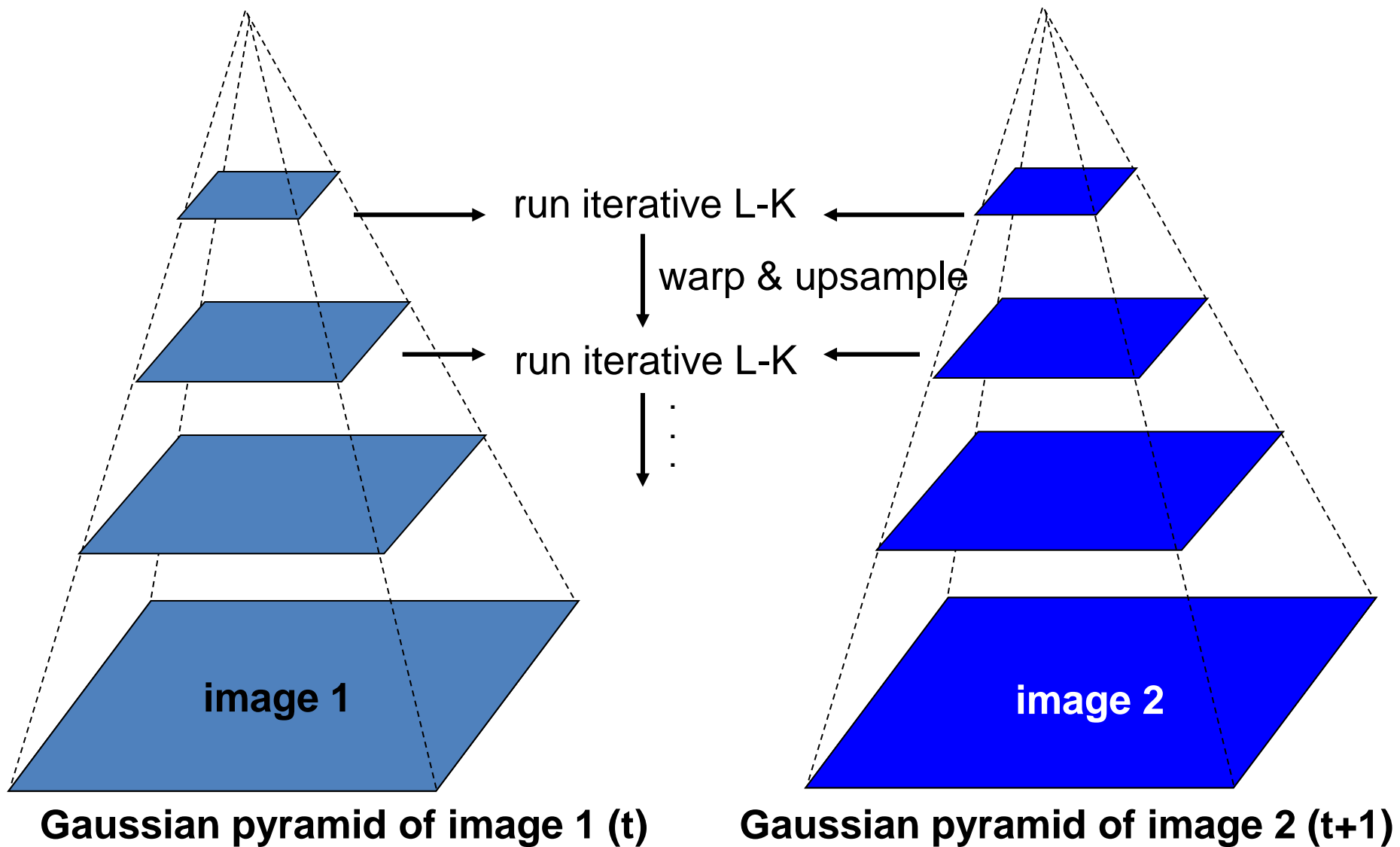
Lucas-Kanade Optical Flow

- Same as Lucas-Kanade feature tracking, but for each pixel
 - As we saw, works better for textured pixels
- Operations can be done one frame at a time, rather than pixel by pixel
 - Efficient

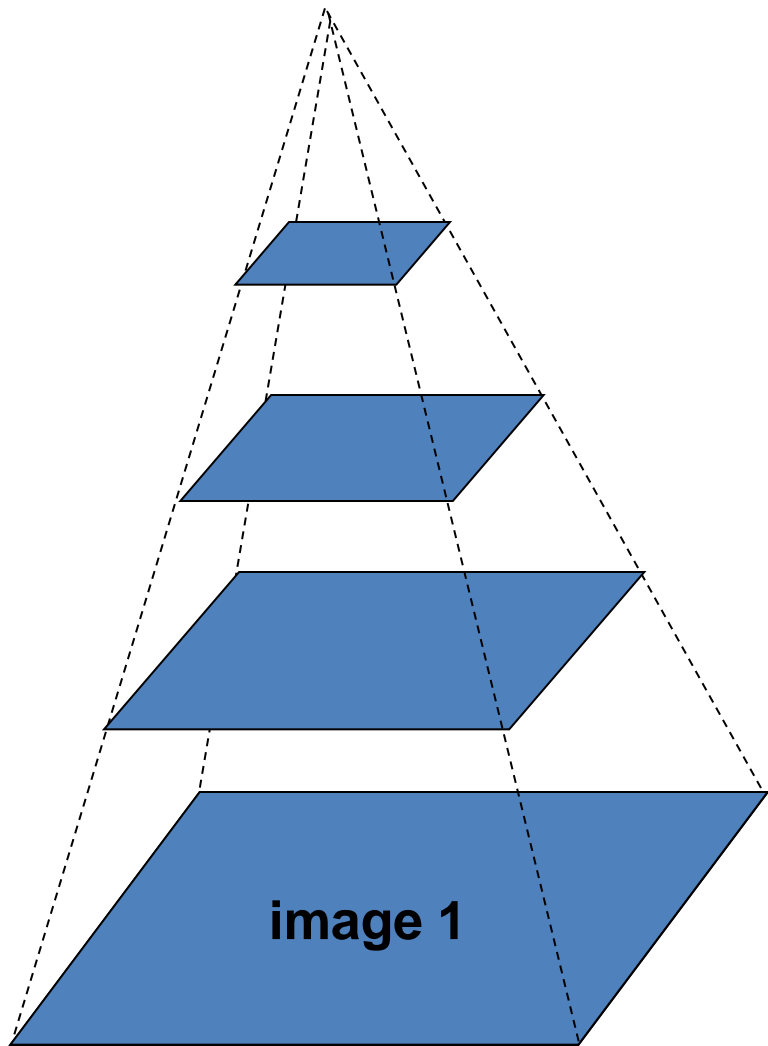
Iterative Refinement

- Iterative Lukas-Kanade Algorithm
 1. Estimate displacement at each pixel by solving Lucas-Kanade equations
 2. Warp $I(t)$ towards $I(t+1)$ using the estimated flow field
 - Basically, just interpolation
 3. Repeat until convergence

Coarse-to-fine optical flow estimation



Coarse-to-fine optical flow estimation



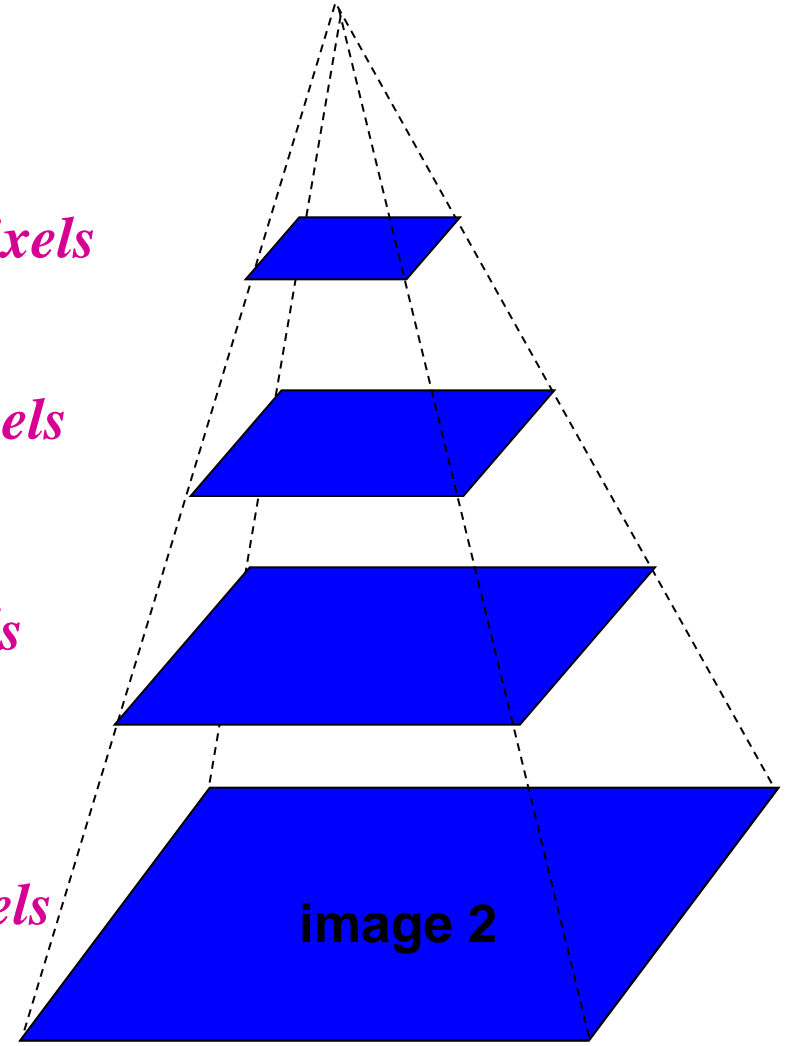
Gaussian pyramid of image 1

$u=1.25$ pixels

$u=2.5$ pixels

$u=5$ pixels

$u=10$ pixels

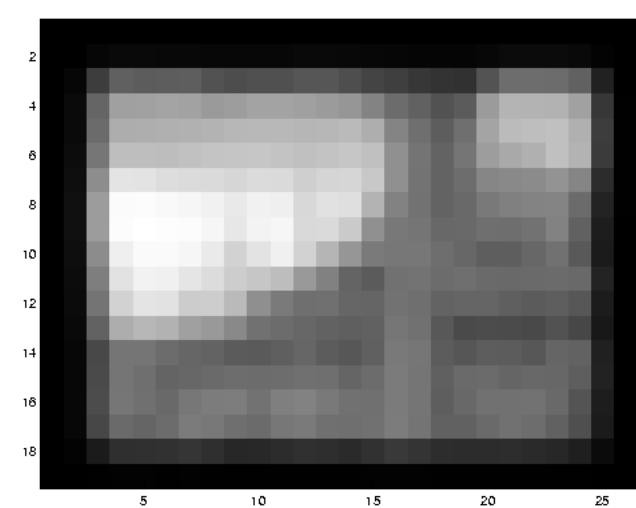
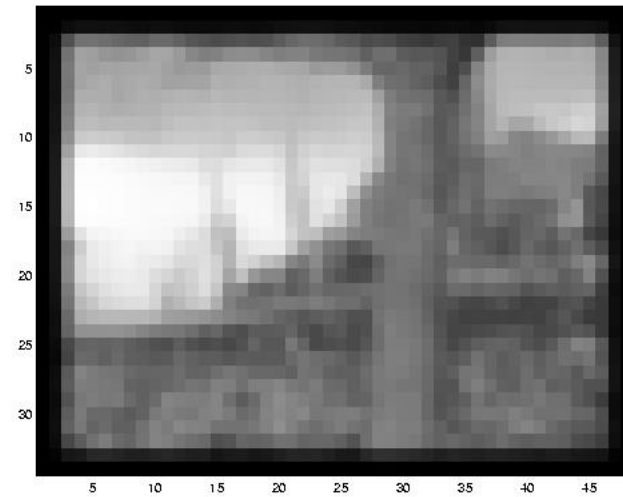
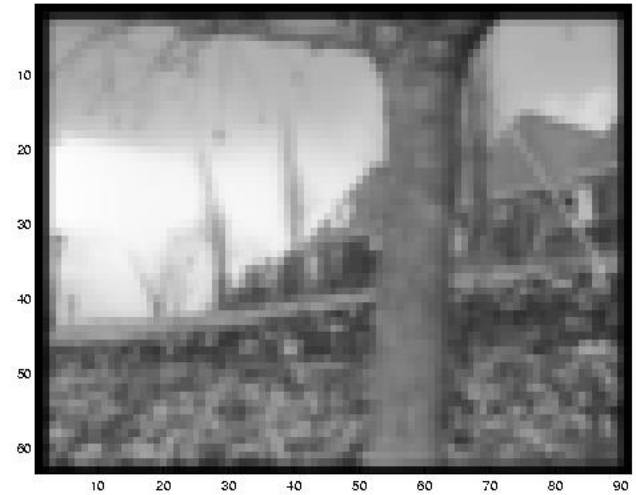


Gaussian pyramid of image 2

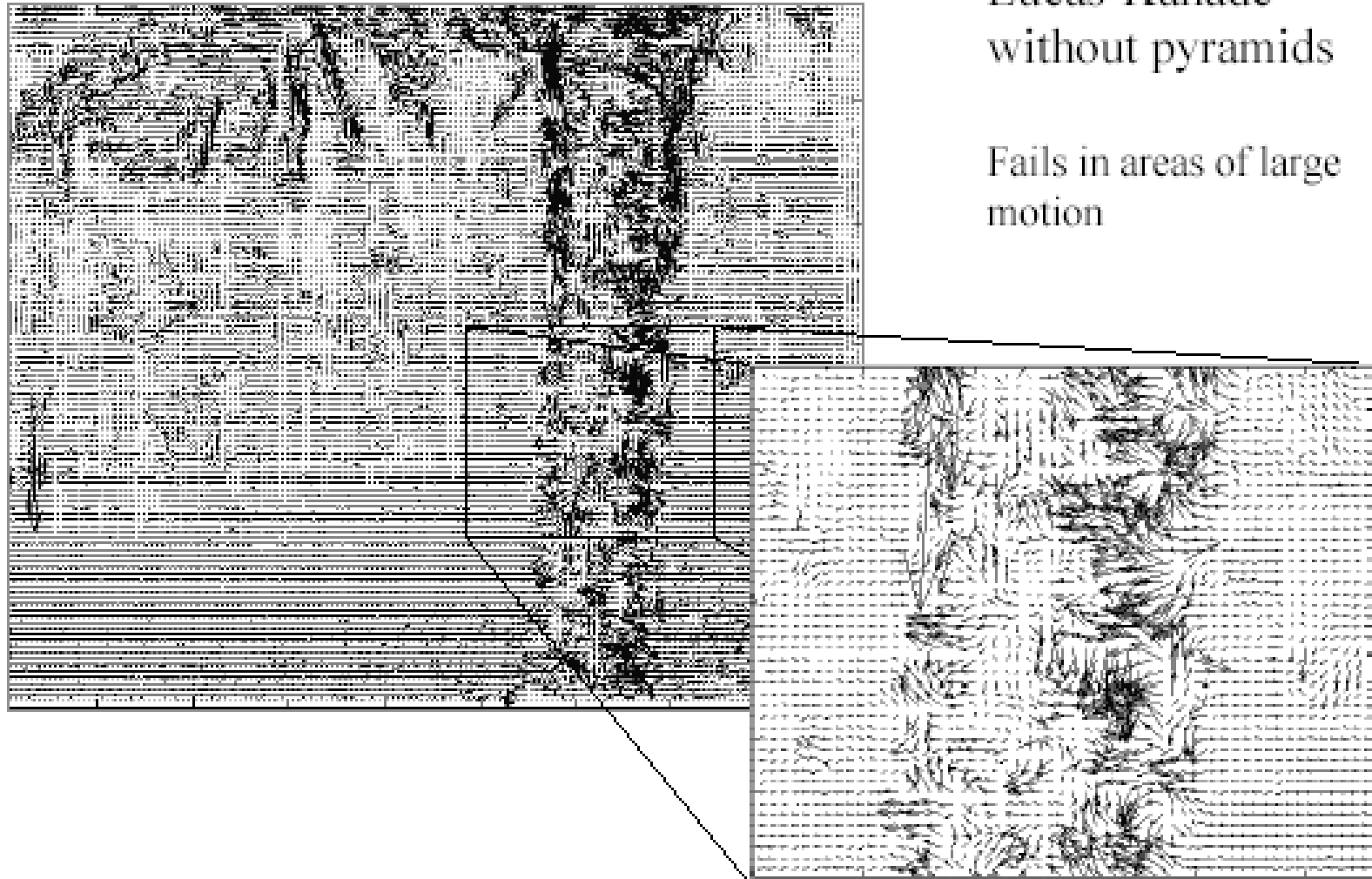
Example



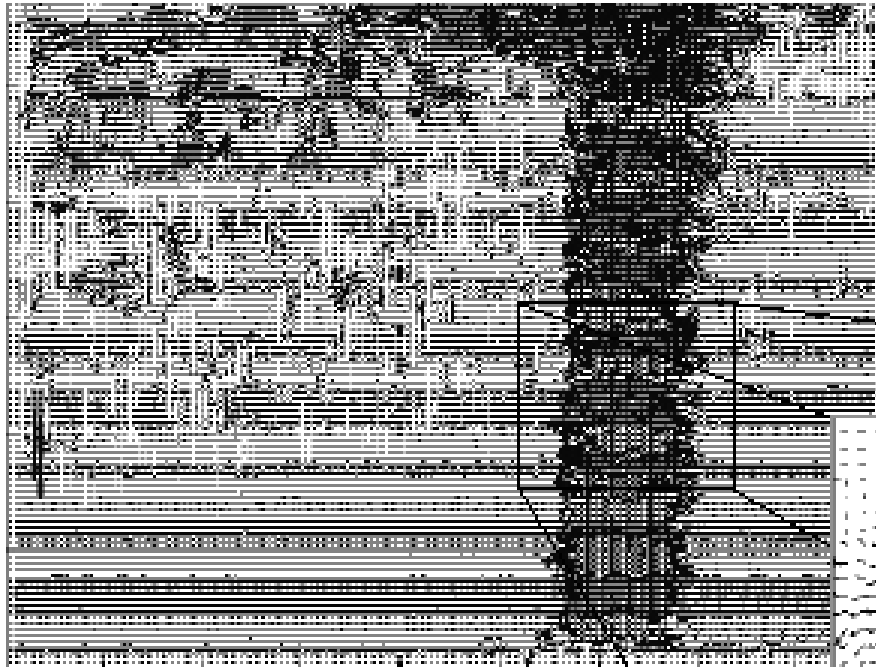
Multi-resolution registration



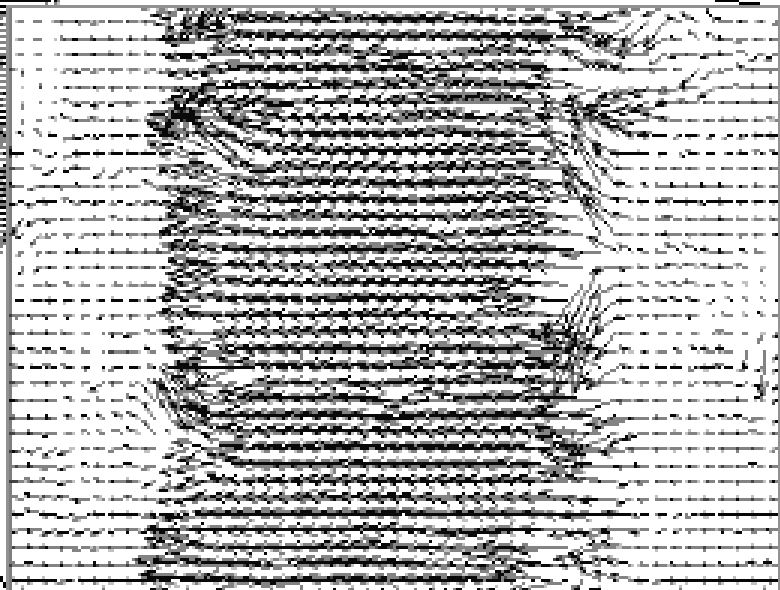
Optical Flow Results



Optical Flow Results



Lucas-Kanade with Pyramids



Errors in Lucas-Kanade

- The motion is large
 - Possible Fix: Keypoint matching
- A point does not move like its neighbors
 - Possible Fix: Region-based matching
- Brightness constancy does not hold
 - Possible Fix: Gradient constancy

Summary

- Major contributions from Lucas, Tomasi, Kanade
 - Tracking feature points
 - Optical flow
- Key ideas
 - By assuming brightness constancy, truncated Taylor expansion leads to simple and fast patch matching across frames
 - Coarse-to-fine registration