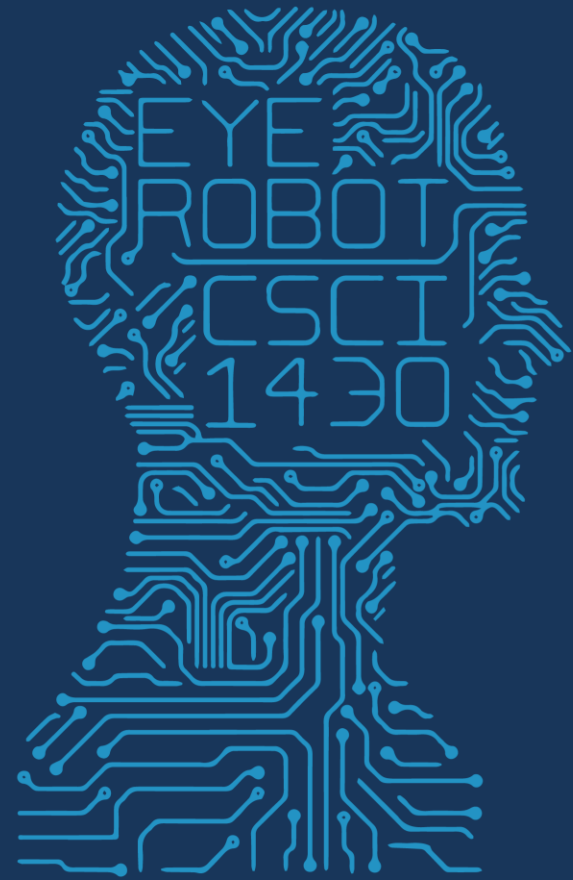




1950

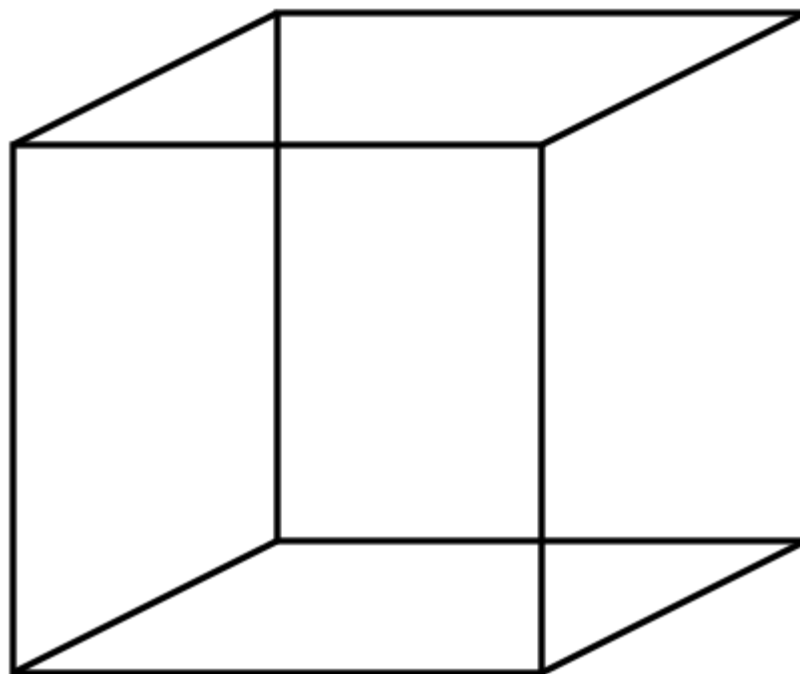
FUTURE VISION



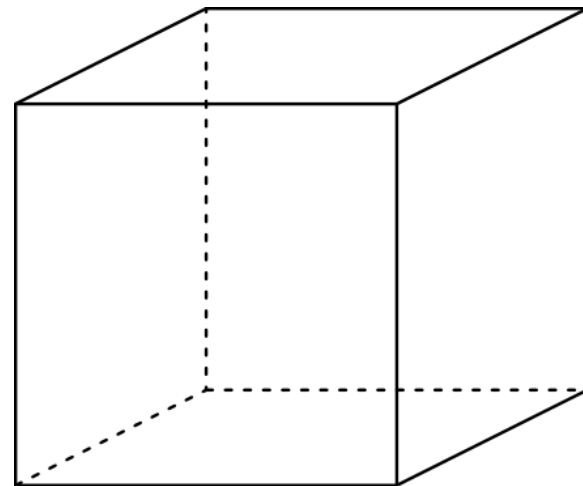
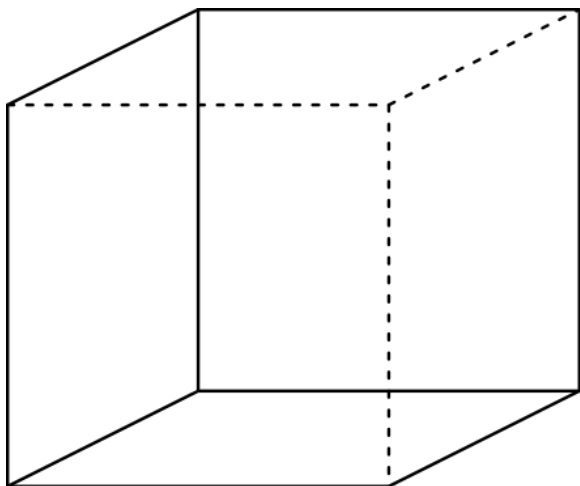
17 APRIL 2019

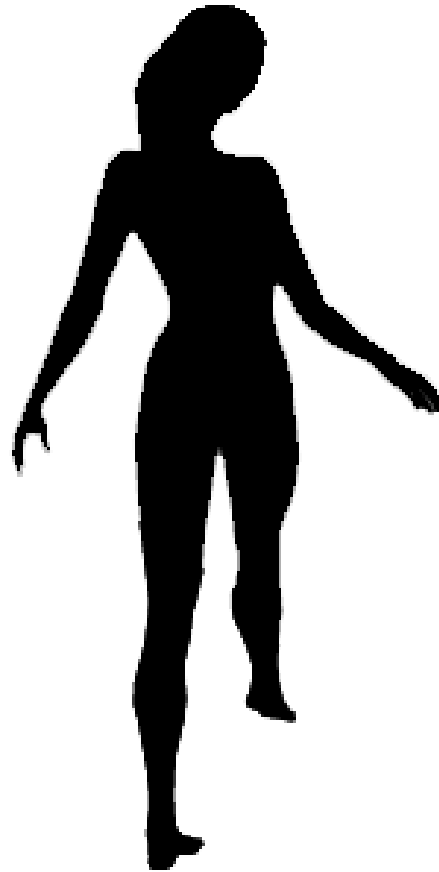
COMPUTER VISION

Multi-stable Perception

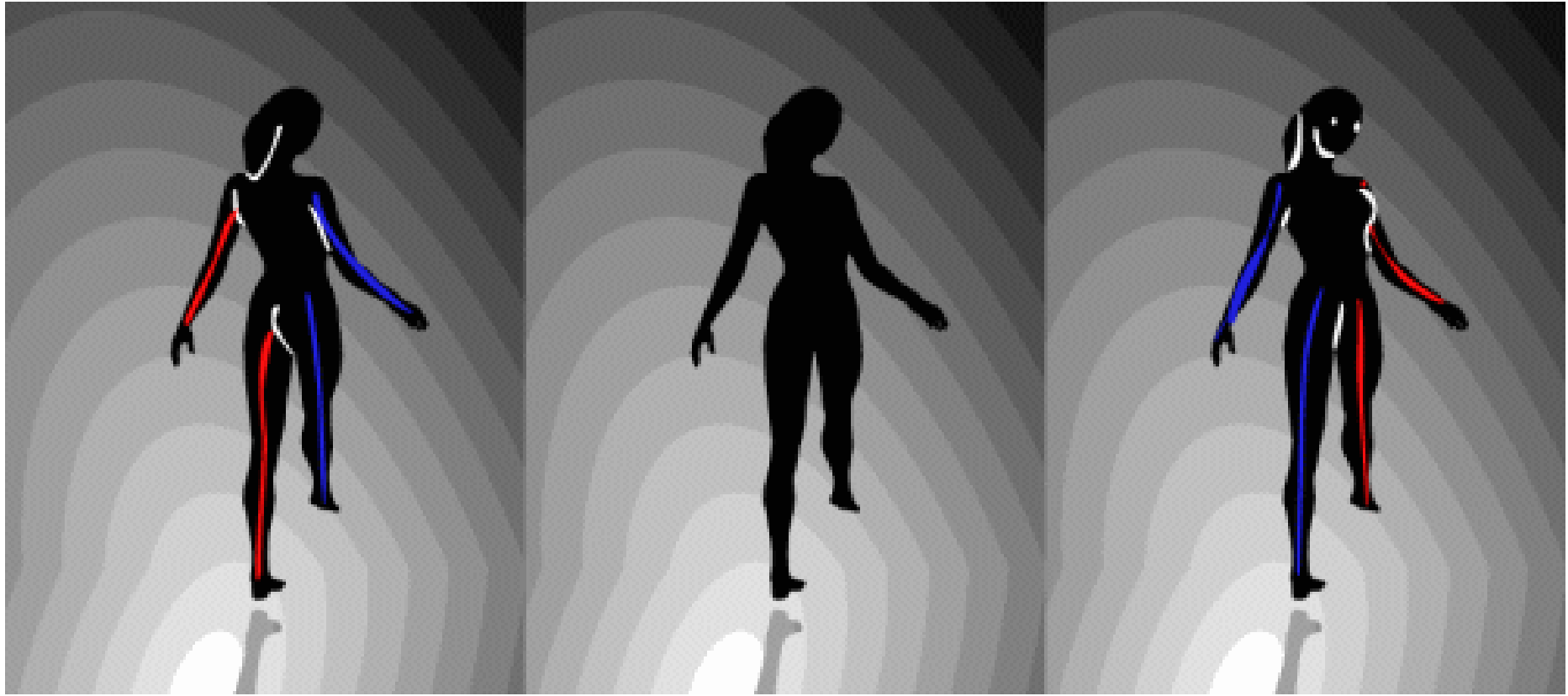


Necker Cube





Spinning dancer illusion, Nobuyuki Kayahara

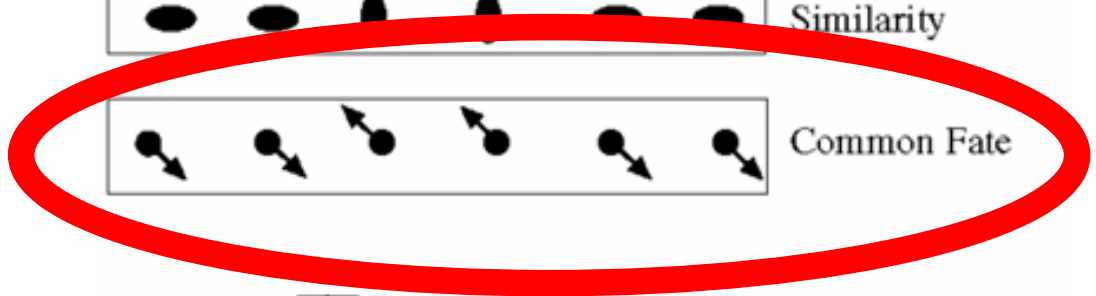
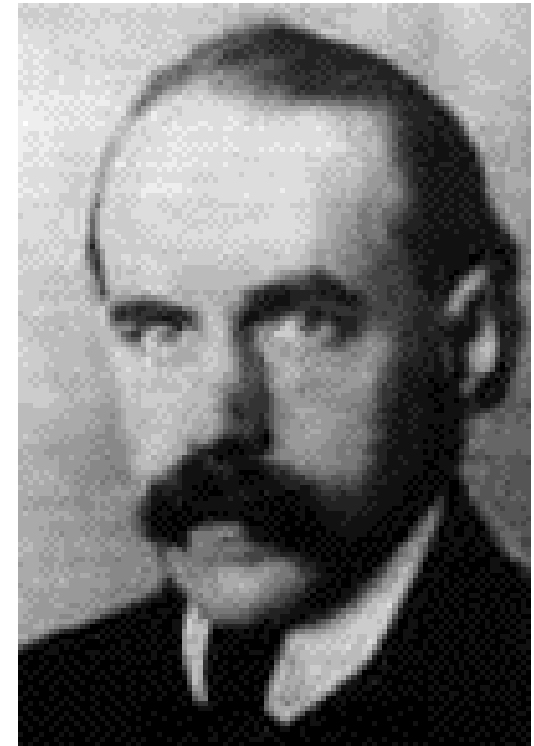


Motion and Gestalt laws of grouping



Gestalt psychology
(Max Wertheimer,
1880-1943)

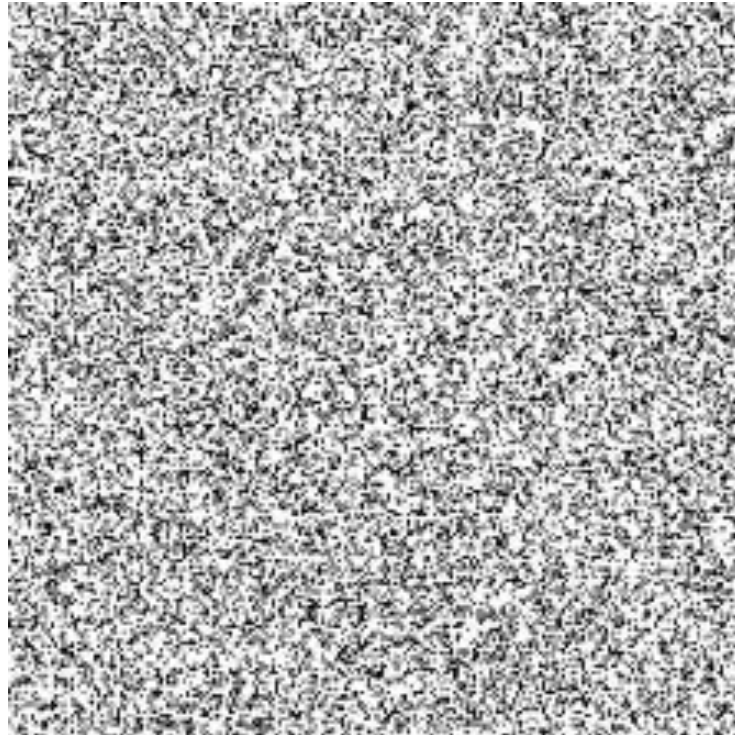
Motion and perceptual organization



Gestalt psychology
(Max Wertheimer,
1880-1943)

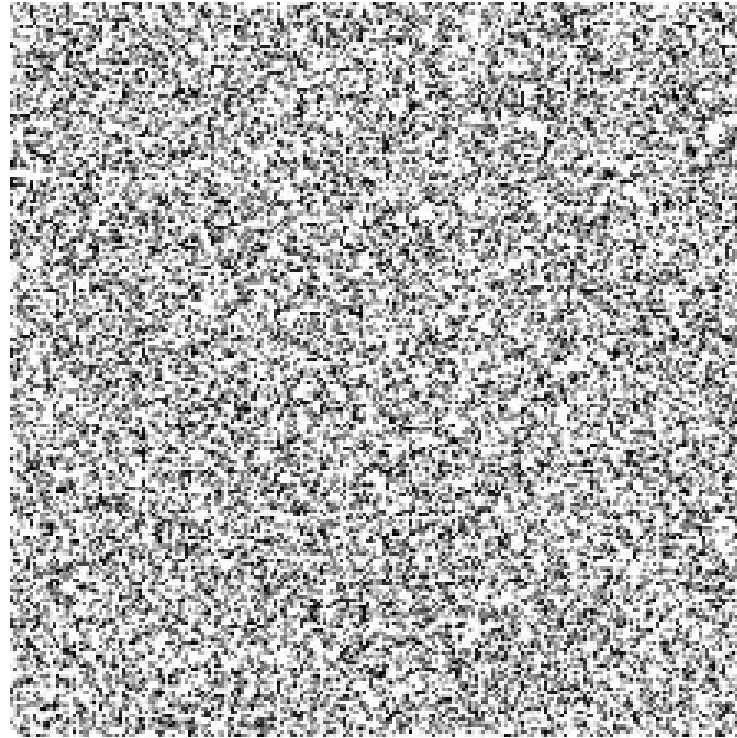
...plus closure, continuation, 'good form'

Motion and perceptual organization



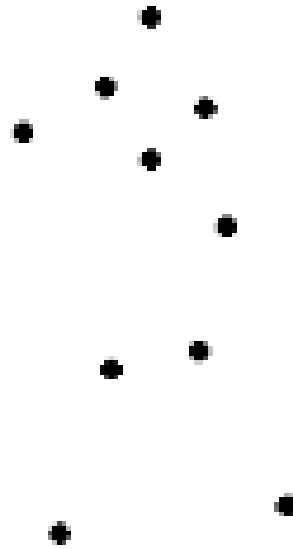
Motion and perceptual organization

- Sometimes motion is the only cue...



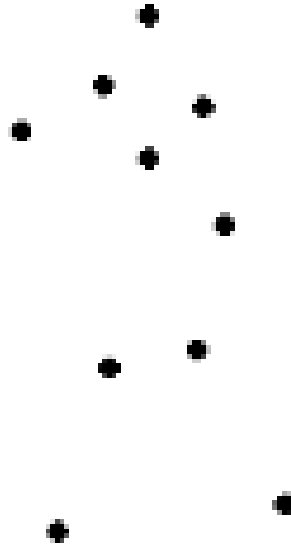
Motion and perceptual organization

Even “impoverished” motion data can evoke a strong percept



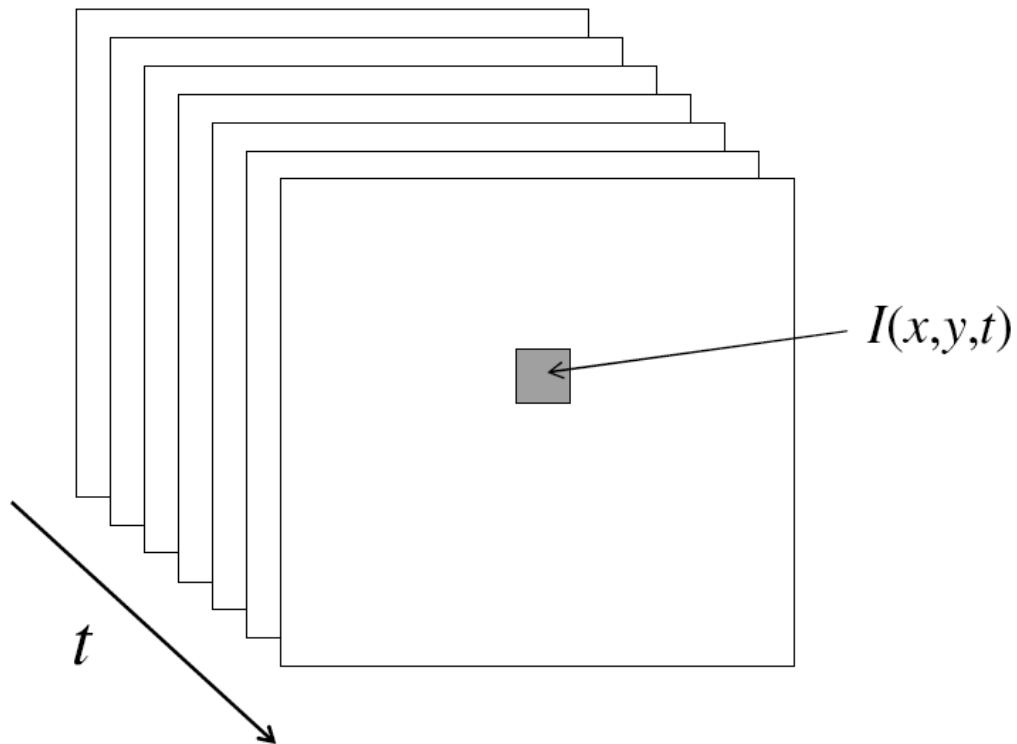
Motion and perceptual organization

Even “impoverished” motion data can evoke a strong percept



Video

- A video is a sequence of frames captured over time
- A 'function' of space (x, y) and time (t)



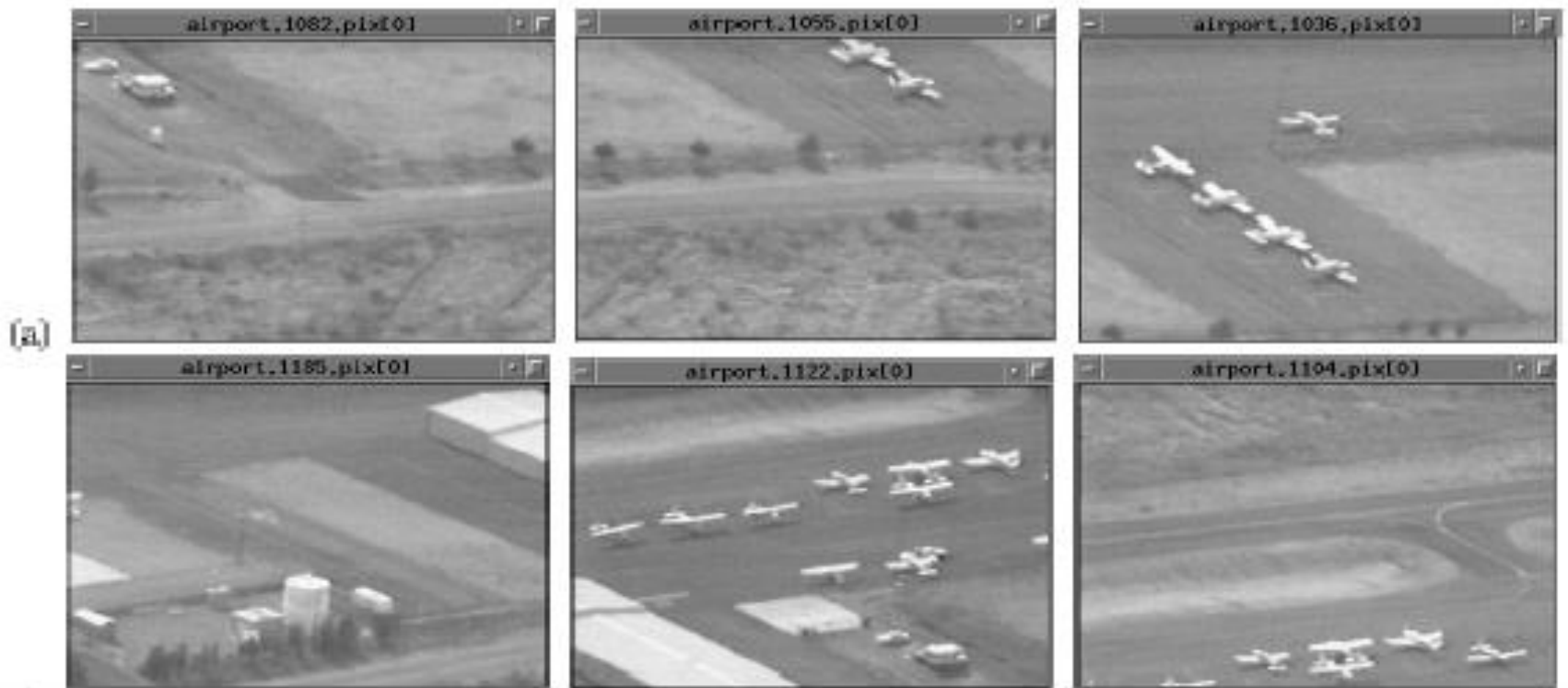
Motion Applications

- Background subtraction
- Shot boundary detection
- Motion segmentation
 - Segment the video into multiple coherently moving objects



© HENNIE LACOCK/CATERS NEWS AGENCY

Mosaicing



(Michal Irani, Weizmann)

Mosaicing



1. Static background mosaic of an airport video clip.

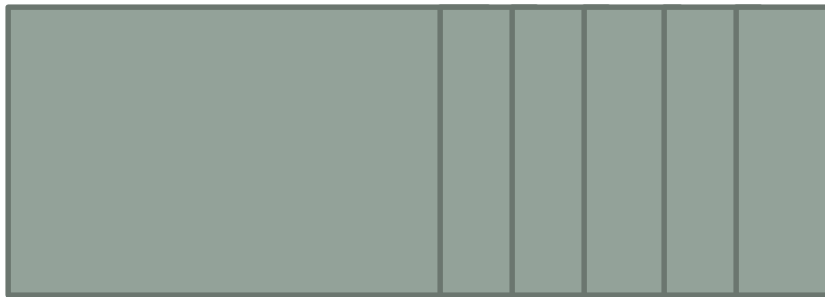
(a) A few representative frames from the minute-long video clip. The video shows an airport being imaged from the air with a moving camera. The scene itself is static (i.e., no moving objects). (b) The static background mosaic image which provides an extended view of the entire scene imaged by the camera in the one-minute video clip.

(Michal Irani, Weizmann)

Mosaicing for Panoramas on Smartphones

Left to right sweep of video camera

Frame t $t+1$ $t+3$ $t+5$



Estimate motion frame to frame, but compare small overlap for efficiency.

Mosaicing for Panoramas on Smartphones



Mosaicing for Panoramas on Smartphones



Mosaicing for Panoramas on Smartphones



Mosaicing for Panoramas on Smartphones

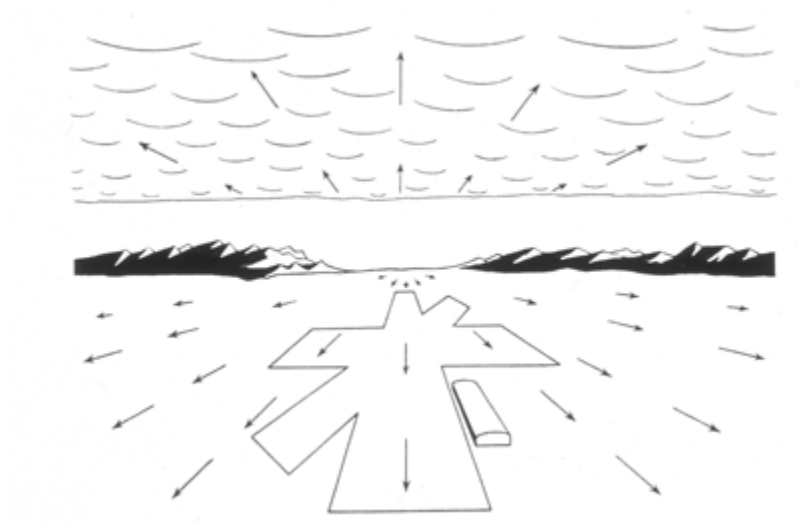


Motion estimation techniques

- 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)
- Direct, dense 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
 - *Optical flow!*

Computer Vision

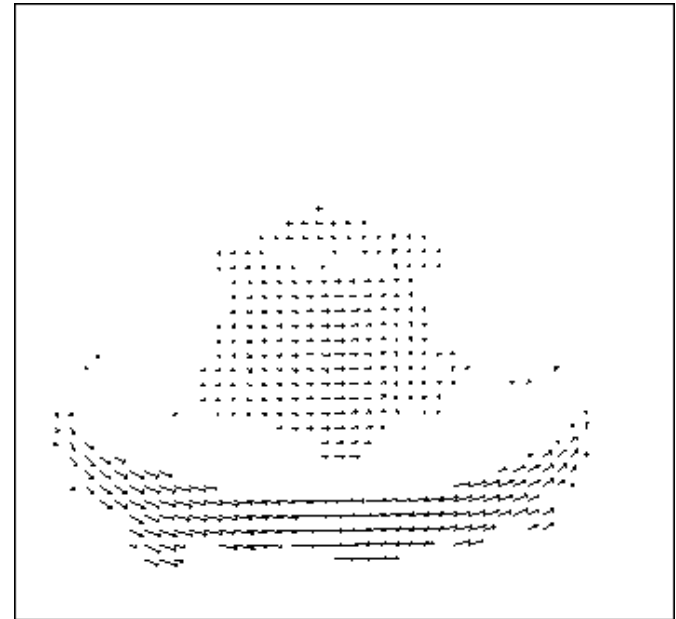
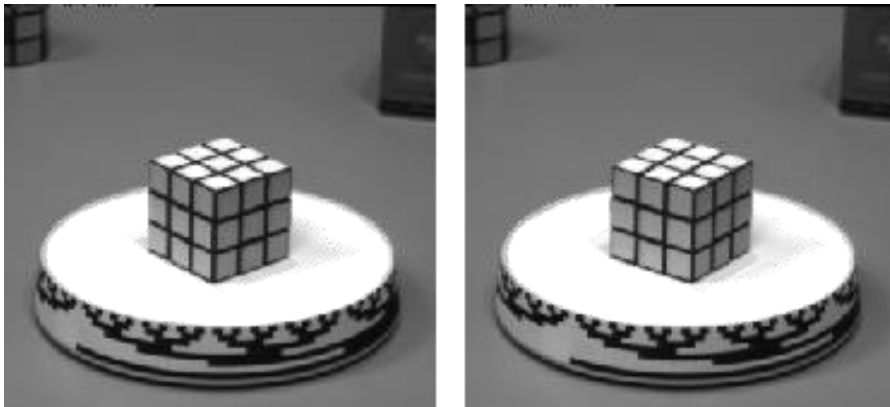
Motion and Optical Flow



Many slides adapted from J. Hays, S. Seitz, R. Szeliski, M. Pollefeys, K. Grauman and others...

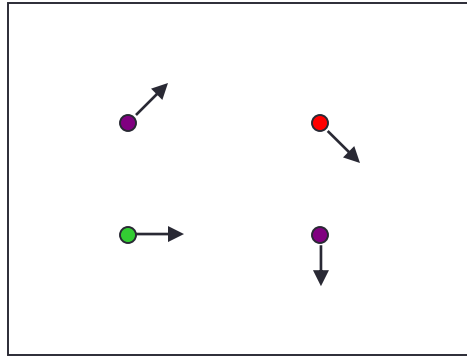
Motion estimation: Optical flow

Optic flow is the **apparent** motion of objects or surfaces

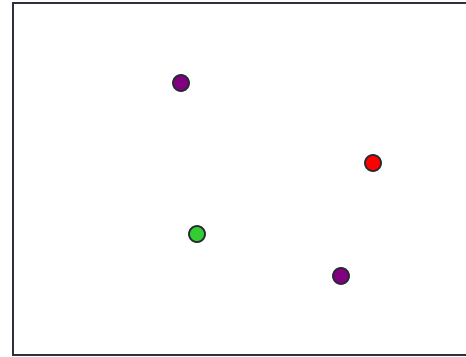


Will start by estimating motion of each pixel separately
Then will consider motion of entire image

Problem definition: optical flow



$I(x, y, t)$



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

How to estimate pixel motion from image $I(x, y, t)$ to $I(x, y, t + 1)$?

No problem!

```
for x1 in I(x1,y1,t)
```

```
  for y1 in I(x1,y1,t)
```

```
    for x2 in I(x2,y2,t+1)
```

```
      for y2 in I(x2,y2,t+1)
```

```
        Patch sum of squared differences
```

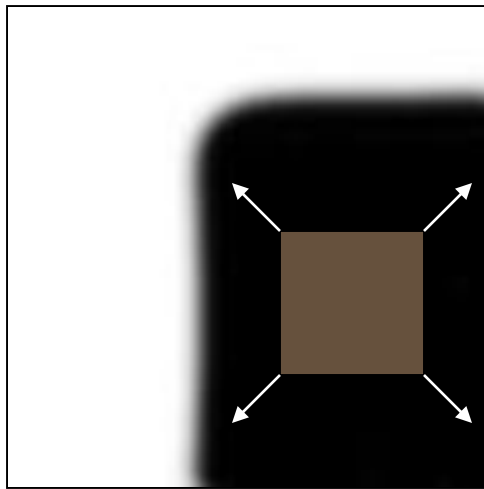
Slow.

Back in early
lectures...

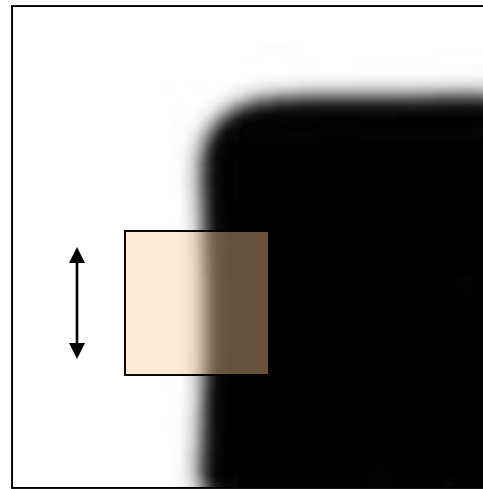
Corner Detection: Basic Idea

Recognize corners by looking at small window.

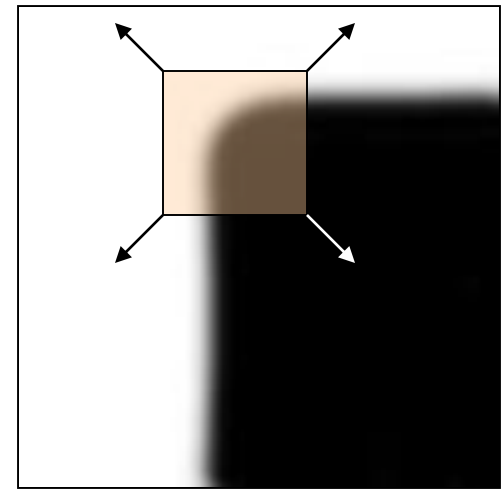
We want a window shift in *any direction* to give a *large change* in intensity.



“Flat” region:
no change in
all directions



“Edge”:
no change
along the edge
direction



“Corner”:
significant
change in all
directions

Corner Detection by Auto-correlation

Change in appearance of window $w(x,y)$ for shift $[u,v]$:

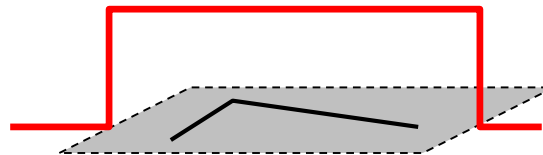
$$E(u, v) = \sum_{x, y} w(x, y) [I(x+u, y+v) - I(x, y)]^2$$

Window
function

Shifted
intensity

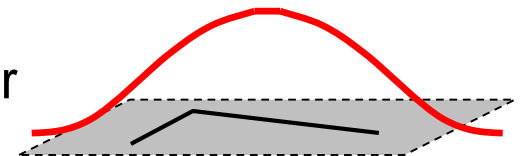
Intensity

Window function $w(x,y) =$



1 in window, 0 outside

or



Gaussian

Corner Detection by Auto-correlation

Change in appearance of window $w(x,y)$ for shift $[u,v]$:

$$E(u, v) = \sum_{x,y} w(x, y) [I(x+u, y+v) - I(x, y)]^2$$

We want to discover how E behaves for small shifts

But this is very slow to compute naively.

$O(\text{window_width}^2 * \text{shift_range}^2 * \text{image_width}^2)$

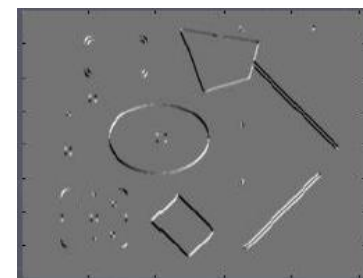
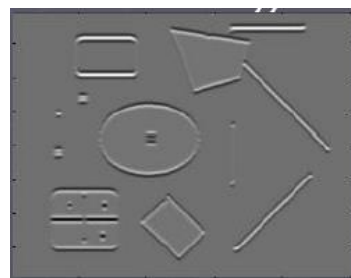
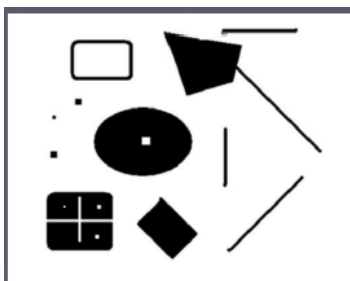
$O(11^2 * 11^2 * 600^2) = 5.2$ billion of these
14.6 thousand per pixel in your image



Corners as distinctive interest points

$$M = \sum w(x, y) \begin{bmatrix} I_x I_x & I_x I_y \\ I_x I_y & I_y I_y \end{bmatrix}$$

2 x 2 matrix of image derivatives
(averaged in neighborhood of a point)



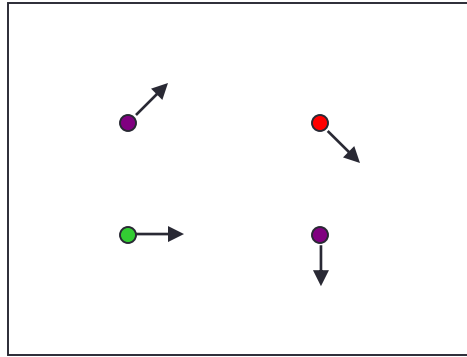
Notation:

$$I_x \Leftrightarrow \frac{\partial I}{\partial x}$$

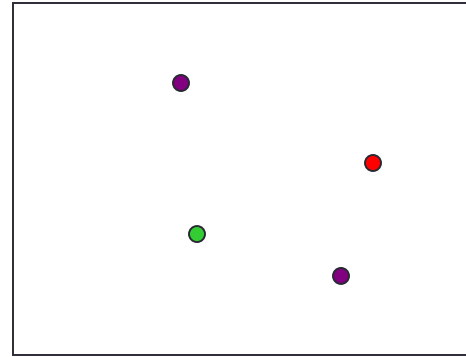
$$I_y \Leftrightarrow \frac{\partial I}{\partial y}$$

$$I_x I_y \Leftrightarrow \frac{\partial I}{\partial x} \frac{\partial I}{\partial y}$$

Problem definition: optical flow



$I(x, y, t)$



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

How to estimate pixel motion from image $I(x, y, t)$ to $I(x, y, t + 1)$?

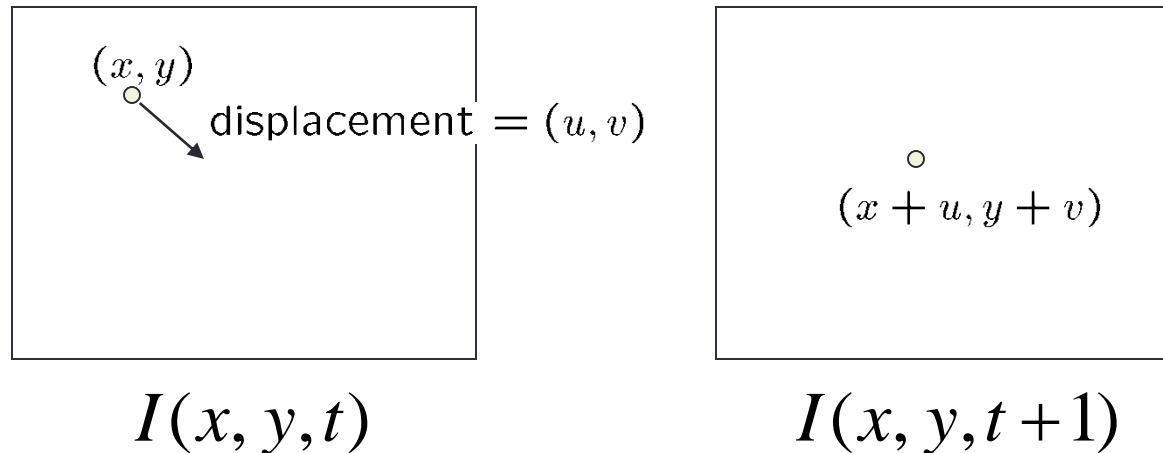
Solve pixel correspondence problem

- Given a pixel in $I(x, y, t)$, look for nearby pixels of the same color in $I(x, y, t + 1)$

Key assumptions

- **Small motion:** Points do not move very far
- **Color constancy:** A point in $I(x, y, t)$ looks the same in $I(x, y, t + 1)$
 - For grayscale images, this is brightness constancy

Optical flow constraints (grayscale images)



Let's look at these constraints more closely

Brightness constancy constraint (equation)

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

Small motion: (u and v are less than 1 pixel, or smoothly varying)

Taylor series expansion of I :

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

Optical flow equation

Combining these two equations

$$0 = I(x + u, y + v, t + 1) - I(x, y, t)$$

(Short hand: $I_x = \frac{\partial I}{\partial x}$
for t **or** $t+1$)

Optical flow equation

- Combining these two equations

$$\begin{aligned}0 &= I(x+u, y+v, t+1) - I(x, y, t) \\ &\approx I(x, y, t+1) + I_x u + I_y v - I(x, y, t) \\ &\approx [I(x, y, t+1) - I(x, y, t)] + I_x u + I_y v \\ &\approx I_t + I_x u + I_y v \\ &\approx I_t + \nabla I \cdot \langle u, v \rangle\end{aligned}$$

(Short hand: $I_x = \frac{\partial I}{\partial x}$
for t **or** $t+1$)

Optical flow equation

- Combining these two equations

$$0 = I(x+u, y+v, t+1) - I(x, y, t)$$

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

(Short hand: $I_x = \frac{\partial I}{\partial x}$
for t or $t+1$)

$$\approx [I(x, y, t+1) - I(x, y, t)] + I_x u + I_y v$$

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

$$\approx I_t + \nabla I \cdot \langle u, v \rangle$$

In the limit as u and v go to zero, this becomes exact

$$0 = I_t + \nabla I \cdot \langle u, v \rangle$$

Brightness constancy constraint equation

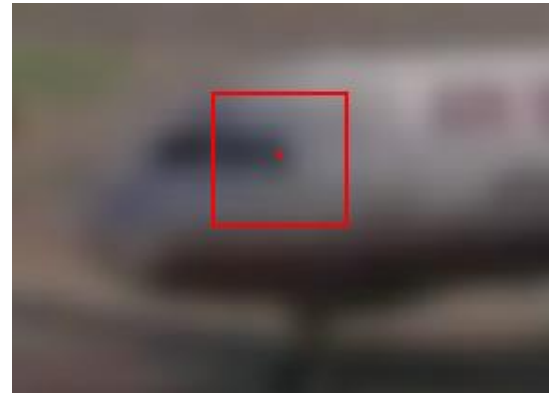
$$I_x u + I_y v + I_t = 0$$

How does this make sense?

Brightness constancy constraint equation

$$I_x u + I_y v + I_t = 0$$

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?

$$0 = I_t + \nabla I \cdot \langle u, v \rangle \quad \text{or} \quad I_x u + I_y v + I_t = 0$$

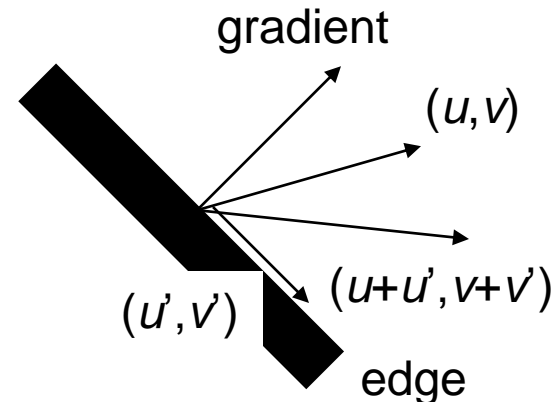
How many equations and unknowns per pixel?

One equation (this is a scalar equation!), two unknowns (u,v)

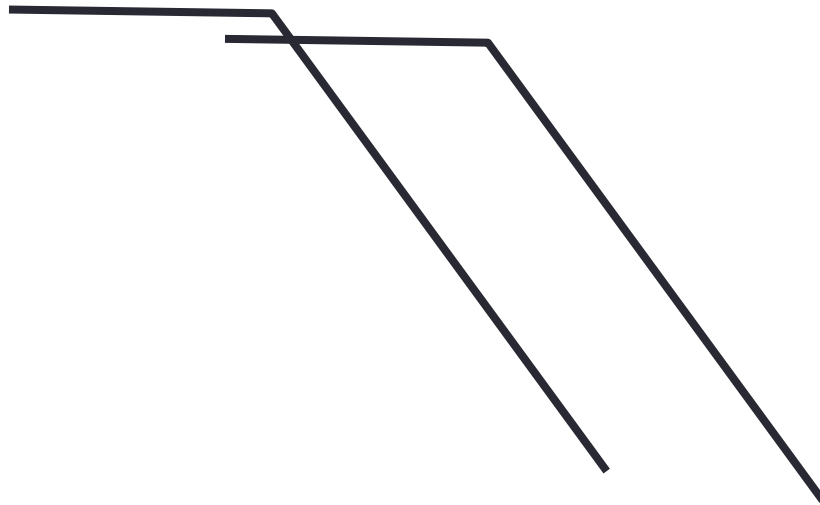
The component of the motion perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

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

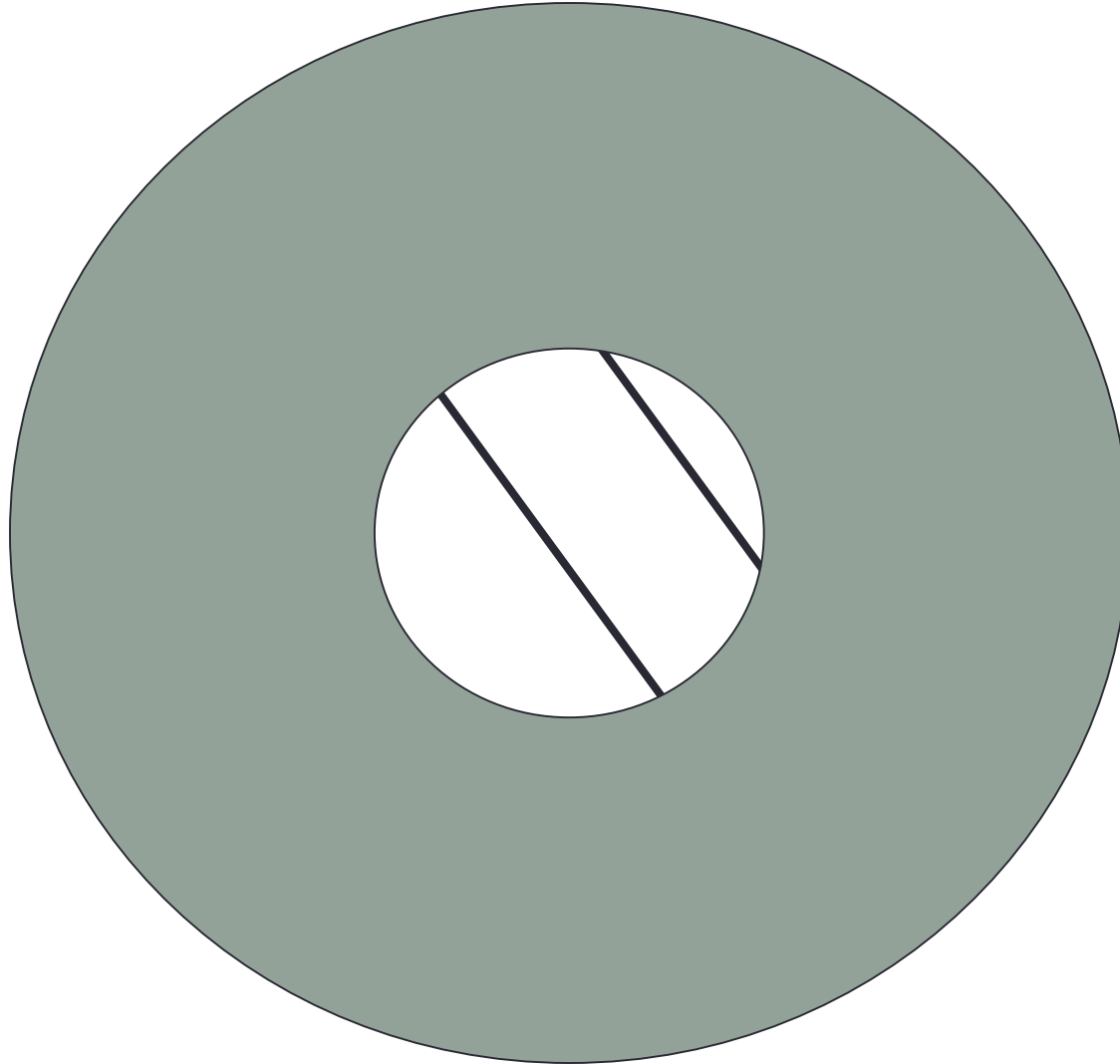
$$\nabla I \cdot [u' \ v']^T = 0$$



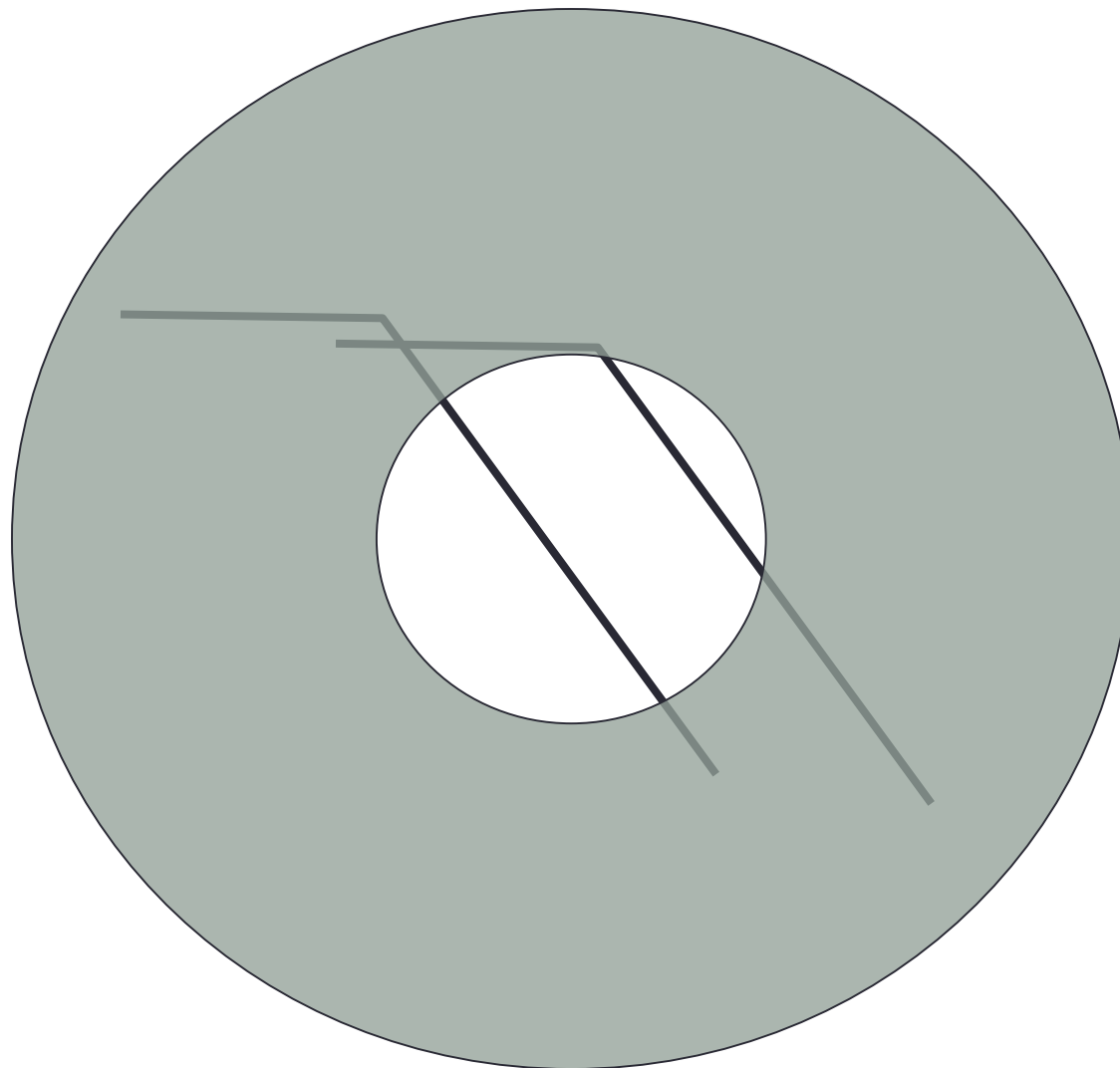
Aperture problem



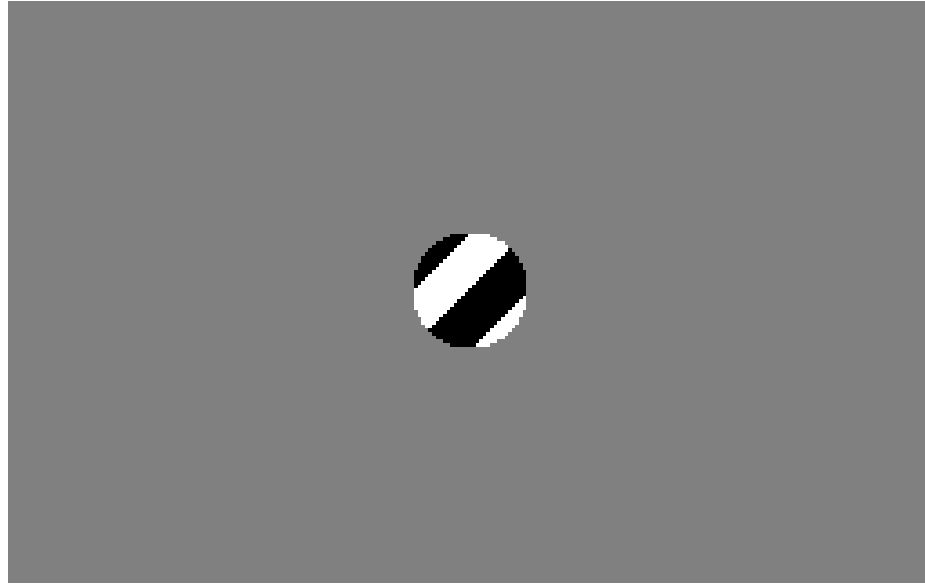
Aperture problem



Aperture problem

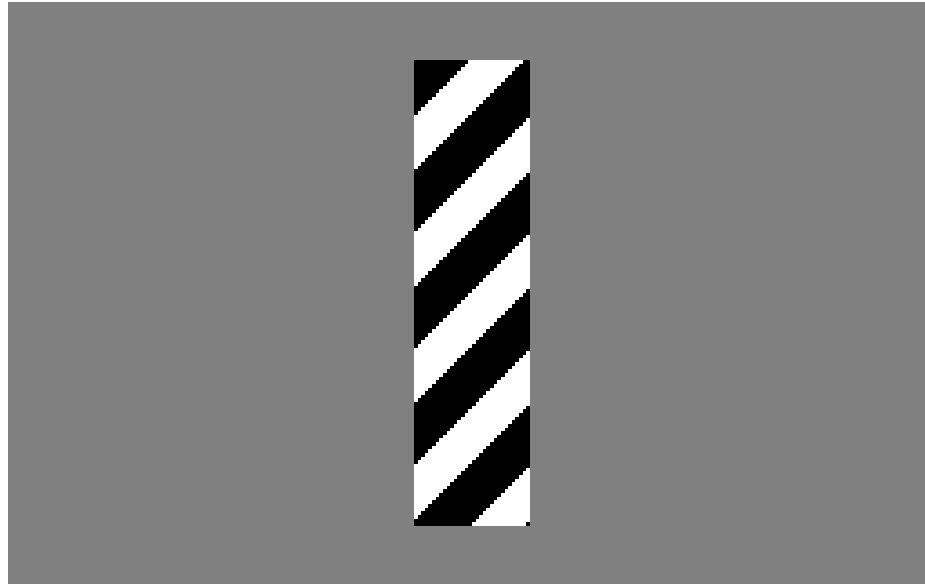


The barber pole illusion



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

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

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

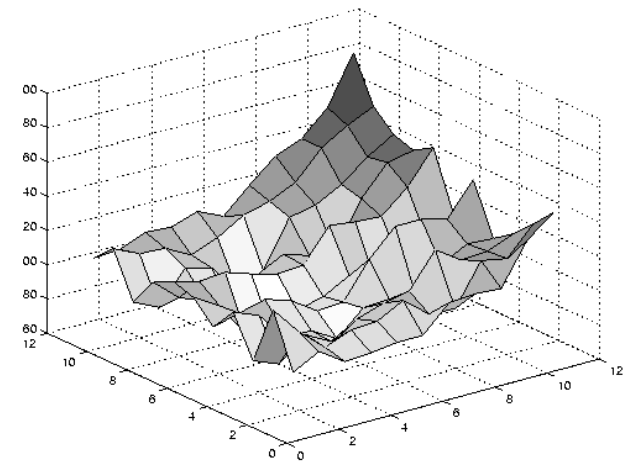
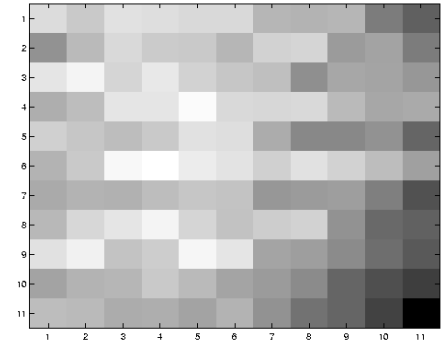
When is this solvable? What are good points to track?

- $A^T A$ 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

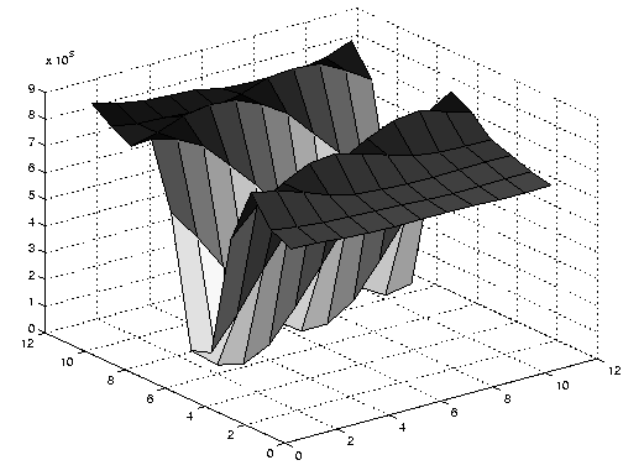
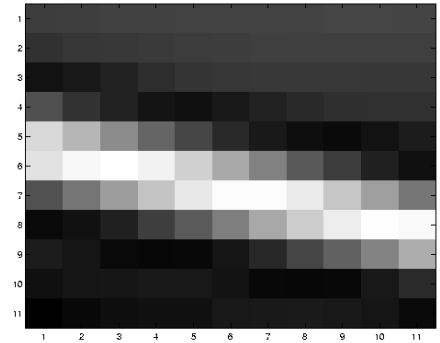
Low texture region



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

- gradients have small magnitude
- small λ_1 , small λ_2

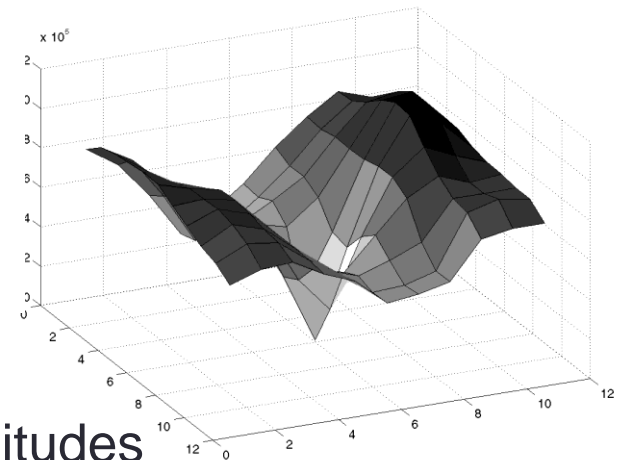
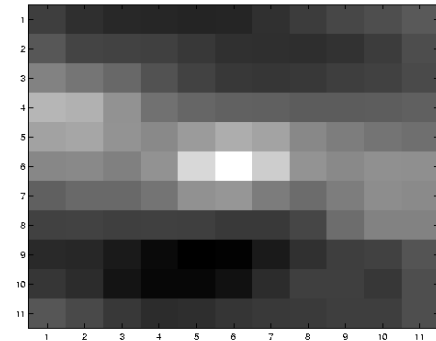
Edge



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

- large gradients, all the same
- large λ_1 , small λ_2

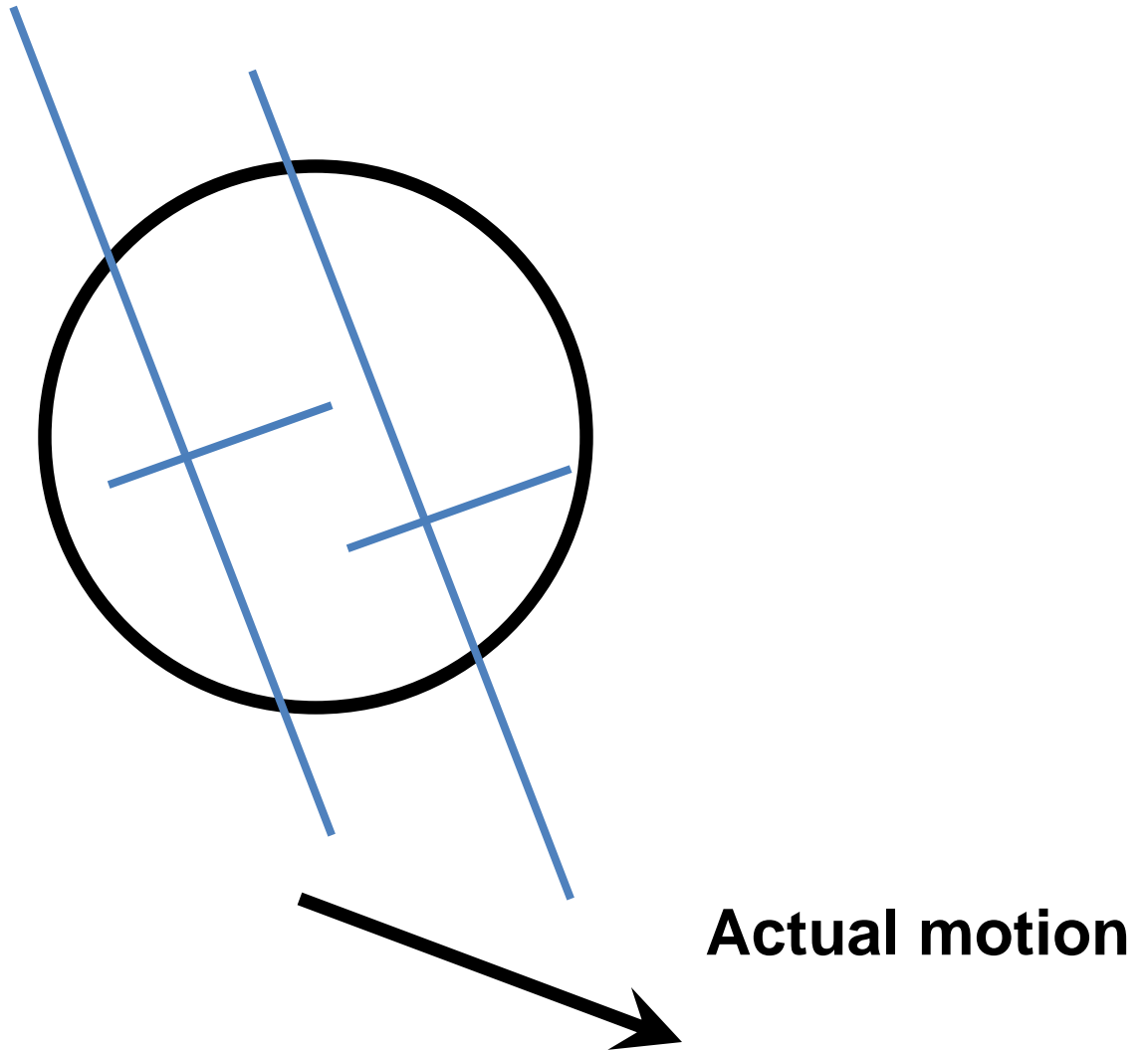
High textured region



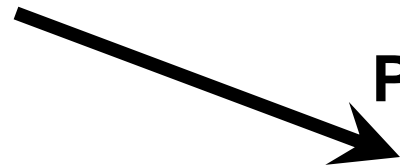
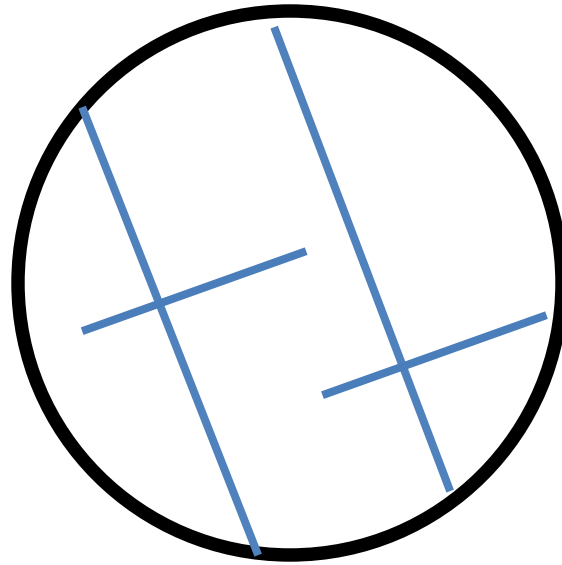
$$\sum \nabla I (\nabla I)^T$$

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

The aperture problem resolved



The aperture problem resolved



Perceived motion

Errors in assumptions

- A point does not move like its neighbors
 - Motion segmentation
- Brightness constancy does not hold
 - Do exhaustive neighborhood search with normalized correlation - tracking features – maybe SIFT – more later....
- **The motion is large (larger than a pixel)**
 1. **Not-linear: Iterative refinement**
 2. **Local minima: coarse-to-fine estimation**

Revisiting the small motion assumption

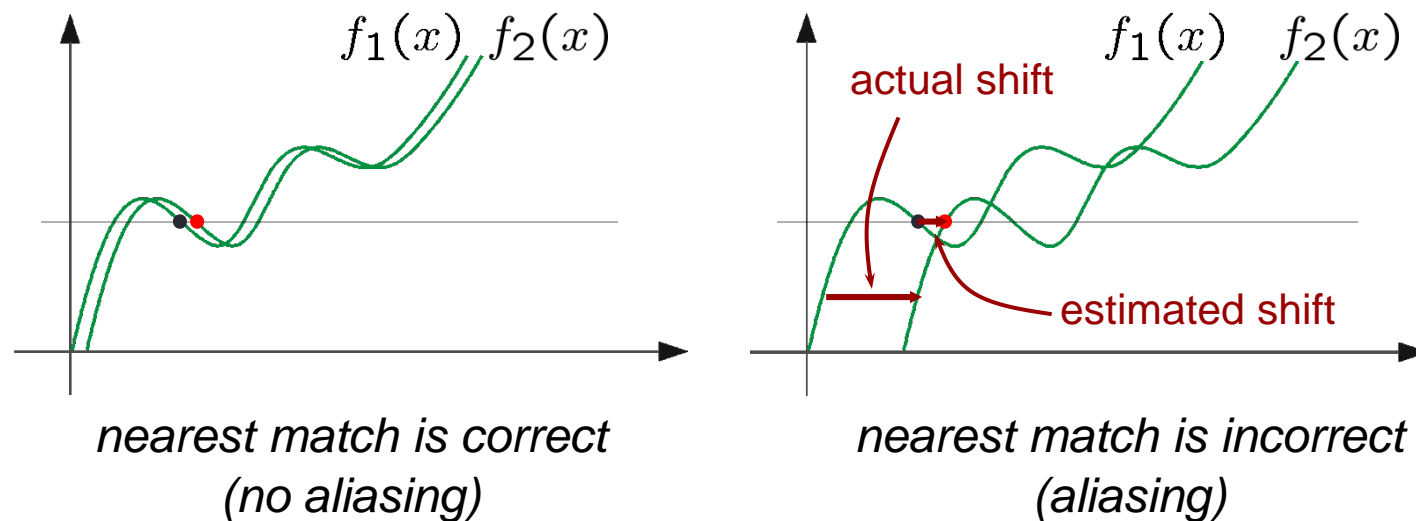


- Is this motion small enough?
 - Probably not—it's much larger than one pixel
 - How might we solve this problem?

Optical Flow: Aliasing

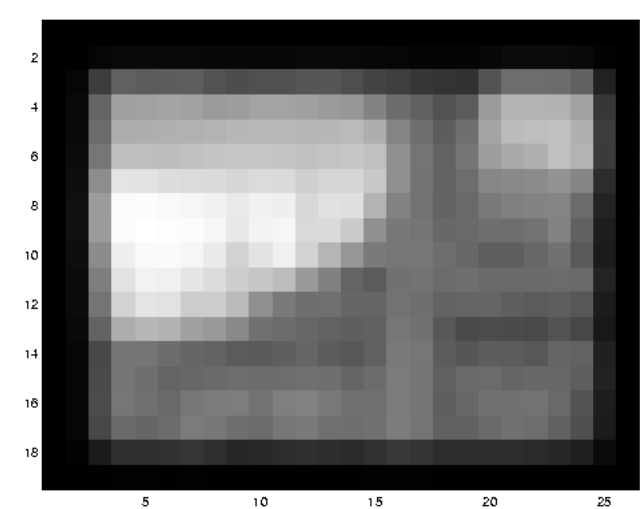
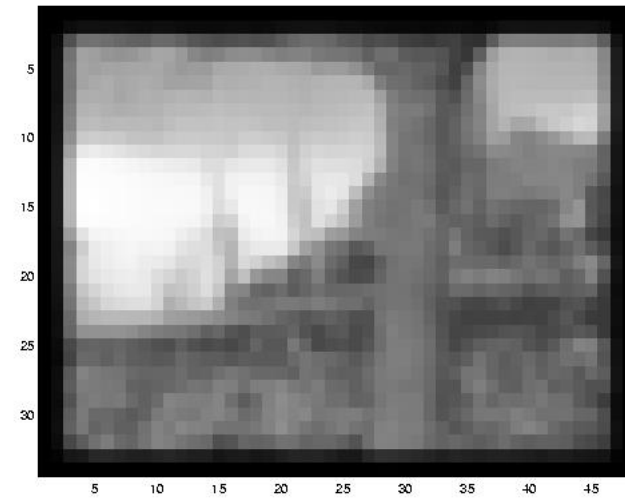
Temporal aliasing causes ambiguities in optical flow because images can have many pixels with the same intensity.

I.e., how do we know which 'correspondence' is correct?

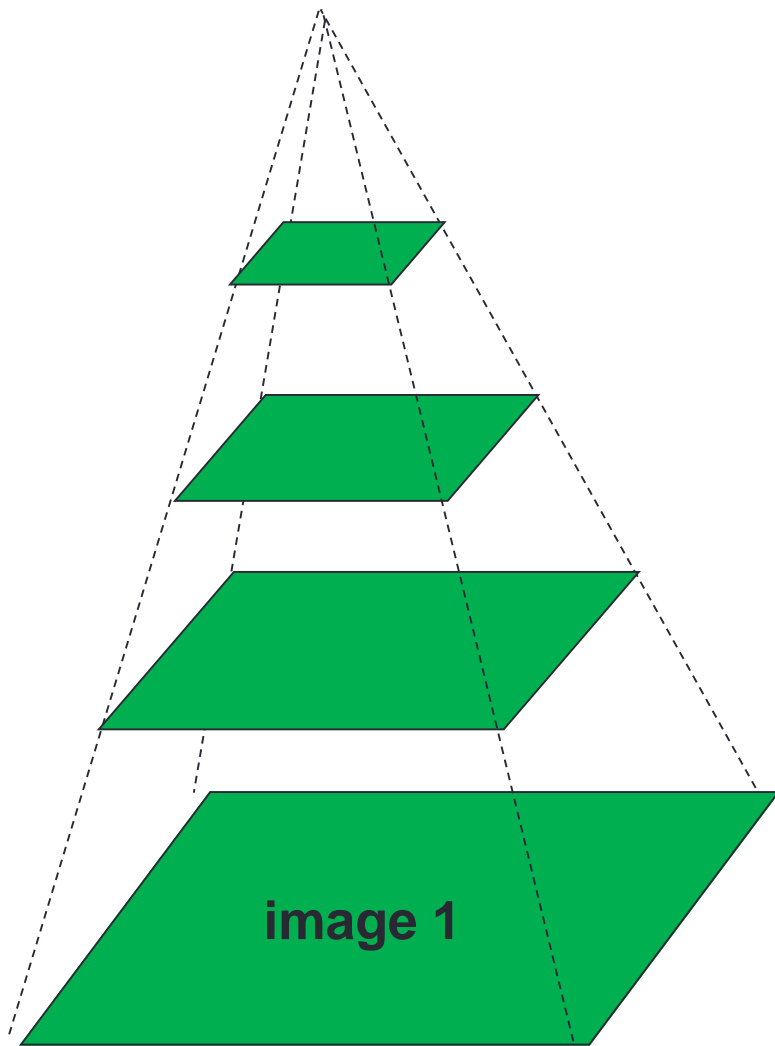


To overcome aliasing: coarse-to-fine estimation.

Reduce the resolution!



Coarse-to-fine optical flow estimation



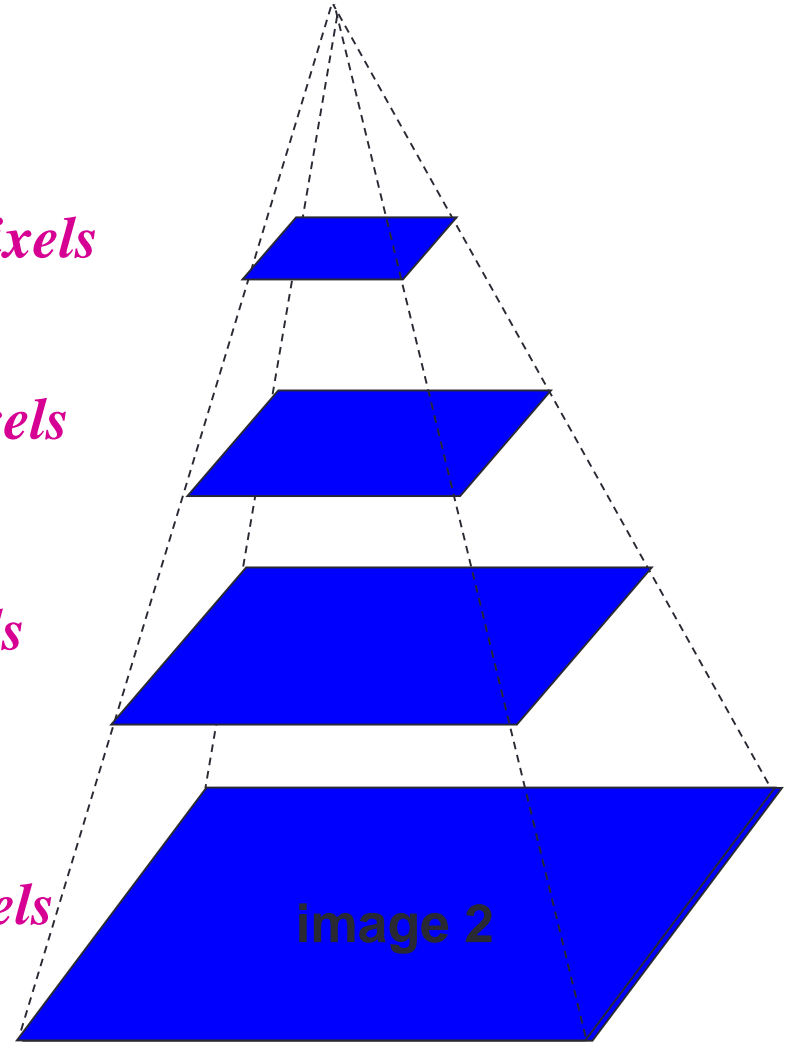
Gaussian pyramid of image 1

$u=1.25$ pixels

$u=2.5$ pixels

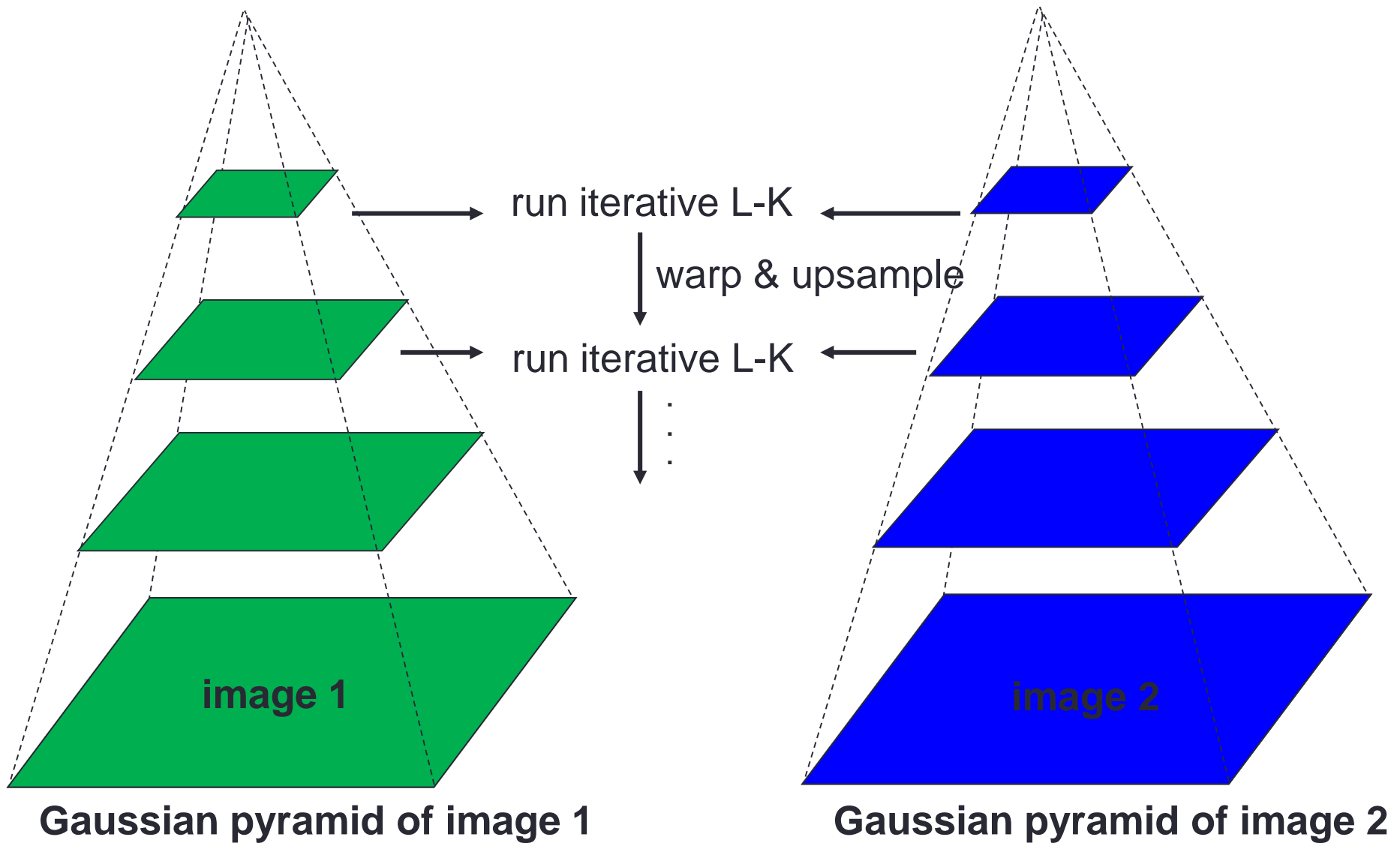
$u=5$ pixels

$u=10$ pixels

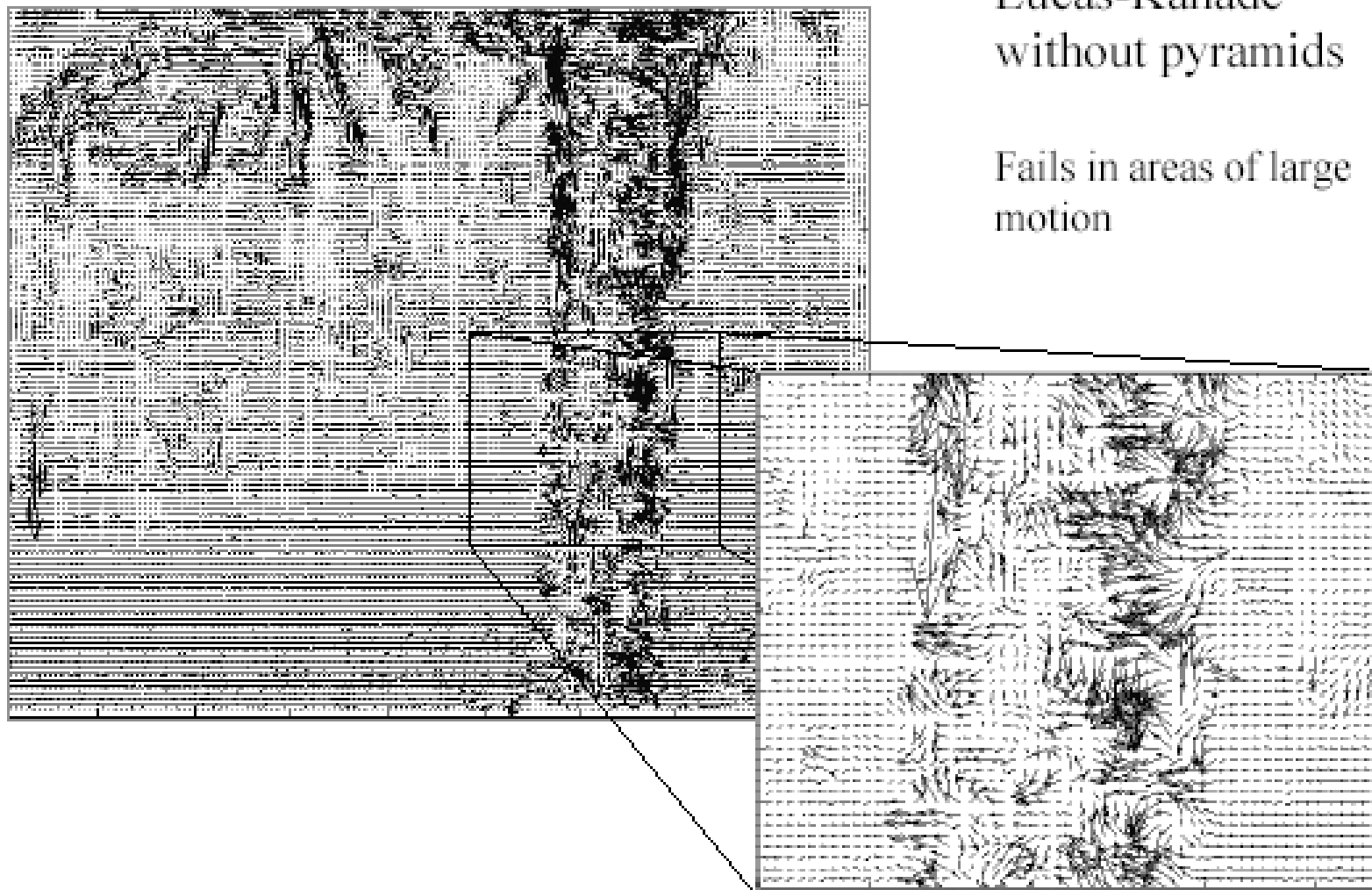


Gaussian pyramid of image 2

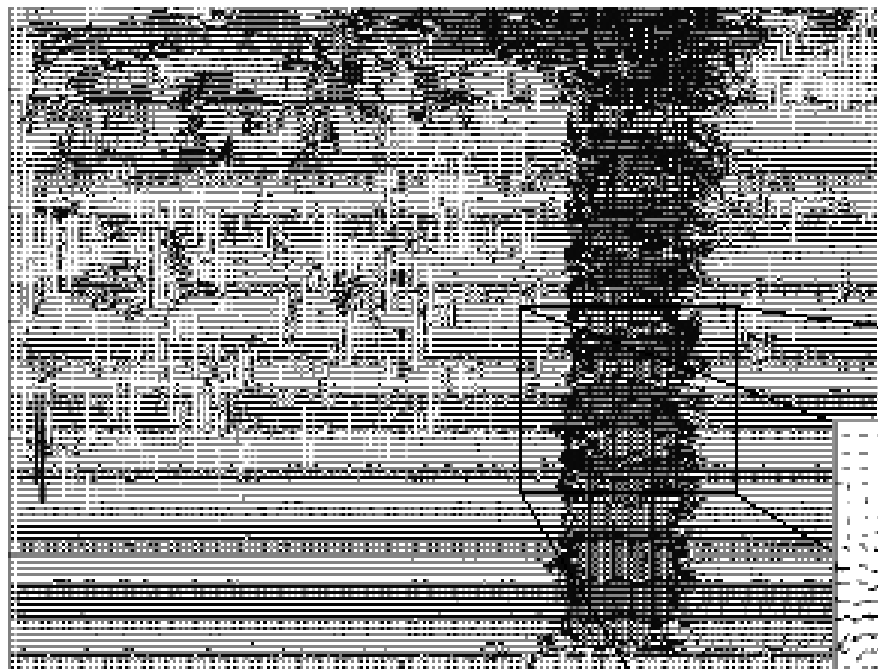
Coarse-to-fine optical flow estimation



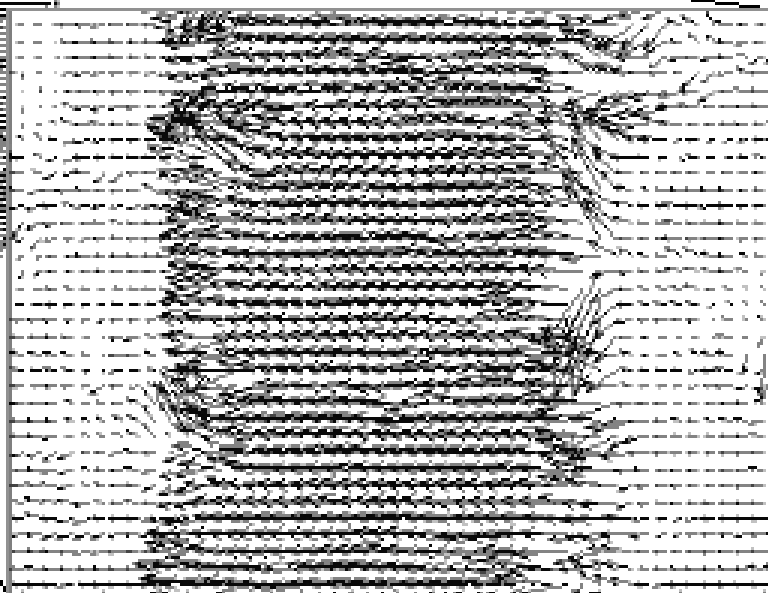
Optical Flow Results



Optical Flow Results



Lucas-Kanade with Pyramids

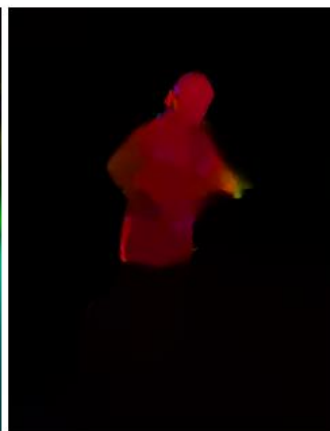
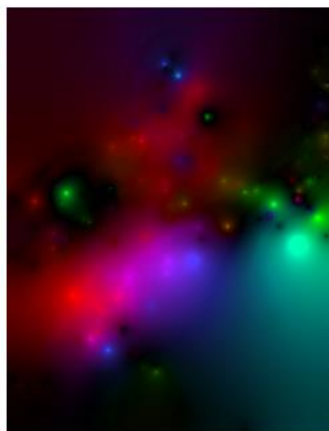
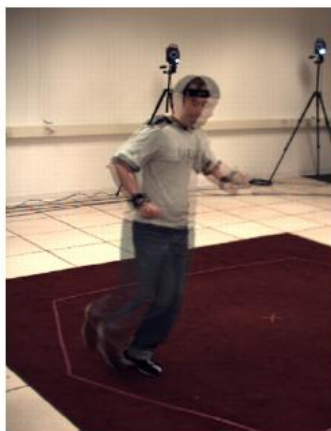


Optical flow

- Definition: the *apparent* motion of brightness patterns in the image
- Ideally, the same as the projected motion field
- Take care: apparent motion can be caused by lighting changes without any actual motion
 - Imagine a uniform rotating sphere under fixed lighting vs. a stationary sphere under moving illumination.

State-of-the-art optical flow, 2009

- Start with something similar to Lucas-Kanade
- + gradient constancy
- + energy minimization with smoothing term
- + region matching
- + keypoint matching (long-range)



Region-based +Pixel-based +Keypoint-based

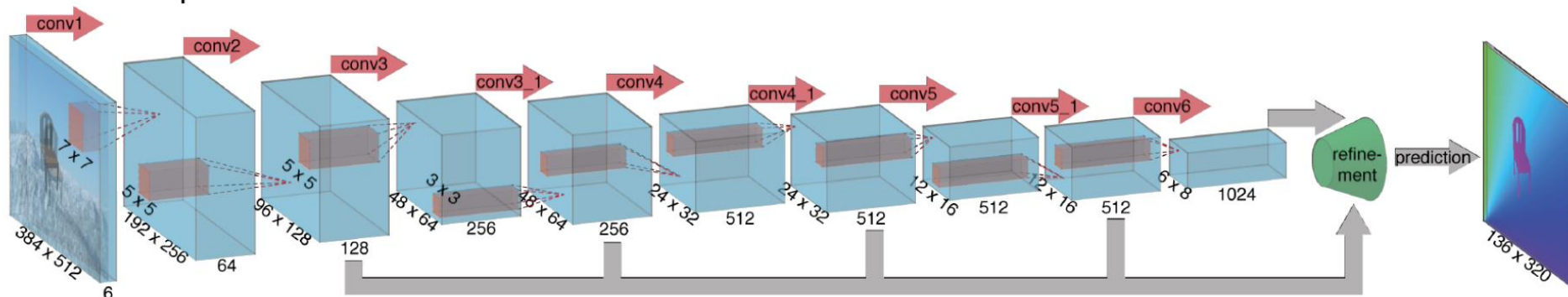
State-of-the-art optical flow, 2015

CNN encoder/decoder

Pair of input frames

Upsample estimated flow back to input resolution

FlowNetSimple



Fischer et al. 2015. <https://arxiv.org/abs/1504.06852>

State-of-the-art optical flow, 2015

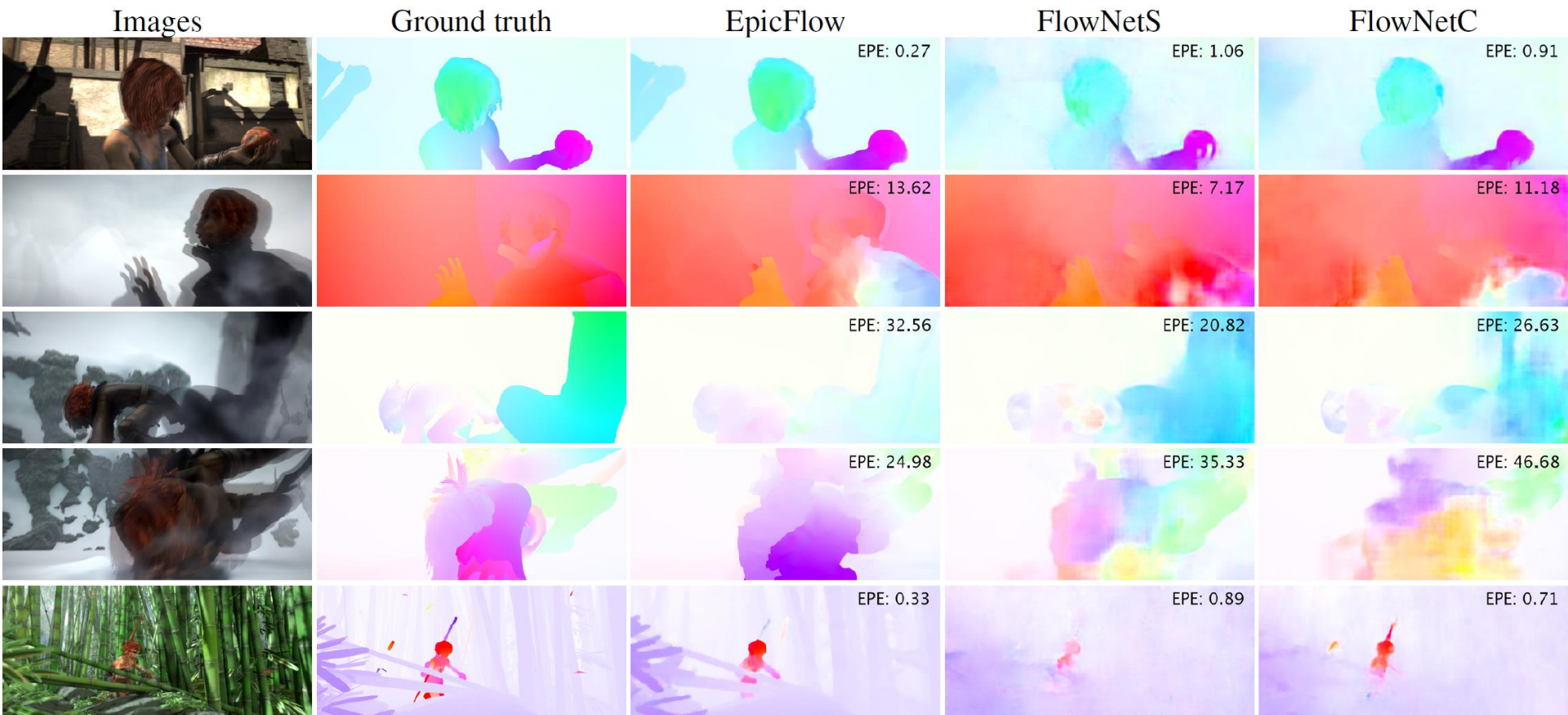
Synthetic Training data



Fischer et al. 2015. <https://arxiv.org/abs/1504.06852>

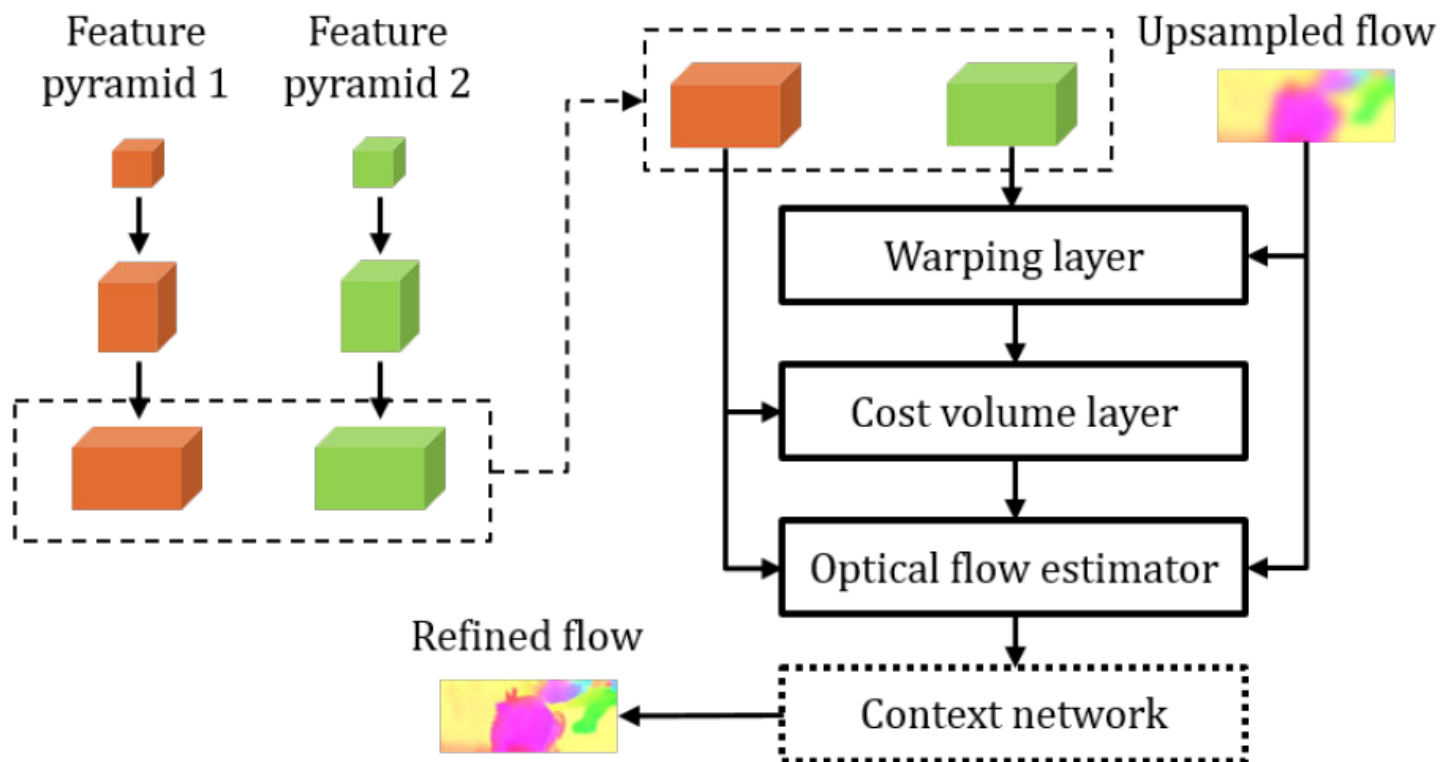
State-of-the-art optical flow, 2015

Results on Sintel

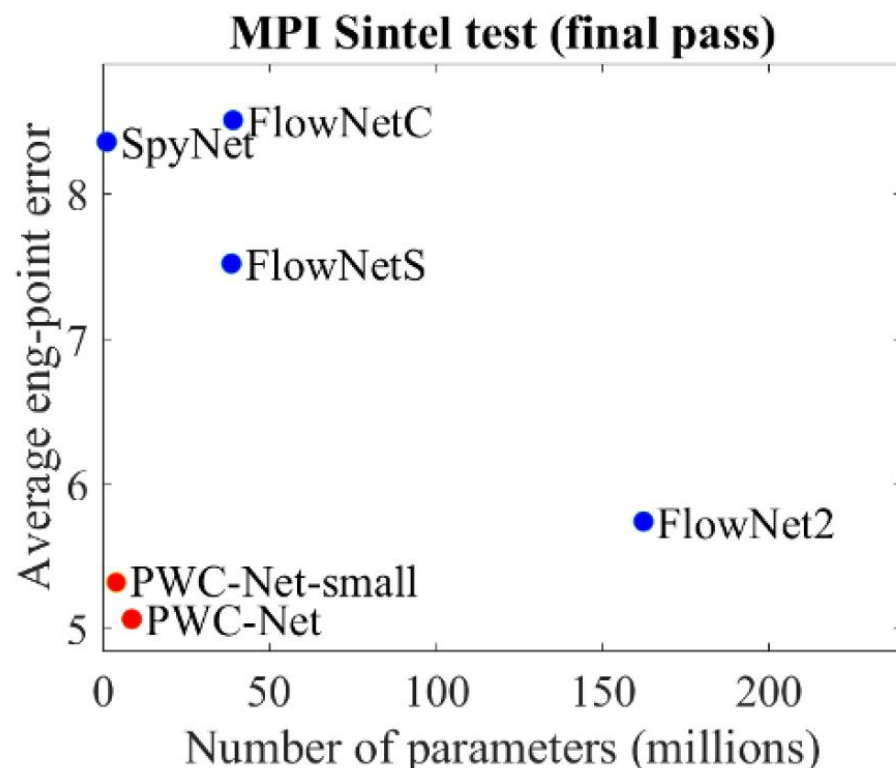
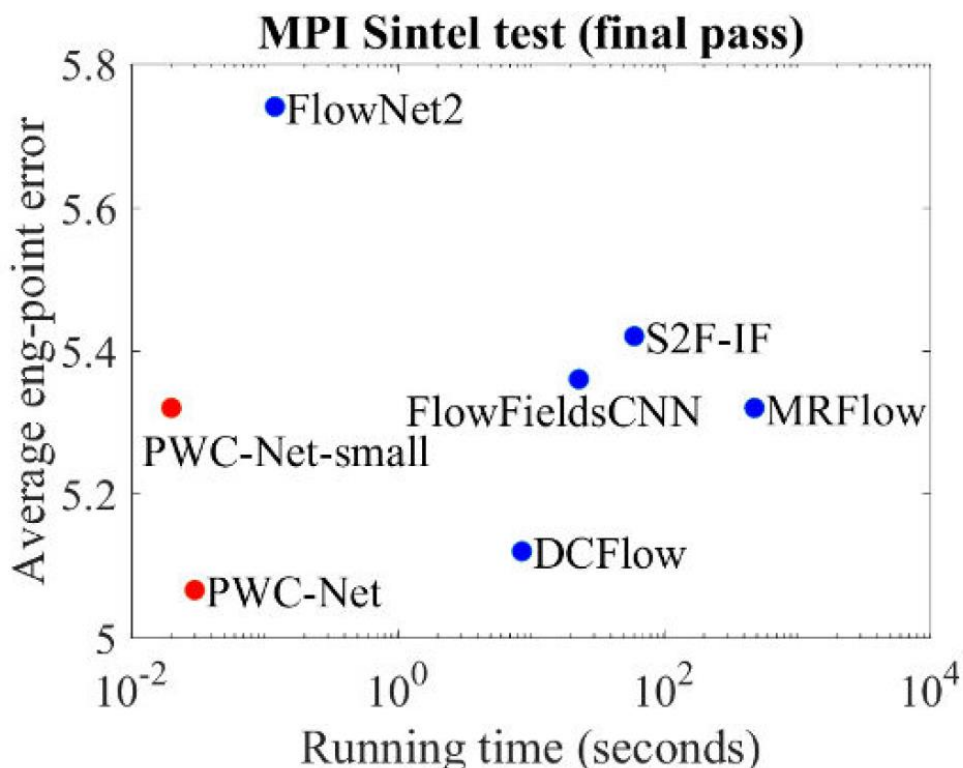


Fischer et al. 2015. <https://arxiv.org/abs/1504.06852>

State-of-the-art optical flow, 2018 (CVPR)



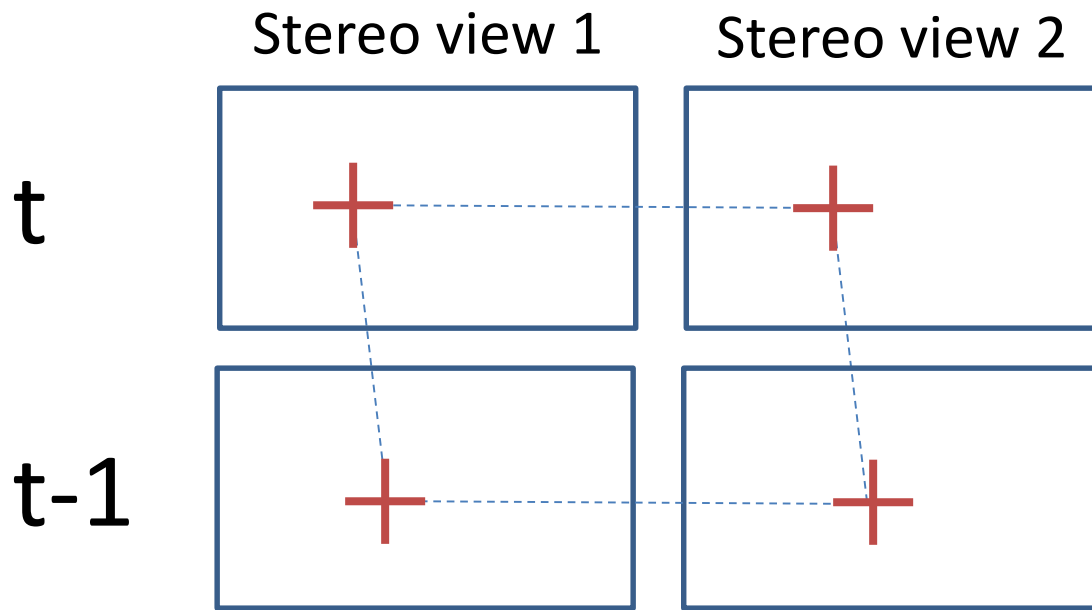
State-of-the-art optical flow, 2018 (CVPR)



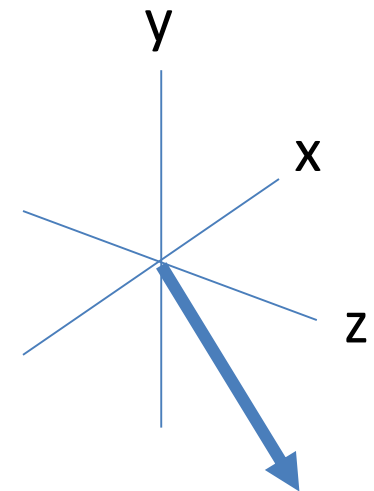
Can we do more? *Scene flow*

Combine spatial stereo & temporal constraints

Recover 3D vectors of world motion

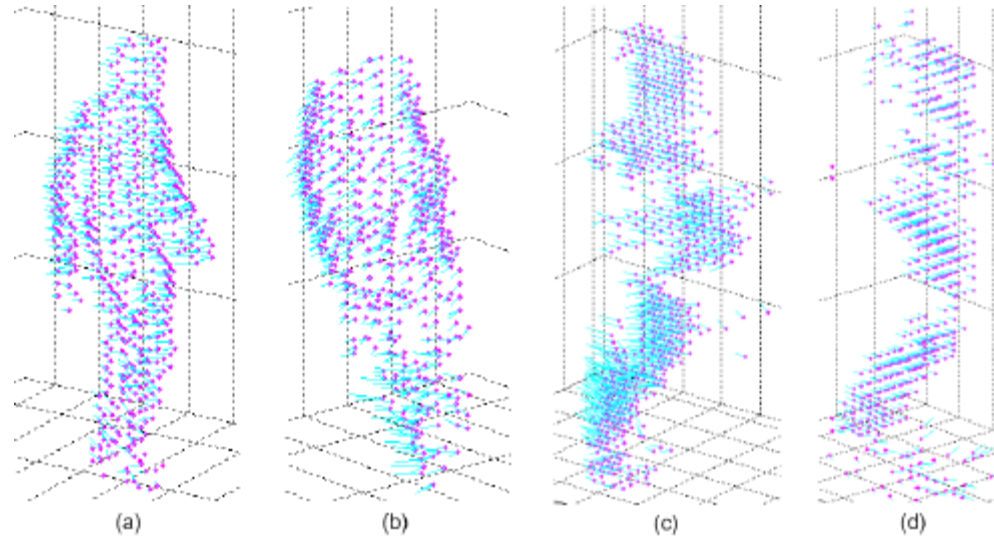


Epipolar constraints
across space and time!



3D world motion
vector per pixel

Scene flow example for human motion



Scene Flow

<https://www.youtube.com/watch?v=RL TK Be6 4>



<https://vision.in.tum.de/research/sceneflow>

[Estimation of Dense Depth Maps and 3D Scene Flow from Stereo Sequences, M. Jaimez et al., TU Munchen]