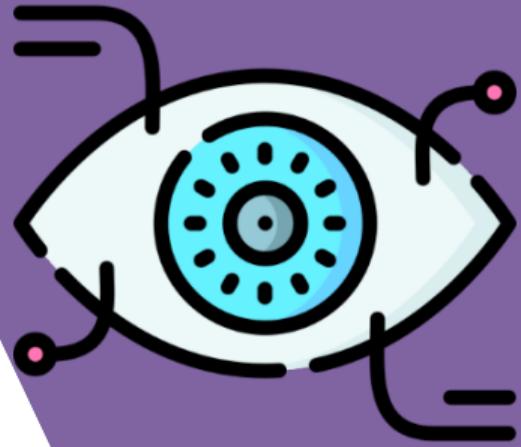


# Computer Vision

A.A. 2024-20245

Lecture 9: Optical Flow



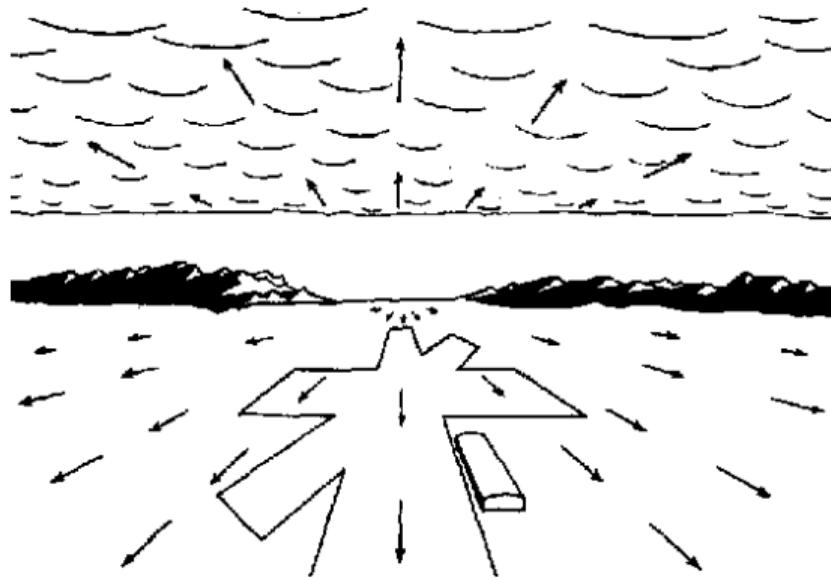
SAPIENZA  
UNIVERSITÀ DI ROMA

ALC<sup>O</sup>R Lab

# References

- Basic reading: Szeliski, Chapter 9.3, 9.4

# Optical Flow



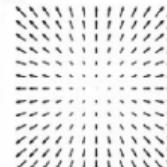
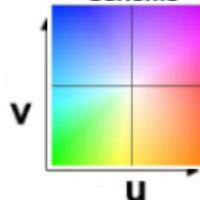
- Optical flow is the **apparent motion** of objects, surfaces, and edges in a visual scene caused by the relative motion between an observer and a scene.

# Optical flow

Optical flow is used to see how every point is moving frame to frame in a video sequence



Encoding Scheme



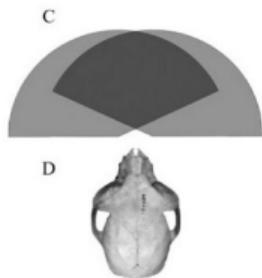
Colour code for visualisation  
HSV encoding scheme: hue, saturation, value

Vector field: displacement, velocity

# Stereo vs. Optical Flow

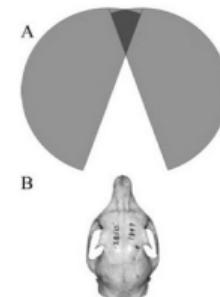
## Stereo

- 2 images at the same time
- Only camera motion
- 1D estimation problem
- Monkeys

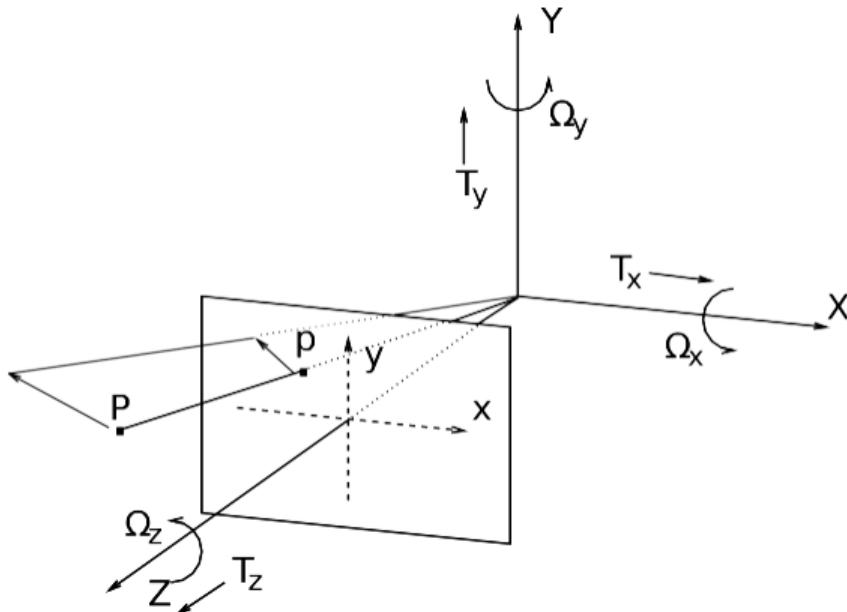


## Optical Flow

- 2 images at 2 time steps
- Camera and object motion
- 2D estimation problem
- Squirrels



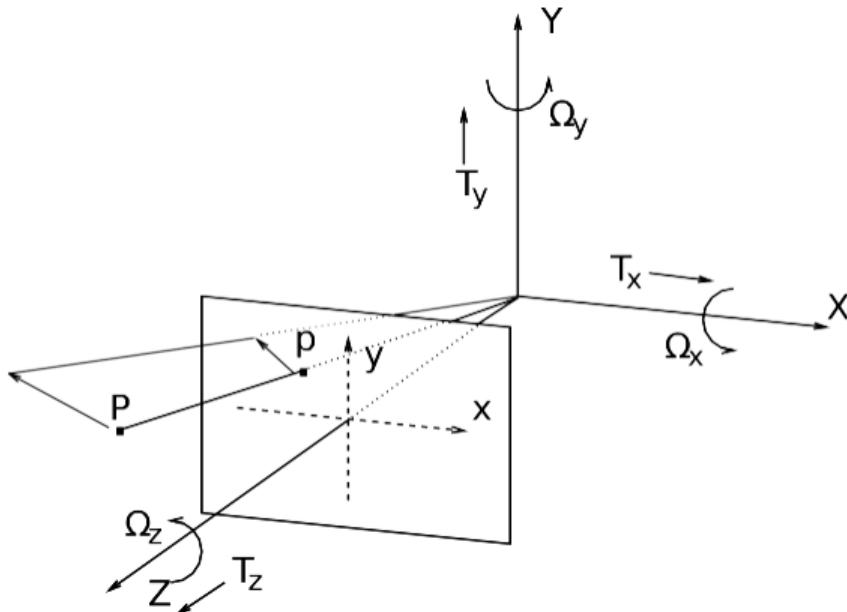
# Optical Flow



## Motion field:

- 2D motion field representing the **projection of the 3D motion** of points in the scene onto the image plane
- Can be the result of camera motion or object motion (or both)

# Optical Flow

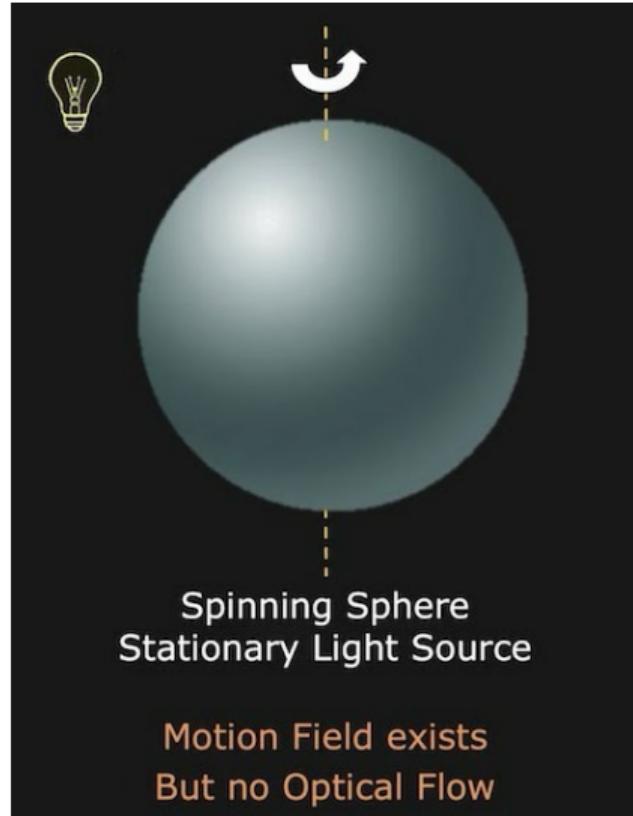


## Optical flow:

- 2D velocity field describing the **apparent motion** in the image (i.e., the displacement of pixels looking “similar”)
- Optical flow  $\neq$  motion field! Why?

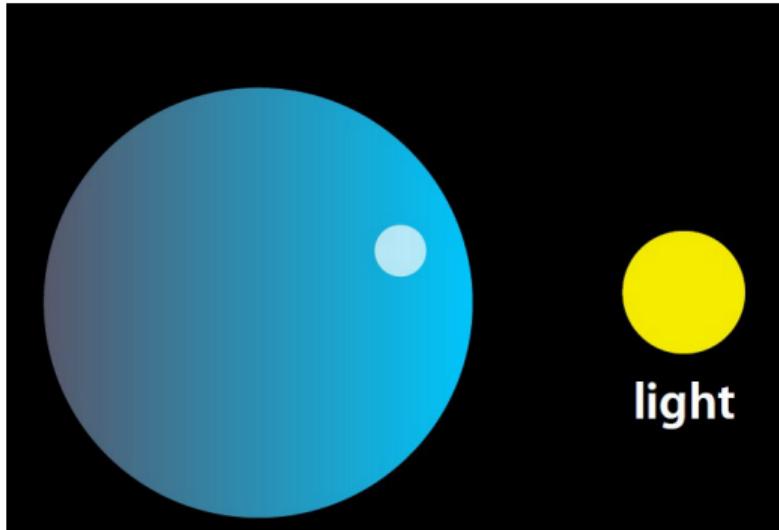
# Thought Experiment

- ▶ Lambertian ball  
**rotating in 3D**
- ▶ What does the 2D  
**motion field** look like?
- ▶ What does the 2D  
**optical flow field** look like?



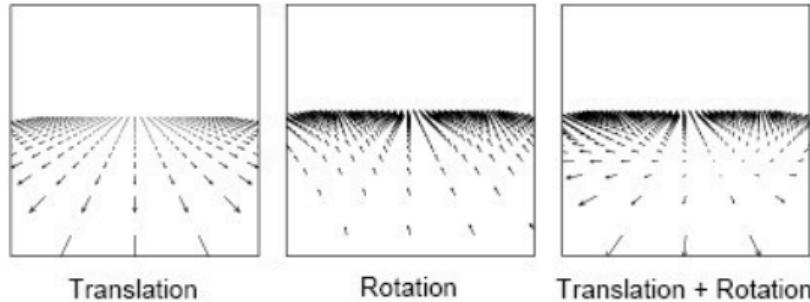
# Thought Experiment

- ▶ Stationary specular ball  
**moving light source**
- ▶ What does the 2D  
**motion field** look like?
- ▶ What does the 2D  
**optical flow field** look like?



Motion field is null, Optical Flow exists

# Optical Flow Field



The optical flow fields tell us something (maybe ambiguous) about:

The **3D structure** of the world

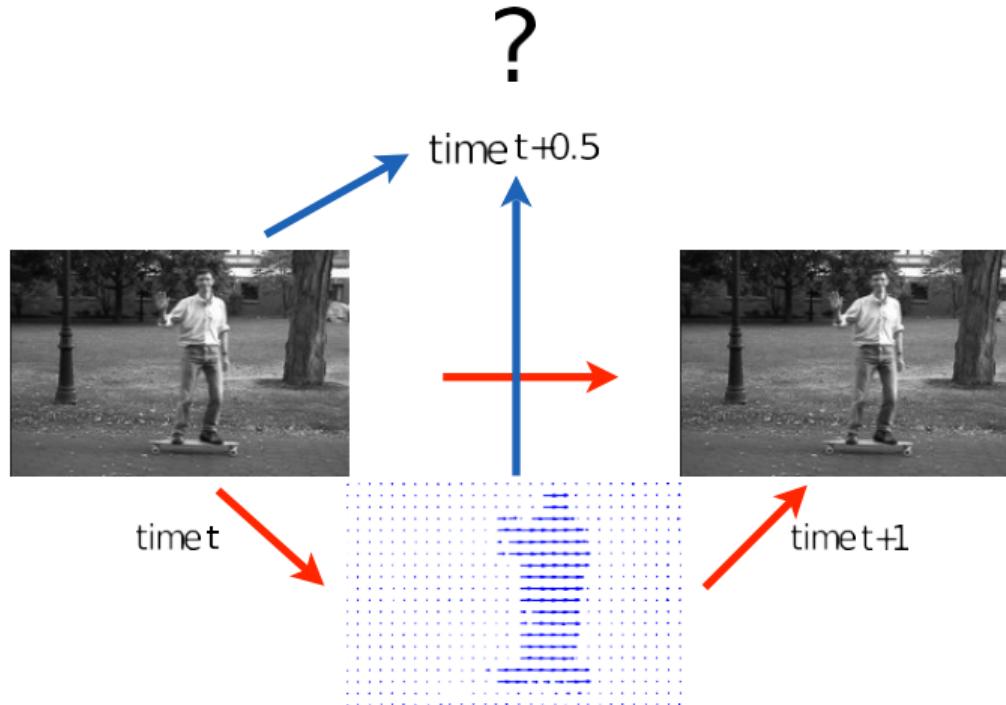
The **motion of objects** in the viewing area

The **motion of the observer** (if any)

In contrast to stereo:

- No epipolar geometry  $\Rightarrow$  **2D estimation problem!**

## Applications: Video Interpolation / Frame Rate Adaption



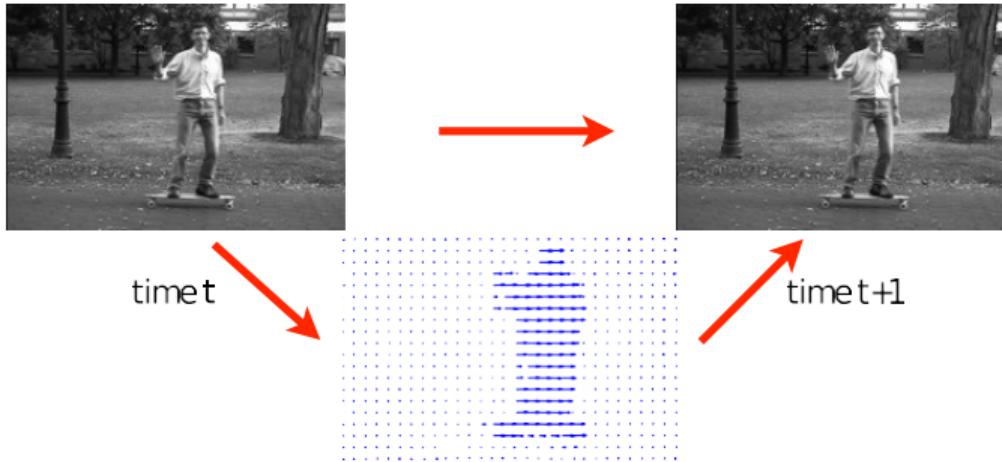
- If we know the image motion we can compute images at intermediate frames

## Applications: Video Interpolation / Frame Rate Adaption



- If we know the image motion we can compute images at intermediate frames

# Applications: Video Compression



- To compress an image sequence, we can predict new frames using the optical flow field and only store how to “fix” the prediction
- Flow fields are smooth, thus easier to compress/store than images!

# Applications: Autonomous Driving



# Uses of motion

The estimation of every pixel in a sequence is a problem with many applications in computer vision

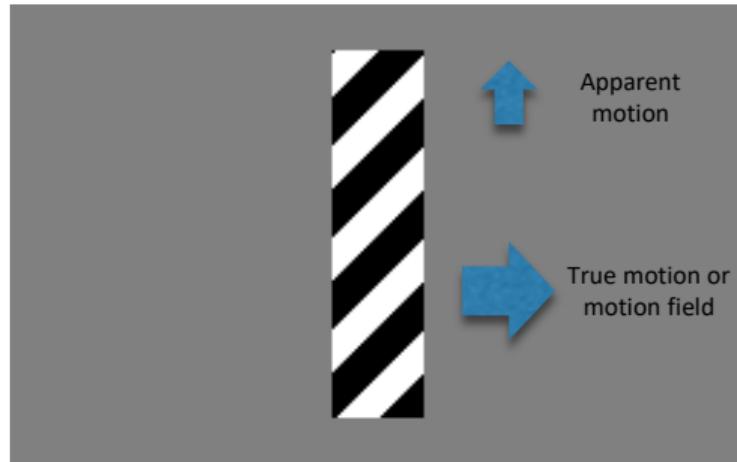
- Improving video quality
  - Motion stabilization
  - Super resolution
- Segmenting objects based on motion cues
- Tracking objects
- Recognizing events and activities

# Aperture Problem



- Aperture problem: A single observation is not enough to determine flow
- Barber Pole: What is the motion field? What is the optic flow field?

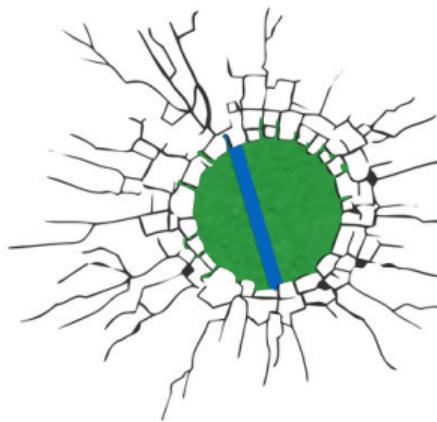
# Barber's pole illusion



# Barber's pole illusion



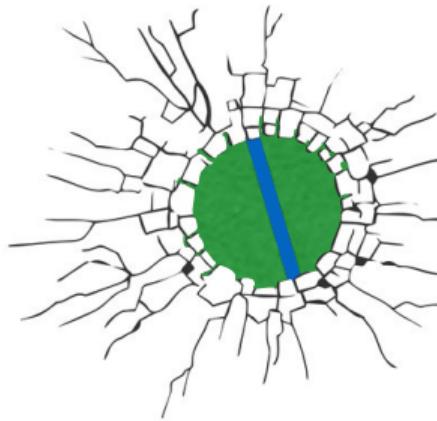
# Aperture Problem



small visible image patch

In which direction is the line moving?

# Aperture Problem

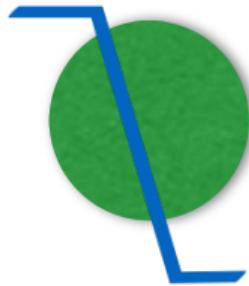


small visible image patch

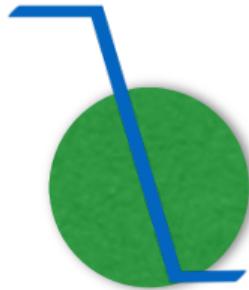
In which direction is the line moving?

# Aperture Problem

Now the full picture

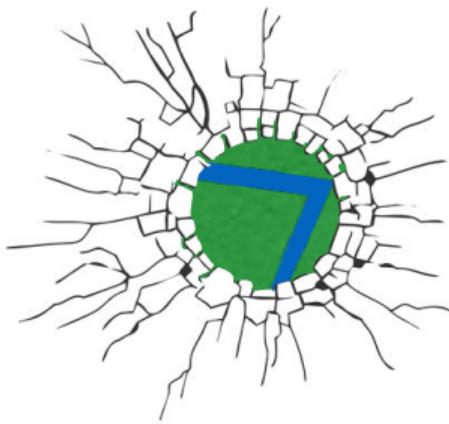


# Aperture Problem



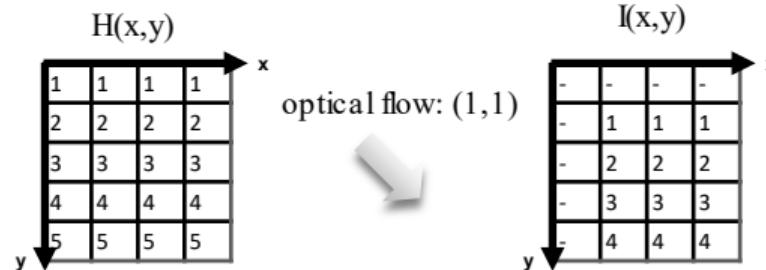


Want patches with different gradients to avoid the aperture problem



Want patches with different gradients to avoid the aperture problem

# Aperture Problem: example



$$\cancel{I_x \dot{a} + I_y v + I_t = 0}$$

Compute gradients

$$I_x(3,3) = 0$$

$$I_y(3,3) = 1$$

$$I_t(3,3) = I(3,3) - H(3,3) = -1$$

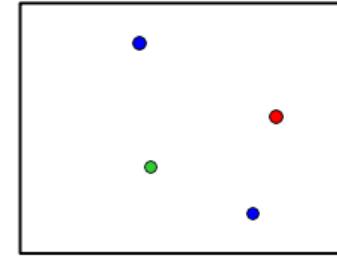
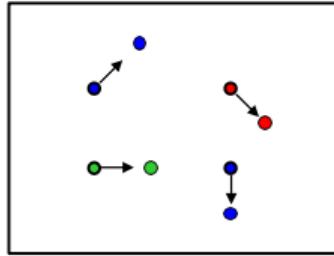
Solution:



$$v = 1$$

We recover the  $v$  of the optical flow but not the  $u$ .

# Estimating optical flow

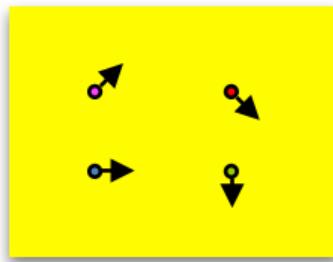


Given two subsequent frames, estimate the apparent motion field  $u(x,y), v(x,y)$  between them.

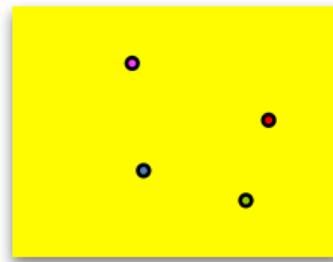
Key assumptions:

- **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

# The approach



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



$$I(x, y, t')$$

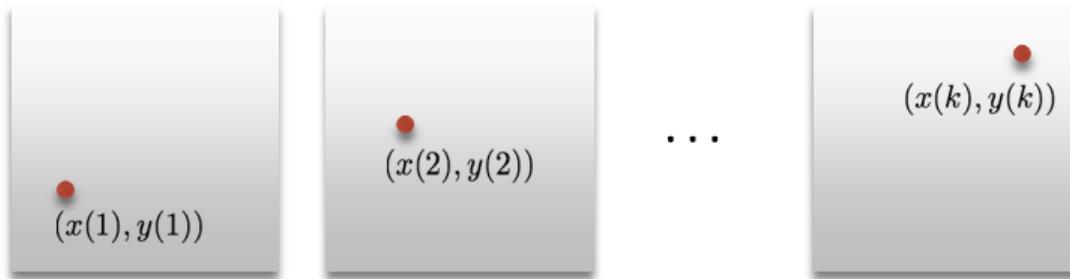
Look for *nearby pixels* with the *same color*

(small motion)

(color constancy)

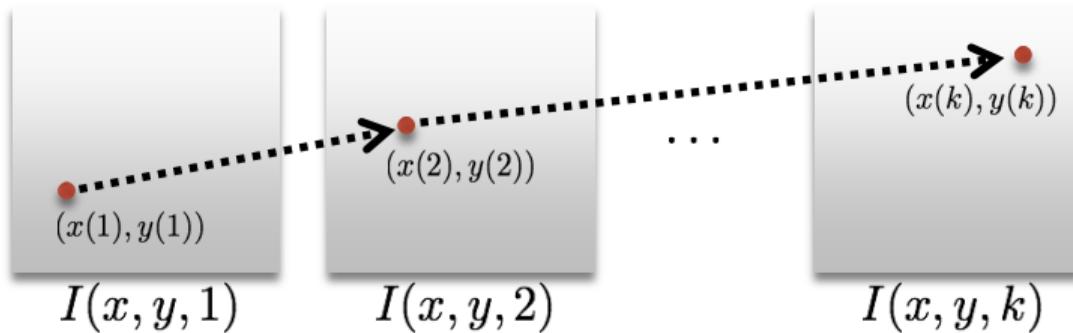
# Brightness constancy Assumption 1

Scene point moving through image sequence



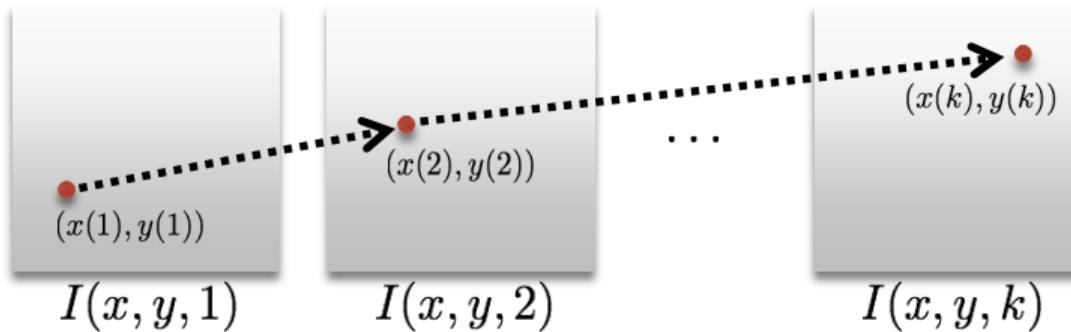
# Brightness constancy Assumption 1

Scene point moving through image sequence



# Brightness constancy Assumption 1

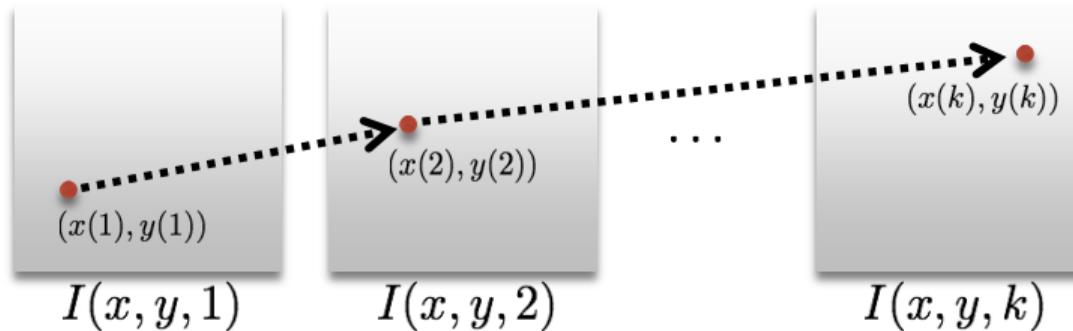
Scene point moving through image sequence



**Assumption: brightness of the point will remain the same**

# Brightness constancy Assumption 1

Scene point moving through image sequence



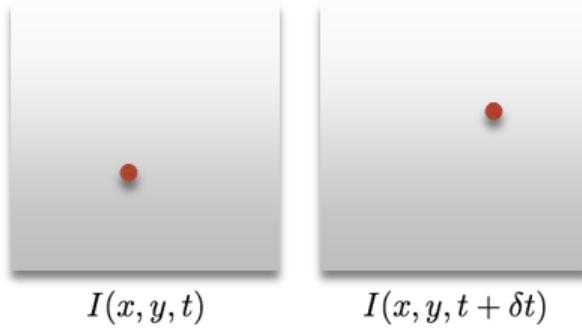
**Assumption: brightness of the point will remain the same**

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

constant

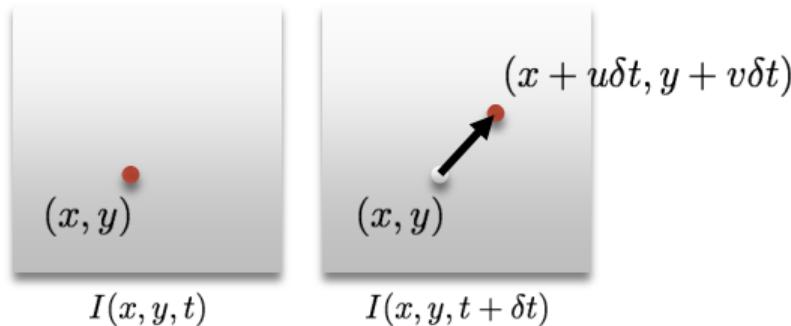
# Small motion

# Assumption 2



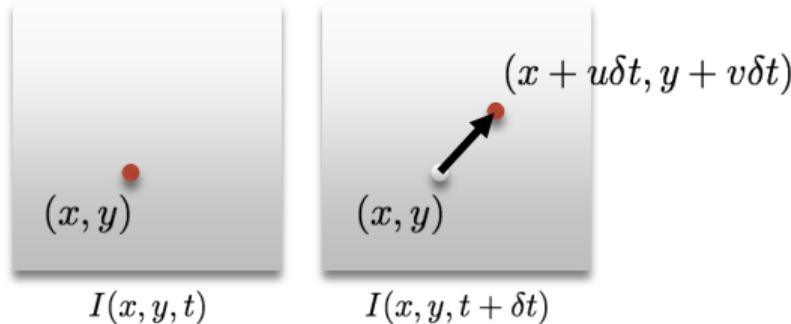
# Small motion

# Assumption 2



# Small motion

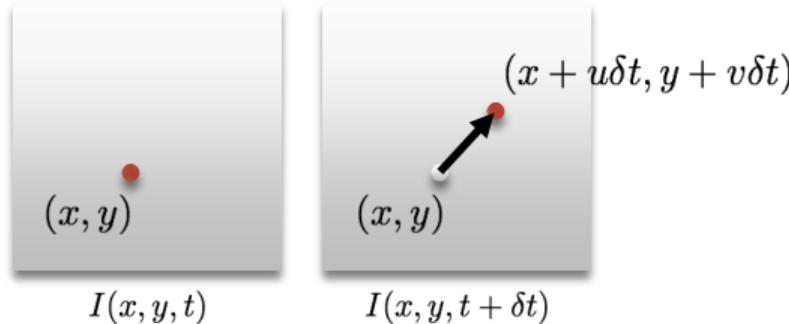
# Assumption 2



Optical flow (velocities):  $(u, v)$       Displacement:  $(\delta x, \delta y) = (u\delta t, v\delta t)$

# Small motion

# Assumption 2



Optical flow (velocities):  $(u, v)$       Displacement:  $(\delta x, \delta y) = (u\delta t, v\delta t)$

For a small space-time step

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

the brightness between two consecutive image frames is the same

# Taylor series expansion

For small space-time step, brightness of a point is the same

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

## Insight:

If the time step is really small, we can *linearize* the intensity function with first order approximation of the Taylor series expansion

## Expand a function as an infinite sum of its derivatives

For one variable

$$f(x + \delta x) = f(x) + \frac{\partial f}{\partial x} \delta x + \frac{\partial^2 f}{\partial x^2} \frac{\delta x^2}{2!} + \dots + \frac{\partial^n f}{\partial x^n} \frac{\delta x^n}{n!}$$

If  $\delta x$  is small:

$$f(x + \delta x) = f(x) + \frac{\partial f}{\partial x} \delta x + O(\delta x^2) \rightarrow \text{Almost Zero}$$

# Taylor series expansion

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



$$I(x, y, t) + \frac{\partial I}{\partial x} \delta x + \frac{\partial I}{\partial y} \delta y + \frac{\partial I}{\partial t} \delta t = I(x, y, t)$$

assuming small motion

fixed point

Partial derivative

$$\frac{\partial I}{\partial x} \delta x + \frac{\partial I}{\partial y} \delta y + \frac{\partial I}{\partial t} \delta t = 0$$

cancel terms

$$\frac{\partial I}{\partial x} \delta x + \frac{\partial I}{\partial y} \delta y + \frac{\partial I}{\partial t} \delta t = 0$$

divide by  $\delta t$   
take limit  $\delta t \rightarrow 0$

$$\frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = 0$$

Brightness Constancy  
Equation

# Brightness constancy equation

$$\frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} = 0$$

**Brightness  
Constancy Equation**

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

(x-flow)                    (y-flow)

shorthand notation

$$\nabla I^\top \mathbf{v} + I_t = 0$$

(1 x 2)                    (2 x 1)

vector form

# Brightness constancy equation

What do the term of the brightness constancy equation represent?

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

flow velocities

↑  
Image gradients  
(at a point p)

↑  
temporal gradient

# Brightness constancy equation

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

$$I_x = \frac{\partial I}{\partial x} \quad I_y = \frac{\partial I}{\partial y}$$

**spatial derivative**

Forward difference

Sobel filter

Derivative-of-Gaussian filter

...

$$I_t = \frac{\partial I}{\partial t}$$

**temporal derivative**

Frame differencing

# Frame differencing

Example of a forward difference

$t$

1	1	1	1	1
1	1	1	1	1
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10
1	10	10	10	10

$t + 1$

1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	10	10	10
1	1	10	10	10

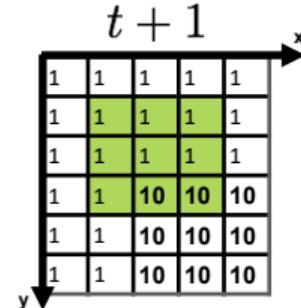
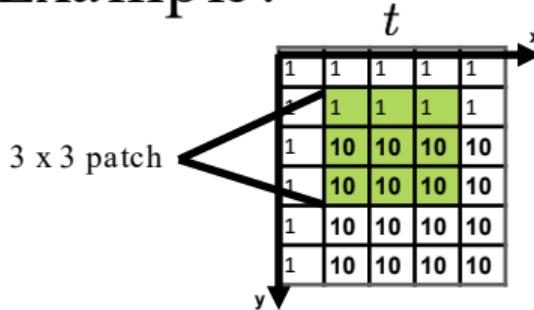
-

$$I_t = \frac{\partial I}{\partial t}$$

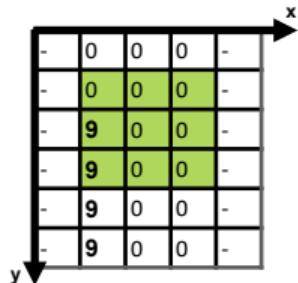
=

0	0	0	0	0
0	0	0	0	0
0	9	9	9	9
0	9	0	0	0
0	9	0	0	0
0	9	0	0	0

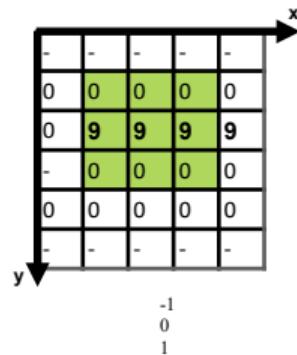
Example:



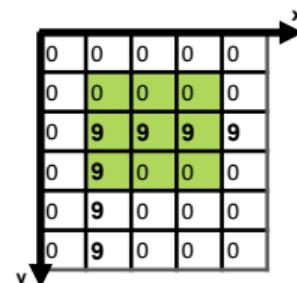
$$I_x = \frac{\partial I}{\partial x}$$



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



$$I_t = \frac{\partial I}{\partial t}$$



-1 0 1

-1  
0  
1

# Brightness constancy equation

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

$$I_x = \frac{\partial I}{\partial x} \quad I_y = \frac{\partial I}{\partial y}$$

**spatial derivative**

Forward difference  
unknown

Sobel filter

Derivative-of-Gaussian filter

...

$$u = \frac{dx}{dt} \quad v = \frac{dy}{dt}$$

**optical flow**

$$I_t = \frac{\partial I}{\partial t}$$

**temporal derivative**

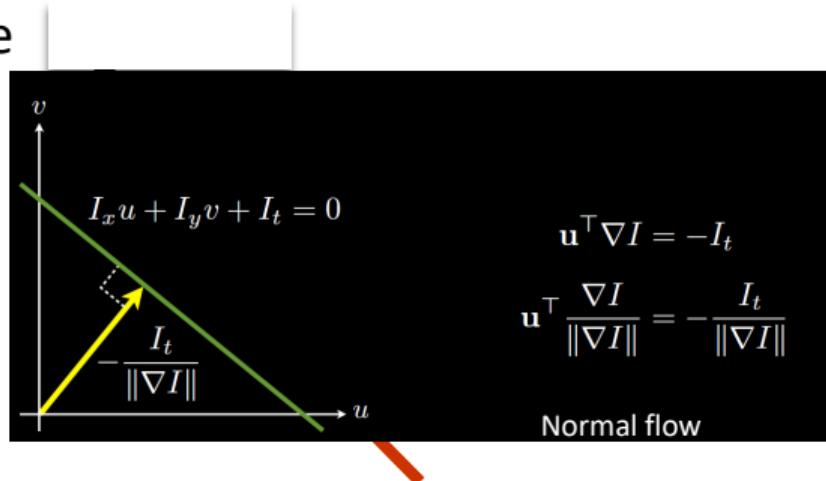
Frame differencing

# Brightness constancy equation

Solution lies on a straight line

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

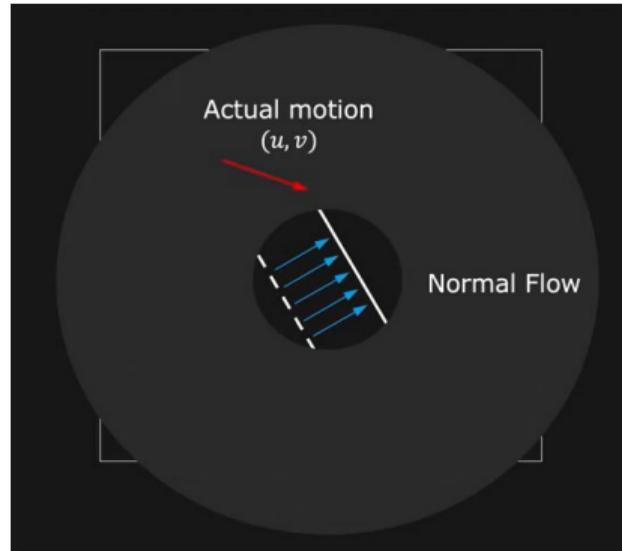
many combinations of  $u$  and  $v$  will satisfy the equality



Optical flow can be split in two components: normal flow and parallel flow

The solution cannot be determined uniquely with a single constraint (a single pixel). The solution is underconstrained  
Where do we get more equations (constraints)?

# Aperture problem



Locally we can only determine normal flow

# Two main methods

- **Lucas-Kanade Optical Flow (1981)**  
method of differences  
**'constant' flow (flow is constant for all pixels)**  
local method (sparse)
- **Horn-Schunck Optical Flow (1981)**  
brightness constancy, small motion  
**'smooth' flow (flow can vary from pixel to pixel)**  
global method (dense)

# Lucas-Kanade Optical Flow

## Assumptions

Assume that the surrounding patch has '**constant flow**'

Neighboring pixels have same displacement: i.e. for each pixel assume motion field, and hence optical flow ( $u, v$ ), is constant within a small neighborhood.

Using a  $5 \times 5$  image patch, gives us 25 equations

$$I_x(\mathbf{p}_1)u + I_y(\mathbf{p}_1)v = -I_t(\mathbf{p}_1)$$

$$I_x(\mathbf{p}_2)u + I_y(\mathbf{p}_2)v = -I_t(\mathbf{p}_2)$$

:

$$I_x(\mathbf{p}_{25})u + I_y(\mathbf{p}_{25})v = -I_t(\mathbf{p}_{25})$$

$n^2$  (25) equations,  
2 unknown: find least  
square solution

# Lukas-Kanade solution

For all points  $(k, l) \in W$ :  $I_x(k, l)u + I_y(k, l)v + I_t(k, l) = 0$

Let the size of window  $W$  be  $n \times n$

In matrix form:

$$\begin{bmatrix} I_x(1,1) & I_y(1,1) \\ I_x(k, l) & I_y(k, l) \\ \vdots & \vdots \\ I_x(n, n) & I_y(n, n) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} I_t(1,1) \\ I_t(k, l) \\ \vdots \\ I_t(n, n) \end{bmatrix}$$

$A$        $\hat{\mathbf{u}}$        $B$   
(Known)    (Unknown)    (Known)  
 $n^2 \times 2$                    $2 \times 1$                    $n^2 \times 1$

$n^2$  equations, 2 unknown: find least square solution

# Least square solution

Solve linear system:  $A\mathbf{u} = B$

$A^T A \mathbf{u} = A^T B$  (Least-Squares using  
Pseudo-Inverse)

In matrix form:

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

Indices  $(k, l)$   
not written  
for simplicity

$\begin{matrix} A^T A & \mathbf{u} & A^T B \\ (\text{Known}) & (\text{Unknown}) & (\text{Known}) \\ 2 \times 2 & 2 \times 1 & 2 \times 1 \end{matrix}$

$$\mathbf{u} = (A^T A)^{-1} A^T B$$

Fast and Easy to Solve

# Least square solution

Solve linear system  $Ax=b$

- that it is equivalent to write (least squares using pseudo-inverse)

Lucas-Kanade Optical Flow

$$\begin{matrix} A^\top A & x & A^\top b \\ \left[ \begin{array}{cc} \sum_{p \in P} I_x I_x & \sum_{p \in P} I_x I_y \\ \sum_{p \in P} I_y I_x & \sum_{p \in P} I_y I_y \end{array} \right] & \left[ \begin{array}{c} u \\ v \end{array} \right] = - \left[ \begin{array}{c} \sum_{p \in P} I_x I_t \\ \sum_{p \in P} I_y I_t \end{array} \right] \\ 2 \times 2 & 2 \times 1 & 2 \times 1 \end{matrix}$$

where the summation is over each pixel  $p$  in patch  $P$

$$x = (A^\top A)^{-1} A^\top b$$

# When does optical flow estimation work?

$$x = (A^\top A)^{-1} A^\top b$$

$A^\top A$  should be invertible ( $\det \neq 0$ )

$A^\top A$  should be well conditioned

$\lambda_1$  and  $\lambda_2$  should not be too small (both are significant enough)

$\lambda_1$  should not be too large respect to  $\lambda_2$  ( $\lambda_1$ =larger eigenvalue)

# When does optical flow estimation work?

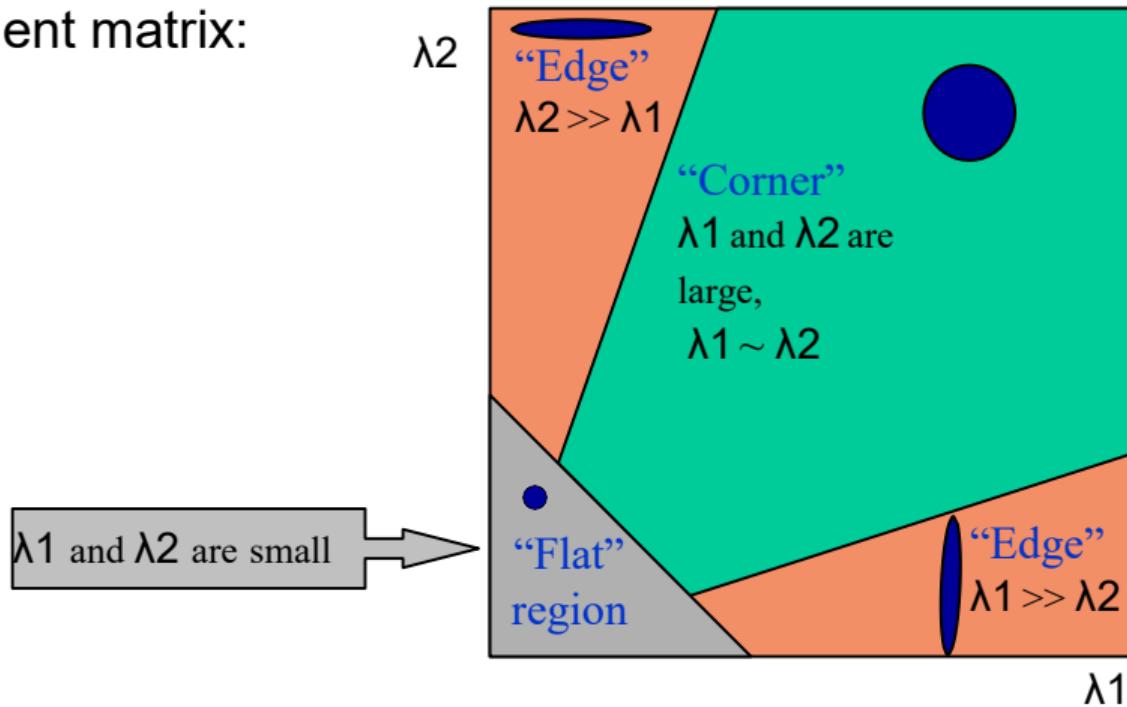
$A^T A$  is the *second moment matrix* (Harris corner detector)

$$A^T A = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix}$$

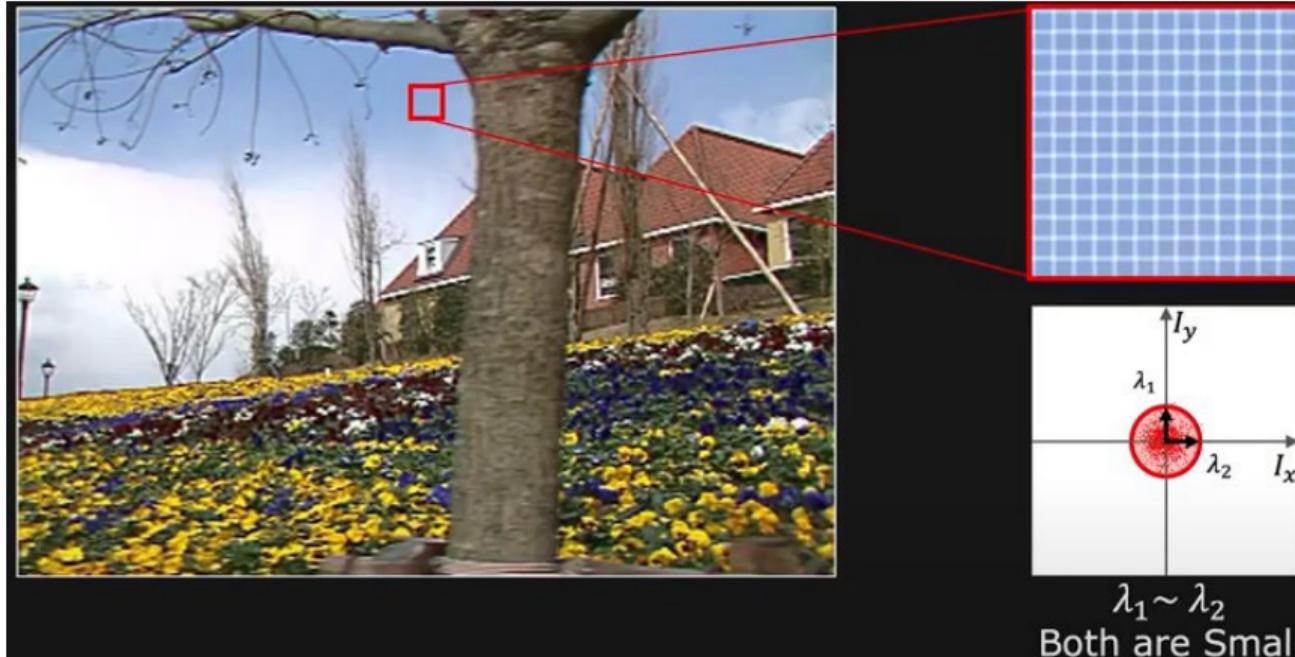
- Eigenvectors and eigenvalues of  $A^T A$  relate to edge direction and magnitude
  - The eigenvector associated with the larger eigenvalue points in the direction of fastest intensity change
  - The other eigenvector is orthogonal to it

# Interpreting the eigenvalues

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

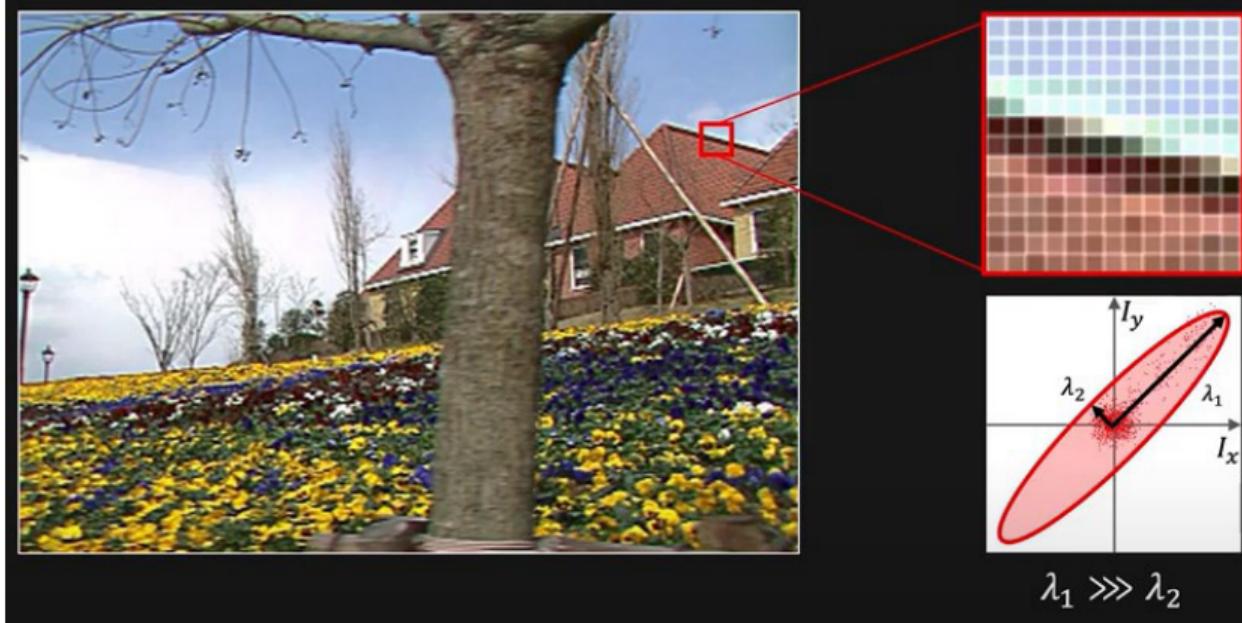


# Low-texture region



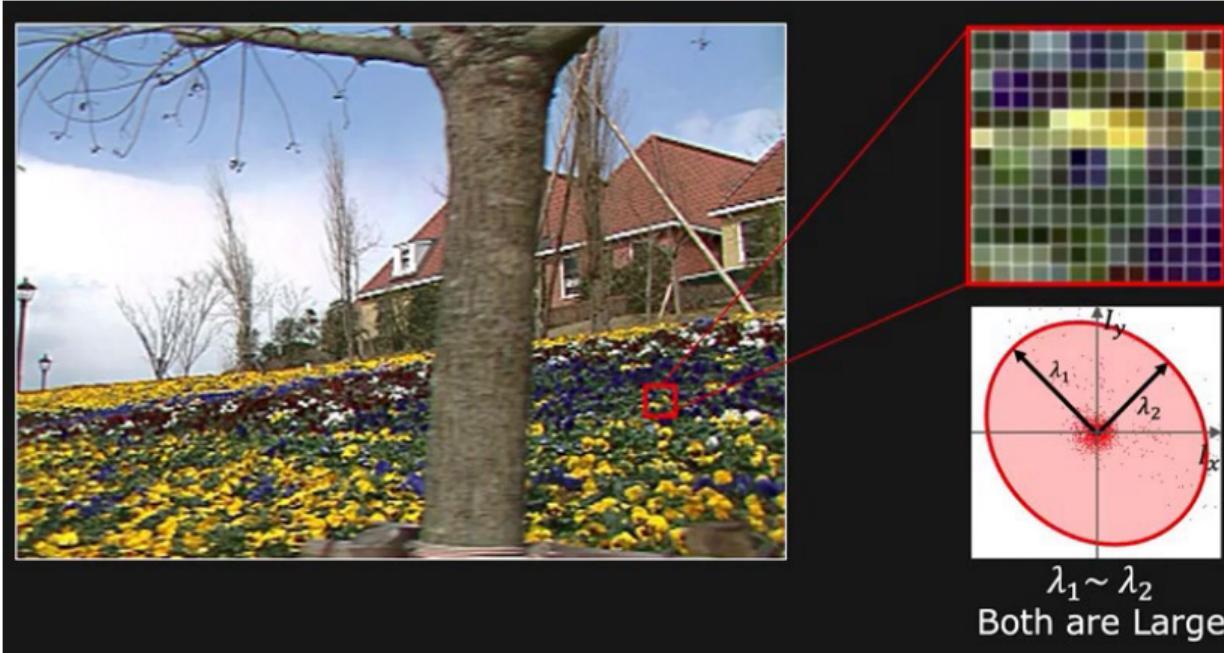
- Equations for all pixels in window are more or less the same
- Gradients have small magnitude
- Can not reliably compute the flow

# Edges



- gradients very large or very small
- prominent gradient in one direction → badly conditioned
- can not reliably compute the flow

# High-texture region



- gradients are different, large magnitudes
- well conditioned
- can reliably compute optical flow

# Implications

Corners and high texture regions are when  $\lambda_1, \lambda_2$  are big; this is also when Lucas-Kanade optical flow works best

Corners and high texture regions are regions with two different directions of gradient (at least)

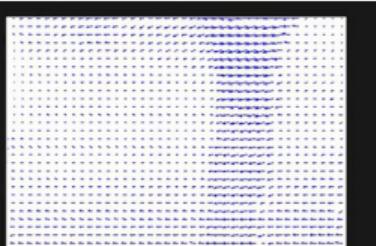
Corners and high texture regions are good places to compute flow

# Lukas-Kanade method

- Decide for local neighborhood of  $W \times W$  pixels and apply uniformly in frames (smoothing frame first with Gaussian filter with a small standard deviation  $\sigma=1.5$ )
- At frame  $t, t+1$  calculate the derivatives  $I_x, I_y, I_t$
- For each couple of frames obtain equational system and solve in the least square sense calculating the eigenvalues
- Plot the optical flow vectors (directions and magnitude)



Image Sequence  
(2 frames)



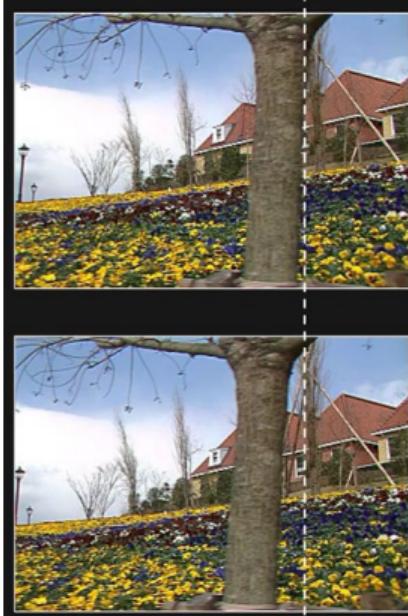
Optical Flow

# Recap

## Key assumptions

- **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



Taylor Series approximation of  
 $I(x + \delta x, y + \delta y, t + \delta t)$  is not valid

Case of large motion

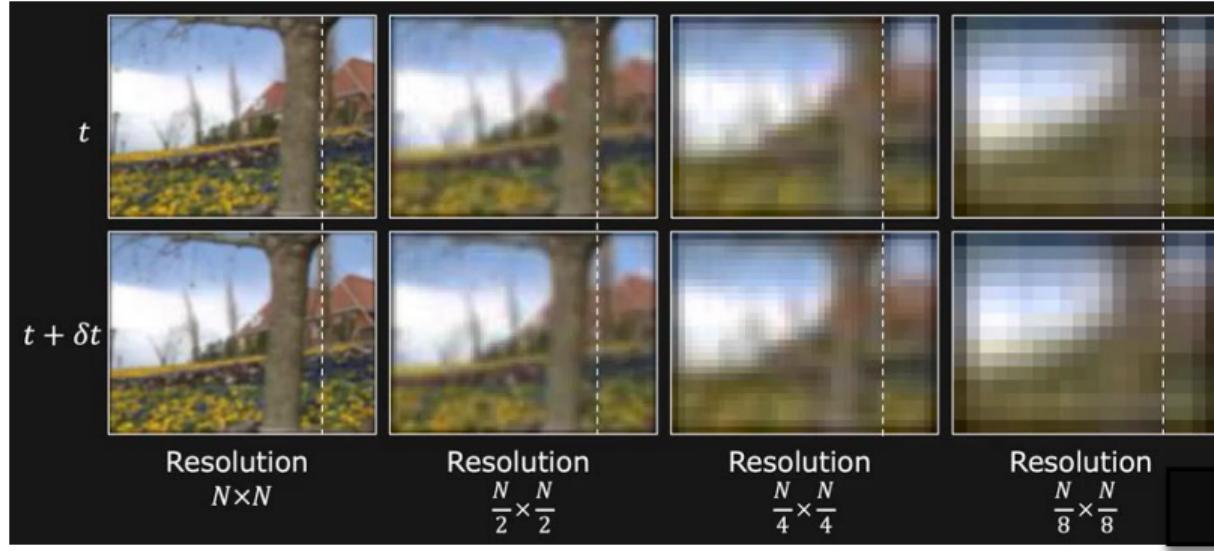
Our simple linear  
constraint equation not valid

$$I_x u + I_y v + I_t \neq 0$$

Is this motion small enough?

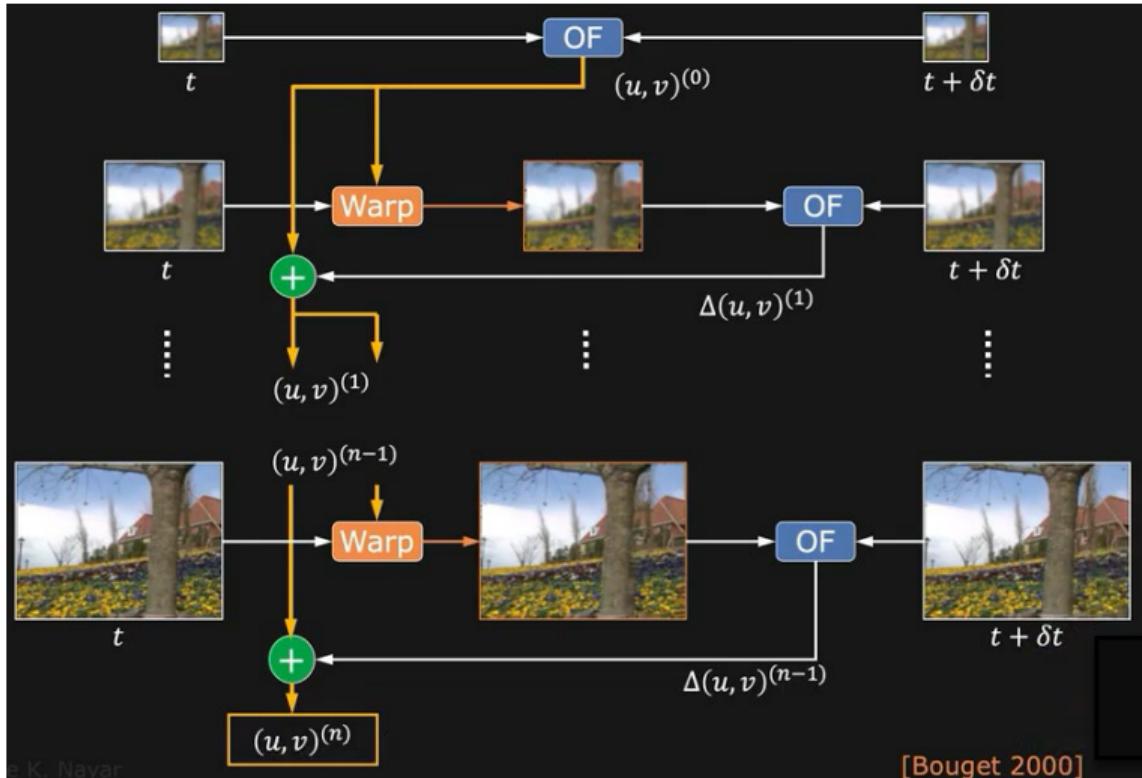
- Probably not: it's much larger than one pixel (2<sup>nd</sup> order terms dominate)
- How might we solve this problem?

# Reduce the resolution

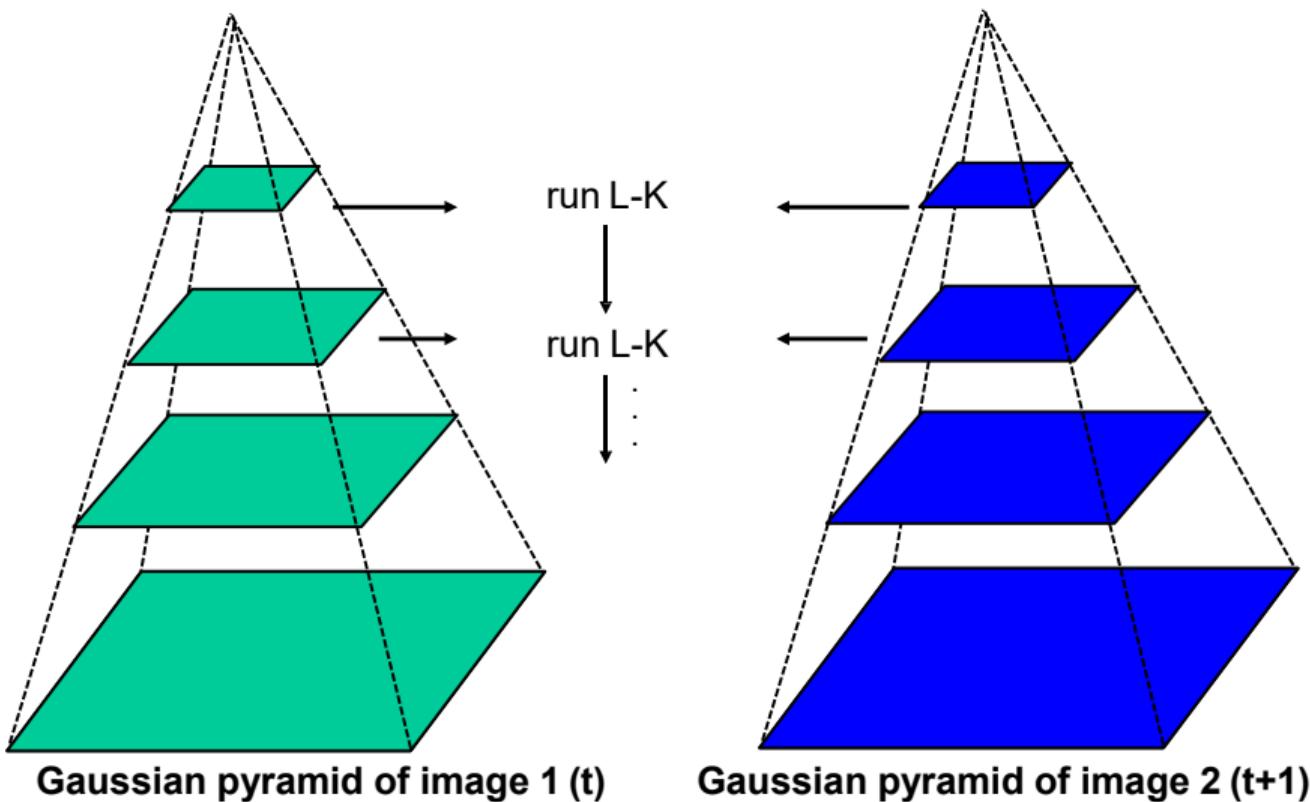


At lowest resolution, motion  $\leq 1$  pixel

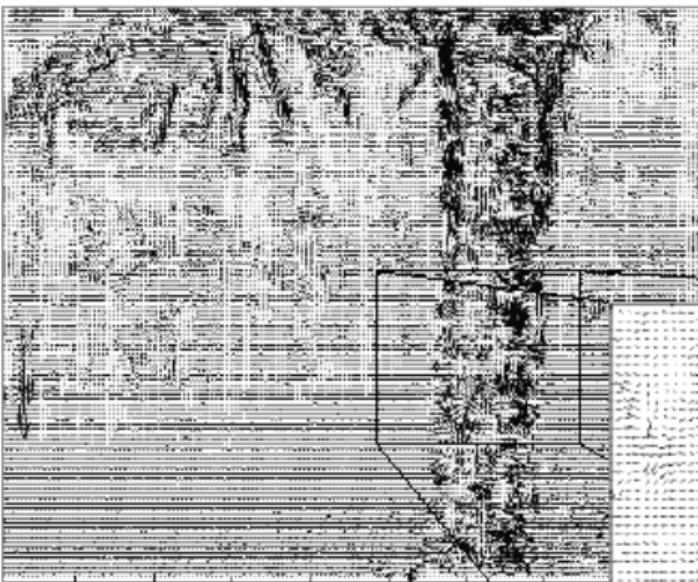
# Coarse-to-fine Flow Estimation



# Coarse-to-fine optical flow estimation

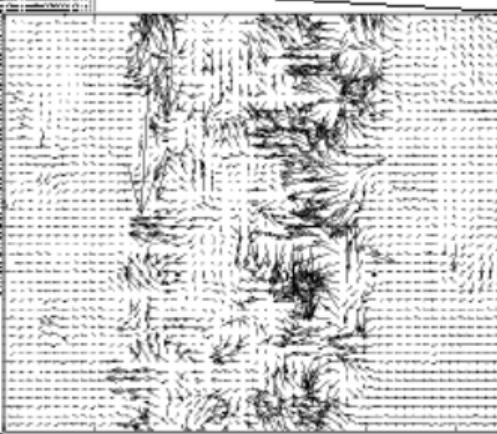


# Optical Flow Results

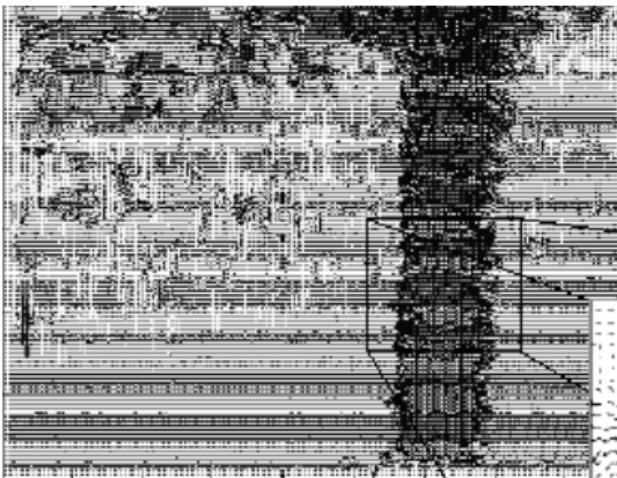


Lucas-Kanade  
without pyramids

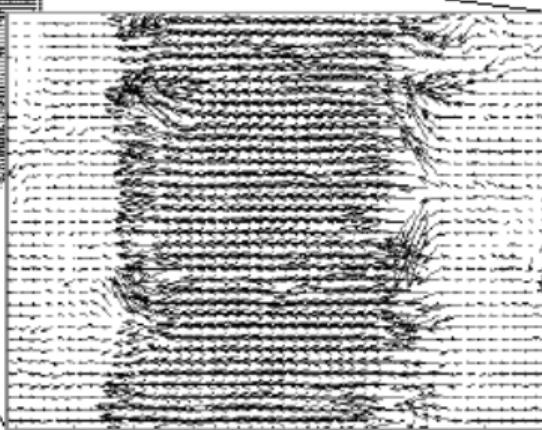
Fails in areas of large motion



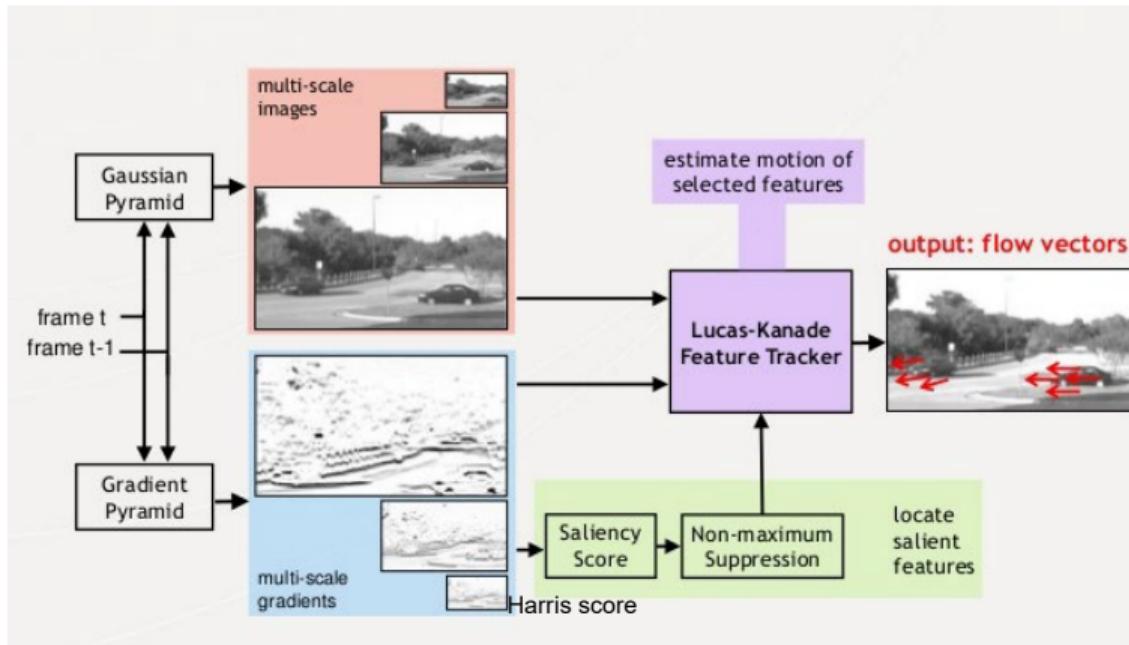
# Optical Flow Results



Lucas-Kanade with Pyramids



# Lucas- Kanade feature tracker



# Motion segmentation

How do we represent the motion in this scene?

- Break image sequence into “layers” each of which has a coherent (affine) motion (since points move like their neighbors for the spatial coherence assumption)



# How do we estimate the layers?

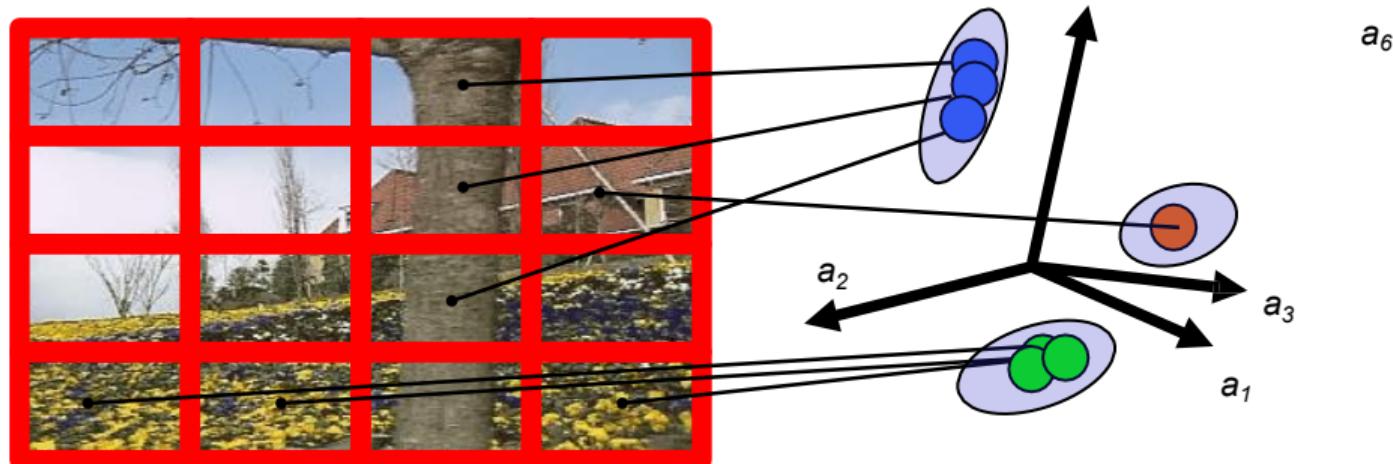
## 1. Obtain a set of initial affine motion hypotheses

- Divide the image into blocks and estimate affine motion parameters in each block by least squares
  - Eliminate hypotheses with high residual error

## 2. Map into motion parameter space

## 3. 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



# How do we estimate the layers?

## 1. Obtain a set of initial affine motion hypotheses

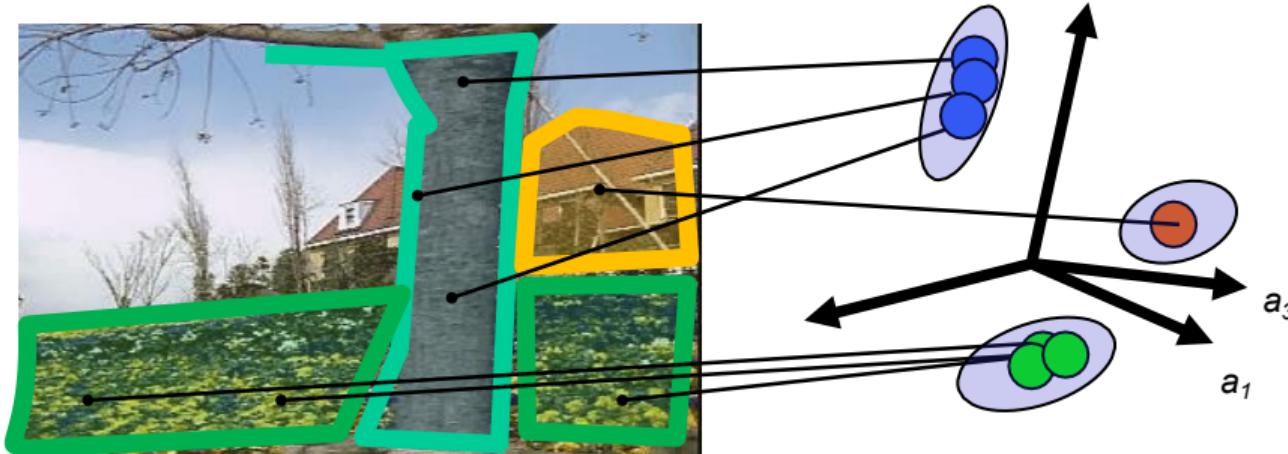
- Divide the image into blocks and estimate affine motion parameters in each block by least squares
  - Eliminate hypotheses with high residual error

## 2. Map into motion parameter space

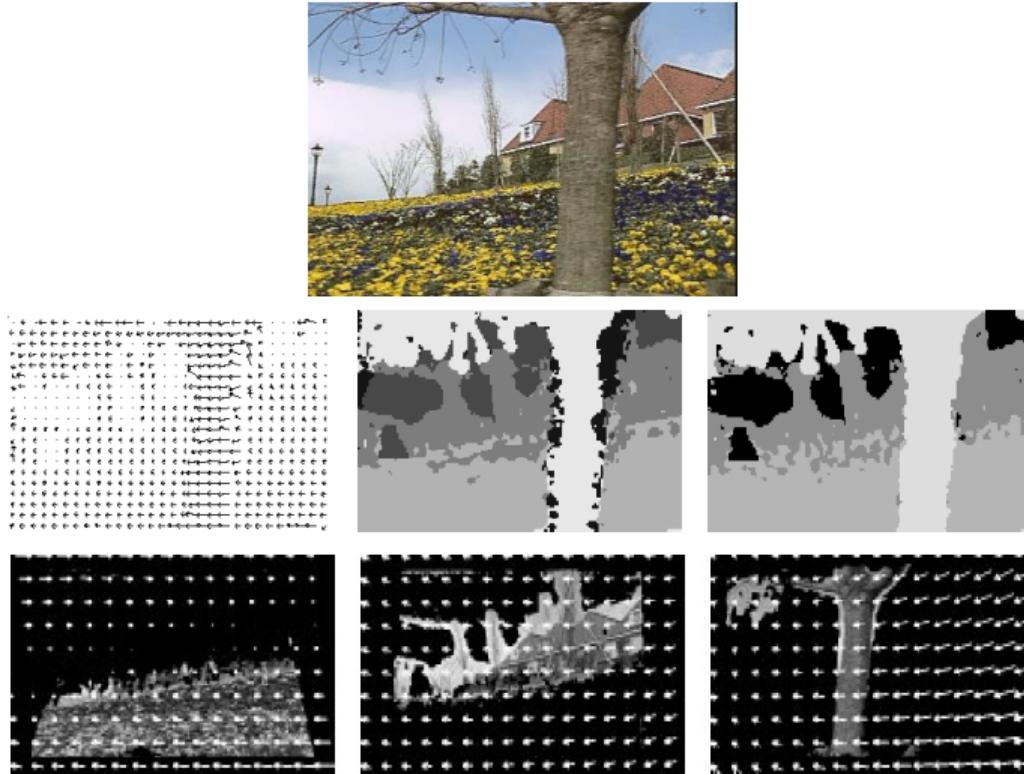
## 3. Perform k-means clustering on affine motion parameters

- 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

## 4. Assign each pixel to best hypothesis --- iterate



# Example result



# Determining Optical Flow

# Horn-Schunck Optical Flow

- ▶ Consider the image  $I$  as a function of continuous variables  $x, y, t$
- ▶ Consider  $u(x, y)$  and  $v(x, y)$  as continuous flow fields (functions)
- ▶ Minimize the following **energy functional**

$$\begin{aligned} E(u, v) &= \iint \underbrace{(I(x + u(x, y), y + v(x, y), t + 1) - I(x, y, t))^2}_{\text{quadratic penalty for brightness change}} \\ &\quad + \lambda \cdot \underbrace{\left( \|\nabla u(x, y)\|^2 + \|\nabla v(x, y)\|^2 \right)}_{\text{quadratic penalty for flow change}} dx dy \end{aligned}$$

with regularization parameter  $\lambda$  and gradient  $\nabla = \left( \frac{\partial}{\partial x}, \frac{\partial}{\partial y} \right)$ .

# Horn-Schunck Optical Flow

$$\begin{aligned} E(u, v) = & \iint (I(x + u(x, y), y + v(x, y), t + 1) - I(x, y, t))^2 \\ & + \lambda \cdot \left( \|\nabla u(x, y)\|^2 + \|\nabla v(x, y)\|^2 \right) dx dy \end{aligned}$$

- ▶ Minimizing this directly is a hard problem because the energy is highly **non-convex** and has many local optima
- ▶ Solution: **linearize** the brightness constancy assumption

# Horn-Schunck Optical Flow

First-Order Multivariable Taylor Series:

$$f(x, y) \stackrel{a,b}{\approx} f(a, b) + \frac{\partial f(a, b)}{\partial x}(x - a) + \frac{\partial f(a, b)}{\partial y}(y - b)$$

Therefore, we have:

$$\begin{aligned} & I(x + u(x, y), y + v(x, y), t + 1) \\ & \stackrel{x,y,t}{\approx} I(x, y, t) + I_x(x, y, t)(x + u(x, y) - x) \\ & \quad + I_y(x, y, t)(y + v(x, y) - y) + I_t(x, y, t)(t + 1 - t) \\ & = I(x, y, t) + I_x(x, y, t)u(x, y) + I_y(x, y, tv(x, y)) + I_t(x, y, t) \end{aligned}$$

# Horn-Schunck Optical Flow

Thus, the optical flow energy functional

$$\begin{aligned} E(u, v) &= \iint (I(x + u(x, y), y + v(x, y), t + 1) - I(x, y, t))^2 \\ &\quad + \lambda \cdot (\|\nabla u(x, y)\|^2 + \|\nabla v(x, y)\|^2) dx dy \end{aligned}$$

is approximated by the following **linearized equation**:

$$\begin{aligned} E(u, v) &= \iint (I_x(x, y, t)u(x, y) + I_y(x, y, t)v(x, y) + I_t(x, y, t))^2 \\ &\quad + \lambda \cdot (\|\nabla u(x, y)\|^2 + \|\nabla v(x, y)\|^2) dx dy \end{aligned}$$

# Horn-Schunck Optical Flow

Spatial discretization of the equation

$$\begin{aligned} E(u, v) = & \iint (I_x(x, y, t)u(x, y) + I_y(x, y, t)v(x, y) + I_t(x, y, t))^2 \\ & + \lambda \cdot (\|\nabla u(x, y)\|^2 + \|\nabla v(x, y)\|^2) dx dy \end{aligned}$$

leads to the following **discretized objective:**

$$\begin{aligned} E(\mathbf{U}, \mathbf{V}) = & \sum_{x,y} (I_x(x, y) u_{x,y} + I_y(x, y) v_{x,y} + I_t(x, y))^2 \\ & + \lambda \cdot ((u_{x,y} - u_{x+1,y})^2 + (u_{x,y} - u_{x,y+1})^2 + (v_{x,y} - v_{x+1,y})^2 + (v_{x,y} - v_{x,y+1})^2) \end{aligned}$$

This objective is quadratic in the flow maps  $\mathbf{U}, \mathbf{V}$  and thus has a unique optimum.

# Horn-Schunck Optical Flow

$$\begin{aligned} E(\mathbf{U}, \mathbf{V}) = & \sum_{x,y} (I_x(x, y) u_{x,y} + I_y(x, y) v_{x,y} + I_t(x, y))^2 \\ & + \lambda \cdot ((u_{x,y} - u_{x+1,y})^2 + (u_{x,y} - u_{x,y+1})^2 + \\ & \quad (v_{x,y} - v_{x+1,y})^2 + (v_{x,y} - v_{x,y+1})^2) \end{aligned}$$

- ▶ Differentiate wrt.  $\mathbf{U}, \mathbf{V}$  and set the gradient to 0
- ▶ Results in a huge but **sparse linear system**
- ▶ Can be solved using **standard techniques** (e.g., Gauss-Seidel, SOR)
- ▶ However, linearization works only for **small motions**
- ▶ Solution: Iterative estimation & warping, coarse-to-fine estimation

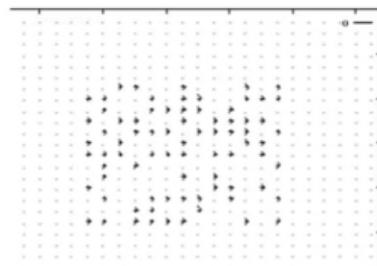
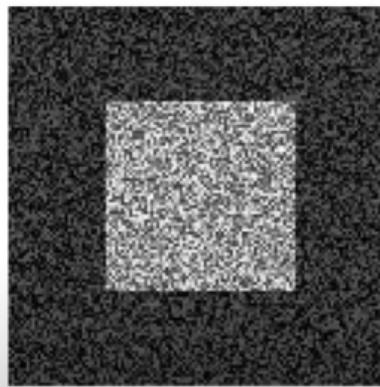
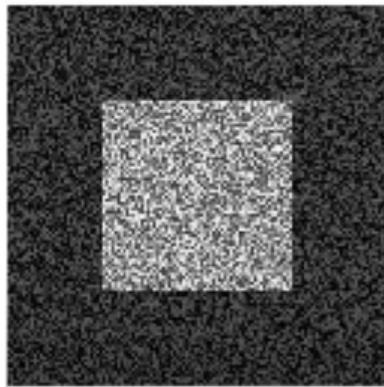
# Horn-Schunck Optical Flow Algorithm

1. Precompute image gradients       $I_y \quad I_x$
2. Precompute temporal gradients     $I_t$
3. Initialize flow field                 $\mathbf{u} = \mathbf{0}$   
 $\mathbf{v} = \mathbf{0}$
4. While not converged

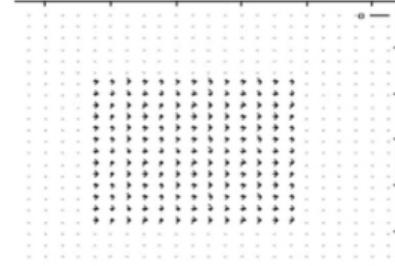
Compute flow field updates for each pixel:

$$\hat{u}_{kl} = \bar{u}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_x \quad \hat{v}_{kl} = \bar{v}_{kl} - \frac{I_x \bar{u}_{kl} + I_y \bar{v}_{kl} + I_t}{\lambda^{-1} + I_x^2 + I_y^2} I_y$$

# Horn – Shunck example



One iteration



10 iterations

# Horn & Schunck Optical Flow

- The HS results are quite plausible already
- However, the flow is very smooth, i.e., to overcome ambiguities we need to set  $\lambda$  to a high value which **oversmooths** flow discontinuities. Why?
- We use a **quadratic penalty** for penalizing changes in the flow
- This does not allow for **discontinuities** in the flow field
- In other words, it penalizes large changes too much and causes oversmoothing

# Robust Estimation of Optical Flow

# Probabilistic Interpretation

The HS optimization problem can be interpreted as **MAP inference in a MRF**:

$$p(\mathbf{U}, \mathbf{V}) = \frac{1}{Z} \exp \{-E(\mathbf{U}, \mathbf{V})\}$$

Maximum A Posteriori inference in Markov Random Fields method to find the most probable configuration of variables given some observed data

with the following **Gibbs energy**:

$$\begin{aligned} E(\mathbf{U}, \mathbf{V}) &= \sum_{x,y} (I_x(x, y) u_{x,y} + I_y(x, y) v_{x,y} + I_t(x, y))^2 \\ &+ \lambda \cdot ((u_{x,y} - u_{x+1,y})^2 + (u_{x,y} - u_{x,y+1})^2 + \\ &\quad (v_{x,y} - v_{x+1,y})^2 + (v_{x,y} - v_{x,y+1})^2) \end{aligned}$$

Remark: As  $\mathbf{U}, \mathbf{V}$  are continuous, we solve inference with gradient descent (not BP)

# Probabilistic Interpretation

The HS optimization problem can be interpreted as **MAP inference in a MRF**:

$$p(\mathbf{U}, \mathbf{V}) = \frac{1}{Z} \exp \{-E(\mathbf{U}, \mathbf{V})\}$$

The corresponding **Gibbs distribution** is Gaussian:

$$\begin{aligned} p(\mathbf{U}, \mathbf{V}) &\propto \prod_{x,y} \exp \left\{ - (I_x(x, y) u_{x,y} + I_y(x, y) v_{x,y} + I_t(x, y))^2 \right\} \\ &\quad \times \exp \left\{ -\lambda (u_{x,y} - u_{x+1,y})^2 \right\} \times \exp \left\{ -\lambda (u_{x,y} - u_{x,y+1})^2 \right\} \\ &\quad \times \exp \left\{ -\lambda (v_{x,y} - v_{x+1,y})^2 \right\} \times \exp \left\{ -\lambda (v_{x,y} - v_{x,y+1})^2 \right\} \end{aligned}$$

Remark: As  $\mathbf{U}, \mathbf{V}$  are continuous, we solve inference with gradient descent (not BP)

# Robust Regularization

**Quadratic penalties translate to Gaussian distributions:**

$$\begin{aligned} p(\mathbf{U}, \mathbf{V}) &\propto \prod_{x,y} \exp \left\{ - (I_x(x,y) u_{x,y} + I_y(x,y) v_{x,y} + I_t(x,y))^2 \right\} \\ &\quad \times \exp \left\{ -\lambda (u_{x,y} - u_{x+1,y})^2 \right\} \times \exp \left\{ -\lambda (u_{x,y} - u_{x,y+1})^2 \right\} \\ &\quad \times \exp \left\{ -\lambda (v_{x,y} - v_{x+1,y})^2 \right\} \times \exp \left\{ -\lambda (v_{x,y} - v_{x,y+1})^2 \right\} \end{aligned}$$

- ▶ Both assumptions are invalid (e.g., discontinuities at object boundaries). Why?
- ▶ Gaussian distributions correspond to squared loss functions
- ▶ Squared loss functions are not robust to outliers!
- ▶ Outliers occur at object boundaries (violation of smoothness/regularizer)
- ▶ Outliers occur at specular highlights (violation of photoconsistency/data term)

# Robust Regularization

Solution: Use **robust data term** and smoothness penalties

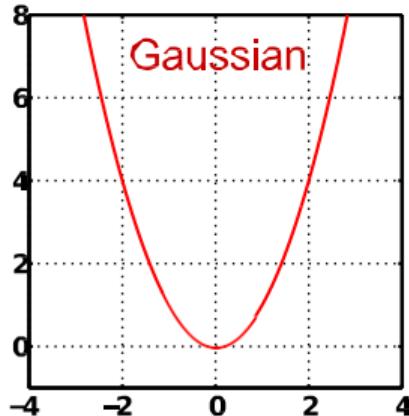
$$\begin{aligned} p(\mathbf{U}, \mathbf{V}) &\propto \prod_{x,y} \exp \left\{ -\rho_D(I_x(x, y) u_{x,y} + I_y(x, y) v_{x,y} + I_t(x, y)) \right\} \\ &\times \exp \left\{ -\lambda \rho_S(u_{x,y} - u_{x+1,y}) \right\} \times \exp \left\{ -\lambda \rho_S(u_{x,y} - u_{x,y+1}) \right\} \\ &\times \exp \left\{ -\lambda \rho_S(v_{x,y} - v_{x+1,y}) \right\} \times \exp \left\{ -\lambda \rho_S(v_{x,y} - v_{x,y+1}) \right\} \end{aligned}$$

- ▶ How to choose  $\rho_D(\cdot)$  and  $\rho_S(\cdot)$ ? We want a prior that allows for discontinuities in the optical flow and a likelihood that allows for outliers and occlusions
- ▶ Replace Gaussian with (heavy-tailed) **Student-t distribution** (=Lorentzian penalty):

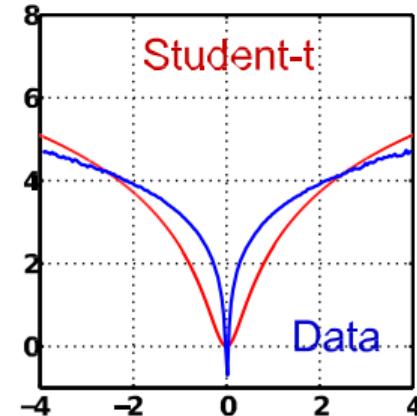
$$p(x) \propto \left( 1 + \frac{x^2}{2\sigma^2} \right)^{-\alpha} \Rightarrow \rho(x) = -\log(p(x)) = \alpha \log \left( 1 + \frac{x^2}{2\sigma^2} \right)$$

# Robust Regularization

Gaussian penalty (squared loss) vs. robust Student-t penalty:



negative  
log-density  
(i.e. energy)



$$p(x) \propto 1 + \frac{x^2}{2\sigma^2}^{-\alpha}$$

$$\rho(x) = \alpha \log 1 + \frac{x^2}{2\sigma^2}$$

- Proposed by [Black-Anandan, ICCV 1993], [Sun et al., CVPR 2010] and others

# Result of Horn & Schunck



# Results of Sun / Black & Anandan



# End-to-End Deep Learning

# Optical flow with deep neural networks

- “Classical” approaches: complex optimization problems, computationally expensive  
→ Not suitable for real-time applications
- First DNN approaches: trade-off between accuracy and size of the model

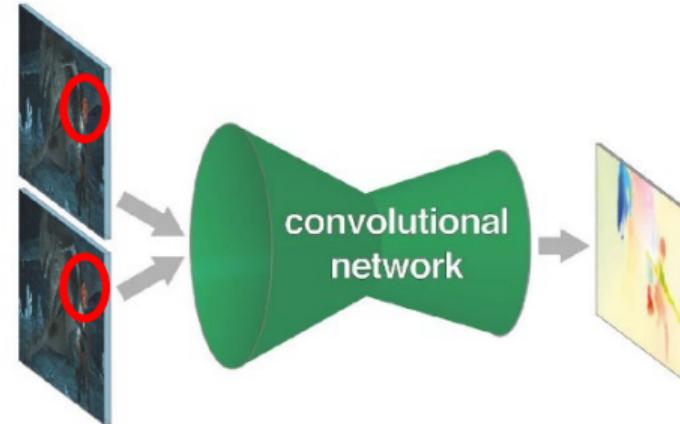
# FlowNet

End to end supervised learning for optical flow estimation:

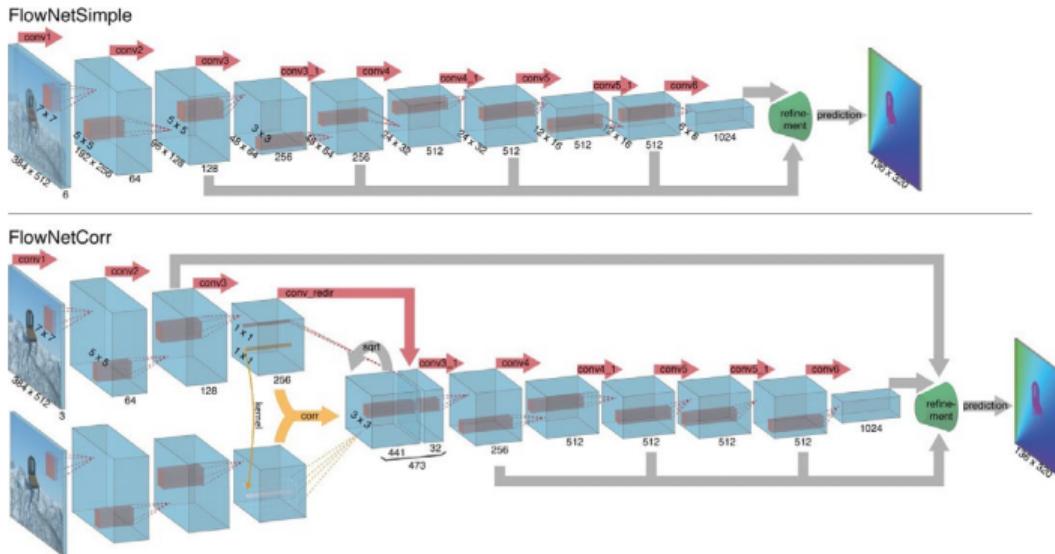
- encoder/decoder architecture

The idea: given two consecutive frames need to match features at different location in input images

The network has to find if an object in the first image is at another location in the second images



# FlowNet

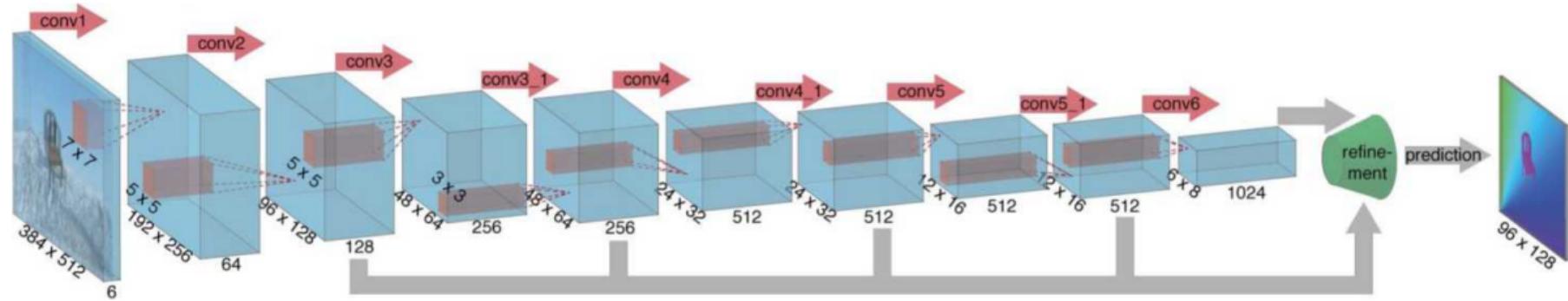


- Encoder: Convolution with stride, decoder: Upconvolution, skip-connections
- Multi-scale loss (EPE in pixels), curriculum learning, synth. training data
- End-Point Error (EPE) is a commonly used metric to evaluate the accuracy of optical flow estimation, measured in pixels. It calculates the Euclidean distance between the predicted flow vectors and the ground truth flow vectors at each pixel location.

# FlowNet (encoder)

Option A: stack both input images together and feed them through a generic network.

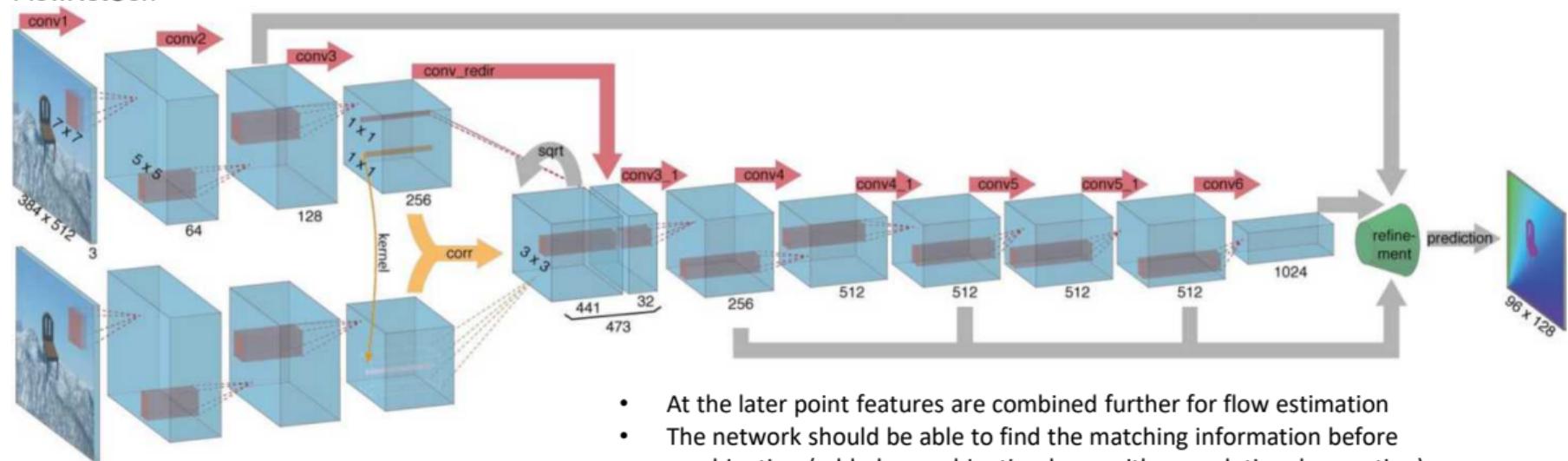
FlowNetSimple



# FlowNet (encoder)

Option B: create two separate, yet identical processing streams for the two images and combine them at a later stage.

FlowNetCorr

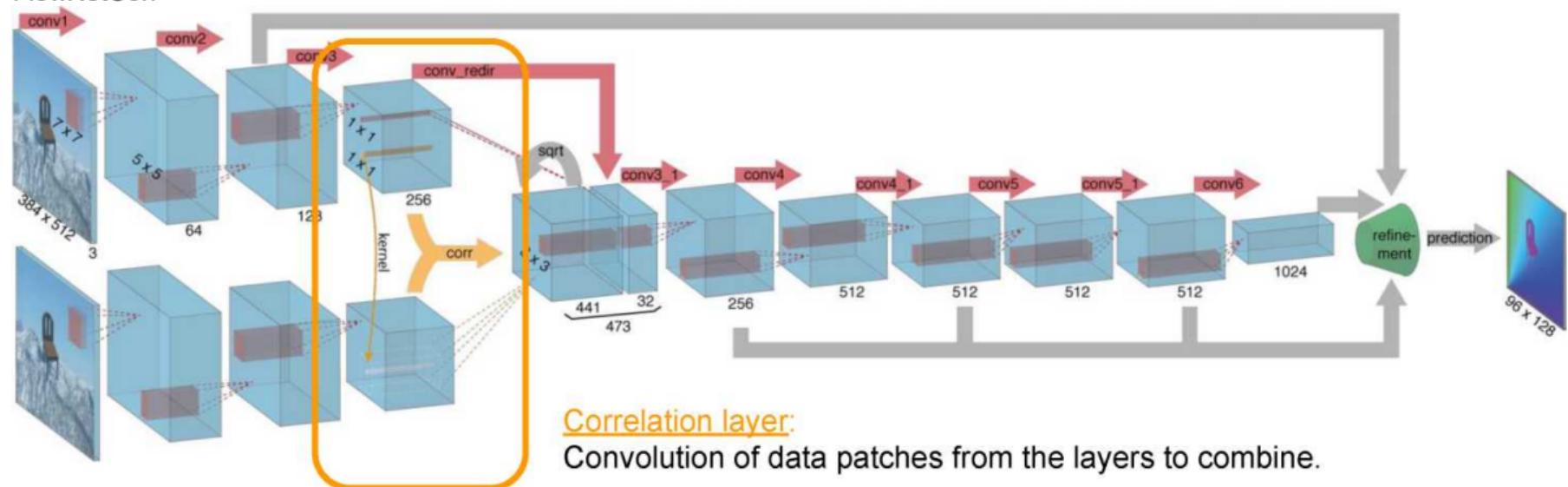


- At the later point features are combined further for flow estimation
- The network should be able to find the matching information before combination (added a combination layer with convolutional operation)
- The features are skip connected to the refinement layer

# FlowNet (encoder)

Option B: create two separate, yet identical processing streams for the two images and combine them at a later stage.

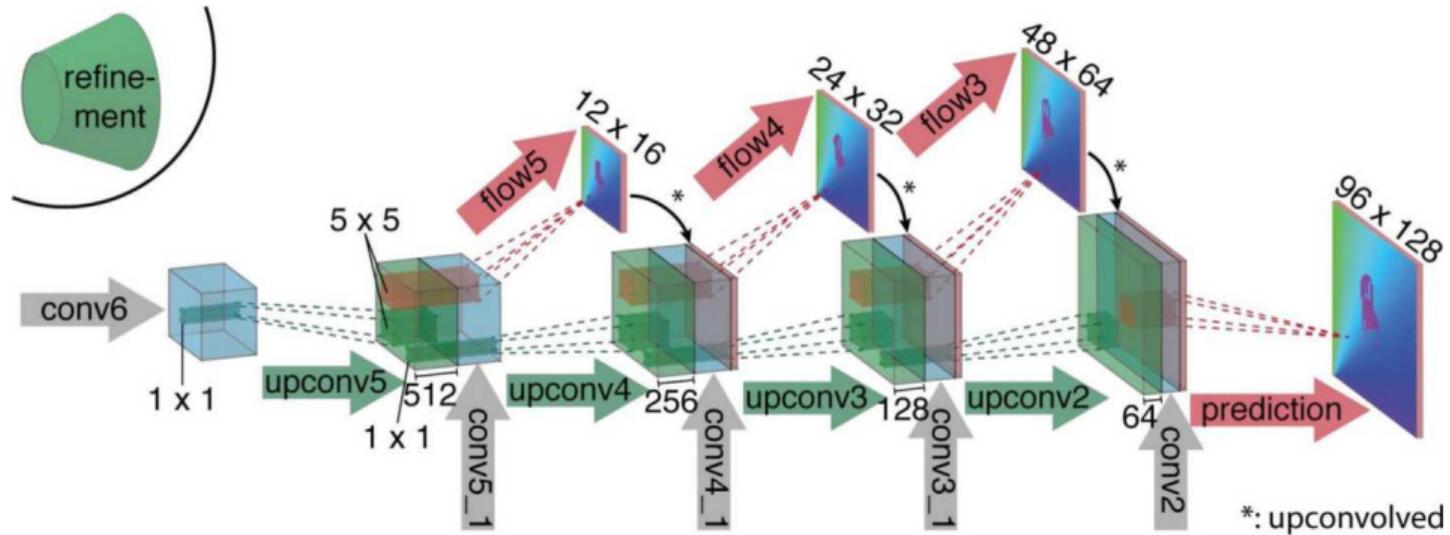
FlowNetCorr



# FlowNet (decoder)

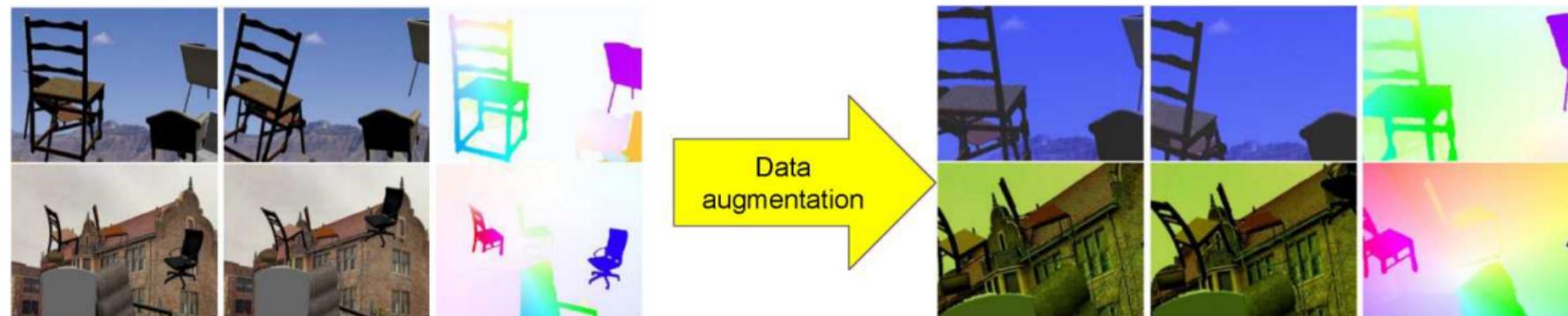
Upconvolutional layers: Unpooling features maps + convolution.

Upconvolutioned feature maps are concatenated with the corresponding map from the contractive part.



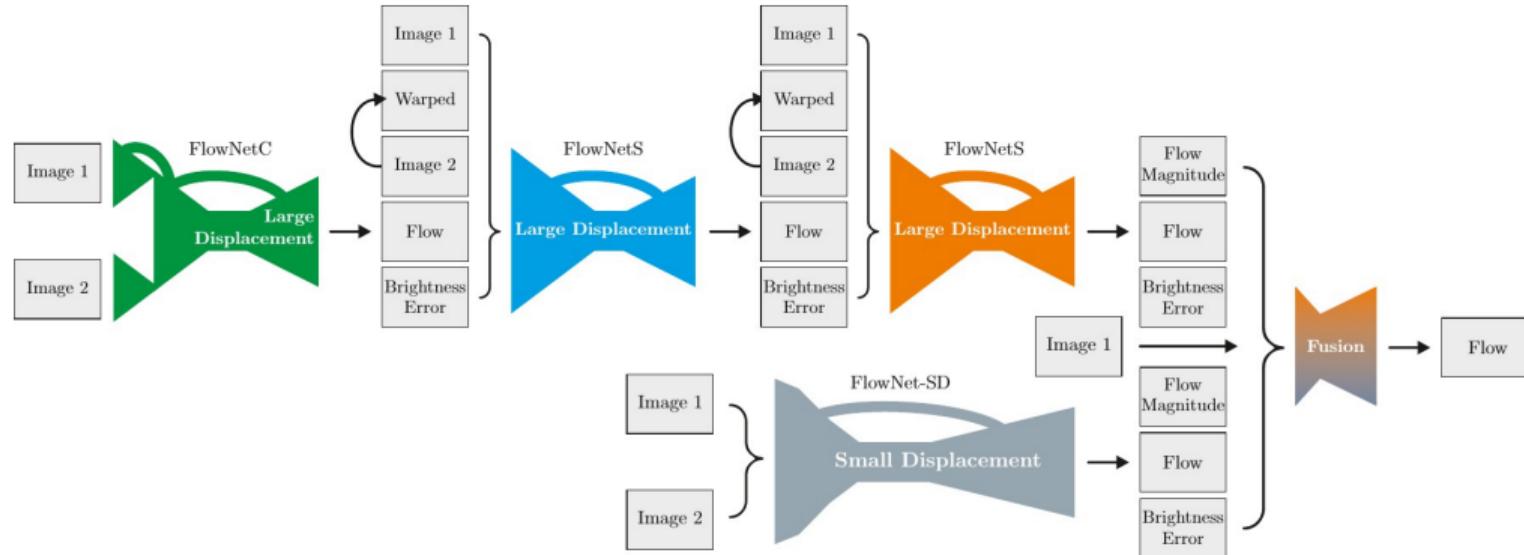
- Features map are scaled to the original resolution
- Upconvolutional layer (unpooling and convolution) doubling the resolution
- Concatenate the 'upconvolution' results with the features from the 'contractive' part of the network
- Bilinear upsampling (variational upsampling) in the last stage

# Optical Flow: FlowNet (synthetic)



Convnets trained on these unrealistic data generalize well to existing datasets such as Sintel and KITTI.

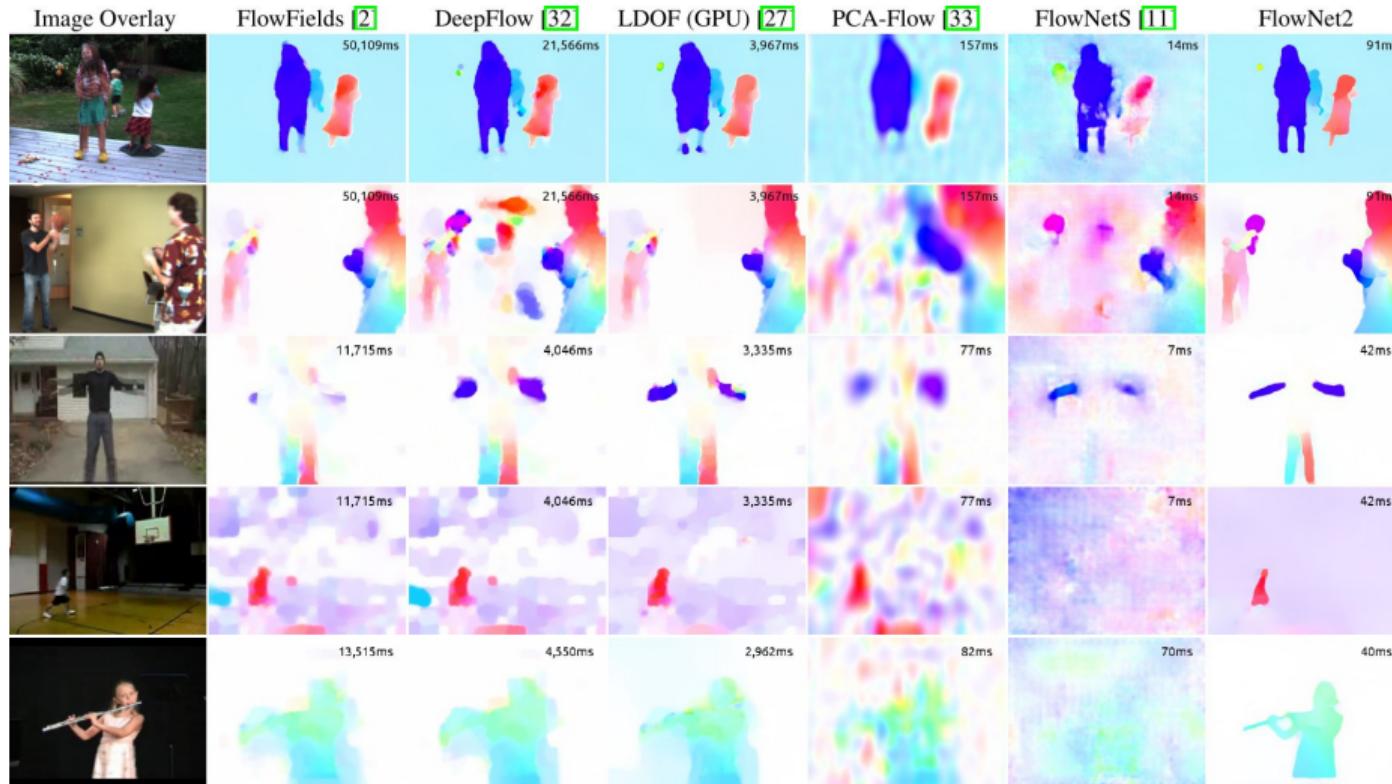
# FlowNet2



[Ilg et al., CVPR 2017]

- Combination (stacking) of multiple FlowNets with warping

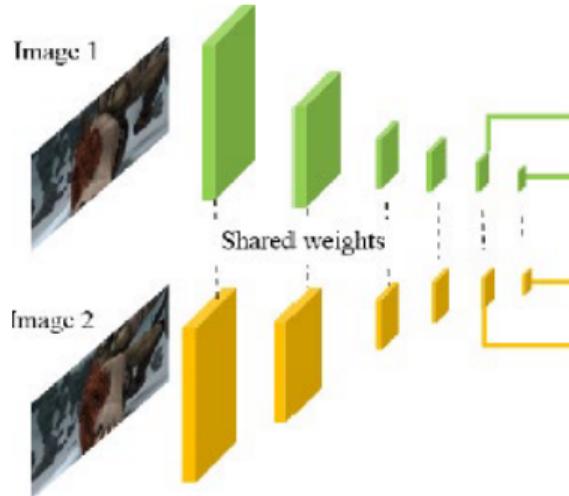
# FlowNet2 – Results



# PWC-Net (Pyramid, Warping, Cost volume)

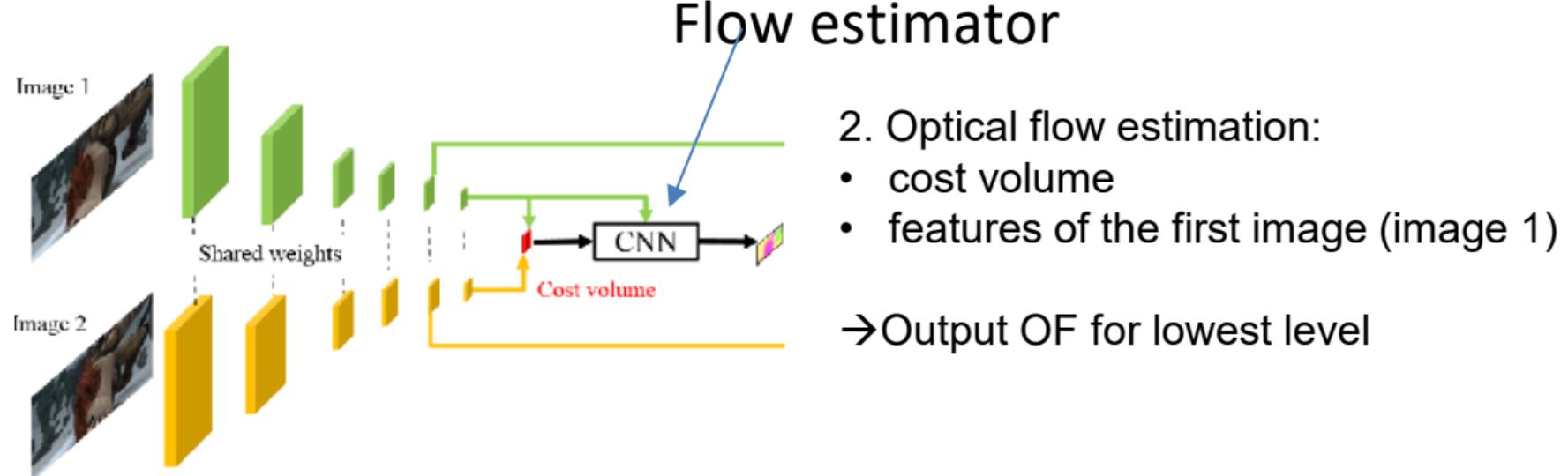
1. Feature extraction from input images with feature pyramid, i.e. convolutional layers
  - Reduction of spatial resolution
2. Optical flow estimation for every level of feature pyramid
  - Start with last convolutional layer, finish on input level
  - Warping and cost volume used in optical flow estimation

# PWC-Net

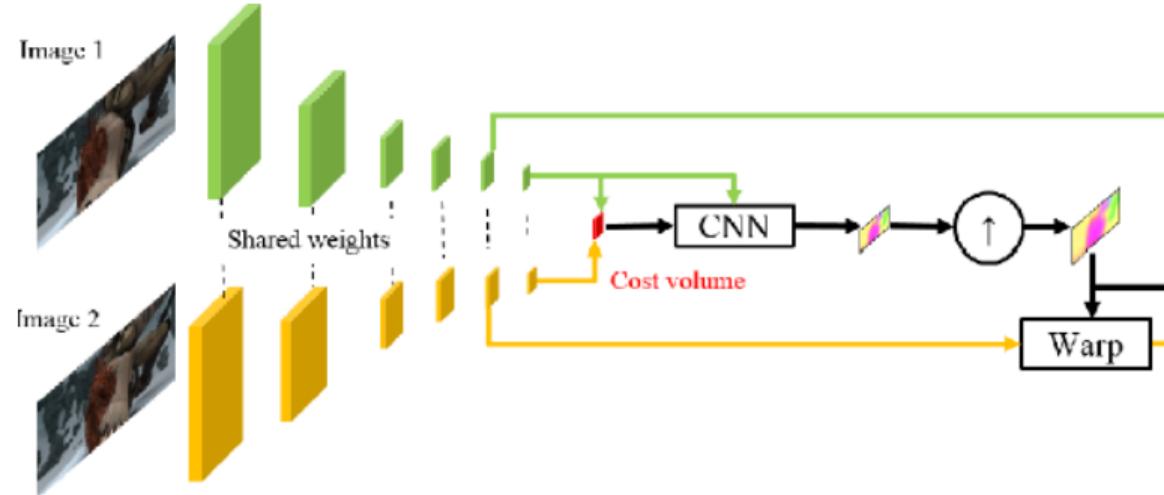


1. Compute cost volume: find most similar pixel in features w.r.t other image
  - We can use normalized cross correlation to get the cost volume of a pixel
  - Invariance to color change, not in scale so pyramidal processing
  - Learning the features pyramid at lower scale we can learn global context to do the matching

# PWC-Net



# PWC-Net

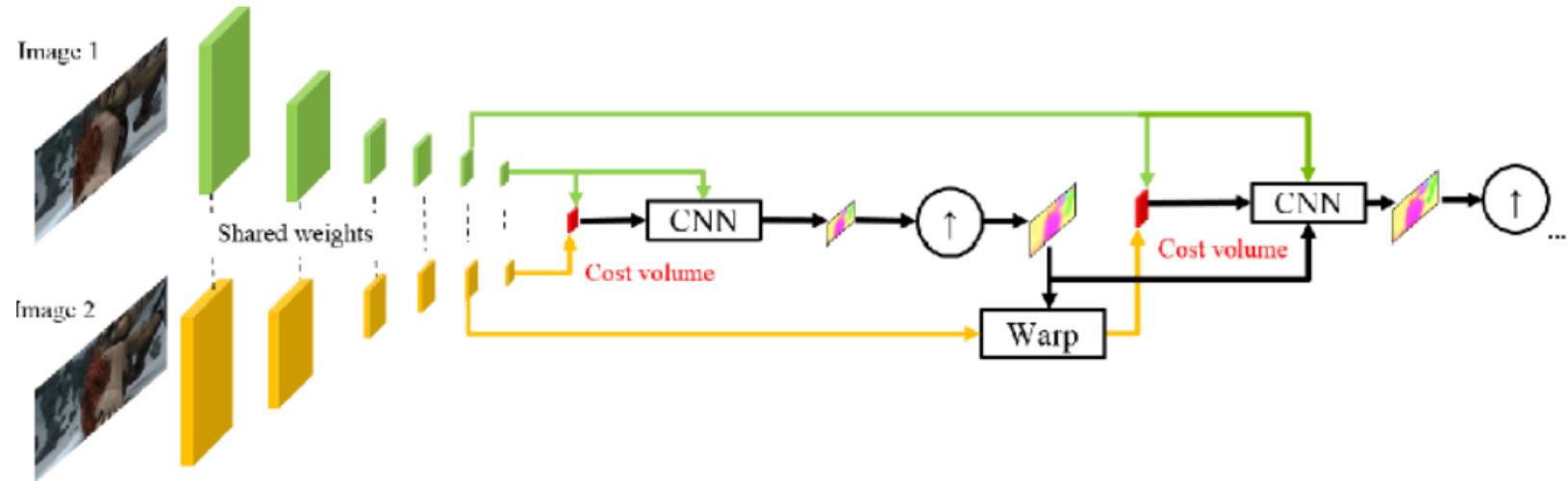


Why warping?

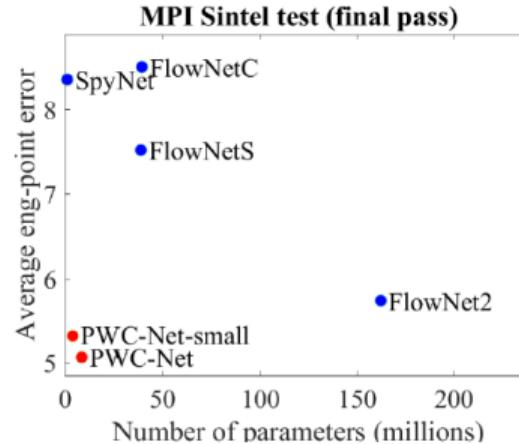
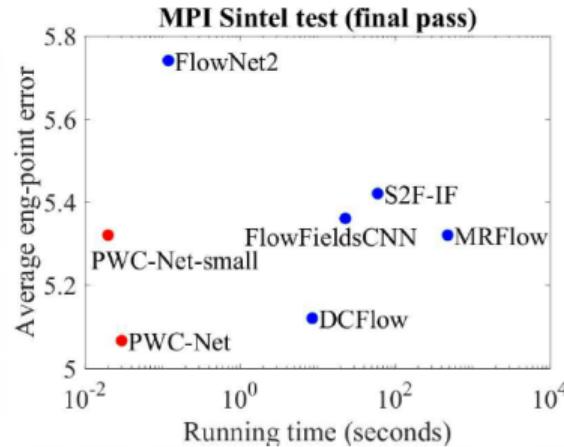
- Second image becomes more similar to first image
- Pixel displacement becomes smaller

1. Upsample OF to match spatial dimensions
2. Warp the features of the 2nd image towards the features of the 1st image (using the estimated flow field)

# PWC-Net



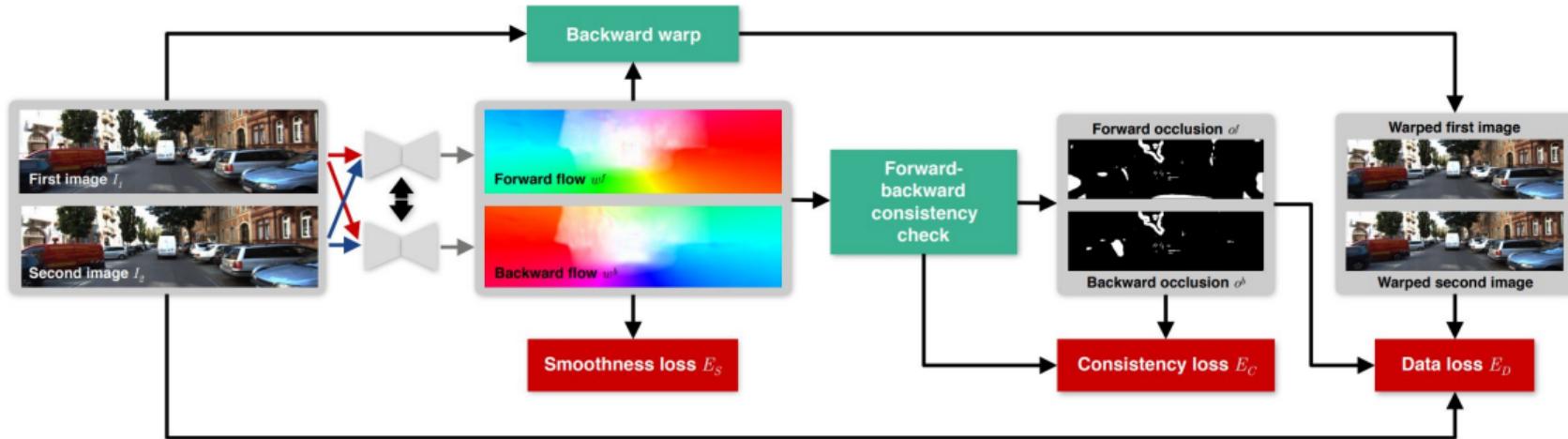
# PWC-Net results



- State-of-the-art accuracy with end-to-end training
- Inference fast enough for real-time application
- PWC-Net-small for mobile applications

Code publicly available: <https://github.com/NVlabs/PWC-Net>

# UnFlow: Unsupervised Optical Flow



- ▶ Learn optical flow without supervision by warping images into the other frame
- ▶ Uses photometric/smoothness terms as loss functions for training the neural net

# Optical Flow Summary

- ▶ Classical OF approaches state-of-the-art until 2016
- ▶ DL based methods on par or better since 2017
- ▶ But require
  - ▶ Big models
  - ▶ Enormous amount of (synthetic) training data
  - ▶ GPU compute time
  - ▶ Sophisticated curriculum learning schedules arranging the data based on measure of difficulty.
- ▶ Top performing DL methods borrow many elements from classical methods (e.g., warping, cost volume, coarse-to-fine estimation, loss functions)

Acknowledgement: some slides and material from Bernt Schiele, Mario Fritz, Michael Black, Bill Freeman, Fei-Fei, Justin Johnson, Serena Yeung, R. Szeliski, Ioannis Gkioulekas, Noah Snavely, Abe Davis, Kris Kitani, Xavier Giró-i-Nieto, Shree Nayar, Andreas Geiger