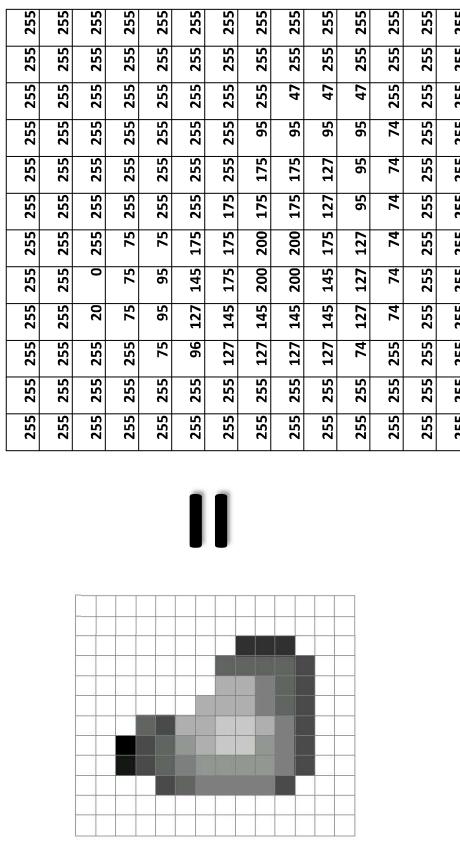


Topics

What is an image?

- A grid (matrix) of intensity values
 - Image
 - Linear filter
 - Cross-correlation
 - Convolution
 - Mean filter
 - Gaussian Filter
 - Image Sharpening
- 
- (common to use one byte per value: 0 = black, 255 = white)

Filters

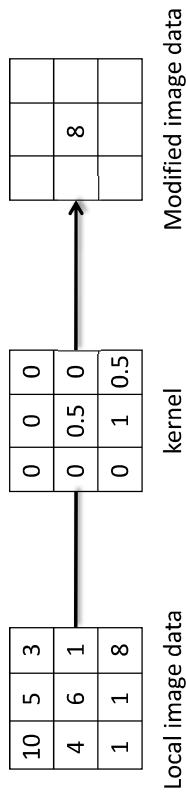
Image filtering

- Modify the pixels in an image based on some function of a local neighborhood of each pixel



Linear filtering

- One simple version of filtering: linear filtering (cross-correlation, convolution)
 - Replace each pixel by a linear combination (a weighted sum) of its neighbors
 - The prescription for the linear combination is called the “kernel” (or “mask”, “filter”)



Source: Zhang et al.

Convolution

- Same as cross-correlation, except that the kernel is “flipped” (horizontally and vertically)

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i - u, j - v]$$

This is called a convolution operation:

$$G = H * F$$

Cross-correlation

- Let F be the image, H be the kernel (of size $2k+1 \times 2k+1$), and G be the output image

$$G[i, j] = \sum_{u=-k}^{\infty} \sum_{v=-k}^{\infty} H[u, v] F[i + u, j + v]$$

This is called a **cross-correlation** operation:

$$G = H \otimes F$$

- Can think of as a “dot product” between local neighborhood and kernel for each pixel

Source: Zhang et al.

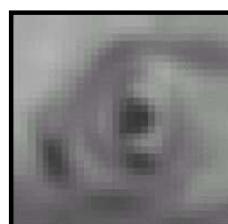
Mean filtering

A 3x3 grid of empty boxes, intended for drawing or writing, with a bounding box of approximately [464, 111, 534, 188].

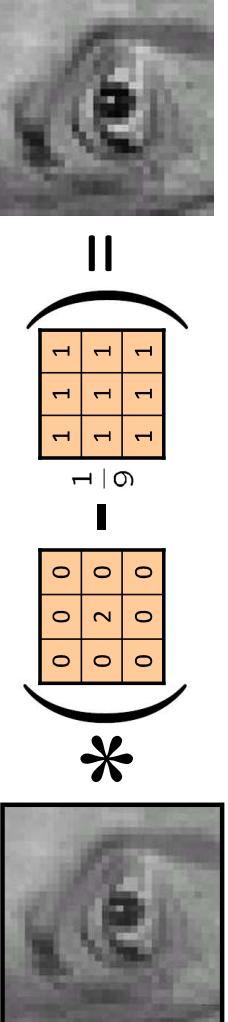
H

E G

Linear filters: Mean filter

$$\text{Original} \quad \quad \quad \text{Blur (with a mean filter)}$$
$$\begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} = \frac{1}{9}$$


Linear filters: examples

$$\text{Original} \quad \quad \quad \text{Sharpening filter}$$
$$\begin{matrix} 0 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 0 \end{matrix} - \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} =$$


Source: D. Lowe

Sharpening

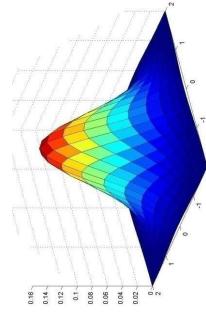
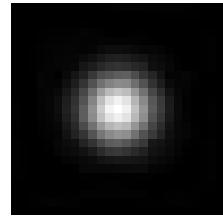
Gaussian Kernel

$$G_\sigma = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$$



before after

Source: D. Lowe



0.003	0.013	0.022	0.013	0.003
0.013	0.059	0.097	0.059	0.013
0.022	0.097	0.159	0.097	0.022
0.013	0.059	0.097	0.059	0.013
0.003	0.013	0.022	0.013	0.003

$5 \times 5, \sigma = 1$

- Constant factor at front makes volume sum to 1 (can be ignored, as we should re-normalize weights to sum to 1 in any case)

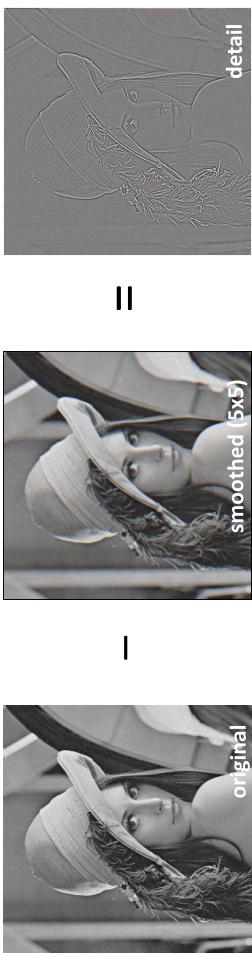
Source: D. Lowe

Source: C. Rasmussen

Gaussian Kernel

Sharpening revisited

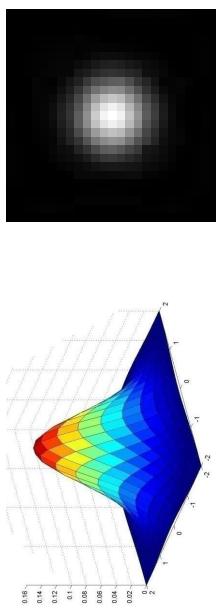
- What does blurring take away?



Let's add it back:



Source: C. Rasmussen



$$G_\sigma = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$$

Topics

- How can we differentiate a *digital* image $F[x, y]$?
 - Option 1: reconstruct a continuous image, f , then compute the derivative
 - Option 2: take discrete derivative (finite difference)
- Edge Detection
 - Image derivatives
 - Image gradient
 - Sobel



Source: S. Lazebnik

Image derivatives

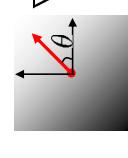
$$\frac{\partial f}{\partial x}[x, y] \approx F[x + 1, y] - F[x, y]$$

Image gradient

Edge detection: smooth first

- The gradient of an image: $\nabla f = \left[\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$

The gradient points in the direction of most rapid increase in intensity

$$\nabla f = \left[\frac{\partial f}{\partial x}, 0 \right]$$


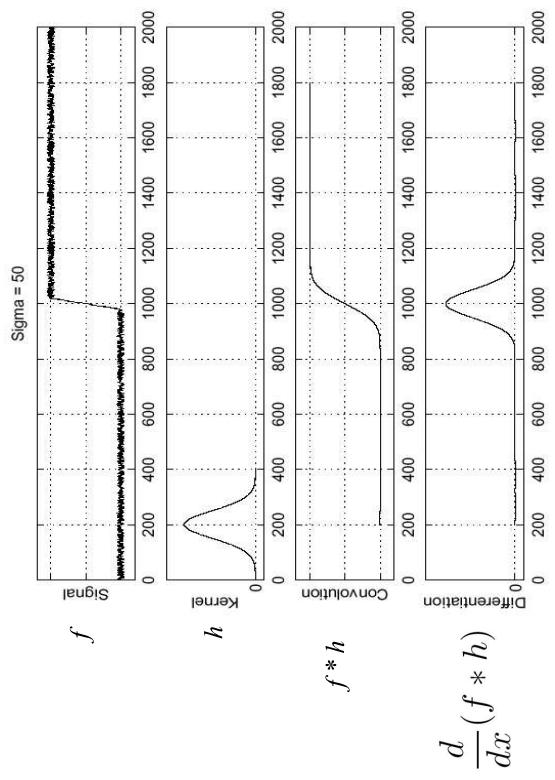
$$\nabla f = \left[0, \frac{\partial f}{\partial y} \right]$$


The edge strength is given by the gradient magnitude:

$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

The gradient direction is given by:

$$\theta = \tan^{-1} \left(\frac{\partial f}{\partial y} / \frac{\partial f}{\partial x} \right)$$



Source: Steve Seitz

Source: S. Seitz

Derivative of Gaussian filter

The Sobel operator

- Common approximation of derivative of Gaussian

$$\begin{array}{c} \frac{1}{8} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} \\ \mathcal{S}_x \quad \mathcal{S}_y \end{array}$$

To find edges, look for peaks in $\frac{d}{dx}(f * h)$

- The standard defn. of the Sobel operator omits the 1/8 term
 - doesn't make a difference for edge detection
 - the 1/8 term is needed to get the right gradient magnitude

Topics

- Image Sampling
 - Gaussian pyramid
 - Image upsampling



Upsampling

- This image is too small for this screen:

- How can we make it 10 times as big?
- Simplest approach:
repeat each row
and column 10 times
- ("Nearest neighbor
interpolation")

- Color and Texture
 - Histogram
 - Edge-based Texture Measures
 - Local Binary Pattern Measure
 - Co-occurrence Matrix Features

Topics

Histograms

Two Edge-based Texture Measures

- A histogram of a gray-tone image is an array $H[*]$ of bins, one for each gray tone.
 - $H[i]$ gives the count of how many pixels of an image have gray tone i .
 - $P[i]$ (the normalized histogram) gives the percentage of pixels that have gray tone i .

1. edgeness per unit area

$$\text{Fedgeness} = \frac{|\{p \mid \text{gradient_magnitude}(p) \geq \text{threshold}\}|}{N}$$

where N is the size of the unit area

2. edge magnitude and direction histograms

$$F_{magdir} = (H_{magnitude}, H_{direction})$$

- where these are the normalized histograms of gradient magnitudes and gradient directions, respectively.

26

Local Binary Pattern Measure

Local Binary Pattern (LBP)

$$LBP_{p,r}(N_c) = \sum_{p=0}^{P-1} g(N_p - N_c) 2^p$$

- N_c : center pixel
- N_p : neighbor pixel
- r : radius (for 3x3 cell, it is 1).

binary threshold function $g(x)$ is,

$$g(x) = \begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$$

- For each pixel p , create an 8-bit number $b_1 \ b_2 \ b_3 \ b_4 \ b_5 \ b_6 \ b_7 \ b_8$, where $b_i = 0$ if neighbor i has value less than or equal to p 's value and 1 otherwise.

- Represent the texture in the image (or a region) by the histogram of these numbers.

	1	2	3	
8	100	101	103	
	40	50	80	4
	50	60	90	
7				5

1	1	1	1	1	1	0	0
---	---	---	---	---	---	---	---

27

Co-occurrence Matrix Features

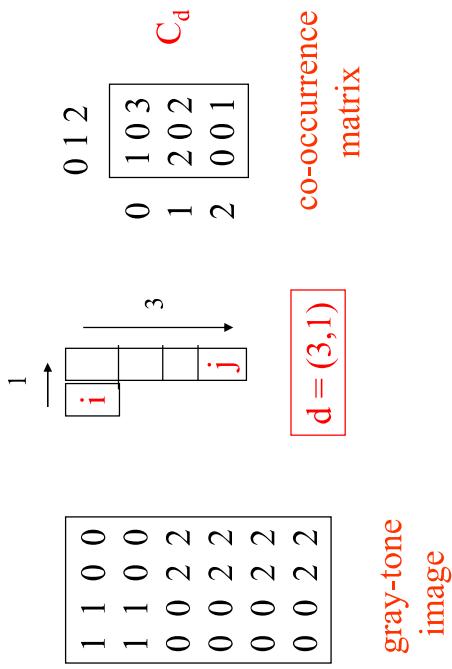
A co-occurrence matrix is a 2D array C in which

- Both the rows and columns represent a set of possible image values.

- $C_d(i,j)$ indicates how many times value i co-occurs with value j in a particular spatial relationship d .

- The spatial relationship is specified by a vector $d = (dr, dc)$.

29



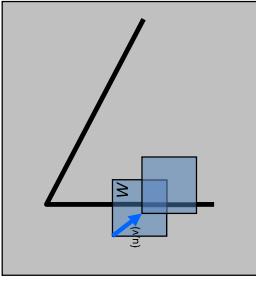
From C_d we can compute N_d , the normalized co-occurrence matrix, where each value is divided by the sum of all the values

30

Topics

- Corners and blobs
 - Harris corner detection
 - Blob detection

Corner detection: the math



From C_d we can compute N_d , the normalized co-occurrence matrix,

- define an SSD "error" $E(u, v)$:

$$\begin{aligned} E(u, v) &\approx \sum_{(x,y) \in W} [I_x u + I_y v]^2 \\ &\approx A u^2 + 2 B u v + C v^2 \end{aligned}$$
$$A = \sum_{(x,y) \in W} I_x^2 \quad B = \sum_{(x,y) \in W} I_x I_y \quad C = \sum_{(x,y) \in W} I_y^2$$

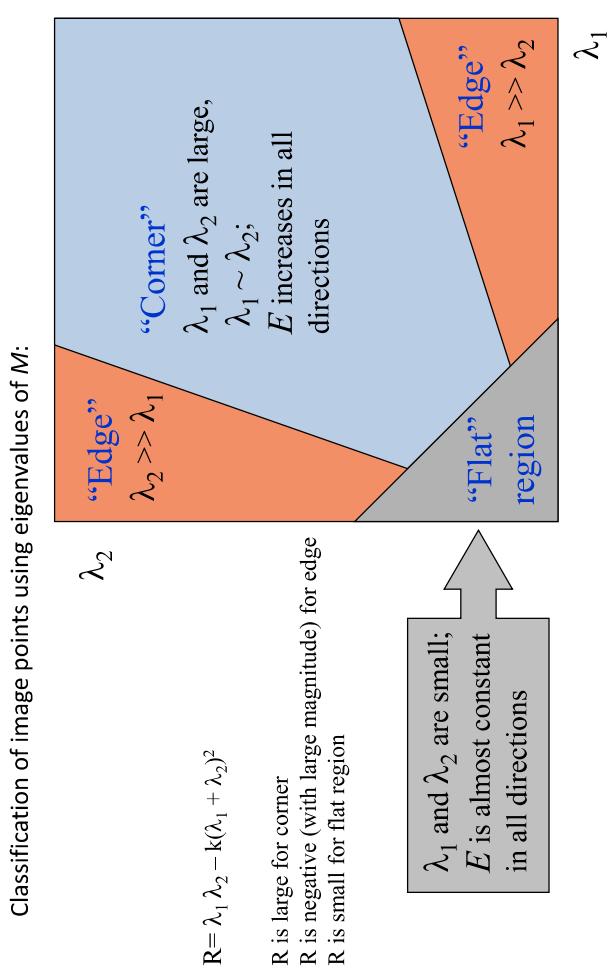
- Thus, $E(u, v)$ is locally approximated as a quadratic error function

Corner detection: the math

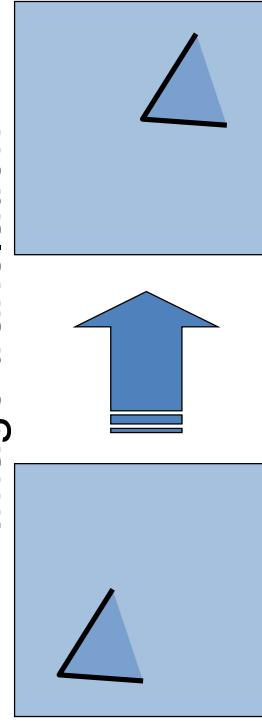
The surface $E(u, v)$ is locally approximated by a quadratic form.

$$\begin{aligned}
 E(u, v) &\approx Au^2 + 2Buv + Cv^2 \\
 &\approx \begin{bmatrix} u & v \end{bmatrix} \begin{bmatrix} A & B \\ B & C \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} \\
 A &= \sum_{(x,y) \in W} I_x^2 \\
 B &= \sum_{(x,y) \in W} I_x I_y \\
 C &= \sum_{(x,y) \in W} I_y^2
 \end{aligned}$$

Interpreting the eigenvalues



Harris detector: Invariance properties
-- Image translation



Other Versions of The Harris operator

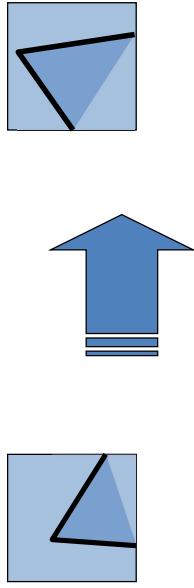
$$\begin{aligned}
 f &= \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2} \\
 &= \frac{\text{determinant}(H)}{\text{trace}(H)}
 \end{aligned}$$

λ_{\min} is a variant of the “Harris operator” for feature detection

- Derivatives and window function are shift-invariant

Harris detector: Invariance properties –

– Image rotation

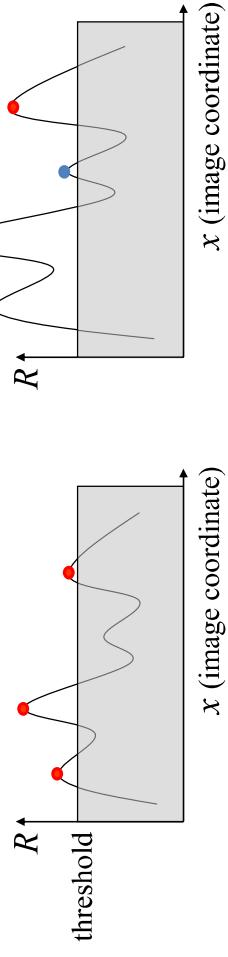


Harris detector: Invariance properties –

Affine intensity change

$$\begin{array}{c} \text{blue square} \\ \rightarrow \\ \text{blue square} \end{array} \quad I \rightarrow aI + b$$

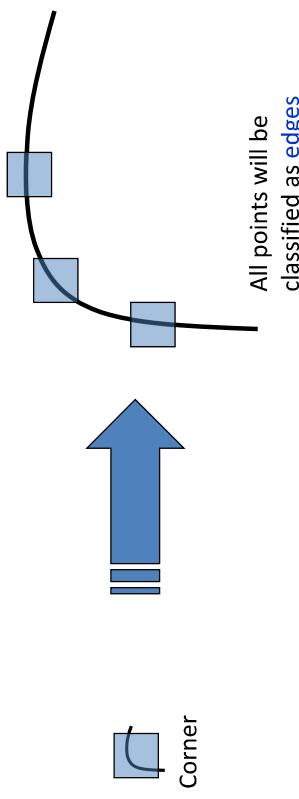
- Only derivatives are used => invariance to intensity shift $I \rightarrow I + b$
- Intensity scaling: $I \rightarrow aI$



Partially invariant to affine intensity change

Harris Detector: Invariance Properties

• Scaling



Not invariant to scaling

Scale Selection

- Instead of computing f for larger and larger windows, we can implement using a fixed window size with a Gaussian pyramid



(sometimes need to create in-between levels, e.g. a $\frac{3}{4}$ -size image)

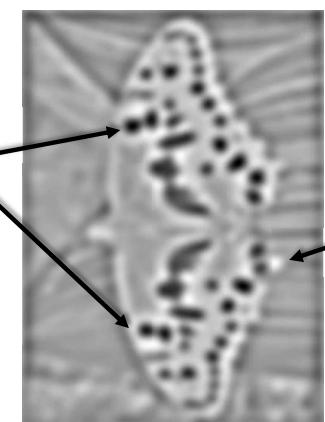
Laplacian of Gaussian

- “Blob” detector

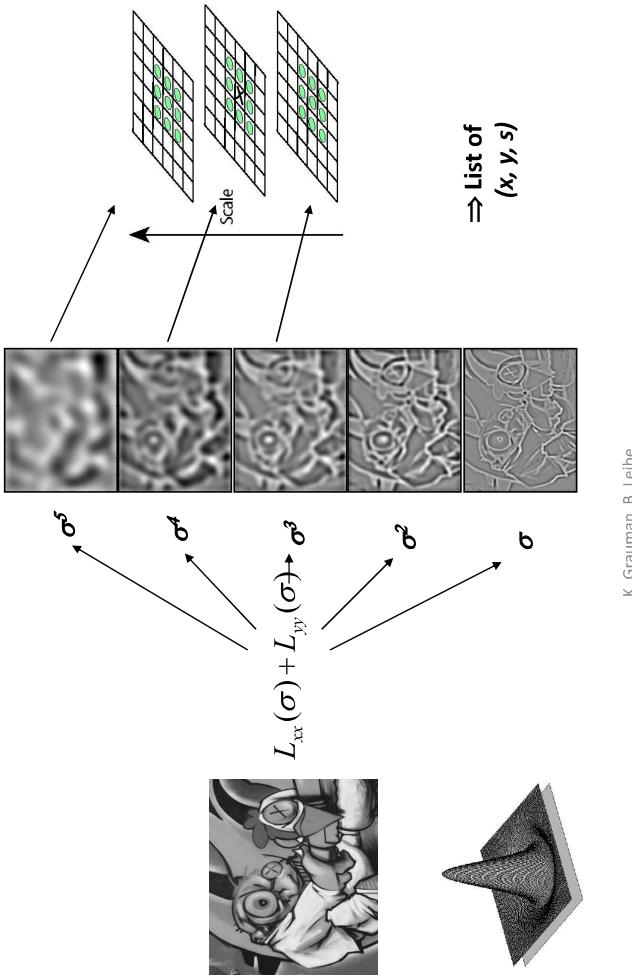


$$* =$$

- Find maxima and minima of Log operator in space and scale



Scale-space blob detector



K. Grauman, B. Leibe

Topics

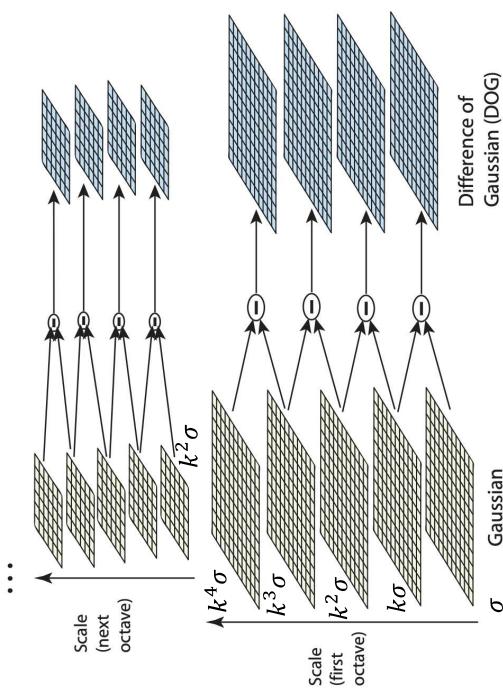
Scale Invariant Feature Transform

- Scale Invariant Feature Transform
 - Building scale-space
 - Interest point detection
 - Orientation assignment
 - SIFT feature descriptor
 - SIFT distance calculation

Scale Invariant Feature Transform

- Scale space peak selection
 - Potential locations for finding features
- Key point localization
 - Accurately locating the feature key points
- Orientation assignment
 - Assigning orientation to the key points
- Key point descriptor
 - Describing the key point as a high dimensional vector (128) (SIFT Descriptor)

Building the Scale Space



Adapted from IICV 2004 by David Lowe

Peak Detection



Adapted from slide by Mubarak Shah

Assignment of the Orientation

Full version

- Divide the 16×16 window into a 4×4 grid of cells (2x2 case shown below)
 - Compute an orientation histogram (8 bin) for each cell (relative orientation and magnitude)
 - $16 \text{ cells} * 8 \text{ orientations} = 128 \text{ dimensional descriptor}$
- The orientation histogram has **36** bins covering the 360 degree range of orientations.
- The samples added to the histogram is weighted by the **gradient magnitude**.
- The **dominate direction** is the peak in the histogram.

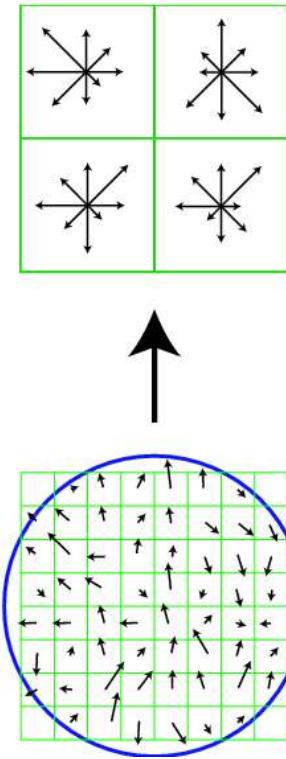


Image gradients

Keypoint descriptor

Adapted from IICV 2004 by David Lowe

Adapted from slide by David Lowe

Feature distance

Feature distance

How to define the difference between two features f_1, f_2 ?

- Simple approach: L_2 distance, $\|f_1 - f_2\|$ (aka SSD)
- can give good scores to ambiguous (incorrect) matches



How to define the difference between two features f_1, f_2 ?

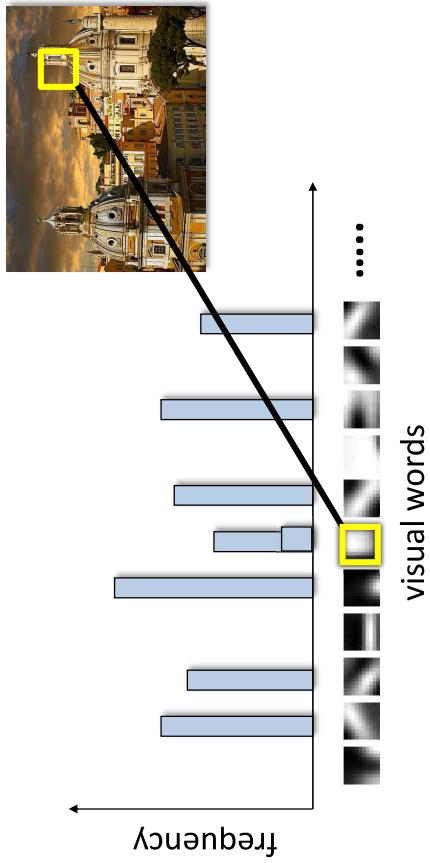
- Better approach: ratio distance = $\|f_1 - f_2\| / \|f_1 - f'_2\|$
- f'_2 is best SSD match to f_1 in l_2
- f'_2 is 2nd best SSD match to f_1 in l_2
- gives large values for ambiguous matches



Topics

Images as histograms of visual words

- Inspired by ideas from text retrieval
 - [Sivic and Zisserman, ICCV 2003]



TF (term frequency)- IDF(inverse document frequency) weighting

- Instead of computing a regular histogram distance, we'll weight each word by it's *inverse document frequency*

- inverse document frequency (IDF) of word j =

$$\log \frac{\text{number of documents}}{\text{number of documents in which } j \text{ appears}}$$

TF-IDF weighting

- To compute the value of bin j in image l :

term frequency of j in l \times inverse document frequency of j

- To compute the value of bin j in image l :

Inverted file

- Each image has $\sim 1,000$ features
- We have $\sim 1,000,000$ visual words
→ each histogram is extremely sparse (mostly zeros)
- Inverted file
 - mapping from words to documents

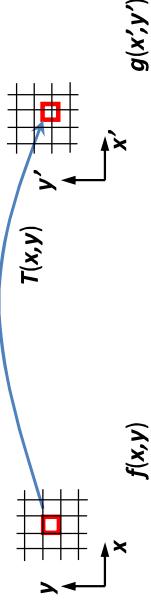
```
"a": [2]
"banana": [2]
"is": [0, 1, 2]
"it": [0, 1, 2]
"what": [0, 1]
```

Topics

- Transformation and Alignment
 - Image warping
 - All 2D Linear Transformations
 - Homogeneous coordinates
 - Affine Transformations
 - RANSAC

Forward Warping

- Send each pixel $f(x, y)$ to its corresponding location $(x', y') = T(x, y)$ in $g(x', y')$
- What if pixel lands “between” two pixels?
- Answer: add “contribution” to several pixels, normalize later
- Can still result in holes



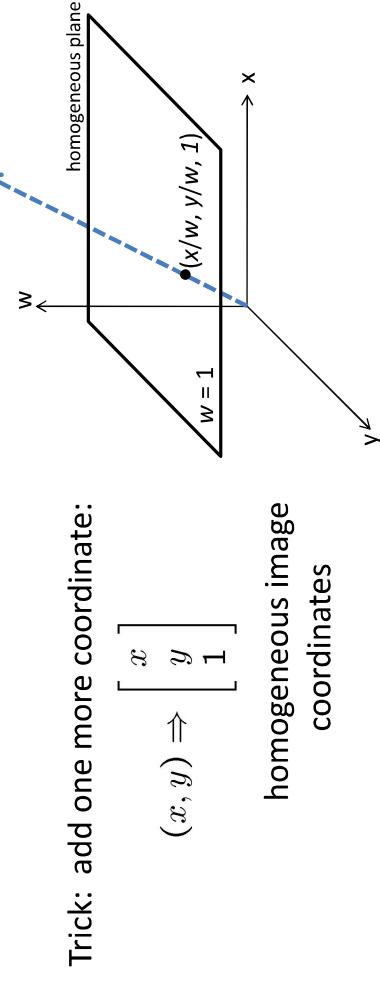
All 2D Linear Transformations

- Linear transformations are combinations of ...

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix} \begin{bmatrix} i & j \\ k & l \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Homogeneous coordinates

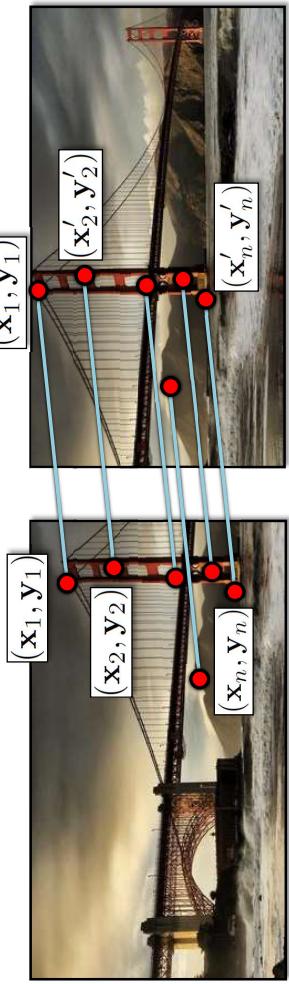


Converting *from* homogeneous coordinates

$$\begin{bmatrix} x \\ y \\ w \end{bmatrix} \Rightarrow (x/w, y/w)$$

Affine Transformations

Simple case: translations



- Affine transformations are combinations of ...

– Linear transformations, and $\begin{bmatrix} x' \\ y' \\ w \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix}$

– Translations $\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

2D *in-plane* rotation

Translate

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & sh_x & 0 \\ sh_y & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Shear

$$(\mathbf{x}_t, \mathbf{y}_t) = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}'_i - \mathbf{x}_i, \frac{1}{n} \sum_{i=1}^n \mathbf{y}'_i - \mathbf{y}_i \right)$$

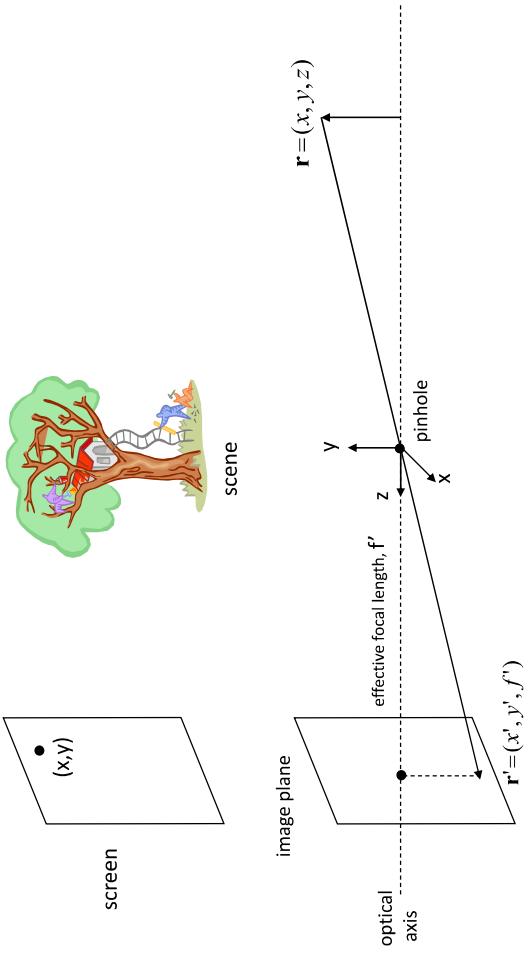
Displacement of match $i = (\mathbf{x}'_i - \mathbf{x}_i, \mathbf{y}'_i - \mathbf{y}_i)$

Topics

- Cameras
 - Pinhole camera
 - Camera parameters
 - Modeling projection
- General version:
 1. Randomly choose s samples
 - Typically $s = \text{minimum sample size that lets you fit a model}$
 2. Fit a model (e.g., line) to those samples
 3. Count the number of inliers that approximately fit the model
 4. Repeat N times
 5. Choose the model that has the largest set of inliers

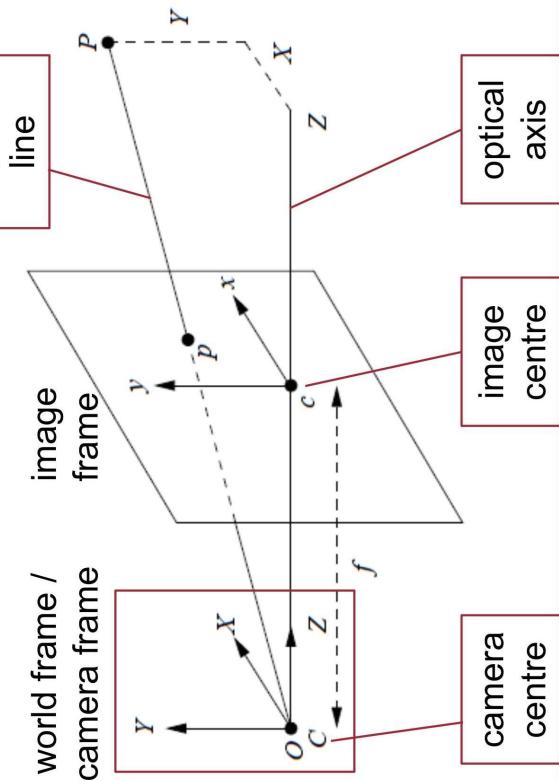
RANSAC

Pinhole and the Perspective Projection



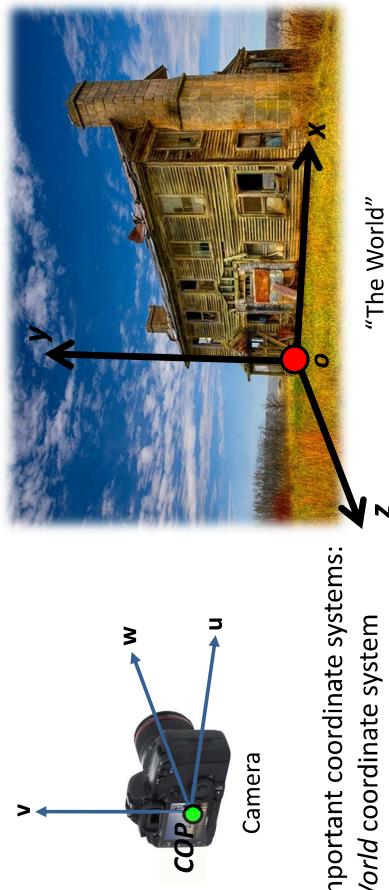
$$\frac{\mathbf{r}'}{f'} = \frac{\mathbf{r}}{z} \quad \Leftrightarrow \quad \frac{x'}{f'} = \frac{x}{z} \quad \frac{y'}{f'} = \frac{y}{z}$$

○ Math representation



Camera parameters

- How can we model the geometry of a camera?



Two important coordinate systems:

1. *World* coordinate system
2. *Camera* coordinate system

Camera parameters

- To project a point (x, y, z) in *world* coordinates into a camera
- First transform (x, y, z) into *camera* coordinates
- Need to know
 - Camera position (in world coordinates)
 - Camera orientation (in world coordinates)
- The formation of image frame
 - Need to know camera *intrinsic*s

Intrinsic Parameters

- In the image frame, denote location of c (*principle point*) image plane as c_x and c_y
 - Image principle point:
 - Intersection between the camera optical axis and image plane
 - Then
- $$P' = (x', y') = (f \frac{x}{z} + c_x, f \frac{y}{z} + c_y)$$

Intrinsic Parameters

- Points in digital image are expressed as in pixels
- Points in image plane are represented in physical measurement (e.g., centimeter)
- The mapping between digital image and image plan can be something like $\frac{\text{pixels}}{\text{cm}}$
- We can use two parameters, k and l , to describe the mapping. If $k = l$, then the camera has “square pixels”.
- The equation now becomes:

$$P' = (x', y') = (f k \frac{x}{z} + c_x, f l \frac{y}{z} + c_y) \\ = (\alpha \frac{x}{z} + c_x, \beta \frac{y}{z} + c_y)$$

Intrinsic Parameters

$$P' = (x', y') = (\alpha \frac{x}{z} + c_x, \beta \frac{y}{z} + c_y)$$

In matrix form:

$$P' = \begin{bmatrix} \alpha & 0 & c_x & 0 \\ 0 & \beta & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = MP$$

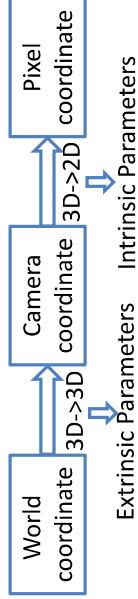
$$P' = MP = \begin{bmatrix} \alpha & 0 & c_x \\ 0 & \beta & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & 1 \end{bmatrix} P = K[I \ 0]P$$

K: Camera matrix (or calibration matrix)

Extrinsic Parameters

- What if the information about the 3D world is available in a different coordinate system?
- We need to relate the points from world reference system to the camera reference system
- Given a point in world reference system P_w , the camera coordinate is computed as

$$P = \begin{bmatrix} R & T \\ 0 & 1 \end{bmatrix} P_w$$



Projection Matrix

- Combining intrinsic and extrinsic parameters, we have

extrinsic parameters

$$P' = K \begin{bmatrix} R & T \end{bmatrix} P_w = M P_w$$

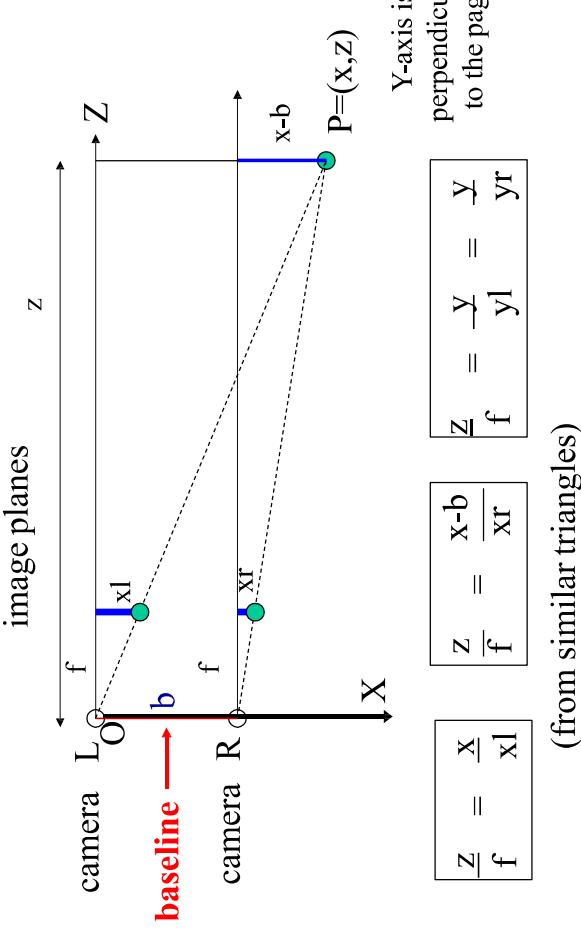
intrinsic parameters

- K changes as the type of camera changes
- Extrinsic parameters are independent of camera

Topics

- Stereo Vision and Structure from Motion
 - Depth and Disparity
 - Epipolar geometry
 - Stereo matching
 - Structure from motion: problem definition

Optic axes of 2 cameras are parallel



(from similar triangles)

$$\frac{z}{f} = \frac{x}{x_l} \quad \frac{z}{f} = \frac{x-b}{x_r}$$

$$\frac{z}{f} = \frac{y}{y_l} = \frac{y}{y_r}$$

Y-axis is perpendicular to the page.

3D from Stereo Images

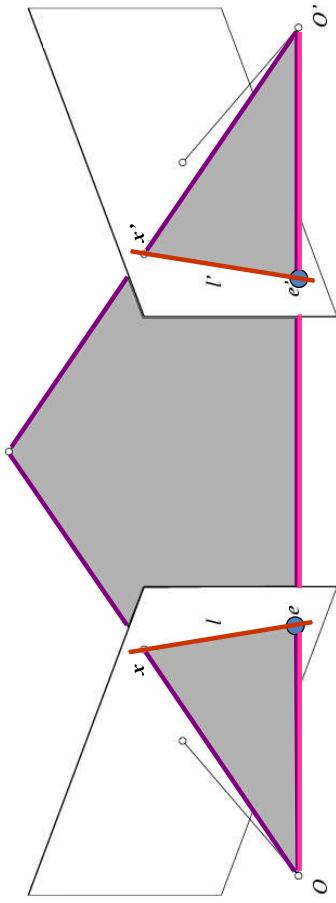
For stereo cameras with parallel optical axes, focal length f , baseline b , corresponding image points (x_l, y_l) and (x_r, y_r) , the location of the 3D point can be derived from previous slide's equations:

$$\text{Depth } z = \frac{f^*b}{d} / (x_l - x_r) = f^*b/d$$

$$x = x_l^*z/f \quad \text{or} \quad b + x_r^*z/f$$

$$y = y_l^*z/f \quad \text{or} \quad y_r^*z/f$$

Note that depth is inversely proportional to disparity



- **Baseline** – line connecting the two camera centers

- **Epipoles**

= intersections of baseline with image planes
= projections of the other camera center

- **Epipolar Plane** – plane containing baseline (1D family)

- **Epipolar Lines** – intersections of epipolar plane with image planes (always come in corresponding pairs)

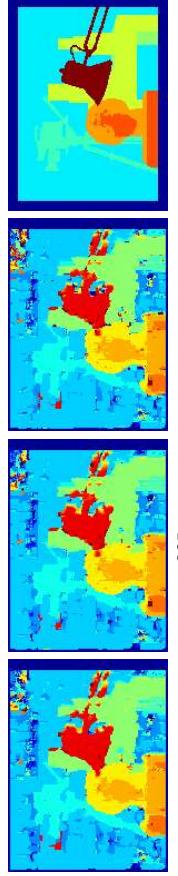
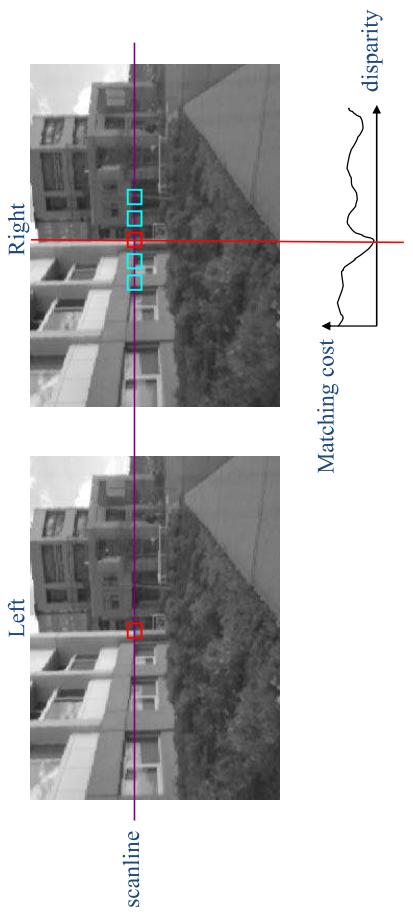
78

Matching windows:

Similarity Measure

Formula	Similarity Measure
$\sum_{(i,j) \in W} I_1(i,j) - I_2(x+i, y+j) $	Sum of Absolute Differences (SAD)
$\sum_{(i,j) \in W} (I_1(i,j) - I_2(x+i, y+j))^2$	Sum of Squared Differences (SSD)
$\sum_{(i,j) \in W} I_1(i,j) - \bar{I}_1(i,j) - I_2(x+i, y+j) + \bar{I}_2(x+i, y+j) $	Zero-mean SAD
$\frac{\sum_{(i,j) \in W} I_1(i,j) \cdot I_2(x+i, y+j)}{\sqrt{\sum_{(i,j) \in W} I_1^2(i,j) \cdot \sum_{(i,j) \in W} I_2^2(x+i, y+j)}}$	Normalized Cross Correlation (NCC)

- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- **Matching cost: SSD, SAD, or normalized cross correlation**



79

Structure from motion

-

Optical Flow

Figure 1 consists of two panels. The top panel shows a square container with four particles: a red dot at the top center, a purple dot at the bottom left, a purple dot at the top right, and a green dot at the bottom center. The bottom panel shows the same square container with the same four particles. Each particle now has a small arrow indicating its initial velocity: the red dot has an arrow pointing up and to the left; the bottom-left purple dot has an arrow pointing down and to the left; the top-right purple dot has an arrow pointing down and to the right; and the green dot has an arrow pointing up and to the right.

- Given two subsequent frames, estimate the point translation

The brightness constancy constraint

Diagram illustrating a coordinate system with origin (x_o, y_o) and a displacement vector (u, v) pointing to $(x_o + u, y_o + v)$. The displacement vector is shown as an arrow originating from the origin.

- Brightness Constancy Equation:

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

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

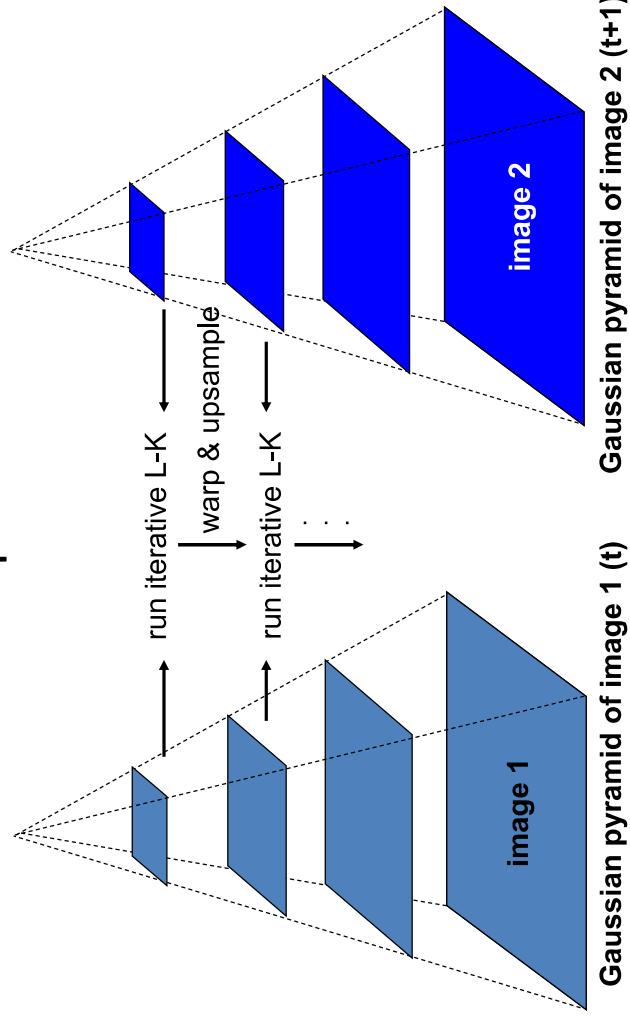
$$I(x+u, y+v, t+1) \approx I(x, y, t) + I_x \cdot u + I_y \cdot v + I_t \cdot \boxed{Difference}$$

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

Iterative Refinement

- **Iterative Lucas-Kanade Algorithm**

1. Estimate velocity at each pixel by solving Lucas-Kanade equations
2. Warp $I(t-1)$ towards $I(t)$ using the estimated flow field
 - use image warping techniques
3. Repeat until convergence



A Few Details

- **Top Level**
 - Apply L-K to get a flow field representing the flow from the first frame to the second frame.
 - Apply this flow field to warp the first frame toward the second frame.
 - Return L-K on the new warped image to get a flow field from it to the second frame.
 - Repeat till convergence.
- **Next Level**
 - Upsample the flow field to the next level as the first guess of the flow at that level.
 - Apply this flow field to warp the first frame toward the second frame.
 - Rerun L-K and warping till convergence as above.
- **Etc.**

Coarse-to-fine optical flow estimation