

## Contents

1	<b>The digital image</b>	9.7
1.1	Image as 2D signal	9.8
1.2	Image sources	9.9
1.3	Sampling	
1.4	Reconstruction	9.10
1.5	Quantization	9.11
1.6	Image Properties	9.12
1.7	Image Noise	9.13
1.8	Colour Images	

2	<b>Image segmentation</b>	10
2.1	Thresholding	10.1
2.2	Region Growing	10.2
2.2.1	Variations	
2.3	Spatial relations	10.3
2.4	Morphological Operations	
2.4.1	Medial Axis Transform (MAT, skeletonization)	

3	<b>Image filtering</b>	3.1
3.1	Linear Shift-Invariant Filtering	

4	<b>Image features</b>	4.1
4.1	Template matching	4.2
4.2	Edge detection	
4.2.1	Canny edge detector	
4.3	Feature detection	4.3.1
4.3.1	Hough transform	
4.3.2	Detecting corner points	
4.3.3	Most accurately localizable patterns	
4.3.4	Lowe's SIFT features	

5	<b>Fourier Transform</b>	5.1
5.1	Aliasing	5.2
5.2	Definition	5.3
5.3	Sampling	5.4
5.4	Image Restoration	

6	<b>Unitary transforms</b>	6.1
6.1	Karhunen-Loeve Transform	6.2
6.2	Basic images and eigenimages (EI)	6.3
6.3	Eigenfaces (EF)	6.4
6.4	Fisherfaces/LDA	6.5
6.5	JPEG image compression	

7	<b>Scale-space representations</b>	7.1
7.1	Image pyramid	7.2
7.2	Applications	7.3
7.3	Pyramids	7.4
7.4	Haar Transform	

8	<b>Optical Flow</b>	8.1
8.1	Applications	8.2
8.2	Brightness constancy	8.3
8.3	Mathematical formulation	8.4
8.4	The aperture problem	8.5
8.5	Optical Flow meaning	8.6
8.6	Regularization: Horn & Schunck algorithm	8.7
8.7	Lucas-Kanade: Integrate over a Patch	8.8
8.8	Gradient-Based Estimation	8.9
8.9	Aperture Problem	8.10
8.10	Iterative Optical Flow Estimation	8.11
8.11	Pyramid/Coarse-to-fine	8.12
8.12	Robust Motion Estimation	8.13
8.13	Motion Models	8.14
8.14	Low-order Parametric Deformations	8.15
8.15	Global Smoothing	8.16
8.16	Probabilistic Formulations	8.17
8.17	Parametric motion models	

9	<b>Video Compression</b>	9.1
9.1	Perception of motion	9.2
9.2	Interlaced video format	9.3
9.3	Why compress video?	9.4
9.4	Lossy video compression	9.5
9.5	Block-matching motion estimation	9.5.1
9.5.1	Determining the best matching block	
9.6	Motion vector and motion vector field	

9.6.1	Practical Half-Pixel Motion Estimation Algorithm	9.7
9.7	Block Matching Algorithm	9.8
9.8	Frame types	9.9
9.9	Summary of Temporal Processing	
9.10	Basic Video Compression Architecture	
9.11	Scalable Video Coding	
9.12	Standards	
9.13	Quality measure	

10	<b>Radon Transform</b>	10.1
10.1	Definition	10.2
10.2	Fourier Slice Theorem	10.3
10.3	Reconstruction	

11	<b>Image denoising</b>	11.1
11.1	Sparse representations for image restoration	11.2
11.2	Sparse linear model	
12	<b>Texture</b>	12.1
12.1	Histogram	12.2
12.2	Texture synthesis	12.2.1
12.2.1	Methods	12.3
12.3	Texture Transfer	

13	<b>Questions</b>	
----	------------------	--

A	<b>Big equations</b>	
---	----------------------	--

## 1 The digital image

Problems of digital cameras sensors

- transmission interference
- compression artefacts
- spilling
- scratches, sensor noise
- bad contrast

### 1.1 Image as 2D signal

**Signal:** function depending on some variable with physical meaning

**Image:** continuous function

**2 variables:**  $xy$ -coordinates

**3 variables:**  $xy$ +time

Brightness is usually the value of the function, but other physical values are: Temperature, pressure, depth...

**What is an image?**

- A picture or pattern of a value varying in space and/or time
- Representation of a function  $f: \mathbb{R}^n \rightarrow S$
- In digital form, e.g.:  
 $I: \{1, \dots, X\} \times \{1, \dots, Y\} \rightarrow S$

- For greyscale CCD images,  $n = 2$ ,  $S = \mathbb{R}^+$

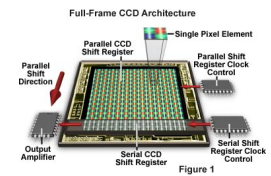
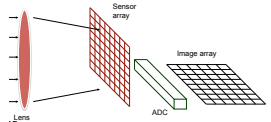
**What is a pixel?**

Not a little square! E.g. gaussian or cubic reconstruction filter.

### 1.2 Image sources

**digital camera (CCD)**

A Charge Coupled Device (CCD).



**analog to digital Conversion**

- The ADC measures the charge and digitizes the result.
- Conversion happens line by line.
- The charges in each photosite move down through the sensor array.

**Blooming** Buckets have finite capacity. Photosite saturation causes blooming.

**bleeding or smearing** during transit buckets still accumulate some charges. Due to tunneling and CCD Architecture. (Influenced by time "in transit" versus integration time. Effect is worse for short shutter times)

**dark current** CCDs produce thermally-generated charge. They give non-zero output even in darkness. Partly, this is the *dark current* and it fluctuates randomly. One can reduce it by cooling the CCD.

**CMOS** Has same sensor elements as CCD. Each photo sensor has its *own amplifier*. This leads to more noise (reduced by subtracting "black" image) and lower sensitivity (lower fill rate). The uses of standard CMOS technology allows to put other components on chip and "smart" pixels.

**CCD vs. CMOS** CCD: mature technology, specific technology, high production cost, high power consumption, higher fill rate, blooming, sequential readout. **CMOS:** recent technology, standard IC technology, cheap, low power, less sensitive, per pixel amplification, random pixel access, smart pixels, on chip integration with other components, rolling shutter (sequential read-out of lines)

### 1.3 Sampling

**ID** Sampling takes a function and returns a vector whose elements are values of that function at the sample points.

**Undersampling** "Missing" things between samples. Information lost

**aliasing** signals "traveling in disguise" as other frequencies. (Can happen in undersampling.)

### 1.4 Reconstruction

Inverse of sampling. Making samples back into continuous function. For output (need realizable method), for analysis or processing (need mathematical method), amounts to "guessing" what the function did in between.

**Bilinear interpolation**

$$f(x, y) = (1-a)(1-b)f[i, j] + a(1-b)f[i+1, j] + abf[i+1, j+1] + (1-a)bf[i, j+1]$$

**Nyquist frequency** Half the sampling frequency of a discrete signal processing system. Signal's max frequency (bandwidth) must be *smaller* than this.

**sampling grids** cartesian sampling, hexagonal sampling and non-uniform sampling

### 1.5 Quantization

real valued function will get digital values - integer values. Quantization is lossy and can't be reconstructed. Simple quantization uses equally spaced levels with  $k$  intervals.

**usual quantization intervals**  
Grayscale image: 8 bit =  $2^8 = 256$  gray-values. **color image RGB** (3 channels): 8 bit/channel =  $2^{24} = 16.7M$  colors. Nonlinear, for example log-scale.

### 1.6 Image Properties

**Image resolution:** Clipped when reduced. **Geometric resolution:** Whole picture but

crappy when reduced. **Radiometric resolution:** Number of colors.

### 1.7 Image Noise

**additive Gaussian Noise** Common model  $I(x, y) = f(x, y) + c$ , where  $c \sim \mathcal{N}(0, \sigma^2)$ . So that  $p(c) = \frac{1}{\sigma\sqrt{2\pi}} e^{-c^2/2\sigma^2}$ .

**Poisson noise** (shot noise)

$$p(k) = \lambda^k e^{-\lambda} / k!$$

**Rician noise:** (appears in MRI)

$$p(I) = \frac{I}{\sigma^2} \exp\left(-\frac{I^2 + f^2}{2\sigma^2}\right) I_0\left(\frac{If}{\sigma^2}\right)$$

**Multiplicative noise:**  $I = f + fc$

**Signal to noise ratio (SNR)**  
 $s = F/\sigma$  is an index of image quality, where

$$F = \frac{1}{XY} \sum_{x=1}^X \sum_{y=1}^Y f(x, y)$$

Often used instead: **Peak Signal to Noise Ratio (PSNR)**  $s_{\text{peak}} = F_{\text{max}}/\sigma$

### 1.8 Colour Images

Consist of red, green and blue channel.

**Prism (with 3 sensors)** Separate light in three beams using dichroic prism. Requires 3 sensors and precise alignment. Gives good color separation.  $\rightarrow$  high-end cameras

**Filter mosaic** Coat filter directly on sensor. "Demosaicing" to obtain full colour & full resolution image.  $\rightarrow$  low-end cameras

**Filter wheel** rotate multiple filters in front of lens. Allows more than 3 colour bands.  $\rightarrow$  static scenes

**new color CMOS sensor, foveon's X3** blue, green, red sensor, one above the other (descending)  $\rightarrow$  better image quality

## 2 Image segmentation

"Segmentation is the ultimate classification problem. Once solved, Computer Vision is solved."

**Interim Summary** Segmentation is hard. It is easier if you define the task carefully: (1) Segmentation task binary or continuous? (2) What are regions of interest? (3) How accurately must the algorithm locate the region boundaries?

**Definition** it partitions an image into *regions of interest*. It is the first stage in many automatic image analysis systems. A *complete segmentation* of an image  $I$  is a finite set of regions  $R_1, \dots, R_N$ , such that

$$I = \bigcup_{i=1}^N R_i \cap R_j = \emptyset, \quad \forall i \neq j$$

**segmentation quality** the quality of a segmentation depends on what you want to do with it. Segmentation algorithms must be chosen and evaluated with an application in mind.

### 2.1 Thresholding

Is a simple segmentation process, produces a binary image  $B$ . It labels each pixel in or out of the region of interest by comparison of the greylevel with a threshold  $T$ :

$$B(x, y) = \begin{cases} 1 & \text{if } I(x, y) \geq T \\ 0 & \text{if } I(x, y) < T. \end{cases}$$

**Choosing  $T$**  By trial and error. Compare results with ground truth. Automatic methods. (ROC curve)

**Chromakeying** Control Lighting!

"Plain" distance measure (e.g.)  
 $I_a = |I - g| > T$

$T \sim 20$ ,  $g = \begin{pmatrix} 0 & 255 & 0 \end{pmatrix}^T$ . Problems: Variation is *not* the same in all 3 channels. Hard alpha mask

$I_{\text{comp}} = I_a I_a + (1 - I_a) I_b$

**Gaussian model per pixel** (Like chromakeying.) mean  $\mu \rightarrow I_{\mu}$ , standard deviation  $\sigma \rightarrow I_{\Sigma}$ .  $I_a = |I - I_b| > T$ ,  $T = \begin{pmatrix} 20 & 20 & 10 \end{pmatrix}$ ,  $I_b =$  background image. Or better (e.g.)

$$I_a = \sqrt{(I - I_b)^T \Sigma^{-1} (I - I_b)}$$

**ROC Analysis** Receiver operating Characteristic. An ROC curve characterizes the performance of a binary classifier. A binary classifier distinguishes between two different types of things.

**Classification error** Binary classifiers make errors. Two types of input to a binary classifier: Positives, negatives. Four possible outcomes in any test: True positive, true negative, false negative, false positive.

**ROC Curve** Characterizes the error trade-off in binary classification tasks. It plots the *true positive fraction*

TP fraction = true positive count/ $P$ ,

$P = TP + FN$  and *false positive fraction*

FP fraction = false positive count/ $N$ ,  $N = FP + TN$ . ROC curve always passes through (0, 0) and (1, 1).

**MAP** (Maximum A Posteriori) detector

**Operating points** choose an *operating point* by assigning relative costs and values to each outcome,  $V_{TN}, V_{TP}, C_{FN}, C_{FP}$ .  $V$  and  $C$  being values and costs. For simplicity, often  $V_{TP}N = V_{TP}P = 0$ .

**Performance Assessment** In real-life, we use two or even three separate sets of test data: (1) A *training set*, for tuning the algorithm, (2) A *validation set* for tuning the performance score, (3) An *unseen test set* to get a final performance score on the tuned algorithm.

**Pixel connectivity** Define neighbors, e.g. (for 2D) 4-neighborhood or 8-neighborhood

**Pixel paths** There are e.g. 4- and 8-connected paths. ( $p_i$  neighbor of  $p_{i+1}$ ).

**Connected regions** A region is 4- or 8-connected if it contains a(n) 4- or 8-connected path between any two of its pixels.

### 2.2 Region Growing

(1) Start from a seed point or region. (2) Add neighboring pixels that satisfy the criteria defining a region. (3) Repeat until we can include no more pixels.

```
function B = RegionGrow(I, seed)
[X,Y] = size(I);
visited = zeros(X,Y);
visited(seed) = 1;
boundary = emptyQ;
boundary.enQ(seed);
while (~boundary.empty())
    nextPoint = boundary.deQ();
    if (include(nextPoint, seed))
        visited(nextPoint) = 2;
        foreach(x,y) in N(nextPoint)
            if (visited(x,y) == 0)
                boundary.enQ(x,y);
            visited(x,y) = 1;
        end
    end
end
```

**end**

### 2.2.1 Variations

**seed selection** (1) One seed point, (2) Seed region, (3) Multiple seeds.

**seed selection** (1) Greylevel thresholding, (2) Greylevel distribution model. E.g. include if  $(I(x, y) - \mu)^2 < (n\sigma)^2$ ,  $n = 3$ . Can update  $\mu$  and  $\sigma$  after every iteration, (3) color or texture information.

**snakes** A snake is an *active contour*. It's a polygon. Each point on contour moves away from seed while its image neighborhood satisfies an inclusion criterion. Often the contour has smoothness constraints, the algorithm iteratively minimizes an energy function:

$$E = E_{\text{tension}} + E_{\text{stiffness}} + E_{\text{image}}$$

### 2.3 Spatial relations

**Markov Random Fields** Markov chains have 1D structure. At every time, there is one state. This enabled use of dynamic programming. **Markov Random fields** break this 1D structure! • Field of sites, each of which has a label, simultaneously. • Label at one site depend on others, no 1D structure to dependencies. • This means no optimal, efficient algorithms, except for 2-label problems. Minimize

$$\text{Energy}(y; \theta, \text{data}) = \sum_i \psi_1(y_i; \theta, \text{data}) + \sum_{i,j} \psi_2(y_i, y_j; \theta, \text{data})$$

•  $i, j$  edges  
**FG-BG segmentation** The code does the following: • background RGB Gaussian model training (from many images) • shadow modeling (hard shadow and soft shadow) • graphcut foreground-background segmentation

### 2.4 Morphological Operations

They are local pixel transformations for processing region shapes. Most often used on binary images. Logical transformations based on comparison of pixel neighborhoods with a pattern.

**8-neighbor erode** (Minkowsky subtraction) Erase any foreground pixel that has one eight-connected neighbor that is background

**8-neighbor dilate** (Minkowsky addition) Paint any background pixel that has one eight-connected neighbor that is foreground. **Applications:** Smooth region boundaries for shape analysis, remove noise and artefacts from an imperfect segmentation, match particular pixel configurations in an image for simple object recognition

**structuring elements** morphological operations take two arguments 1. a binary image 2. a structuring element Compare the structuring element to the neighborhood of each pixel. This determines the output of the morphological operation. The structuring element is also a binary array and has an origin.

$$I_1 \cup I_2 = \{x: x \in I_1 \text{ or } x \in I_2\}, \\ I_2 \cap I_1 = \{x: x \in I_1 \text{ and } x \in I_2\}, \\ I^c = \{x: x \notin I\}, \\ I_1 \setminus I_2 = \{x: x \in I_2 \text{ and } x \notin I_2\}.$$

**Erosion** of binary image  $I$  by the structuring element  $S$  is defined by  
 $I \ominus S = \{z \in E \mid S_z \subset I\}$   
 $S_z$  translation of  $S$  by vector  $z$ .

**Dilation** is  $I \oplus S = \bigcup_{h \in S} I_b$ .

**Opening**  $I \circ S = (I \ominus S) \oplus S$ .  
**Closing**  $I \bullet S = (I \oplus S) \ominus S$ .

To remove holes in the foreground and islands in the background, do both opening and closing. The size and shape of the structuring element determine which features survive. In the absence of knowledge about the shape of features to remove, use a circular structuring element.

**Granulometry** Provides a size distribution of distinct regions or "granules" in the image. We open (opening as above) the image with increasing structuring element size and count the number of regions after each operation. Creates "morphological sieve".

```
function gSpec = granulolo(I, T, maxRad)
% Segment the image I.
B = (I > T);
% Open the image at each structuring element size up
% to a maximum and count the remaining regions.
numRegions(x) = max(max(connectedComponents(O)));
end
gSpec = diff(numRegions);
```

**Hit-and-miss transform**  $H = I \otimes S$  Searches for an exact match of the structuring element. Simple form of template matching.

**Thinning**  $I \odot S = I \setminus (I \otimes S)$

**Thickening**  $I \oslash S = I \cup (I \otimes S)$

**Sequential thinning/thickening** With structuring elements  $S_1, \dots, S_n$  and sequential thinning/thickening •

$I \bullet \{S_i: i = 1, \dots, n\} = ((I \bullet S_1) \bullet \dots \bullet S_n)$   
Several sequences of structuring elements are useful in practice. These are usually the set of rotations of a single structuring element, sometimes called the *Golay alphabet*. See bwmorph in matlab.

**2.4.1 Medial Axis Transform (MAT, skeletonization)**

The skeleton and MAT are stick-figure representations of a region  $X \in \mathbb{R}^2$ . Start a grassfire at the boundary of the region, the skeleton is the set of points at which two fire fronts meet.

**Skeleton** Use structuring element  
 $B = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ .

The  $n$ -th skeleton subset is  
 $S_n(X) = (X \ominus B) \setminus [(X \ominus B) \otimes B]$   
where  $\ominus$  denotes  $n$  successive erosions. The skeleton is the union of all the skeleton subsets  $S(X) = \bigcup_{n=1}^{\infty} S_n(X)$ .

**Reconstruction** can reconstruct region  $X$  from its *skeleton subsets*.  
 $X = \bigcup_{n=0}^{\infty} S_n(X) \oplus B$

**Applications and problems** The skeleton/MAT provides a stick figure representing the region shape. Used in object recognition, in particular, character recognition. **Problems:** Definition of a maximal disc is poorly defined on a digital grid and is sensitive to noise on the boundary. Sequential thinning output sometimes preferred to skeleton/MAT.

### 3 Image filtering

Image filtering is modifying the pixels in an image based on some function of a local neighborhood of the pixels.

**3.1 Linear Shift-Invariant Filtering**

About modifying pixels based on *neighborhood*. Local methods simplest. Linear means *linear combination* of neighbors. Linear methods simplest. *Shift-invariant* means doing the same for each pixel. Same for all is simplest. Useful to: Low-level image processing operations, smoothing and noise reduction, sharpen, detect or enhance features.

**Linear operation**  $L$  is a linear operation if

$$L[\alpha I_1 + \beta I_2] = \alpha L[I_1] + \beta L[I_2]$$

**Output**  $I'$  of linear image operation is a weighted sum of each pixel in the input  $I$

$$I'_j = \sum_{i=1}^N \alpha_{ij} I_i, \quad j = 1 \dots N$$

**Linear Filtering** Linear operations can be written:

$I'(x, y) = \sum_{i,j \in N(x,y)} K(x, y; i, j) I(i, j)$   
 $I$  = input image;  
 $I'$  = output of operation.







$$\Delta = -\frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}.$$

**Remarks**

- 1 Coupled PDE solved using iterative methods and finite differences
- 2  $\hat{x} = \Delta \hat{x} - \lambda \left( \frac{\partial I}{\partial x} \hat{x} + \frac{\partial I}{\partial y} \hat{y} + i \right) \frac{\partial I}{\partial x}$ , where  $g(\mathbf{x})$  is a weighting function that determines the *support* of the estimator (the region within which we combine constraints). It is common to let  $g(\mathbf{x})$  be Gaussian in order to weight constraints in the center of the neighborhood more highly, giving them more influence. The 2D velocity  $\hat{\mathbf{u}}$  that minimizes  $E(\mathbf{u})$  is the least squares flow estimate.
- 3 More than two frames allow a better estimation of  $I$ .
- 4 Information spreads from corner-type patterns.
- 5 Errors at boundaries
- 6 Example of *regularisation*: selection principle for the solution of illposed problems.

## 8.7 Lucas-Kanade: Integrate over a Patch

The Lucas-Kanade method assumes that the displacement of the image contents between two nearby instants (frames) is small and approximately constant within a neighborhood of the point  $p$  under consideration, thus the optical flow equation can be assumed to hold for all pixels within a window centered at  $p$ . Namely, the local image flow (velocity) vector  $(\hat{x}, \hat{y})$  must satisfy

$$\frac{\partial I(q_k)}{\partial x} \hat{x} + \frac{\partial I(q_k)}{\partial y} \hat{y} = -\frac{\partial I(q_k)}{\partial t}$$

for  $k = 1, \dots, n$  and  $q_k$  the pixels inside the window. These equations can be written in matrix form

$$A\mathbf{v} = \mathbf{b} \quad (8.1)$$

where  $\mathbf{x} = (x \ y)^T$ ,  $\mathbf{v} = (\hat{x} \ \hat{y})^T$  and

$$A_{ij} = \frac{\partial I(q_i)}{\partial x_j}, \quad \mathbf{b}_i = -\frac{\partial I(q_i)}{\partial t}$$

Eq. 8.1 is overdetermined, so do compromise solution by the least squares principle Eq. A.3

## 8.8 Gradient-Based Estimation

Assume *brightness constancy*. Let  $f_1(\mathbf{x})$  and  $f_2(\mathbf{x})$  be 1D signals (images) at two time instants. Let  $f_2 = f_1(\mathbf{x} - \delta)$ , where  $\delta$  denotes translation.

$$\begin{aligned} &\sim f_1(\mathbf{x}) - f_2(\mathbf{x}) \\ &= \delta f_1'(\mathbf{x}) + O(\delta^2) \\ &\sim \delta = \frac{f_1(\mathbf{x}) - f_2(\mathbf{x})}{f_1'(\mathbf{x})} \approx \delta \end{aligned} \quad (8.2)$$

Assume displaced image well approximated by first-order Taylor series

$$\begin{aligned} I(\mathbf{x} + \mathbf{u}, t+1) \\ \approx I(\mathbf{x}, t) + \mathbf{u} \cdot \nabla I(\mathbf{x}, t) + I_t(\mathbf{x}, t) \end{aligned} \quad (8.3)$$

Insert Eq. 8.2 in Eq. 8.3 to get

$$\nabla I(\mathbf{x}, t) \cdot \mathbf{u} + I_t(\mathbf{x}, t) = 0. \quad (8.4)$$

This is called the *gradient constraint equation*.

**Intensity Conservation** Tracking points of constant brightness can also be viewed as the estimation of 2D paths  $\mathbf{x}(t)$  along which intensity is conserved:

$$\begin{aligned} I(\mathbf{x}(t), t) &= c, \\ \frac{d}{dt} I(\mathbf{x}(t), t) &= 0 \\ \nabla I \cdot \dot{\mathbf{x}} + I_t &= 0. \end{aligned} \quad (8.5)$$

$\sim$  gradient-constraint equation Eq. 8.4.

**Least-Squares estimation** One cannot recover  $\mathbf{u}$  from one gradient constraint since Eq. 8.4 is one equation with two unknowns,  $u_1$  and  $u_2$ . The intensity gradient constrains the flow to a one parameter family of velocities along a line in *velocity space*. One can see from Eq. 8.4 that this line is perpendicular to  $\nabla I$  and its perpendicular distance from the origin is  $|I_t| / \|\nabla I\|$ .

One common way to further constrain  $\mathbf{u}$  is to use gradient constraints from nearby pixels, assuming they share the same 2D velocity. With many constraints there may be no velocity that simultaneously satisfies them all, so instead we find the ve-

locity that minimizes the constraint errors. The least-squares (LS) estimator minimizes the squared errors:

$$E(\mathbf{u}) = \sum_{\mathbf{x}} g(\mathbf{x}) [\mathbf{u} \cdot \nabla I(\mathbf{x}, t) + I_t(\mathbf{x}, t)]^2, \quad (8.6)$$

where  $g(\mathbf{x})$  is a weighting function that determines the *support* of the estimator (the region within which we combine constraints). It is common to let  $g(\mathbf{x})$  be Gaussian in order to weight constraints in the center of the neighborhood more highly, giving them more influence. The 2D velocity  $\hat{\mathbf{u}}$  that minimizes  $E(\mathbf{u})$  is the least squares flow estimate.

The minimum of  $E(\mathbf{u})$  can be found from its critical points, where its derivatives with respect to  $\mathbf{u}$  are zero; i.e.,

$$\begin{aligned} \frac{\partial E(u_1, u_2)}{\partial u_1} &= \sum_{\mathbf{x}} g(\mathbf{x}) [u_1 I_x^2 + u_2 I_x I_y + I_x I_t] \\ &= 0 \\ \frac{\partial E(u_1, u_2)}{\partial u_2} &= \sum_{\mathbf{x}} g(\mathbf{x}) [u_2 I_y^2 + u_1 I_x I_y + I_y I_t] \\ &= 0 \end{aligned}$$

These equations may be rewritten in matrix form:

$$M\mathbf{u} = \mathbf{b} \quad (8.7)$$

$$M = \begin{pmatrix} \sum g I_x^2 & \sum g I_x I_y \\ \sum g I_x I_y & \sum g I_y^2 \end{pmatrix}, \quad \mathbf{b} = - \begin{pmatrix} \sum g I_x I_t \\ \sum g I_y I_t \end{pmatrix}.$$

When  $M$  has rank 2, then the LS estimate is  $\hat{\mathbf{u}} = M^{-1}\mathbf{b}$ .

**Implementation Issues** Usually we wish to estimate optical flow at every pixel, so we should express  $M$  and  $\mathbf{b}$  as functions of position  $\mathbf{x}$ , i.e.,  $M(\mathbf{x})\mathbf{u}(\mathbf{x})$ . Note that the elements of  $M$  and  $\mathbf{b}$  are local sums of products of image derivatives. An effective way to estimate the flow field is to first compute derivative images through convolution with suitable filters. Then, compute their products ( $I_x^2$ ,  $I_x I_y$ ,  $I_y^2$ ,  $I_x I_t$  and  $I_y I_t$ ), as required by Eq. 8.7. These quadratic images are then convolved with  $g(\mathbf{x})$ , to obtain the elemnts of  $M(\mathbf{x})$  and  $\mathbf{b}(\mathbf{x})$ .

In practice, the image derivatives will be approximated using numerical differentiation. It is important to use a consistent approximation scheme for all three directions.

## 8.9 Aperture Problem

When  $M$  in Eq. 8.7 is rank deficient one cannot solve for  $\mathbf{u}$ . This is often called the aperture problem as it invariably occurs when support  $g(\mathbf{x})$  is sufficiently local. However, the important issue is not the width of the image structure. However, the important issue is not the width of support, but rather the dimensionality of the image structure. Even for large regions, if the image is one-dimensional then  $M$  will be singular. When each image gradient within a region has the same spatial direction, it is easy to see that rank  $M = 1$ . Moreover, note that a single gradient constraint only provides the normal component of  $\mathbf{u}$ ,

$$\hat{\mathbf{u}} = \frac{-I_t}{\|\nabla I\|} \frac{\nabla I}{\|\nabla I\|}$$

## 8.10 Iterative Optical Flow Estimation

Equation 8.7 provides an optimal solution, but not to our original problem. Higher-

$$|\delta - \hat{\delta}| = \left| \frac{\delta^2}{f_1''(\mathbf{x})} \right| 2 \left| f_1'(\mathbf{x}) \right| + O(\delta^3)$$

For a sufficiently small displacement, and bounded  $|f_1''|/|f_1'|$ , we expect reasonably accurate estimates. This suggests a form of Gauss-Newton optimization in which we use the current estimate to *undo* the motion, and then we reapply the estimator to the *warped* signals to find the residual motion. This continues until the residual motion is sufficiently small.

In 2D, given an estimate of the optical flow field  $\mathbf{u}^0$ , we create a *warped* image sequence  $I^0(\mathbf{x}, t)$ :

$$I^0(\mathbf{x}, t + \delta t) = I(\mathbf{x} + \mathbf{u}^0 \delta t, t + \delta t), \quad (8.8)$$

where  $\delta t$  is the time between consecutive frames. Assuming that  $\mathbf{u} = \mathbf{u}^0 + \delta \mathbf{u}$ , from brightness constancy and Eq. 8.8 we get

$$I^0(\mathbf{x}, t) = I^0(\mathbf{x} + \delta \mathbf{u}, t + 1)$$

If  $\delta \mathbf{u} = 0$ , then clearly  $I^0$  would be constant through time (assuming brightness constancy). Otherwise, we can estimate the residual flow using

$$\delta \hat{\mathbf{u}} = M^{-1} \mathbf{b} \quad (8.9)$$

where  $M$  and  $\mathbf{b}$  are computed by taking spatial and temporal derivatives (differences) of  $I^0$ . The refined optical flow estimate then becomes

$$\mathbf{u}^1 = \mathbf{u}^0 + \delta \hat{\mathbf{u}}.$$

In an iterative manner, this new flow estimate is then used to rewrap the original sequence, and another residual flow can be estimated.

This iteration yields a sequence of approximate objective functions that converge to the desired objective function. At iteration  $j$ , given the estimate  $\mathbf{u}^j$  and the warped sequence  $I^j$ , our desired objective function is

$$\begin{aligned} E(\delta \mathbf{u}) &= \sum_{\mathbf{x}} g(\mathbf{x}) \left[ I(\mathbf{x}, t) - I(\mathbf{x} + \delta \mathbf{u}, t+1) \right]^2 \\ &= \sum_{\mathbf{x}} g(\mathbf{x}) \left[ I^j(\mathbf{x}, t) - I^j(\mathbf{x} + \delta \mathbf{u}^j, t+1) \right]^2 \\ &=: \tilde{E}(\delta \mathbf{u}). \end{aligned} \quad (8.10)$$

The gradient approximation to the difference to  $E(\delta \mathbf{u})$  gives the approximate objective function  $\tilde{E}$ . From  $\left| \tilde{d} - \delta \right|$  one can show that  $\tilde{E}$  approximates  $E$  to the second-order in the magnitude of the residual flow,  $\delta \mathbf{u}$ . The approximation error vanishes as  $\delta \mathbf{u}$  is reduced to zero. The iterative refinement with rewrapping reduces the residual motion at each iteration so that the approximate objective function converges to the desired objective function, and hence the flow estimate converges to the optimal LS estimate  $E(\delta \mathbf{u})$ .

The most expensive step at each iteration is the computation of image gradients and the matrix inverse in 8.9. One can, however, formulatio the problem so that the spatial image derivatives used to form  $M$  are taken at time  $t$ , and as such, do not depend on the current flow estimate  $\mathbf{u}^j$ . To see this, note that the spatial derivatives are computed at time  $t$  which leads to  $I(\mathbf{x}, t) = I^j(\mathbf{x}, t)$ . Of course  $\mathbf{b}$  in 8.9 will always depend on the warped image sequence and must be recomputed at each iteration. In practice, when  $M$  is not recomputed from the warped sequence then the spatial and temporal derivatives will not be centered at the same location in  $(\mathbf{x}, y, t)$  and hence more iterations may be needed. =====

$$\nabla I(\mathbf{x}, t) \cdot \mathbf{u} + I_t(\mathbf{x}, t) = 0.$$

This is called the *gradient constraint*

equation.

## 8.11 Pyramid/Coarse-to-fine

Limits of the (local) gradient method:

- 1 Fails when intensity structure within window is poor
- 2 Fails when displacement is large (typical operating range is motion of 1 pixel per iteration!). Linearization of brightness is suitable only for small displacements.
- 3 Brightness is no strictly constant in images. Actually less problematic than it appears, since we can pre-filter images to make them look similar.

In practice, our images have temporal sampling rates lower than required by the sampling theorem to uniquely reconstruct the continuous signal. AS a consequence, temporal aliasing is a common problem in motion estimation.

The spectrum of a translating signal is confined to a plane through the origin in the frequency domain. That is, if we construct a space-time signal  $f(\mathbf{x}, t)$  by translating a 2D signal  $f_0(\mathbf{x})$  with velocity  $\mathbf{u}$ , i.e.,  $f(\mathbf{x}, t) = f_0(\mathbf{x} - \mathbf{u}t)$ , one can show that the space-time Fourier transform of  $f(\mathbf{x}, t)$  is given by

$$\begin{aligned} F(\omega_x, \omega_y, \omega_t) &= F_0(\omega_x, \omega_y) \cdot \delta(\omega_1 \omega_x + \omega_2 \omega_y + \omega_t), \end{aligned} \quad (8.11)$$

where  $F_0$  is the 2D Fourier transform of  $f_0$ . Eq. 8.11 shows that the spectrum is nonzero only on a plane, the orientation of which gives the velocity. When the continuous signal is sampled in time, replicas of the spectrum are introduced at intervals of  $2\pi/T$  radians, where  $T$  is the time between frames, it is easy to see how this causes problems; i.e., the derivative filters may be more sensitive to the spectral replicas at high spatial frequencies than to the original spectrum on the plane through the origin.

Optical flow can be estimated at the coarsest scale of a Gaussian pyramid, where the image is significantly blurred, and the velocity is much slower (due to subsampling). The coarse-scale estimate can be used to warp the next (finer) pyramid level to *stabilize* its motion. Since the velocities after warping are slower, wider low-pass frequency band will be free of aliasing. One can therefore use derivatives at the finer scale to estimate the residual motion. This coarse-to-fine estimation continues until the finest level of the pyramid (the original image) is reached. Mathematically, this is identical to iterative refinement except that each scale's estimate must be up-sampled and interpolated before warping the next finer scale.

While widely used, coarse-to-fine methods have their drawbacks, usually stemming from the fact that fine-scale estimates can only be as reliable as their coarse-scale precursors; a poor estimate at one scale provides a poor initial guess at the next finer scale, and so on. That said, when aliasing does occur, one must use some mechanism such as coarse-to-fine estimation to avoid local minima in the optimization.

**8.12 Robust Motion Estimation**

The LS estimator is optimal when the gradient constraint errors, i.e.,

$$e(\mathbf{x}) := \mathbf{u} \cdot \nabla I(\mathbf{x}, t) + I_t(\mathbf{x}, t)$$

are mean-zero Gaussian, and the errors in different constraints are independent and identically distributed (IID.). Not surprisingly, this is a fragile assumption, for example, brightness constancy is often violated due to changing surface orientation, specularities reflections, or time-varying shadows. When there is significant depth variation in the scene , the constant motion

model is extremely poor, especially at occlusion boundaries. LS estimators are not suitable when the distribution of gradient constraint errors is heavy-tailed, as they are sensitive to small numbers of measurement outliers. It is therefore often crucial that the quadratic estimation in Eq. 8.6 be replaced by a robust estimator,  $\rho(\cdot)$ , which limits the influence of constraints with larger errors:

$$E(\mathbf{u}) = \sum_{\mathbf{x}} g(\mathbf{x}) \rho(e(\mathbf{x}, \sigma) \cdot \sigma).$$

For example the redescending Geman-McClure estimator

$$\rho(e, \sigma) = e^2 / \left( e^2 + \sigma^2 \right),$$

where  $\sigma^2$  determines the range of constraint errors for which influence is reduced.

## 8.13 Motion Models

Thus far we have assumed that the 2D velocity is constant in local neighbourhoods. Nevertheless, even for small regions this is often a poor assumption. We now consider generalizations to more interesting motion models.

**Affine Model** General first-order affine motion is usually a better model of local motion than a translational model. An affine velocity field centered at location  $\mathbf{x}_0$  can be expressed in matrix form as

$$\mathbf{u}(\mathbf{x}; \mathbf{x}_0) = A(\mathbf{x}; \mathbf{x}_0) \mathbf{c}, \quad (8.12)$$

where  $c_1 \ c_2 \ c_3 \ c_4 \ c_5 \ c_6$  is the motion model parameters, and

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

From Eq. and we get the gradient constraint equation

$$\nabla I(\mathbf{x}, t) A(\mathbf{x}; \mathbf{x}_0) \mathbf{c} + I_t(\mathbf{x}, t) = 0,$$

for which the LS estimate for the neighbourhood has the form

$$\hat{\mathbf{c}} = M^{-1} \mathbf{b} \quad (8.13)$$

where now  $M$  and  $\mathbf{b}$  are given by

$$\begin{aligned} M &= \sum_{\mathbf{x}} g A^T \nabla I^T \nabla I A, \\ \mathbf{b} &= - \sum_{\mathbf{x}} g A^T \nabla I^T I_t. \end{aligned}$$

When  $M$  is rank deficient there is insufficient image structure to estimate the six unknowns. Affine models often require larger support than constant models, and one may need a robust estimator instead of the LS estimator.

Iterative refinement is also straightforward with affine motion models. Let the optimal affine motion be  $\mathbf{u} = A\mathbf{c}$ , and let affine estimate at iteration  $j$  be  $\mathbf{u}^j = A\mathbf{c}^j$ . Because the flow is linear in the motion parameters, it follows that  $\delta \mathbf{u} = \mathbf{u} - \mathbf{u}^j$  and  $\delta \mathbf{c} = \mathbf{c} - \mathbf{c}^j$  satisfy

$$\delta \mathbf{u} = A \delta \mathbf{c}.$$

Accordingly, defining  $I^j(\mathbf{x}, t)$  to be the original sequence  $I(\mathbf{x}, t)$  warped by  $\mathbf{u}^j$  as in Eq. 8.8 we use the same LS estimator as in Eq. 8.13, but with  $I$  and  $\hat{\mathbf{c}}$  replaced by  $I^j$  and  $\hat{\delta \mathbf{c}}$ .

## 8.14 Low-order Parametric Deformations

There are many other polynomial and rational deformations that make useful motion models. *Similarity deformations*, comprising translation ( $d_1, d_2$ ), 2D rotation  $\theta$ , and uniform scaling by  $s$  are a special case of the affine model, but still very useful in practice. In a neighborhood centered at  $\mathbf{x}_0$  it has the form as Eq. 8.13, but with  $\mathbf{c} = (d_1, d_2, s \cos \theta, s \sin \theta)^T$  and

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

## 8.15 Global Smoothing

While area-based regression is commonly used, some of the earliest formulations of optical flow estimation assumed smooth-

ness through nonparametric motion models, rather than an explicit parametric model in each local neighbourhood. One such energy functional was proposed by Horn and Schunck in Eq. A.2. A key advantage of global smoothing is that it enable propagation of information over large distances in the image. In image regions of nearly uniform intensity, such as a blank wall or tabletop, local methods will often yield singular (or poorly conditioned) systems of equations. Global methods can *finn* in the optical flow from nearby gradient constraints.

The equation above can be minimized directly with discrete approximations to the integral and the derivatives. This yields a large system of linear equations. The main disadvantage of global methods is computational efficiency. Another problem is in the setting of the *regularization parameter*  $\lambda$  that determines the amount of desired smoothing.

## 8.16 Probabilistic Formulations

One problem with the above estimators is that, although they provide useful estimates of optical flow, they do not provide confidence bounds. Nor do they show how to incorporate any prior information one might have bout motion to further constrain the estimates. As a result, one may not be able to propagate flow estimates from one time to the next, nor know how to weight them when combining flow estimates from different information sources. These issues can be addressed with a probabilistic formulation.

The cost function 8.10 has a simple probabilistic interpretation. Up to normalization constants, it corresponds to the log likelihood of a velocity under the assumption that intensity is conserved up to gaussian noise.

$$I(\mathbf{x}, t) = I(\mathbf{x} + \mathbf{u}, t + 1) + \eta.$$

If we assume that the same velocity  $\mathbf{u}$  is shared by all pixels within a neighbourhood, that  $\eta$  is white Gaussian noise with standard deviation  $\sigma$ , and uncorrelated at different pixels, we obtain the conditional density

$$p(I | \mathbf{u}) \propto e^{E(\mathbf{u})/2\sigma^2}.$$

=====

## 8.17 Parametric motion models

Global miton models offer:

- 1 More constrained solutions than smoothness (Horn-Schunck)
- 2 Integration over a large area than a translation-only model can accomodate (Lucas-Kanade)

## 9 Video Compression

### 9.1 Perception of motion

Perception of motion: Human visual system is specifically sensitive to motion. Eyes follow motion automatically. Some distortions are not as perceivable as in image coding (would be if we froze frame). No good psycho-visual model available. Vusal perception is limited to < 24Hz. Asussession of images will be perceived as continuous if frequency is sufficiently high. Cinema 24/24Hz, TV 25 Hz or 50 Hz. We still nee to avoid aliasing (wheel effect). High-rendering frames desired in computer games (needed due to absence of motion blur). Flicker can be perceived up to > 60Hz in particular in periphery. Issue addressed by 100 Hz TV.

### 9.2 Interlaced video format

Two temporally shifted half images, increase of frequency 25 Hz → 50 Hz. Reduction of spatial resolution. Full image representation: progressive.

### 9.3 Why compressed video?

Raw HD TV signal 720p @50 Hz:  
1280 · 720 · 50 · 24 bits/s  
= 1 105 920 000 bits/s > 1 Gb/s  
Only 20 Mb/s HDTV channel bandwidth requires compression of factor of 60 (0.4 bits/pixel on average)

### 9.4 Lossy video compression

Take advantage of redundancy. Spatial correlation between neighboring pixels. Temporal correlation between frames. Drop perceptuall unimportant details.

**Temporal Redundancy** Take advantage of similarity between successive frames

**Temporal processing** Usually high frame rate: Significant temporal redundancy. Possible representations along temporal dimension:

- 1 Transform/subband methods: Good for textbook case of constant velocity uniform global motion. Inefficient for nonuniform motion, i.e. real-world motion. Requires large number of frame stores which leads to delay. (Memory cost may also be an issue.) Is ineffective for many scene changes or high motion.
- 2 Predictive methods: Good performance using only 2 frame stores. However, simple frame differencing is not enough...

**Goal** Exploit the temporal redundancy

**Predict current frame** based on previously coded frames

#### Types of coded frames:

- 1 I-frame: Intra-coded frame, coded independently of all other frames.
- 2 P-frame: Predictively coded frame, coded based on previously coded frame I or P. Can send motion vector plus changes.
- 3 B-frame: Bi-directionally predicted frame, coded based on both previous and future coded frames I and P. In case something is uncovered.

#### Motion-compensated prediction

Simple frame differencing *fails* when there is motion. Must account for motion. → Motion-compensated (MC) prediction. MC-prediction generally provides significant improvements. Questions: How can we estimate motion? How can we form MC-prediction?

#### Ideal situation

- 1 Partition video into moving objects
- 2 describe object motion → Generally very difficult

### Practical approach

**Block-Matching** Motion Estimation:

- 1 Partition each frame into blocks, e.g. 16 × 16 pixels
- 2 Describe motion of each block → No object identification required and good, robust performance.

## 9.5 Block-matching motion estimation

**Assumptions**

- 1 Translational motion within block:

$$f(n_1, n_2, k_{cur}) = f(n_1 - mv_1, n_2 - mv_2, k_{ref}).$$

**ME Algorithm**

- 1 Divide current frame into non-overlapping  $N_1 \times N_2$  blocks.
- 2 For each block, find the best matching block in reference frame.

### 9.5.1 Determining the best matching block

For each block in the current frame, search for best matching block in the reference frame.

for determining "best match":

$$\begin{aligned} \text{Metrics} &= \sum_{n_1, n_2 \in \text{Block}} [f(n_1, n_2, k_{cur}) - f(n_1 - mv_1, n_2 - mv_2, k_{ref})]^2. \\ \text{MAE} &= \sum_{n_1, n_2 \in \text{Block}} |f(n_1, n_2, k_{cur}) - f(n_1 - mv_1, n_2 - mv_2, k_{ref})|. \end{aligned}$$

**Candidate blocks** All blocks in, e.g.  $(\pm 32, \pm 32)$  pixel area

**Strategies for searching** candidate blocks





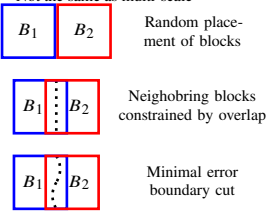
Reaction-diffusion systems. It is believed that these systems show behaviors which are qualitatively equivalent to real processes (Morphogenesis) found in the nature, such as animal markings (shells, fish, wild cats. . .).

Efros & Leung '99 extended

Observation: neighbor pixels are highly correlated. Idea: Unit of synthesis is block

- Exactly the same but now we want  $P(B|N(B))$ .

- Much faster: Synthesize all pixes in a block at once
- Not the same as multi-scale



- Algorithm
- Pick size of block and size of overlap
  - Synthesize blocks in raster order
  - Search input texture for bolck that satisfies overlap constraints (above and left)
  - Paste new block into resulting texture. Use dynamic programming to compute minimal error boundary cut

12.3 Texture Transfer

Take the texture from one object and "paint" it onto another object. This requires spearating texture and shape. That's hard, but we can cheat. Assume we can capture shape by baundary and rough shading. Then, just add another constraint when sampling: similarity to underlying image at that spot.

**Shape from textures** Possibility to re-cover normals both for deterministic and probabilistic textures.

13 Questions

- $I_{\text{comp}} = I_{\alpha} I_a + (1 - I_{\alpha}) I_b$
- MAP, Maximum a posteriori detector.
- graph cuts
- Solve MRFs with graph cuts
- impulse response  $r(-x, -y)$
- Canny nonmaxima suppression
- Entropy Coding (Huffman code)
- Aperture problem: normal flow
- Lucas-Kanade: Iterative refinement/local gradient method
- Coarse-to-fine-estimation
- SNR scalability EI, EP frame
- MPEG Structure

- Radon transform: Just FT with first Fourier coefficient?
- Sparse representations for image restoration: Inpainting, demosaick-ing

A Big equations

$$\mathcal{F}[h](u, v) = \frac{1}{2\ell} \int_{-\ell}^{\ell} dx_1 \exp(-i2\pi u x_1) \cdot \underbrace{\int_{-\infty}^{\infty} dx_2 \delta(x_2) \exp(-i2\pi v x_2)}_{=1}$$
$$= \text{sinc}(2\pi u \ell)$$
(A.1)

$$E = \iint dx dy \left[ \left( \frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{\partial y}{\partial t} + \frac{\partial I}{\partial t} \right)^2 + \alpha^2 (\|\nabla \dot{x}\| + \|\nabla \dot{y}\|)^2 \right]$$
(A.2)

$$\mathbf{v} = \left( \frac{\sum_i w_i I_x(q_i)^2}{\sum_i w_i I_x(q_i) I_y(q_i)} \frac{\sum_i w_i I_x(q_i) I_y(q_i)}{\sum_i w_i I_y(q_i)^2} \right)^{-1} \cdot \begin{pmatrix} -\sum_i w_i I_x(q_i) I_t(q_i) \\ -\sum_i w_i I_y(q_i) I_t(q_i) \end{pmatrix}$$
(A.3)