

## Contents

### Contents

1	The digital image	9.7	Block Matching Algorithm
1.1	Image as 2D signal	9.8	Frame types
1.2	Image sources	9.9	Summary of Temporal Processing
1.3	Sampling	9.10	Basic Video Compression Architecture
1.4	Reconstruction	9.11	Scalable Video Coding
1.5	Quantization	9.12	Standards
1.6	Image Properties	9.13	Quality measure
1.7	Image Noise		
1.8	Colour Images		

2	Image segmentation		
2.1	Thresholding		
2.2	Region Growing		
2.2.1	Variations		
2.3	Spatial relations		
2.4	Morphological Operations		
2.4.1	Medial Axis Transform (MAT, skeletonization)		

3	Image filtering		
3.1	Linear Shift-Invariant Filtering		

4	Image features		
4.1	Template matching		
4.2	Edge detection		
4.2.1	Canny edge detector		
4.3	Feature detection		
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	Fourier Transform		
5.1	Aliasing		
5.2	Definition		
5.3	Sampling		
5.4	Image Restoration		

6	Unitary transforms		
6.1	Karhunen-Loeve Transform		
6.2	Basic images and eigen-images (E)		
6.3	Eigenfaces (EF)		
6.4	Fisherfaces/LDA		
6.5	JPEG image compression		

7	Scale-space representations		
7.1	Image pyramid		
7.2	Applications		
7.3	Pyramids		
7.4	Haar Transform		

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

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

**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

**1D** 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.

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

**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  $c \sim N(0, \sigma^2)$ . So that  $p(c) = (2\pi\sigma^2)^{-1} e^{-c^2/2\sigma^2}$ .

**Poisson noise** (shot noise)

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

**Rician noise** (appears in MRI)

$$p(I) = \frac{1}{\sigma^2} \exp\left(-\frac{I^2 + f^2}{2\sigma^2}\right) I_0\left(\frac{If}{\sigma^2}\right)$$

**Multiplicative noise**  $I = f + f \cdot c$

**Signal to noise ration (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.

being values and costs. For simplicity, often  $V_T N = V_T P = 0$ .

**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 models. 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 info, for example.

**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:

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

### 2.2.3 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 is 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

$$\begin{aligned} \text{Energy}(y; \theta, \text{data}) &= \sum_i \psi_1(y_i; \theta, \text{data}) + \sum_{\langle i, j \rangle} \psi_2(y_i, y_j; \theta, \text{data}) \\ \psi &= i, j \text{ edges} \end{aligned}$$

**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

$S_n(X) = (X \ominus_n B) \setminus [(X \ominus_n B) \circ B]$  where  $\ominus_n$  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*.

$$X = \bigcup_{n=0}^{\infty} S_n(X) \oplus_n B$$

**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:

$I'(x, y) = \sum_{i,j \in N(x,y)} K(x, y; i, j) I(i, j)$  where  $I$  = input image;  $I'$  = output of operation.  $K$  is *kernel* of the operation.  $N(m, n)$  is a neighbourhood of  $(m, n)$ .

**Correlation** e.g. template matching. Linear operation:  $I' = K I$

$I'(x, y) = \sum_{i,j \in N(x,y)} K(i, j) I(x + i, y + j)$   $A = I(1)$ ,  $A + C = I(3)$ ,  $A + B = I(2)$ ,  $A + B + C + D = I(4)$ .

**Edge** The filter window falls off the edge of the image, we need to extrapolate, methods: (1) clip filter (black) (2) wrap around (3) copy edge (4) reflect across edge (5) vary filter near edge

**Filter at boundary** (1) ignore, copy or truncate. No processing of boundary pixels. Pad image with zeros (matlab). Pad image with copies of edge rows/columns (2) truncate kernel (3) reflected indexing (4) circular indexing

**Separable kernels** Separable filters can be written

$$K(m, n) = f(m)g(n)$$

for a rectangular neighbourhood with size  $(2M + 1) \times (2N + 1)$ ,  $I'(m, n) = f * (g * I(N(m, n)))$ ,  $I''(m, n) = \sum_{j=-N}^N g(j) I(m, n - j)$

$$I'(m, n) = \sum_{i=-M}^M f(i) I''(m - i, n) = \sum_{x,y=-\infty}^{\infty} [s(x, y) - t(x - p, y - q)]^2$$

$$= \sum_{x,y=-\infty}^{\infty} |s(x, y)|^2 + |t(x, y)|^2 - 2 \cdot \sum_{x,y=-\infty}^{\infty} s(x, y) \cdot t(x - p, y - q)$$

$$= s(p, q) * t(-p, -q) \leq \sqrt{\sum |s(x, y)|^2} \cdot \sqrt{\sum |t(x, y)|^2}$$

**Gaussian Kernel** Idea: Weight contributions of neighboring pixels:

$$N_{\mu=0, \sigma}(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{x^2+y^2}{2\sigma^2}}$$

Smoothing with a Gaussian instead of a box filter removes the artefact of the verti-

inequality was used. Equality  $\Rightarrow s(x, y) = \alpha \cdot t(x - p, y - q)$  with  $\alpha = N(x)/N(y)$ . Amount of smoothing depends on  $\sigma$  and window size. Width  $> 3\sigma$ .

**Scale space** Convolution of a Gaussian with  $\sigma$  with itself is a gaussian with  $\sigma\sqrt{2}$ . Repeated convolution by a Gaussian filter produces the scale space of an image.

**Gaussian filter top-5** (1) Rotationally symmetric. (2) Has a single lobe.  $\rightarrow$  Neighbor's influence decreases monotonically. (3) Still one lobe in frequency domain.  $\rightarrow$  No corruption from high frequencies (4) Simple relationship to  $\sigma$  (5) Easy to implement efficiently

**Differential filters** • Prewitt operator:  $\begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$ . • Sobel operator:  $\begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}$ . • Laplacian operator:  $\begin{pmatrix} 0 & 1 & 0 \\ -1 & 4 & -1 \\ 0 & 1 & 0 \end{pmatrix}$ . • High-pass filter:  $\begin{pmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{pmatrix}$ .

**High-pass filters** • Laplacian operator:  $\begin{pmatrix} 0 & 1 & 0 \\ -1 & 4 & -1 \\ 0 & 1 & 0 \end{pmatrix}$ . • High-pass filter:  $\begin{pmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{pmatrix}$ .

**Differentiation and convolution**  $\frac{\partial f}{\partial x} = \lim_{\varepsilon \rightarrow 0} \left( \frac{f(x + \varepsilon, y) - f(x, y)}{\varepsilon} \right)$

**Laplacian operator** Detects discontinuities by considering second derivative  $\nabla^2 f(x, y) = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$ . Isotropic (rotationally invariant) operator, zero-crossings mark edge location, discrete-space approximation by convolution with  $3 \times 3$  impulse response  $\begin{pmatrix} 0 & 1 & 0 \\ -1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ , or  $\begin{pmatrix} 0 & 1 & 0 \\ 1 & -8 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ .

**Laplacian of Gaussian** The Laplacian operator is very sensitive to fine detail and noise, so blur it first with Gaussian.  $\rightarrow$  do it in one operator Laplacian of Gaussian (LoG)

$\text{LoG}(x, y) = -\frac{1}{\pi\sigma^4} \left[ 1 - \frac{x^2 + y^2}{2\sigma^2} \right] \cdot e^{-\frac{x^2 + y^2}{2\sigma^2}}$

**Integral images** integral images (also known as summed-area tables) allow to efficiently compute the convolution with a constant rectangle  $I'(x, y) = I + \alpha [K * I]$ , where  $K$  is a high-pass filter kernel and  $\alpha \in [0, 1]$ .

**Integral images** integral images (also known as summed-area tables) allow to efficiently compute the convolution with a constant rectangle  $I'(x, y) = \int_0^x dx' \int_0^y dy' I(x', y')$

**4.2.1 Canny edge detector** (1) Smooth image with a gaussian filter (2) Compute gradient magnitude and angle (Sobel, Prewitt, ...)

$M(x, y) = \sqrt{(\partial_x f)^2 + (\partial_y f)^2}$ ,  $\alpha(x, y) = \arctan(\partial_y f / \partial_x f)$

(3) Apply nonmaxima suppression to gradient magnitude image (4) Double thresholding to detect strong and weak edge pixels (5) Reject weak edge pixels not connected with strong edge pixels

**Canny nonmaxima suppression** Quantize edge normal to one of four directions: horizontal,  $-45^\circ$ , vertical,  $45^\circ$ . If  $M(x, y)$  is smaller than either of its neighbors in edge normal direction  $\rightarrow$  suppress; else keep

**Double-thresh. of grad. magn.** **strong edge:**  $M(x, y) \geq \theta_{\text{high}}$  **weak edge:**  $\theta_{\text{high}} > M(x, y) \geq \theta_{\text{low}}$

Typical setting:  $\theta_{\text{high}}, \theta_{\text{low}} = 2, 3$ . Region labeling of edge pixels. Reject regions without strong edge pixels.

**Robustness of Harris corner detector** (1) Invariant to brightness offset:  $f(x, y) \rightarrow f(x, y) + c$  (2) Invariant to shift and rotation (3) Not invariant to scaling

**4.3.4 Lowe's SIFT features** Recover features with position, orientation and scale.

**Position** (1) Look for strong responses of DoG filter, (2) only consider local maxima.

**Scale** (1) Look for strong responses of DoG filter over scale space. (2) only consider local maxima in both position and scale. (3) Fit quadratic around maxima for subpixel accuracy.

**Orientation** (1) Create histogram of local gradient directions computed at selected scale. (2) Assign canonical orientation at peak of smoothed histogram (3) Each key specifies stable 2D coordinates  $(x_0, y_0, r)$

**SIFT description** (1) Thresholded image gradients are sampled over  $16 \times 16$  pixel



array of locations in space. (2) Create an array of orientation histograms (3) 8 orientations  $\times 4 \times 4$  histogram array = 128 dimensions

## 5 Fourier Transform

### 5.1 Aliasing

One can't shrink an image by taking every second pixel. If we do, characteristic errors appear. Typically, small phenomena look bigger; fast phenomena can look slower. Common phenomena (1) Wagon wheels rolling the wrong way in movies. (2) Checkerboards misrepresented in ray tracing (3) Striped shirts look funny on color television.

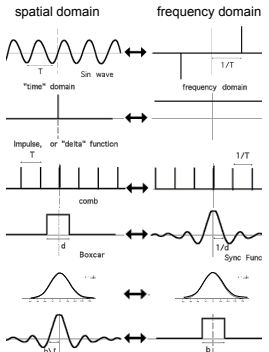
### 5.2 Definition

Represent function on a new basis. Basis elements have the form  $e^{-i2\pi(u x + v y)}$ . The Fourier transform is

$$\hat{f}(u, v) = \int_{\mathbb{R}^2} dx dy f(x, y) e^{-i2\pi(u x + v y)} = \frac{1}{2\ell} [\theta(x_1 + \ell) - \theta(x_1 - \ell)] \delta(x_2)$$

Basis functions of Fourier transform are eigenfunctions of linear systems.

#### Important functions



**Convolution theorem** (1) The Fourier transform of the convolution of two functions is the product of their Fourier transforms

$$\hat{f} \cdot \hat{g} = \widehat{f * g}$$
$$\widehat{f \cdot g} = \mathcal{F}(f \cdot g)$$

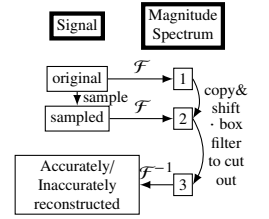
### 5.3 Sampling

Go from continuous world to discrete world, from function to vector. Samples are typically measured on regular grid. We want to be able to approximate integrals sensibly -> Delta function

$$\delta_{2D}(f(x, y)) = \sum_{i,j=-\infty}^{\infty} f(x, y) \delta(x - i, y - j) = f(x, y) \sum_{i,j=-\infty}^{\infty} \delta(x - i, y - j)$$

with  $S$  = Sample operator.

#### FT of sampled signal



In the figure above the accuracy depends on the overlapping wave functions in "2". The box filter then can't cut out appropriately the magnitude spectrum to get a proper result in "3". This leads to an inaccurately reconstructed signal.

**Proper sampling** To avoid this effect, this is the procedure: original signal -> lp filtering -> lp filt. sign -> sample -> sampl.sign. -> reconstr. -> reconstr.sign

**Smoothing as low-pass filtering** The message of the FT is that high

frequencies lead to blurring with sampling. Solutions: press high frequencies before sampling. A filter whose FT is a box is bad, because the filter kernel has infinite support. Common solution: use a Gaussian.

### 5.4 Image Restoration

**Nyquist sampling theorem** Nyquist theorem: The sampling frequency must be at least twice the highest frequency.  $\omega_s \geq 2\omega$ . If this is not the case, the signal needs to be bandlimited before sampling, e.g. with a low-pass filter.

#### Pixelization

Possibilities: Square pixels, Gaussian reconstruction filter, Bilinear interpolation, perfect reconstruction filter.

**Motion blurring** Each light dot is transformed into a short line along the  $x_1$ -axis:  $h(x_1, x_2)$

$$h(x_1, x_2) = \frac{1}{2\ell\sigma^2} \exp\left(-\frac{x_1^2 + x_2^2}{2\sigma^2}\right)$$

**Noise** Gaussian blurring kernel:

$$\widehat{h(x)} = \frac{1}{2\ell\sigma^2} \exp\left(-\frac{x_1^2 + x_2^2}{2\sigma^2}\right)$$

**Problem**  $f(x) \xrightarrow{h(x)} g(x) \xrightarrow{\hat{h}(x)} f(x)$ . The "inverse" kernel  $\hat{h}(x)$  should compensate the effect of the image degradation  $h(x)$ , i.e.,

$$(\hat{h} * h)(x) = \delta(x)$$

$\hat{h}(x)$  may be determined more easily in Fourier space:

$$\mathcal{F}[\hat{h}](u, v) \cdot \mathcal{F}[h](u, v) = 1$$

To determine  $\mathcal{F}[\hat{h}]$ , we need to estimate (1) the distortion model  $h(x)$  (point spread function) or  $\mathcal{F}[h](u, v)$  (modulation transfer function) (2) the parameters of  $h(x)$ , e.g. for deblurring

**Motion Blur FT** Eq. 14 Problem:  $\mathcal{F}[\hat{h}](u) = 1/\hat{h}(u)$ . sinc has many zeroes and these frequencies can't be recovered! Solution: Regularized reconstruction filter

$$\hat{F}[\hat{h}](u, v) = \frac{\mathcal{F}[h]}{\|\mathcal{F}\|^2 + \varepsilon}$$

Singularities are avoided by the regularization  $\varepsilon$ .

**Space-time super-resolution** One can put two movies of the same thing and merge their frames for space and time super-resolution.

**Spatial super-resolution** • lens + pixel = low-pass filter (edised to avoid aliasing) • Low-res images =  $D * H * G * S$  (desired high-res image). D:decimate, H:lens+pixel, G: Geometric warp • Simplified case for translation:  $LR = (D * G) * (H * HR)$ . G is shift-invariant and commutes with H. First compute H HR, then deconvolve HR with H. • Super-resolution needs to restore attenuated frequencies. Many images improve  $S/N$  ratio  $\sim \sqrt{n}$ , which helps. Eventually Gaussian's double exponential always dominates.

## 6 Unitary transforms

Digital image as a matrix:

$$f = \begin{bmatrix} f(0,0) & \dots & f(N-1,0) \\ \vdots & & \vdots \\ f(0,L-1) & \dots & f(N-1,L-1) \end{bmatrix}$$
$$= f_{yx}$$

or as a vector

$$\mathbf{f} = \begin{pmatrix} f(0,0) \\ \vdots \\ f(N-1,L-1) \end{pmatrix}$$

**General approach** (1) Sort samples  $f(x, y)$  of an  $M \times N$  image (or rectangular block in the image) into column vector of length  $M \times N$ . (2) Compute transform coefficients  $\mathbf{c} = \mathbf{A}^T \mathbf{f}$  where  $\mathbf{A}$  is a matrix of size  $(MN)^2$ . (3) Transform  $\mathbf{A}$  is unitary, iff  $\mathbf{A}^{-1} = \mathbf{A}^*$  (4) If  $\mathbf{A}$  is real-valued, i.e.  $\mathbf{A} = \mathbf{A}^T$ , transform is orthonormal.

**Energy conservation**  $\|\mathbf{c}\|^2 = \|\mathbf{f}\|^2$

**Image collection**  $f_i$  one image,  $F = f_1, \dots, f_n$ .

#### Auto-correlation function

$$R_{ff} = E[f_i f_i^*] = FF^* / n$$

**Energy distribution** Energy is conserved, but often will be unevenly distributed among coefficients. Autocorrelation matrix:

$$\mathbf{R}_{cc} = E[\mathbf{c}\mathbf{c}^*] = E[\mathbf{A}\mathbf{f}\mathbf{f}^* \mathbf{A}^*] = \mathbf{A} \mathbf{R}_{ff} \mathbf{A}^*$$

Mean squared values ("average energies") of the coefficients  $c_i$  are on the diagonal of  $\mathbf{R}_{cc}$ .

$$E[c_i^2] = [\mathbf{R}_{cc}]_{ii} = [\mathbf{A} \mathbf{R}_{ff} \mathbf{A}^*]_{ii}$$

**Eigenmatrix of autocorrelation matrix** Definition: Eigenmatrix  $\Phi$  of autocorrelation matrix  $\mathbf{R}_{ff}$ . (1)  $\Phi$  is unitary (2) The columns of  $\Phi$  form a set of eigenvectors of  $\mathbf{R}_{ff}$ , i.e.,

$$\mathbf{R}_{ff} \mathbf{f} = \Phi \Lambda$$

where  $\Lambda = \text{diag}(\lambda_0, \dots, \lambda_{MN-1})$ . (3)  $\mathbf{R}_{ff}$  is symmetric nonnegative definite, hence  $\lambda_i \geq 0$  for all  $i$  (4)  $\mathbf{R}_{ff}$  is normal matrix, i.e.  $\mathbf{R}_{ff}^* \mathbf{R}_{ff} = \mathbf{R}_{ff} \mathbf{R}_{ff}^*$ , hence unitary eigenmatrix exists.

### 6.1 Karhunen-Loeve Transform

Strongly correlated samples with equal energies  $\xrightarrow{A}$  uncorrelated samples, most of the energy in first coefficient.

**Properties** (1) Unitary transform with matrix  $\mathbf{A} = \Phi^*$  where the columns of  $\Phi$  are ordered according to decreasing eigenvalues. (2) Transform coefficients are pairwise uncorrelated

$\mathbf{R}_{cc} = \mathbf{A} \mathbf{R}_{ff} \mathbf{A}^* = \Phi^* \mathbf{R}_{ff} = \Phi^* \Phi \Lambda = \Lambda$  (3) Energy concentration property: No other unitary transform packs as much energy into the first  $J$  coefficients, where  $J$  is arbitrary. Mean squared approximation error by choosing only first  $J$  coefficients is minimized.

**Optimal energy concentration** (1) To show optimum energy concentration property, consider the truncated coefficient vector  $\mathbf{b} = \mathbf{I}_J \mathbf{c}$ , where  $\mathbf{I}_J$  contain ones on the first  $J$  diagonal positions, else zeros. (2) Energy in first  $J$  coefficients for arbitrary transform  $\mathbf{A}$

$$E = \text{tr}(\mathbf{R}_{bb}) = \text{tr}(\mathbf{I}_J \mathbf{R}_{cc} \mathbf{I}_J^T)$$
$$= \text{tr}(\mathbf{I}_J \mathbf{A} \mathbf{R}_{ff} \mathbf{A}^* \mathbf{I}_J) = \sum_{k=0}^{J-1} a_k^T \mathbf{R}_{ff} \bar{a}_k$$

where  $a_k^T$  is the  $k$ -th row of  $\mathbf{A}$ . (3) Lagrangian cost function to enforce unit-length basis vectors

$$L = E + \sum_{k=0}^{J-1} \lambda_k (1 - a_k^T \bar{a}_k)$$
$$= \sum_{k=0}^{J-1} a_k^T \mathbf{R}_{ff} \bar{a}_k + \sum_{k=0}^{J-1} \lambda_k (1 - a_k^T \bar{a}_k)$$

Differentiating  $L$  with respect to  $a_j$  yields necessary condition

$$\mathbf{R}_{ff} \bar{a}_j = \lambda_j \bar{a}_j, \quad \forall j < J$$

### 6.2 Basis images and eigenimages (EI)

For a unitary transform, the inverse transform  $\mathbf{f} = \mathbf{A}^* \mathbf{c}$  can be interpreted in terms of the superpositions of "basis images" (columns of  $\mathbf{A}^*$ ) of size  $MN$ . If the transform is a KL transform, the basis images, which are the eigenvectors of the autocorrelation matrix  $\mathbf{R}_{ff}$ , are called "eigenimages". If energy concentration works well, only a limited number of eigenimages is needed to approximate a set of images with small error. These eigenimages form an optimal linear subspace of dimensionality  $J$ .

**EI for recognition** To recognize complex patterns (e.g., faces), large portions of an image (say of size  $MN$ ) might have to be considered. High dimensionality of "image space" means high computational burden for many recognition techniques. Transform  $\mathbf{c} = \mathbf{W} \mathbf{f}$  can reduce dimensionality from  $MN$  to  $J$  by representing the image by

$\mathbf{c}$  coefficients. Idea: tailor a KLT to the specific set of images of the recognition task to preserve the salient features.

**Simple recognition** Simple Euclidean distance (SSD) between images. Best match wins

$$\text{argmin}_i D_i = \|\mathbf{I}_i - \mathbf{I}\|$$

Computationally expensive, i.e. requires presented image to be correlated with every image in the database!

**Eigenspace matching** Let  $\mathbf{I}_i$  be the input image,  $\mathbf{I}$  the database. The "character" of the face  $\hat{\mathbf{J}} = \mathbf{J} - \langle \mathbf{J} \rangle$ , with  $\mathbf{J}$  being any image (set). Do KLT (aka PCA) transformation

$$\hat{\mathbf{I}}_i \mapsto p_i, \quad E^* \hat{\mathbf{I}}_i = p_i$$
$$\hat{\mathbf{I}}_i \sim \hat{\mathbf{I}}_i \approx E p_i$$
$$\|\hat{\mathbf{I}}_i - \mathbf{I}\| \approx \|p_i - p\|$$

with closest rank-k approximation property of SVD. Approximate

$$\text{argmin}_i D_i = \|\mathbf{I}_i - \mathbf{I}\| \approx \|p_i\|$$

### 6.3 Eigenfaces (EF)

Concatenate face pixels into "observation vector",  $\mathbf{x}$ .

**EI for recognition** (1) Input image, (2) normalize, (3) subtract mean face, (4) KLT, (5) Find most similar  $p_i$ , (6) similarity measure, (7) rejection system, (8) result of identification.

**Limitations of EFs** Differences due to varying illumination can be much larger than differences between faces!

### 6.4 Fisherfaces/LDA

Training data: For eigenfaces distance of difference of illumination are within individual variance. Key idea: Find directions where ratio of between/within individual variance are maximized. Linearly project to basis where dimension with good signal to noise ratio are maximized.

**Fisher linear discriminant analysis** Eigenimage method maximizes "scatter" within the linear subspace over the entire image set - regardless of classification task

$$E_{\text{opt}} = \text{argmax}_E (\det(E R E^*))$$

Fisher linear discriminant analysis: Maximize between-class scatter, while minimizing within-class scatter,

$$F_{\text{opt}} = \text{argmax}_F \left( \frac{\det(F R_B W^*)}{\det(F R W F^*)} \right)$$

$$\mathbf{R}_B = \sum_{i=1}^c N_i (\mathbf{\mu}_i - \mathbf{\mu})(\mathbf{\mu}_i - \mathbf{\mu})^*$$

$$\mathbf{R}_W = \sum_{i=1, \dots, c} (\Gamma_i \mathbf{\ell} - \mathbf{\mu}_i)(\Gamma_i \mathbf{\ell} - \mathbf{\mu}_i)^*$$

$\Gamma_i \mathbf{\ell} \in \text{Class}(i)$

$N_i$  are the samples in class  $i$  and  $\mathbf{\mu}_i$  is the mean in class  $i$ . Solution: Generalized eigenvectors  $\mathbf{w}_i$  corresponding to the  $k$  largest eigenvalues  $\{\lambda_i | i = 1, \dots, k\}$ , i.e.

$$\mathbf{R}_B \mathbf{w}_i = \lambda_i \mathbf{R}_W \mathbf{w}_i, \quad i = 1, \dots, k$$

Problem: within-class scatter matrix  $\mathbf{R}_W$  at most of rank  $L - c$ , hence usually singular. Apply KLT first to reduce dimension of feature space to  $L - c$  (or less), proceed with Fisher LDA in low-dimensional space.

**Eigenfaces vs. Fisherfaces** Eigenfaces conserve energy but the two classes e.g. in 2D are no longer distinguishable. FLD (Fisher LDA) separates the classes by choosing a better 1D subspace. Fisher faces are much better in varying illuminations.

**Varying illumination (FF)** All images of same Lambertian surface with different illumination (without shadows) lie in a 3D linear subspace. Single point source at infinity

$f(x, y) = a(x, y) \left( \mathbf{\ell}^T \mathbf{n}(x, y) \right) L$ ,  $a(x, y)$  surface albedo,  $L$  light source intensity. Superposition of arbitrary num-

ber of point sources at infinity is still in same 3D linear subspace, due to linear superposition of each contribution to image. Fisher images can eliminate within-class scatter.

#### Appearance manifold approach

For every object, (1) sample the set of viewing conditions (2) use these images as feature vectors (3) apply a PCA over all the images (4) keep the dominant PCs (5) sequence of views for one object represent a manifold in space of projections (6) what is the nearest manifold for a given view?

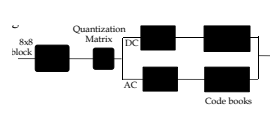
**Object-pose manifold** Appearance changes projected on PCs (1D pose changes). Sufficient characterization for recognition and pose estimation.

### 6.5 JPEG image compression

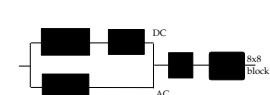
We don't resolve high frequencies too well...let's use this to compress images...JPEG!

**Concept** Block-based discrete cosine transform (DCT)

**JPEG Encoding and Decoding** Encoding:



Decoding:



#### DCT

A variant of discrete fourier transform: Real numbers, fast implementation. Block sizes: (1) small block: faster, correlation exists between neighboring pixels (2) better compression in smooth regions The first coefficient  $B(0,0)$  is the DC component, the average intensity. The top-left coefficients represent low frequencies, the bottom right high frequencies.

#### Entropy Coding (Huffman code)

symbol	prob.	code	binary fraction
Y	0.5	1	0.1
Z	0.25	01	0.01
X	0.125	001	0.001
W	0.125	000	0.000

The code words, if regarded as a binary fraction, are pointers to the particular interval being coded. In Huffman code, the code words point to the base of each interval. The average code length is  $H = -\sum p(s) \log_2 p(s) \rightarrow \text{optimal}$ .

## 7 Scale-space representations

From an original signal  $f(x)$  generate a parametric family of signals  $f^t(x)$ , where fine-scale information is successively suppressed.

### 7.1 Image pyramid

Level 0:  $1 \times 1$ , Level 1:  $2 \times 2$ , Level 2:  $4 \times 4$ , Level  $J-1$ :  $N/2 \times N/2$ , Level  $J$  (base):  $N \times N$ .

### 7.2 Applications

(1) Search for correspondence: look at coarse scales, then refine with finer scales (2) Edge tracking: a "good" edge at a fine scale has parents at a coarser scale (3) Control of detail and computational cost in matching: e.g. finding stripes; terribly important in texture representation

### 7.3 Pyramids

**Gaussian pyramid** Smooth with gaussians, because

"gaussian" = another gaussian. Progressively blurred and subsampled versions of the image. Adds scale invariance to fixed-size algorithms.

**Laplacian Pyramid** <+> Shown the information added in gaussian pyramid at each spatial scale. Useful for noise reduction & coding.

**Wavelet/QMRF** Bandpassed representation, complete, but with aliasing and some non-oriented subbands. Recursive application of a two-band filter bank to the low-pass band of the previous stage yields octave band splitting.

**Stereable pyramid** Shown components at each scale and orientation separately. Non-aliased subbands. Good for texture and feature analysis.

### 7.4 Haar Transform

Two major sub-operations: (1) Scaling captures info at different frequencies (2) Translation captures info at different locations Can be represented by filtering and downsampling. Relatively poor energy compaction.

### 8 Optical Flow

In Visual Computing, people seem like to use

$$\mathbf{I}_\bullet = \frac{\partial \mathbf{I}}{\partial \mathbf{x}}, \quad u = \frac{\partial x}{\partial t}, \quad v = \frac{\partial y}{\partial t}$$

where the subindex means a derivative if and only if we are talking about  $\mathbf{I}$ .

### 8.1 Applications

1 tracking 2 structure from motion 3 stabilization 4 compression 5 Mosaicing

### 8.2 Brightness constancy

**Definition of Optical Flow** "Apparent motion of brightness patterns". Ideally, the optical flow is the projection of the three-dimensional velocity vectors on the image.

**Caution required** 1 Uniform, rotating sphere  $\mathcal{OF} = 0$  2 No motion, but changing lighting  $\mathcal{OF} \neq 0$

### 8.3 Mathematical formulation

$$I(x, y, t) = \text{brightness at } (x, y) \text{ at time } t$$

$$I\left(\frac{dx}{dt} \delta t, y + \frac{dy}{dt} \delta t, t + \delta t\right) = I(x, y, t)$$

**Optical flow constraint equation:**

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

**8.4 The aperture problem**

The motion of an edge seen through an aperture is (in some cases) inherently ambiguous. E.g. the edge is physically moving upwards, but the edge motion alone is consistent with many other possible motions, and in this case the edge e.g. appears to move diagonally.

### 8.5 Optical Flow meaning

Estimate of observed projected motion field. Not always well defined! Compare: 1 Motion Field (or Scene Flow), projection of 3-D motion field 2 Normal Flow: observed tangent motion 3 Optic Flow: Apparent motion of the brightness pattern: Apparent motion of the brightness pattern (hopefully equal to motion field) 4 Consider barber pole illusion

**Planar motion** Ideal motions of a plane,  $X, Y$  being the horizontal and vertical direction and  $Z$  normal to the image plane: 1 translation in  $X$  2 translation in  $Z$  3 rotation around  $Z$  4 rotation around  $Y$

### 8.6 Regularization: Horn & Schunck algorithm

The Horn-Schunck algorithm assumes smoothness in the flow over the whole image, thus, it tries to minimize distortions in flow and prefers solutions which show more smoothness. The flow is formulated as a global energy functional which is the sought to be minimized. This function is given for two-dimensional image streams

$$\text{as Eq. 15: The associated ELE are}$$
$$\frac{\partial L}{\partial \bar{x}} - \frac{\partial}{\partial x} \frac{\partial L}{\partial \bar{x}} - \frac{\partial}{\partial y} \frac{\partial L}{\partial \bar{y}} = 0,$$
$$\frac{\partial L}{\partial \bar{y}} - \frac{\partial}{\partial x} \frac{\partial L}{\partial \bar{x}} - \frac{\partial}{\partial y} \frac{\partial L}{\partial \bar{y}} = 0.$$

$$\text{this gives}$$
$$\frac{\partial L}{\partial x} \left( \frac{\partial I}{\partial x} \bar{x} + \frac{\partial I}{\partial y} \bar{y} + \frac{\partial I}{\partial t} \right) - \alpha^2 \Delta \bar{x} = 0,$$
$$\frac{\partial L}{\partial y} \left( \frac{\partial I}{\partial x} \bar{x} + \frac{\partial I}{\partial y} \bar{y} + \frac{\partial I}{\partial t} \right) - \alpha^2 \Delta \bar{y} = 0.$$

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

**Remarks** 1 Coupled PDE solved using iterative methods and finite differences

$$\bar{x} = \Delta \bar{x} - \lambda \left( \frac{\partial L}{\partial x} \bar{x} + \frac{\partial L}{\partial y} \bar{y} + I \right) \frac{\partial L}{\partial \bar{x}},$$

2 More than two frames allow a better estimation of  $\mathbf{I}$ . 3 Information spreads from corner-type patterns. 4 Errors at boundaries 5 Example of regularization: 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  $(\dot{x}, \dot{y})$  must satisfy

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

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

$$\mathbf{A} \mathbf{v} = \mathbf{b}$$

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

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

Eq. 1 is overdetermined, so do compromise solution by the least squares principle Eq. 16

### 8.8 Gradient-Based Estimation

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

$$\sim f_1(x) - f_2(x) = \delta f_1'(x) + O(\delta^2)$$
$$\sim \delta = \frac{f_1(x) - f_2(x)}{f_1'(x)} \approx \delta$$

Assume displaced image well approximated by first-order Taylor series



consider generalizations to more interesting motion models.

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

$$\mathbf{u}(\mathbf{x}; \mathbf{x}_0) = A(\mathbf{x}; \mathbf{x}_0) \mathbf{c}, \quad (12)$$

where  $(c_1 \ c_2 \ c_3 \ c_4 \ c_5 \ c_6)^T$  are 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

$$\hat{\mathbf{c}} = M^{-1} \mathbf{b} \quad (13)$$

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

$$M = \sum_{\mathbf{x}} g A^T \nabla I^T \nabla I A,$$

$$\mathbf{b} = - \sum_{\mathbf{x}} g A^T \nabla I^T I_t.$$

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 folows 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 we use the same LS estimator as in Eq. 13, but with  $I$  and  $\hat{\mathbf{c}}$  replaced by  $I^j$  and  $^j\hat{\mathbf{c}}$ .

### 8.14 Low-order Parametric Deformations

There ary 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 centred 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 optiact flow estimation assumed smoothness 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. 15. 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 euations. 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 esti-

mates from different information sources. These issues can be addressed with a probabilistic formulation.

The cost function 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

$$p(I(\mathbf{x}, t) | \mathbf{u}) =_{\propto} e^{E(\mathbf{u})/2\sigma^2}.$$

#### 8.17 Parametric motion models

Globa miton models offer:

- ➊ More constrained solutions than smoothness (Horn-Schunck)
- ➋ 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 < 24 Hz. Asuccession of images will be perceived as continuous if frequency is sufficiently high. Cinema 2424 Hz, TV 25 Hz or 50 Hz. We still nee to avoid aliasing (wheel effect). High-rendering frame-rates 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 temporarily shifted half images, increase of frequency 25 Hz → 50 Hz. Reduction of spatial resolution. Full image representation: progressive.

#### 9.3 Why compress video?

Raw HD TV signal 720*p* @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:

- ➊ 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.
- ➋ Productive 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:**

- ➊ I-frame: Intra-coded frame, coded independently of all other frames.
- ➋ P-frame: Predictively coded frame, coded based on previously coded frame I or P. Can send motion vector plus changes.

- ➌ 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**

- ➊ Partition video into moving objects
- ➋ describe object motion → Generally very difficult

**Practical approach** Block-Matching Motion Estimation:

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

#### 9.5 Block-matching motion estimation

**Assumptions:** ➊ Translational motion within block:

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

**ME Algorithm** ➊ Divide current frame into non-overlapping  $N_1 \times N_2$  blocks.

- ➋ 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.

**Metrics** for determining “best match”:

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

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

**Candidate blocks** All blocks in, e.g. (±32, ±32) pixel area

**Strategies for searching** candidate blocks for best match.

- ➊ Full search: Examine all candidate blocks
- ➋ Partial (fast) search: Examine a carefully selected subset.

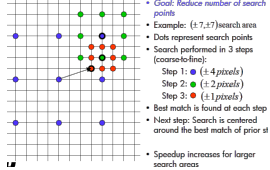
**Motion vector** Estimate of motion for best matching block.

#### 9.6 Motion vector and motion vector field

**Motion vector** Expresses the *relative horizontal and vertical offsets*  $(mv_1, mv_2)$ , or motion, of a given block from one frame to another.

**Motion vector field** Collection of motion vectors for all the blocks in a frame.

**Example of fast motion estimation search**



**Motion Vector Presision**

- ➊ Motivation: Motion is not limited to integer-pixel offsets. However, video is only known at discrete pixel locations. To estimate sub-pixel motion, frames must be spatially interpolated.
- ➋ Fractional MVs are used to represent the sub-pixel motion.
- ➌ Improved performance (extra complexity is worthwhile)
- ➍ Half-pixel ME used in most standards: MPEG-1/2/4
- ➎ Why are half-pixel motion vectors better? They can capture half-pixel motion. Averaging effect (from spatial interpolation) reduces prediction error → Improved prediction.

MC-prediction usually performs well; In compression have a second chance to recover when it performs badly.

MC-prediction yields

- ➊ Motion vectors
- ➋ MC-prediction error or residual → Code error with conventional image coder
- ➌ Sometimes MC-prediction may *perform badly*
- ➍ Examples: complex motion, new imagery (occlusions)
- ➎ Approach: 1. Identify frame or individual blocks where prediction fails

2. Code without prediction

### 9.10 Basic Video Compression Architecture

- ➊ Exploiting the redundancies:

- ➋ Temporal: MC-prediction (P and B frames)
- ➌ Spatial: Block DCT
- ➍ Color: color space conversion
- ➎ Scalar quantization of DCT coefficients
- ➏ Zigzag scanning, runlength and Huffman coding of the nonzero quantized DCT coefficients

**Disadvantages**

- ➊ Assumes translational motion model → Breaks down for more complex motion.
- ➋ Offer produces blocking artifacts (OK for coding with Block DCT)

**Bidirectional MC prediction** is used to estimate a block in the current frame from a block in:

- ➊ Previous frame
- ➋ Future frame
- ➌ Average of ablock from the previous frame and a block from the future frame
- ➍ Neither, i.e. code current block without prediction

*Example:* Prediction with P- and B-frames

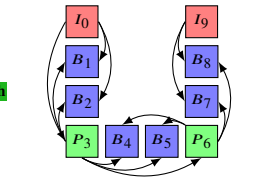
- ➊ Motion compensated prediction: Predict the current frame based on reference frame(s) while compensating for the motion.
- ➋ Examples of block-based motion-compensated prediction (P-frame) and bi-directional prediction (B-frame).

#### 9.8 Frame types

Main addition over image compression: Exploit the temporal redundancy. Predict current frame based on previously coded frames. Three types of coded frames:

- ➊ *I-frame*: Intra-coded frame, coded independently of all other frames
- ➋ *P-frame*: Predictively coded frame, coded based on previously coded frame
- ➌ *B-frame*: Bi-directionally predicted frame, coded based on both previous and future coded frames.

**MPEG: Group of Pictures (GOP)**



Starts with an I-frame, ends with frame right before next I-frame. “Open” ends in B-frame, “closed” in P-frame. MPEG Encoding a parameter, but “typical”:

**IBBPBBPBBBI, IBBPBBPBBPBBBI.**

Periodic I-frames enable random access into the coded bitstream. Parameters:

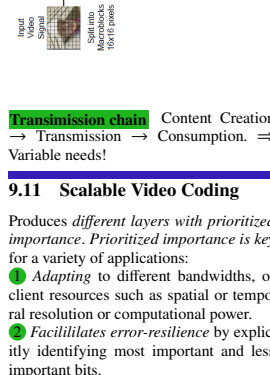
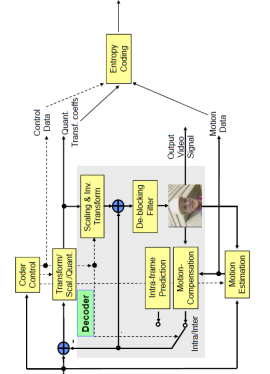
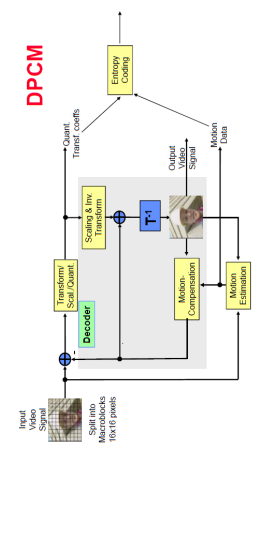
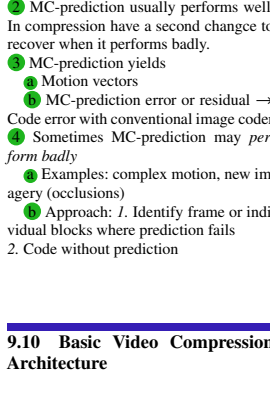
- ➊ Spacing between I frames
- ➋ Number of B frames between I and P frames

**Example compression performance**

I:  $\frac{1}{7}$ , P:  $\frac{1}{20}$ , B:  $\frac{1}{50}$ , Average:  $\frac{1}{27}$ .

#### 9.9 Summary of Temporal Processing

- ➊ Use MC-prediction (P and P frames) to reduce temporal redundancy.



- ➊ Decompose video into *multiple layers of prioritized importance*
- ➋ Code layers into *base and enhancement* bitstreams
- ➌ Progressively combine *one or more bitstreams* to produce *different levels of video quality*.

Example of scalable coding with base and two enhancement layers: Can produce three different qualities:
➊ Base layer
➋ Base + Enh1 layers
➌ Base + Enh1 + Enh2 layers
Scalability with respect to: Spatial or temporal resolution, bit rate, computation, memory.

- ➍ **Example**
  - Encode image/video into three layers: Base, Enh1, Enh2
  - Low-bandwidth receiver: Send only Base layer.
  - Medium-bandwidth receiver: Send Base & Enh1 layers
  - High-bandwidth receiver: Send all three layers Base, Enh1, Enh2.
- Can adapt to different clients and network situations
- Three basic types of scalability (refine video quality along three different dimensions):

- ➊ Temporal scalability → Temporal resolution
- ➋ Spatial scalability → Spatial resolution
- ➌ SNR (quality) scalability → Amplitude resolution
- Each type of scalable coding provides scalability of one dimension of the video signal. Can combine multiple types of

scalability to provide scalability along multiple dimensions

**Temporal Scalability** based on the use of *B-frames* to refine the *temporal resolution* B-frames are dependent on other frames. However, no other frame depends on a B-frame. Each B-frame may be discarded without affecting other frames.

**Spatial scalability** Based on refining the *spatial resolution* *Base layer* is *low resolution* version of video. *Enh1* contains *coded difference* between upsampled base layer and original video. Also called pyramid coding.

**SNR scalability** Based on refining the *amplitude resolution* *Base layer* uses a *coarse quantizer*. *Enh1* applies a *finer quantizer* to the difference between the original DCT coefficients and the coarsely quantized base layer coefficients.

#### 9.12 Standards

**Goal** Ensuring *interoperability*: Enabling communication between devices made by different manufacturers. Promoting a technology or industry. Reducing costs.

**Scope of standardization** Not the encoder, not the decoder. Just the *bitstream syntax* and the *decoding process* (e.g. use IDCT but not how to implement the IDCT) This enables improved encoding and decoding strategies to be employed in a standard-compatible manner.

#### 9.13 Quality measure

**objective: PSNR**

- Error for one pixel, difference between original and decoded value

$$e(v, h) = \tilde{x}(v, h) - x(v, h)$$

- Mean-squared-Error, MSE e.g. over an image

$$e\text{MSE} = \sqrt{\frac{1}{N \cdot M} \sum_{v,h=1}^{v=N, h=M} e^2(v, h)}$$

- Peak-Signal-to-Noise-Ratio

$$\text{PSNR} = [\max x]^2 / e^2_{\text{MSE}}$$

E.g.  $x = 2^K$  or 255. One can use a log-scale like dB.

#### A Big equations

$$\begin{aligned} \mathcal{F}[h](u, v) &= \frac{1}{2\ell} \int_{-\ell}^{\ell} dx_1 \exp(-i2\pi u x_1) \cdot \underbrace{\int_{-\infty}^{\infty} dx_2 \delta(x_2) \exp(-i2\pi v x_2)}_{=1} \\ &= \text{sinc}(2\pi u \ell) \end{aligned} \quad (14)$$

$$E = \iint dx dy \left[ \left( \frac{\partial I}{\partial x} \frac{dx}{dt} + \frac{\partial I}{\partial y} \frac{dy}{dt} + \frac{\partial I}{\partial t} \right)^2 + \alpha^2 (\|\nabla \hat{x}\| + \|\nabla \hat{y}\|)^2 \right] \quad (15)$$

$$\mathbf{v} = \left( \frac{\sum_i w_i I_x(q_i)^2}{\sum_i w_i I_x(q_i) I_y(q_i)} \frac{\sum_i w_i I_x(q_i) I_y(q_i)}{\sum_i w_i I_y(q_i)^2} \right)^{-1} \cdot \begin{pmatrix} -\sum_i w_i I_x(q_i) I_t(q_i) \\ -\sum_i w_i I_y(q_i) I_t(q_i) \end{pmatrix} \quad (16)$$

#### 10 Questions

- $I_{\text{comp}} = I_a I_a + (1 - I_a) I_b$
- MAP, Maximum a posteriori detector.
- graph cuts
- Solve MRFs with graph cuts
- impulse response  $t(-x, -y)$
- 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