

Contents

Contents

1	The digital image	9.4	Lossy video compression
1.1	Image as 2D signal	9.5	Block-matching motion estimation
1.2	Image sources	9.5.1	Determining the best matching block
1.3	Sampling		
1.4	Reconstruction		
1.5	Quantization		
1.6	Image Properties		
1.7	Image Noise		
1.8	Colour Images		

2	Image segmentation	9.6	Motion vector and motion vector field
2.1	Thresholding	9.6.1	Practical Half-Pixel
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	9.7	Block Matching Algorithm
3.1	Linear Shift-Invariant Filtering	9.8	Frame types
		9.9	Summary of Temporal Processing

4	Image features	9.10	Basic Video Compression Architecture
4.1	Template matching	9.11	Scalable Video Coding
4.2	Edge detection	9.12	Standards
4.2.1	Canny edge detector	9.13	Quality measure

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		

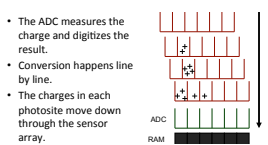
5	Fourier Transform	5.1	Aliasing
		5.2	Definition
		5.3	Sampling
		5.4	Image Restoration

6	Unitary transforms	6.1	Karhunen-Loeve Transform
		6.2	Basis images and eigen-images (EI)
		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?



Blooming Buckets have finite capacity. Photoisate 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

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.

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

crappily 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 $p(c) = (2\pi\sigma^2)^{-1}e^{-c^2/2\sigma^2}$.

Poisson noise (shot noise)

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

Rician noise (appears in MRI)

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

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. → high-end cameras

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

Filter wheel rotate multiple filters in front of lens. Allows more than 3 colour bands. → static scenes

new color CMOS sensor, foveon's X3 blue, green, red sensor, one above the other (descending) → 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

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

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 grayscale 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 $\sigma \rightarrow I_{\Sigma}$. $I_a = |I - I_{bg}| >$

$$T, T = \begin{pmatrix} 20 & 20 & 10 \end{pmatrix}, I_{bg} = \text{background image. Or better (e.g.)}$$

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

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*

$$TP \text{ fraction} = \text{true positive count} / P,$$

$$P = TP + FN \text{ and false positive fraction}$$

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

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 $V_{TN} = V_{TP} = 0$.

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

$$I \odot S = (I \odot S) \odot S.$$

To remove holes in the foreground and islands in the background, do both opening and closing. Theseize 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".

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 dependend on others, no 1D structure to dependencies. • This means no optimal, efficient algorithms, except for 2-label problems. Minimize

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

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.

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

Erosion of binary image I by the structuring element S is defined by

$$I \odot S = \{z \in E \mid S_z \subset I\}$$

S_z translation of S by vector z .

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

Opening $I \odot S = (I \odot S) \oplus S$.

$$\text{Closing } I \bullet S = (I \oplus S) \odot S.$$

To remove holes in the foreground and islands in the background, do both opening and closing. Theseize 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 = granulo(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. for x=1:maxRad O = imopen(B, strel('disk',x)); numRegions(x) = max(max(connectedComponents(O))); end gSpec = diff(numRegions);

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

$$\text{Thinning } I \odot S = I \setminus (I \odot S)$$

$$\text{Thickening } I \odot S = I \cup (I \odot S)$$

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, theskeleton 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 $S_n(X) = (X \ominus B) \setminus [(X \ominus B) \odot B]$ 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 provies a stick figure representing the region shap. Used in object recognition, in particula, 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.

Diagram of template matcher

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. k is *kernel* of the operation. $N(m, n)$ is a neighbourhood of (m, n) .

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

Convolution e.g. point spread function $I'(x, y) = \sum_{i,j \in N(x,y)} K(i, j) I(x + i, y + j)$

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 $(2M + 1) \times (2N + 1)$,

$$I'(m, n) = f * (g * I(N(m, n))),$$

$$I''(m, n) = \sum_{i=-M}^N g(j) I(m, n - j),$$

$$I'(m, n) = \sum_{i=-M}^M f(i) I''(m - i, n) \rightarrow (2M + 1) + (2N + 1) \text{ operations!}$$

Smoothing kernels (low-pass filters)

$$\bullet \text{ Mean filter: } \frac{1}{9} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \bullet \text{ Weighted smoothing filters: } \frac{1}{10} \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{pmatrix},$$

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 vertical and horizontal lines. Gaussian smoothing Kernel is *separable*! $N(x, y) = N(x)N(y)$. Amount of smoothing depends on σ and window size. Width > 3σ .

Scale space Convolution of a Gaussian with σ 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. → Neighbor's influence decreases monotonically. (3) Still one lobe in frequency domain. → No corruption from high frequencies. (4) Simple relationship to σ

Diagram of template matcher

(5) Easy to implement efficiently

Differential filters • Prewitt operator: $\begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}$, • Sobel operator: $\begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}$.

<

$r(x, y) \rightarrow r(-x, -y)$
 search object
 peak(s) \rightarrow location(s) (p, q) Remove
 mean before template matching to avoid
 bias towards bright image areas.

4.2 Edge detection

Idea (continuous-space): Detect local gradient

$\| \nabla f(x, y) \| = \sqrt{(\partial_x f)^2 + (\partial_y f)^2}$
 Digital image: Use finite differences instead:

difference $(-1 \ 1)$, **central difference** $(-1 \ 0 \ 1)$: **Prewitt** $\begin{pmatrix} -1 & 0 & 1 \\ 0 & [0] & 0 \\ -1 & 0 & 1 \end{pmatrix}$, **Sobel** $\begin{pmatrix} -1 & 0 & 1 \\ 0 & [0] & 0 \\ -2 & [0] & 2 \end{pmatrix}$, **Roberts** $\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$, $\begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$.

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

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)

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

4.2.1 Canny edge detector

(1) Smooth image with a gaussian filter (2) Compute gradient magnitude and angle (Sobel, Prewitt,...)

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

(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° , vertical, 45° . 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.

4.3 Feature detection

4.3.1 Hough transform

Problem: fit a straight line (or curve) to a set of edge pixels. Hough transform (1962): generalized template matching technique. (1) Consider detection of straight lines $y = mx + c$. (2) draw a line in the parameter space m, c for each edge pixel x, y and increment bin counts along line. Detect peak(s) in (m, c) plane. (3) Alternative parametrization avoids infinite-slope problem $x \cos \theta + y \sin \theta = \rho$

circle detection find circles of fixed radius r . For circles of undetermined radius,

use 3D Hough transform for parameters (x_0, y_0, r)

4.3.2 Detecting corner points

Many applications benefit from features localized in (x, y) . Edges well localized only in one direction \rightarrow detect corners. Desirable properties of corner detector: (1) Accurate localization, (2) invariance against shift, rotation, scale, brightness change, (3) robust against noise, high repeatability

4.3.3 Most accurately localizable patterns

Local displacement sensitivity

$S(\Delta x, \Delta y) = \sum_{x, y \in \text{window}} [f(x, y) - f(x - \Delta x, y - \Delta y)]^2$

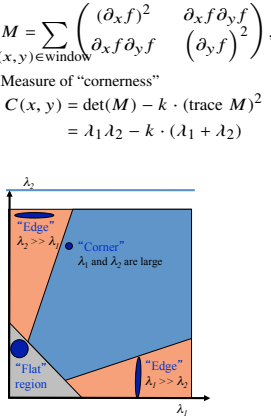
Linear approximation for small $\Delta x, \Delta y$

$f(x + \Delta x, y + \Delta y) \approx f(x, y) + \partial_x f(x, y) \Delta x + \partial_y f(x, y) \Delta y$
 $S(\Delta x, \Delta y) \approx \sum_{(x, y) \in \text{window}} \begin{pmatrix} \partial_x f & \partial_y f \end{pmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} \mathbf{M} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix}$

Feature point extraction

$\text{SSD} \approx \Delta^T \mathbf{M} \Delta$
 Find points for which the following is large
 $\min \Delta^T \mathbf{M} \Delta$
 for $\| \Delta \| = 1$, i.e. maximize eigenvalues of \mathbf{M} .

Keypoint detection Often based on eigenvalues λ_1, λ_2 of \mathbf{M} ("structure matrix"/"normal matrix"/"second-moment matrix")
 $\mathbf{M} = \sum_{(x, y) \in \text{window}} \begin{pmatrix} (\partial_x f)^2 & \partial_x f \partial_y f \\ \partial_x f \partial_y f & (\partial_y f)^2 \end{pmatrix}$,
 Measure of "cornerness"
 $C(x, y) = \det(\mathbf{M}) - k \cdot (\text{trace } \mathbf{M})^2 = \lambda_1 \lambda_2 - k \cdot (\lambda_1 + \lambda_2)^2$



(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

Corner importance weight Give more importance to central pixels by using Gaussian weighting function
 $M = \sum_{x, y \in \text{window}} G(x - x_0, y - y_0, \sigma)$

$\hat{f} \cdot \hat{g} = \widehat{f \cdot g}$
 $\begin{pmatrix} (\partial_x f)^2 & \partial_x f \partial_y f \\ \partial_x f \partial_y f & (\partial_y f)^2 \end{pmatrix}$

Compute subpixel localization by fitting parabola to *cornerness function*

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 $(x, y, \text{scale}, \text{orientation})$

SIFT description (1) Thresholded image gradients are sampled over 16×16 array of locations in scale space. (2) Create array of orientation histograms (3) 8 orientations $\times \frac{1}{4} \times 4$ histogram array = 128 dimensions

5 Fourier Transform

5.1 Aliasing

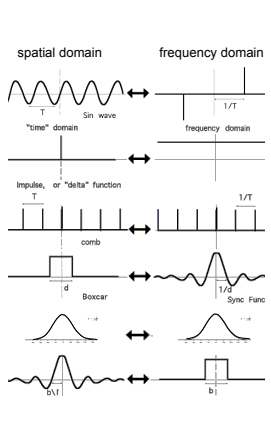
One can't shrink an image by taking every second pixel. If we do, characteristic errors appear. Typically, small phenomena look bigger; fast phenomena can look slower. Common phenomena (1) Wagon wheels rolling the wrong way in movies. (2) Checkboards misrepresented in ray tracing (3) Striped shirts look 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 $\hat{f}(u, v) = \iint_{\mathbb{R}^2} dx dy f(x, y) e^{-i2\pi(u x + v y)}$

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
 $\hat{f} \cdot \hat{g} = \widehat{f \cdot g}$

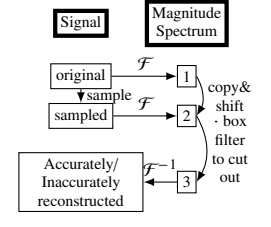
The Fourier transform of the product of two functions is the convolution of the Fourier transforms
 $\hat{f} \cdot \hat{g} = \widehat{f \cdot g}$

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 \rightarrow Delta function
 $S_{2D}(f(x, y))$

$\hat{f} \cdot \hat{g} = \widehat{f \cdot g}$
 $\hat{f} \cdot \hat{g} = \widehat{f \cdot g}$
 with $\mathcal{S} =$ 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 $\xrightarrow{\text{lp filtering}}$ lp filt. sign. $\xrightarrow{\text{lp filt. sign.}}$ reconstr. sign.

Smoothing as low-pass filtering The message of the FT is that high frequencies lead to trouble with sampling. Solutionspress 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.

Nyquist sampling theorem Nyquist theorem: The sampling frequency must be at least twice the highest frequency. $\omega_s \geq 2\omega_c$. If this is no the case, the signal needs to be bandlimited before sampling, e.g. with a low-pass filter.

5.4 Image Restoration

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:
 $h(x_1, x_2) = \frac{1}{2\ell} [\theta(x_1 + \ell) - \theta(x_1 - \ell)] \delta(x_2)$

Noise Gaussian blurring kernel:
 $h(x_1, x_2) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{x_1^2 + x_2^2}{2\sigma^2}\right)$

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 defocussing
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

$\hat{F}[\hat{h}](u, v) = \frac{\mathcal{F}[h]}{\|\mathcal{F}\|^2 + \epsilon}$
 Singularities are avoided by the rugelarization ϵ .

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 (ediseder to avoid aliasing) • Low-res images = $D * H * G *$ (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 ration $\sim \sqrt{n}$, 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 $\mathbf{c} = \mathbf{A} \mathbf{f}$ 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. $\mathbf{A} = \bar{\mathbf{A}}$, transform is orthonormal.

Energy conservation $\| \mathbf{c} \|^2 = \mathbf{c}^* \mathbf{c} = \mathbf{f}^* \mathbf{A}^* \mathbf{A} \mathbf{f} = \| \mathbf{f} \|^2$
Image collection f_i one image, $F = [f_1, \dots, f_n]$.

Auto-correlation function

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

energy distribution Energy is conserved, but often will be unevenly distributed among coefficients. Autocorrelation matrix:
 $R_{cc} = E[\mathbf{c} \mathbf{c}^*] = E[\mathbf{A} \mathbf{f} \mathbf{f}^* \mathbf{A}^*] = \mathbf{A} R_{ff} \mathbf{A}^*$
 Mean squared values ("average energies") of the coefficients c_i are on the diagonal of R_{cc} .

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

Eigenmatrix of autocorrelation matrix Definition: Eigenmatrix Φ of autocorrelation matrix R_{ff} . (1) ϕ is unitary (2) The columns of Φ form a set of eigenvectors of R_{ff} , i.e.,
 $R_{ff} \Phi = \Phi \Lambda$,

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

6.1 Karhunen-Loeve Transform

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

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

(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 contains ones on the first J diagonal positions, else zeros.

(2) Energy in first J coefficients for arbitrary transform \mathbf{A}
 $E = \text{tr}(R_{\mathbf{b}\mathbf{b}}) = \text{tr}(\mathbf{I}_J R_{cc} \mathbf{I}_J)$

$= \text{tr}(\mathbf{I}_J \mathbf{A} R_{ff} \mathbf{A}^* \mathbf{I}_J) = \sum_{k=0}^{J-1} a_k^T R_{ff} \bar{a}_k$
 where a_k^T is the k -th row of \mathbf{A} . (3) Lagrangian cost function to enforce unit-length basis vectors

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

Differentiating L with respect to a_j yields necessary condition
 $R_{ff} \bar{a}_j = \lambda_j \bar{a}_j, \quad \forall j < J$

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 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 J 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

$\arg \min_i D_i = \| \mathbf{I}_i - \mathbf{I} \|$

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 $\hat{J} = J - \langle J \rangle$, with J being any image (set). Do KLT (aka PCA) transformation

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

with closest rank-k approximation property of SVD. Approximate

$\arg \min_i D_i = \| \mathbf{I}_i - \mathbf{I} \| \approx \| p_i - p \|$

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 ar maximized.

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

$E_{\text{opt}} = \arg \max_E \left(\det \left(E R E^* \right) \right)$,

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

$F_{\text{opt}} = \arg \max_F \left(\frac{\det \left(F R_B W^* \right)}{\det \left(F R_W F^* \right)} \right)$,

$R_B = \sum_{i=1}^c N_i \left(\mathbf{\mu}_i - \bar{\mathbf{\mu}} \right) \left(\mathbf{\mu}_i - \bar{\mathbf{\mu}} \right)^*$,

$R_W = \sum_{i=1, \dots, c} \left(\Gamma_{\ell} - \mathbf{\mu}_i \right) \left(\Gamma_{\ell} - \mathbf{\mu}_i \right)^*$,
 $\Gamma_{\ell} \in \text{Class}(i)$

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

$R_B \mathbf{w}_i = \lambda_i R_W \mathbf{w}_i, i = 1, \dots, k$

Problem: within-class scatter matrix 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) ile in a 3D linear subspace. Single point source at infinity

$f(x, y) = a(x, y) \left(\ell^T n(x, y) \right) L$,

$a(x, y)$ surface albedo, L light source intensity. Superposition of arbitrary number 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 withi-class scatter.

Appearance manifold approach

For everyobject, (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?

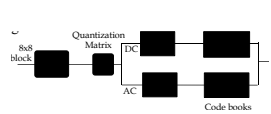
Object-pose manifold Appearance changes projected on PCs (1D pose changes). Sufficient characterization for recognition and pose estimation.

6.5 JPEG image compression

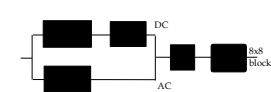
We don't resolve high frequencies too well...let's use this to compress images ...JPEG!

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
Z	0.5	1	0.1
Y	0.25	01	0.01
X	0.125	001	0.001
W	0.125	000	0.000

● Coupled PDE solved using iterative methods and finite differences

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

$$\dot{\mathbf{y}} = \Delta \dot{\mathbf{y}} - \lambda \left(\frac{\partial I}{\partial \mathbf{x}} \dot{\mathbf{x}} + \frac{\partial I}{\partial \mathbf{y}} \dot{\mathbf{y}} + \dot{\mathbf{I}} \right) \frac{\partial I}{\partial \mathbf{y}}.$$

● More than two frames allow a better estimation of $\dot{\mathbf{I}}$. ● Information spreads from corner-type patterns.

● Errors at boundaries ● 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 \mathbf{p} under consideration, thus the optical flow equation can be assumed to hold for all pixels within a window centered at \mathbf{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

$$\mathbf{A} \mathbf{v} = \mathbf{b} \quad (1)$$

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

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

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

8.8 Gradient-Based Estimation

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

$$\begin{aligned} &\leadsto f_1(\mathbf{x}) - f_2(\mathbf{x}) \\ &= \delta f_1'(\mathbf{x}) + \mathcal{O}\left(\delta^2\right) \\ &\leadsto \tilde{\delta} = \frac{f_1(\mathbf{x}) - f_2(\mathbf{x})}{f_1'(\mathbf{x})} \approx \delta \end{aligned} \quad (2)$$

Assume displaced image well approximated by first-order Taylor series

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

Insert Eq. 2 in Eq. 3 to get

$$\nabla I(\mathbf{x}, t) \cdot \mathbf{u} + I_t(\mathbf{x}, t) = 0. \quad (4)$$

This i s 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:

$$\begin{aligned} &I(\mathbf{x}(t), t) = c, \\ &\leadsto \frac{d}{dt} I(\mathbf{x}(t), t) = 0 \\ &\leadsto \nabla I \cdot \dot{\mathbf{u}} + I_t = 0. \end{aligned} \quad (5)$$

\leadsto gradient-constraint equation Eq. 4.

● **Least-Squares estimation** One cannot recover \mathbf{u} from one gradient constraint since Eq. 4 is one equation with two unknowns, u_1 and u_2 . The intensity gradient instead constrains the flow to a one-parameter family of voelcities along a line in *velocity space*. One can see from Eq. 4 that this line is perpendicular to ∇I and its perpendicular distance from the origin is

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 velocity that minimizes the constraint errors. The least-squares (LS) estimator mini-

mizes the squared errors:

$$E(\mathbf{u}) = \sum_{\mathbf{x}} g(\mathbf{x}) [\mathbf{u} \cdot \nabla I(\mathbf{x}, t) + I_t(\mathbf{x}, t)]^2, \quad (6)$$

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.,

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

These equations may be rewritten in matrix form:

$$\mathbf{M} \mathbf{u} = \mathbf{b} \quad (7)$$

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

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

● **Implementation Issues** Usually we wish to estimate optical flow at every pixel, so we should express \mathbf{M} and \mathbf{b} as functions of position \mathbf{x} , i.e., $\mathbf{M}(\mathbf{x})\mathbf{u}(\mathbf{x})$. Note that the elements of \mathbf{M} and \mathbf{b} are local sums of products of image derivatives. An effective way to estimate the flow field ist 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. 7. These quadratic images are then convolved with $g(\mathbf{x})$, to obtain the elemnts of $\mathbf{M}(\mathbf{x})$ 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 \mathbf{M} in Eq. 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 \mathbf{M} will be singular. When each image gradient within a region has the same spatial direction, it is easy to see that rank $\mathbf{M} = 1$. 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 7 provides an optimal solution, but not to our original problem. Higher-order terms ignored!

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

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 $I^0(\mathbf{x}, t)$:

$$I^0(\mathbf{x}, t + \delta t) = I\left(\mathbf{x} + \mathbf{u}^0 \delta t, t + \delta t\right), \quad (8)$$

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

$$\delta \hat{\mathbf{u}} = \mathbf{M}^{-1} \mathbf{b} \quad (9)$$

where \mathbf{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

$$\begin{aligned} E(\delta \mathbf{u}) &= \sum_{\mathbf{x}} g(\mathbf{x}) \left[I(\mathbf{x}, t) \right. \\ &\quad \left. - I\left(\mathbf{x} + \hat{\mathbf{u}} + \delta \mathbf{u}, t + 1\right) \right]^2 \\ &= \sum_{\mathbf{x}} g(\mathbf{x}) \left[I^j(\mathbf{x}, t) \right. \\ &\quad \left. - I^j\left(\mathbf{x} + \delta \mathbf{u}^j, t + 1\right) \right]^2 \\ &=: \tilde{E}(\delta \mathbf{u}). \end{aligned} \quad (10)$$

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 9. One can, however, formulatio the problem so that the spatial image derivatives used to form \mathbf{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 9 will always depend on the warped image sequence and must be recomputed at each iteration. In practice, when \mathbf{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 i s called the *gradient constraint equation*.

8.11 Pyramid/Coarse-to-fine

Limits of the (local) gradient method:

- Fails when intensity structure within window is poor
- Fails when displacement is large (typical operating range is motion of 1 pixel per iteration!). Linearization of brightness is suitable only for small displacements.
- 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

$$\begin{aligned} F(\omega_x, \omega_y, \omega_t) &= F_0(\omega_x, \omega_y) \\ &\cdot \delta\left(u_1 \omega_x + u_2 \omega_y + \omega_t\right), \end{aligned} \quad (11)$$

where F_0 is the 2D Fourier transform of f_0 . Eq. 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 filtersl 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 sub-sampling). 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 scal'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 prevents 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 gra-dient 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 will be extremely poor, especially at occlusion boundaries.

LS estimators are not suitable when the distribution of gradient constraint errors is heavy-tails, as they are sensitive to small numbers of measurement outliers. It is therefore often crucial that the quadratic estimation in Eq. 6 be replaced by a robust estimator, $\rho(\cdot)$, which limits the influence of constraints with larger errors:

$$\begin{aligned} E(\mathbf{u}) &= \sum_{\mathbf{x}} g(\mathbf{x}) \rho(e(\mathbf{x}), \sigma) \\ \rho(e, \sigma) &= e^2 / \left(e^2 + \sigma^2 \right), \end{aligned}$$

where σ^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

$$\begin{aligned} \mathbf{u}(\mathbf{x}; \mathbf{x}_0) &= A(\mathbf{x}; \mathbf{x}_0) \mathbf{c}, \quad (12) \\ \text{where } &(\epsilon_1 \ 2 \ \epsilon_3 \ \epsilon_4 \ \epsilon_5 \ \epsilon_6)^T \text{ are the motion model parameters, and} \\ A(\mathbf{x}; \mathbf{x}_0) &= \begin{pmatrix} 1 & 0 & x-x_0 & y-y_0 & 0 & 0 \\ 0 & 1 & 0 & 0 & x-x_0 & y-y_0 \end{pmatrix}. \end{aligned}$$

From Eq. and we get the gradient constraint equation

$$\begin{aligned} \nabla I(\mathbf{x}, t) A(\mathbf{x}; \mathbf{x}_0) \mathbf{c} + I_t(\mathbf{x}, t) &= 0, \\ \text{for which the LS estimate for the neigh-} &\text{bourhood has the form} \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{c}} &= \mathbf{M}^{-1} \mathbf{b} \quad (13) \\ \text{where now } \mathbf{M} \text{ and } \mathbf{b} \text{ are given by} \\ \mathbf{M} &= \sum_{\mathbf{x}} g A^T \nabla I^T \nabla I A, \\ \mathbf{b} &= - \sum_{\mathbf{x}} g A^T \nabla I^T I_t. \end{aligned}$$

When \mathbf{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

$$\begin{aligned} \delta \mathbf{u} &= A \delta \mathbf{c}. \\ \text{Accordingly, defining } I^j(\mathbf{x}, t) &\text{ to be the original sequence } I(\mathbf{x}, t) \text{ warped by } \mathbf{u}^j \text{ as in Eq. 8 we use the same LS estimator as in Eq. 13, but with } I \text{ and } \hat{\mathbf{c}} \text{ replaced by } I^j \text{ and } \hat{\delta \mathbf{c}}. \end{aligned}$$

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 θ , 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 $(d_1, d_2, s \cos \theta, s \sin \theta)^T$ 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 optiocl 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* λ 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 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 η is white Gaussian noise with standard deviation σ , and uncorrelated at different pixels, we obtain the conditional density

$$\begin{aligned} p(I \mid \mathbf{u}) &= \propto e^{E(\mathbf{u})/2\sigma^2} \\ &= \end{aligned}$$

8.17 Parametric motion models

Global miton models offer:

- More constrained solutions than smoothness (Horn-Schunck)
- 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. Visual perception is limited to < 24Hz. Asusession of images will be perceived as continuous if frequency is sufficiently high. Cinema 2424 Hz, TV 25 Hz or 50 Hz. We still need 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 100Hz 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?

$$\begin{aligned} \text{Raw HD TV signal } 720p @ 50 \text{ Hz:} \\ 1280 \cdot 720 \cdot 50 \cdot 24 \text{ bits/s} \\ = 1\,105\,920\,000 \text{ bits/s} > 1 \text{ Gb/s} \end{aligned}$$

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:

- 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.
- 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:**

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

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

9.5 Block-matching motion estimation

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

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

9.5 Block-matching motion estimation

● **Assumptions:** ● Translational motion within block:

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

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

- 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”:

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

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

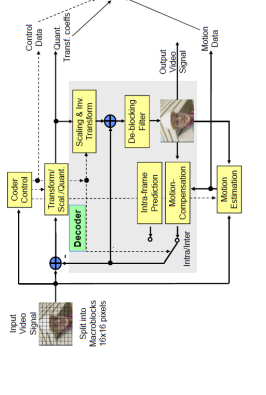
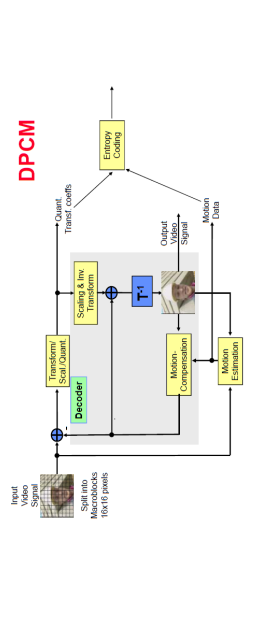
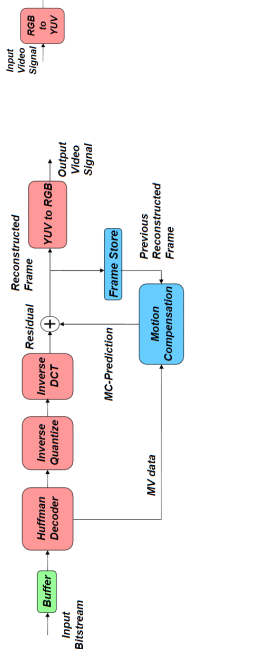
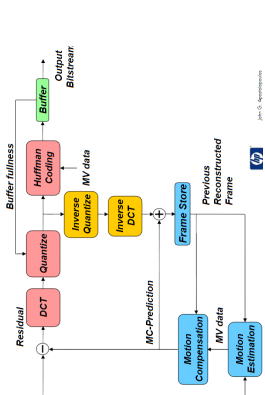
● **Strategies for searching** candidate blocks for best match.

- Full search: Examine all candidate blocks
- 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

- 2 Scalar quantization of DCT coefficients
- 3 Zigzag scanning, runlength and Huffman coding of the nonzero quantized DCT coefficients



Transmission chain Content Creation
→ Transmission → Consumption. ⇒
Variable needs!

9.11 Scalable Video Coding

Produces *different layers with prioritized importance*. *Prioritized importance is key* for a variety of applications:

- 1 *Adapting* to different bandwidths, or client resources such as spatial or temporal resolution or computational power.
- 2 *Facilitates error-resilience* by explicitly identifying most important and less important bits.

Procedure:

- 1 Decompose video into *multiple layers of prioritized importance*
- 2 Code layers into *base and enhancement* bitstreams
- 3 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: 1 Base layer 2 Base + Enh1 layers 3 Base + Enh1 + Enh2 layers Scalability with respect to: Spatial or temporal resolution, bit rate, computation, memory.

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):

- 1 Temporal scalability → Temporal resolution
- 2 Spatial scalability → Spatial resolution
- 3 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
 $e(v, h) = \hat{x}(v, h) - x(v, h)$
• 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

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

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

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

A Big equations

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

$$= \text{sinc}(2\pi u \ell)$$

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

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