
Compact Lecture

Multimedia Coding: Methods & Applications

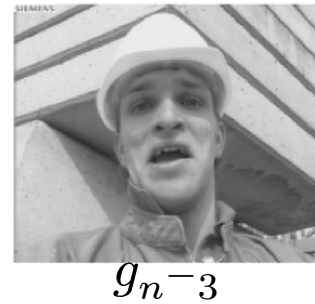
Part 4: Video Coding Fundamentals

4.1: Motion Estimation and Compensation

Dr. Klaus Illgner

Dr. Uwe Rauschenbach

What is „Video“?



**Video is sequence of images $\{g\}$,
where the images have a ordered relationship in time**

Key Feature:

Difference between images mainly
caused by motion



What can be done for Efficient Coding?

Resolution of standard TV:

720 x 576, 25 Hz, 4:2:0 → 165,9 Mbps (90 min → 112 GB)
still image compression 10:1 → 16,6 Mbps (90 min → 11,2 GB)

→ amount of data even for SDTV too large for transmission and storage

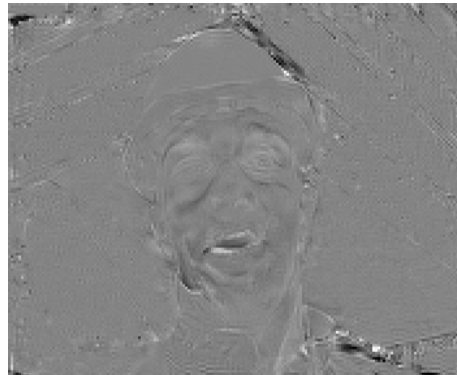
Approach for coding:

→ transmit only modified image areas

→ extend still image coding into temporal domain (kind of “3D”)



No compensation
H = 6.4bit



Motion compensated
H = 4.4bit

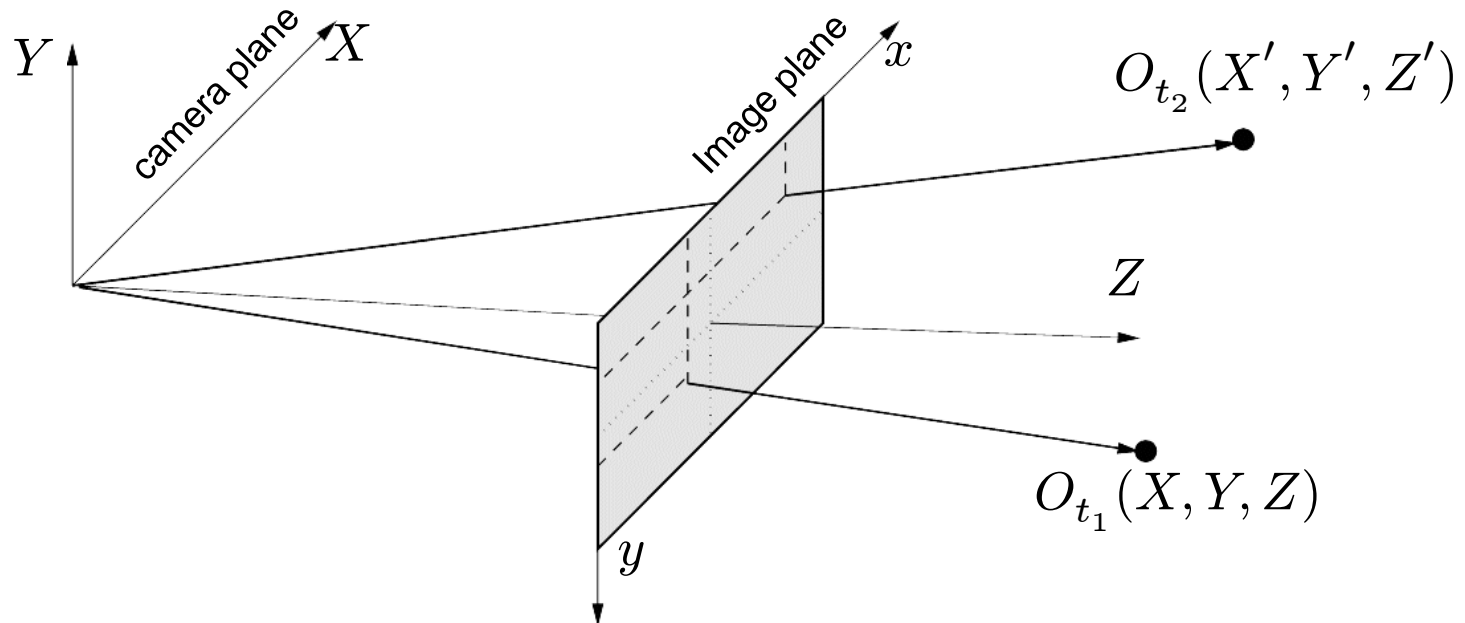
Approach:

Estimate motion
(reason for changes of the image)

Problem:

How to describe “motion”?

Image Generation



Mapping 3D world \rightarrow 2D image plane:

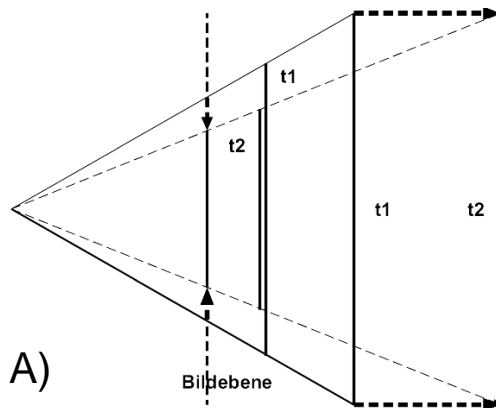
Geometrical optics for modeling

Motion:

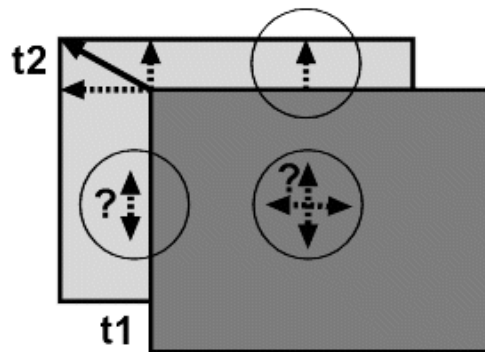
- Projection onto image plane is time variable
- 3D object movement \rightarrow moving of 2D regions

Mapping of Motion

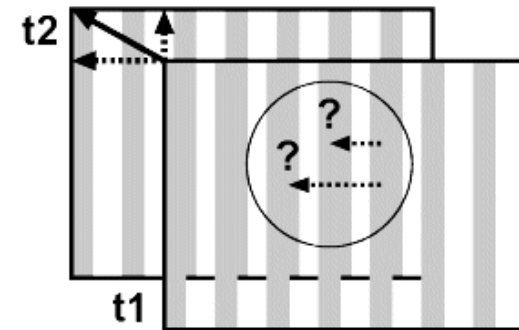
Problem: motion in the image plane is not unique
(no one-to-one mapping between 2D and 2D world)



A)



B)



C)

A) change of size caused by shortening (lengthening), change of depth, rotation

B) aperture problem → locale motion description

C) correspondence problem, in particular for periodic structures → aliasing

→ a unique description requires to assume a certain model

Modeling Motion

- **Consistency of objects:**

opaque, diffuse reflecting, geometrical form

- **Motion of objects:**

translation, rotation, deformation

a) estimating physical parameters

→ motion analysis

model: parametric description

b) finding correspondences

→ coding

model: displacement

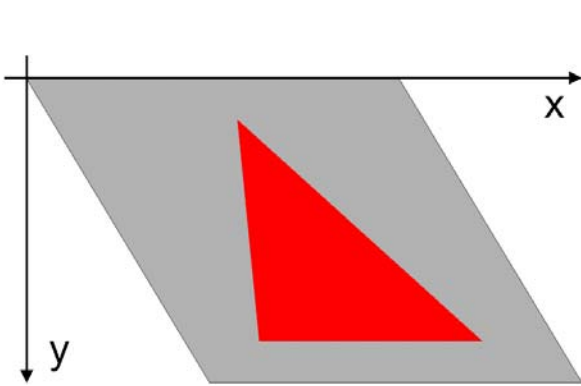
- **Movements of the camera**

zoom, pan, rotation

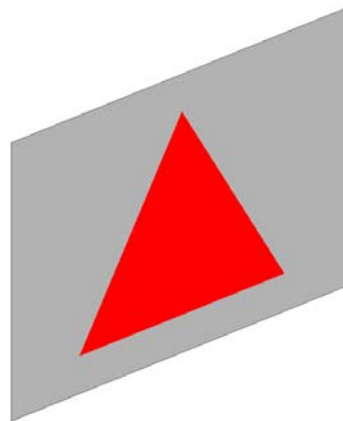
2D Affine Mapping

Transforming coordinates and coordinate systems

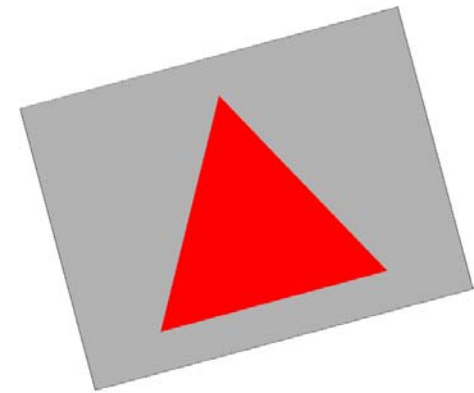
$$\begin{pmatrix} x' \\ y' \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0.5 \\ 0 & 1 \end{pmatrix}$$



$$\begin{pmatrix} 1 & 0 \\ -0.5 & 1 \end{pmatrix}$$



$$\begin{pmatrix} \cos(\alpha) & -\sin(\alpha) \\ \sin(\alpha) & \cos(\alpha) \end{pmatrix}$$

Parametric Motion Model (1)

3D Objektbewegung (3D affin) $\mathbf{X}' = \mathbf{A}\mathbf{X} + \mathbf{X}_0$

\mathbf{A} \rightarrow rotation, deformation

$\mathbf{X}_0 \in \mathbb{R}^3$ \rightarrow translation

\mathbf{X}, \mathbf{X}' \rightarrow coordinates in the 3D space

Mapping a point of an object assuming entral projection

$$x' = x \frac{Z}{Z'} + X_0 \frac{F}{Z'} \quad y' = y \frac{Z}{Z'} + Y_0 \frac{F}{Z'}$$

Resulting 2D motion of 3D moving of a rigid plane in space (3D):

$$x' = \frac{a_1x + a_2y + a_3}{a_7x + a_8y + 1} \quad y' = \frac{a_4x + a_5y + a_6}{a_7x + a_8y + 1}$$

Parametric Motion Model (2)

Modell	Parameter	Object form	Object motion	Projection
2D translation	2	arbitrary	2D translational	parallel
2D affine	6	planar	3D affine	parallel
3D affine	8	planar	3D affine	central
2D flexible	2N	2D linear in sections	2D flexible in sections	arbitrary

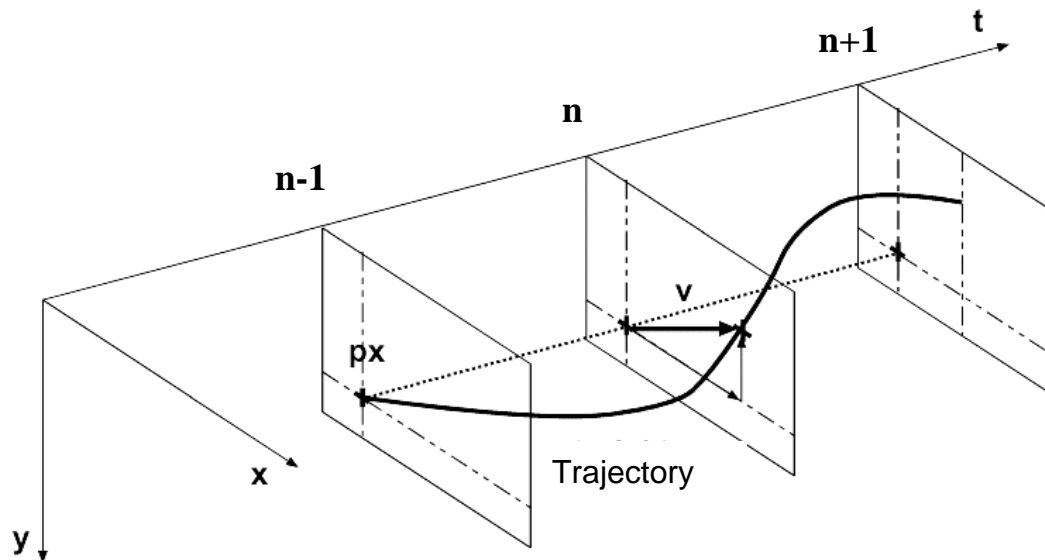
Estimating the parameters:

- directly via feature points
 - requires to identify N measurement points (minimum 1 point / parameter)
- indirectly via a displacement vector field

Describing Motion as Displacement

Assumption: $g(\mathbf{x}) \mapsto O(\mathbf{X}), \quad \forall \mathbf{x} \in \mathbb{R}^2, \mathbf{X} \in \mathbb{R}^3$

Moving of $O(\mathbf{X}) \rightarrow$ trace in the image sequence \rightarrow motion trajectory



Displacement vector $v(\mathbf{x})$: describes motion in the 2D image plane
Motion estimation: estimating the displacement $v(\mathbf{x})$

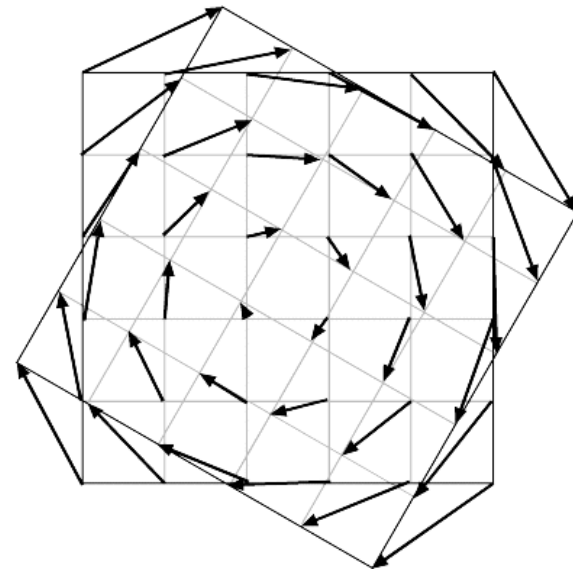
Displacement (Motion) Vector Field (MVF)

Vector field is discrete in space

$$\mathbf{v}_n = \{\mathbf{v}_n[\mathbf{x}] \in^2, \mathbf{x} \in^2\}$$

Characteristics:

- Each vector only describes translational motion
- Dynamic change of vectors
- Describing other types of motion (incl. object deformation) exploiting the change of vectors over time AND the context



Implementation:

- Describing motion discrete in space (location)
 - MVF is samples in space on a sub-grid of the image grid
 - MVF is termed dense, if there is a vector for each pixel
- Quantization and thresholding of the vector amplitude

Approaches for Motion Estimation

Motion Estimation (ME) $\mathbf{v}_n = \text{ME}(g_n, g_{n-1})$

Assumptions:

- Unique assignment of motion to 2D plane $g[x] \leftrightarrow O(X)$
- Intensity of pixel remains unaltered over time $\implies g_n[\mathbf{x}] = g_{n-1}[\mathbf{x} - \mathbf{v}[\mathbf{x}]]$

Algorithms

- Matching approach
minimizing an error criteria / maximizing a similarity criterion
- Gradient approach
evaluate the continuity equation \rightarrow solving an equation system
- Statistical approach
MVF is the realization of a random process; maximizing the probability

Correspondences in Space

Assuming dense MVF

$$\forall_{\mathbf{y} \in \mathcal{N}_{\mathbf{x}}^2} \mathbf{v}[\mathbf{x}] \approx \mathbf{v}[\mathbf{y}]$$

Motion model

2D motion can approximated as

Locally translational

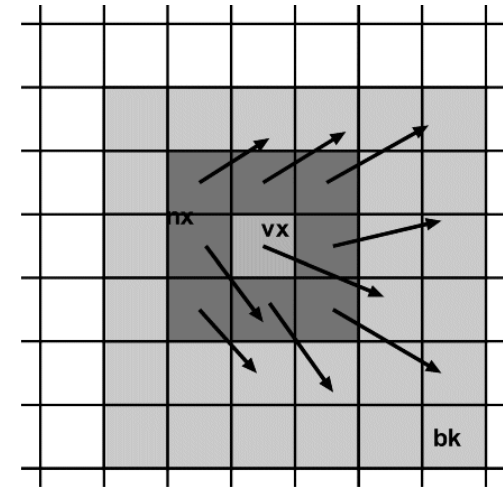
$$\forall_{\mathbf{y} \in \mathcal{N}_{\mathbf{x}}^2} \mathbf{v}[\mathbf{x}] = \mathbf{v}[\mathbf{y}]$$

„Matching“ principle

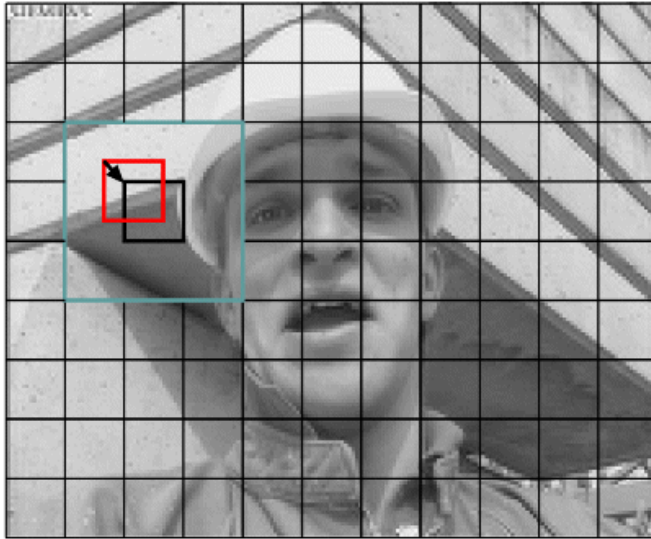
- Partitioning the image into regions, e.g.. blocks
- Minimizing a distance criterion $e()$ for each region

$$e(\mathbf{x}_k, \mathbf{v}) = \sum_{\mathbf{x} \in \mathcal{B}_k} f\left(g_n[\mathbf{x}], g_{n-1}[\mathbf{x} - \mathbf{v}[\mathbf{x}]]\right) \rightarrow \min$$

$f()$: function to weight each pixel error



Block Matching



g_{n-3}



g_n

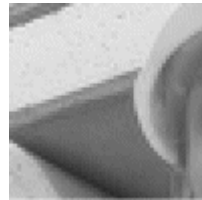
Error criterion:

$$e(\mathbf{x}_k, \mathbf{v}) = \sum_{\mathbf{x} \in \mathcal{B}_k} \left| g_n[\mathbf{x}] - g_{n-3}[\mathbf{x} - \mathbf{v}] \right|^\alpha, \quad \alpha \in \{1, 2\}$$

Estimation criterion:

$$\mathbf{v}[\mathbf{x}_k] = \operatorname{argmin}_{\mathbf{v}_i \in \mathcal{V}} \{e(\mathbf{x}_k, \mathbf{v}_i)\} \quad \mathcal{V}: \text{set of test vectors}$$

Progression of the Error Criterion (Example)



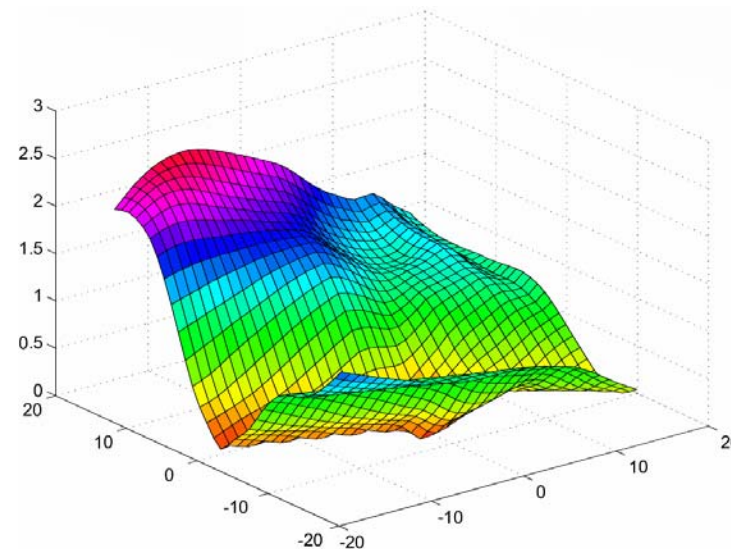
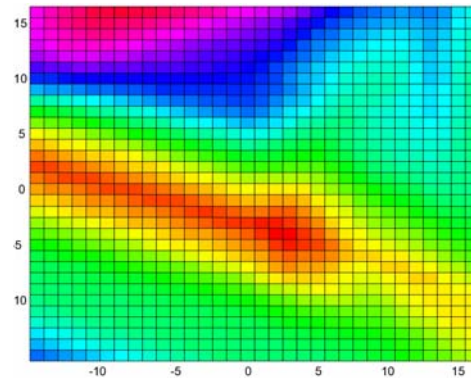
Search area



Reference block
from current image



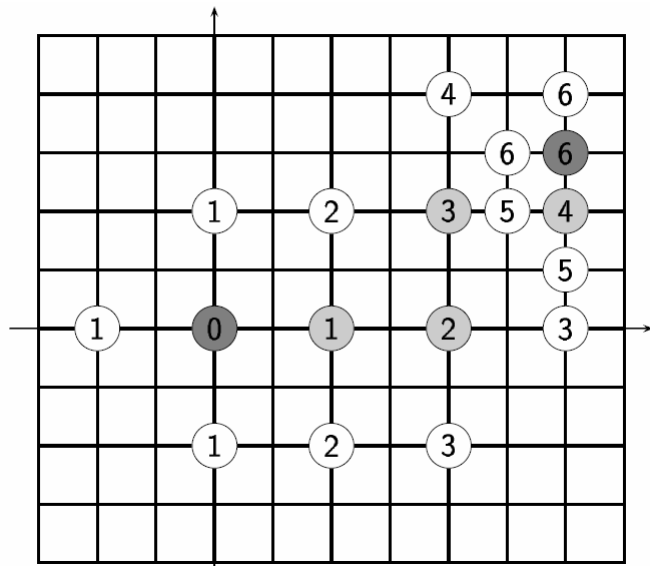
Most similar block
from previous image



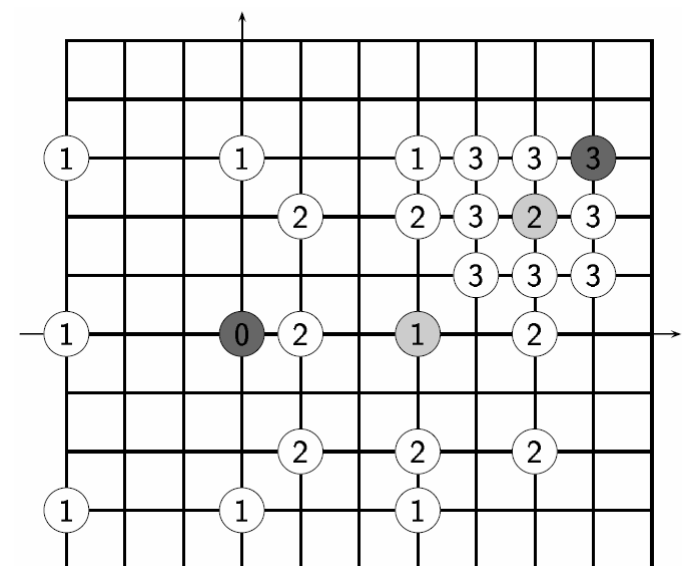
Progression of quadratic error $e()$ over the displacement v

Search Strategies

1. full search: computational expensive, guarantees a minimal error
2. logarithmic search (assumption: convex error progression)
3. search in multiple steps (three step with decreasing step size)
4. fast search based on Schwarz Inequality



Logarithmic search

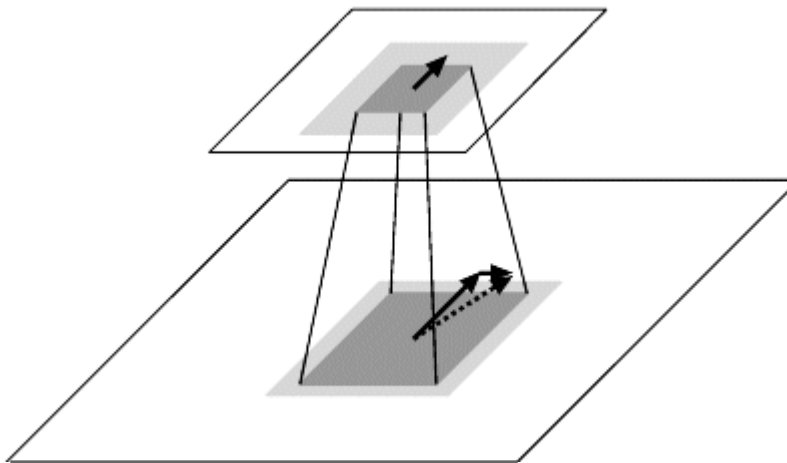


Multiple step search

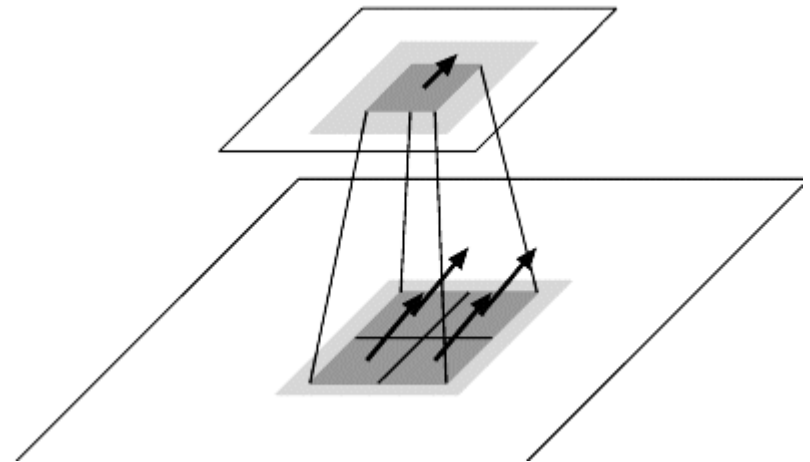
Hierarchical Search

4. Hierarchical Search

- a) approximation on image of reduced resolution
using search algorithm out of 1....4
- b) successive refinement on images of higher resolution
reduced search area \rightarrow reduced computational complexity



Identical image area per block
 \rightarrow Enlarged block size



Identical block size
 \rightarrow Refined description

Extensions (1)

Displacement vector amplitudes are quantized to resolution of image grid:

$$\mathbf{v}[\mathbf{x}] \in \Lambda$$

Increasing the amplitude resolution: $k \cdot \mathbf{v}_n[\mathbf{x}] \in \mathbb{Z}^2$

→ motion estimation on interpolated images

$$\mathbf{v}_n = \text{ME}\left(\text{E}(g_n), \text{E}(g_{n-1})\right)$$

Interpolation operator:

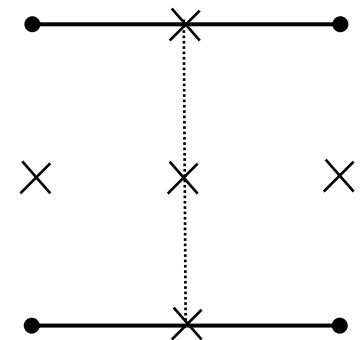
$$\text{E} : \tilde{g} = \text{E}(g) = [g]_{\uparrow k} * h_E$$

Typical interpolation by factor 2 (half-pel) or 4 (quarter-pel)

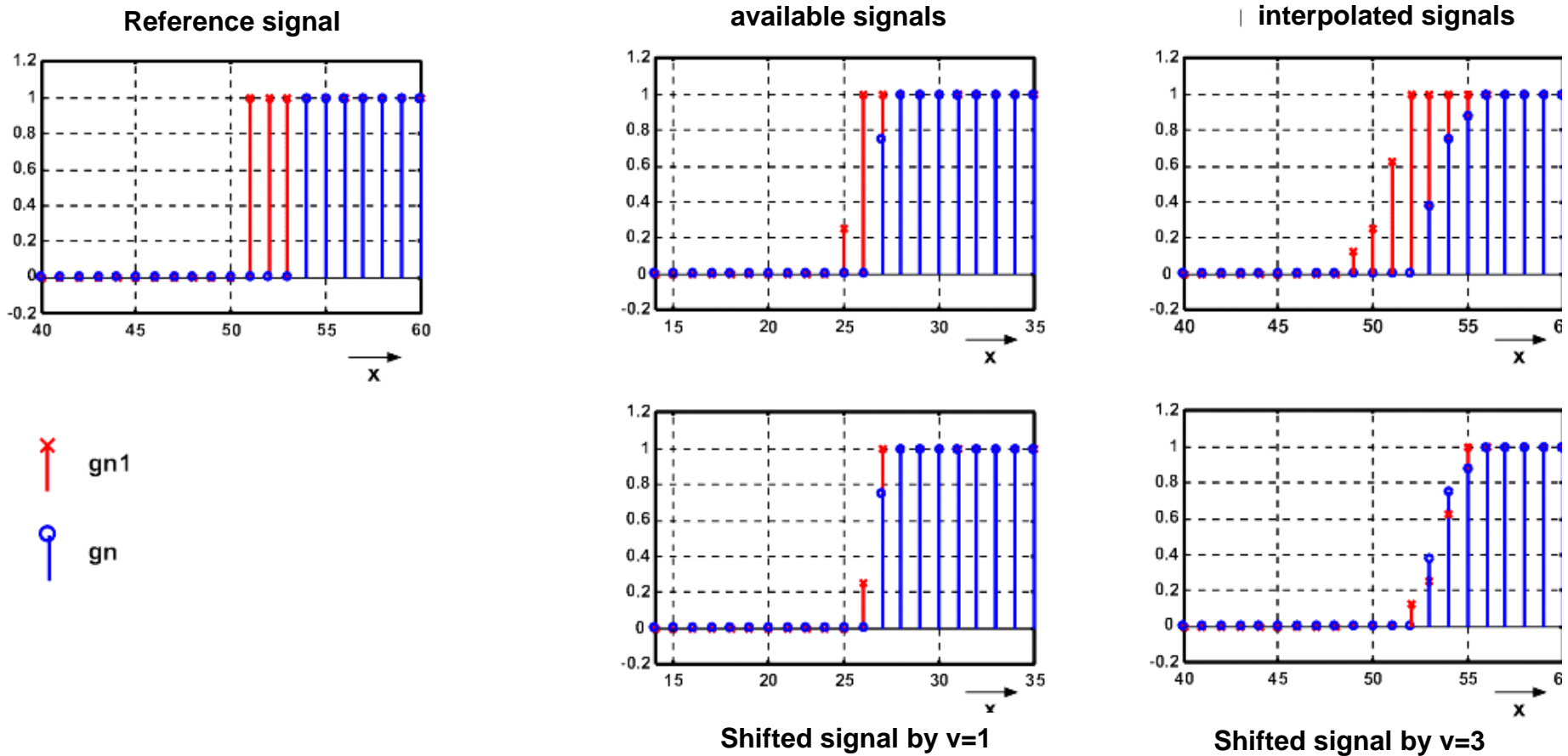
Example:

- $k = 2$
- bilinear interpolation (separable filter kernel)

$$h_E[x] = \frac{1}{2}\delta[x-1] + \delta[x] + \frac{1}{2}\delta[x+1]$$



Example for Half-pel Motion Estimation



Extensions (2)

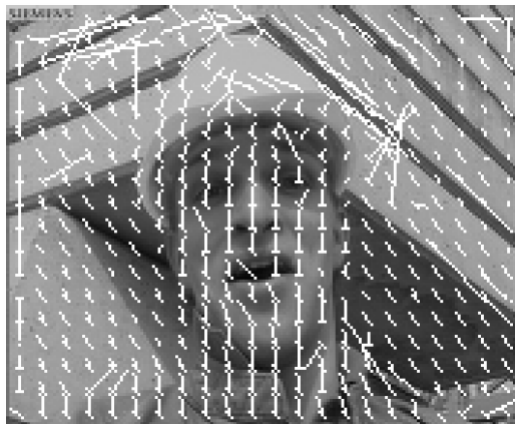
Regular vector fields → smoothness constraint:

Assumption: neighbored vectors describe similar motion
(homogeneous motion, rigid bodies)

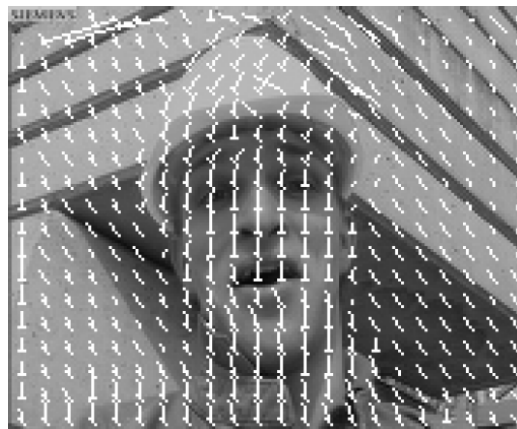
→ Motion estimation taking neighbored vectors into account

$$\mathbf{v}[\mathbf{x}_k] = \operatorname{argmin}_{\mathbf{v}_i \in \mathcal{V}} \left\{ \Psi(e(\mathbf{x}_k, \mathbf{v}_i)) + \lambda \sum_{\mathbf{y} \in \mathcal{N}_{\mathbf{x}_k}} |\mathbf{v}_i - \mathbf{v}[\mathbf{y}]| \right\}, \quad \lambda \in \mathbb{R}$$

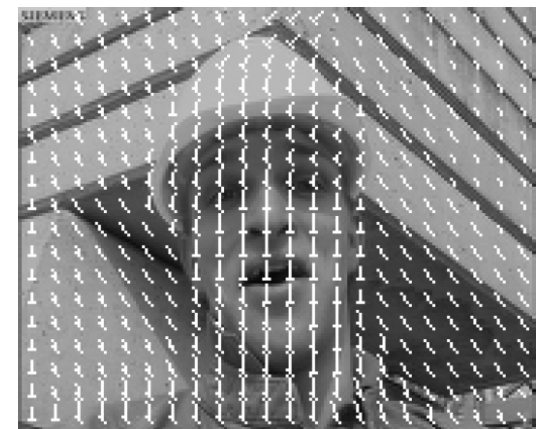
$\Psi(\cdot)$ Weight function



$\lambda = 0$



$\lambda = 0.05$



$\lambda = 0.5$

Extensions (3)

Increasing the spatial resolution of MVF:

- using smaller blocks
estimation is less reliable due to aperture effects and noise
→ hierarchical block matching
- interpolating motion vector fields
interpolation requires adaptation according to motion model and
consideration of motion discontinuities at boundaries



4x4 independent



16x16



4x4 hierarchical

Other Matching Approaches (1)

Correlation as measure for similarities of functions:

$$\begin{aligned}\varphi_{g_n g_{n-1}}(\mathbf{v}) &= (g_n * g'_{n-1})(\mathbf{v}) \\ &= c \cdot \sum_{\mathbf{x}} g_n[\mathbf{x}] \cdot g_{n-1}[\mathbf{x} - \mathbf{v}]\end{aligned}$$

$$g'[x] = g[-x] \quad \text{normalization: } c^{-1} = \sum_{\mathbf{x}} g_n[\mathbf{x}] \cdot \sum_{\mathbf{x}} g_{n-1}[\mathbf{x}]$$

Criterion for motion estimation

→ Maximizing the cross correlation function

$$\mathbf{v} = \operatorname{argmax}_{\mathbf{v}_i \in \mathcal{V}} \{ \varphi_{g_n g_{n-1}}(\mathbf{v}_i) \}$$

- high computational effort
- robust against illumination changes
 - reduced constraints for motion model

Other Matching Approaches (2)

Correspondences in the frequency domain

$$\begin{aligned} g[\mathbf{x}] & \circ \bullet G(\mathbf{f}_x) \\ g[\mathbf{x} - \mathbf{v}] & \circ \bullet G(\mathbf{f}_x) \cdot \exp(-j2\pi \langle \mathbf{v}, \mathbf{f}_x \rangle) \end{aligned}$$

Interpretation: motion results in characteristic shifts of the phase

Characteristics:

- high computational complexity
- phase signal typically have significant high frequency components
→ Estimation by matching unreliable
- displacements present in the entire image can be identified
- motion can not be assigned to local regions

Region Oriented Motion Estimation (1)

Partitioning into regions:

- Homogeneity of features (texture, motion)
- Correspondence of objects (*a priori* knowledge)

→ Estimating of the region form required

- Describing the contour
- Approximating as set of elementary regions such as blocks or triangles

Ill-posed problem: motions versus regions form

Region Oriented Motion Estimation(2)

Approaches for solving the problem:

- **Alternating estimation**

- Estimate motion based on existing region partitions (segmentation)
- Update the segmentation constraint to the feature **motion**

- **Simplified version**

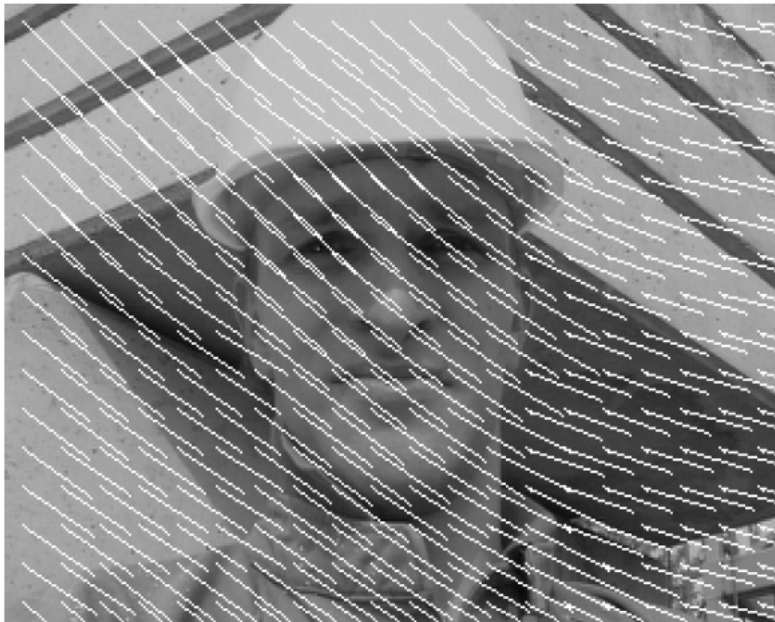
- Estimate motion using **elementary regions**
- Create segmentation by *Split and/or Merge*

- **Criterion for motion estimation includes a region model**

- Smoothness constraint:
 - motion and texture within regions homogeneous
 - contours of regions are smooth
- Discontinuities of features (in particular motion) at region boundaries

→ very high computational complexity

Examples 1:

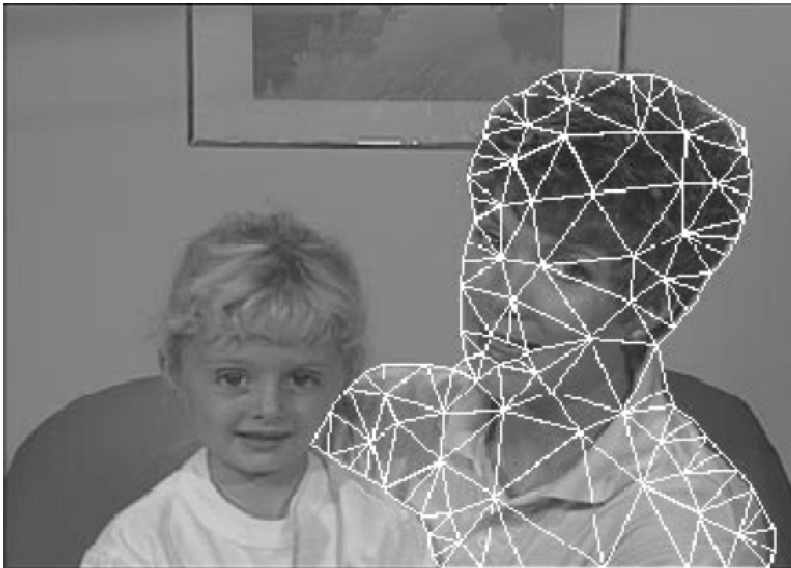


Vector field, calculated by
Block matching and
regularization

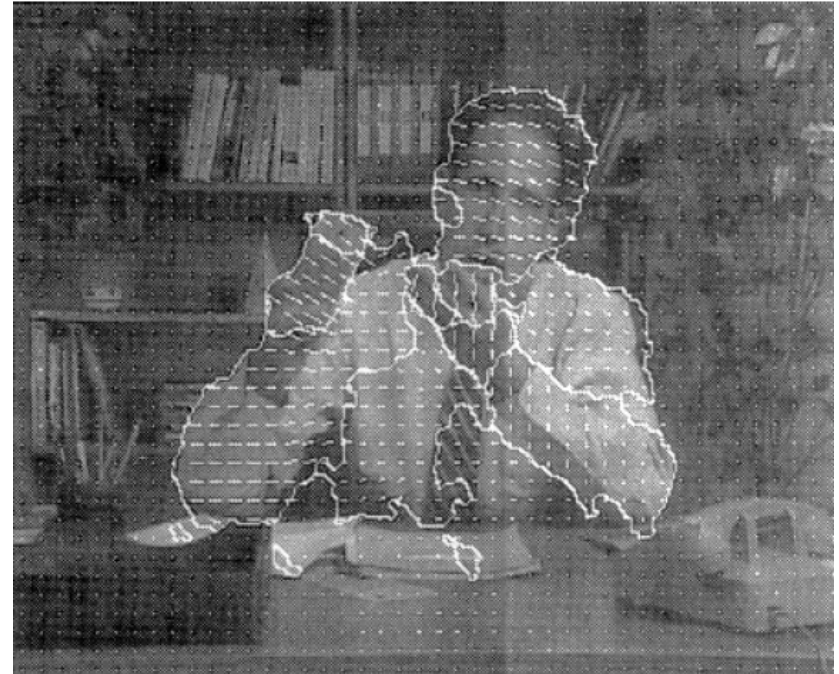


Segmentation derived from the
displacement vector field

Example 2:



Estimate „motion“ of
nodes of a triangular mesh



Estimate jointly the
Segmentation and motion
Based on a statistical approach

Motion Compensation

Goal: Coding the difference between images

→ compensation of motion

→ predict an image based on estimated motion

$$\hat{g}_n[\mathbf{x}] = \text{MC}(g_{n-1}, \mathbf{v}_n) = g_{n-1}[\mathbf{x} - \mathbf{v}[\mathbf{x}]], \quad \mathbf{x} \in \Lambda$$

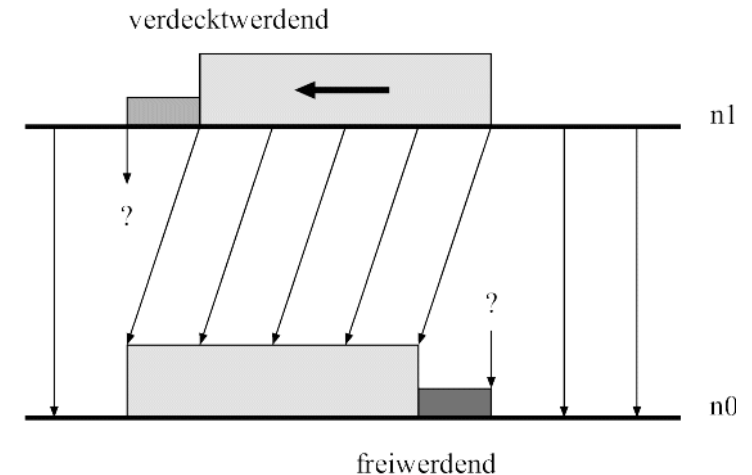
mit $|g_n - \hat{g}_n|^2 \rightarrow \min$

Error free prediction in reality not possible:

- images are sampled on a grid
- motion model is just an assumption
- VVF sampled and quantised
- border effects
- uncovered background

→ Coding of the **prediction error**

$$d_n = g_n - \hat{g}_n$$



Extensions (1)

Amplitude resolution higher than image grid of image to be compensated

- e.g. block matching with half-pel resolution
- Affine motion parameters

→ subpel compensation required:

$$\hat{g}_n = \text{RD} \left(\text{MC} \left(\text{EX}(g_{n-1}), \mathbf{v}_n \right) \right)$$

Reduction operator R:

$$\text{RD} : g = \text{RD}(\tilde{g}) = [g * h_R]_{\downarrow k}$$



8x8
Pixel-grid
accurate
resolution



8x8
Half-pel accurate
resolution

Extensions (2)

Blocking artifacts due to discontinuities at block boundaries

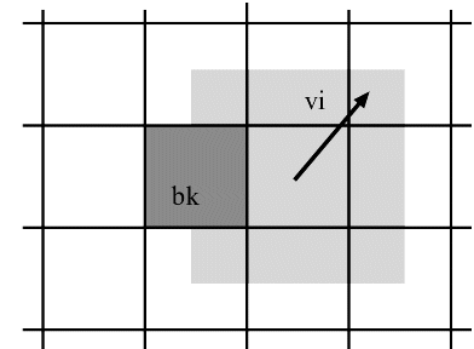
- post-processing filter
- In-loop filter (Deblocking-filter → H.264)
- Prediction with overlapping blocks (OBMC)

$$\hat{g}_n[\mathbf{x}] = \sum_{\mathbf{v}_i \in \mathcal{N}_{\mathbf{x}}^2} w(i) g_{n-1}[\mathbf{x} - \mathbf{v}_i]$$

w : 2D weighting function with $\sum_{x,y} w[x,y] = 1$

Optimization problem:

- displacement vectors depends on weight window
 - weight window depends on displacements vector
- iterative approach



Pixel accurate resolution with
Overlap compensation