







Virtuelle Realität Interaktion per Computer Vision

G. Zachmann
Clausthal University, Germany
cg.in.tu-clausthal.de

Was ist Computer-Vision?

- "The target problem ... is computing *properties of the 3-D world* from one or more *digital images*."
[Trucco, Verri: "Introductory Techniques for 3-D Computer Vision"]
- "...computer vision involves computers *interpreting images*." [dito]
- "Computer vision is the science and technology of *machines that see*. ... artificial systems that obtain *information from images*."
[Wikipedia]



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 2


Motivation

- Aufgaben:
 - Tracking von Objekten (Körper, Hand, Kopf, Augen)
 - Shape-Erkennung (Geste, Körperhaltung, Gesichtsausdruck)
 - Erkennung eines Bewegungsablaufs (z.B. dynamische Gesten)
- Langfristig der richtige Weg
- Aber: robuste und schnelle Erkennung nicht trivial!
 - Verschiedene Techniken werden eingesetzt, jede hat ihre Vor- und Nachteile

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 3

Background Subtraction

- Aufgabe: Bewegliches Interaktionsobjekt (Vordergrund) vom Hintergrund extrahieren
- Idee: Subtrahiere aktuell aufgenommenes Bild von einem Referenzhintergrund



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 4

■ Naiver Ansatz:

- Erzeuge Referenzbild $\mathcal{R} : [1, \text{width}] \times [1, \text{height}] \rightarrow [0, 255]^3$
 - nehme Hintergrundbild auf
- Subtrahiere pixelweise Eingabebild \mathcal{B}_t zum Zeitpunkt t vom Referenzbild

$$P(x, y) = \|\mathcal{R}(x, y) - \mathcal{B}_t(x, y)\|_2$$
- Ergebnisbild P ist ein Grauwertbild
- Pixel mit Differenz größer einem Schwellwert τ gehört zum Vordergrund

$$P(x, y) = \begin{cases} \text{Vordergrund} & : p(x, y) > \tau \\ \text{Hintergrund} & : p(x, y) \leq \tau \end{cases}$$

■ Problem: Schon leichte Variation des Hintergrundes (z.B. Beleuchtungsvariation) stört das Ergebnis

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 7

■ Etwas bessere Idee

- Farbverteilung für Hintergrund **pro Pixel** aufbauen
 - Hintergrund über längeren Zeitraum aufnehmen
 - Für jedes Pixel eigene Farbverteilung
 - Die einfachste Möglichkeit ist eine **unimodale Funktion**, d.h. eine Funktion mit nur einem Extremum
 - Häufigsten verwendete unimodale Funktion in der CV-Community ist die **Gaussfunktion**

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \mu, \Sigma) = \frac{1}{(2\pi)^{3/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$

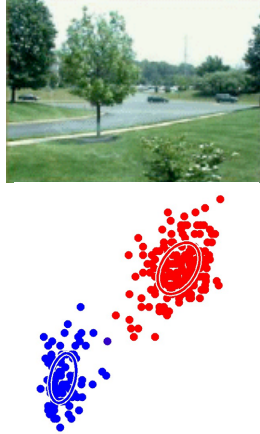
- μ ist Mittelwert; $\mathbf{u}_1, \mathbf{u}_2$ die Eigenvektoren von Σ und λ_1, λ_2 die Eigenwerte mit $\lambda_1 > \lambda_2 > 0$
- Funktioniert ganz gut z.B. bei Beleuchtungsvariation

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 8

- Unimodale Fkt versagt schon bei leicht beweglichem Hintergrund (z.B. wedelnde Bäume)
- Verbesserung: **Mixture of Gaussians** (Parameter z.B. durch Clustering-Verfahren bestimmen)

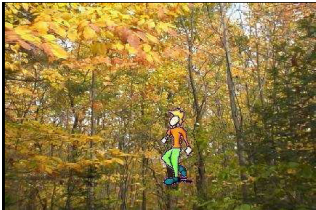
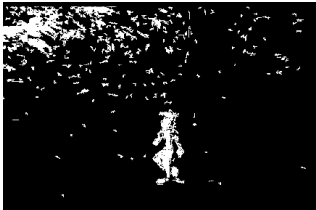
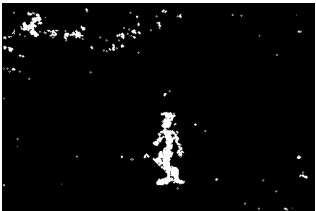

$$p(\mathbf{x}) = \sum_{k=1}^K \omega_k \mathcal{N}(\mathbf{x} | \mu, \Sigma)$$
 - Im Falle des Baumes, K=2, eine Gaussfkt für den grünen Baum, die anderen für den weiss-blauen Himmel
- **Non-parametric Model** bei stärker variierendem Hintergrund. z.B. Kernel Density Estimator

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N k\left(\frac{\mathbf{x} - x_n}{h}\right)$$



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 9

Beispiel

	
Source	Einfaches Thresholding
	
Non-parametric Model	Mixture of Gaussians

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 10

Grenzen des Verfahrens

- Variation des Hintergrundes nur in begrenztem Umfang
 - Objekte/Personen (die ich nicht tracken will) im Hintergrund dürfen sich nicht bewegen
 - Kameraposition statisch
- Tracking über längeren Zeitraum kaum möglich
 - Problem des sich stetig verändernden Hintergrunds, z.B. Beleuchtungsvariation

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 11

Farbsegmentierung

- Aufgabe: Bestimme alle Pixel im Bild die zu dem zu detektierenden Objekt gehören
- Voraussetzung: homogene Farbe des Zielobjektes
- Idee:
 - Farbverteilung des Objektes erstellen
 - Segmentierung der Vordergrundpixel über Farbraum
- Vorteil: Funktioniert auch bei beweglicher Kamera und variierendem Hintergrund



Robocup 2008



Farbhistogram
[Kai Uwe Barthel, FHTW Berlin]

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 12

Notationen

- Gegeben die Zufallsvariablen
 - X ein 3-dim Zufallsvektor, der Farbwerte repräsentiert
 - Y eine Zufallsvariable, welche die Segmentierung repräsentiert, d.h. sie kann die Werte fg (foreground) und bg (background) annehmen
 - rgb repräsentiert einen konkreten Farbwert z.B. (255,0,0)
 - $P(X=rgb)$ gibt die Wahrscheinlichkeit für das Vorkommen des Farbwertes rgb in einem Bild an
 - $P(Y=fg)$ ist die Wahrscheinlichkeit, dass ein Pixel im Bild zum Vordergrund gehört, $P(Y=bg)$ analog
 - $P(X=rgb | Y=fg)$ ist die bedingte Wahrscheinlichkeit, dass der Farbwert rgb zum Vordergrund gehört, analog für bg
 - Im Folgenden kürzen wir $P(X=rgb | Y=fg)$ mit $P(rgb | fg)$ ab, analog für bg

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 13

Segmentierung

- Pixel wird als Vordergrund klassifiziert, wenn

$$\frac{P(fg|rgb)}{P(bg|rgb)} > \tau$$
- Berechnung der Wahrscheinlichkeiten über *Regel von Bayes*

$$\frac{P(fg|rgb)}{P(bg|rgb)} = \frac{P(rgb|fg)P(fg)P(\overline{rgb})}{P(rgb|bg)P(bg)P(\overline{rgb})}$$
- Satz: *Regel von Bayes*

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 14

Ground-Truth-Daten erstellen

- Wahrscheinlichkeiten auf der rechten Seite der Gleichung

$$\frac{P(fg|rgb)}{P(bg|rgb)} = \frac{P(rgb|fg)P(fg)}{P(rgb|bg)P(bg)}$$
 sind einfach zu ermitteln z.B. durch manuelles Labeln und Histogram Counting
- Training: Zielobjekt aus verschiedenen Viewpoints und unter mehreren Beleuchtungsbedingungen aufnehmen


$$P(rgb|fg) = \frac{fg(rgb)}{\#Pixel_{fg}} \quad P(rgb|bg) = \frac{bg(rgb)}{\#Pixel_{bg}}$$

$$P(fg) = \frac{\#Pixel_{fg}}{\#Pixel_{fg} + \#Pixel_{bg}} \quad P(bg) = \frac{\#Pixel_{bg}}{\#Pixel_{fg} + \#Pixel_{bg}}$$

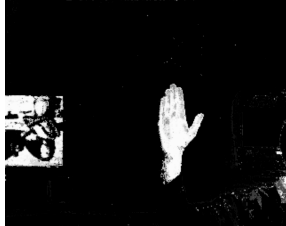
G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 15

Grösse der Trainingsdaten

- Ziel: **robuste** Verteilung mit wenigen Bildern; z.B. durch
 - Mixture of Gaussians
 - Kernel Density Funktion
- Histogramm zwar genauer, **aber** viel größerer Trainingsdatensatz erforderlich
- Beispiel: Hautsegmentierung



Originalbild



Hautsegmentierung

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 16

Threshold

- Wahl des **Schwellwertes** ist Kompromiss zwischen "false positives" und "false negatives"

Comparison of histogram and mixture models

Probability of correct detection

Probability of false detection

gut schlecht

gut schlecht

1. Histogram model using histogram with 32³ bins
2. Mixture of gaussian model

$\tau \ll$

$\tau \gg$

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 17

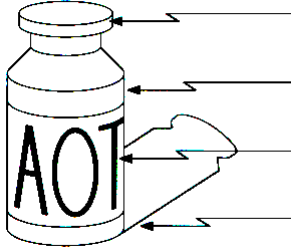
Grenzen des Verfahrens

- Hintergrundabhängig
 - Objektfarbe darf nicht zu häufig im Hintergrund vorkommen
- Beleuchtungsabhängig
 - Helligkeit der Lichtquelle(n)
 - Farbe der Lichtquelle(n)
- Speziell bei Hautsegmentierung: Personenabhängig
 - verschiedene Menschen haben verschiedene Hautfarbe

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 18

Bildkanten

- Bisher Objektfarbe als Feature
- Kanten beschreiben Form des Objektes
- Kanten entstehen durch mehrere Faktoren



Diskontinuität der Oberflächen Normale

Unterschied im Tiefenwert

Farbunterschied

Beleuchtungsvariation

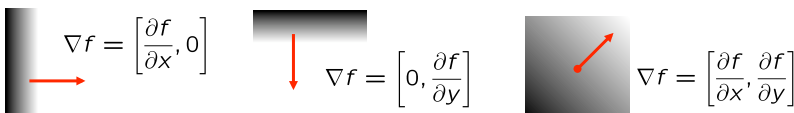
G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 19

Extraktion

- Kanten sind Farbunterschiede pro Bildraumabstand, wird also durch Bildgradient repräsentiert

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

- Gradient zeigt in Richtung der größten Änderung



- Richtung und Intensität des Gradienten

$$\theta = \tan^{-1} \left(\frac{\partial f / \partial y}{\partial f / \partial x} \right) \quad \|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x} \right)^2 + \left(\frac{\partial f}{\partial y} \right)^2}$$

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 20

- Analytische Funktion eines Bildes nicht bekannt
- Deshalb diskrete Ableitung

$$\frac{\partial f}{\partial x}(x, y) \approx f(x + 1, y) - f(x, y)$$

- Bekannteste Kantenoperatoren:



$\begin{bmatrix} +1 & 0 \\ 0 & -1 \end{bmatrix}$	$\begin{bmatrix} -1 & 0 & +1 \\ -1 & 0 & +1 \\ -1 & 0 & +1 \end{bmatrix}$	$\begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix}$
$\begin{bmatrix} 0 & +1 \\ -1 & 0 \end{bmatrix}$	$\begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ +1 & +1 & +1 \end{bmatrix}$	$\begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix}$
Roberts	Prewitt	Sobel

- Implementierung als Faltung

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 21

Canny-Edge-Detector

- Für viele CV-Algorithmen wird ein binäres Kantenbild benötigt
- Wir können bisher nur Kantenintensitätsbild berechnen

- Guter Algorithmus: **Canny Edge Detector**

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 22

- Canny's Zielsetzung:
 - **Gute Kantendetektion:** Der Algo soll so viele echte Kanten finden wie möglich
 - **Gute Lokalisierung:** Von Algo gefundene Kanten sollten so nah wie möglich an tatsächlichen Kanten liegen
 - **Gute Kantenantwort:** Eine Kante sollte nicht mehrfach gefunden werden
- Lösung kann durch Ableitung der Gaussfunktion approximiert werden

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 23

Algorithmus

- Gaussian Smoothing Filter auf Bild anwenden, z.B.

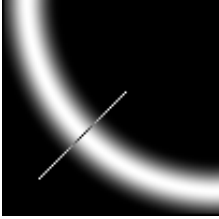
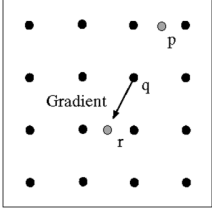
$$\frac{1}{159} \begin{bmatrix} 2 & 4 & 5 & 4 & 2 \\ 4 & 9 & 12 & 9 & 4 \\ 5 & 12 & 15 & 12 & 5 \\ 4 & 9 & 12 & 9 & 4 \\ 2 & 4 & 5 & 4 & 2 \end{bmatrix}$$
- Ableitung des resultierenden Bildes berechnen, z.B. Sobel Op.
- Non-Maximum-Suppression Algo anwenden (nächste Folie)
- Hysteresis Thresholding:
 - Sei K das Kantenbild und T_1, T_2 zwei Schwellwerte mit $T_2 > T_1$
 - markiere alle Pixel $K(x, y) \geq T_2$ als Kante
 - markiere alle Pixel $K(x, y) < T_1$ als nicht Kante
 - Markiere alle Pixel $T_1 \leq K(x, y) < T_2$ als Kante, wenn es einen Pfad von einem Kantenpixel nach (x, y) gibt, der nur Pixel mit Wert $\geq T_1$ hat

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 24

Non-Maximum-Suppression


- Idee: Wähle lokales Maximum als Pixel und lösche alle anderen Kantenpixel
- **Algorithmus:**
 - Betrachte jedes Pixel q im Kantenbild K
 - *Teste für benachbarte Pixel p und r :*

$$K(p) < K(q) \wedge K(r) < K(q)$$
 - *Falls ja $\rightarrow q$ ist lokales Maximum ; setze q als Kantenpixel*
Sonst; setze q als Hintergrundpixel





G. Zachmann Virtuelle Realität und Simulation - WS 08/09
Interaktion per Computer-Vision 25


Beispiel




Kantenoperator
→



↓ Non-max-suppression



← Hysteresis Thresholding



G. Zachmann Virtuelle Realität und Simulation - WS 08/09
Interaktion per Computer-Vision 27

Gesichtserkennung

- Unterscheide zwischen vier verschiedenen Problemstellungen
 1. **Lokalisation**: finde Position aller Gesichter in einem Bild
 2. **Erkennung**: bei gegebener Position des Gesichtes, ordne Gesicht einer bestimmten Person zu
 3. **Emotion/Ausdruck**: erkennen ob eine Person lacht, traurig ist etc.
 4. **Tracking**: Gegeben initiale Position, verfolge Gesicht über einen längeren Zeitraum
- Wir behandeln primär Lokalisation
- Tracking kann als Lokalisation pro Frame formuliert werden;
 - Ab Frame 2 ist grobe Position des Gesichtes bekannt
 - Daher ist die abzusuchende Bildregion deutlich kleiner (unter der Voraussetzung kontinuierlicher Bewegung)

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 28

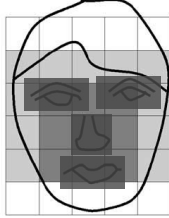
Video



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 29

Template-Matching-Ansatz

- Idee:
 - Erstelle ein Muster des Gesichtes.
 - Suche im Bild nach diesem Muster.
 - Melde Position des Gesichtes bei hinreichend guter Übereinstimmung
- Vorgehensweise
 - Extraktion **invarianter Features**, d.h. Features die möglichst robust gegenüber Störungen (*Beleuchtung, Hintergrund, unterschiedliche Gesichtsform*) sind
 - Berechnung der Wahrscheinlichkeit für das Vorhandensein eines Gesichtes durch **Kombination** der Features



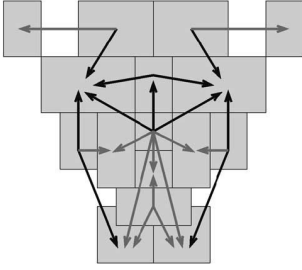
G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 30

- Extraktion invarianter Features:
 1. Extraktion von **low-level Features** (z.B. Kanten, Hautfarbe)
 2. Basierend auf den low-level Features, berechne durch Template Matching potentiellen Positionen von **Gesichtsfeatures** wie
 1. *Augen(-brauen)*
 2. *Nase(-nlöcher)*
 3. *Mund/Lippen*
 4. *Wangen*
 5. *Gesichtsrand*

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 31

Konkretes Verfahren

- Idee: Helligkeitswerte im Gesicht variieren zwar von Pixel zu Pixel, aber Relationen im Großen bleiben erhalten
- Erstellen des Gesichts-Templates:
 - Definiere Template als Menge von rechteckigen Teilregionen des Gesichtes
 - Jede Region entspricht einem signifikanten Gesichtsfeature wie Augen, Mund, Nase, Wangen, Stirn
 - Feature sind durchschnittliche Helligkeitswert für jede dieser Regionen
 - kann entweder über Ground-Truth-Datensatz berechnet werden oder
 - durch Expertenwissen (d.h. manuell festlegen)
 - Speichere zu jedem paar von Regionen (mit Pfeilen markiert) Quotient der Helligkeitswerte



Das Diagramm zeigt ein schematisches Gesichtstemplate, das in rechteckige Regionen unterteilt ist. Pfeile verbinden diese Regionen in einem Netzwerk, was die Beziehungen zwischen den verschiedenen Gesichtsteilen (wie Augen, Nase, Mund, Wangen, Stirn) darstellt, die für die Berechnung von Helligkeitsquotienten genutzt werden.

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 32

- Lokalisierung des Gesichtes im Bild:
 - Berechne im Eingabebild für alle möglichen Positionen die Regionen
 - Ein Regionenpaar matcht mit dem Template, wenn der Helligkeitsquotient einen bestimmten Schwellwert überschreitet
 - Bildposition wird als Gesicht klassifiziert, wenn mehr als eine bestimmte Anzahl an Regionenpaaren passen

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 33

Gestenerkennung

- Aufgabe: Gegeben n Handgesten, entscheide ob bzw. welche Geste im Bild zu sehen ist
- Anwendungen:
 - Gesten als Shortcuts in Programmen
 - Simple Navigations in VR
 - Automatische Übersetzung von Gebärdensprachen in Wort/Schrift
- Zwei Grundsätzlich verschiedene Herangehensweisen
 - Modellbasierte Methoden
 - Bildbasierte Methoden

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 42

Video



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 43



Modellbasierte Methoden

- Basis ist ein Modell der Hand bzw. der Handgesten
 - Als Modell kann z.B. eine künstliche/gerenderte Hand dienen
- Ziel:
 - Vergleiche die Gesten, generiert durch das Handmodell mit dem Eingabebild
 - Bestimme das Modell, das am besten zum Bild passt
 - Kantenabstandsmass (*Chamfer, Hausdorff*, direkter Vergleich)
 - Überlappung der Handregionen von hautsegmentiertem Eingabebild und Template
 - Passt die Modellgeste hinreichend gut, wird die Geste vom Verfahren erkannt, sonst wird Bild als Hintergrund erkannt

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 44

Ein modellbasierter Ansatz

- Gegeben: Eine Handgeste G und Eingabebild I

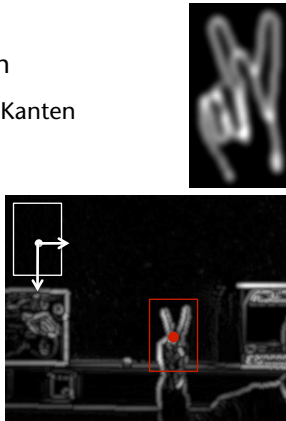



- Ziel: Finde die Position im Eingabebild I , die am besten zur Handgeste G passt
- Verwende Kanten als Feature zum Vergleich von Eingabebild und Handgeste

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 45

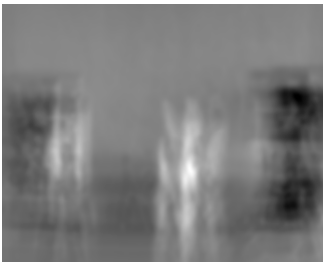
Confidence Map

- Template T_G zur Handgeste G generieren
 - Rendere Handgeste G und extrahiere daraus Kanten
 - Wende Smoothing auf Kantenbild an
- Matching Algorithmus:
 - Für alle Pixel im Eingabebild
 - Extrahiere Kanten aus Eingabebild I ,
Ergebnis ist Kantenbild I_K
 - Berechne "Ähnlichkeit" zwischen Template
und lokale Bildumgebung zum Pixel durch
pixelweise Multiplikation von Template und Kantenbild
- Dies beschreibt eine Faltung des Kantenbildes I_K mit dem Template T_G (genauer mit Spiegelung des Templates)



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 46

- Ergebnis nennt man *Confidence Map*
 - Gleiche Auflösung wie Eingabebild
 - Lokale Maxima geben beste Kandidaten für Position der Handgeste an
- Aufwand der Faltung ist $O(n*m*w*h)$
 - nxm ist Auflösung von I_K und wxh Auflösung von T_G
- Analogie zwischen Faltung im Ortsraum und Multiplikation im Fourierraum (nächste Folie)
- Aufwand für *Fast-Fourier-Transformation(FFT)*: $O(n*m*\log n*\log m)$
- Aufwand für Multiplikation im Fourierraum $O(n*m)$
- Faltung über Fourierraum lohnt wenn $\log n*\log m < c*w*h$



G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 47

Video

Original image

Template index: our approach; Position: hand labeled

Chamfer based approach

Our approach

G. Zachmann Virtuelle Realität und Simulation - WS 08/09 Interaktion per Computer-Vision 49