

TrabHCI

Lifelong

Learning

Programme

Technologies to Reduce the Access Barrier in Human Computer Interaction Erasmus Intensive Programme 2012-1-FI1-ERA10-09684

Image processing for gesture recognition: from theory to practice

Ivan Bernabucci, Michela Goffredo, Maurizio Schmid University Roma TRE ivan.bernabucci@uniroma3.it

- 1. Image and video processing for posture and gesture recognition (today)
- 2. Object oriented programming (tomorrow, with Ivan)
- 3. Kinect programming and SDK (tomorrow, with Antti)
- 4. Project work... (starting wednesday)

TrabHC



Image and video processing for posture and gesture recognition

Ivan Bernabucci, Michela Goffredo, Maurizio Schmid Roma Tre University maurizio.schmid@uniroma3.it

Images and video are everywhere!



Personal photo albums



Movies, news, sports

flick

REAKING NEWS



Broadcast Yourself





What's Computer Vision?







• **Data transformation** from a still or video camera into either a decision or a new representation for achieving some particular **goal**.

- The input data may include some **contextual information** such as "the camera is mounted in a car" or "one person is in the scene".
- The **decision** might be "the person is still" or "there are 5 cars on the road"...





Perceive the "world behind the picture"



In a machine vision system, a computer receives a grid of numbers from the camera or from disk: that's a *digital image*.





Perceive the "world behind the picture"

Moreover, data is corrupted by noise and distortions:

- from variations in the world (weather, lighting, reflections, movements);
- imperfections in the lens and mechanical setup, finite integration time on the sensor (motion blur);
- electrical noise in the sensor or other electronics;
- compression artifacts...

Additional contextual knowledge can often be used to work around the limitations imposed on us by visual sensors.

<u>General rule: the more constrained a computer vision context is,</u> <u>the more we can rely on those constraints to simplify the</u> <u>problem and the more reliable our final solution will be.</u>





Origins of computer vision

(a) Original picture.



(b) Differentiated picture.

L. G. Roberts, Machine Perception of Three Dimensional Solids, Ph.D. thesis, MIT Department of Electrical Engineering, 1963.





Vision as a source of semantic information



Vision as a source of semantic information



Scene and context categorization



Qualitative spatial information



Challenges: viewpoint variation



Challenges: illumination



image credit: J. Koenderink





Challenges: scale



Bio



Challenges: deformations



Xu, Beihong 1943





Challenges: occlusions



Magritte, 1957





Challenges: background clutter





Emperor shrimp and commensal crab on a sea cucumber in Fiji Photograph by Tim Laman



onal Geographic Society. All rights reserved.

Challenges: object intra-class variation



Challenges: local ambiguity







Challenges or opportunities?

- Images may be confusing, but they also reveal the structure of the world through numerous cues
- e.g. Linear perspective, texture gradient,...







Challenges or opportunities?

Shape and lighting cues



Challenges or opportunities?

Grouping cues: Similarity (color, texture, proximity, shape...)







Applications

Real-time stereo



NASA Mars Rover



Structure from motion



Multi-view stereo for community photo collections





Pollefeys et al.

Goesele et al. ROMA TRE Biomedical Engineering Laborator

Applications



Factory inspection

Surveillance

TrabHCI



Reading license plates, checks, ZIP codes



Monitoring for safety (Poseidon)



Autonomous driving, ----robot navigation----



ROMA TRE BioLab³ Hold

Applications



Assistive technologies



Entertainment (Sony EyeToy)



Movie special effects





Digital cameras (face detection for setting focus,

[Face priority AE] When a bright part of the face is too bright





Visual search (MSR Lincoln)





- WHY Image and video processing for posture and gesture recognition for Human-computer Interaction (HCI)?
- Computing, communication and display technologies progress quicker & quicker.
- Mechanical devices (keyboard & mouse) for HCI are the bottleneck in the effective utilization of available progresses.
- The future is "natural", i.e. be inspired by the natural human to-human communication modalities:
 - Speech
 - Gestures



Human-computer Interaction

Human-computer interfaces inspired by human-human communication.

In 1991, Myron Krueger wrote a pioneering book, "Artificial reality", where it is reported:

"Natural interaction means voice and gesture. New interface technologies requires tools and features that mimic the principles of human communication."



Gesture recognition

- Gestures are useful for computer interaction since they're the most primary and expressive form of human communication.
- Webster's dictionary defines gestures as

"...the use of motions of the limbs or body as a means of expression; a movement usually of the body or limbs that expresses or emphasizes an idea, sentiment or attitude"



20 years of entertainment waiting for it...









Gestural taxonomy



Gesture recognition

- No single method for automatic gesture recognition is suitable for every application: each gesture-recognition algorithm depends on:
 - user cultural background;
 - application domain;
 - environment.







• System requirements vary depending on the aim of the application (i.e. entertainment system vs surgical system..)





Vision-based gesture recognition

- Visual interpretation of:
 - body pose
 - body gesture
 - hand pose
 - hand gesture



Using imaging devices: still/video cameras



In a nutshell



Vision-based gesture recognition

- Advantages of video-based gestures applications over conventional human-machine interaction:
 - Access information while maintaining total sterility: touchless interfaces; healthcare environments.
 - Easily explore large and complex data.
 - Provide a source of expressiveness.

TrabHCI

36

• Overcome social handicaps associated with impairments: people with disabilities; elderly.





Differences

Do you know the difference between handicap & disability?

Disability: inability to execute some class of movements, or pick up sensory information of some sort, or perform some cognitive function, that typical unimpaired humans are able to execute or pick up or perform.

Social handicap: inability to accomplish something one might want to do because of the environment/situation.





Differences

Do you know the difference between handicap & disability?

We can be handicapped, even when we are not disabled.

Italians who do not speak Japanese will be handicapped when they visit Tokyo, because while most people will be able to gather important information by reading signs on buildings, they will not.

And one can have disabilities, without being handicapped relative to many tasks, if the proper tools and supporting structures are provided.





Vision-based gesture recognition

- Gesture interpretation:
 - **1. Definition** of gestures and dictionary
 - **2. Temporal modelling** of gestures: set of temporal parameters
 - **3. Spatial modelling** of gestures: characterization of spatial properties of limb's trajectories



In a nutshell



- **Image**: 2D function, **I**(x,y), where x and y are spatial coordinates, and the amplitude of f at any pair (x, y) is called the intensity or gray level of the image at that point.
- If x, y, and the amplitude values of **I** are all finite, discrete quantities, we call the image a **digital image**.



TrabHCI

41



- **Digital image processing**: processing digital images by means of a digital computer.
- A digital image is composed of a finite number of elements (pixels) defined by:
 - location (x,y)

42

TrabHCI

• intensity value



An image may be continuous with respect to:

- the x- and y-coordinates;
- in amplitude.

...but to convert it into a digital form, we will need to sample that function in both domains (coordinates and amplitude).

- **sampling**: digitizing the image coordinate values.
- **quantization**: digitizing the image amplitude (black2white)



Sampling

Spatial resolution is the smallest discernible detail in an image.



Four representations of the same image, with variation in the number of pixels used:

a) 256 x 256;
b) 128 x 128;
c) 64 x 64;
d) 32 x 32.



Quantization

Gray-level resolution refers to the smallest discernible change in grav level.



Four representations of the same image, with variation in the number of grey levels used:



Color imaging

Most real-world images are not monochrome, of course, but full color!



Images from digital cameras are usually in Red-Green-Blue (RGB) colorspace.

A number of color spaces or color models have been suggested and each one of them has a specific color coordinate system and each point in the color space represents only one specific color.

Each color model may be useful for specific applications.





• What is it possible to do with color video frames for motion detection?

Example: skin detection









Image processing

- We said that **Computer vision** is the transformation of data from a still or video camera into either a decision or a new representation.
- **Image processing** is part of Computer Vision and aim at transforming the image so that information can be extracted.



48



Global operator: threshold

- Simple and useful image processing method for getting information is segmenting pixels with respect to their values, i.e. **segmenting objects**.
- The basic global threshold algorithm aims at
 - scanning the image pixel by pixel;

TrabHC

- labelling each pixel whether its value greater or less than a threshold value T.
- If the pixel value is >= T, then it's set to a maximum value M (usually 255)
- If the pixel value is < T, then it's set to a minimum value m (usually 0)



Global operator: threshold

• Of course, the result depends on the choice of T...







 Choosing the threshold looking at the histogram of the pixel intensity values



TrabHCI

51



Global operator: threshold

- Sometimes pixels belonging to the object of interest are in the range between 2 values.
- In this case, we need 2 thresholds
- Remember Beckham?





Thresholds in HSV colorspace: H(0-20); S (30-150); V(80-255)





• After image threshold, we have a **binary image** where white pixels (255) correspond to the object of interest (+ noise)



 Image processing include algorithms which allow connecting isolated pixels sufficiently close to other and/or deleting isolated pixels (Morphology operators)

TrabHCl 53



- The basic morphological transformations are called **dilation and erosion**, and they aim at:
 - removing noise;
 - isolating individual elements;
 - joining disparate elements

in an image.

Morphology is a local image operator.



Dilation is a convolution of an image A with a kernel B and causes <u>bright</u> regions within an image to grow.

• As the kernel B is scanned over the image, we compute the <u>maximal</u> pixel value overlapped by B and replace the image pixel under the central point with that maximal value.



 The kernel can be any shape or size (usually small solid square or disk)



Erosion is a convolution of an image A with a kernel B and causes <u>dark</u> regions within an image to grow.

• As the kernel B is scanned over the image, we compute the <u>minimum</u> pixel value overlapped by B and replace the image pixel under the central point with that minimum value.



 The kernel can be any shape or size (usually small solid square or disk)

TrabHCI

56



How and when do we need to use them?

 The erode operation is often used to eliminate "speckle" noise in an image. The idea here is that the speckles are eroded to nothing while larger regions that contain visually significant content are not affected.

TrabHCI

57





How and when do we need to use them?

TrabHCl

58

 The dilate operation is often used when attempting to find connected components (i.e., large discrete regions of similar pixel color or intensity). The utility of dilation arises because in many cases a large region might otherwise be broken apart into multiple components as a result of noise, shadows, or some other similar effect. A small dilation will cause such components to "melt" together into one.





Edge detection

- Tipically, image processing prefer to have information about **features**, rather than object detection.
- The **detection of edges** is a fundamental tool in image processing, in the areas of feature detection and feature extraction, which aim at identifying points in a digital image at which the image brightness changes sharply or, more formally, has discontinuities.
- AIM: to capture important events and changes in properties of the world assuming that **discontinuities in image brightness** are likely to correspond to:
 - discontinuities in depth and surface orientation,

TrabHCI

59

 changes in material properties and variations in scene illumination.



Edge detection

Set of connected curves that indicate the boundaries of objects.

• Applying an edge detection algorithm to an image may significantly reduce the amount of data to be processed and may therefore filter out information that may be regarded as less relevant, while preserving the important structural properties of an image.



What are discontinuities?

Consider the problem of detecting edges in the following 1D signal. We may intuitively say that there should be an edge between the 4th and 5th pixels.

| 5 | 7 | 6 | 4 | 152 | 148 | 149 |
|---|---|---|---|-----|-----|-----|
| | | | | | | |

BUT here?

| 5 | 7 | 6 | 41 | 113 | 148 | 149 |
|---|---|---|----|-----|-----|-----|
| | | | | | | |

Not trivial at all... **TrabHCl** 61



Edge detection

Hopefully...



Mr. Canny



Mr. Sobel



Mr. Roberts



...& Miss Lena





Edge detection



Lena for Mr. Canny

63

TrabHCI



Lena for Mr. Sobel



Lena for Mr. Roberts



Image processing: basics

• **Image processing** aims at **transforming the image** so that information can be extracted.



In a nutshell



Video camera

- Monochrome
- color
- IR IR

- ▶ 1 camera
- ▶ 2 or more cameras
- Stereoscopic view

- CCD
- CMOS

- Depth cameras
- Spatial resolution (#pixels)
- Temporal resolution (fps)





