# Max, our Agent in the Virtual World
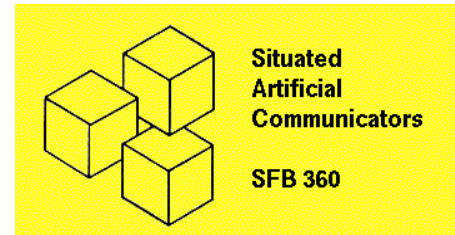
A Machine that Communicates with Humans

## Ipke Wachsmuth

University of Bielefeld

---

## University of Bielefeld

*Collaborative Research Center SFB 360 and
Artificial Intelligence & Virtual Reality Lab*

Situated
Artificial
Communicators

SFB 360

Technische Fakultät

Labor für
Künstliche Intelligenz
& Virtuelle Realität

sound check

---

## Collaborative Research Center SFB 360

**SFB 360 Thematic fields**

- Speech and Visual Perception
- Perception and Reference
- Knowledge and Inference
- Speech-Action Systems

*started in July 1993, overall funding
by Deutsche Forschungsgemeinschaft
Directors:  Prof. Gert Rickheit
                 Prof. Ipke Wachsmuth*
`www.sfb360.uni-bielefeld.de`

### Disciplines involved

- Linguistics
- Psycholinguistics
- Psychology
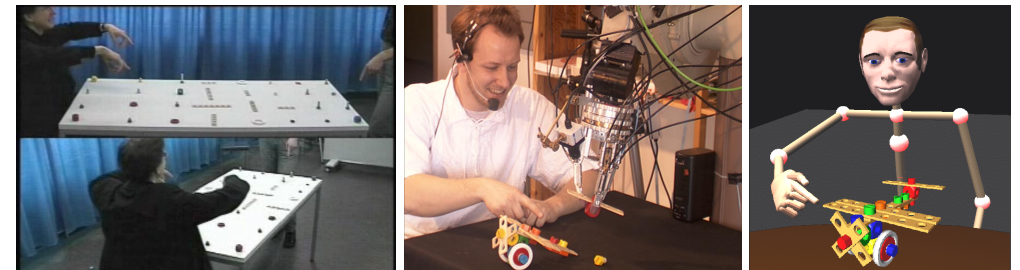- Informatics
- Neuroinformatics
- Artificial Intelligence

---

## Leading research questions

Situated
Artificial
Communicators
SFB 360

How do humans communicate in a cooperative task robustly and successfully?

What can be learned from this about particular features of human intelligence?

Can we transfer communication abilities to artificial systems of robotics and virtual reality?
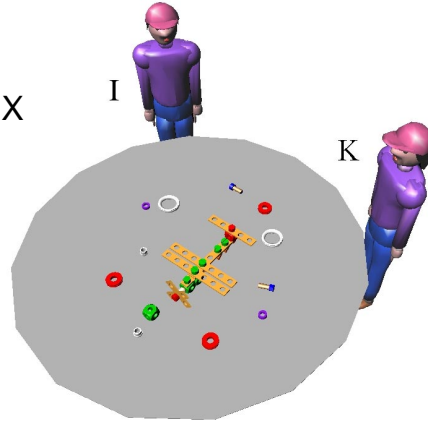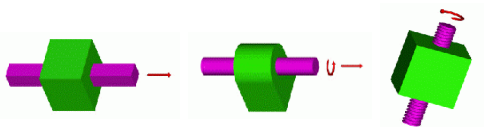
## Scenario for investigation

As most cognitive abilities are decisively situated, a specific reference situation serves to investigate task-oriented discourse.

For illustration the assembly of a model aeroplane from the BAUFIX construction kit is used.

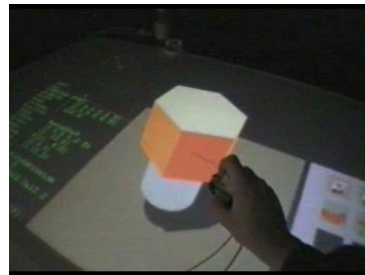A human instructor (I) and an artificial constructor (K) cooperate by way of an „Instructor-Constructor Dialog"

---

## Situated communication

I: Mount it at the right.
K: You mean here?

---

## Virtual Constructor

- „*lower kinematic pairs"* can be modeled (uncoupled or coupled)
- based on Roth's 1994 book: „Konstruieren mit Konstruktionskatalogen"



prismatic pair        cylindrical pair        helical pair

revolute pair        spherical pair        planar pair

---

## Virtual Constructor

*Everything buildable with the 'BAUFIX' kit can be built in virtual reality.*

Structural descriptions adapted dynamically:

- object descriptions are updated to comply with current situation
- make actual conceptualization available for dialogue



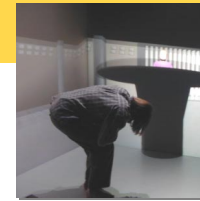( e.g., 'bar' gets to be 'tail unit' )

COAR representation formalism
*Artificial Intelligence Review 10(3-4),* 1996

## Artificial Intelligence & Virtual Reality



### Lab Research Mission

- AI methods used to establish an intuitive communication link between humans and multimedia

- Highly interactive Virtual Reality by way of multi-modal input and output systems (gesture, speech, gaze)

- Scientific enquiry and engineering of information systems closely interwoven (cognitive modeling approach)

---

## New lab inaugurated 15 July '02

… and Max

- 3-sided Cave-like display
- 6 D-ILA projectors
- Passive stereo, circular-polarisation filters
- Gesture tracking: marker-based infrared-camera system
- precision hand posture tracking by two wireless datagloves
- 8-channel spatial sound system

---

## Applications

**Embodied Conversational Agents** [Cassell et al. 2000]

Computer-generated characters that demonstrate human-like properties in „face-to-face" communication. Three aspects:

**Multimodal Interfaces**
with natural modalities like speech, facial displays, hand gestures, and body stance

**Software Agents**
that represent the computer in an interaction with a human or represent their human users in a digital environment (as „avatars")

**Dialog Systems**
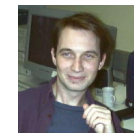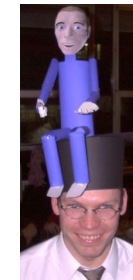where verbal as well as nonverbal devices advance the human-machine dialog

---

## „Brains"

Labor für Künstliche Intelligenz & Virtuelle Realität

**Mit »Max« wird der Computer menschlicher**

| Ipke Wachsmuth | Bernhard Jung | | Marc Latoschik | Timo Sowa |

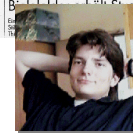| Ian Voß | Peter Biermann | Stefan Kopp | Alf Kranstedt | Nadine Leßmann |

## Computers

- 'Artabel Fleye 160' Linux Cluster for application and rendering
- 5 server nodes (double-Pentium III-class PCs)
- 8 graphic nodes (single-Pentium IV-class PCs) with NVIDIA GeForce 3 graphics
- nodes linked via 2GBit/s Myrinet-network for distributed OpenGL rendering

  – thanks for €s to DFG –



FLEYE
POWERED BY ARTABEL

---

## Agent MAX

*An artificial communicator situated in virtual reality*



Research into fundamentals of communicative intelligence:

- PHYSIS – the body system (especially gestures)
- COGNITION – the knowledge system
- EMOTION – the valuation system
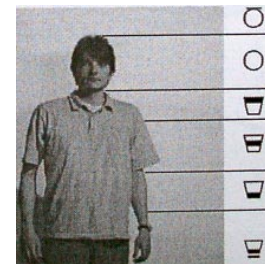
---

## PHYSIS: Articulated body



Hand animated by key framing
Body animated by model-based animation
Motion generators running concurrently and synchronized

Kinematic skeleton with 53 degrees of freedom (DOF) in 25 joints for the body and 25 DOF for each hand

---

## HamNoSys for gesture form description

(„Hamburg Notation System" – Institut für Deutsche Gebärdensprache, Hamburg)

| Symbol | ASCII-equivalent | Description |
|--------|------------------|-------------|
| ⊐ | BSifinger | indexfinger stretched |
| ▲ | EFinA | extended ahead |
| ○ | PalmL | palm orientated left |
| ⊔ | LocShoulder | location shoulder height |
| ⌐→ | LocStretched | fully stretched out |
| ↑ | MoveA | hand move ahead |
| → | MoveR | hand move right |
| <etc.> | ... | ... |
| ( ) | ( ) | executed in parallel |
| [ ] | [ ] | executed in sequence |

## Outlining a (roughly) rectangular shape

HamNoSys + movement constraints + timing constraints
(selected)          {STATIC, DYNAMIC}          {Start, End, Manner}
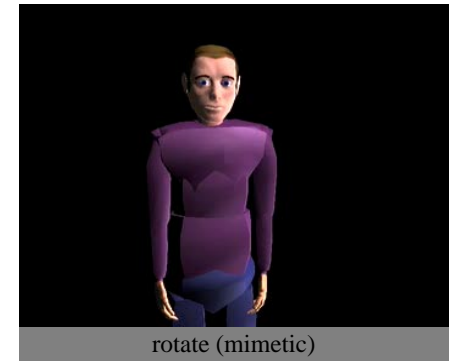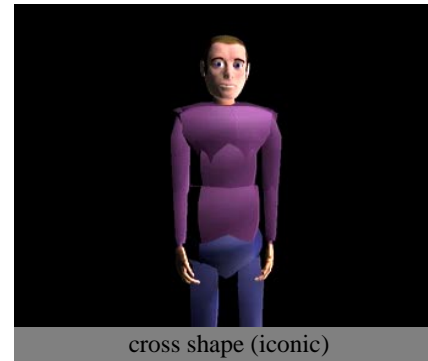
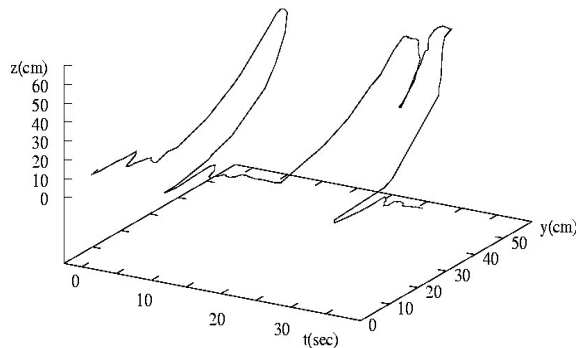**Articulated Communicator**

Gesture mappings

DrawRect
  (PARALLEL (Start 0.9, 0)(End 2.9, 0))
    (SEQUENCE (Start 0.9, 0)(End 2.9, 0))
      (PARALLEL (Start 0.9, 0)(End 1.8, 0))
        (DYNAMIC (Start 0.9, 0)(End 1.8, 0)(HandLocation ((LocShoulder LocCenter LocNorm)(LocShou
        (STATIC (Start 0.9, 0)(End 1.8, 0)(PalmOrientation (PalmD)))
      (PARALLEL (Start 1.9, 0)(End 2.1, 0))
        (DYNAMIC (Start 1.9, 0)(End 2.1, 0)(HandLocation ((LocShoulder LocLeftBeside LocNorm)(LocC
        (STATIC (Start 1.9, 0)(End 2.1, 0)(PalmOrientation (PalmR)))
      (PARALLEL (Start 2.2, 0)(End 2.9, 0))
    (STATIC (Start 0.9, 0)(End 2.9, 0)(HandShape (BSifinger)))
Handaction_1

---

## More form gestures

cross shape (iconic)          rotate (mimetic)
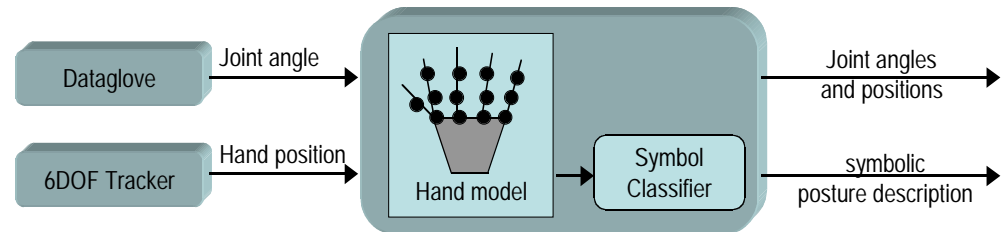
---

## Measuring gestures

**Segmentation cues**

- strong acceleration of hands, stopps, rapid changes in movement direction
- strong hand tension
- symmetries in two-hand gestures

z(cm) 60 50 40 30 20 10 0

y(cm) 50 40 30 20 10 0

t(sec) 0 10 20 30

---

## Analyzing gestures

Dataglove → Joint angle → Hand model

6DOF Tracker → Hand position → Hand model → Symbol Classifier

Hand model → Joint angles and positions

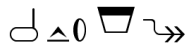Symbol Classifier → symbolic posture description

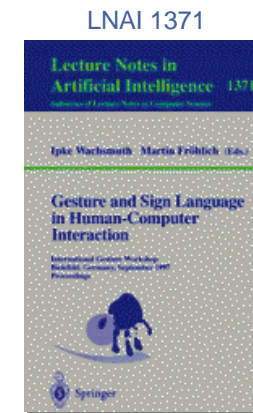Symbolic classification of gesture *shape* (HamNoSys)

# Gesture imitation game

- Human displays gestures, Max imitates them
- Parsing of gesture input: HamNoSys
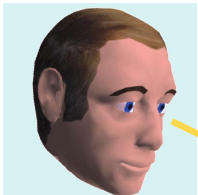- HamNoSys for specification of gesture output

Real time!

# 2 Gesture books (1998, 2002)

LNAI 1371

LNAI 2298
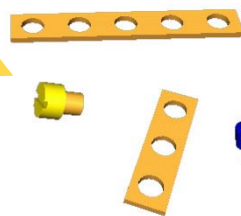
# COGNITION: Analyzing language

**I: Steck die gelbe Schraube in die lange Leiste.**

- **speech recognition**
- **syntactic-semantic parsing**
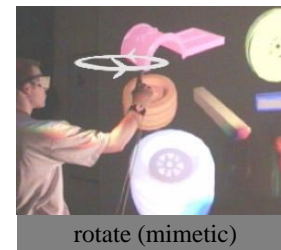- **reference to perceived scene**

Insert the yellow bolt into the long bar.

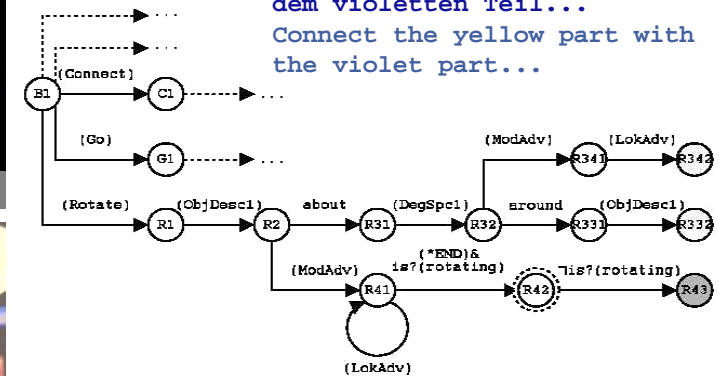| | | |
|---|---|---|
| steck | COMMAND | CONNECT |
| die | DET | |
| gelbe | COLOR | YELLOW |
| Schraube | OBJECTTYPE | BOLT |
| in | PREP | IN |
| die | DET | |
| lange | SIZE | LARGE |
| Leiste | OBJECTTYPE | BAR |

```
(select x (OBJECTTYPE(x)= BOLT and COLOR(x)= YELLOW))
(select y (OBJECTTYPE(y)= BAR  and SIZE(y) = LARGE))
```
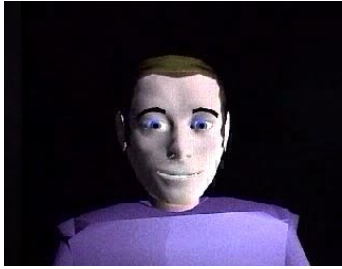
# Multimodal Analysis: tATN

**Verbinde das gelbe Teil mit dem violetten Teil...**
**Connect the yellow part with the violet part...**

point/select (deictic)

rotate (mimetic)

- Integration of speech and gesture
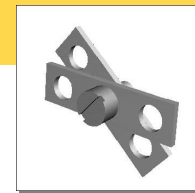- Interpretation in application context

## Lip-synchronous speech



Text-to-Speech:

TXT2PHO (IKP Uni Bonn), MBROLA
Phoneme transcription is the basis for automatic generation of visemes.

(Concept-to-Speech: TO DO)

Historical: Zemanek-Vocoder



- one viseme for *M, P, B*
- one viseme for *N, L, T, D*
- one viseme for *F, V*
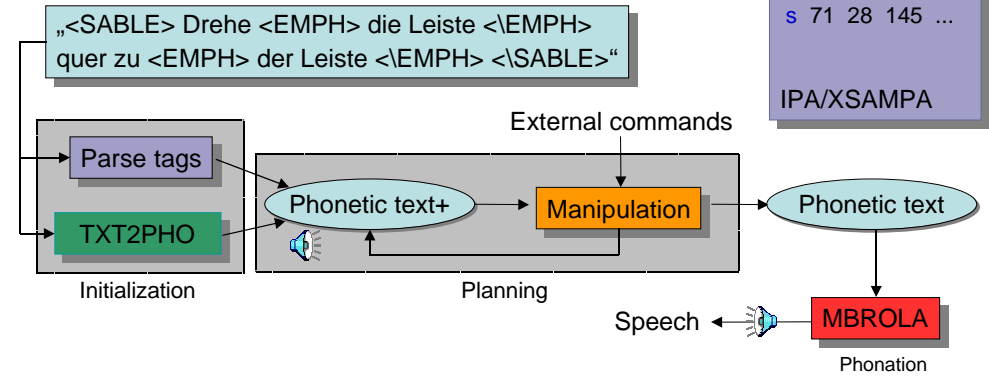- one viseme for *K, G*
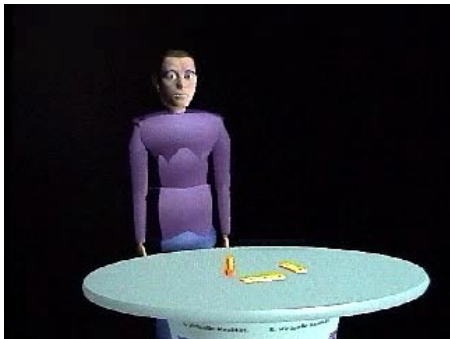- plus visemes for the vowels

## Speech and accentuation



🔊 Drehe die Leiste quer zu der Leiste.
Turn this bar crosswise to that bar.

🔊 Drehe *die* Leiste quer zu *der* Leiste.
Turn *this* bar crosswise to *that* bar.

Phonetic text:

s  105  18  176 ...
p  90  8  153
a:  104  4  150 ...
s  71  28  145 ...

IPA/XSAMPA

„<SABLE> Drehe <EMPH> die Leiste <\EMPH>
quer zu <EMPH> der Leiste <\EMPH> <\SABLE>"

External commands

Parse tags

TXT2PHO → Phonetic text+ → Manipulation → Phonetic text

Initialization     Planning

Speech ← MBROLA

Phonation

## Uttering speech and gesture



**And now take this bar and make it this big.**

MURML: XML-based markup language for multimodal utterance representations

```
<utterance>
  <specification>
    Und jetzt nimm <time id="t1"/> diese Leiste
    <time id="t2" chunkborder="true"/>
    und mach sie <time id="t3"/> so gross. <time id="t4"/>
  </specification>
  <behaviorspec id="gesture_1">
   <gesture>
    <affiliate onset="t1" end="t2"/>
    <constraints>
     <parallel>
      <static slot="HandShape" value="BSifinger"/>
      <static slot="ExtFingerOrientation"
          value="$object_loc_1" mode="pointTo"/>
      <static slot="GazeDirection" value="$object_loc_1"
          mode="pointTo"/>
     </parallel>
```
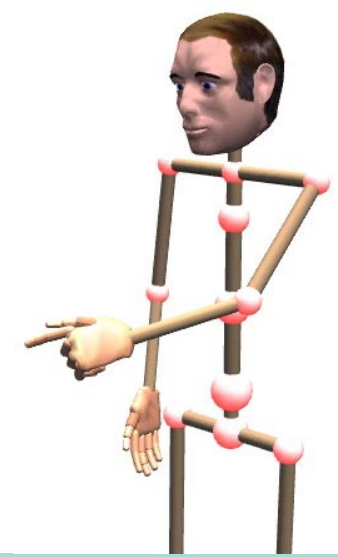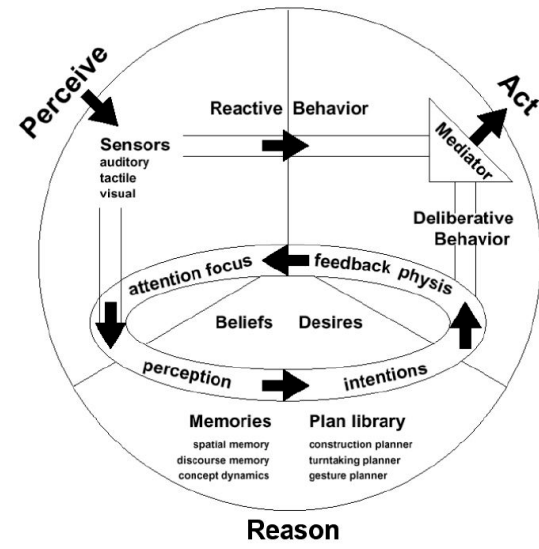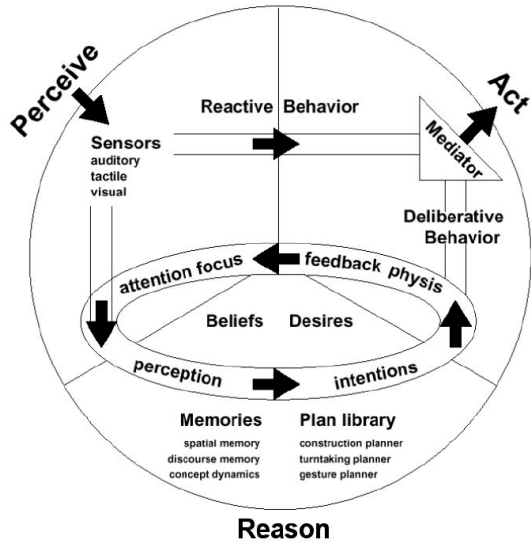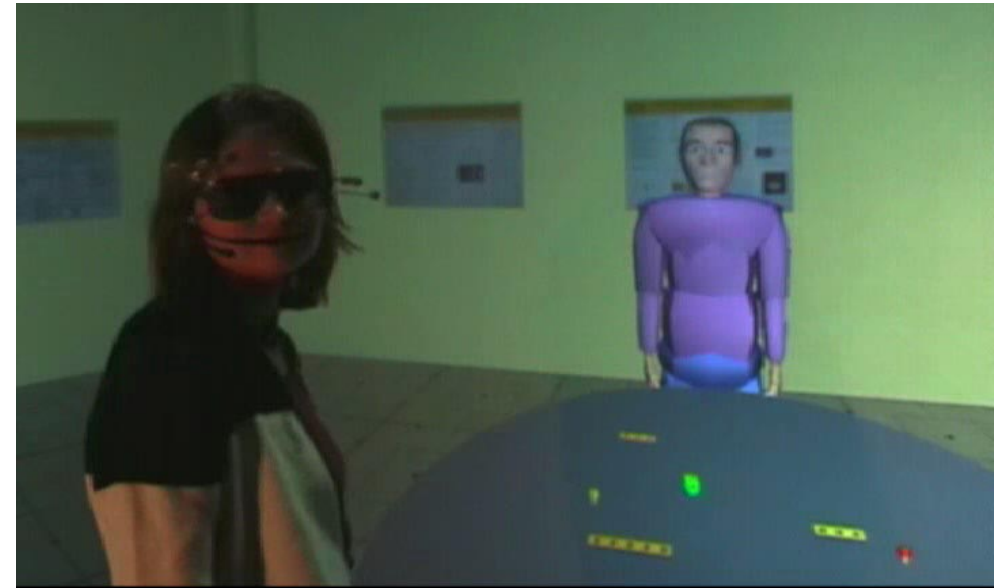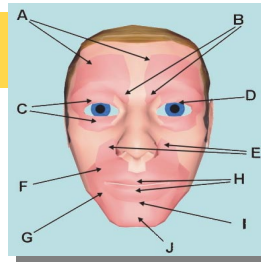
## Cognitively motivated architecture

Situated Artificial Communicators SFB 360
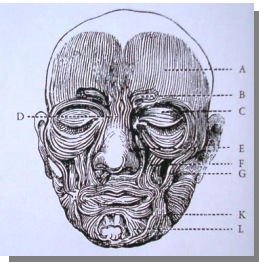
## Cognitively motivated architecture

- Perceive, Reason, Act running concurrently
- parallel processing by a reactive and a delibera-tive system
- information feedback in a cognitive loop
- BDI kernel with selfcon-tained dynamic planners
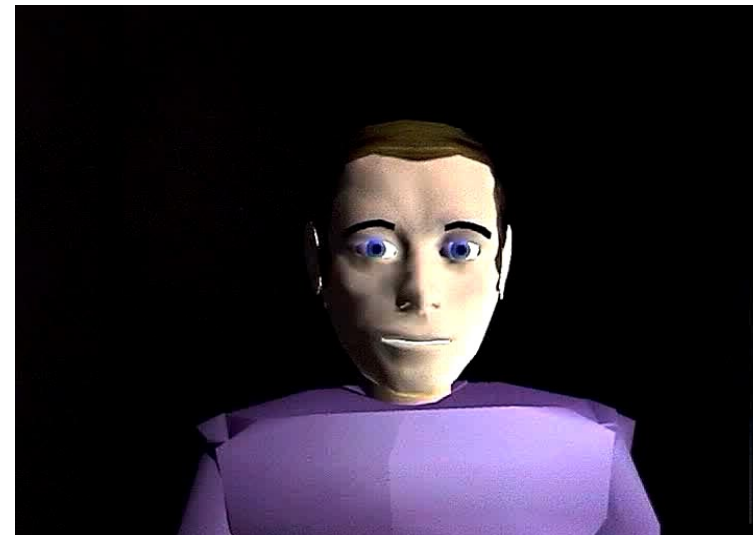- account for embodiment (physis) of the agent, multimodality

## Communicating with Agent Max…

## Facial expression



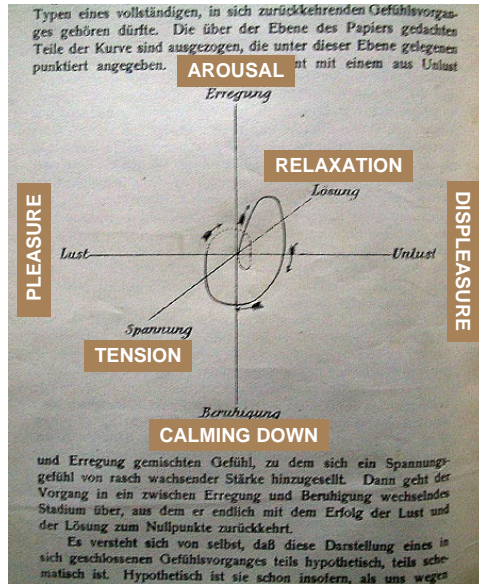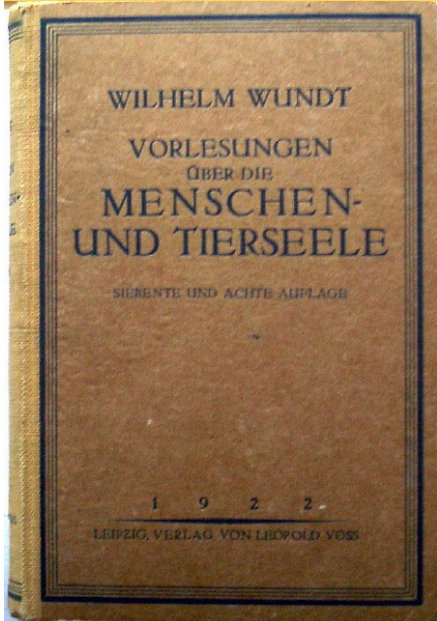| Muskeln des linken Bildes (von Sir Ch. Bell) | | Virtuelle Muskeln des rechten Bildes (MAX) | |
|---|---|---|---|
| A | Stirnmuskel | A | Stirnmuskel |
| B | Augenbrauenrunzler | B | Augenbrauenrunzler |
| C | Augenringmuskel | C | Augenringmuskel |
| D | Pyramidenmuskel der Nase | D | Augenlidmuskel |
| E | Heber der Oberlippe u d. Nasenflügels | E | Heber der Oberlippe u d. Nasenflügels |
| F | eigentlicher Lippenheber | F | Jochbeinmuskel u. Mundwinkelheber |
| G | Jochbeinmuskel | G | Mundwinkelherabzieher |
| K | Mundwinkelherabzieher | H | Ringmuskel des Mundes |
| L | Viereckiger Kinnmuskel | I | Unterlippenherabzieher |
| | | J | Unterkiefer |

## Expression of EMOTION



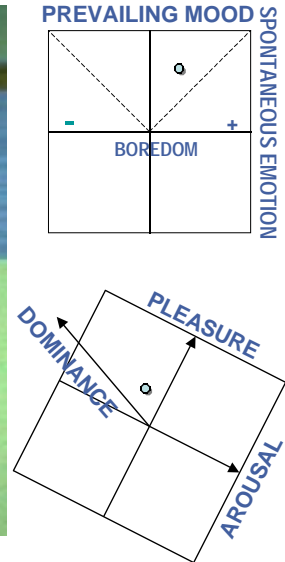Coordinated control of face muscles based on Action Units (Ekman/Friesen)

Student project (Körber Prize!)

Emotive system under development

## Wundt Emotion Dynamics

## Affect – Dynamic Emotion Space



EmoMax-Demo
25.02.2003

PREVAILING MOOD    SPONTANEOUS EMOTION
BOREDOM
PLEASURE
DOMINANCE
AROUSAL

## Embodied Communication



Anthropomorphic appearance
– humanoid body
– personality
– facial expression
– gesture
– spoken language
– emotional features

Intentionality
– knowledge / beliefs
– desires / motivations
– intentions
– commitments
– emotions...

e.g., BDI architecture ++
*(Beliefs - Desires - Intentions)*



THE END