

# Language Technology II: Language-Based Interaction

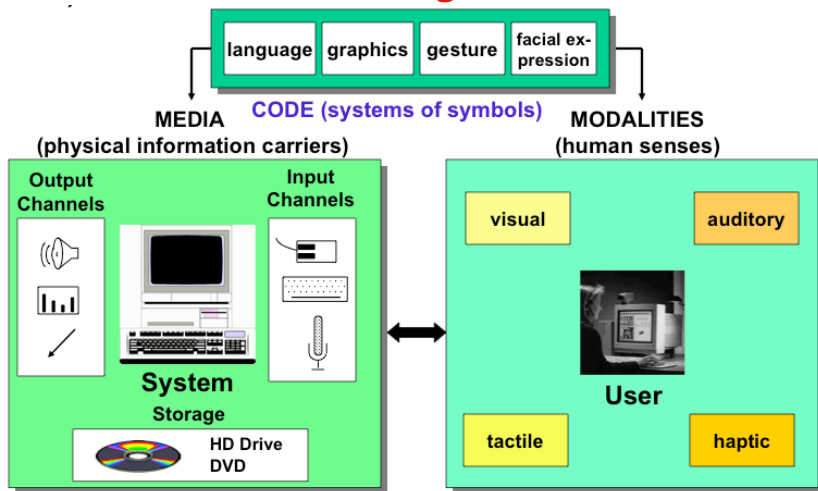
## Beyond Spoken Dialogue Systems

Ivana Kruijff-Korbayová  
korbay@coli.uni-sb.de  
www.coli.uni-saarland.de/courses/late2/

*I have reused slides from presentations of W. Wahlster, M. Johnston and J. Cassell*



## Interaction Using More Senses



(Wahlster, 2003)



## Outline

- Multimodality and multimedia
- Multimodal Dialogue Systems for Mobile Environments
- Embodied Conversational Agents
  - Talking heads
  - Virtual animated characters
  - Talking robots



## Multimodal and Multimedia Interfaces

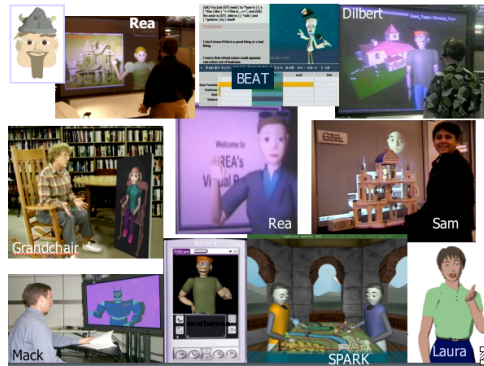
- A multimodal/multimedia interface supports more than one mode/medium through which the users and applications programs can communicate
  - E.g., natural language, graphics, sketching, animation, gestures, menus, video, sounds ...
- **Medium** = physical carrier of information
  - E.g., acoustic vs. optical
- **Mode** = particular sign system
  - E.g., language vs. graphics
- Examples:
  - Text + graphics output on display = **multimodal presentation**
  - Speech + text output on display = **monomodal multimedia presentation**
  - Speech + DataGlove input = **multimedia input**

(Wahlster, 2004)

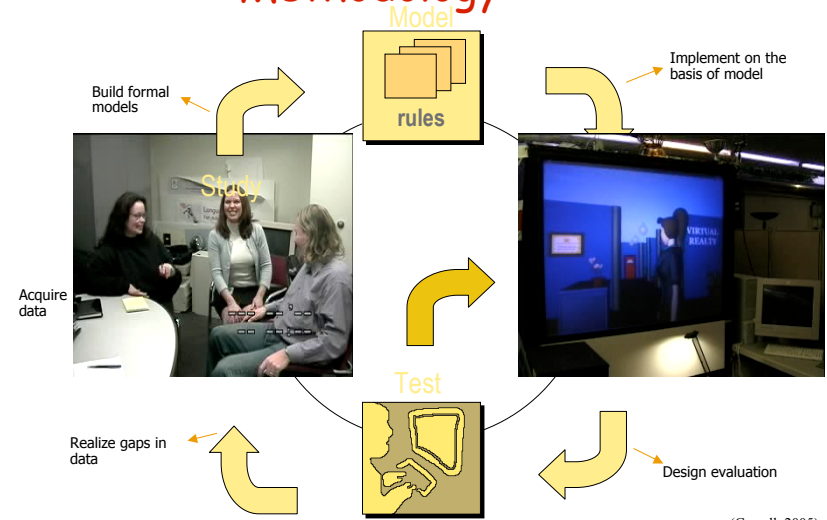


# Anthropomorphic Interfaces

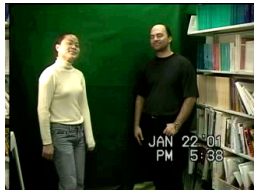
- = interfaces which have a "persona", i.e. at least a face or a whole body often also called Embodied Conversational Agents (ECA)
  - Talking heads
  - Virtual animated characters
- Added aspects of social interaction



# Methodology

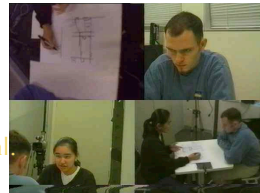


(Cassell, 2005)



Posture Shifts mark the beginning of new discourse segments (Cassell et al., '01)

Looks towards the listener indicate that further grounding is needed (Nakano, et al. '02)



Gestures are more likely to occur with rhematic material than thematic material (Cassell et al. '94)

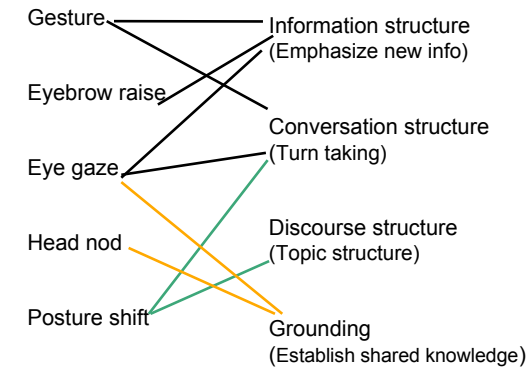
Small talk occurs before face-threatening discourse moves (Bickmore & Cassell, '02)



(Cassell, 2005)

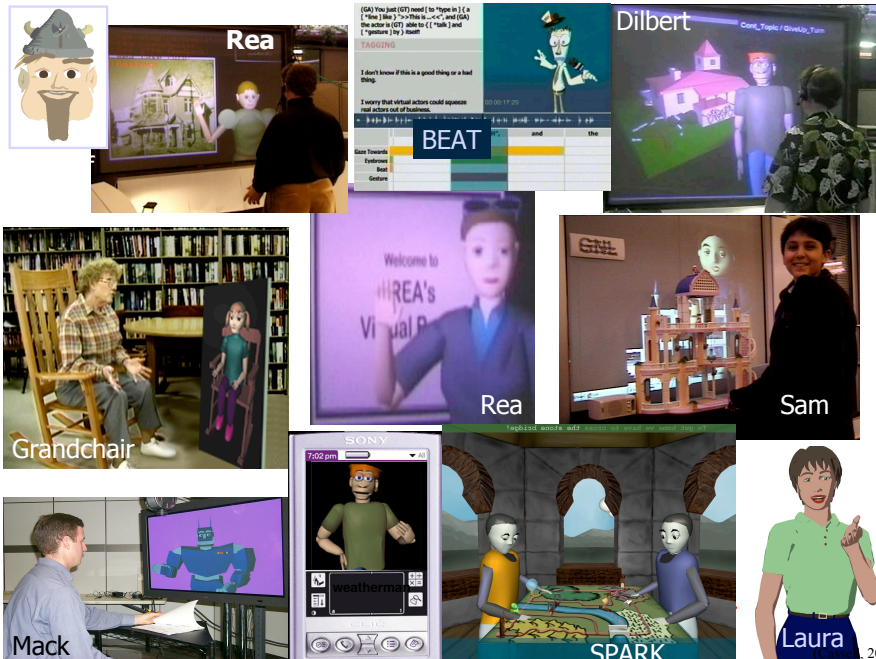


# Relationship between Linguistic Structure & Behavioral Cues

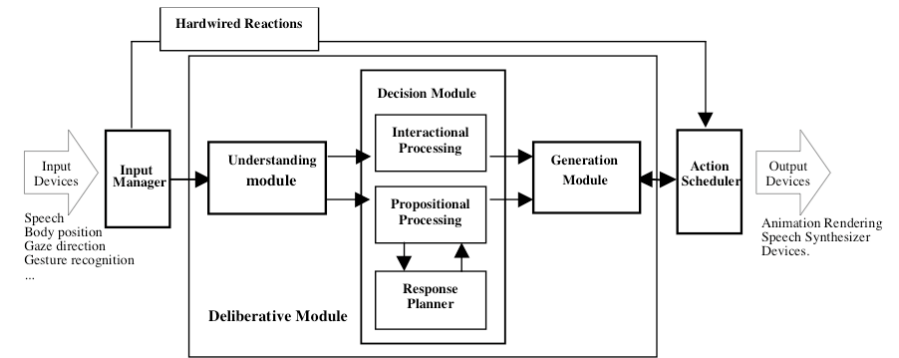


(Cassell, 2005)





## REA's Architecture



## Intelligent MM Interfaces

*Coexistence of input and output in different media and modes*  $\neq$  *Effective user interface*



## Limitations of Current MM Technology

- Although multimedia technology is so popular, that virtually everyone is using it in some form, the technology is still in its infancy:
  - Canned presentation segments
  - Hardwired presentation sequences
  - Very limited user adaptation
  - Inadequate media coordination
  - No unanticipated presentation possible
  - No follow-up questions possible
  - Very limited user interaction
  - No high-level authoring tools
  - No content-based reuse strategies
  - No inference services
  - No representation of the content of the presentation
  - No formal semantics/pragmatics of presentations

## MM Interfaces for Public Kiosks

- Since the introduction of ATMs in the 70's, public kiosks have been deployed to provide users with a broad range of information and services
- BUT: Majority have rigid system-initiative graphical interfaces with user input by touch or keypad
  - Thus, they can only support simple tasks for able-bodied users
- To support more complex tasks for a broader range of users, kiosks will need to provide a more flexible and natural user interface
  - MM interfaces provide naturalness and flexibility
  - E.g., Gustafson et al. 1999 (August), Narayan et al. 2000 (MVPQ), Raisamo 1998, Lamel et al. 2002 (MASK), Wahlster et al. 2003 (SmartKom Public), Cassell et al. (MACK)

(Johnston, 2004)



## Burning Research Issues in Multimodal Dialogue Systems

- **Multimodality:** from alternate modes of interaction towards mutual disambiguation and synergistic combinations
- **Discourse Models:** from information-seeking dialogs towards argumentative dialogs and negotiations
- **Domain Models:** from closed world assumptions towards the open world of web services
- **Dialog Behaviour:** from automata models towards a combination of probabilistic and plan-based models



## Intelligent MM Interfaces

*Coexistence of input and output in different media and modes*  $\neq$  *Effective user interface*

- From alternate modes of interaction to **composite** multimodality
- Careful coordination of different media and modes in a coherent and cooperative dialogue is required



## Composite Multimodality: Input

- Composite input:
  - Enabling users to provide a single contribution (turn) which is optimally distributed over the available input modes e.g., speech + ink "zoom in here"
- Motivation
  - Naturalness
  - Certain kinds of content within a single communicative act are best suited to particular modes, e.g.,
    - Speech for complex queries or constraints, reference to objects currently not visible or intangible
    - Ink/gesture for selection, indicating complex graphical features
  - Empirical studies
    - Task performance and user preference advantages (Oviatt et al. 1999)
    - Compensation for errors (Oviatt 1999; Bangalore & Johnston n 2000)

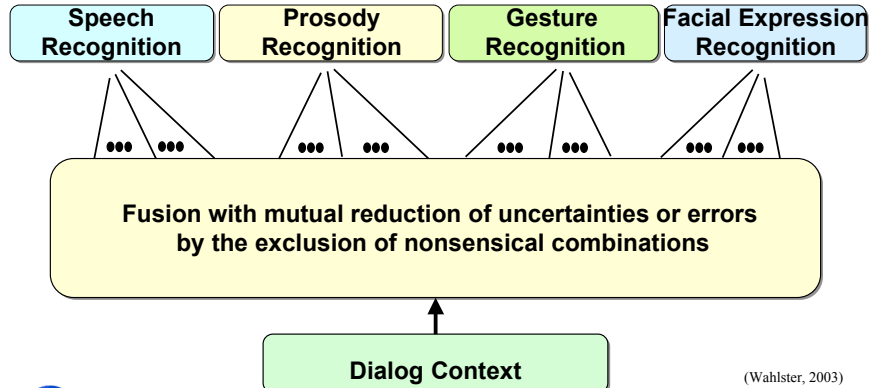


(Johnston, 2004)



# Composite Multimodality: Input Fusion

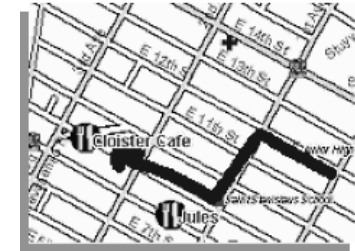
Mutual disambiguation and synergistic combinations: semantic fusion of multiple modalities in dialog context helps to reduce ambiguity and errors



(Wahlster, 2003)

# Composite Multimodality: Output

- Composite output:
  - Allowing for system output to be optimally distributed over the available output modes, e.g.,
    - High level summary in speech, details in graphics: "Take this route across town to the Cloister Café"
    - Multimodal help providing examples for the user: "To get the phone number for a restaurant, circle one like *this* and say or write *phone*." (Hastie et al. 2002)
  - Output should be dynamically tailored to be maximally effective given the situation and user preferences
- Same motivation as for multimodal input

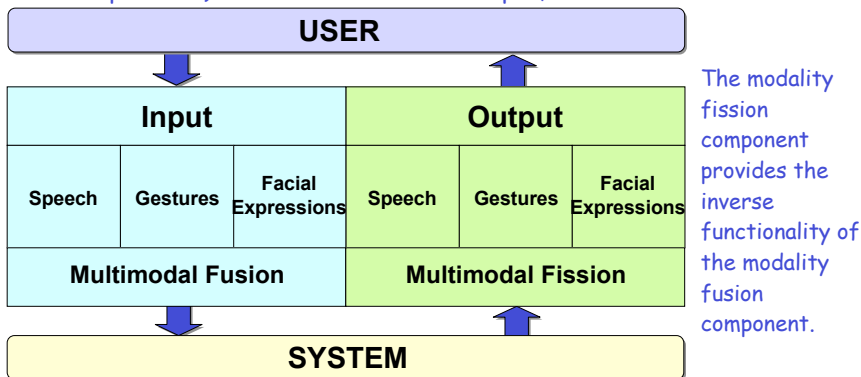


(Johnston, 2004)



# Full Symmetric Multimodality

Symmetric multimodality means that all input modes (speech, gesture, facial expression) are also available for output, and vice versa.



**Challenge:** A dialogue system with symmetric multimodality must not only understand and represent the user's multimodal input, but also its own.

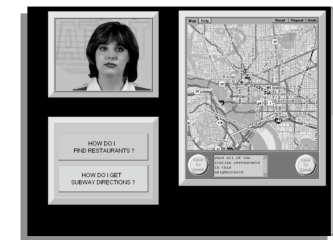
(Wahlster, 2003)



# MATCH:

## Multimodal Access to City Help

- Interactive city guide and navigation for information-rich urban environments
  - Finding restaurants and points of interest, getting info, subway routes for New York and Washington, D.C.
- Composite input and output
  - Speech, ink, graphics
- Mobile (standalone on a tablet or distributed WLAN)
- MATCHkiosk (deployed at AT&T visitor center in DC)
  - Social interaction
  - Also printed output




(Johnston, 2004)



# MATCH

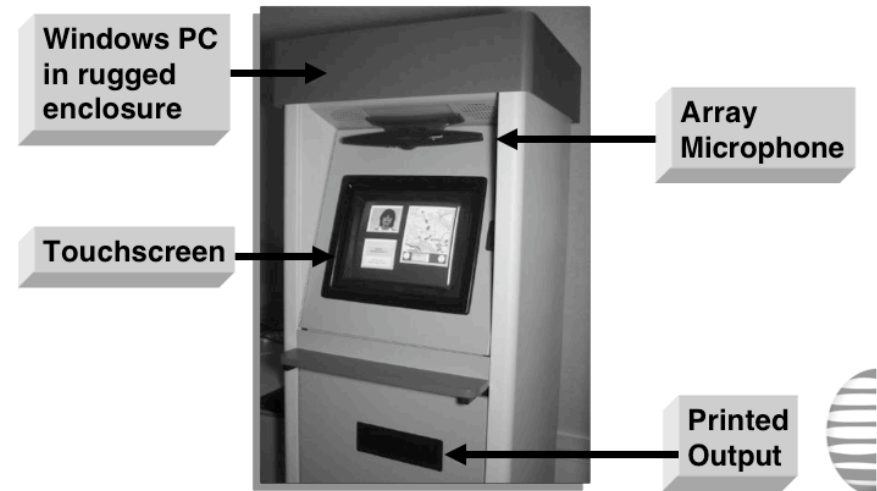


- Finding restaurants
  - Speech: "show inexpensive italian places in chelsea"
  - Multimodal: "cheap italian places in this area"
  - Pen: 
- Getting info: "numbers for these"
- Subway routes: "how do I get here from Broadway and 95th street"
- Pan/zoom map: "Zoom in here"

(Johnston, 2004)



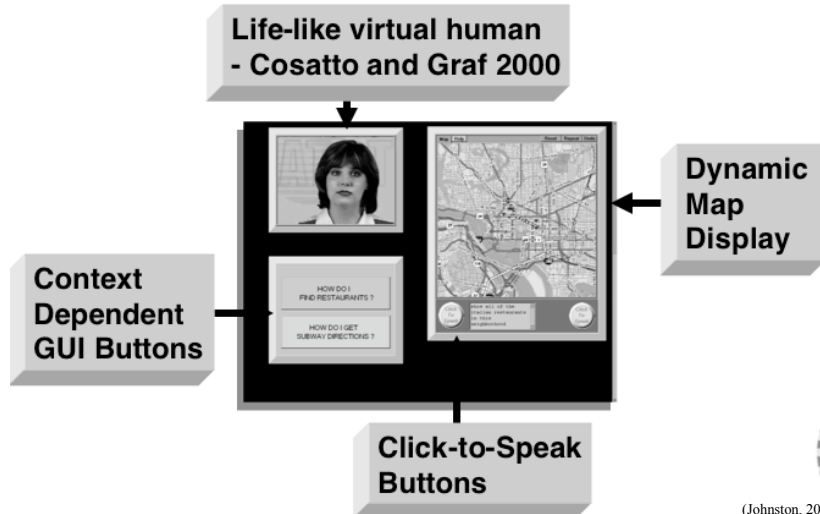
# MATCHKiosk



(Johnston, 2004)



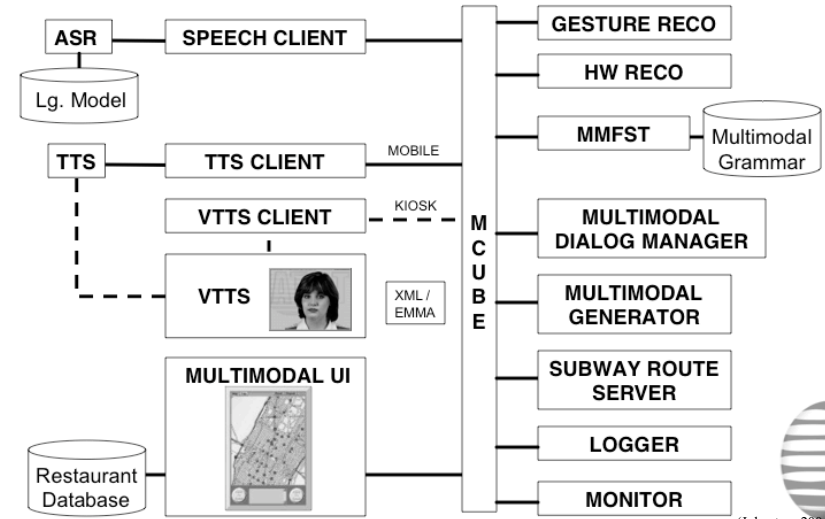
# MATCHKiosk



(Johnston, 2004)



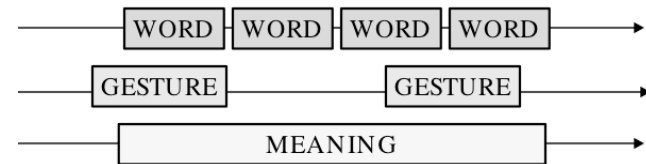
# MATCH Architecture



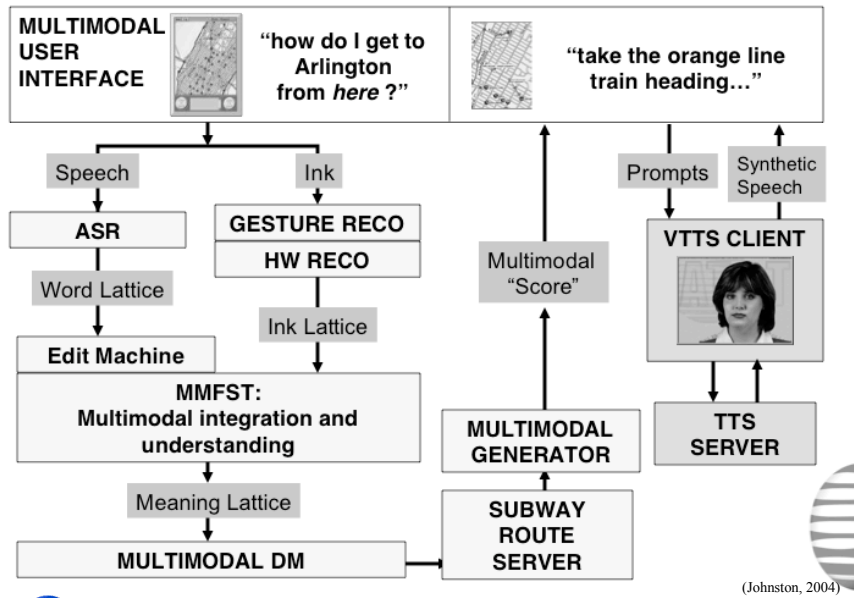
(Johnston, 2004)



# Multimodal Understanding



- Associate word sequence + gesture sequence with meaning
  - Late integration: first compute meaning of word sequence and meaning of gesture sequence, then "merge" the meanings, e.g., (Johnston 1998), SmartKom: overlay (Alexandersson&Becker 2001; Pflieger 2002)
  - Early integration: compute meaning of a composite word+gesture sequence: MMFST (Johnston&Bangalore 2000,2004)
    - Enables compensation for errors in individual modes

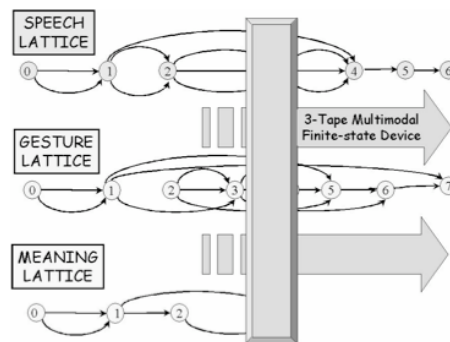


(Johnston, 2004)



## MMFST: Early Multimodal Integration

- Speech and gesture parsing, multimodal integration, and understanding in single MM grammar model
  - (Johnston&Bangalore 2000,2004)
  - Compiled to efficient finite state device
    - G:W transducer aligns speech and ink
    - G\_W:M transducer takes a composite alphabet of speech and gesture symbols and outputs meaning
  - Compiled from a declarative multimodal CFG (terminals are triples W:G:M = Words:Gestures:Meaning)
- Robust, efficient
- Enables compensation for errors



(Johnston, 2004)



## MMFST: Example

"Phone numbers for these three restaurants"



S → eps:eps:<cmd> CMD eps:eps:</cmd>  
 CMD → phone:eps:<phone> numbers:eps:eps  
 for:eps:eps DEICTICNP  
 eps:eps:</phone>  
 DEICTICNP → DDETPL eps:area:eps eps:selection:eps  
 NUM RESTPL eps:eps:<restaurant>  
 eps:SEM:SEM eps:eps:</restaurant>  
 DDETPL → these:G:eps  
 RESTPL → restaurants:restaurant:eps  
 NUM → three:3:eps

Johnston et al. (2004)



## User-Tailored Generation

"compare these restaurants"

- User-tailored summaries, comparisons or recommendations can be generated using a model of user preferences



*Compare-A:* Among the selected restaurants, the following offer exceptional overall value. Uguale's price is 33\$. It has excellent food quality and good decor. Da Andrea's price is 28\$. It has very good food quality an good decor. John's Pizzeria's price is 20\$. It has very good food quality and mediocre decor.

*Compare-B:* Among the selected restaurants, the following offer exceptional overall value. Babbo's price is 60\$. It has superb food quality. Il Mulino's price is 65\$. It has superb food quality. Uguale's price is 33\$. It has excellent food.

Johnston et al. (2004)



## The Need for Going Mobile

Broadband mobile Internet access technologies via UMTS or mobile hotspots pave the way for a wide spectrum of added-value web services.



but: the user must input more and more complex commands to specify his information needs.

PDA's and smartphones with tiny keyboards and mice are useless for mobile settings.

➔ Multimodal Dialogue Systems for Mobile Systems

(Wahlster, 2003)



## Multimodal Mobile Devices

Today's Cell Phone

Third Generation UMTS Phone



**Verbmobil**  
Speech only

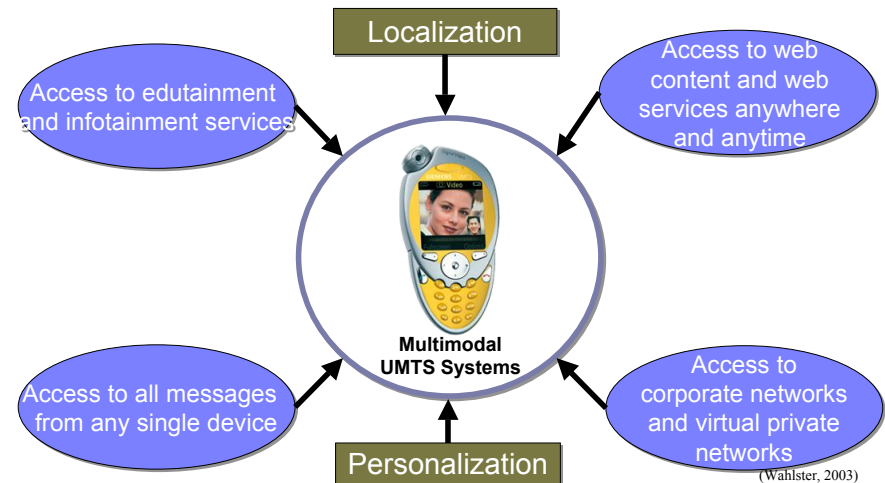


**SmartKom**  
Speech, Graphics and Gesture

(Wahlster, 2003)



## Intelligent Interaction with Mobile Internet Services



(Wahlster, 2003)



## Mobile Access to an Edutainment System



Mobile Dialogue with a Virtual Tourist Guide for the Heidelberg Castle

Location-adaptive Query Interpretation



Cooperation: DFKI, EML, DFG, FhG

(Wahlster, 2003)



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

33

## Spoken Dialogues with the Car Navigation System



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

34

## Pervasive Computing

- conversational access to information and services anytime / anywhere
- event-driven interaction rather than traditional "well-formed" dialogues
  - task-switching
  - interruption
- deployment in robust environments
- adaptive and multimodal
  - needs of user
  - different environments



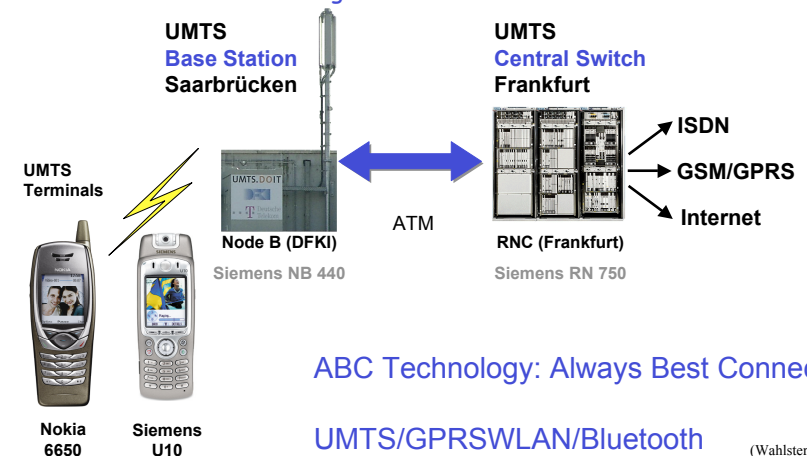
6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

35

## Going Mobile

DFKI operates the first German Test Site for Innovative 3G/4G Broadband Mobile Technologies



ABC Technology: Always Best Connected

UMTS/GPRS/WLAN/Bluetooth

(Wahlster, 2003)



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

36



(Wahlster, 2003)

## CityShopper's Added-Value Mobile Service

Please let me know, when I pass a shop selling batteries.

- CityShopper sends a note to the user or activates an alarm as soon as the user approaches a shop that offers of an item on his active shopping list.
- CityShopper's spatial alarm can be combined with:
  - route planning and navigation
  - temporal and spatial optimization of the shopping tour



(Wahlster, 2003)



## Spatial Alarm by CityShopper

A Location-based Active Shopping List based on UMTS  
(Cooperation of DFKI with Telekom)

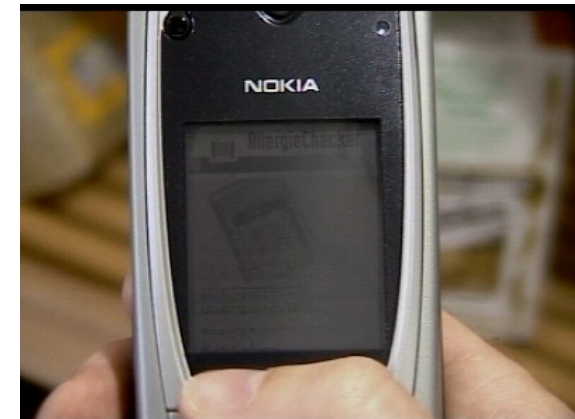


(Wahlster, 2003)



## Mobile AllergyChecker

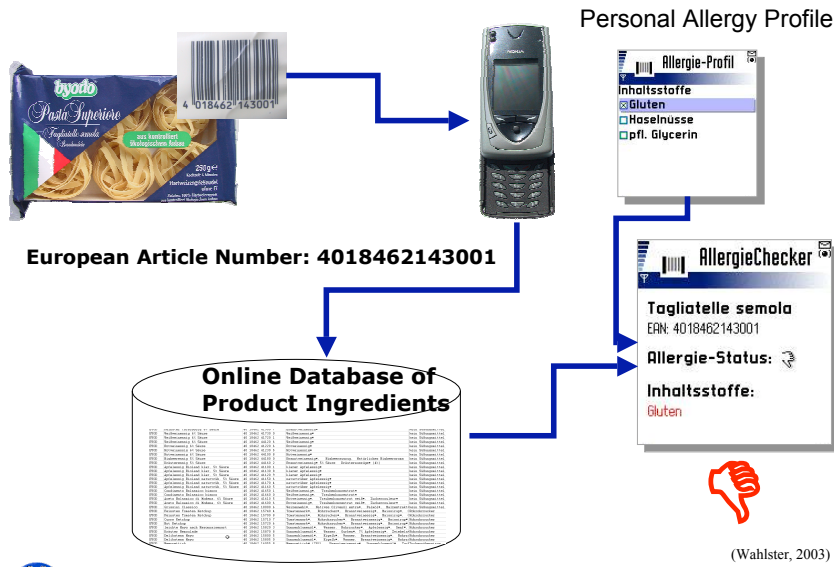
Mobile Bar Code Analyzer and Broadband Personalized Allergy Checker  
(cooperation of DFKI with Mineway)



(Wahlster, 2003)

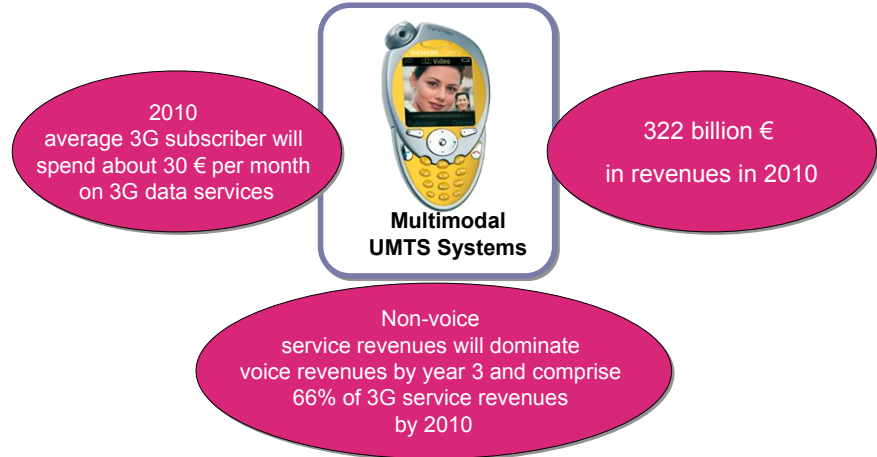


Real-Time Video-based Product Identification and Personalized Allergy Check

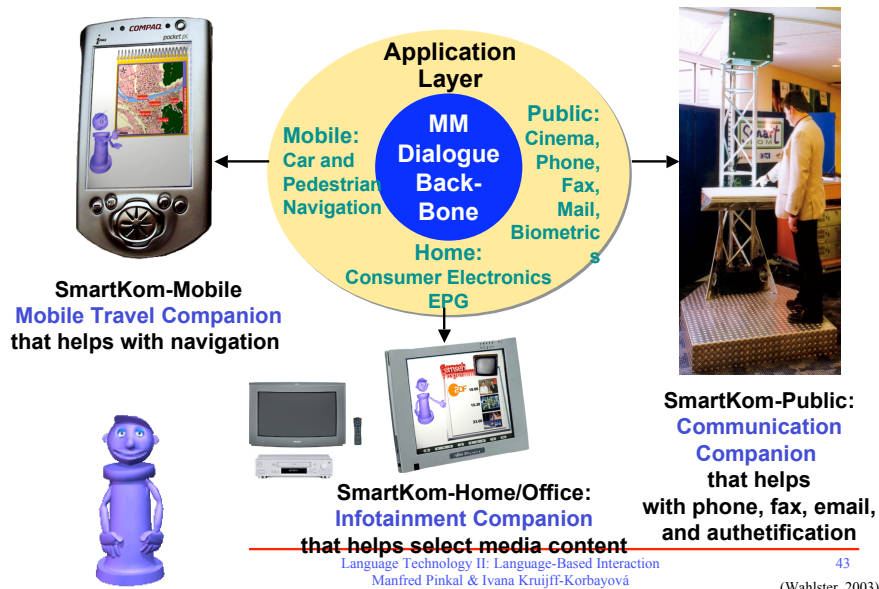


Market studies predict:

Cumulative revenues of almost 1 trillion € from launch of mobile 3G services until 2010



SmartKom



SmartKom



SmartKom-Home: DFKI Collaboration with Sony und Philips

# SmartKom-Public (Kiosk Version)

Multimodal Control of TV-Set

Multimodal Control of VCR/DVD Player

3 dual Xeon 2.8 Ghz processors with 1.5 GB main memory

Infrared Camera for Gestural Input, Tilting CCD Camera for Scanning, Video Projector

Microphone

Camera for Facial Analysis

Projection Surface

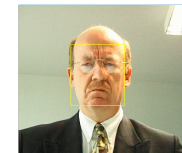
Speakers for Speech Output

(Wahlster, 2003)

# Using Facial Expression Recognition for Affective Personalization

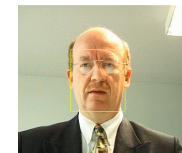
Processing ironic or sarcastic comments  
 (1) Smartakus: Here you see the CNN program for tonight.

(2) User: *That's great.*



(3) Smartakus: I'll show you the program of another channel for tonight.

(2') User: *That's great.*



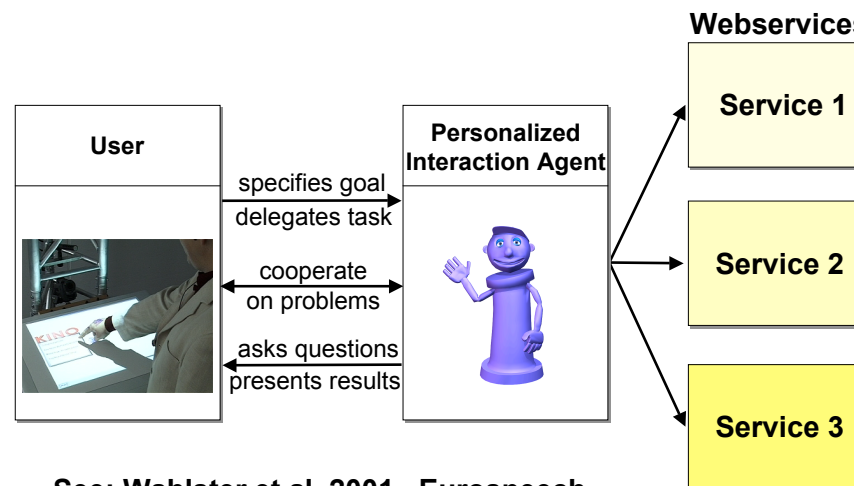
(3') Smartakus: Which of these features do you want to see? (Wahlster, 2003)

# SmartKom

	Input by the User	Output by the Presentation agent
Speech	+	+
Gesture	+	+
Facial Expressions	+	+

(Wahlster, 2003)

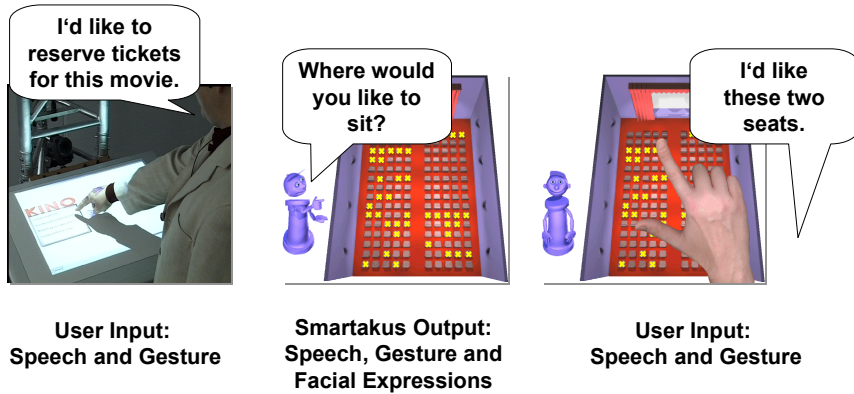
# SmartKom



See: Wahlster et al. 2001 , Eurospeech

(Wahlster, 2003)

# SmartKom



User Input:  
Speech and Gesture

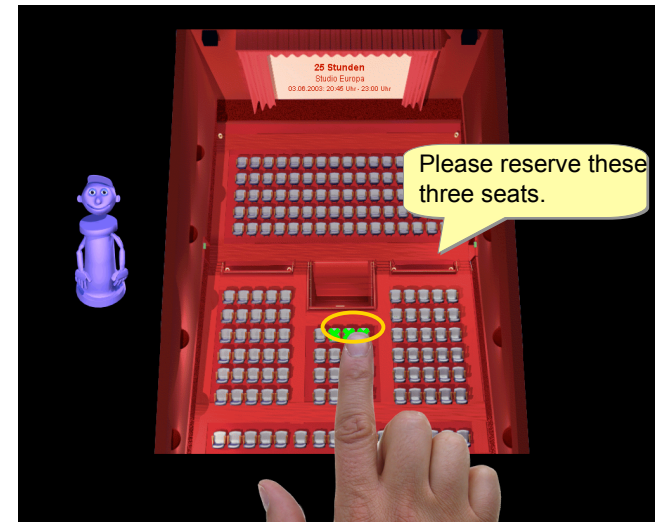
Smartakus Output:  
Speech, Gesture and  
Facial Expressions

User Input:  
Speech and Gesture

(Wahlster, 2003)



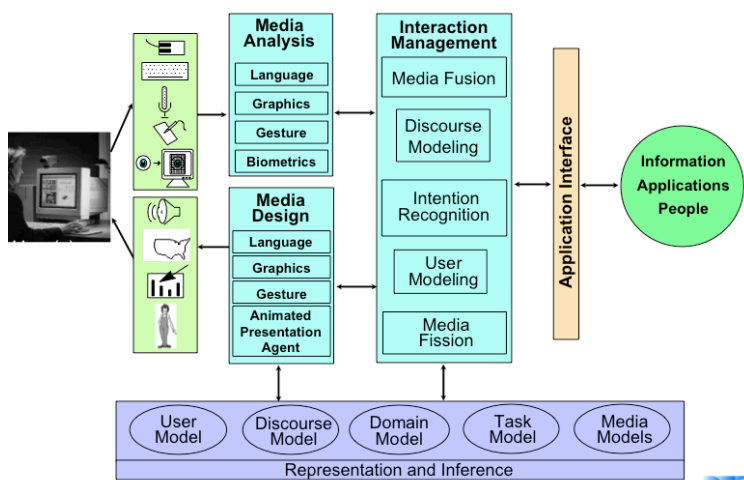
# SmartKom



(Wahlster, 2003)



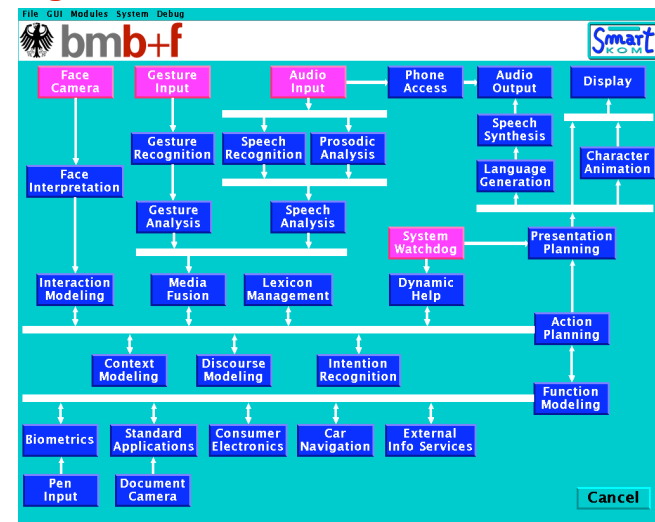
# SmartKom Architecture



© W. Wahlster



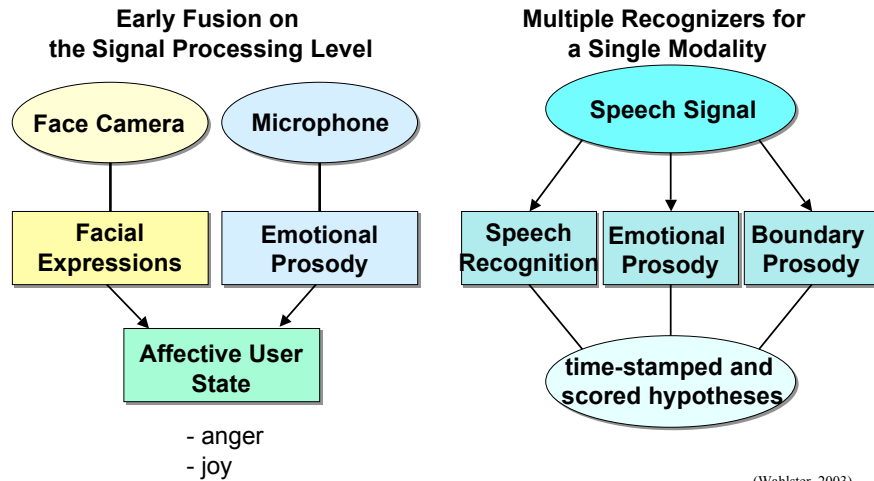
# The High-Level Control Flow of SmartKom



(Wahlster, 2003)



# Fusing Symbolic and Statistical Information in SmartKom



# Computational Mechanisms for Modality Fusion and Fission

Modality Fusion	Modality Fission
Unification	Planning
Overlay Operations	Constraint Propagation
<b>Ontological Inferences</b>	
<b>M3L: Modality-Free Semantic Representation</b>	

# M3L Representation of an Intention

**Lattice Fragment** I would like to know more about this

```

<intentionLattice>
[... ]
<hypothesisSequences>
<hypothesisSequence>
<score>
<source> acoustic </source>
<value> 0.96448 </value>
</score>
<source> gesture </source>
<value> 0.99791 </value>
</score>
<source> understanding </source>
<value> 0.91667 </value>
</score>
<hypothesis>
<discourseStatus>
<discourseAction> set </discourseAction>
<discourseTopic><goal> epg_info </goal></discourseTopic>
[... ]
<event id="dim868">
<informationSearch id="dim869">
<pieceOfInformation>
<broadcast id="dim863">
<avMedium>
<avMedium id="dim866">
<avType> featureFilm </avType>
<title> Enemy of the State </title>
[... ]
</pieceOfInformation>
</informationSearch>
[... ]
</hypothesisSequence>
[... ]
</hypothesisSequences>
</intentionLattice>
    
```

Confidence in the Speech Recognition Result (acoustic 0.96448)

Confidence in the Gesture Recognition Result (gesture 0.99791)

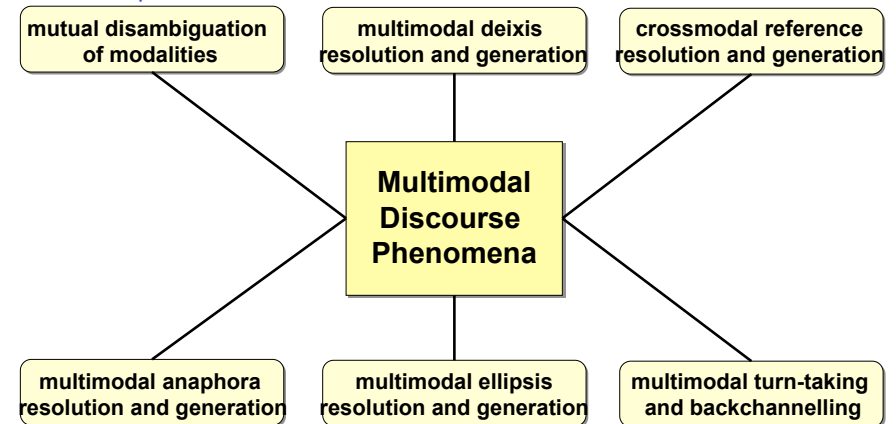
Confidence in the Speech Understanding Result (understanding 0.91667)

Planning Act (informationSearch id="dim869")

Object Reference (title: Enemy of the State)

# Multimodal Discourse Phenomena

Symmetric multimodality is a prerequisite for a principled study of many discourse phenomena.



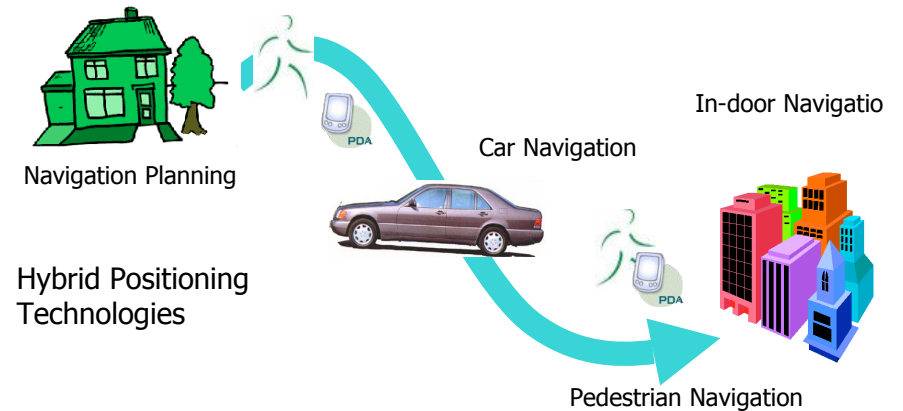


# Adaptive Perceptual Feedback on the System State



(Wahlster, 2003)

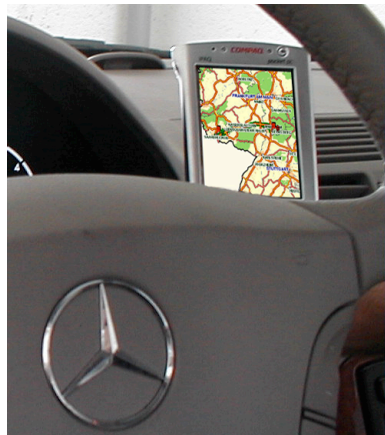
# Seamless Navigation Services



(Wahlster, 2003)

# Getting Driving and Walking Directions via SmartKom

SmartKom can be used for Multimodal Navigation Dialogues in a Car



**User:** I want to drive to Heidelberg.

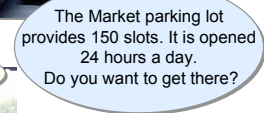
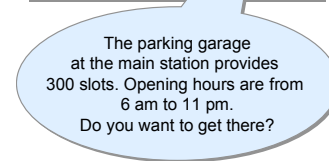
**Smartakus:** Do you want to take the fastest or the shortest route?

**User:** The fastest.

**Smartakus:** Here you see a map with your route from Saarbrücken to Heidelberg.

(Wahlster, 2003)

# Spoken Navigation Dialogues with SmartKom



(Wahlster, 2003)

## Salient Characteristics of SmartKom

- **Seamless integration** and **mutual disambiguation** of multimodal input and output on semantic and pragmatic levels
- Situated understanding of possibly **imprecise, ambiguous, or incomplete** multimodal input
- **Context-sensitive interpretation** of dialog interaction on the basis of **dynamic discourse and context models**
- Adaptive generation of **coordinated, cohesive and coherent** multimodal presentations
- Semi- or fully automatic **completion of user-delegated tasks** through the integration of information services
- **Intuitive personification** of the system through a presentation agent

(Wahlster, 2003)



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

65

## SmartKom Keypoints

- Various types of unification, overlay, constraint processing, planning and ontological inferences are the fundamental processes involved in SmartKom's **modality fusion and fission** components.
- The key function of modality fusion is the reduction of the overall uncertainty and the mutual disambiguation of the various analysis results based on a **three-tiered representation of multimodal discourse**.
- A multimodal dialogue system must not only understand and represent the user's input, but its **own multimodal output**.

(Wahlster, 2003)

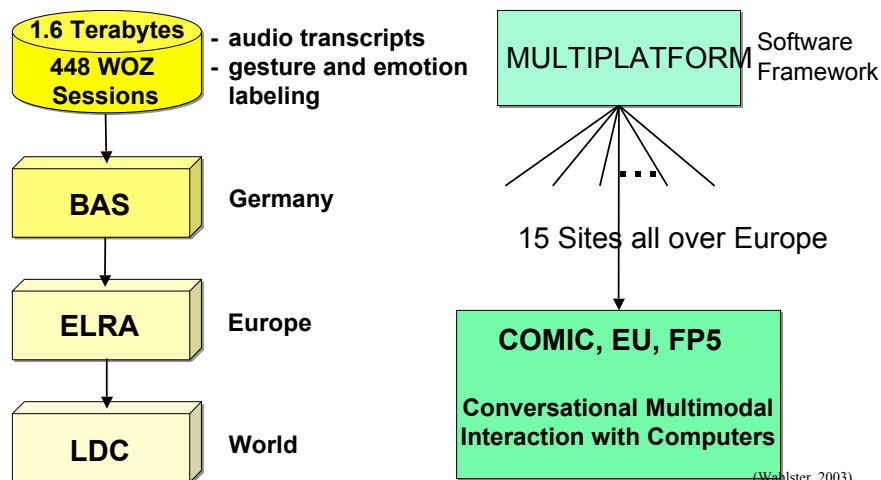


6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

66

## SmartKom's Impact on Software Tools and Resources for Research on Multimodality



(Wahlster, 2003)



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

67

## Ambient Intelligence

=  
**Embedding Artificial Intelligence in Everyday Objects and Environments**  
Key Characteristics of Ambient Intelligence:

- Embedded** Many networked devices are integrated into the environment
- Situated** These devices can recognize situational context
- Personalized** They can be tailored towards your needs and affects
- Adaptive** They can change in response to you and your task
- Pro-active** They try to anticipate your plans and intentions

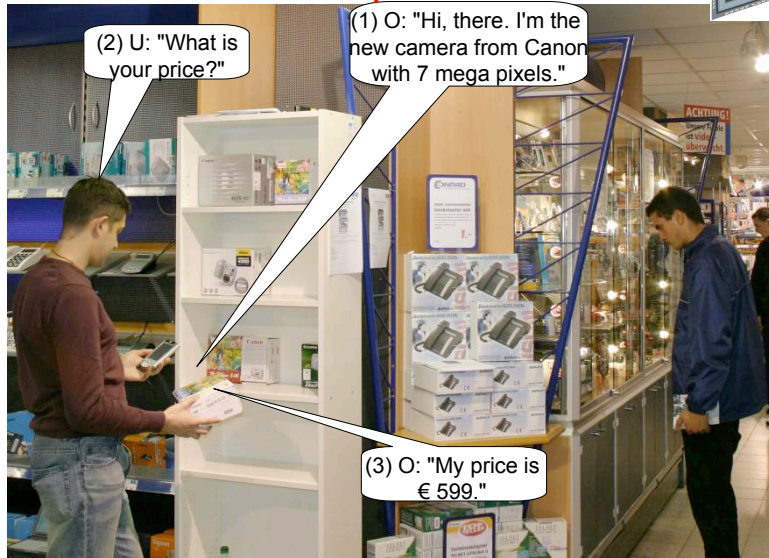


6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

68

## Mobile Shop Assistant



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

69

## Automatic Comparison Shopping with Ambient Intelligence



### Multimodal Fusion of Speech, Gesture and Physical Actions (Intra- and Extra-Gestures)



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

70

## Conversational Robots

- Some examples:
  - WITAS (CSLI)
    - Unmanned aerial vehicle: start, land, locate & follow objects
    - Collaborative activities, multitasking
  - Godot (Edinburgh)
    - Navigation, learning natural language descriptions of objects
  - MEL (MERL)
    - Hosting
    - Engagement behaviors (initiate, maintain and disengage)
    - Modeling of look-tracking (head movement)
  - COSY project (DFKI)
- Language + vision + other sensory input + gestures
- Activities situated in physical world



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

71

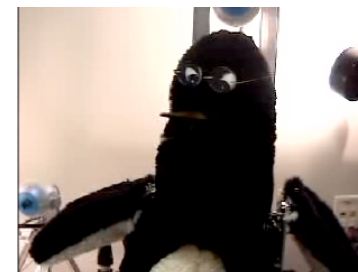
## Conversational Robots



WITAS



GODOT



MEL



6/28/05  
Beyond Spoken...

Language Technology II: Language-Based Interaction  
Manfred Pinkal & Ivana Kruijff-Korbayová

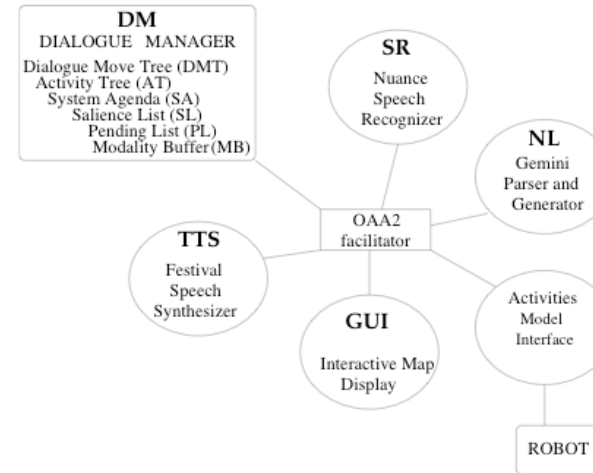
72

## WITAS

Multi-modal Utterances	Dialogue Move
Operator (O) : <i>Our job is to look for a red car</i>	Command (Joint Activity)
UAV (U) : <i>Ok. I am looking for one.</i>	Report (Confirm Activity)
O : <i>Fly here please [+click on map]</i>	Command (Deictic)
U : <i>Okay. I will fly to waypoint one</i>	Report (Confirm Activity)
U : <i>Now taking off and flying there.</i>	Report (Current Activity)
O : <i>Stop that. Go to the tower instead.</i>	Command, Revision
U : <i>I have cancelled flying to waypoint one. I will fly to the tower.</i>	Report (Activity status)
O : <i>What are you doing ?</i>	Wh-question (Current Activity)
U : <i>I am searching for a red car and flying to the tower</i>	Answer (Current Activity)
O : <i>What will you do next ?</i>	Wh-question (Planned Activity)
U : <i>I have nothing planned.</i>	Answer(Planned Activity)
U : <i>I see a red car on main street [display on map, show video images], Is this the right car ?</i>	Report, Yn-question (Activity)
O : <i>Yes, that's the right car</i>	Yn-answer (Positive)
U : <i>Okay. I am following it .</i>	Report (Current activity)



## WITAS

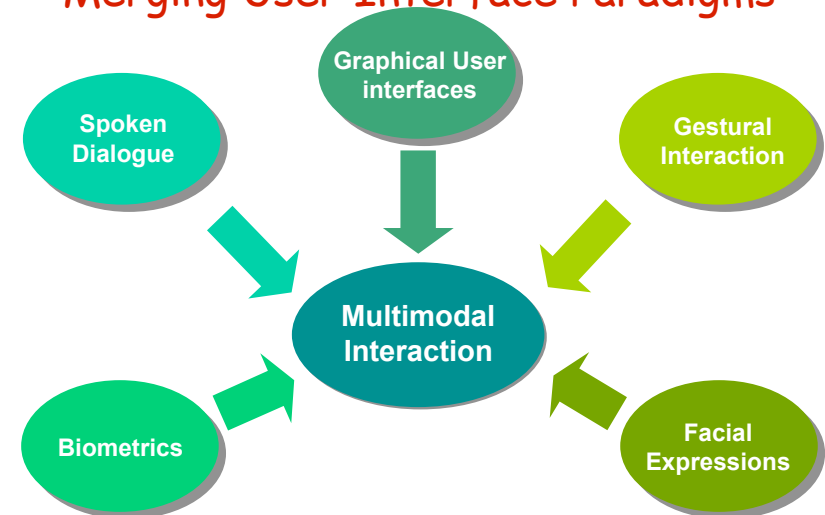


## Wrapping Up

- Multimodal and anthropomorphic interfaces
- Symmetric multimodality
- Composite multimodality
  - Fusion
  - Fission
- Adaptivity, personalization, collaboration (activity & pro-activity)
- Anthropomorphic objects
- Robots



## Merging User Interface Paradigms



(Wahlster, 2003)



## References

- M. Johnston et al. "MATCH: An architecture for Multimodal Dialogue Systems." In Proc. Of the 40th Annual Meeting of ACL. pp. 376-383. 2002.
- N. Pflieger et al. "Robust Multimodal Discourse Processing." In Proc. Of DiaBruck. pp. 107-114. 2003.  
SmartKom website: <http://www.smartkom.org/>
- J. Cassell. "More than just another pretty face: Embodied conversational interface agents." Communications of the ACM 43(4): 70-78.2000.
- Oliver Lemon, Alexander Gruenstein, and Stanley Peters, "Collaborative Activities and Multi-tasking in Dialogue Systems." Traitement Automatique des Langues (TAL), special issue on dialogue, 43(2): 131 - 154, 2002.
- Sidner et al. "Egagement by Looking:Behaviors for Robots when collaborating with people." In Proc. Of DiaBruck. pp. 123-130. 2003.

