# BITS Workshop

## Oct. 8/9 2002, Munich

**Tilman Becker**

**German Research Center for Artifical
Intelligence, DFKI**

**D–66123 Saarbrücken, Germany**

***Tilman.Becker@dfki.de***

- **from**

  **speech corpora**

- **to**

  **speech synthesis**

- **to**

  **speech generation**


- **Natural Language Generation (NLG) is the customer of the customer of the corpus**

# SmartKom: A Transportable Interface Agent



**Application Layer**

**MM Dialogue Back-bone**

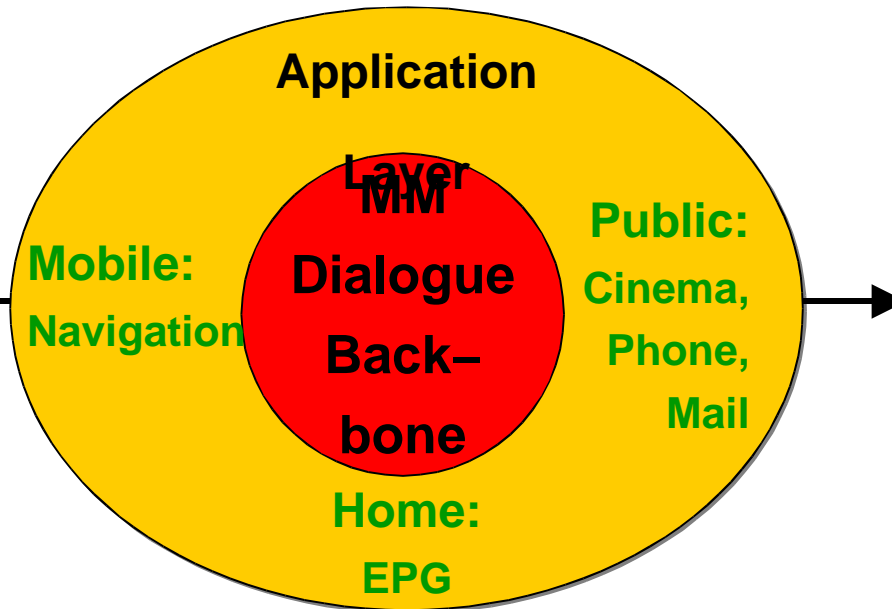**Mobile:** Navigation

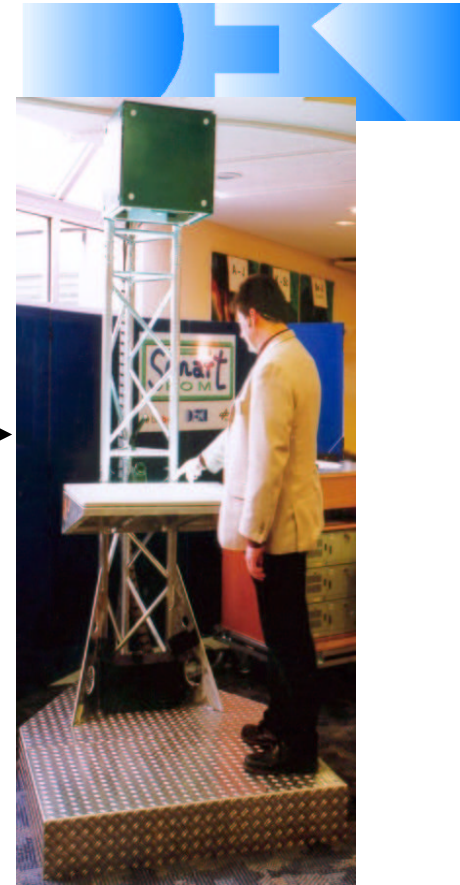**Public:** Cinema, Phone, Mail

**Home:** EPG

**SmartKom−Mobile**:
A Handheld
Communication
Assistant

**SmartKom−Home/Office**:
Multimodal Portal
to Information Services

**SmartKom−Public**:
A Multimodal
Communication
Kiosk

# SmartKom

**Scenario:**

public (mobile, home)

**Application:**

movie information

(EPG, email, phone, fax,

address book,

tv and vcr control,
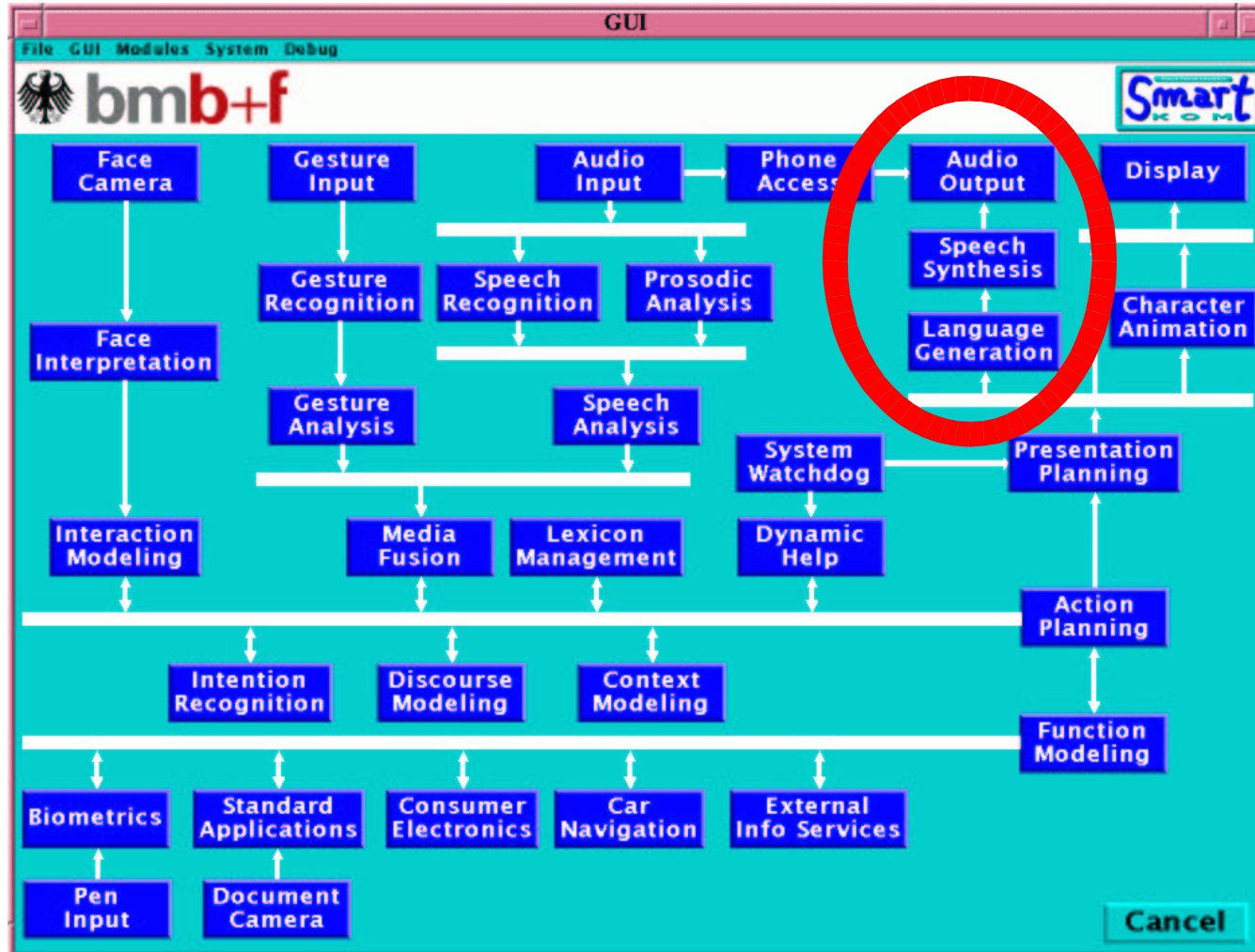
routing/tourist info)



**U:** *I want to make a reservation in (R) this movie theater*

**S:** This theater does not take reservations

**U:** *Then a different one, (R) this one perhaps*

# Control GUI in SmartKom

```
:syntaxElement  case="acc"  argumentStatus="Object"  syntaxCategory="NP">
 <syntaxElement  syntaxCategory="Det">
  <lexicalElement  partOfSpeechTag="ART">
   <text>              die                </text>
  </lexicalElement>
 </syntaxElement>
 <syntaxElement  syntaxCategory="N">
  <lexicalElement    partOfSpeechTag="NN">
   <text>       Anfangszeiten    </text>
   </lexicalElement>
  </syntaxElement>


taxElement>
ence>
urseElement>
```
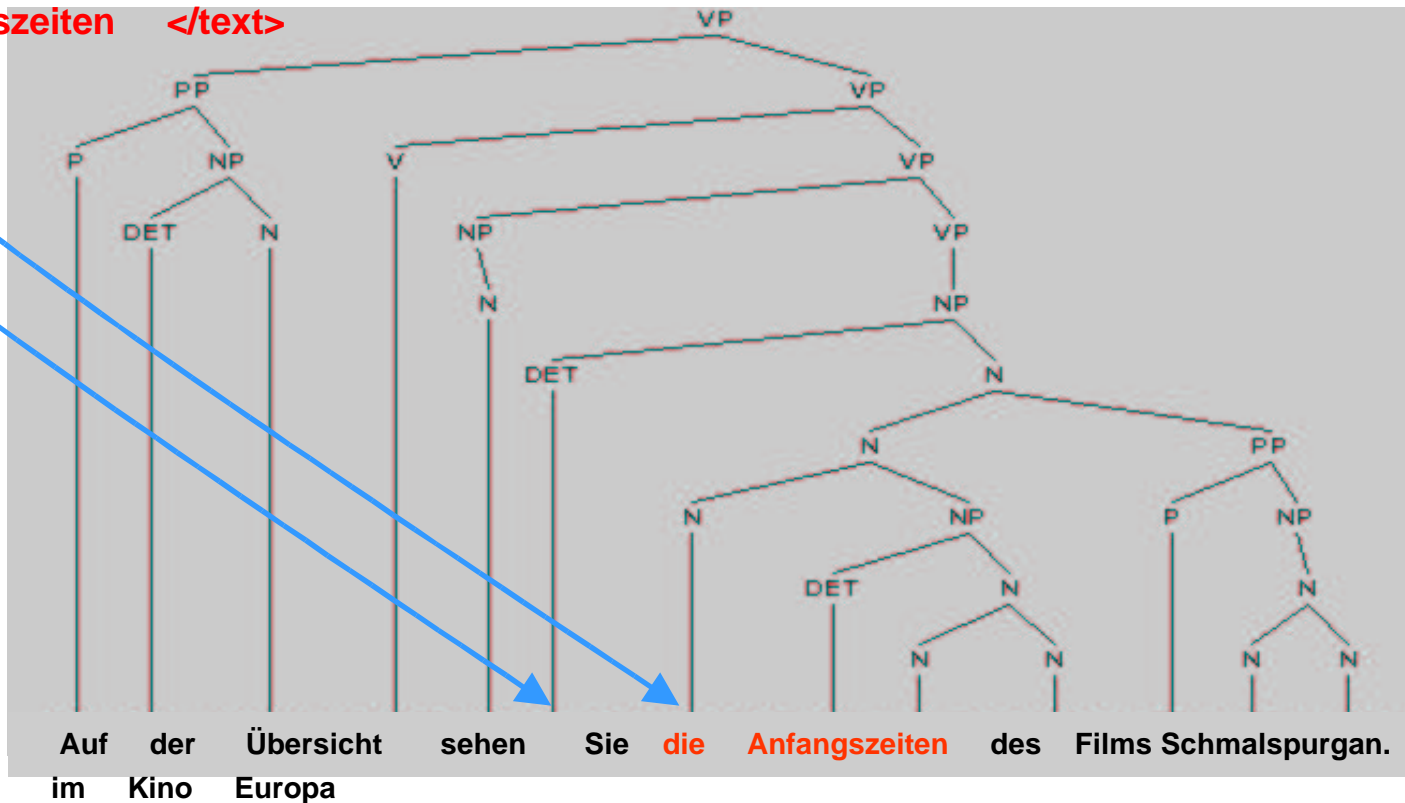
**XML Format**

**graphical**

**representation**

6

# CTS: Available Information

- **Discourse Structure**

  – **e.g., elaboration, contrast**

- **Discourse Status**

  – **given/new, topic/focus**

- **Semantics?**

- **Syntax:**

  – **boundaries: syntactic structure**

  – **sentence type, POS, morphology, ...**

- **Attitude, Emotion,**

- **Multimodal Aspects:**

  – **coordination with gesture, posture, facial expressions, ...**

# But what is relevant?

- **Many different types of data available but**

- **Which should be part of a corpus:**

    – **automatic annotation**

    – **standardized formats**

    – **relevance to synthesis**

# Annotation Formats

- **Information structure**

  – **RST**

- **Attitude?**

- **Syntax**

  – **Penn Treebank**

  – **Chunks, topological fields, ...**

- **POS**

  – **STTS**