# Chapter 6. Data Collection in a Wizard-of-Oz Experiment

**in Reinforcement Learning for Adaptive Dialogue Systems
by: Rieser & Lemon.**

## Course: Autonomous Machine Learning

Hanan Al Rahbi
hanan.alrahbi@snu.ac.kr

Interdisciplinary Program in Cognitive Science
Seoul National University

# The Turk

https://youtu.be/RdT4yG8wczQ

# Contents

# What is a WOZ?

- A research method in which a human being simulates the intelligent behavior of a machine
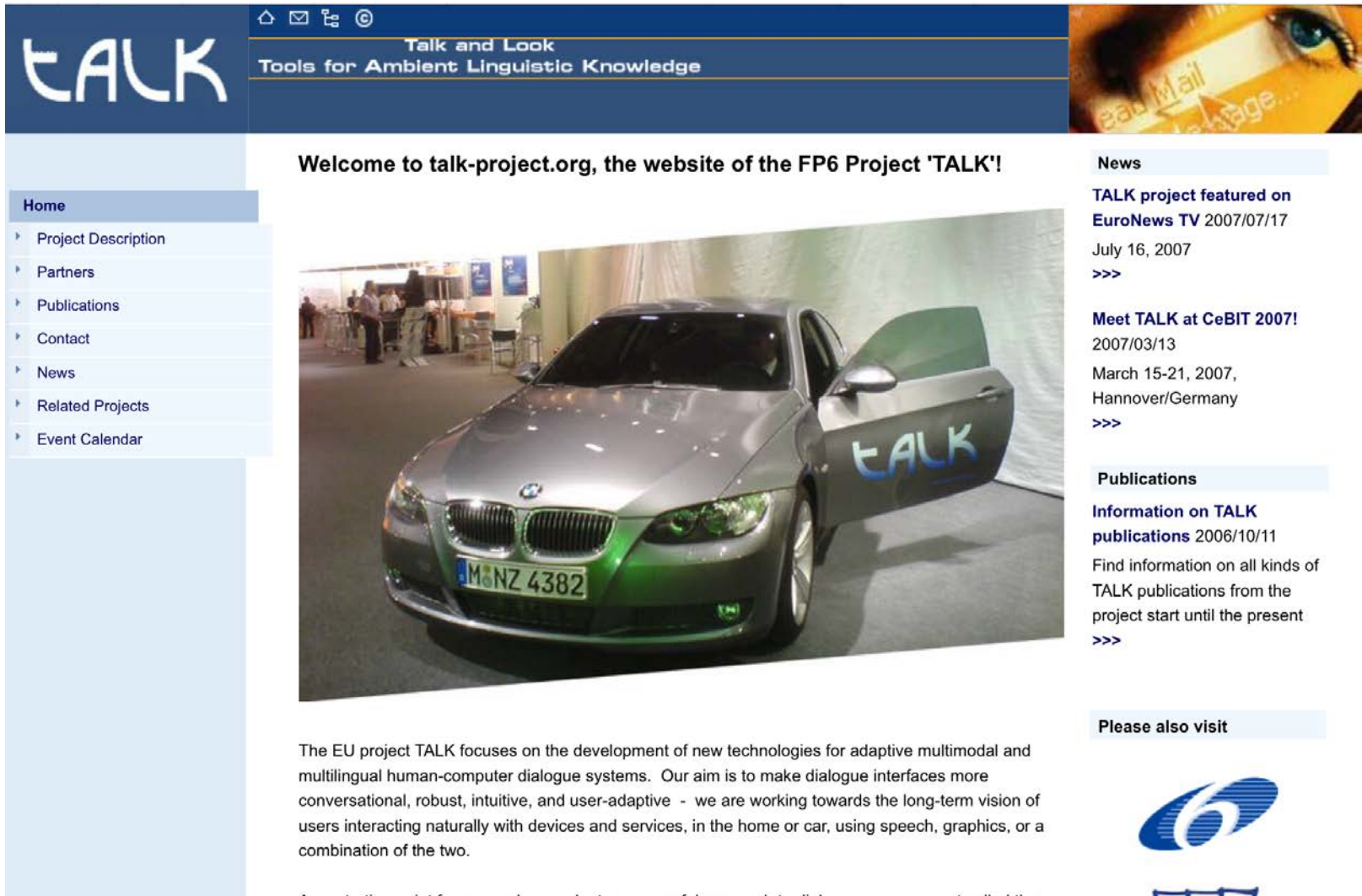- Before one is able to build a full working system

# What is it used for?

- Collecting initial data before a system is designed
- Producing more intelligent behavior by current machines

# Why is it important?

- Allows data-driven development for domains with no available prototypes

- Helps creating effective policies with the highest rewards

- Saves effort and time

# The Experiment

# Experimental setup



Multimodal Wizard-of-Oz data collection setup for an in-car music player application

# Recruited Subjects: Wizards & Users

- This experiment focuses on the <u>behavior</u> of <u>users</u> and <u>wizards</u>

|  | Quota | Age | Language |
|---|---|---|---|
| Wizards | 5 (2F, 1M) | 20~35 | German: Native<br>English: Good |

**No Experience in dialogue systems**

€

|  | Quota | Age | Field of Study |
|---|---|---|---|
| Users | 21 (11F, 10M) | 20~30 | Social Science:23.8%<br>Languages:23.8%<br>Natural Sciences:28.6%<br>Arts:17% |

# Experimental Procedure and Task Design

1. Training wizards (database, interaction with users)
2. User and wizard placed in separate rooms
3. User received sheet of instructions upon arrival
4. Introducing the user to the driving simulator (tested)
5.  User could solve the tasks in any order they preferred
6. After each task user filled task-specific questionnaire
7. User interviewed by experiment leader

a) Simple text-message conveying how many results were found
b) Output of a list of just the name (album, song, or artist)
c) A table of complete search results
d) A table of complete search results but only displaying a subset of of columns.

# Experimental Procedure and Task Design

- Designed 10 task sets
- Every task set was used at least twice
- Each set contains 4 tasks of 2 different types:
  - Search for a specific title/ album
  - Build a playlist

# Noise Simulation

- HCI vs. WOZ

- Related work:
  - Skauntze (2003,2005)
  - Stuttle et al (2004); Williams and Young (2004a)
    - Even with high noise wizards are able to interpret the ASR output well and assimilate contextual knowledge about what user actions are likely to follow

# Noise Simulation

- ## Method:

  - To approximate speech recognition errors, a tool was used to randomly delete parts of the transcribed utterances

  - Wizards also build up their own hypotheses about what the user really said (misunderstandings)

  - Word deletion rate of the text varied:

    - 20% weakly corrupted = deletion rate of 20%

    - 20% strongly corrupted= deletion rate of 50%

    - In 60% of the cases the wizard saw the transcribed speech uncorrupted

# Noise Simulation

**Example 6.2.1**

*uncorrupted:  Zu dieser Liste bitte Track 'Tonight' hinzufügen.*
  *[ Add track 'Tonight' to this list.]*
*weakly corrupted:  Zu dieser Liste bitte Track Tonight . . . .*
  *[. . . track 'Tonight' to this list.]*
*strongly corrupted:  Zu . . . Track Tonight . . . .*
  *[. . . track 'Tonight' to . . . .]*

# Results and Discussion

- 30% of the corrupted utterances had a noticeable effect on the interaction

- 7% of all user turns lead to a communication error (much lower than the current WER for spoken dialogue systems {around 30%})

- On the other hand, the error rate is higher than for human-human communication

# Results and Discussion

- Shortcomings of the deletion method:
  - Deleting words is a rather crude simulation of real-world acoustic problems (justified)
  - Time delay introduced by transcribing the utterances (both of user and wizard)

- This method is not suitable for studying detailed error, however, it can be sufficient in order to study natural presentation strategies under the presence of noise.

# Corpus Description

- 21 sessions, containing 72 dialogue, with about 1600 turns were gathered

- Data for each session includes video and audio recordin, questionnaire data, transcripts, and a log file

- The logging information per session consists of OAA messages in chronological order

- Corpus is marked up and annotated using Nite XML Toolkit (NXT)

# Analysis

- Results of corpus analysis for multimodal presentation strategies
- Qualitative measures:
    - 22.3% of the 793 wizard turns were annotated as presentation strategies, resulting in 177 instances for learning
    - 48% screen output
        - 78.6% the table option
        - 17% the list
        - 0.04% text only

    - Verbal presentation only present 1.6 items on average
        - Where wizard summarized the results by presenting the options for the most distinctive feature to the user.

# Analysis

- Did the Wizards apply significantly different strategies? It is important to compare! (data will be used for learning)
  - Dialogue length is about the same with very slight differences between wizards
  - Most wizards were equally successful in completing tasks, only one was better with 100% task success, where another one scored 78% task success
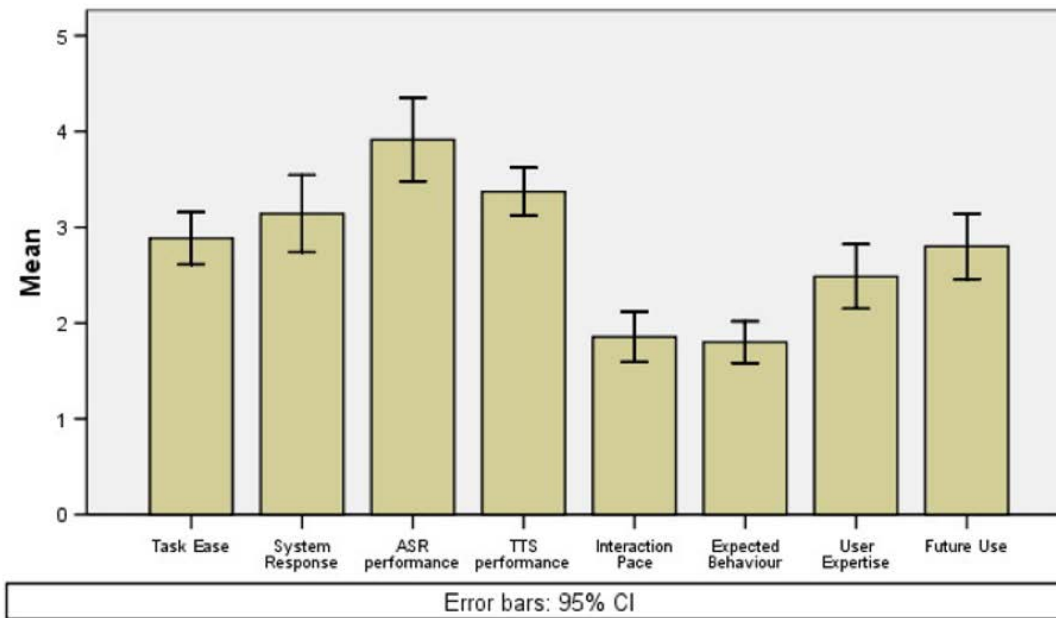
Therefore we can say they applied similar strategies (this doesn't mean they react the same way)

However, multimodal behavior of wizards is very limited
  - Only 3 users selected an item by clicking

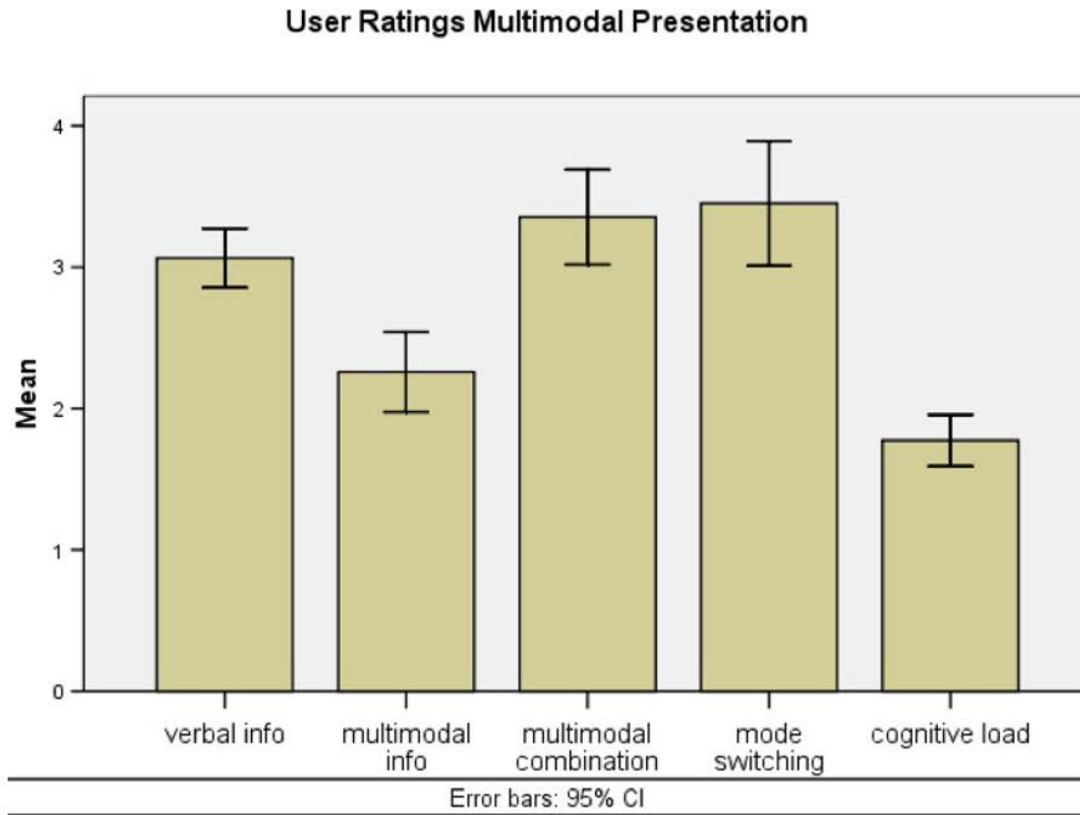# Analysis

- Subjective Ratings from the User Questionnaire



**Fig. 6.6** User Ratings on the PARADISE questions on a 5-point Likert Scale

# Analysis

- Subjective Ratings from the User Questionnaire



**Fig. 6.7** User Ratings on the multimodal presentation questions on a 5-point Likert Scale

# Discussion

- ## Common mistakes (the wizard either):
  - Either wizard displayed too much information On the screen
  - Or fail to present results early enough

  Screen output should display <u>appropriate</u> amount of information

  - There is a need for a strategy which decides <u>how many</u> database search results to present to the user, <u>when</u>, and <u>which modality to use</u> in an adaptive optimal matter
  - Also a strategy to help minimize the large lists displayed, cut the length of the dialogue, as well as the noise
  - Include information about users driving performance is very important
  - There should be a better and more realistic in-car simulation (the screen size)