

Experimenting with grounding strategies in dialogue

Volha Petukhova¹, Harry Bunt², Andrei Malchanau¹, Ramkumar Aruchamy¹

¹ Saarland University, Spoken Language Systems, Germany; ² Tilburg University, Netherlands
{v.petukhova, andrei.malchanau, ramkumar.aruchamy}@lsv.uni-saarland.de harry.bunt@uvt.nl

Abstract

This paper discusses empirically grounded strategies for the generation of feedback acts by a dialogue system in a way that supports a natural communication style and therefore leads to higher user acceptance. User evaluation of an implemented prototype system shows that an appropriate strategy can be generated by rules that are based on an analysis of human-human dialogue behaviour for a given task and domain.

1 Introduction

While conversational speech-based applications have recently begun penetrating the mass market, commercial dialogue systems are still limited to a rather restricted communication behaviour modelled on information providing tasks. Some systems developed for research purposes allow for more natural conversations, but they are often limited to a narrow domain with manually crafted domain models and pre-baked dialogue strategies. Alternatively, dialogue strategies can be adapted through reinforcement learning, but this requires large amounts of training data, while offering only a limited range of dialogue actions.

In this paper we show how a relatively small amount of 'Wizard-of-Oz' (WoZ) data and focused analysis of the phenomena related to grounding can help to design various strategies and communicative styles in order for a dialogue system to exhibit behaviour that is more natural to its users.

2 Observed grounding behaviour

To simulate user's information-seeking and system's information-providing behaviour we designed a set of quiz games. Data has been collected in a WoZ setting with the Wizard holding the facts about a famous person's life, and a

player guessing his/her identity by asking questions of various types. 338 dialogues were collected (16 hours comprising about 6.000 speaking turns, 18 turns per dialogue), transcribed and annotated with ISO 24617-2 dialogue acts.¹

For an interactive system it is important to know that its contributions are understood and accepted (i.e. grounded) by the user. In our quiz scenario, if the answer is understood and accepted by the player, he continues with his next question. However, we do not just observe question-answer pairs. Players very often signal their understanding and acceptance of the previous system utterance by repeating or rephrasing (part of) it, known as 'implicit verification', or accepting answers with inarticulate positive feedback like 'Okay', 'mm-mhm', 'yeah', 'right', etc. This allows the user to verify the correctness of the system's recognition of the preceding utterance, and gives the user the possibility to correct mistakes on the fly (allo-feedback). In case of positive feedback from a player, the Wizard often explicitly acknowledges it, and in case of negative feedback always reacts to it.

We analysed the data for the occurrence of sequences of Questions, Answers, positive/negative Auto- and AlloFeedback acts. Table 1 presents the frequencies of the patterns that were observed. These patterns were used to construct a decision tree for feedback generation, weighting possible transitions from one state to another. It may be observed that the simple Question-Answer sequence is the most frequent pattern, however explicit positive Auto-Feedback occurs quite often.

A dialogue system that provides positive auto-feedback after every user contribution would exhibit a style of communicative behaviour is unnatural and even annoying. It is therefore interesting

¹For the ISO 24617-2 dialogue act annotation standard see Bunt et al., 2012; for details on the data collection and the annotation see Petukhova et al., 2014.

| Observed sequence | Frequency (in %) |
|---------------------------------------------------------------------------------------------------------------|------------------|
| P:Question1 - W:Answer - P:Question2 | 47.1 |
| P:Question1 - W:pos. AutoFeedback - W:Answer - P:Question2 | 28.6 |
| P:Question1 - W:neg. AutoFeedback(execution: answer not found) - P:pos. AutoFeedback - P:Question2 | 7.6 |
| P:Question1 - W:neg. AutoFeedback - P:Repeat/rephrase Question - W:pos. AutoFeedback - W:Answer - P:Question2 | 6.3 |
| P:Question1 - W:pos. AutoFeedback - W:Answer - P:pos. Allo/AutoFeedback - P:Question2 | 4.9 |
| P:Question1 - W:neg. AutoFeedback(execution: answer not found) - P:Question2 | 2.8 |
| P:Question1 - W:neg. AutoFeedback - P:Repeat/rephrase Question - W:Answer - P:Question2 | 2.7 |

Table 1: Observed sequences of player-system acts ranked according to relative frequencies.

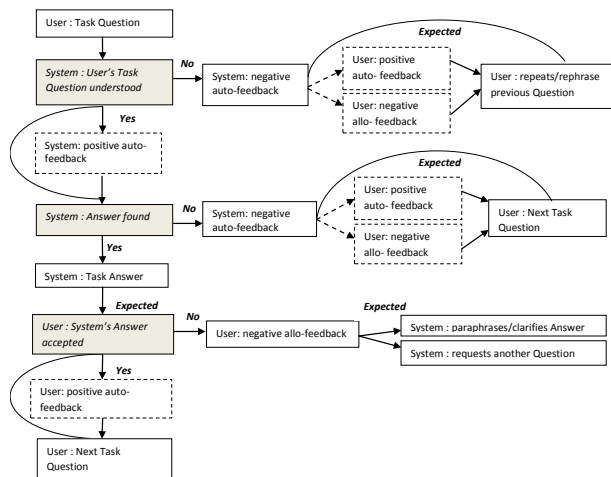


Figure 1: Decision tree for the generation of dialogue acts by the system. Dashed boxes present optional actions; gray boxes represent actual or expected processing states.)

to consider strategies where positive feedback is generated *regularly* when the dialogue reaches a crucially important state, and only *occasionally* in other situations, e.g. when it is not vital for task performance, and *regularly*, e.g. when the dialogue reaches a crucially important state. For our scenario and dialogue setting we designed a decision tree that incorporates the observations and analysis of our data for the generation of various types of feedback, see Figure 1.

To evaluate to what extent users feel that grounding strategies modelled in such a way lead to natural and flexible interactive behaviour, three question-answering system prototypes were implemented using the NPCEditor tool², extending the dialogue management strategies defined in the NPCEditor in order for the system to show more complex interactive behaviour beyond question-answering, by adding more dialogue manager states and a wider variety of positive and negative auto- and allo-feedback act types.

For the evaluation we investigated user satisfaction using a questionnaire filled in after interacting with the system. A within-subject evaluation was performed with 6 users who played a game us-

²<https://confluence.ict.usc.edu/display/VHTK/NPCEditor>

ing three different system prototypes: (1) minimal query - response (MQR) setting; (2) system always generating explicit auto-feedback to player's query (AEFR); and (3) system generating explicit feedback according to the decision tree shown in Figure 1 (DEFR).

We tested user satisfaction by asking subjects to rate their level of agreement on the following parameters: (i) learnability (the ease to use the system, e.g. rules well explained); (ii) ability to get the requested information; (iii) correctness of answers; (iv) frequency and type of system feedback; (v) speed of responses; (vi) naturalness of the interaction and (vii) overall attitude, e.g. likability and engagement. For each parameter we obtained the agreement scores. Responses for each question were summed up and divided by the number of participants to calculate the level of agreement in terms of average Likert scores.³ The results show that players in general appreciate explicit feedback, and when the system generated feedback acts according to the decision tree it received the highest score on all criteria without exception: MQR was rated 3.4 on 5-point Likert scale; AEFR - 3.6; and DEFR - 4.5.

This exploratory study left some unexplored and/or not implemented options. For instance, the behaviour in other dimensions than the feedback dimensions such as Turn-, Time-, Own- and Partner Communication Management, and Discourse Structuring deserves attention; findings there may well lead to more interesting and flexible behaviour on the part of the system.

References

- Bunt, H. et al. 2012 ISO 24617-2: A semantically-based standard for dialogue annotation. In: *Proceedings of LREC 2012*, Istanbul, Turkey
- Petukhova, V. et al. 2014 The DBOX corpus collection of spoken human-human and human-machine dialogues. In: *Proceedings of LREC 2014*, Reykjavik, Iceland

³Statistical tests were not performed, since the number of raters (6) was too low to draw statistically significant conclusions. Our goal was rather to obtain first impressions of the acceptance or rejection of different grounding strategies by human players.