

From Proceedings of the 13th International Conference on Artificial Intelligence in Education, AIED 2007, July 9-13, Marina del Rey, California. This is the authors' own copy that is made available from the authors' web page according to the copyright agreement.

The Influence of Learner Characteristics on Task-Oriented Tutorial Dialogue

Kristy Elizabeth BOYER, Mladen A. VOUK, James C. LESTER
*Department of Computer Science, North Carolina State University
Raleigh, NC 27695, USA
{keboyer, vouk, lester}@ncsu.edu*

Abstract. Tutorial dialogue has been the subject of increasing attention in recent years, and it has become evident that empirical studies of human-human tutorial dialogue can contribute important insights to the design of computational models of dialogue. Students with particular characteristics may have specific dialogue profiles, and knowledge of such profiles could inform the design of tutorial dialogue systems whose strategies leverage the characteristics of the target population and address the communicative needs of those students. This paper reports on a study that was conducted to investigate the influence of learner characteristics (performance levels, self-efficacy, and gender) on the structure of task-oriented tutorial dialogue. A tutorial dialogue corpus was gathered from interactions transpiring in the course of problem-solving in a learning environment for introductory computer science. Analyses of the annotated dialogues suggest that the dialogue structure of (1) low-performing students differs significantly from that of high-performing students, (2) students with low self-efficacy differs significantly from that of students with high self-efficacy, and (3) males differs significantly from that of females.

1. Introduction

Providing intelligent tutoring systems with the ability to engage learners in rich natural language dialogue has been a goal of the AI & Education community since the inception of the field. With the investigation of tutorial dialogue in a number of systems devised to support a broad range of conversational phenomena (e.g., CIRCSIM [1], BEETLE [2], the GEOMETRY EXPLANATION TUTOR [3], WHY2/ATLAS [4], ITSPOKE [5], SCOT [6], ProPL [7] and AUTOTUTOR [8]), we have begun to see the emergence of a core set of foundational requirements and functionalities for mixed-initiative natural language interaction. Moreover, recent years have witnessed the appearance of corpus studies empirically investigating speech acts in tutorial dialogue [9], dialogues' correlation with learning [10, 11, 12, 13], student uncertainty in dialogue [14, 15], and comparing text-based and spoken dialogue [5].

While all learners may engage in some universal form of tutorial dialogue, it may be the case that different populations of learners engage in qualitatively different forms of

dialogue. It seems plausible that students with particular characteristics may have specific dialogue profiles, and knowledge of such profiles could inform the design of tutorial dialogue systems whose strategies leverage the characteristics of the target population and address the communicative needs of those students. This paper reports on a study investigating the influence of learners' achievement levels, self-efficacy, and gender on task-oriented tutorial dialogue.

Given that human-human tutorial dialogue offers a promising model for effective communication [16], an experiment was conducted to study naturally occurring tutorial dialogues in a task-oriented learning environment. A text-based dialogue interface was incorporated into a learning environment for introductory computer science. In the environment, students undertook a programming task and conversed with human tutors while designing, implementing, and testing Java programs. To ensure that only natural language was used for communication and to eliminate the possibility of non-verbal communication such as gesture and body language, tutors were physically separated from students in an adjoining lab. The tutors' interface included a real-time synchronized view of the students' problem-solving workspace. Dialogues were logged and the resulting corpus was then manually annotated with tutorial dialogue acts. Analyses of the annotated dialogues suggest that the dialogue structure of low-performing students differs significantly from that of high-performing students, that the dialogue structure of students with low self-efficacy differs significantly from that of students with high self-efficacy, and that the dialogue structure of males differs significantly from that of females.

2. Task-Oriented Tutorial Dialogue Corpus and Dialogue Acts

While all tutorial dialogue is undertaken in support of learning tasks, one genre of tutorial dialogue is directly situated in the task at hand: these dialogues emerge as a result of the creation of learning artifacts such as designs, proofs, or computer programs. The domain investigated in this study, which has also been studied in the ProPL tutorial dialogue project [7], is that of computer programs. Here, students design, implement, and test programs (in the case of the study, Java programs) to meet a given specification. In the course of constructing the artifact, tutors and students pose questions to one another, tutors offer advice and feedback, and students make statements about the artifacts.

The Java Corpus was gathered by logging text-based dialogues between tutors and novice computer science students. The learning task was to complete a programming problem that required students to apply fundamental concepts such as iteration, modularization, and sequential-access data structures. Table 1 presents two sample annotated dialogue excerpts from the Java Corpus. In Dialogue Excerpt A, the tutor interacts with a low performing student, Student A, whose pre-test score was well below the median. The structure of Dialogue A illustrates many features commonly seen with low performing students, such as seeking to establish confirmation of a proposed plan before proceeding to implementation. In contrast to Dialogue A, Dialogue B illustrates some common characteristics of dialogues with high performing students. Student B seeks tutorial advice, then proceeds directly to implementation with no pre-emptive request for tutor feedback.

Table 1: Dialogue Excerpts

Dialogue Excerpt A	Dialogue Excerpt B
Tutor: Do you know how to do that? [TQ]	Student B: So I need an if for each digit? [TQ]
Student A: Not really. [A]	Tutor: One if should suffice, since it will be called in each iteration. [A]
Tutor: Well we first need a new String that will hold zipCode's string value. [HA]	Tutor: You just need to know which element to reference. [A]
Student A: So String z = zipCode? [RF]	Tutor: This would be done in the inner loop. [HA]
Tutor: Close. [PLF]	Student B: Ok. [ACK]
Tutor: Then you can set that string equal to ""+zipcode. [HA]	<i>(Student works for 2.5 minutes.)</i>
Student A: Ok so String z="" + zipCode [RF]	Tutor: You've got the right idea. [UPF]
Tutor: Yeah. [PPF]	Student B: Yeah, I had programmers block. [EX]
Student A: Then what? [TQ]	<i>(Student works for 3 minutes.)</i>
	Tutor: Perfect. [UPF]

The Java Corpus consists of 5034 dialogue acts: 3075 tutor turns and 1959 student turns. The corpus was manually annotated with a set of tutorial dialogue acts designed to capture the salient characteristics of task-oriented tutorial dialogues. The coding scheme (Table 2) draws on a scheme devised for tutorial dialogue on qualitative physics problems [10]. While most of the acts in this scheme are present in the Java corpus as well, the particular dialogues in the Java corpus made it difficult to make judgements about short answer questions versus deep answer questions and to make fine-grained distinctions between hinting levels. The four-category scheme [9] and a more expansive non-tutorial dialogue act catalogue [17] also contributed common acts.

The entire corpus was annotated by a single annotator. In an agreement study to evaluate the consistency of the coding scheme and its application to the corpus, a second annotator labelled a subset of 969 acts (of the total 5034 acts in the corpus). This yielded a 0.75 Kappa between the two annotators, indicating a reasonable inter-rater reliability.

3. Experimental Design

Subjects were students enrolled in an introductory computer science course and were primarily freshman or sophomore engineering majors in disciplines such as mechanical, electrical, and computer engineering.

The corpus was gathered from tutor-student interactions between 35 students and 6 tutors during a one-week study. Tutors and students were completely blind to each other's characteristics as they worked together remotely from separate labs. Tutors observed student problem-solving actions (e.g., programming, scrolling, running programs) in real time. The tutors consisted of four graduate students and two undergraduates in computer science; five were male and one was female. Tutors all had some level of tutoring experience, and were not instructed about specific tutorial strategies.

Subjects first completed a pre-survey including items about self-efficacy, attitude toward computer science, and attitude toward collaboration. Subjects then completed a ten item pre-test over specific topic content. The tutorial session was controlled at 50 minutes for all subjects, after which subjects completed a post-survey and post-test containing variants of the items on the pre- versions. Any subject whose session was interrupted due to technical difficulties or external factors, or who completed the task early, was omitted from the data set for analysis ($n_{\text{omitted}}=6$).

Table 2: Dialogue Acts

	Act	Description	Examples
Student/Tutor	Task Question (TQ)	Questions about goals to pursue, ordering of goals, and the specific problem being solved.	"Where should we start?" "Should we use an array?"
	Concept Question (CQ)	Questions about domain elements, concepts, or facts that would apply over many different problems.	"How do I declare an array?" "I don't know how to write a loop."
	Answer (A)	Answers to a task or conceptual question.	"No." or "Yes." "We need to give it an index."
	Acknowledgement (ACK)	Positive acknowledgement of a previous statement.	"Okay." or "Yeah." "Alright."
	Extra Domain (EX)	A statement not related to the computer science discussion.	"Hello." or "Sorry." "Nice working with you."
Student	Request Feedback (RF)	A request for evaluative feedback on completed or proposed problem solving steps.	"So should I do array[0] = 1?" "Does that look good?"
	Signal Non-Understanding (SNU)	An indication that a previous statement by the tutor is not clear.	"Kind of makes sense." "Not really." or "I'm confused."
	Statement (S)	Assertion of fact.	"I am going to use a for loop." "We need to initialize that variable."
	(Un)Prompted Positive Feedback (UPF/PPF)	Positive feedback regarding problem solving action.	"Good job." or "Looks great." "Yep."
Tutor	(Un)Prompted Lukewarm Feedback (ULF/PLF)	Partly positive, partly negative feedback regarding student problem solving action.	"The first part is right, but..." "You're close." or "Well, almost."
	(Un)Prompted Negative Feedback (UNF/PNF)	Negative feedback regarding student problem solving action.	"No." "Actually, that won't work."
	Hint/Advice (HA)	Problem solving or conceptual hint or advice not in answer to a direct question.	"Each digit is represented by 5 bars." "Let's move on."
	Request to Confirm Understanding (RCU)	A request for student to confirm or disconfirm understanding.	"Does that make sense?" "Are you with me?"

4. Results and Discussion

To compare dialogue structure based on learner characteristics, three partitioning criteria were applied to the student population: incoming performance level, self-efficacy rating, and gender. After briefly noting overall learning effectiveness, this section reports on dialogue structure characteristics for each student sub-population based on each of the three partitioning criteria.

For each student, learning gain was gauged by the difference between pre and post-test scores. On average, students scored 13% higher on the post-test than the pre-test. A pair-wise difference t-test indicates that the difference is statistically significant ($p < 0.05$).

4.1 Dialogue Profile Analyses

For each student dialogue session, the relative frequency of each dialogue act was computed as the ratio of the number of occurrences of that dialogue act to the total number of dialogue acts in the session. The relative frequency of dialogue acts was then computed for high-performing and low-performing students, for high-efficacy and low-efficacy students, and for female and male students. To determine whether intra-group differences in means were significant, t-tests were performed. Table 3 summarizes the relative frequency results, omitting all dialogue acts for which there was no significant difference by learner characteristics. Bolded values indicate statistically significant differences ($p \leq 0.05$). It should be noted that the three partitions are not independent. For example, high performing students were more often in the high self-efficacy group, and most females were in the low performing group. Despite these confounds, we draw meaningful conclusions by examining each learner characteristic individually.

Students were divided into low performing and high performing groups based on the median pre-test score. Analyses yielded the following findings: 1) High performing students made more acknowledgements, requested feedback less often, and made more declarative statements than low performing students. 2) Tutors paired with low performing students made more extra-domain statements, gave more prompted feedback, and made more requests for confirmation of understanding than tutors paired with high performing students.

Following an instrument devised by Bandura to measure domain-specific self-efficacy [18], students were asked to rate their confidence in being able to complete a programming assignment under a variety of circumstances. Because the problem used for this study was drawn from a standard problem set for the course, students had an experiential basis on which to judge their ability to complete the problem. Statistically significant differences in dialogue structure emerged when students were grouped by their confidence level regarding whether they could complete a simple laboratory assignment on their own. Analyses yielded the following findings: 1) Students in the high self-efficacy group made more declarative statements, or assertions, than students in the low self-efficacy group. 2) Tutors paired with low self-efficacy students gave more negative feedback and made fewer acknowledgements than tutors paired with low self-efficacy students.

Although females comprised a small number of our subjects, some statistically significant results emerged. 1) Women made more requests for feedback and fewer declarative statements than men. 2) Tutors paired with women gave more positive feedback and made more requests to confirm understanding than tutors paired with men.

4.2 Discussion

These findings extend those of previous studies investigating tutorial dialogue and learning effectiveness which have found correlations of dialogue structure and content with learning [11, 12, 13]. Of particular interest is a large spoken tutorial dialogue study conducted as part of the ITSPPOKE project [10]. The ITSPPOKE study found that student

utterances exhibiting reasoning and reasoning-oriented questions posed by the tutor were positively correlated with learning in a human-computer corpus, as were the introduction of new concepts in the dialogue by students in a human-human corpus. The Java Corpus study reported on here found that learner characteristics appear to significantly affect the structure of tutorial dialogue, and that both tutor and student dialogue acts appear to be affected by these differences. Tutors more often engaged in more extra-domain conversation, provided additional feedback, and more frequently engaged in discussions to gauge students' level of understanding when conversing with low performing, low efficacy, or female students. These same groups of students tended to request more feedback, make fewer declarative statements, and make fewer acknowledgements. It seems likely that learner characteristics affect (and are affected by) tutorial dialogue issues analogous to those bearing on help-seeking behaviors [19] and self-explanation [20].

These findings suggest that it may be possible to devise tutorial dialogue strategies that address the specific communicative needs of different groups of learners. Putting gender differences aside because of the limited data, several design implications could be considered for tutorial dialogue systems:

- **Encouraging Reflection:** If a student with low incoming performance initiates few requests for feedback, the system should consider taking remedial action such as asking task questions or concept questions to assess student understanding.
- **Giving Adequate Feedback:** Systems should be prepared to give prompted feedback more often when working with low-performing or low-efficacy students.

Table 3: Relative Frequency Results

Dialogue Act	Pre-test Performance		Self-efficacy Level		Gender	
	$n_{high}=17$, $n_{low}=18$		$n_{high}=19$, $n_{low}=16$		$n_{female}=7$, $n_{male}=28$	
Student Acknowledgement (ACK)	High	10.3%	High	9.1%	Female	8.3%
	Low	6.3%	Low	7.3%	Male	8.2%
Tutor Acknowledgement (ACK)	High	2.3%	High	2.5%	Female	1.2%
	Low	1.5%	Low	1.2%	Male	2.1%
Tutor Extra Domain (EX)	High	6.9%	High	8.5%	Female	8.4%
	Low	9.4%	Low	7.8%	Male	8.2%
Student Request Feedback (RF)	High	1.2%	High	1.6%	Female	2.8%
	Low	2.2%	Low	2.0%	Male	1.5%
Student Statement (S)	High	3.5%	High	3.5%	Female	1.2%
	Low	1.6%	Low	1.4%	Male	2.8%
Tutor Prompted Positive Feedback (PPF)	High	0.5%	High	0.9%	Female	1.7%
	Low	1.5%	Low	1.3%	Male	0.9%
Tutor Prompted Lukewarm Feedback (PLF)	High	0.1%	High	0.2%	Female	0.6%
	Low	0.6%	Low	0.5%	Male	0.3%
Tutor Prompted Negative Feedback (PNF)	High	0.1%	High	0.0%	Female	0.3%
	Low	0.3%	Low	0.4%	Male	0.1%
Tutor Request Confirmation of Understanding (RCU)	High	0.1%	High	0.3%	Female	1.4%
	Low	0.9%	Low	0.8%	Male	0.3%

- **Making Acknowledgements:** When interacting with high-efficacy students, systems might give more acknowledgements than in their default setting; this may more accurately reflect the interaction expected when working with a human tutor.
- **Maintaining Conversational Comfort:** When interacting with students who have been deemed to be low-performing prior to the tutoring session, systems should consider making slightly more extra-domain statements, which could create a more conversational setting in which weaker students might feel more at ease.

5. Conclusion

Tutorial dialogue exhibits structural regularities that cut across learning tasks and domains. However, learner characteristics may profoundly affect the structure of tutor-student conversations. Analyses of task-oriented tutorial dialogues indicate that students' incoming performance levels, self-efficacy, and gender significantly influence the structure of dialogue. The findings suggest that learner characteristics may be considered in designing tutorial dialogue strategies that more effectively target the specific needs of students with particular characteristics.

The study reported here represents a first step toward understanding how learner characteristics affect the structure of tutorial dialogue. Several directions for future work appear promising. First, it will be important to explore the influence of learner characteristics on tutorial dialogue in the presence of surface level information about students' utterances. This line of investigation is of particular interest given recent results indicating that lexical cohesion in tutorial dialogue with low-performing students is found to be highly correlated with learning [21]. Second, a comparative analysis of alternate tutoring strategies on the effectiveness and efficiency of student learning for students with targeted characteristics will yield a clearer picture of the space of tutorial dialogue. Third, students' motivation and frustration undoubtedly influence (and are influenced by) the structure and content of tutorial dialogue, so developing a better understanding of these interrelationships will contribute to more effective tutorial dialogue management.

Acknowledgements

The authors wish to thank Tiffany Barnes and the members of the IntelliMedia Center for Intelligent Systems at NC State University for their insightful discussions on this research, along with the software development team of August Dwight and Taylor Fondren. Special thanks to Scott McQuiggan for assistance in framing the efficacy questions of this work and for assistance in preparing the manuscript. This work was supported in part by National Science Foundation Grants CNS-0540523 and REC-0632450. Any opinions, findings, conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

References

- [1] M. Evens and J. Michael, *One-on-One Tutoring by Humans and Computers*. Mahwah, New Jersey: Lawrence Erlbaum Associates, 2006.
- [2] C. Zinn, J. D. Moore and M. G. Core, "A 3-tier planning architecture for managing tutorial dialogue," *Intelligent Tutoring Systems, Sixth International Conference (ITS 2002)*, 2002.
- [3] V. Alevan, K. Koedinger and O. Popescu, "A Tutorial Dialog System to Support Self-explanation: Evaluation and Open Questions," *Proceedings of the 11th International Conference on Artificial Intelligence in Education*, pp. 39-46, 2003.
- [4] K. VanLehn, P. W. Jordan, C. P. Rose, D. Bhembe, M. Bottner, A. Gaydos, M. Makatchev, U. Pappuswamy, M. Ringenberg and A. Roque, S. Siler, R. Srivastava "The architecture of Why2-Atlas: A coach for qualitative physics essay writing," *Proceedings of the 6th International Conference on Intelligent Tutoring Systems*, pp. 158-167, 2002.
- [5] D. J. Litman, C. P. Rosé, K. Forbes-Riley, K. VanLehn, D. Bhembe and S. Silliman, "Spoken Versus Typed Human and Computer Dialogue Tutoring," *International Journal of Artificial Intelligence in Education*, vol. 16, pp. 145-170, 2006.
- [6] H. Pon-Barry, K. Schultz, E. O. Bratt, B. Clark and S. Peters, "Responding to Student Uncertainty in Spoken Tutorial Dialogue Systems," *International Journal of Artificial Intelligence in Education*, vol. 16, pp. 171-194, 2006.
- [7] H. C. Lane and K. VanLehn, "Teaching the Tacit Knowledge of Programming to Novices with Natural Language Tutoring," *Computer Science Education*, vol. 15, pp. 183-201, 2005.
- [8] A. Graesser, G. Jackson, E. Mathews, H. Mitchell, A. Olney, M. Ventura, P. Chipman, D. Franceschetti, X. Hu, M.M. Louwerse, N.K. Person, and the Tutoring Research Group, "Why/AutoTutor: A Test of Learning Gains from a Physics Tutor with Natural Language Dialog," *Proceedings of the Twenty-Fifth Annual Conference of the Cognitive Science Society*, pp. 474-479, 2003.
- [9] J. Marineau, P. Wiemer-Hastings, D. Harter, B. Olde, P. Chipman, A. Karnavat, V. Pomeroy, S. Rajan, A. Graesser, and the Tutoring Research Group, "Classification of speech acts in tutorial dialog," in *Proceedings of the Workshop on Modelling Human Teaching Tactics and Strategies of ITS-2000*, pp. 65-71, 2000.
- [10] K. Forbes-Riley, D. Litman, A. Huettner and A. Ward, "Dialogue-learning correlations in spoken dialogue tutoring," *Proceedings of the 12th International Conference on Artificial Intelligence in Education (AIED 2005)*, pp. 225-232, 2005.
- [11] M. G. Core, J. D. Moore and C. Zinn, "The role of initiative in tutorial dialogue," *Proceedings of the Tenth Conference on European Chapter of the Association for Computational Linguistics*, pp. 67-74, 2003.
- [12] C. Rose, D. Bhembe, S. Siler, R. Srivastava and K. VanLehn, "The Role of Why Questions in Effective Human Tutoring," *Proceedings of Artificial Intelligence in Education*, 2003.
- [13] S. Katz, D. Allbritton and J. Connelly, "Going Beyond the Problem Given: How Human Tutors Use Post-Solution Discussions to Support Transfer," *International Journal of Artificial Intelligence in Education*, vol. 13, pp. 79-116, 2003.
- [14] J. Liscombe, J. Hirschberg and J. Venditti, "Detecting Certainty in Spoken Tutorial Dialogues," *Proceedings of Interspeech*, 2005.
- [15] K. Forbes-Riley and D. Litman, "Using Bigrams to Identify Relationships Between Student Certainty States and Tutor Responses in a Spoken Dialogue Corpus," *Proceedings of the 6th SIGdial Workshop on Discourse and Dialogue*, 2005.
- [16] M. T. H. Chi, S. A. Siler, H. Jeong, T. Yamauchi and R. G. Hausmann, "Learning from human tutoring," *Cognitive Science*, vol. 25, pp. 471-533, 2001.
- [17] A. Stolcke, K. Ries, N. Coccaro, E. Shriberg, R. Bates, D. Jurafsky, P. Taylor, R. Martin, C. Van Ess-Dykema and M. Meteer, "Dialogue act modeling for automatic tagging and recognition of conversational speech," *Computational Linguistics*, vol. 26, pp. 339-373, 2000.
- [18] A. Bandura, "Guide for constructing self-efficacy scales," in *Self-Efficacy Beliefs of Adolescents*, F. Pajares and T. Urdan, Eds. Greenwich, Connecticut: Information Age Publishing, 2006, pp. 307-337.
- [19] V. Alevan, B. McLaren, I. Roll and K. Koedinger, "Toward tutoring help seeking: applying cognitive modeling to meta-cognitive skills," *Proceedings of the 7th International Conference on Intelligent Tutoring Systems*, pp. 227-239, 2004.
- [20] M. T. H. Chi, N. Leeuw, M. H. Chiu and C. LaVancher, "Eliciting Self-Explanations Improves Understanding," *Cognitive Science*, vol. 18, pp. 439-477, 1994.
- [21] A. Ward and D. Litman, "Cohesion and learning in a tutorial spoken dialog system," *Proceedings of the 19th International FLAIRS (Florida Artificial Intelligence Research Society) Conference*, 2006.