

A New Perspective of Negotiation-Based Dialog to Enhance Metacognitive Skills in the Context of Open Learner Models

Raja M. Suleman¹ · Riichiro Mizoguchi¹ · Mitsuru Ikeda¹

Published online: 12 August 2016

© International Artificial Intelligence in Education Society 2016

Abstract Negotiation mechanism using conversational agents (chatbots) has been used in Open Learner Models (OLM) to enhance learner model accuracy and provide opportunities for learner reflection. Using chatbots that allow for natural language discussions has shown positive learning gains in students. Traditional OLMs assume a learner to be able to manage their own learning and already in an appropriate affective/behavioral state that is conducive for learning. This paper proposes a new perspective of learning that advances the state of the art in fully-negotiated OLMs by exploiting learner's affective & behavioral states to generate engaging natural language dialogues that train them to enhance their metacognitive skills. In order to achieve this, we have developed the NDLtutor that provides a natural language interface to learners. Our system generates context-aware dialogues automatically to enhance learner participation and reflection. This paper provides details on the design and implementation of the NDLtutor and discusses two evaluation studies. The first evaluation study focuses on the dialogue management capabilities of our system and demonstrates that our dialog system works satisfactorily to realize meaningful and natural interactions for negotiation. The second evaluation study investigates the effects of our system on the self-assessment and self-reflection of the learners. The results of the evaluations show that the NDLtutor is

This paper is an extended version of the paper (Suleman et al. 2015) published at AIED2015.

✉ Raja M. Suleman
suleman@jaist.ac.jp

Riichiro Mizoguchi
mizo@jaist.ac.jp

Mitsuru Ikeda
ikeda@jaist.ac.jp

¹ School of Knowledge Science, Japan Advanced Institute of Science & Technology, Nomi, Ishikawa, Japan

able to produce significant improvements in the self-assessment accuracy of the learners and also provides adequate support for prompting self-reflection in learners.

Keywords Intelligent tutoring system · Open learner model · Affect & behavior modelling · Metacognition · Interest-based negotiation

Introduction

The paradigm of Open Learner Models (OLM) was introduced in Intelligent Tutoring Systems in order to involve the learner in the overall learning experience (Bull and Pain 1995; Dimitrova 2003). OLMs generate the Learner Model (LM) of a learner by diagnosing their knowledge during their interactions with the system (VanLehn 1988). This is achieved by evaluating the learner's answers to a series of questions on a particular topic or domain. Previous LMs were hidden from the learners and were only accessible to the system. OLMs externalize the contents of the LM to promote independent learning. This is done in order to provide transparency and increase learner's trust in the system (Bull and Kay 2010). Negotiated OLMs (Bull and Vatrapsu 2012; Bull 2016) achieve this by maintaining separate belief bases for both the learner and the system. The learner is allowed to inspect (view) and edit their own belief base however they can only inspect (view) the belief base of the system. Negotiation mechanisms are used to resolve any conflict (difference) that might occur between the learner's belief base and that of the system. The result of this negotiation is used to update the LM accordingly.

Allowing the learner to edit their belief base results in scenarios where the learner's belief about their own knowledge is different from that of the system. Such events trigger an interaction where the system tries to negotiate the changes made by the learner in their belief base in an effort to remove the difference of beliefs between the learner's belief base and the system's belief base. The aim of this negotiation is to increase the accuracy of the system's LM and enhance the role of the learner in the construction and maintenance of their LMs, which increases learner reflection (Bull and Pain 1995; Kerly et al. 2008b; Dimitrova 2003).

Different approaches to negotiation have been deployed by previous fully negotiated OLMs which include menu-based interfaces (Bull and Pain 1995) and conceptual graphs (Dimitrova 2003). Conversational agents or chatbots were introduced to allow for more flexible and naturalistic negotiations (Kerly and Bull 2006). The natural language interface provided by a chatbot (Kerly et al. 2008b) improves the quality of dialogues by easing the communication between the learner and the system. The use of a chatbot yielded positive learning gains and was successful in increasing self-assessment accuracy (Kerly and Bull 2008). Through a successful trial with different age groups the research was able to identify the novelty and effectiveness of using a chatbot to discuss the LM content with the learners in the context of OLMs.

Research has shown that expert human-tutors are successful as they try to engage students according to their affective and behavioral states, which provides a sense of empathy and encourages learner involvement (Lepper et al. 1993). We believe current OLM implementations can be largely enhanced by explicit use of the information regarding such states of a learner to control the flow of the dialogue.

Improving the metacognitive abilities of the learner has always been a key role of OLMs (Bull and Kay 2013) and these systems have shown to be successful in promoting self-reflection. However, there is no explicit mechanism in current OLMs to scaffold the metacognitive processes. Self-reflection is implied implicitly, i.e. how the learner is reflecting or evaluating themselves is left on the part of the learner. The system does not explicitly involve the learner into a discussion that can motivate them to practice these skills more actively.

A conflict may occur because the learner may be confused about their knowledge, or simply have a misconception which leads them to change their LMs. The system challenges the change made by the learner and requires them to justify himself. This creates an interesting prospect to involve the learner into a discussion about their belief and what led them to believe so. Humans become stronger advocates of their beliefs once they are challenged, and are intrinsically motivated to defend their beliefs (Gal and Rucker 2010). This provides an excellent opportunity to involve an intrinsically motivated learner in a deep learning dialogue which not only discusses the domain knowledge but also encourages them to assess the discussion to promote self-reflection. In order to capture this opportunity and make use of the context, we propose a paradigm of Negotiation-Driven Learning (NDL).

Learning is maximized by proactive participation of learners; we believe that such a context is ideal to engage a learner in a dialogue that explicitly targets the metacognitive skills of the learner and provides them the scaffolding to utilize and enhance these skills. Research on the effects of using learner's affective and behavioral states to shape negotiations has shown a positive impact on the overall learning gains (du Boulay et al. 2010; Fredrickson 1998). This has not been previously studied in the context of OLMs. In NDL we aim to exploit the utility created by the occurrence of a conflict by engaging a learner in a natural language dialogue according to their affective and behavioral states and promote metacognitive skills in them through reflective dialogues and self-assessments.

The rest of the paper is organized as follows; the next section introduces the paradigm of Negotiation-Driven Learning. Here we provide the outline of the system architecture and the details of the design of dialogues in NDL. We then describe the Wizard-of-Oz experiment which is used for selecting learner's emotional states for our system as well as generating system libraries for rules of dialog management and NLP matters. This is followed by the discussion on the objectives achieved during the WoZ experiment. Next we provide a description of the different phases of NDL and then provide an example dialogue to illustrate how our envisioned system interacts with the learner. The next section introduces our implementation of the NDLtutor which is followed by two evaluation studies. The first evaluation study evaluates the dialogue management capabilities and validates the emotional states that were selected for our system. The second evaluation study explores the effects of our system on the self-reflection and self-assessment skills of the learners. In the end we give an overview of the related work and finally conclude the paper with a few concluding remarks.

Negotiation-Driven Learning

This paper proposes a new learning paradigm of Negotiation-Driven Learning which aims at *enhancing* the role of negotiations in OLMs to facilitate constructive learning.

When a learner is involved in a learning exercise, they are not only learning something new, but they are also implicitly involved in learning how to learn. More often than not they are more inclined towards executing well-practiced strategies rather than monitoring themselves. NDL aims at encouraging learners to use these metacognitive skills more actively and effectively.

NDL acts as a component of the ITS which is triggered when a conflict between the beliefs of the system and the learner occur. During its interaction with the learner the system tries to understand why the learner holds a certain belief (cause of the conflict) and tries to help them understand why it might not be true. The system uses the information about the learner's affective and behavioral states to control the flow of the dialogue to ensure maximum engagement. An NDL dialogue session is concluded when the learner is able to defend their claim, or shows an understanding of their incorrect belief by accepting the system's justification/proposal. The system's LM is updated with the outcome of the dialogue and the ITS resumes the normal course of tutoring.

Generating Dialogues in NDL

NDL allows learners to interact with the system in a natural language interface. In order to accomplish this, the system follows the negotiation protocol proposed in (Miao 2008) to allow the learner to provide justification of their change. This protocol is consistent with other protocols that have been defined and used in previous versions of OLMs (Bull and Pain 1995; Dimitrova 2003; Van Labeke et al. 2007). The system asks the learner to justify the changes they make to their belief base. If the justification provided by the learner contains an incorrect idea, the system rejects it. If the justification provided by the learner contains an "assertion", the system can ask for more information to accept it or provide a proposal to the learner to continue the dialogue further. The system initiates a reasoning process which is used to understand the motivation behind the change made by the learner. The system and the learner have equal rights to ask for further information; accept or reject a justification provided by the other party; therefore the system needs to be capable of deploying an alternative strategy in case a learner rejects its proposal/justification.

Facilitating Metacognitive Skills

Facilitating metacognitive skills has been the core of recent research on ITSs and OLMs (Bull and Kay 2007; Mitrovic and Martin 2002; Mitrovic and Martin 2007). Learners who are good at using their metacognitive skills perform better than those who are unable to use such skills actively (Garner and Alexander 1989; Schraw and Dennison 1994). NDL emphasizes the importance of actively using and enhancing these skills during an interaction between the learner and the system. Figure 1 shows the dialogue session after a few dialogue moves encompassing domain-specific reasoning. Once the learner is able to answer the domain specific questions to an acceptable standard, the system requires them to summarize their answers and reflect upon their discussion with the system. This is done to reinforce their understanding and encourage self-assessment.

The dialogue session in Fig. 1 highlights a major feature of NDL that distinguishes our approach from the current implementations of OLMs. The system engages in a

T: Great! Now can you summarize your answers in one paragraph?

S: A Stack is a Data Structure that uses the First In Last Out scheme. It allows us to put an element on the top of the stack by using the Push function. We can only remove the Top element of the stack by using the Pop function. When an element is popped it reduces the size of the stack, so it has one less element.

T: Good job! Do you see the difference between the flow of logic in your initial answer and the last answer?

S: Yes! My first answer was not complete. I think I missed out some important details. The last answer looks much more complete and I think it covers the important points about a stack. I think it is a proper definition of a stack.

T: Why do you think you were not able to come up with a complete answer the first time?

S: I thought that was a complete answer. But now I realize that it wasn't. I think I was not sure about the requirements of your question. But as you asked me further questions, I realized that I had only given a partial answer. I knew what stacks are, but I was not able to synchronize the details. When you asked me to summarize my answer, I realized that I had known all the details about a stack, but I was not able to combine them into a complete answer.

Fig. 1 Sample NDL dialogue (Reflection Phase)

domain discussion if the learner is unable to justify the change they made in their belief base. The domain discussion phase is used to analyze how much the student knows about a specific topic. If a learner is more knowledgeable or has improved/increased their knowledge they are able to answer the question within the first attempt. This provides the system with the information about their knowledge level in the topic. For less knowledgeable students who are not able to answer the question according to the defined standard (criteria), the system engages in a series of funneling questions in order to understand their level of understanding/knowledge of the topic. For such students, at the end of the domain dialogue session, the system explicitly encourages them for self-assessment by asking them to reflect upon the past interaction and evaluate how the discussion helped them formulate their final answers.

Identifying Learner's States

All ITSs aim to engage learners to maximize learning; however a learner's engagement highly depends upon the affective and behavioral state they are in (Lehman et al. 2008). If a learner is in some sub-optimal state, the system needs to diagnose such states in order to help a learner move into an optimal state that is more conducive to learning. When a learner is in an optimal state of learning, they are more focused and learn better. Hence the system needs to ensure that such a state is maintained. There is an abundance of literature on modeling affect and motivation with varied views (Afzal and Robinson 2011; Bursell and Picard 2007; Conati and Maclaren 2009; D'Mello and Graesser 2012; Woolf et al. 2010). However it is agreed that an exact estimation of such states is not required in practice as the main focus of an ITS is to improve the cognitive state of a learner, and the knowledge about these states support the system in its reasoning process (du Boulay et al. 2010).

The process of learning requires the learner to be interested, motivated and confident to engage in a productive discussion with the system. Table 1 shows a list of Affective & Behavioral states that were selected to be used in NDL to model the affective/behavioral state of the learner. These states have been selected from previous research on the subject (Lehman et al. 2008; du Boulay et al. 2010), and they provide a good approximation of the learner's mental state. How these states were shortlisted will be discussed in the experiment section of the paper. The precision of modeling these states is not of principal importance, but an approximation of these states can allow the system to engage the learner more actively.

Affective states are related to emotions or feelings and therefore are more prominent during the domain-independent discussions where learner responses are generally influenced by how they are feeling. On the other hand behavioral states are related to the interaction of the learner and hence domain-dependent discussions are mostly influenced by the behavioral states of the learner. Metacognitive states of a learner are more difficult to gauge as they are implicit in nature and are used subconsciously. However, understanding the *context* of a dialogue can help in estimating the approximate metacognitive state of the learner. Further discussion about these states will be continued in the Wizard of Oz experiment section.

System Architecture

We propose the use of Interest-Based Negotiations (IBN) (Fisher 1983) in NDL. IBN aims at exploring underlying interests of the parties rather than their negotiating positions and considers negotiating parties as allies working together for mutual gain, which is the essence of the negotiation process.

Since negotiation is a process of understanding, we make use of IBN to generate the dialogues in NDL. To realize the envisioned interactions in our system we extend the computational model proposed in (Tao et al. 2006) on the automation of IBN. Our system consists of the following functional components as shown in Fig. 2:

State Reasoner: handles all the state-related tasks. It generates the State Model (SM) for the learner by translating learner inputs to the corresponding affective and behavioral states. The State Updater (SU) updates all these state in real-time with each transaction. It also stores previously held states of the learner to understand learner progression.

Table 1 List of selected affective & behavioral states of learner in NDL

Affective states	
Confused	Poor comprehension of material, attempts to resolve erroneous belief
Frustrated	Difficulty with the material and an inability to fully grasp the material
Engaged	Emotional involvement or commitment
Behavioral states	
Confident	The feeling or belief in one's abilities or qualities
Interested	Wanting to know or learn more about something
Motivated	Having a motive or incentive to perform an action

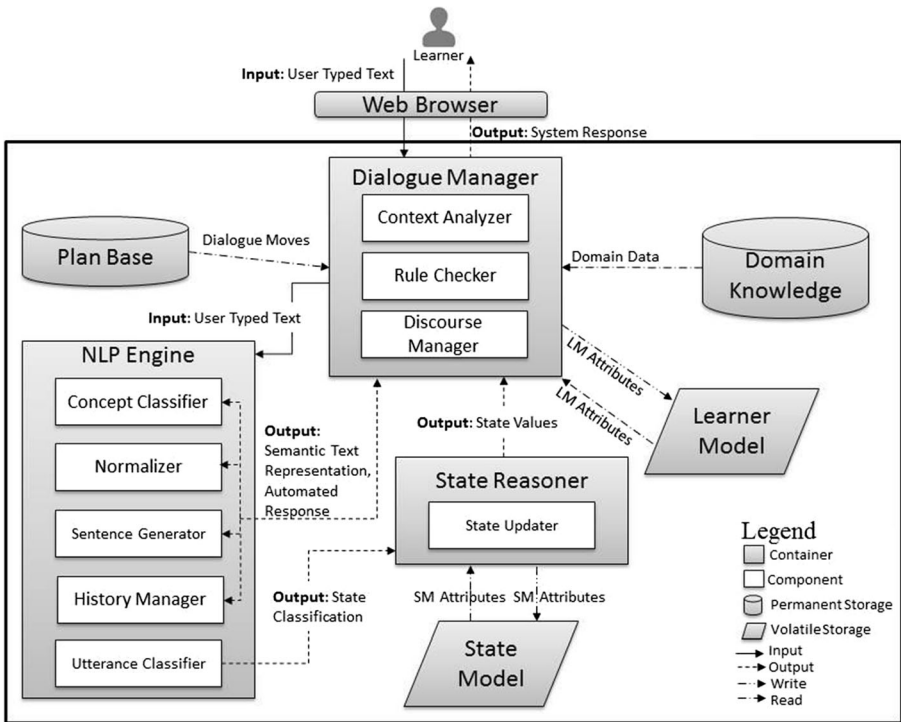


Fig. 2 NDL system architecture

Dialogue Manager: consists of the Rules Checker (RC) which is an inference engine and uses the information from the SM in conjunction with the LM in order to select the next system move with the maximum utility according to the current context. The Context Analyzer (CA) submodule uses the information from the SR and the NLPE in order to articulate the current context. It also consists of the Discourse Manager (DiM) that controls the flow of the overall dialogue.

NLP Engine: this is the core module for providing a Natural Language interface to the learner. NDL does not require a complete NLP understanding as we are interested in the concept-level cognition of the learner’s input. To accomplish this, the NLPE consists of submodules which include:

Concept Classifier: uses a Normalized Distance Compression algorithm to return a list of concept identifiers that most closely match the learner input.

Normalizer: manages stemming and spell checking for the learner input.

Sentence Generator: uses the concepts identified along with the current context to generate a list of possible utterances of the system. These possibilities are matched with the response library and the best matching phrase is selected to generate sentences automatically.

History Manager: stores information about the concepts used by the system and the concepts expressed by the learner. This information is passed to the RE, which uses it to classify the current context.

Utterance Classifier: uses a Cosine Similarity Index algorithm to return a list of state classifiers that are identified from the learner input.

Plan Base: holds the different negotiation moves available to the system according to a specific context. The information regarding the consequences of using a move in a specific context and state are used to update a move's adequacy to that context in the PB.

Designing Dialogues for NDL

Realizing interactions such as the one shown in Fig. 1 requires that the system not only understands the learner's characteristics but is able to comprehend their answers to provide a proper response. In NDL we wanted to introduce a more flexible, open, and natural method of interaction between the learner and the system. The use of chatbots has been documented to ease the negotiation process and improve engagement levels (Kerly and Bull 2006; Kerly et al. 2008b). In light of these previous studies on the use of chatbots in OLMs we put forward the following questions for ourselves:

- Q1. Can a conversational agent provide a more natural and flexible negotiation interface to the learner than a menu-based system?
- Q2. What kind of dialogue moves would be required to facilitate such a negotiation?
- Q3. What will be the challenges of implementing such a chatbot?
- Q4. Which emotional states of a learner we need to pay attention to for realizing usable IBN-based dialogues?

To find the answers to these questions, we conducted a Wizard-of-Oz (WoZ) experiment. Natural language dialogue is complex in nature and the interaction patterns differ from learner to learner. Such inconsistencies were to be faced in negotiating the LM with learners; therefore, we required empirical data in order to support our system design. The WoZ approach has been shown to be valuable for collecting data in scenarios which require complex interactions between the users and the systems (Dahlbäck et al. 1993). Our experiment design was based on the structure and guidelines on conducting a WoZ experiment provided by the previous study on CALMsystem (Kerly et al. 2008b). Building upon the findings of the previous studies our experiment design included a self-annotation mechanism for students to annotate their input according to the option they think best describes their current emotional state. We also used the findings of the previous study to generate a list of possible outcomes/markers that could be related compared afterwards. Since in the WoZ experiments, users are under the impression that they are interacting with a system, many application-specific characteristics of a textual dialogue can be elicited.

For this experiment we created an independent OLM. The domain of "Data Structures" was used for this experiment. The system gave a multiple-choice questions test to capture their understanding. These test scores were used to analyze the performance of each student and the results were used to generate the learner model. At the end of the test, the learner was allowed to update their belief base about their knowledge in the corresponding topic. This allowed for the wizard to initiate a dialogue

in the case of a conflict occurring between the system's set of beliefs and the learner's set of belief. Ensuring a mixed-initiative dialogue system, the participants were also allowed to initiate a dialogue with the system by themselves at any time. During their interaction with the system, the participants typed their inputs and were required to annotate each input according to a drop-down list of states provided to them (self-annotation). They had the liberty to select multiple states which they thought best represented their mental state or they could provide a new/different state not available in the list.

Wizard of Oz Experiment

The study was conducted with the students of Bahria University, Islamabad, Pakistan. A total of 45 students from the fourth semester of the Software Engineering course participated in the experiment. All participants had completed the compulsory courses of computer programming (C++, OOP, and Data Structures) as a course requirement. The first author acted as a secondary experimenter while the experiment was conducted and supervised by the local instructor (Senior Lecturer in the SE department). The author was available via an online connection throughout the duration of the experiment. The participants were given an introduction to ITSs and OLMs by the secondary experimenter through a Skype video conferencing session. The session included an introduction to the aims and objectives of ITSs and their real-life applications. The participants were also introduced to the different categories of OLMs and were shown the interfaces and interaction possibilities provided by some OLMs, specifically Mr. COLLINS (Bull and Pain 1995), STyLE-OLM (Dimitrova 2003) and CALMsystem (Kerly et al. 2008b). An initial survey was conducted to understand their expectations from such a system.

The participants were provided with a web interface to interact with the system. All interactions between the system and the participants were logged and the interaction transcripts were stored for future analysis. Once the participants had completed their sessions with the system, another survey was conducted to get their feedback about the system and the interaction possibilities it provided.

The participants were randomly divided into three groups; one uncontrolled group and two controlled groups. This was done in order to ensure that the system responses generated during each phase would be valid enough for a diverse group of learners. The experiment was conducted in three phases where in the first phase with the uncontrolled group, there was no negotiation protocol set for the wizard. The wizard conducted open-ended dialogues with the participants without following any set of rules. The dialogue scenarios captured in these interactions were translated into IF/THEN clauses in order to generate the initial 'rules library'. The interaction logs were also used to generate a corpus for system responses that constituted the first version of the response library. Figure 3 shows a screenshot of the response library available to the wizard.

In order to generate the response library, the protocol discussed previously was used to classify system utterances. The strategies used are:

1. ASK for JUSTIFICATION: ask to justify a response/claim.
2. GIVE JUSTIFICATION: provide justification for the last utterance/action.
3. ACCEPT JUSTIFICATION: accept the claim if it is justified.
4. REJECT JUSTIFICATION: reject a claim if it is not justified.

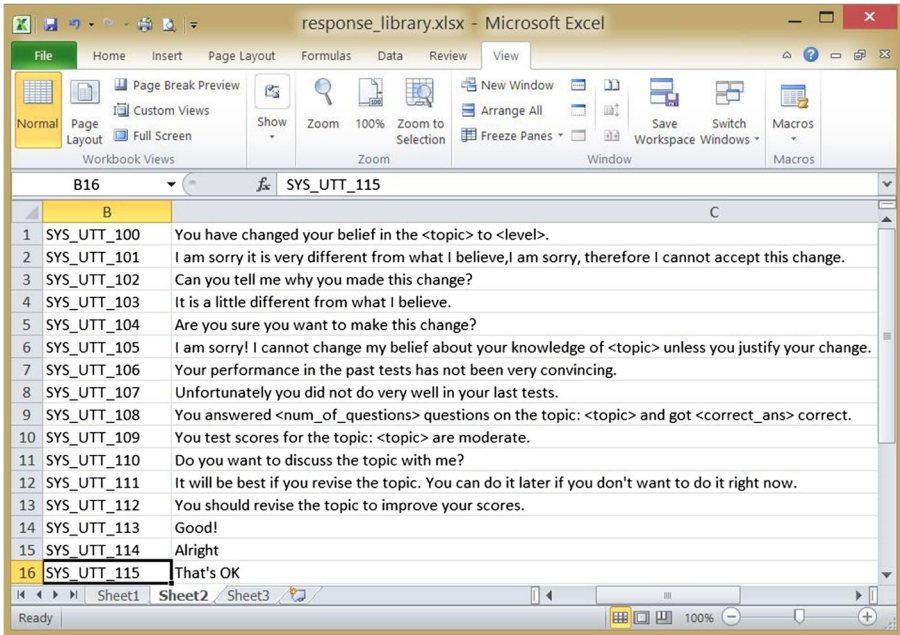


Fig. 3 Response library

5. GIVE PROPOSAL: propose an alternative solution
6. ACCEPT PROPOSAL: accept a proposal.
7. REJECT PROPOSAL: reject a proposal.
8. PROVIDE FEEDBACK: provide feedback corresponding to the last action.

Both the rules library and the response library were saved in MS Excel file for quick access to an appropriate response to the learner. Each system response was given unique identifier SYS_UTT_#, where '#' was a unique numerical value. This allowed the wizard to only select and copy/paste the corresponding system utterance in the next phase.

The second phase with the controlled group-1 was conducted under controlled conditions where the wizard used the rules and response libraries generated from the analysis of the interactions in the previous phase under the protocol guidelines to respond to student inputs. During this phase there were certain scenarios which did not occur in the previous phase and hence had no corresponding rules in the rules library to select an appropriate response from the response library. In such situations, the wizard had the liberty to improvise the response and such a situation was highlighted for future analysis of the dialogues.

The third phase was conducted with the controlled group-2. The interaction logs of the first two phases were used to update the rules and response libraries. The analysis of the first two phases allowed for the improvement of the rules and response libraries for the wizard by including missing rules and responses for new dialogue scenarios. The third phase of the experiment was almost completely automated with 85 % of the wizard responses being generated by using the rules library. The results from this phase were again used to update the libraries to accommodate missing rules or responses.

Table 2 shows an example of a rule used by the wizard in order to select a corresponding system response.

The students were divided randomly in three equal groups for the three phases of the experiment. This meant that for each phase, we had 15 students interacting with the wizard. Each group had a single interaction session with the wizard. The interactions in the first phase were the longest with an average interaction time of 33 min. As there was no set negotiation protocol, this meant that the students and the wizard indulged in a very open discussion. The interaction times of the second phase were considerably shorter as a negotiation protocol was introduced and the discussion was more directed. The average interaction time in this phase was 20 min. The third phase saw the shortest interactions as it used the formalized rules and response library. Average interaction time for the third phase was 16 min. All of the interactions were concluded successfully with the student either accepting the wizard's proposal or retaining their initial stance about their knowledge level.

Results

The interaction logs and the conversation transcripts from the WoZ experiment were transcribed and analyzed in order to understand the kind of dialogues the participants engaged in with the system. In the 45 conversations between the student's and the wizard there were a total of 195 negotiation fragments. The number of user initiated conversations was 80. The mean interaction time was 27.4 min. Off-topic discussions or small talk constituted 13.4 % of all conversations. 45.6 % of the conversations were related to domain-specific discussions while the remaining 41 % conversations constituted the inputs used to approximate learner characteristics.

While off-topic conversation during a tutoring session may be seen as counter-productive to the construction of knowledge, it has been found to be an effective strategy to keep the learners engaged. Expert tutors utilize off-topic conversations in scenarios where the learner seems to be disengaged or frustrated. It is seen as useful strategy to build a sense of trust and empathy using a dialogue that does not require the learner to recall domain or task-oriented knowledge. Having the ability to engage at a certain level of small talk allows the system to provide responses to user inputs that are not related to the domain or the task at hand. This gives the system the ability to hold more naturally flowing dialogues with the learners.

Classifying Student's Affective and Behavioral States

Affect relates to the emotional reaction (feeling) one has towards an attitude object (learning task). For example, if a student is confused about a mathematical concept

Table 2 Sample Rule for wizard to select system response

IF	User has changed their belief in topic and the difference between their belief value and the system's belief value is greater than 2
THEN	Highlight User Change: {SYS_UTT_100} REJECT CHANGE: {SYS_UTT_101} ASK for JUSTIFICATION: {SYS_UTT_102}

(attitude object), whenever they are exposed to a problem related to that concept, they feel confused. Behavior relates to how one behaves when exposed to an attitude object. Considering the previous example, if the student is confused about a concept, they are most likely to avoid it and be less interested in taking on the problem.

There are many unknown categories of learner’s mental states and an in-depth evaluation of all these states was out of the scope of our study. For the initial classification of the participant’s affective states we used Ekman’s six “basic” emotions (Ekman 1973) and a set of learning-focused affective states identified in (Graesser et al. 2006; D’Mello et al. 2007) for our study. The list of affective states that was used for this study include: *confusion, engagement, frustration, curiosity, eureka, surprise, anger, fear and sadness*. Similarly for the classification of behavioral states, we used only the “states” identified in (De Vicente and Pain 2002) based on the theories of motivation in education (Malone and Lepper 1987; Keller 1983). Choosing between different states is not a trivial task; therefore, we concentrated on the states that would have a deeper impact on the outcome of an interaction. We limited our study to the states that characterize a student’s behavior while interacting with a human tutor which include: *confidence, interest, satisfaction, effort and motivation* along with their negative dimensions. The occurrence frequencies of the states shown in Figs. 4 and 5 were used as the measure of acceptance which narrowed the affective states list to; *confused, frustrated, and engaged*. Whereas the behavioral states selected were; *confident, interested and motivated*.

The interaction logs generated during the experiment consist of self-annotated typed input by the participants. There is no gold standard for understanding and detecting the mental state of a learner from an interaction log. To this effect we employ the Multiple-judge strategy (Graesser et al. 2006) to manually annotate the interaction logs. The judges included the participants (self-annotations) and two expert judges (assistant professors) and two intermediate judges (lecturers). One of the expert judges was a professor of psychology while one of the intermediate judges was a lecturer in linguistics. This selection of judges provided us with a diverse pool of experience which was very helpful during the discussions over the annotated utterances. The judges were provided with the

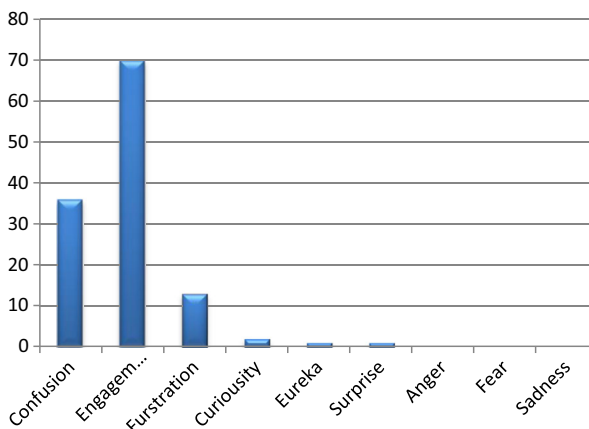


Fig. 4 Occurrences for each affective state

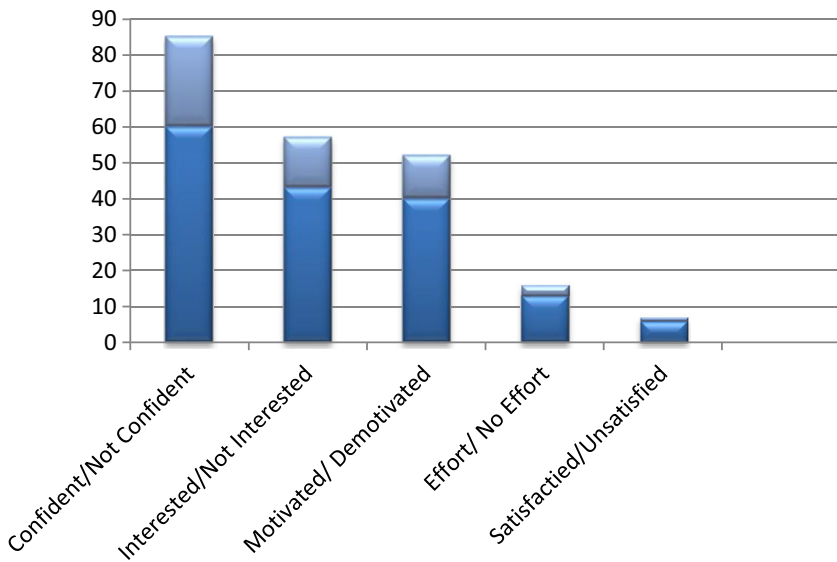


Fig. 5 Occurrences for each behavioral state

learner interactions along with the list of affective and behavioral states classified for this study. They were also given the liberty to add a new state if they deemed necessary in order to capture the approximation of the participant’s mental state. We are aware of the subjective nature of this classification scheme which might not reflect the true mental state of a learner. However, we have previously emphasized that an approximation scheme can be considered sufficient to control the flow of the dialogues. An incorrect classification of a learner state does not drastically impede the dialogue course as the system uses the context and dialogue history to ensure an effective flow of the dialogue. We will discuss this topic in the evaluation section below.

The judges were provided specific guidelines for annotating the transcripts. They were required to highlight any markers in the student’s input that might point towards a specific attribute of their mental state. For example using “*Ummm...*” in the beginning of an utterance was classified by tutors as a sign of “low confidence” or “guessing”. A similar “vocal” sound is associated to a thinking person. However, it was noticed during the experiments that when the students were thinking, they did not type “*ummm*”, but rather made the vocal sound. Another important aspect of annotation was the consideration of “context” while annotating the transcripts. Context plays an important role in helping to decipher the rationale behind a specific utterance and in most cases the thought process involved. For example, if a student is asked a question related to the domain and they answer;

“I don’t know.... But I think it is”

This input from the student is treated as “confused” and their answer is “not confident” but he is considered to be “interested” as he is trying to answer the question. Similarly, the very basic utterance “OK” can have multiple meanings which can be

elicited if the context in which the utterance occurs is known. The strategy to highlight markers in text and convey a context was very helpful in fine-tuning the rules in the library.

The annotated transcripts from the judges were compared with each other to find the matching and conflicting annotations. The list of conflicting annotations was discussed with all the expert judges in order to reach a consensus regarding a specific learner utterance and its relation to a specific affective or behavioral state.

The self-annotated lists of the participants were then matched with the agreed upon judge's annotated list in order to generate a list of student utterances classification according to the affective and behavioral states. A list of utterances with no matches, or mismatches was also generated during this process. These lists were deliberated upon by the judges in order to remove any discrepancy between the annotated values. As mentioned previously, the panel of judges included an expert tutor of psychology and a lecturer in linguistics. This diversity of experience helped the panel to annotate utterances mismatching annotations to generate a complete list.

An interesting observation during the analysis of the inputs was the positive and negative dimensions of the specified states and how they affected the course of the dialogue. It was observed that in case of affective states, a negative affective state required more system involvement than a positive affective state. For example, if a learner was confused (negative state), the system had a better opportunity to help him realize his confusion than when he was not confused (positive state), in which case the system intervention was minimum. Contrary to this, the dimensions of behavioral states played a much greater role in the interactions between the learner and the system. A confident learner reacted differently than a learner who was not confident. It was observed that both positive and negative dimensions of behavioral states impacted the system's interactions with the learner.

Inputs Related to Affective States

In Fig. 4 we can see the distribution of the occurrence of the affective states in the learner inputs. These occurrences were calculated by comparing the tutor's annotated data from the experiment with the self-annotated data of the learners during their discussion with the system. The findings were consistent with previous study (D'Mello et al. 2007). The most often occurring affective states were selected and identified as *confusion*, *engagement* and *frustration*. The rest of the affective states were almost non-existent in both tutor and participant annotations. Table 3 shows a list of a few learner inputs and their corresponding affective state.

Inputs Related to Behavioral States

In Fig. 5 we can see the distribution of the occurrence of the behavioral states in the learner inputs. The states with the highest occurrence frequencies i.e. *confident*, *interested* and *motivated* were selected for the classification of the learner utterances. Frequencies of the corresponding negative states were also added to the chart to show the number of occurrences along both positive & negative dimensions. Table 4 shows a list of learner inputs and corresponding behavioral states.

Table 3 Examples of affective states corresponding to user inputs

User Input	Affective State
I don't understand	Confused
No! I still don't understand	Confused, Frustrated
I don't know	Confused, Frustrated
I don't need your help!	Frustrated
What is this?	Confused
How?	Confused
I can't do this	Frustrated
Wow I did it!	Engaged
Yes, I think I got it	Engaged
I know it	Engaged

Revisiting Questions Set for the Experiment

At the end of the experiment, we analyzed the user interactions with the system, the observations made by the authors during the experiment and the discussions with the tutors' panel while annotating the learner utterances, in order to answer the questions we set for ourselves before the experiment. The first question we put forward was:

- Q1. Can a conversational agent provide a more natural and flexible negotiation interface to the learner than a menu-based system?

The participants had never used an ITS before and therefore they did not have a hands-on experience of using a menu-based OLM. However as mentioned earlier, in the pre-experiment setup, the participants were given an introductory lecture on OLMs and the interaction possibilities provided by a few OLMs. They were shown interfaces

Table 4 Examples of behavioral states corresponding to user inputs

User Input	Behavioral State
Yes I know	Confident
Ok, Yes, Yeah, Yeah sure, Sure (Context Dependent)	Confident, Motivated
I want to discuss this	Motivated, Interested
No (Context Dependent)	Uninterested
I'm not sure	Not confident
I don't think so	Not confident
I don't want to ...	Uninterested
I can't do this	Demotivated
I want to solve this	Motivated
Can you help me?	Interested
Let's talk about something else	Uninterested
Not now	Uninterested

and interaction fragments of previous system in order to familiarize them with the concept and applications of OLMs. In the post-experiment survey the participants noted that the natural language dialogue conducted by the wizard was a very natural and realistic form of interaction as it closely related to some form of chat messaging provided by many SNS and SMS application they use to in their daily lives. Majority of the participants were of the opinion that using a natural language negotiation approach would allow the student to interact with the system more openly. On the question of replacing the natural language interface with a menu-based interface, most of the participants answered in the negative as they thought it would make them feel controlled and confined. The authors are aware that the results from this question do carry a bias as the participants had no prior experience of a menu-based system. However, previous WoZ experiments have concluded that students did prefer a chatbot over a menu-based interaction system (Kerly et al. 2007). Now we consider the second question:

Q2. What kind of dialogue moves would be required to facilitate such a negotiation?

The answer to this question was investigated during the analysis of the results of the experiment. The negotiation protocol that was provided to the wizard proved to be sufficient in handling the course of negotiations from different participants. It was noted that apart from following the negotiation protocol, the system also needs to be able to handle a fair amount of off-topic discussions or small talk. This was in-line with the findings and guidelines provided by previous WoZ experiment to study the use of a conversational agent in an OLM (Kerly et al. 2008b). This became more evident in the interaction of less interested/motivated participants. However it was also noted that almost all participants did engage in some form of small talk with the wizard during their interactions. Therefore, the discourse manager not only needs to follow a negotiation protocol, but also needs to be able to deal with small talk initiated by the learner, or in some cases initiated by the system itself in order to engage the learner and keep continuity in the discussion. Another finding that resonated with results of previous WoZ experiment is consideration that the system should be able to keep track and control the level of small talk during an interaction. This will be essential to ensure that the learner does not spend too much time off-topic.

Q3. What will be the challenges of implementing such a chatbot?

This question relates to the challenges we could foresee for our system after conducting and analyzing the experiment. Our findings were in line with previous work on the use of chatbots in OLMs (Kerly and Bull 2006; Kerly et al. 2008a; Kerly et al. 2007). The more prominent challenge was the implementation of a natural language interface that will be able to handle vast array of user inputs. The research on Natural Language Processing (NLP) has been continuing for years and there is no single, best NLP approach that can be used to generate a 100 % realistic dialogue environment. Keeping this limitation in mind, we needed to decide upon the tradeoff between the usability of an NLP technique and its complexity. Spending too much effort and time on implementing the NL interface would negatively affect the scope of

the project. Hence it was decided to keep the complexity of the NLP to a minimum and with each development iteration try to improve upon it.

The second challenge that was acknowledged was the complex nature of learner states and identifying such states automatically. Since we will not be using any sensory information and only use the typed input to generate an approximation of the learner's state, this will make the task simpler but the accuracy of the resulting states will remain questionable. Further research will be required in this respect in order to maximize the usability of the state model at acceptable cost.

The third challenge identified by the analysis of the experiment was related to the user experience. Learners with different knowledge level use the system different and their interaction patterns also vary significantly. A chatbot in a learning environment needs to be able to adapt to this change in character and keep the learner engaged and on topic. As identified earlier, small talk can act as a good strategy to bring back the learner who loses interest, but the system needs so ensure that the small talk should not hinder the learning process.

Lastly the experiment also gave insight to the problem of authoring a chatbot script from scratch. Most of the current chatbot implementations use specialized script formats that increase the learning curve and require some time to generate. This is normally due to the fact that domain-dependent and domain-independent dialogue fragments are merged into the same script. To minimize this complexity a scheme that separates the domain-dependent and domain-independent utterances and uses a mechanism to merge them at runtime would allow tutors to concentrate more on the domain-dependent section of the chatbot. This would result in faster development times and maintenance tasks would be more simplified.

Q4. Which emotional states of a learner we need to pay attention to for realizing usable IBN-based dialogues?

The Wizard of Oz experiment aided us to short-list the affective and behavioral states that were most prominent during the interactions between the participants and the wizard. As discussed in detail in the previous sections, the deliberation between the experts over the annotated logs allowed us to finalize a list of three affective and three behavioral states that will be used to control the flow of the dialogue in our system. The sufficiency of these six emotional states will be evaluated in the first evaluation study discussed below.

Phases of NDL

NDL acts as a component of an ITS which is normally triggered by an event such as a conflict between the learner's belief in their knowledge and the system's belief about the same. If necessary, an NDL dialogue may also be initiated explicitly by the learner by using the <DISCUSS > command followed by the "topic" they wish to discuss with the system. Every NDL dialogue is comprised of three phases: Initialization, Domain Discussion, and Reflection. Dividing a dialogue into three phases makes it easier to handle the different inputs expected during each phase, hence the rules can be specifically written for a phase.

Initialization Phase

The first phase of an NDL dialogue is the initialization phase. This phase is used to initialize the values of the State Model for the current dialogue as well as set a foundation for the dialogue to follow. In the initialization phase, the system tries to understand what triggered the dialogue. If a conflict is the cause of a dialogue being triggered, then the system tries to identify the reason behind the learner's change.

During the initialization phase, the system asks the learner for a justification for their action. A weak justification is challenged by the system until the learner is able to justify himself. The values from the State Model are used to select the next move of the system during this phase. The initialization phase can lead to two possible outcomes for the learner; domain discussion or take a test to prove their claim. A learner who chooses to take a test is still encouraged to discuss the topic with the system before they take the test.

Domain Discussion Phase

Once a learner agrees to discuss the topic with the system, the domain discussion phase is initiated. This phase is directly related to the domain knowledge of the learner and discusses the selected topic with the learner using the natural language interface. The discussion starts with a focal question about the topic. The learner responses are classified according to the Utterance Classifiers (UT) and annotated with the State Model (SM) values. A complete description of how this is done will be discussed in a separate paper. Once a response has been classified, the system uses this information to generate the system response accordingly. Each focal question in the domain discussion phase has a list of attributes related to it, which include:

- List of Correct Answers according to the degree of completeness i.e. (EXPERT, INTERMEDIATE, and NOVICE).
- List of concepts related to the topic that constitutes a good answer.
- List of common misconceptions related to the topic/concept.

When a learner's input is classified as an answer, it is firstly matched with the misconception list. If the learner's answer contains any misconception, the system initiates a remedial dialogue which focuses on the identified misconception and the funnels through the related concepts in order to identify the cause of the misconception. If the answer does not contain any misconception, then the system matches the content of the input with the list of correct answers to score the degree of completeness. Finally the learner's input is tokenized in order to match the concepts related to the questions. The scores of the learner's answer and concept coverage are then used to calculate a final score. This score is used by the RC in order to match the corresponding rule in the rule-base to generate the next system response/move. As mentioned before, the implementation details are not a part of this paper and will be discussed in a separate paper. The domain discussion phase is completed once the learner is able to provide an answer with an acceptable degree of completeness along with a medium to high concept coverage.

Reflection Phase

The last phase of the NDL dialogue is one of the core features that distinguish our system from the current OLMs. The reflection phase of NDL is initiated at the end of the domain discussion. This phase is utilized for the explicit reflection for the learner and does not discuss the domain rather the dialogue which just occurred between the learner and the system. During this phase the system engages the learner into a dialogue that encourages self-assessment in the learner. The system discusses the learner's final answer with respect to his initial answer and the discourse that led them to it. The learner is asked open-ended questions that make them compare their answers and assess how they were able to improve upon them. There is no correct or wrong answers in this phase, however the system does keep a track of the learner responses and uses the dialogue history and the learner's verbosity in order to advance the dialogue.

This phase allows the learner to reflect upon their discussion with the system. They are encouraged to identify the causes of confusion and how they were able to clarify them. This phase offers the learner a series of questions explicitly targeted towards self-reflection and evaluation. Moreover the learner is also encouraged to evaluate their learning strategy and how they can improve it. The discussions from this phase are saved and are available to the learner as assessment logs they can access at any time.

Learner Interaction in NDL

The data collected in the WoZ experiment helped us in identifying the different state-transitions that are likely to occur during a dialogue session in NDL. As mentioned previously, the third phase of the experiment was almost completely automatic. To show the kind of dialogues the system is expected to hold with the learners, we will show an example of an interaction between a learner and the wizard from the final phase of the experiment.

Example: A confident student who is confused about his knowledge level but is interested to discuss his beliefs with the learner

The learner in this example completed the topic of Stack and gave a test where he could not perform very well. The system updated his LM to reflect a “Medium” understanding of the topic. The learner then changed his belief to “Very High”, which was very different from that of the system. Therefore the wizard interrupted the proceeding by initiating the following dialogue:

Initialization Phase

-
- [I.1] Tutor You have changed your belief in the topic of Stacks to Very High.
I am sorry but this is very different from what I believe, therefore I cannot accept this change.
Can you tell me why you made this change?
- [I.2] Student I think I know more now
- [I.3] Tutor That's great!

However your performance in the topic of Stack has not been very convincing.
Do you want me to show you an overview of your performance in the topic of Stack?

[I.4] Student Ok

[I.5] Tutor You answered 10 questions on the topic of Stacks and got 03 answers correct.

[I.6] Student Yeah... but I'm better prepared now

[I.7] Tutor That's great!

In order to justify your change, you will have to take a Test to improve your score.
Before you take the test would you like to discuss the topic of Stack with me?

[I.1] When the system identified a conflict between the belief values of the learner and the system, the chatbot interrupted the proceedings by initializing the dialogue. The first utterances of the system are a direct cause of the user action i.e. change made to their belief base. The system highlights the value set by the learner which caused the conflict. The system then asks the learner to justify their claim.

[I.2] The learner response is classified by the system as a “weak_confident_justification”. The learner claims that they know more which is classified as a “confident act” but the use of “I think” in the start of their utterance is classified as a “weak act”.

[I.3] since the learner showed some confidence, the wizard starts with a positive feedback to reinforce their confidence. The learner's justification has been marked as weak and their claim cannot be justified the system rejects their justification. As the current behavioral state of the learner has a positive value, the system gives the proposal to view their performance.

[I.6] The learner's utterance is classified as confident and interested as he continues to show interest in defending his claim.

[I.7] as the system already provided the learner an overview of their performance and the learner did not revert back his change, the system gives him the proposal to discuss the topic in order to gauge his knowledge level. Discussing the topic instead of directly giving the test allows the system to highlight the weaknesses, confusions and misconceptions the learner might have about the topic during the discussion.

Domain Discussion Phase

[DD.1] Tutor What do you know about a Stack?

[DD.2] Student It is a data structure that we use to hold elements. We can take elements in and out of it

[DD.3] Tutor Can you elaborate on that?

[DD.4] Student It is a data structure that holds elements. We use operations to take elements in and out of it.

[DD.5] Tutor Ok.

A stack is an abstract data type that stores elements in reverse order.
What can you tell me about the principle behind storing elements in a stack?

[DD.6] Student The first element we put in the stack is the first element that we take out.

[DD.7] Tutor So you think the first element we put in the stack is the first element that we take out.
I'm afraid that's not correct.

Can you explain the order in which elements are stored in a stack?

[DD.8] Student On top of each other

[DD.9] Tutor Good!

Which of the following is an example of such an ordering scheme:

- A) A pile of books
- B) A line at a bank counter

[DD.10] Student A

Misconception funneling continues...

[DD.17] Tutor That's right!

Question So let's revisit the question:

What can you tell me about the principle behind storing elements in a stack?

Missing concepts funneling continues...

[DD.25]. Tutor Great! Can you summarize your answer now?

[DD.26] Student A stack is an abstract data type that stores elements in reverse order. The principle of storing elements is called LIFO. It states that the LAST element INSERTED in the stack is the FIRST element that is taken OUT of it. It has 2 major operations Push and Pop. Push operation is used to insert an element on the stack and Pop operation is used to remove an element from the stack.

[DD.3] This system utterance is a result of the following rule being triggered:

IF user utterance is an answer and the answer does not contain any misconception and the answer score is low and the concept coverage is low and it is the first try by the learner

THEN prompt the learner to elaborate on their answer.

The above rule is intended to allow the learner to rethink about their answer and add more details to it. This encourages the learner to think about the knowledge they have and decide what new information they can add to their previous answer.

[DD.5] This system utterance is a result of the following rule:

IF user utterance is an answer and the answer does not contain any misconception and the answer score is low and the concept coverage is low and number of tries by the learner is greater than 1

THEN give the learner a neutral feedback and provide a partial statement about the answer start funneling the missing concepts

As described previously, each topic has a set of concepts related to it. When the learner's provides an answer, their answer is tokenized and matched with the list of topics which generates a set of missing concepts. The system then starts to funnel through these concepts one by one in order to help the learner build a concrete answer.

[DD.7] The rule behind this system utterance is as follows:

IF user utterance is an answer and the answer contains a misconception

THEN repeat the learner's answer and provide negative feedback and funnel misconception

Just as every topic has a list of concepts associated with it; it also has a list of common misconceptions. When a learner's answer is matched with a misconception,

the system highlights the misconception by repeating the learner's answer. The system then start an inner loop of funneling through the concepts related to the misconception.

[DD.17] once all of the concepts related to the misconception have been discussed/covered, the system then asks the learner the same question which triggered the misconception. This is done in order to reinforce that the learner is able to answer the same question without the misconception.

[DD.25] As with the case of a misconception, once the system and the learner have discussed all the missing concepts in the learner's initial answer, the system asks the learner to summarize their answer. This gives the opportunity to the learner to collect all the information into one concise answer. The learner's answer is again scored and checked for any misconceptions or missing concepts. If any such irregularities are found, these are highlighted by the system and a corresponding feedback and proposal are generated depending upon the previous dialogue moves.

Reflection Phase

-
- [R.1] Tutor Let's reflect upon the discussion we just had in order to understand what we gained from it. How does your initial answer compare with your final answer?
- [R.2] Student My first answer was not complete. I missed out important details. The last answer looks much more complete and I think it covers the important points about a stack. I think it is a proper definition of a stack.
- [R.3] Tutor What was most confusing for you?
- [R.4] Student I didn't find anything confusing
- [R.5] Tutor You did hold a misconception regarding the following question:
What can you tell me about the principle behind storing elements in a stack?
Were you able to clarify your confusion?
- [R.6] Student Yes
- [R.7] Tutor What confusions do you still have that you need to clarify
- [R.8] Student I think I don't have any more confusions about stack
- [R.9] Tutor What worked well for you?
- [R.10] Student When you broke down the problem to highlight my misconception. Thinking about what I know about stacks and imagining it in my mind helped me highlight my misconception and clear out my confusion.
- [R.11] Tutor What did not work so well for you that you should change?
- [R.12] Student I need to have a more logical flow in my answers. I should think about what I know and learn to relate the knowledge I have with the question at hand.
- [R.13] Tutor Good job!
I have updated my belief base about your knowledge in the topic Stack to HIGH.
If you still want to change it to VERY HIGH, please take a TEST to improve your score.
Would you like to ACCEPT my proposal or take a TEST?
-

[R.1] This system utterance initiates the reflection phase. The system asks the learner to compare their initial answer to their final answer. There is no correct or wrong answer for this phase. However the system does react to the input given by the user in this phase.

[R.5] this utterances shows how the system reacts to the learner's answers in this phase. As the misconception flag had been raised during the domain discussion phase,

the system expected the learner to talk about their confusion of “storing principle in stacks”. However when the learner claims to have had no confusions, the system highlights the question which generated the misconception so that learner can evaluate their input.

[R.13] after the dialogue is completed, the system updates the LM accordingly. If the learner is able to provide a suitable answer to the initial question by the end of discussion, the system proposed to update its belief base by one point increment. However it also allows the learner to take a test if they still want to justify their change.

NDLtutor

Implementation

To evaluate the feasibility of the architecture we defined for NDL, the validity of the affective/behavioral states and the effects on the self-reflection and self-assessment skills of the learners, we implemented a system called NDLtutor. It is an environment that provides the learner with a natural language interface to discuss their LM with the system in the data structures domain. Since NDL can act as a component of an ITS, we implemented NDLtutor as an independent OLM. The NDLtutor diagnoses the knowledge gaps of the learner during discussions about the learner’s beliefs and promotes metacognition by using reflective dialogue strategies. One of the goals of the implementation was to develop an open source system. To accomplish this, the backend of the NDLtutor has been implemented using PHP and MySQL whereas the frontend (user interface) has been designed using HTML5 and jQuery. The backend database consists of:

- *Domain Knowledge*: The domain knowledge is stored as plain text which is divided into topics and sub-topics. Each topic has two sets of questions; 1) *Multiple-Choice Questions* (MCQs) to assess the learner’s performance 2) *Domain Discussion Questions* (DDQs) that are used to discuss the topic with the learner during the conflict resolution phase. Each DDQ has the following format:
 - *Question Text*: The question that is given to the learner.
 - *Expected Answer List*: A list of answers according to expertise level (Expert, Intermediate, and Novice).
 - *Concept List*: A list of concepts that should be part of the learner’s answer
 - *Misconception List*: a list of common misconceptions related to the question.

Misconception Funneling Questions: A list of questions to funnel through each misconception.

- *State Model*: The state model is stored as a list of attributes (states).
- *Learner Model*: The learner model is an overlay of learner’s knowledge upon the domain and is constructed/updated by using the learner’s performance evaluation during the MCQ questions as well as the discussion sessions with the system.
- *The Reflection Log*: The database also stores the learner’s responses during the reflection phase and this act as a self-assessment log for the learner to review at any time.

Interface and Basic Functionality of NDLtutor

Figure 6 shows the interface of the NDLtutor. The interface is divided into three columns. The left column contains the learner's own belief base and the LM generated by the system. The learner is allowed to change their belief base using a drop-down list. The middle column contains the section which provides the MCQ tests to assess the knowledge level of a learner in a topic. These results are used to generate the LM of the learner. The right column provides the learner with the chatbot interface. The chatbot provides the following modes of interaction:

1. *Conflict resolution*: This form of interaction is initiated by the system when the learner's change generates a conflict between the belief base of the system and the learner.
2. *Discussion*: The chatbot allows the learner to initiate a discussion about a topic by using the DISCUSS keyword. The result of this discussion is reflected in the system's LM.
3. *Help*: The learner can also ask for quick explanations using the HELP keyword. This functionality allows the learner to use the chatbot to quickly search for terms/concepts they want to know more about. The system in help mode acts like a glossary and provides the basic definition for that term.

NLPE Class

It provides the Natural Language Processing functions for the NDLtutor. It defines the following member functions:

- `Tokenizer($string)`: breaks the learner's input text into words (tokens).
- `Normalizer($string)`: stems the input using the Porter stemmer algorithm.
- `Utterance_Classifier($string)`: uses Normalized Compression Distance (NCD) algorithm to match the learner's input with the `User_Utterance_Library`. It returns a set of

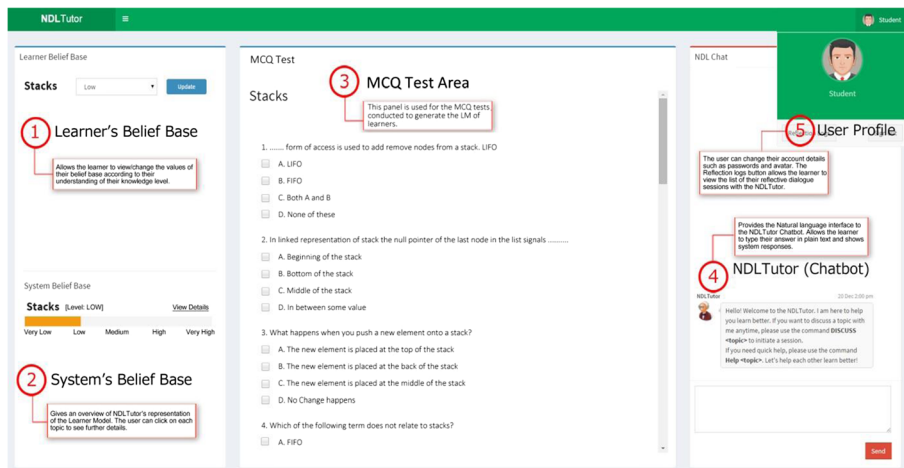


Fig. 6 NDLtutor interface

learner states and a unique utterance classifier. We have tested the `Utterance_Classifier()` with different threshold settings and benchmarked 60 % as the minimum score to qualify for a match.

- `Answer_Scorer($string)`: this returns the score of the learner's utterance when it is classified as an answer. The `Answer_Scorer()` uses the NCD to match the answer to the list of answers (expected/misconceptions/bad) and returns the score of the highest match. Again a threshold of 60 % is set as the minimum score for a match.
- `Concept_Classifier(array)`: matches the tokens in the learner's input text with the list of concepts associated with the expected answer and return a list of missing concepts.
- `Sentence_Generator()`: Discourse Manager (DiM) calls this function to populate a template response from the `System_Utterance_Library` according to the selected system move.

Dialogue_Manager Class

It controls the dialogue capabilities of the system. It uses the following member functions to achieve this task:

- `Context_Analyzer()`: this function constructs and maintains the *session* variables in order to generate the *current context*. The current context is used for user utterance classification as well as selecting next system move tasks.
- `Rule_Checker()`: accesses the `Rules_Library` to match rule conditions according to current values of parameters such as user utterance type, number of tries, misconception identified, concept coverage etc.
- `Discourse_Manager()`: this function uses the system move selected by the `Context_Analyzer` and the `Rule_Checker` to generate the system response by calling the `Sentence_Generator()` functionality of the NLPE class.

Evaluation Study 1

The first evaluation study of the NDLtutor was conducted to assess the dialogue management capabilities of the system, use of affective and behavioral states to control dialogue flow and using a natural language interface as the communication medium. This evaluation focused on:

1. Quality of dialogues produced by the system.
2. Completion of meaningful dialogues.
3. Use of affective and behavioral states to control the flow of dialogues
4. Use of reflection dialogues as a means to promote metacognition and self-assessment.

Participants

The participants for this evaluation were 20 students from the undergraduate Software Engineering program at Bahria University Islamabad, Pakistan. These students were at

the time enrolled in the data structures course and had just recently been introduced to the topic of stacks. The students had no previous experience of using an ITS system.

Method

Before the start of the session the students were given an overview of the system and the functionality available to them by the first author through a video conference session on Skype. They were introduced to the interface and the possible modes of interaction they could use. They were encouraged to inspect/change their belief base whenever they felt necessary. An initial LM was generated using the test scores of the students in their class exam and the lecturer's personal feedback about each student. The LMs were intentionally altered to show the student's knowledge level to be less than their original knowledge level. This was done to motivate the students to challenge the system's representation. The experiment was conducted in the computer lab of the Software Engineering department and a local instructor (Senior Lecturer) was present at the time of the experiment. The first author was also virtually present via Skype to answer any question.

Each student logged into the system for an individual session which typically lasted between 15 to 20 min. Individual logins were provided so that logs of individual interactions could be recorded in the database. The students could view the system's representation of their LM which was a simple skill-meter. To reinforce interaction symmetry the students were allowed to invoke the chatbot directly by using the DISCUSS <topic > command. Hence a negotiation session could be initiated by the system when a student made a conflicting change to their belief base or it could be initiated by the student by using the DISCUSS command. Out of 20 negotiation sessions recorded in the experiment, 18 (90 %) sessions were initiated by the system whereas two (10 %) sessions were student-initiated. A post-experiment survey was conducted to get the student's feedback about the system. Self-reflection dialogue logs were also used to analyze the learner's interest and reaction to the dialogue itself.

Learner Interactions

During the course of the evaluation, different interactions were seen depending upon different characteristics of the students. The major characteristics that influenced a session include:

Knowledge level – the difference between the knowledge levels of the students had a major impact on the interaction. The interaction time of the more knowledgeable students was considerably shorter than that of the less knowledgeable students. This was an obvious observation since the more knowledgeable students were able to justify the change they had made by discussing the topic with the system and required minimal or no support from the system. This observation was in line with the findings of previous research on learner's prior knowledge in multimedia learning environments (Lawless and Brown 1997). Their answers were more concrete and well-formed which left little room for the NDLtutor to continue the dialogue. Figure 7 shows such an interaction log. The student's answer score is high and the concept coverage is high as well, therefore the system does not need

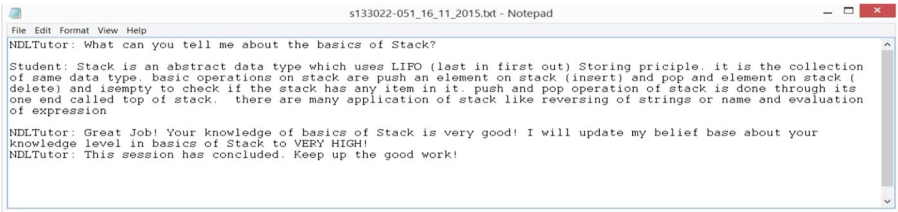


Fig. 7 Knowledgeable learner interaction with NDLtutor

to deploy any funneling strategy. Another important observation from this interaction is that the reflection phase is also influenced by the depth of the discussion. Since the depth of dialogue is so shallow that the system cannot engage the learner in a reflective dialogue regarding their discussion.

Contrary to this, the interaction sessions of the less knowledgeable students were longer and provided more insight to the evaluation. These interactions followed different paths depending upon the student's reply and therefore the NDLtutor needed to make more strategic dialogue decisions. The basic markers for such students were the low answer score and low concept coverage. This provided more room for discussion as the system could ask a series of funneling questions in order to cover the topic. This category of students was the main focus of our study as they allowed us to test our system's dialogue management capabilities. The dialogue fragment in Fig. 8 shows an interaction log of a student with low knowledge of stacks. As seen the student's answer is not complete and allows the system to engage in a funneling discussion about the topic. Such an interaction also provided a gradual transition into the reflection phase.

Affective and Behavioral states – one of the main research issues we are investigating in this study is the impact of using affective and behavioral state of a learner to make dialogue control decisions. Recognizing and responding to the emotional states of the learners have been shown to promote engagement and learning gains (Woolf et al. 2009). The influence of such states was clearly observed during the review of the interaction logs. Figure 9 shows an excerpt of such a dialogue where the NDLtutor is able to identify a specific state and use this information to control the

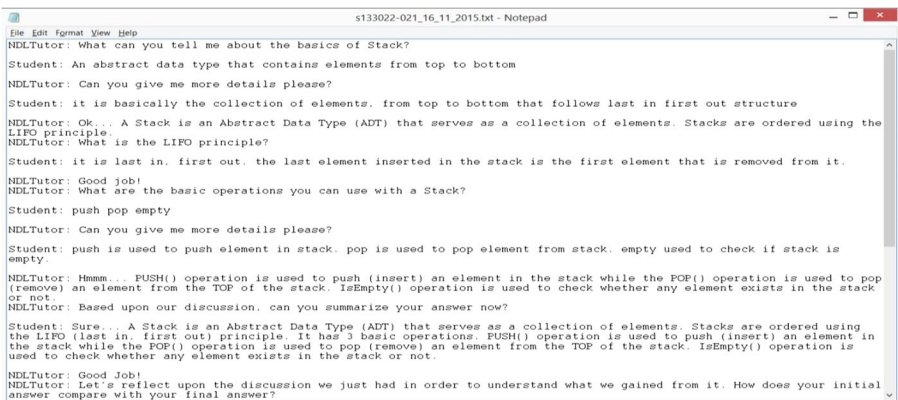


Fig. 8 Less knowledgeable learner Interaction with NDLtutor

```

s133022-019_16_11_2015.txt - Notepad
File Edit Format View Help
NDLTutor: What can you tell me about the basics of Stack?
Student: An abstract data type that contains elements from top to bottom
NDLTutor: Can you give me more details please?
Student: Can you give me more details please?
NDLTutor: I would be happy to help you. But I need you to work with me. Give it a try. I am sure you can do better!
Student: I don't know more about stack
NDLTutor: Ok here's a start... A Stack is an Abstract Data Type (ADT) that serves as a collection of elements.
Stacks are ordered using the LIFO principle.
NDLTutor: What is the LIFO principle?
Student: i know this...its last in first out
NDLTutor: Great work! You have already started answering my questions.
NDLTutor: Stacks are ordered using the LIFO principle. The Last element Inserted in the stack is the First element
taken Out from the stack.
NDLTutor: Let's keep working!
NDLTutor: What are the basic operations you can use with a Stack?
Student: give me a hint?
NDLTutor: Why don't you try it yourself first? I will help you with the correct answer.
Student: i don't know
NDLTutor: HINT: how would you INSERT something in a small jar?
Student: Oh yes! push to insert and pop to remove!
NDLTutor: You are doing great!
NDLTutor: PUSH() operation is used to push (insert) an element in the stack while the POP() operation is used to pop
(remove) an element from the TOP of the stack. IsEmpty() operation is used to check whether any element exists in
the stack or not.
NDLTutor: Based upon our discussion. can you summarize your answer now?

```

Fig. 9 NDLtutor dialogue excerpt showing system's adaptation to the learner's response patterns

dialogue flow. The student is not confident about his knowledge of stacks but shows interest to interact with the NDLtutor and remains on topic. The NDLtutor provides maximum scaffolding to the student as they appear to try harder with every question answered. The student tries to ask for help repeatedly and this is caught by the NDLtutor. To cope with this, the system encourages the student to try to answer by himself before he could receive help/hint. In the future development iteration, the use of hints/help feature would be further formalized so that the NDLtutor can ensure maximum input from the student before providing assistance on domain knowledge.

The student's confidence was found to be more of a personality trait and not directly associated with their knowledge level as we observed less knowledgeable students to show confidence in their interactions as well. However interest and engagement levels were found to be more influenced by the student's knowledge level. Students with very low knowledge of the topic inclined to show less interest in the discussion and repeatedly asked the system for help. This highlighted an important caveat that was not fully taken into consideration during the initial analysis of the system. The act of *gaming* the system was seen in some interactions where the less knowledgeable students were uninterested in the domain discussion and repeatedly asked the system to provide them with help. This was expected as off-task behavior has been linked to students with low motivation and prior knowledge in previous research (Baker et al. 2004). The students used the system's answers during the domain discussion phase and copy-pasted them as their final answer to receive a high score. These new insights were recorded for the next development iteration of the system.

Results and Discussion

As stated earlier the main focus of the evaluation was the *dialogue management capabilities* of the NDLtutor. The results collected from the experiment consisted of

two parts; the interaction logs and the post-experiment survey. Table 5 shows the results of the survey conducted at the conclusion of the experiment phase. The findings were in line with previous researches on tutorial dialogue and learning effectiveness (Core et al. 2003; Rosé et al. 2003; Katz et al. 2003).

While analyzing the results of the survey, the most prominent discovery was the high rate of acceptance from the students. In our understanding, a major factor leading to this outcome was the “Asian culture” influence. We had actually discovered this in one of our earliest survey’s for another study. Asian students tend to be very respectful and polite in their interactions with their tutors. This is a major factor that influences their reactions and it was again prominent in the results of this survey. Having highlighted this, we do recognize that the students were actually very interested and impressed by their interactions with the NDLtutor. They were intrigued by the idea of discussing a topic with a computer tutor in a natural language setting. The authors received multiple emails and Facebook comments from students showing interest in the NDLtutor and volunteering for future experiments. The interaction logs were analyzed in the light of four major criteria set for the experiment.

Quality of dialogues produced by the system – the first criterion was related to the quality of the dialogues generated automatically by the system. It is imperative that the system is able to generate dialogue that engage and motivate students. The analysis of the interaction logs revealed that the system was indeed able to initiate and conduct fruitful dialogues with the student. The user utterance classification scheme that was defined in the earlier section was validated by reviewing the interaction logs and further supported by the survey results where 90 % of the students agreed that the system was indeed able to understand their inputs. In the case of a mismatch the system asked the students to rephrase what they had said which proved to be a good strategy to improve

Table 5 Post-experiment survey results

	<strongly agree.....strongly disagree>					Mean
	(5)	(4)	(3)	(2)	(1)	
Do you think discussing a topic with the chatbot was a good way of justifying your proficiency in that topic?	16	3	1	0	0	4.75
Do you think discussing a topic with the chatbot helped you improve your understanding?	13	4	2	1	0	4.45
Was the chatbot able to correctly understand what you wanted to say?	16	2	0	1	1	4.55
Were the system’s reactions to your inputs valid?	14	4	1	0	1	4.5
Did the chatbot make the negotiation process easy?	14	2	3	1	0	4.45
Did the use of off-topic discussion/small talk make dialogue feel realistic/natural?	4	7	6	2	1	3.61
Did you find the reflection dialogue beneficial?	16	1	3	0	0	4.65
Would you be interested to use a similar system in the future as a study resource?	18	1	1	0	0	4.85

the system's understanding of the inputs. All the sessions were completed successfully which showed the robustness of the system's dialogue management capabilities.

Completion of meaningful dialogues – as discussed above, all dialogue sessions terminated successfully with the mutual agreement between the student and the NDLtutor. The inclusion of small talk in the system corpus proved to be a valuable decision during the system design phase. The post-experiment survey showed that students thought that a minimum amount of small talk made the system feel more realistic and natural. The students also appreciated the misconception funneling functionality of the system and found it to be really helpful in correcting their erroneous beliefs. Moreover it provided them with a chance to discuss the topic in more detail which promoted deeper learning.

Use of affective and behavioral states to control the flow of dialogues – as seen in Fig. 9 above, the use of affective and behavioral states to control the flow of dialogues allowed the system to be more flexible and naturalistic in its responses to the students than the negotiation mechanisms of other existing OLMs. One interesting observation from the interaction logs was that in the case of the system identifying the student's mental state erroneously, the impact on the dialogue was not drastic. This was due to the fact that the system used the information about the student's states in conjunction with the current context of the dialogue. An example of such an occurrence is seen in Fig. 10. The system identifies the student's behavior as "not confident" and raises this point to confirm its classification. The student reacts by reaffirming their belief in what they had said. Since their answer was correct, the system accepts their justification and proceeds to the next dialogue move.

Use of reflection dialogues as a means to promote metacognition and self-assessment – the survey results in Table 5 confirm that the students found the reflection phase to be very helpful in promoting self-reflection. Irrespective of the fact that the more knowledgeable students did not have the reflection phase in their interaction, all the participants unanimously agreed to the usefulness of having a reflective dialogue at the end of the domain discussion. The option of viewing the reflection logs was also welcomed by all the participants. The students accepted that a reflection log would allow them to reflect upon their learning periodically.

Due to the limited empirical data, the question about how the students may use such reflection logs is out of the scope of this evaluation. This will be an interesting prospect to investigate and therefore will be a part of the future evaluations of the system.

```

s133022-036_16_11_2015.txt - Notepad
File Edit Format View Help
NDLTutor: Can you give me more details please?
Student: it is the collection of elements. from top to bottom. we use it in programming.
NDLTutor: Ok... A Stack is an Abstract Data Type (ADT) that serves as a collection of elements. Stacks are ordered using the LIFO principle.
NDLTutor: What is the LIFO principle?
Student: i think its last in first out
NDLTutor: Are you sure? You don't seem to be confident about your answer.
Student: Yes... LIFO is last-in first-out. Last element inserted in the stack is the first element removed from the stack.
NDLTutor: Good Job!
NDLTutor: What are the basic operations you can use with a Stack?
  
```

Fig. 10 NDLtutor dialogue excerpt showing system's confirmation of the student's confidence in his response

Evaluation Study 2

The second evaluation study of the NDLtutor was conducted to assess the pedagogical implications of the NDLtutor. This evaluation focused on:

1. Improvement in Self-assessment accuracy.
2. Effects (if any) of the reflection phase on the Self-reflection skills of learners.

Participants

The participants for this evaluation were 20 students from the undergraduate Software Engineering program at Bahria University Islamabad, Pakistan. 15 students had participated in the first evaluation while the remaining 5 had no previous experience of using an ITS system.

Method

As with the first evaluation, before the start of the session the students were again given an overview of the system and the functionality available to them by the first author via Skype video conferencing session. This was done to accommodate the 5 new students who had volunteered for the study. They were introduced to the interface and the possible modes of interaction they could use. The domain was extended to include the topics of Queues and Linked Lists in addition to the topic of Stacks. The students were asked to concentrate on one topic per session. A single topic was selected per session to ensure maximum concentration and engagement of the students.

For this evaluation the system implementation was updated so that the students had to make an initial self-assessment for each topic after logging into the system. The self-assessment scores were divided into a 5 confidence bands namely; Very Low, Low, Moderate, High, Very High. Each of these bands had a corresponding numerical value assigned to it as follows; Very Low =0, Low =1, Moderate =2, High =3 and Very High =4. Once the students completed the self-assessment they were provided with the option of taking the MCQ test. The MCQs for the topic of Stack were updated from the previous version of the system in order to generate fresh results. The system's learner model for a topic was updated once the student completed the MCQ test for that specific topic. Once the student completed the MCQ test for a selected topic, the system's learner model was updated for that topic. The system then asked the student to confirm their initial self-assessment or update it if they deemed necessary. As the students confirmed/updated their belief base, conflicts occurred between the belief base of the learner and that of the system and at this point the system initiated a dialogue session for the corresponding topic. At the end of the dialogue session the system either accepted the student's change (system's belief base changed) or rejected it (system's belief base remained unchanged). When the student logged off from the system, they were alerted about any discrepancies between the belief bases as a last resort to encourage them to review their belief base in contrast to that of the system.

It is worth mentioning here that in the first evaluation study, we intentionally manipulated (reduced) the system's belief score about the learner's knowledge level

to motivate the learners to challenge the system in order to evaluate the dialog management capacity of NDLtutor for maximum dialogue interactions. However in the second evaluation study, no such manipulations were made to ensure a natural dialog activity of learners in the normal context. As a consequence not all of the students engaged in a dialogue with the system for every topic.

To analyze the effects on the self-reflection of the students, the reflection phase was updated to include a scoring mechanism. The reflection phase consisted of five questions allotted one point, hence five points per reflection session for each topic. These questions were structured specifically so that the answers could be quantified in terms of a numerical value. For example, the learner was asked to score their initial answer on a scale of 0 to 10 (0: minimum, 10: maximum). The student’s answer was then compared with the system’s score of their initial answer. This was done to test whether after completing the domain discussion phase, the student would be able to evaluate their initial answer better. If the student’s scoring of their initial answer matched with that of the system, they were awarded a single point. Details about these measures are described together with the results presented below.

Results

This evaluation focused on the effects of the NDLtutor on the self-assessment and self-reflection skills of the students. To gauge the effects on the self-assessment of the students, we used two discrepancy measures introduced in a previous study on the evaluation of CALMsystem (Kerly et al. 2008b). The selected measures are:

1. Self-assessment accuracy

The self-assessment scores for a student were calculated as the numerical sum of the student’s belief across all three topics. Hence the highest possible self-assessment score for a student could have a value of 12. The self-assessment error was calculated for two cases; a) Before Negotiation and b) After Negotiation. Figure 11 shows the results of the self-assessment evaluation. The mean self-assessment error before negotiation for all the 20 students was 1.6 with a standard deviation of 0.860. The mean self-assessment error for all the 20 students was reduced to 0.65 after negotiation with a standard deviation of 0.653. Hence

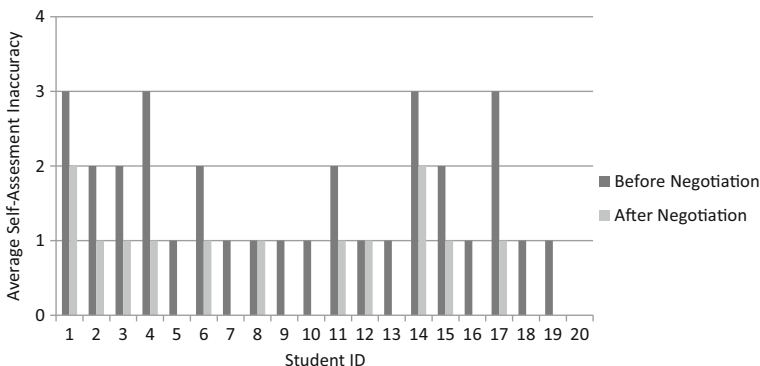


Fig. 11 Self-assessment inaccuracy before & after negotiation with NDLtutor

significant improvements ($t = 3.83, p < 0.0005$) in self-assessment were made by the students after negotiating with the NDLtutor. Figure 11 shows that the students did change their self-assessments after negotiating with the NDLtutor and their final self-assessments at the end of the evaluation study, more closely matched with the system's assessment about their knowledge. Out of the 19 students that engaged in a dialogue with the system, 17 (89.4 %) students made changes to their belief base that resulted in the reduction of the self-assessment error whereas two (10.5 %) students did not make any changes to their belief base after negotiation. The belief bases of 8 (42 %) students matched completely with that of the system at the end of the experiment.

2. No. of Topics with discrepancy

The second discrepancy measure adopted from the previous study on CALMsystem was the reduction in the No. of Topics with discrepancy. This measure was calculated as the difference of the number of topics where the student's belief base value was different from that of the system before negotiation, to number of topics where the student's belief base value was different from that of the system after negotiation. The mean number of topics with discrepancy before negotiation was 1.45 across the three topics for all the students. The mean number of topics with discrepancy after negotiation reduced to 0.65 indicating that there was significant reduction ($t = 3.72, p < 0.0006$) in the number of discrepancies after negotiating the topics with the NDLtutor. Figure 12 shows the number of topics with discrepancy for each individual student. Out of the 19 students that engaged in negotiation with the system, the number of topics with discrepancy reduced for 15 (78.9 %) students whereas the number of topics with discrepancy did not change for 4 (21 %) students at the end of the experiment. The high percentage of students with reduction in the number of topics with discrepancy indicates that the students did in fact reassess (review) their belief bases after negotiating with the NDLtutor.

3. Effects in Self-Reflection

Promoting metacognitive skills of the learners has always been one of the major objectives of OLMs. Opening up the learner model to the learner was intended to

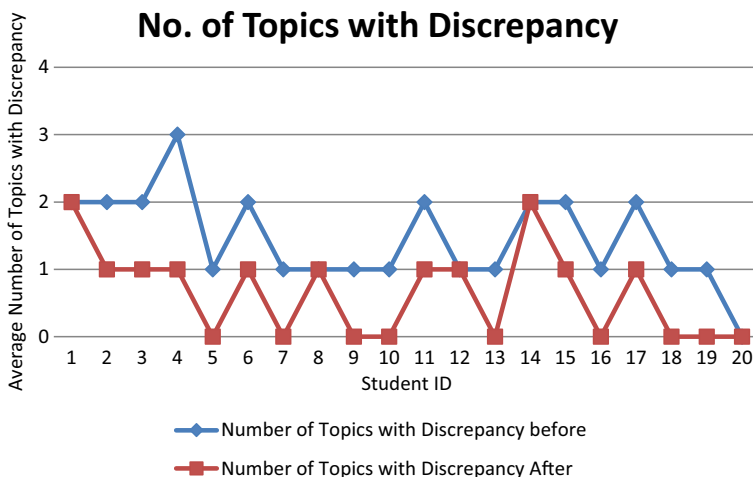


Fig. 12 Number of topics with discrepancy before & after negotiation with NDLtutor

maximize learner participation as well as promote self-reflection in learners (Bull and Vatraru 2012). Learner's self-assessment of their belief in their knowledge level is considered as a reflective activity. An improvement in self-assessment has been used as an indicator for promotion of self-reflection (Kerly et al. 2008b; Dimitrova 2003). How the learner is reflecting is mostly implicit as self-reflection is domain/task-independent. This implicitness of self-reflection skills of a learner and their ability to use such skills makes formally analyzing and assessing such skills a difficult task. It has been argued in the research on assessing and explicitly promoting self-reflection in learners that the system should focus on providing the learners with the tools to engage in some form of reflective activity such as; self-assessment of their belief-base, skill diaries (Long and Alevan 2013) and self-explanations (Gross et al. 2015).

Based on the concept of skill dairies we introduced a reflection phase at the end of each dialogue session in the NDLtutor. The idea is to encourage the learner to reflect upon their discussion with the system. Our aim is to provide the learner with support in a domain-independent form of interaction that helps them in analyzing/realizing how they answered the system's questions during the domain discussion, what were the problems they encountered, what concepts they missed or what misconceptions were highlighted during discussion. To enable the system to analyze the learner's input during this phase, we introduced an informal formative assessment that uses five questions which can be quantified by the system to generate a reflection score at runtime. The reflective score for each student is calculated for each individual session by the system. The five questions carrying one (1) point each are as follows:

- Q1. *On a scale of 0 to 10 (0: Minimum, 10: Maximum) how would you rate your first answer?* – This question is used to analyze the learner's ability to evaluate their initial answer. The value of scale provided by the learner is converted into a percentage value and compared with the system's evaluation of the learner's initial answer i.e. (answer score + concept coverage). If the learner's evaluation score matches the evaluation score of the system (permitted variance: $\pm 15\%$), the learner is awarded one point, otherwise zero point.
- Q2. *How is your last summarized answer different from your first answer?* – The answer to this question is tested for learner's verbosity and their ability to identify the incompleteness of their initial answer. The learner's answer is analyzed for statements relating to incompleteness of their initial answer as well as missing details. If the learner's answer includes these markers, the system awards one point.
- Q3. *What were the concepts that you missed?* – This question is used to check whether the learner is able to recall the concepts they missed in their initial answer. If the learner had missed some concepts during the domain discussion, then the system asks them to list these concepts. The learner is awarded one point if he is able to list all the concepts he missed during the domain discussion. If the learner did not miss any concepts, the system accepts "No" as an answer and awards one point.
- Q4. *Did you encounter any misconceptions?* – Similar to Q3, the learner is asked to state the misconceptions (if any) that were encountered during the domain

discussion. If no misconceptions were encountered the system accepts “No” as an answer and awards one point.

- Q5. *Did you improve your understanding/knowledge on the topic?* – Whether the learner’s belief about their understanding of the topic changed after their interaction with the system. The learner’s answer is analyzed with respect to the score of their final summarized answer in the domain discussion phase. If the learner’s final answer score is higher than their initial answer score then the expected answer to this question is “Yes”, which earns the learner one point. On the other hand, if the learner provides “No” as an answer to this question and their final answer score is low, the system allocates one point and asks them to elaborate on the reasons that might have hindered their learning.

It is necessary to state that this reflective score is not intended to be used as a formal assessment measure; instead we argue that such a score can be used to study the different correlations between the learner’s self-assessment beliefs, their performance in the domain dialogue, and their responses in the reflection phase over a period of time. Figure 13 shows the reflection scores of all the 20 students that participated in the evaluation. The scores are shown for each reflection session a student engaged in. Here is it important to point out that the reflection phase was only initiated for a student who was unable to provide a high scoring answer to the initial Domain Discussion Question in the domain discussion phase. This means that not all of the students engaged in a reflection session across all the three topics. This explains the empty columns of the students for some sessions in Fig. 13. That is to say, an empty column does not show a 0 reflection score, it only indicates that the student did not engage in a reflection phase for the specific session. The mean reflection score across all the three sessions was 3.82 which show that most students were able to get high scores during the reflection phase. 5 (25 %) students engaged in the reflective dialogue of Session-I (Stacks). This number increased in Session-II (Queues) to 12 (60 %). A similar number of students 12 (60 %) engaged in the reflective dialogue for Session-III (Linked Lists). The mean reflection score of Session-I was 2.8 whereas the mean reflection scores for Session-II and Session-III were 3.83 and 4.25, respectively. The stats reveal that more students engaged in the reflective dialogue than those who

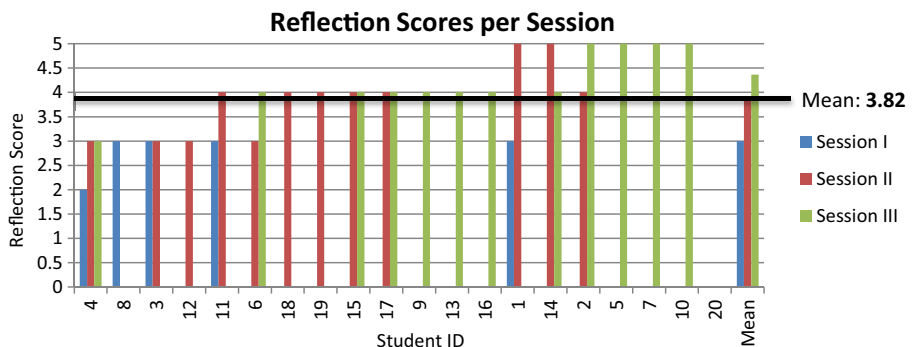


Fig. 13 Self-Reflection scores (in ascending order) of individual students across all sessions

engaged in the initial session, and as the complexity of the domain topic increased. Some further observations are as follows:

1. Not all students engaged in reflective dialogues for all the three sessions. As defined earlier, a reflective dialogue is only conducted after the student has completed the domain discussion dialogue with the system. If there is no discrepancy between the learner's belief and that of the system, or if the learner accepts the system's belief value and updates their own belief base to match the system's belief without challenging the system, in such a scenario no dialogue session is conducted and hence the student does not engage in a reflection dialogue. For instance, in Fig. 13, Student#1, has no reflection score for session-III as this student did not challenge the system and accepted the system's inference about his knowledge level on the topic of Linked Lists. Only Student#4 engaged in all the three reflection sessions.
2. When there is a discrepancy between the system's belief and that of the student, then the system initiates a dialogue, if the student is able to answer the domain question in their first attempt to an acceptable standard, then the system does not have any room for a reflective discussion. Hence for students with high level of knowledge, the possibility of engaging in a reflection session is minimal. This can be seen in Fig. 13 as student # 20 has no reflection scores. This student was able to prove their knowledge during the MCQ tests for all topics, so that the system had no rationale to challenge his beliefs.
3. We found two encouraging suggestions. The first one is that the reflection score for each student remained neutral or positive and did not decline over multiple sessions except Student #14. The second is that the average of the reflection scores for each session increases as shown in the "Mean" column of Fig. 13. Although we cannot claim it with statistical significance in this experiment, whether interacting with the system multiple times had an effect on the learner's answers or did it play any role in training the learners to answer better is an interesting topic and it would be worth exploring in future studies.
4. Another interesting observation is that the reflection scores of the learners suggest a direct correlation to their confidence in their knowledge level of the topic. To find the correlation between the learner's confidence in their knowledge and their reflection score, we calculated the Spearman's rank correlation coefficient for *tied* data. The average correlation coefficient between the learner's confidence in their knowledge and their reflection score across all three topics was found to be 0.674, which shows a positive correlation that is statistically significant at the 0.01 level (for $n = 20$) between the learner's confidence in a topic and the reflection scores. In fact, students who chose low confidence values i.e. "Very Low" and "Low" in their belief bases tended to accept the system's inference without challenging the system. These students were also observed to generally have a below average reflection score in the reflection phase (Student#1, Student#3, Student#4, Student#8 and Student#11 in Fig.13). However, the students who were more confident about their knowledge level and chose "High" or "Very High" values in their belief bases challenged the system more and also scored higher in the reflection phase. For instance, all the students in Session-III had chosen "High" or "Very High" as their belief base value for the topic of Linked List. The only

- exception was Student#4 who had chosen “Moderate” and this student scored below average in the reflection dialogue. Whether or not there is a direct and strong correlation between the student’s confidence in their knowledge level and their reflection scores is an interesting observation. However, this requires a greater number of interaction sessions to be further investigated.
5. While the above observation presents the correlation between learner’s confidence of his/her knowledge level and the reflection score, here we discuss about the correlation that was observed between the learner’s *actual* knowledge level which is in the system’s belief base and his/her reflection score. The correlation between the learner’s knowledge as assessed by the system and their reflection scores was also calculated using the Spearman’s rank correlation coefficient for *tied* data. The average correlation coefficient between the learner’s knowledge as assessed by the system and their reflection scores across all three topics was found to be 0.68, which shows a positive correlation that is statistically significant at the 0.01 level (for $n = 20$) between the learner’s knowledge level in a topic and their reflection scores. From further analysis of the results we were able to define three broad categories of students; 1) Below Average, 2) Average and 3) Above Average according to their knowledge level as assessed by the system during the MCQ tests. Students in each category shared similar characteristics. The Below Average students showed the tendency to only challenge the system’s beliefs on the basic/easier topics i.e. Stacks or in some cases Queues. These students were mostly unable to defend their claim during the domain discussion phase and also scored below average in the reflection phase (Blue bars in Fig. 13). The Average category of the students can be considered as the ideal candidates for the system as these students demonstrated an average level of knowledge and high confidence in their assessments. These students challenged the system’s beliefs across all the topics and were mostly successful in defending their beliefs. Their reflection scores (Red bars in Fig. 13) were also closer to the overall mean score. The Above Average students were the ones who had a high or very high level of knowledge and were also confident about their beliefs. Such students mostly engaged in the dialogue with the system for the advanced topic i.e. Linked List. The Green bars in Fig. 13 show that almost all the students who engaged in the dialogue related to the advanced topic had an above average score. This observation is in line with previous research that students who have better metacognitive skills perform much better than students who have weak metacognitive skills (Swanson 1990; Schraw and Dennison 1994).

Related Work

Figure 14 shows the research themes that motivate and influenced the research on the Negotiation-Driven Learning paradigm. This section provides an overview of these research areas and how they contribute to the development of NDL. Open Learner Models emphasize the active involvement of the learner in the process of improving the accuracy of the Learner Model. OLMs utilize different strategies of negotiation in order

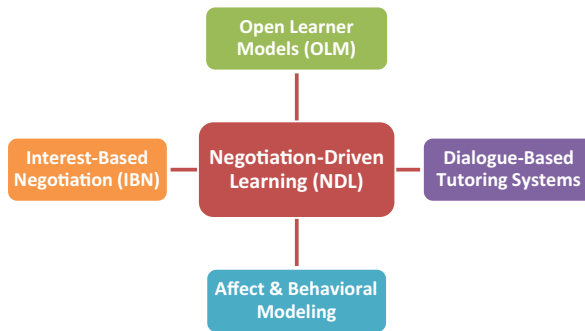


Fig. 14 Research themes in the NDL paradigm

to allow the learner to discuss their LM with the system. These strategies are mainly differentiated on the amount of control the learner and the system have over the course of the dialogue. Fully Negotiated OLMs allow a larger degree of control to the learner as compared to other negotiation strategies by deploying interaction symmetry that provides the same dialogue moves to both the system and the learner (Bull and Vatrupu 2012). Allowing the learner to change their belief base gives them a sense of control over the process while having the ability to defend their beliefs against the system and ask for justification from the system inculcates a sense of trust in them.

Allowing the learner to interact with the system about his LM, opens up new doorways of interaction possibilities and diagnosis. The interactive nature of the dialogues provides an opportunity to promote reflective thinking in the learners. A very important aspect of OLMs has always been the active promotion of metacognitive skills in the learner (Bull and Kay 2013). It has been documented that students who have better metacognitive skills perform much better than students who have weak metacognitive skills (Swanson 1990; Schraw and Dennison 1994). Since metacognitive skills do not have any observable manifestation, such skills are hard to acquire and gauge. However continuous stimuli can lead the learner into learning to use such skills more actively so that these skills are automatically used by the learner while they are learning. Improvement in the metacognitive skills of the learners has mostly been implied implicitly. Externalization has been considered as one of the major sources of self-reflection in learners. When they are able to view their LM and reflect upon their knowledge level. However, this self-reflection remains implicit and OLMs do not provide a clear platform to the learner to keep a track of their metacognitive abilities.

ITS systems are modelled to replicate expert or semi-expert tutors, since expert tutors have shown to have the maximum learning gain in learners. An important aspect of the expert tutors teaching tactics is the ability to react to the student's emotional and motivational states (Chi et al. 2001). A learner's affective and behavioral states play a vital role in the outcome of their interaction with the system as a confident, interested, and motivated learner would interact very differently from a learner who is not confident, uninterested or demotivated. For an automated system to be able to replicate an expert tutor's empathy, it needs to be able to classify the learner's interactions as a possible outcome of a mental state. Intensive research on this effect would contribute to advancement of the field of OLMs.

Learner Models

Intelligent Tutoring Systems use Learner Models to provide adaptive and personalized content to the learners (Self 1998). The system diagnoses the learner's knowledge during its interactions with the learner and uses this information to infer the corresponding learner model (VanLehn 1988). The learner model represents the current state of the learner's knowledge. The basic tasks for any learner model include (Wenger 1987):

1. Storing information about the learner's knowledge and expertise about a particular domain. This information allows the system to compare the learner's knowledge with that of an expert module to generate evaluations and highlight area of weakness.
2. Representation of the learner's knowledge level that allows for an insight into incorrect knowledge and misconceptions held by the learner.
3. Accounting for data by analyzing the information available to the system to generate the diagnosis for the learner. Such diagnostic process can vary depending upon the kind and amount of information available to the system.

Traditionally the learner model was encapsulated from the learner and only visible and available to the system for adaptive tutoring. It has been argued and that involving the learner in the process of constructing and maintaining their learner model not only promotes learner engagement but also has positive effects on their metacognitive skills (Bull and Pain 1995; Kerly et al. 2008b; Dimitrova 2003).

Classes of Open Learner Models

We can identify different classes of OLMs according to the level of control they provide to the learner over the LM. The learner's level of control can be defined as the learner's capability to change the contents of the LM. According to this specification, OLMs can be classified as:

1. **Inspectable:** An inspectable OLM can be considered as a read-only or view only OLM. The LM is completely controlled by the system and is only available to the learner for viewing. The learner has no right to change the contents of the LM directly. The learner can answer questions related to the domain in order to have their model updated. The externalization of the LM has shown to increase learner involvement and promote self-reflection and planning skills (Bull and Kay 2010). All OLM implementations are considered inspectable since they allow the learner to view their learner models in one form or another.
2. **Co-operative:** These models allow the learner and system to jointly construct the learner model. The system asks the learners to provide complementary information required for the modeling process (Beck et al. 1997).
3. **Challenge:** These OLMs allow the learner to challenge the model generated by the system. EI-OSM (Zapata-Rivera et al. 2007) is one such system based on Toulmin's model of argumentation (Toulmin 1958). EI-OSM uses claims, data, warrants, backing and rebuttal to allow learner to add new arguments with supporting evidence. A teacher has the authority to determine which evidence

has the highest strength and the evidence supported by the teacher is considered stronger than the unapproved evidence provided by the learner. Another OLM that allows the learner to challenge the system is xOLM (Van Labeke et al. 2007). The learner is allowed to view the model and select the topic for discussion. The system provides justification for the topic and the learner are provided with three options; 1) agree 2) disagree and 3) move on, to continue the interaction. If the learner agrees with the system, the system's beliefs are reinforced. In case of a disagreement the learner has to provide further information which is used to diagnose the model. Move on allows the learners to end the discussion with the system.

4. Add-Evidence: These OLMs allow the learner to provide additional evidence to be considered in the modeling process. ELM-ART (Weber and Brusilovsky 2001) is an OLM that allows the learners to inspect and edit the contents of their learner model. ELM-ART is implemented as an adaptive interactive textbook where the learner informs the system about their knowledge by providing evidence to support their claim. Evidence can be in the form of answering questions, taking tests or performing tasks. Another OLM that allows the learner to provide evidence is TAGUS (Paiva et al. 1995). The learner can inform (*tell*) the system about the new evidence which is then analyzed by the system to take appropriate action.
5. Editable: Learners have full responsibility and control in editable OLMs. They are allowed to edit their learner model when they deem necessary without the intervention of the system. The system may offer some information regarding its belief base which can be neglected or overridden by the learner. The changes made by the learner are directly reflected in the system's belief base which alters their learner model. Some examples of OLMs in this class are; C-POLMILE (Bull and McEvoy 2003), SASY (Czarkowski et al. 2005) and Flexi-OLM (Mabbott and Bull 2006).
6. Persuaded: Persuaded OLMs also allow the learner to change their learner models but they are required to demonstrate their competency before the system can agree with the changes they made. The system uses questioning techniques to analyze the learner's knowledge level and validate their claim. If the learner is not able to justify the change they made, their changes are rejected by the system and the learner model remains unchanged. Flexi-OLM (Mabbott and Bull 2006) is an OLM that falls in this category.
7. Negotiated: Negotiated OLMs allow for a more collaborative approach towards constructing and maintaining the OLM. Negotiated OLMs use a separate set of beliefs for the learner and the system. The negotiation process is used to resolve the conflicts (discrepancies) between these sets of beliefs. There is an interaction symmetry which provides both the learner and the system with equal rights of interaction. The basic negotiation protocol allows for; ask for justification, provide justification, challenge justification, reject justification, provide proposal, accept proposal or reject proposal. Mr. Collins (Bull et al. 1995) is the first fully negotiated LM which focuses on the discussion of the LM between the learner and the system. Mr. Collins uses a menu-based discussion which allows learners to challenge and respond to the system. While Mr. Collins has been shown to promote learner reflection, the negotiation method used can be considered as restrictive. STyLE-OLM (Dimitrova 2003) is another fully-negotiated system that allows learners to discuss their LM with the system. STyLE-OLM is proposed based on

the idea that interaction is a stimulus for reflection. The dialog is constructed as a conceptual graph that allows the learner to see the explicit connections between the different arguments. However, some learners might find using the graphical interface difficult or distracting. CALMsystem (Kerly et al. 2008b) addresses the problem of using menu selections and conceptual graphs for young learners. In order to provide an easier way to communicate with the system, CALMsystem proposes the use of natural language dialogue. CALMsystem follows the negotiation options provided by Mr. Collins and uses a chat-bot to provide a natural language dialogue. CALMsystem utilized the Lingubot™ (Creative Virtual 2007) technology to build the chatbot. Domain-independent utterances do not affect the course of the dialogue which can be restrictive in a natural language dialogue system. CALMsystem laid the foundations of using natural language conversational agents in the context of OLMs. The background and guidelines provided by CALMsystem formed the basis of the research on NDL and guided the implementation of the NDLtutor. As mentioned earlier, CALMsystem utilizes the negotiation options provided by Mr. Collins, whereas NDLtutor divides the dialogue into three phases so that the interactions in each phase can be handled separately. Where CALMsystem allows the student to take a test to prove their claim, NDLtutor discusses the domain topic and uses funneling questions to elicit missing concepts and misconceptions. This means that if a misconception is encountered during the discussion, the system highlights it and asks funneling questions to help learner realize and remove their misconception in the same session. The learner's model is not updated after each question as in the case of CALMsystem. The NDLtutor only updates the learner model (if necessary) once the dialogue session has terminated i.e. all the phases have been completed. Another distinguishing feature of NDLtutor is the use of the reflection phase to explicitly involve the learner in self-reflection dialogue once the domain (topic) discussion has been completed.

Theories of Automated Negotiation

Negotiation is a vital form of human interaction which ranges from basic information exchange to more complex cooperation or coordination activities. As computer systems evolve to become autonomous agents, it was inevitable for such systems to be able to conduct a negotiation of their own in an automated way. Automated negotiation has found much interest and success in the field of e-commerce where automated autonomous agents negotiate over resources (tangible assets).

Interest-Based Negotiation (IBN) (Fisher 1983) has gained attention from the research community since it provides a good alternative to Position-Based negotiation where all agents are considered adversaries. It is also known as win-win negotiation as all parties try to create a mutual gain. IBN allows the parties to reveal their underlying interest by specifying new information during the course of the dialogue. This information can be used to decide an alternate strategy in real-time which makes IBN more responsive. Since learning is a process of exchanging ideas and understanding problems, IBN seems much more suited for educational systems, as proposed in (Miao 2008). There are no current implementations of OLMs that have tried to utilize IBN as the main negotiation approach.

Affect & Behavioral Modeling

Research has shown that expert human tutors have a higher impact on learning than novice tutors and ITSs (Lehman et al. 2008). This is not only due to the pedagogical strategies employed by such expert tutors but is also deeply rooted in the emotional (affective) and motivational (behavioral) strategies such tutors employ to engage the learners in learning (du Boulay et al. 2010). Affect and behavior are closely entwined in a bi-directional relationship. Moreover a learner may not only experience a positive affective or behavioral state, but also a negative state. Such a negative state might even be necessary for a learner to be engaged in the process of learning. Understanding the state the learner is in can allow a system to be more empathetic towards them which leads to higher levels of engagement. It has been argued that while an exact estimation of a specific state might not be possible or even required, an approximation of these states can be as helpful in continuing the learning process. The terminology of “caring systems” encompasses such systems which are meta-affectively and meta-cognitively aware. NDLtutor aspires to inherit such attributes to provide adequate support the learners to promote their cognitive and meta-cognitive skills.

Dialogue-Based Tutoring Systems

ITS systems have come a long way from having simple human-computer interfaces to adopting conversational interfaces. Apart from the conventional text display and graphics such systems employ an automated conversational agent that is able to speak to the student using synthesized speech accompanied by facial expressions and gestures. This makes the learner’s experience more interactive and has also been shown to increase engagement.

Dialogue-based tutoring systems have deployed different forms of strategies to maximize learning. Knowledge construction dialogues (KCD) were used to encourage students to infer or construct the target knowledge in the ATLAS system (Freedman 1999). KCDs connect principles and relate them to common sense knowledge to help students to discuss their knowledge. ATLAS was originally developed for CIRC SIM tutor and also provides a natural language interface to the learners. Immediate feedback strategy was employed in ANDES (Gertner and VanLehn 2000; VanLehn 1996) to help college and high-school physics students to do their homework problems. ANDES highlighted the use of real-time hints and feedback to help student solve given tasks. A similar approach has been adopted in the NDLtutor in order to ensure that the learners are provided real-time feedback and hints to help them answer questions. NDLtutor uses the hints function to also analyze the learner’s help-seeking patterns. This allows the system to ensure that the learners are only provided help when they have made an effort to answer the system’s questions. This is done to prevent learners from gaming the system.

One of the most successful systems in this category has been AutoTutor (Graesser et al. 1999; Person et al. 2003). It is an ITS that provides a natural language dialogue to interact with the learner. AutoTutor provides the learner with an interactive agent that speaks out the question in addition to displaying the text on the screen. AutoTutor engages the learner in a deep reasoning dialogue which requires the learner to provide comprehensive explanations. Autotutor’s strength lies in its ability to handle learner

responses during the course of the dialogue. Autotutor uses advanced statistical NLP techniques such as Latest Semantic Analysis (Graesser et al. 2000) to analyze learner response and classify responses into corresponding speech acts. NDLtutor follows a similar approach; however it only employs basic NLP techniques such as Normalized Compression Distance (Cilibrasi and Vitanyi 2005) to classify and analyze learner inputs. The main reason for the difference in the complexities of the NLP techniques used in these systems is that the aim of AutoTutor is to act as a teacher and teach/construct knowledge. Contrary to this, NDLtutor does not adopt the role of a teacher, but only reinforces and discusses what the learner knows to improve their understanding.

Concluding Remarks

Open Learner Models maximize learner involvement by engaging them in a process of collaboratively constructing and maintaining their learner model (Bull et al. 1995; Dimitrova 2003; Kerly et al. 2008b). This research on OLMs has shown to produce significant learning gains. Negotiated OLMs utilize different interaction strategies to enhance self-assessment and promote self-reflection in learners. Conversational agents have been used in this regard to provide naturalistic mode of interaction between the learner and the system. This not only eased the communication process but also improved the self-assessment accuracy of the learners. Following the success of using chatbots in OLMs, this study investigates the possibilities of enhancing the capabilities of such chatbots and their implications on the learner's learning. This study introduces the paradigm of Negotiation-Driven Learning which uses a chatbot employing Interest-Based Negotiation strategy to discuss the learner model with the students. We discuss the use of approximations of a learner's affective & behavioral states in order to control the flow of dialogues. Such a scheme enables a more reactive and responsive dialogue between the learner and the system and yields significant self-assessment improvement in learners. We also highlight the explicit reflection phase of NDL for the promotion of metacognitive skills using a reflective dialogue at the end of every session which can also be used as a self-reflection log by the learners.

This paper provides the details of the architecture of our system, the design and implementation, and presents the discussion on the results of two evaluation studies. Our system consists of 5 main components that interact with each other to provide an open-ended, natural language dialogue interface to the learners. We have discussed in details the Wizard-of-Oz experiment that was conducted to collect the data to support our system design. The data that we collected during the experiments was analyzed to select three affective and three behavioral states used to control the dialogue in NDL. The results from the last phase of the WoZ experiment showed that the data we collected and the resulting rules and response libraries allowed the wizard to conduct negotiations with the learners in the domain of Data Structures almost automatically. We then discussed the implementation details of our system the NDLtutor and presented two evaluation studies of our system which has been developed using the architecture obtained by the WoZ experiment. We evaluated the interactions between the students and the NDLtutor to highlight the potential benefits of using our approach. We have argued that our approach provides new insights into combining best practices that have only been used separately in existing OLMs to develop an intelligent tutoring

system that is capable of engaging learners in dialogues that promote metacognitive skills in them.

Providing a Natural Language interface to learners can ease the communication process but adds to the overall complexity. NLP is a research field in its own right and a complete understanding of learner inputs is out of the scope of this study. In order to minimize such complexity we use the Normalized Compression Distance (NCD) and Cosine Similarity Index algorithms to find the matching utterances in the classification process. Natural language discourse requires the system to be capable of handling some amount of off-task behavior or small talk. Being able to cope with this allows the dialogue to flow more naturally. During the WoZ experiment, we collected a significant amount of data that related to small talk or off-topic discussions. This allowed us to classify user utterances as small talk and generate corresponding system utterances. The system is able to carry on a certain amount of off-topic discussion before encouraging the learner to get back on the task at hand. In case of non-responsive students, the system uses the same small talk in an attempt to engage the learner in learning. The current implementation of NDLtutor utilizes basic NLP techniques to manage the dialogue between the learner and the system. It is needless to say that the system would definitely benefit from the use of more advanced NLP techniques. To this effect, further investigations into the use of NLP techniques that do not outweigh their usability with their complexity for the said task are underway to improve the classification and scoring results.

We presented an evaluation framework based on the previous work on using chatbots in OLMs (Kerly et al. 2008b) and evaluated our system accordingly. Our findings were consistent with the previous research and showed a significant improvement in learner's self-assessment abilities after negotiating with the NDLtutor. We also evaluated the explicit reflection phase we introduced in our system in order to support and promote self-reflection in learners. The informal formative assessment of the reflection phase provided evidence that a reflection dialogue can engage learners to analyze and assess their understanding of their knowledge. The results of the evaluation also highlighted characteristics in the interaction patterns that were common in different categories of learners. While our findings demonstrate that the NDLtutor does provide adequate support to the learners to promote reflective thinking, the scope of our results were confined by the scale of the study. Further investigation is needed to make an in-depth analysis into the forms of reflection supported by our system.

NDL follows the notion that learning is maximized by participation in the learning process and negotiation provides an excellent opportunity to challenge the learners which promotes metacognitive skills by motivating them to think more objectively about their learning. NDL finds its roots in the theory of repetition in learning. We believe that continuously engaging learners in dialogue that encourage them to utilize their metacognitive abilities allows them to use such abilities more efficiently over time and our proposed approach has the potential to achieve these desired results.

Acknowledgments The authors are greatly indebted to Vania Dimitrova, Cristina Conati and Amali Weerasinghe for their invaluable suggestions and kind guidance. The authors are grateful to the faculty of Bahria University Islamabad, Pakistan for their kind cooperation in conducting the Wizard-of-Oz experiment. We also thank all the participants who took part in the experiment. This research is partially supported by JSPS KAKENHI Grant Number 26240033 and by the Ministry of Education, Culture, Sports, Science & Technology (MEXT), Japan.

References

- Afzal, S., & Robinson, P. (2011). Natural affect data: Collection and annotation. In *New perspectives on affect and learning technologies* (pp. 55–70). Springer New York.
- Baker, R. S., Corbett, A. T., Koedinger, K. R., & Wagner, A. Z. (2004). Off-task behavior in the cognitive tutor classroom: when students game the system. In *Proceedings of the SIGCHI conference on Human factors in computing systems* (pp. 383–390). ACM.
- Beck, J., Stem, M., & Woolf, B. P. (1997). Cooperative Student Models. In *Artificial Intelligence in Education, 1997: Knowledge and Media in Learning Systems: Proceedings of AI-ED 97, World Conference on Artificial Intelligence in Education, Kobe, Japan* (Vol. 39, p. 127). IOS Press.
- Bull, S. (2016). Negotiated learner modelling to maintain today's learner models. *Research and Practice in Technology Enhanced Learning*, 11(1), 1–29.
- Bull, S., & Kay, J. (2007). Student models that invite the learner in: The SMILI© Open learner modelling framework. *International Journal of Artificial Intelligence in Education*, 17(2), 89–120.
- Bull, S., & Kay, J. (2010). Open learner models. In *Advances in intelligent tutoring systems* (pp. 301–322). Springer Berlin Heidelberg.
- Bull, S., & Kay, J. (2013). Open learner models as drivers for metacognitive processes. In *International handbook of metacognition and learning technologies* (pp. 349–365). Springer, New York.
- Bull, S., & McEvoy, A. T. (2003). An Intelligent Environment with an Open Learner Model for the Desktop PC and Pocket PC. *Artificial Intelligence in Education: Shaping the Future of Learning through Intelligent Technologies*, 97, 389.
- Bull, S., & Pain, H. (1995). “Did I Say what I Think I Said, and Do You Agree with Me?”: *Inspecting and Questioning the Student Model*. University of Edinburgh, Department of Artificial Intelligence.
- Bull, S., & Vatrupu, R. (2012). Negotiated learner models for today. In *Proceedings of 20th International Conference of Computers in Education (ICCE2012)*, pp. 96–100.
- Bull, S., Pain, H., & Brna, P. (1995). Mr. Collins: A collaboratively constructed, inspectable student model for intelligent computer assisted language learning. *Instructional Science*, 23(1–3), 65–87.
- Burleson, W., & Picard, R. (2007). Evidence for gender specific approaches to the development of emotionally intelligent learning companions. *IEEE Intelligent Systems*, 22(4), 62–69.
- Chi, M. T., Siler, S. A., Jeong, H., Yamauchi, T., & Hausmann, R. G. (2001). Learning from human tutoring. *Cognitive Science*, 25(4), 471–533.
- Cilibrasi, R., & Vitanyi, P. (2005). Clustering by compression. *IEEE Transactions on Information Theory*, 51(4), 1523–1545.
- Conati, C., & Maclaren, H. (2009). Empirically building and evaluating a probabilistic model of user affect. *User Modeling and User-Adapted Interaction*, 19(3), 267–303.
- Core, M. G., Moore, J. D., & Zinn, C. (2003). The role of initiative in tutorial dialogue. In *Proceedings of the tenth conference on European chapter of the Association for Computational Linguistics-Volume 1* (pp. 67–74). Association for Computational Linguistics.
- Creative Virtual. (2007). Creative Virtual UK web site. www.creativevirtual.com.
- Czarkowski, M., Kay, J., & Potts, S. (2005). Web framework for scrutable adaptation. In *Workshop on Learner Modelling for Reflection. AI in Education* (pp. 11–18).
- D’Mello, S., & Graesser, A. (2012). AutoTutor and affective AutoTutor: Learning by talking with cognitively and emotionally intelligent computers that talk back. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 2(4), 23.
- D’Mello, S., Picard, R., & Graesser, A. (2007). Towards an affect-sensitive Autotutor. *IEEE Intelligent Systems*, 22(4), 53–61.
- Dahlbäck, N., Jönsson, A., & Ahrenberg, L. (1993). Wizard of Oz studies: why and how. In *Proceedings of the 1st international conference on Intelligent user interfaces* (pp. 193–200). ACM.
- De Vicente, A., & Pain, H. (2002). Informing the detection of the students’ motivational state: an empirical study. In *International Conference on Intelligent Tutoring Systems* (pp. 933–943). Springer Berlin Heidelberg.
- Dimitrova, V. (2003). STyLE-OLM: interactive open learner modelling. *International Journal of Artificial Intelligence in Education*, 13(1), 35–78.
- du Boulay, B., Avramides, K., Luckin, R., Martínez-Mirón, E., Méndez, G. R., & Carr, A. (2010). Towards systems that care: a conceptual framework based on motivation, metacognition and affect. *International Journal of Artificial Intelligence in Education*, 20(3), 197–229.
- Ekman, P. (1973). Universal facial expressions in emotion. *Studia Psychologica*, 15(2), 140–147.
- Fisher, R. U. W. (1983). *Getting to Yes: Negotiating Agreement without giving in*. New York: Penguin books.

- Fredrickson, B.L. (1998). What good are positive emotions? *Review of General Psychology*, 2(3), 300–319.
- Freedman, R. (1999). Atlas: A plan manager for mixed-initiative, multimodal dialogue. In *AAAI-99 workshop on mixed-initiative intelligence* (pp. 1–8).
- Gal, D., & Rucker, D. D. (2010). *When in doubt, shout! Paradoxical influences of doubt on proselytizing*. Washington: *Psychological Science*, 21(7), 1701–1707.
- Garner, R., & Alexander, P. A. (1989). Metacognition: Answered and unanswered questions. *Educational Psychologist*, 24(2), 143–158.
- Gertner, A. S., & VanLehn, K. (2000, June). Andes: A coached problem solving environment for physics. In *International conference on intelligent tutoring systems* (pp. 133–142). Springer Berlin Heidelberg.
- Graesser, A. C., Wiemer-Hastings, K., Wiemer-Hastings, P., Kreuz, R., & the Tutoring Research Group (1999). AutoTutor: A Simulation of a Human Tutor. *Journal of Cognitive Systems Research*, 1(1), 35–51.
- Graesser, A. C., Wiemer-Hastings, P., Wiemer-Hastings, K., Harter, D., Tutoring Research Group, T. R. G., & Person, N. (2000). Using latent semantic analysis to evaluate the contributions of students in AutoTutor. *Interactive Learning Environments*, 8(2), 129–147.
- Graesser, A. C., McDaniel, B., Chipman, P., Witherspoon, A., D’Mello, S., & Gholson, B. (2006). Detection of emotions during learning with AutoTutor. In *Proceedings of the 28th Annual Meetings of the Cognitive Science Society* (pp. 285–290).
- Gross, S., Mokbel, B., Hammer, B., & Pinkwart, N. (2015). Learning feedback in intelligent tutoring systems. *KI-Künstliche Intelligenz*, 29(4), 413–418.
- Katz, S., Allbritton, D., & Connelly, J. (2003). Going beyond the problem given: How human tutors use post-solution discussions to support transfer. *International Journal of Artificial Intelligence in Education*, 13(1), 79–116.
- Keller, J. M. (1983). Motivational Design of Instruction. *Instructional Design Theories and Models: An Overview of their Current Status*. Ed. CM Reigelruth.
- Kerly, A., & Bull, S. (2006). The potential for chatbots in negotiated learner modelling: a wizard-of-oz study. In *International Conference on Intelligent Tutoring Systems* (pp. 443–452). Springer Berlin Heidelberg.
- Kerly, A., & Bull, S. (2008). Children’s interactions with inspectable and negotiated learner models. In *International Conference on Intelligent Tutoring Systems* (pp. 132–141). Springer Berlin Heidelberg.
- Kerly, A., Hall, P., & Bull, S. (2007). Bringing chatbots into education: Towards natural language negotiation of open learner models. *Knowledge-Based Systems*, 20(2), 177–185.
- Kerly, A., Ahmad, N., Bull, S. (2008a). Investigating Learner Trust in Open Learner Models Using a ‘Wizard of Oz’ Approach. *International Conference on Intelligent Tutoring Systems* (pp. 722–724). Springer, Berlin.
- Kerly, A., Ellis, R., & Bull, S. (2008b). CALMsystem: A Conversational Agent for Learner Modelling. *Knowledge-Based Systems*, 21(3), 238–246.
- Lawless, K. A., & Brown, S. W. (1997). Multimedia learning environments: Issues of learner control and navigation. *Instructional Science*, 25(2), 117–131.
- Lehman, B., Matthews, M., D’Mello, S., & Person, N. (2008). What are you feeling? Investigating student affective states during expert human tutoring sessions. In *International Conference on Intelligent Tutoring Systems* (pp. 50–59). Springer Berlin Heidelberg.
- Lepper, M. R., Woolverton, M., Mumme, D. L., & Gurtner, J. (1993). Motivational techniques of expert human tutors: Lessons for the design of computer-based tutors. *Computers as Cognitive Tools*, 1993, 75–105.
- Long, Y., & Alevan, V. (2013). Skill diaries: Improve student learning in an intelligent tutoring system with periodic self-assessment. In *International Conference on Artificial Intelligence in Education* (pp. 249–258). Springer Berlin Heidelberg.
- Mabbott, A., & Bull, S. (2006, June). Student preferences for editing, persuading, and negotiating the open learner model. In *International Conference on Intelligent Tutoring Systems* (pp. 481–490). Springer Berlin Heidelberg.
- Malone, T. W., & Lepper, M. R. (1987). Making learning fun: A taxonomy of intrinsic motivations for learning. *Aptitude, learning, and instruction*, 3(1987), 223–253.
- Miao, Y. (2008). An intelligent tutoring system using interest based negotiation. In *2008 10th International Conference on Control, Automation, Robotics and Vision*, (pp. 2219–2223).
- Mitrovic, A., & Martin, B. (2002, May). Evaluating the effects of open student models on learning. In *International Conference on Adaptive Hypermedia and Adaptive Web-Based Systems* (pp. 296–305). Springer, Berlin Heidelberg.
- Mitrovic, A., & Martin, B. (2007). Evaluating the effect of open student models on self-assessment. *International Journal of Artificial Intelligence in Education*, 17(2), 121–144.
- Paiva, A., Self, J., & Hartley, R. (1995, August). Externalising learner models. In *Proceedings of World Conference on Artificial Intelligence in Education* (pp. 509–516).

- Person, N. K., Graesser, A. C., Kreuz, R. J., & Pomeroy, V. (2003). Simulating human tutor dialog moves in AutoTutor. *International Journal of Artificial Intelligence in Education (IJAIED)*, 12, 23–39.
- Rosé, C. P., et al. (2003) The role of why questions in effective human tutoring. *Proceedings of the 11th International Conference on AI in Education*.
- Schraw, G., & Dennison, R. S. (1994). Assessing metacognitive awareness. *Contemporary Educational Psychology*, 19(4), 460–475.
- Self, J. (1998). The defining characteristics of intelligent tutoring systems research: ITSs care, precisely. *International Journal of Artificial Intelligence in Education (IJAIED)*, 10, 350–364.
- Suleman, R. M., Mizoguchi, R., & Ikeda, M. (2015). Negotiation-driven learning. In *International Conference on Artificial Intelligence in Education* (pp. 470–479). Springer, International Publishing.
- Swanson, H. L. (1990). Influence of metacognitive knowledge and aptitude on problem solving. *Journal of Educational Psychology*, 82(2), 306.
- Tao, X., Miao, Y., Shen, Z., Miao, C., & Yelland, N. (2006). Interest based negotiation automation. In *International Conference on Intelligent Computing* (pp. 211–222). Springer, Berlin Heidelberg.
- Toulmin, S. E. (1958). *The uses of argument*. Cambridge: Cambridge University Press.
- Van Labeke, N., Brna, P., & Morales, R. (2007). Opening up the interpretation process in an open learner model. *International Journal of Artificial Intelligence in Education*, 17(3), 305–338.
- VanLehn, K. (1988). Student modeling. *Foundations of Intelligent Tutoring Systems*, 55–78.
- VanLehn, K. (1996). Conceptual and meta learning during coached problem solving. In *International Conference on Intelligent Tutoring Systems* (pp. 29–47). Springer Berlin Heidelberg.
- Weber, G., & Brusilovsky, P. (2001). ELM-ART: An adaptive versatile system for Web-based instruction. *International Journal of Artificial Intelligence in Education (IJAIED)*, 12, 351–384.
- Wenger, E. (1987). *Artificial intelligence and tutoring systems: computational and cognitive approaches to the communication of knowledge*. California: Morgan Kaufmann Publisher.
- Wolf, B., Bursleson, W., Arroyo, I., Dragon, T., Cooper, D., & Picard, R. (2009). Affect-aware tutors: recognising and responding to student affect. *International Journal of Learning Technology*, 4(3–4), 129–164.
- Wolf, B. P., Arroyo, I., Muldner, K., Bursleson, W., Cooper, D. G., Dolan, R., & Christopherson, R. M. (2010). The effect of motivational learning companions on low achieving students and students with disabilities. In *International Conference on Intelligent Tutoring Systems* (pp. 327–337). Springer Berlin Heidelberg.
- Zapata-Rivera, D., Hansen, E., Shute, V. J., Underwood, J. S., & Bauer, M. (2007). Evidence-based approach to interacting with open student models. *International Journal of Artificial Intelligence in Education*, 17(3), 273–303.