# Highlights

**Exploring Emerging Technologies for Requirements Elicitation Interview Training: Empirical Assessment of Robotic and Virtual Tutors**

Binnur Görer,Fatma Başak Aydemir

- We propose the extensible architecture REIT for requirements elicitation training systems and demonstrate its applicability by implementing two instances of it with different agent structures and interaction modalities.

- We empirically evaluate the two interview training systems – RoREIT with an embodied robotic agent and VoREIT with a virtual voice agent – with the students of a graduate level requirements engineering course.

- The participants rated VoREIT more favorably for the ease of use and engagement while RoREIT yielded a notably higher learning gain compared to VoREIT.

- Our findings indicate that each system has its own distinct advantages and weaknesses. Software engineering educators can customize REIT based on their needs and available resources.

- We publicly share our system implementation and study materials in Görer and Aydemir (2023a).

# Exploring Emerging Technologies for Requirements Elicitation Interview Training: Empirical Assessment of Robotic and Virtual Tutors

Binnur Görer[a,*], Fatma Başak Aydemir[a]

[a]*Boğaziçi University, Istanbul, Türkiye*

## ARTICLE INFO

## Abstract

Requirements elicitation interviews are a widely adopted technique, where the interview success heavily depends on the interviewer's preparedness and communication skills. Students can enhance these skills through practice interviews. However, organizing practice interviews for many students presents scalability challenges, given the time and effort required to involve stakeholders in each session. To address this, we propose REIT, an extensible architecture for Requirements Elicitation Interview Training system based on emerging educational technologies. REIT has components to support both the interview phase, wherein students act as interviewers while the system assumes the role of an interviewee, and the feedback phase, during which the system assesses students' performance and offers contextual and behavioral feedback to enhance their interviewing skills. We demonstrate the applicability of REIT through two implementations: RoREIT with a physical robotic agent and VoREIT with a virtual voice-only agent. We empirically evaluated both instances with a group of graduate students. The participants appreciated both systems. They demonstrated higher learning gain when trained with RoREIT, but they found VoREIT more engaging and easier to use. These findings indicate that each system has distinct benefits and drawbacks, suggesting that REIT can be realized for various educational settings based on preferences and available resources.

## 1. Introduction

Requirements elicitation is a set of activities to gather stakeholders' needs and desires for a system-to-be (Zowghi and Coulin, 2005) and is crucial for the success of software projects. Software developers can reduce project risks, improve team communication, and deliver high-quality software products that meet customer needs and desires by investing time and effort in requirements elicitation (van Lamsweerde, 2009). Among various elicitation techniques, interviews are the most popular and effective (Davis et al., 2006).

Requirements elicitation interviews require a combination of theoretical knowledge of interview techniques and soft skills such as interview management, behavioral control, and confidence (Hadar et al., 2014). It is essential to practice in a real interview setting to improve these soft skills and overcome interview nervousness (Andrews et al., 2006; Powell et al., 2021). A popular teaching strategy for requirements elicitation interviews is role-playing, which allows students to improve their abilities while playing the roles of stakeholders and requirements engineers. However, the excessive human effort required to plan and oversee these activities often makes them impractical in a regular classroom context (Debnath and Spoletini, 2020). Technological tools such as games, domain expert systems, and simulations offer an alternative approach to teaching requirements elicitation interviews by incorporating a digital interview partner. These tools can provide students with an interactive and engaging learning experience, allowing them to practice their skills in a controlled environment (Daun et al., 2021). Additionally, these tools can provide a standardized experience for all students, regardless of their partners' performance, ensuring that each student has an equal opportunity to develop their interviewing skills.

Emerging technologies have revolutionized the education ecosystem by creating new and innovative learning methods. For example, the rise of online learning platforms has enabled students to access educational content anywhere and anytime, making the learning experience more flexible (Alraimi et al., 2015). Social robots have become a prevalent source of help in educational activities where success depends on regular practice and attentive

---

*Corresponding author
✉ binnur.gorer@boun.edu.tr (B. Görer); basak.aydemir@boun.edu.tr (F.B. Aydemir)
ORCID(s): 0000-0001-9153-9244 (B. Görer); 0000-0003-3833-3997 (F.B. Aydemir)

supervision (Belpaeme et al., 2018). Gamification techniques have been used to motivate and engage students by incorporating game design elements into education contexts (Caponetto et al., 2014). Emerging technologies such as chatbots (Paschoal et al., 2018), serious games (Vega et al., 2009; Hainey et al., 2011; Yasin et al., 2018; García et al., 2020; Ibrahim et al., 2019; Garcia et al., 2019), simulators (Debnath and Spoletini, 2020) have also been used to support requirements engineering education.

This paper explores the application of emerging technologies in requirements engineering education, aimed at resolving scalability concerns arising from the considerable human resources required for conducting elicitation interview practices. We propose REIT as an extensible architecture for requirements elicitation interview training systems.

Systems implementing REIT feature two primary phases to help students practice and improve their elicitation interviewing skills. In the first phase, known as the interview session, the student takes the role of a requirements engineer while the system acts as a project stakeholder. Using a predetermined scenario, the system presents multiple choices for the following question for the student to ask at each turn of the interview and expects the student to choose the proper one. After the interview phase, the system passes to the tutoring phase, where it assesses the student's performance in the interview and provides feedback to improve their interview skills. The system presents the errors based on the choices of the student and allows the student to revisit the incorrect sections. It also offers contextual feedback to reinforce their learning. At the end of the feedback session, the system provides a behavioral analysis as potential soft skill improvement areas for the student.

REIT is a modular architecture that allows easy customization with alternative agent structures to accommodate the diverse needs and requirements of students and educators. We implement and evaluate two variants of REIT: RoREIT and VoREIT. RoREIT employs an embodied robotic agent, enabling audio-visual interaction capabilities to better emulate a human-like interviewing context, albeit at a higher cost due to the required robotic hardware. In contrast, VoREIT adopts a virtual voice-based agent providing voice-only interaction, a cost-effective alternative offering a reduced human-like setup.

We conducted a user study with the students of a graduate level requirements engineering course to evaluate RoREIT and VoREIT. The participants rated both RoREIT and VoREIT favorably for perceived attitudes and usefulness without a significant difference. Regarding perceived ease of use, the participants' scores for VoREIT were higher than RoREIT, and this disparity is statistically significant. VoREIT is also found to be more engaging by the participants, although both of the systems rated with scores more than moderate levels. RoREIT yielded a notably higher learning gain compared to VoREIT, where the difference in learning gains between the two systems is significant. These results indicate that neither system dominates the other; rather, each system has its own distinct advantages and weaknesses.

The main contributions of this work are as follows:

- We propose an extensible architecture for requirements elicitation training systems and demonstrate its applicability by implementing two instances of it with different agent structures and interaction modalities.

- We empirically evaluate the implemented interview training systems – RoREIT with an embodied robotic agent and VoREIT with a virtual voice agent – with the students of a graduate level requirements engineering course.

- We publicly share our system implementation and study materials in Görer and Aydemir (2023a).

REIT builds upon our previous work (Görer and Aydemir, 2023c). Yet, it is designed to have a higher level of flexibility to support diverse agent configurations instead of a single agent configuration presented in our previous work, and it incorporates a behavioral feedback analysis component to support the soft skill improvement of students that is missing from our previous work. The user study to evaluate systems implementing REIT is conducted with different subjects and research questions after their implementation.

The rest of the paper is organized as follows. Section 2 reviews the relevant related work for the use of emerging technologies in education in general and in requirements engineering (RE) education specifically. In Section 3, we outline the REIT architecture together with the system descriptions of RoREIT and VoREIT. In Section 4, we present our research questions, the design, and the findings of the user study. Section 5 covers the discussion points, the limitations of the study, and the implications of our system on RE education. Section 6 describes the threats to validity and how we address them. Finally, Section 7 provides the concluding remarks and future work.

## 2. Related Work

Implementing innovative technologies in education encompasses diverse digital learning methods like e-learning, game-based learning, and artificial intelligence (AI) assisted learning (Liaw, 2008; Sitzmann, 2011; Ciolacu et al., 2018). Although these technologies have undergone continuous development over the last decade, the COVID-19 pandemic has expedited the adoption of technological tools in education. Remote learning has necessitated that these technologies should serve all students from various backgrounds and age groups worldwide. Consequently, having specialized solutions created with students' needs in mind has become imperative amidst this rapid digital educational transformation (Zhao and Watterston, 2021). Moreover, the long-term applicability and accessibility of the systems turned out to be critical factors for institutions and instructors to adopt any technology in their education programs (Pelletier et al., 2022). As a result, assessing a technological system's effectiveness in facilitating students' learning experience and outcome and its affordability by the institutions became essential to ensure long-term sustainability and effectiveness (Lu et al., 2022).

Requirements Engineering Education and Training (REET) plays a pivotal role in equipping students with the necessary skills to confront the industry's challenges. To provide students with practical exposure in real-world scenarios, it is imperative to merge academic knowledge with hands-on training (Daun et al., 2021). However, conventional pedagogical approaches often necessitate instructors' or peers' active involvement during practical training, potentially constraining the duration and repeatability of projects. Although several research has suggested using technical solutions to enhance the effectiveness of REET to address this challenge, they are relatively scarce given the advances in educational technology (Daun et al., 2021). We briefly outline the studies that propose using emerging technologies in REET. We only included the studies with an implemented solution and excluded proposal-only research that lack any quantitative or qualitative evaluation.

The majority of earlier research recommends software solutions to help REET learners standardize and validate the requirements stated during the elicitation processes. To assist in the creation of accurate and consistent software requirements specifications, Garbers and Periyasamy (2006) implement an interactive editing tool that enables the user to develop an IEEE-compliant requirements document. Using the tool, they target to teach how to improve requirements quality. A small number of students evaluated the tool and offered their feedback on its usability though no detailed analysis was provided on the efficacy of the tool. The issue of verifying the accuracy of the elicited requirements is addressed in Ogata and Matsuura (2012) and Kakeshita and Yamashita (2015). Ogata and Matsuura (2012) implement an automated tool that takes model-based requirements prepared using Unified Modelling Language (UML) and creates a prototyped product model. Their goal was to help instructors confirm the validity of students' requirements more easily. A group of students tested the tool and created more specific requirements. However, the missing details about the experiment, like the studied business case and the number of created requirements, limit the validity of the experiment results. A more sophisticated tool is offered in Kakeshita and Yamashita (2015) to help students determine whether their requirements analysis models are correct by comparing them to the one the instructor provided. This approach helps the students learn how to create more accurate models. The authors present a very rough analysis of survey results to measure the perceived usability and usefulness of the tool. Liang and De Graaf (2010) remark on the need for remote collaboration of developers and customers while creating a requirements document. To teach this RE activity to the students, they develop a web-based wiki platform to encourage remote cooperation on the same requirements document. They present the feedback from students on the tool's usability but do not provide any evidence about its pedagogical efficacy.

The use of digital games in REET suggests that games can be an effective and engaging method for teaching RE concepts and skills. These games are designed to simulate the challenges and complexities of RE tasks, such as eliciting and prioritizing requirements, and provide learners with opportunities to practice RE in a safe and supportive environment. A digital game is proposed in Vega et al. (2009) to teach the stages of elicitation through the symbolic presentation of software requirements workshop activities. They aim to bring real-life elements of RE to the virtual world and offer an effective learning environment by playing. However, the study was limited in prototyping, and no user study was conducted. Hainey et al. (2011) develop a game as a motivating and engaging platform to teach requirements collection and analysis procedures. As the pioneering study in REET to show a comparative evaluation of games to the paper-based traditional educational methodologies, they conducted a user study with a large and diverse group of students to assess the learning effectiveness of the game and students' perceptions of it. Likewise, an educational game is proposed in Yasin et al. (2018) to teach the security requirements analysis process. Their user

study primarily assesses the learning experience of the students in terms of the game's usability and utility as well as its motivational impact on the students.

A three-dimensional (3D) serious game is created by Garcia et al. to improve students' understanding and application of requirements engineering methodologies (Garcia et al., 2019; García et al., 2020). The learning task in the study is to improve a set of pre-prepared requirements by inquiring more about them via elicitation interviews. The students were given the templates and asked to follow the game's guides to elicit better requirements. They conducted a large-group user study to measure the game's impact on student factors like motivation and satisfaction and the improvement of students' skills in preparing requirements with having a control group that did the same task without the game. Nonetheless, they did not present any statistical analysis of the game's effectiveness results. Another similar serious game was developed in Ibrahim et al. (2019) to create an interesting and engaging learning platform, consequently aiming to improve students' understanding of the RE concepts. They aim to offer diverse experiences with multiple scenarios by having multiple levels within the game. However, only one level of the game is implemented, and pilot testing is applied with a few students to obtain feedback on the game's usability. The utilization of virtual reality technology for UML modeling in a 3D environment is demonstrated in Ochoa and Babbit (2019). However, due to the lack of comprehensive user research, the suitability of the system for REET could not be assessed.

In a pioneering study outlined in Nakamura et al. (2014), the importance of incorporating a simulated stakeholder in requirements elicitation training is emphasized. The researchers designed a domain expert system that could monitor students' chat messages as they worked on eliciting requirements for a given project. The system intervened as necessary, encouraging students to ask questions and providing answers by conducting keyword-based searches on its extensive domain knowledge. Unfortunately, the effectiveness of the system was not adequately demonstrated in the user study because the scenarios used were too simplistic and did not require the students to seek assistance from the domain expert system. Paschoal et al. (2018) developed a prototype chatbot intending to aid students in improving their ability to elicit requirements. The chatbot was designed to adapt its responses to the student's expertise level, thereby appropriately customizing the complexity of the answer. Despite conducting an evaluation study to showcase the system's effectiveness compared to an out-of-context chatbot, the researchers utilized an unremarkable measurement that may threaten the study's validity. Another chatbot-based interview simulator is developed in Laiq and Dieste (2020). The authors use cloud-based AI systems that can recognize the questions of the users provided in natural language and answer the questions based on the provided context. However, the authors did not include a comprehensive experimental study to provide more information about the simulator's operation and efficacy. Debnath and Spoletini (2020) propose an interview simulator with a multi-modal conversational agent to allow students to practice elicitation interviews. The simulator can also evaluate users' responses and provide a report at the end of the interview. While their approach is a significant contribution to the literature in addressing the characteristics of real requirements elicitation interviews, the prototype they have developed falls short of their original proposal in terms of interaction and dialogue generation. Their preliminary analysis shows that the participants who used the simulator before interviewing a human fictional customer made fewer mistakes than those who did not use the tool. In Konlog and Spoletini (2023), a web-based application is proposed to help with the creation of efficient training plans that are in accordance with the resources the instructor has. The application can accommodate different training programs, like our training schemes, and make them all available for the user to choose from in accordance with their needs.

The existing literature shows that although there have been certain advancements in the incorporation of technology within REET, such efforts remain relatively limited when contrasted with other educational domains, such as K-12 and higher education (Leoste et al., 2021; Timotheou et al., 2023), both in terms of the quantity of conducted studies and the diversity of explored technologies. Most studies in REET perform pilot user studies with students to assess suggested solutions, revealing insightful information about the students' preferences. Nonetheless, despite the importance of proving the efficacy of the suggested solutions through controlled experiments, many studies frequently lack this component, as noted in Daun et al. (2017). Furthermore, none of them explain the long-term applicability and utility of the suggested technology in REET-related education programs. In our study, we offer reproducible systems and evaluate them in terms of both user preferences and system effectiveness. Our purpose is to make it simple for instructors to assess the presented systems by taking into account their benefits and drawbacks in relation to the education objectives and available resources of the institute.

# 3. Interview Training System

We propose a modular and extensible system architecture for Requirements Elicitation Interview Training, named REIT. The architecture is designed to support requirements elicitation interview training activities which employ two main phases: interview practicing and interview performance evaluation. In the first phase, the agent acts as a stakeholder and allows the interviewer to practice an elicitation interview. Following a pre-determined scenario, the system presents multiple options for the interviewer's next question at each interview stage. The scenario progresses based on the interviewer's responses. For each interview turn, the agent tracks and analyzes the interviewer's responses and facial expressions. In the interview performance evaluation phase, the agent acts as a tutor and revisits each erroneous turn, highlighting the incorrectly selected options and the reasoning behind them. The interviewer is given the opportunity to correct their answers.

REIT is built using the Robot Operating System (ROS) framework (Quigley et al., 2009). This approach simplifies the creation of a modular architecture, allowing different modules to communicate synchronously through the publisher/subscriber messaging protocol. REIT enables the use of physical or virtual agents. The system controller can incorporate various agent features and behaviors. In this study, we implement two versions of this architecture, each with different agents and associated features: RoREIT is built as a multimodal interactive robotic system, whereas VoREIT is implemented with a virtual voice-only agent[1].

In this section, we first introduce the components of REIT. We then describe its customized versions designed with a robotic agent RoREIT and a voice agent VoREIT. Following, we outline the interaction flow of the interview training process and each step's functionality.

## 3.1. System Architecture

The modular interview training architecture REIT is comprised of eight independent modules as shown in Figure 1, namely *Database, Speech Recognizer, Interaction Engine, Dialogue Displayer, Stream Recorder, Facial Expression Analyzer, Feedback Evaluator* and *Trainer Agent*. The computation load is distributed across computing resources to enable real-time application and prevent any delays. We used two standard consumer laptops with modest memory and computing power (16 GB memory, 2.6 GHz 6-core CPU). The entire system is built to function autonomously, except for the speech recognizer module.
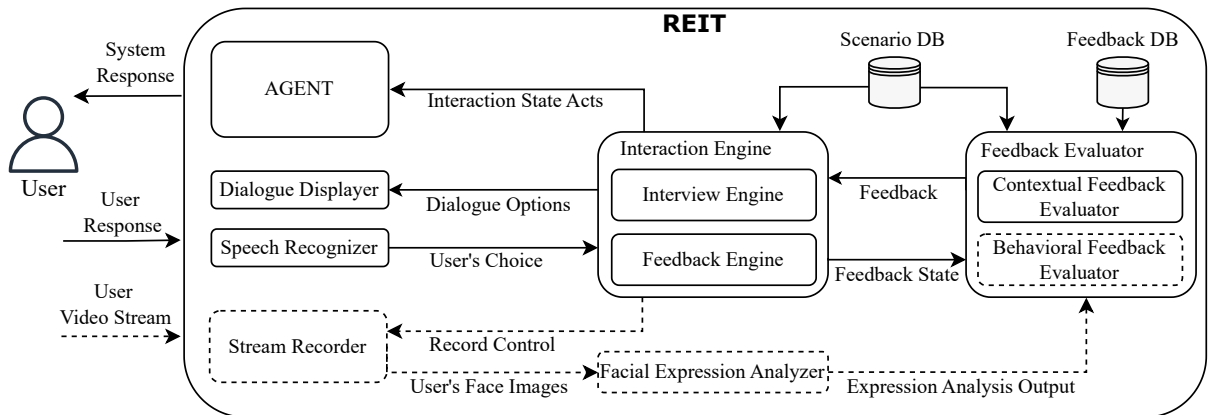


**Figure 1:** The architecture REIT for requirements elicitation interview training system. The dashed lines indicate optional inputs/outputs and modules.

*Database* includes scenario and feedback databases. The *scenario database* is a collection of scenario files. An interview scenario is set up as a dialogue in which the requirements engineer and the stakeholder alternate responses. The requirements engineer has three options to choose from at each turn: two of them are manipulated with some mistakes that are commonly encountered in real elicitation interviews as described in Bano et al. (2019), and the other is the correct response. The selected option of the requirements engineer determines the stakeholder response

---

[1]The code repository link for the systems' implementations is available in Görer and Aydemir (2023a).

and subsequent options for the requirements engineer for the following turn. As a result, the responses from the requirements engineer determine how the conversation flows across all possible dialogue paths in the scenario. The scenarios are created using the Twine tool[2] and exported to HTML format for deployment. The *feedback database* comprises feedback texts associated with each distinct mistake type. For example, the system provides the feedback "Please do not use technical jargon the stakeholder may not be familiar with" if the selected option contains technical jargon. We devise multiple feedback texts for each mistake type by rephrasing the text to avoid repeating identical content in each occurrence of a certain mistake.

*Speech recognizer* module converts speech inputs into text format. It captures user responses during both the interview and feedback sessions. In our study, this module was executed in a Wizard of Oz fashion, where a human operator was included in the process to convert the user's unstructured speech input (*User Response*) into the expected system input (*User's Choice*). The human operator handles participant speech in its original form and converts it into the intended system input only when it corresponds with one of the available dialogue options without adding their interpretations or judgments.

*Interaction engine* controls the interaction flow by iterating over the interaction steps. The overall interaction, described in Section 3.3, is modeled as a finite state machine. Each state has its defined actions and next-state transition logic. The engine decides the next state based on the current state and the relevant system inputs and variables defined in the state's transition logic. The interaction engine has two internal sub-engines devoted to managing the two main phases of the interview training:

- *Interview Engine* manages the interaction in the interview session. It queries the scenario with the given user's response (*User's Choice*) to retrieve the stakeholder response text, which is then sent to the agent module to be synthesized and spoken to the user (via *Interaction State Acts*). Meanwhile, the following question is collected from the scenario. As the agent finishes speaking the stakeholder response, the next question (*Dialogue Options*) is passed to the dialogue displayer, where the user will be shown the possible dialogue options.

- *Feedback Engine* controls the flow of the feedback session. The engine iterates over each incorrect interview turn. To remind the user of the interview stage, the preceding stakeholder text and dialogue options are sent to the dialogue displayer together with the previously chosen option. Simultaneously, the contextual feedback text (*Feedback*) is received from the feedback evaluator module. It is then delivered to the agent controller to be uttered by the agent. After the user makes a second attempt, the feedback state (*Feedback State*), which contains the user's updated response, is delivered to the feedback evaluator. The returned feedback (*Feedback*), indicating whether the user's second attempt is correct or incorrect, is communicated to the agent module and the dialogue displayer in order to notify the user both verbally and visually.

*Dialogue displayer* is a graphical user interface to communicate text-based input of REIT to the user, presented on the system's display. At each interview turn, it shows the dialogue options to the user to select one from (see Figure 7a). When the user makes a selection, the option is highlighted with a yellow background to let them know that the system understood their selection, as shown in Figure 7b. The tool is also used in the feedback session to display the revisited interview turns with incorrect responses. For each incorrect interview turn, the stakeholder's prior response and the corresponding options are displayed, with the incorrect option highlighted in red (see Figure 8a). The user's choice is highlighted in yellow following their second attempt to reassure the user that the system has recognized their input. The yellow background of the selected option will change to green if it is chosen correctly this time (see Figure 8b). Otherwise, it is turned red, and the correct option is shown with a green background to inform the user. At the end of the overall session, the feedback evaluator module evaluates the user's behavioral and contextual performance during the interview and provides a comprehensive analysis to the user. The analysis's findings are verbally communicated through supplementary visualizations. The dialogue displayer is utilized to present these visualizations to the user.

*Stream Recorder* captures the video stream from the user (*User Video Stream*), which is then processed by the facial expression analysis module. The recorded content includes both the user's speech and frontal face images. OBS Studio[3], an open-source recording software, is utilized for capturing the video content of the interview session. This

---

[2]https://twinery.org/
[3]https://github.com/obsproject/obs-studio

tool effectively captures, encodes, and records video content and provides a Python API for controlling the program inside REIT. During each interview turn, the recording starts as soon as the dialogue options are presented to the user and continues through the user's evaluation of the available option, the utterance of the selected option, and receiving the corresponding response of the agent. The recording then stops, saving the data captured during that turn for later analysis to provide feedback on the user's behavioral performance. The start and stop requests (*Record Control*) are managed by the interaction engine.

*Facial Expression Analyzer* evaluates the recorded face images of the user to determine their emotional state during each interview turn. This information is then used to provide feedback to the user, helping them to improve their soft skills. The goal of the emotional analysis is to give the user an understanding of their emotional state during the interview and highlight the significance of emotional control during requirements elicitation interviews. For the representation of the emotional states, we use Russell's circumplex model of emotions (Russell, 1980). This model describes emotions in terms of valence and arousal dimensional spaces. Valence refers to the degree of positivity or negativity associated with a particular emotion, while arousal refers to the level of excitement or calmness experienced. The FaceChannel library (Barros et al., 2020) is utilized for the automated analysis of emotional states in arousal and valence dimensions. The library performs real-time analysis of facial expressions, making it an efficient and practical solution for our experimental setup with limited computational resources. It predicts the user's emotional state from the given image frames. The predicted values are represented on a continuous scale ranging from -1 to 1, with negative values indicating a low level of arousal or valence and positive values indicating a high level.

*Feedback Evaluator* assesses the user's performance throughout the interview practicing. It has two components:

- *Contextual Feedback Evaluator* evaluates the user's responses for each of the interview turns. If the user's choice is incorrect, that interview turn is saved with the preceding stakeholder text, user dialogue options, and the selected incorrect option. The mistake in the incorrect option and the corresponding feedback text for that mistake type are retrieved from the scenario and feedback databases, respectively. The evaluation results are communicated to the feedback engine in *Feedback* message. During the feedback session, the feedback engine goes through each incorrect turn by displaying the interview turn and requests the agent module to mention the associated feedback verbally.

  The contextual feedback evaluator module is also utilized to evaluate the accuracy of the user's second attempt (sent by *Feedback State*) provided upon feedback. If the second attempt is also inaccurate, no additional feedback is given now. The evaluation result, indicating whether the user's second attempt is correct or incorrect, is communicated to the feedback engine, which notifies the user both verbally and visually through the agent and the dialogue displayer.
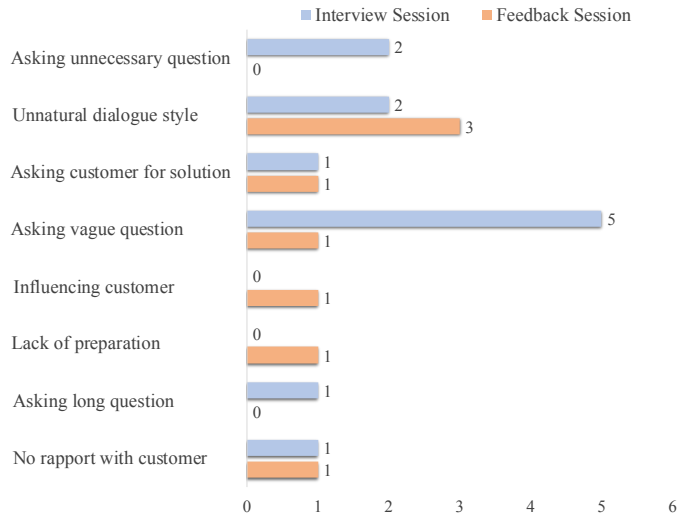
  At the end of the feedback session, an overall analysis of the contextual performance is presented to the user. The accuracy of the user's choice per each interview turn is visualized and presented to the user as shown in Figure 2a. Blue ticks indicate the right choice, while red exclamation marks indicate the wrong choice. The second attempt's results are also displayed to notify the user if an erroneous turn is fixed during the feedback session. Subsequently, the number of mistakes per category is shown for the interview and feedback session as in Figure 2b. The incorrect choice in each turn is manipulated with mistakes belonging to one or two categories. Hence, the total number of mistakes per category is greater or equal to the number of incorrect turns. Our goal with this analysis is to help the user to identify the mistake categories with which they struggle and monitor their development upon obtaining feedback. The accompanying speech texts, along with the visualizations in Figure 2, are sent to the feedback engine to be transferred to the agent module and the dialogue displayer.

- *Behavioral Feedback Evaluator* takes the output from the facial expression analyzer and uses it to evaluate the user's behavioral performance. The module processes the results, which include the user's emotional state represented in terms of pleasure and activation, and computes the median values for each interview round. These values are depicted on the 2D circumplex model as in Figure 3a. The categorical emotions are placed on the circumplex to help the user better understand the representation of emotions in the dimensions of pleasure and activation. While conducting a requirements elicitation interview, it is important to avoid communicating negative emotions to the stakeholder. These emotions, such as nervousness, stress, upset, sadness, depression, and boredom, can negatively impact the interaction and decrease overall performance. Hence, the regions

| Turn ID | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Interview Session | ✓ | ✓ | ! | ! | ✓ | ✓ | ! | ✓ | ! | ✓ | ! | ✓ | ! | ✓ | ! | ! | ! | ! | ✓ | ✓ |
| Feedback Session | ○ | ○ | ! | ✓ | ○ | ! | ○ | ✓ | ○ | ! | ○ | ! | ○ | ✓ | ! | ! | ✓ | ○ | ○ |

(a) The state of user choice for each interview turn is displayed as correct or incorrect by blue ticks and red exclamation marks, respectively. The accuracy of user's second attempt in the feedback session is also displayed for the incorrect turns.
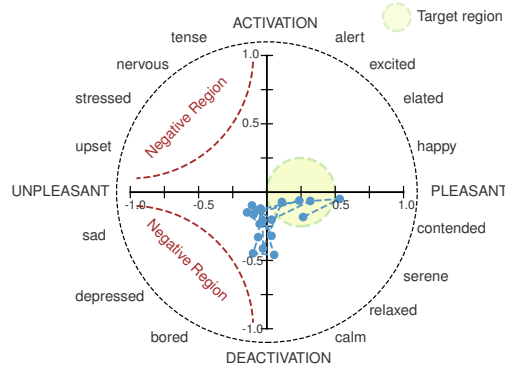


(b) Number of user mistakes per mistake class in the interview session (in their first attempts) and feedback session (in their second attempts).

**Figure 2:** Visual presentation of the overall contextual performance analysis of a sample user at the end of the session.
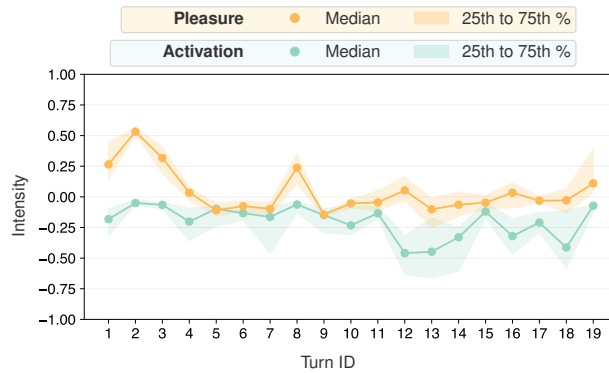
containing the negative emotions are indicated specifically on the circumplex. Emotional control is important for maintaining a professional and pleasant interaction during the interviews. It is desired for the analyst to display expressions with slightly positive emotions around neutral activation levels. To make it easier for the user to understand the agent's discussions around this topic, the desired emotional states are visually emphasized on the circumplex by highlighting the target region.

In order to provide a more informative analysis across the turns, the descriptive statistics of pleasure and activation values are calculated and visualized as shown in Figure 3b. The median value provides a measure of central tendency, and the interquartile range ($25^{th}$ to $75^{th}$ percentile) shows the spread of the data. This helps the user to understand how their emotional expressions changed during the interview and identify any specific turns where the emotional expressions deviated from their typical pattern. This analysis could provide insights into the user's emotional control during the interview and helps them to identify possible shortcomings for improvement.

*Trainer Agent* module is a versatile and adaptable component designed to interact with users in various ways. Its primary objective lies in providing users with an intuitive interactive experience while efficiently communicating system actions. In this study, the module's implementation extends to both embodied and virtual forms, accommodating physical robotic agents as well as voice-based conversational agents. However, the module is designed to be flexible to support other sorts of agents. For instance, it can readily integrate with virtually embodied 3D characters within virtual reality platforms, permitting customization to cater to distinct user experiences. This module's extensible design further facilitates the incorporation of additional components, such as speech-driven gesture generation and user-adaptive emotional expressions, to enhance the agent's interactive capabilities. These additional features can be orchestrated by

(a) User's facial expressions during each interview turn are displayed on the 2D circumplex model with dimensions of pleasure and activation. The median value is used as the central value to represent the overall emotional state during each turn. The blue circles show the median values of emotional intensity per each interview round.



(b) The visual representation of the descriptive statistics of the emotional expressions in the dimensions of pleasure and activation. The intensity value in each dimension ranges from -1 to 1. The median values and 25% to 75% percentile range are provided across all the interview turns.

**Figure 3:** Visual presentation of the overall behavioral performance analysis of a sample user at the end of the session.

the interaction engine through state actions. In our study, we implement two versions of the trainer agent module with different levels of interaction capabilities.

- *Robotic Agent* is implemented with a physically embodied robotic platform, the Nao robot, which is shown in Figure 4. As the most widely used humanoid bipedal robot for academic research, Nao is effective, affordable, and simple to program (Gouaillier et al., 2008). Its appealing design, expansive sensing, and acting capabilities make it sound for social robot-human interaction. It is 58 cm tall, mobile, and has auditory, visual, and tactile senses. Nao has two video cameras in the forehead and mouth, providing images with resolution up to 1280x960 at 30 frames per second. In our study, we used the $4^{th}$ version of Nao, which has the embedded Linux-based operating system OpenNAO and middleware NAOqi 2.1.4[4]. Naoqi is the main program of Nao which provides a low-level programming interface to control the joints of the robot, adjust LEDs color and intensity, and manage the built-in text-to-speech component.

  The robot controller module is depicted in Figure 5a. It is responsible for controlling the actuators on the robotic platform in order to carry out the actions requested by the interview engine (*Interaction State Acts*). The module manages the robot's speech synthesizer and the various sub-controllers of the robot's actuators, such as LEDs and motors. The robot controller communicates with the interaction engine to receive information about the current interaction state actions and converts them to the appropriate actions to be taken by the robot. The module
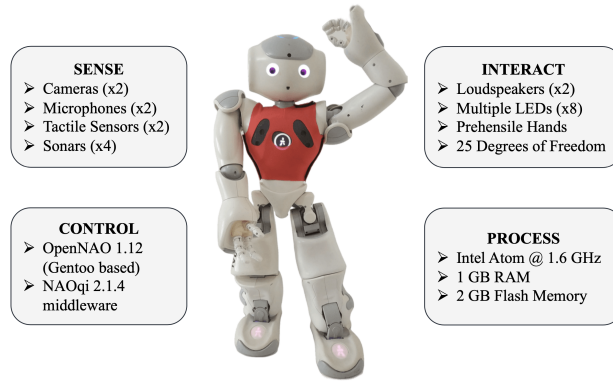
---

[4]http://doc.aldebaran.com/2-1/index.html

**Figure 4:** The Nao robot used in the robotic agent module.

is designed to provide seamless integration between the interview engine and the robotic platform, ensuring smooth and effective interactions between the user and the robot. The robot controller module is run on top of Naoqi middleware. Naoqi serves as a messaging interface between the Nao robot and the laptop and therefore is executed on both devices.

For the speech synthesizer module, we used the built-in functionality of the Nao robot. This module converts text into speech and plays it through the robot using the selected voice with customizable control commands such as pitch and speed. Besides speech, the robot has been programmed to use non-verbal cues such as head movements, gaze direction, and body gestures to convey social cues and enhance the overall interaction experience with the user. Emphatic gestures are created to express joy and despair in response to the users' correct and incorrect choices in their second attempts during the feedback session. Likewise, eye blinking is implemented using the eye LEDs of Nao to increase perceived liveness. These accompanying social skills aim to make the robot's communication more natural and engaging, leading to a better user experience and improved performance during the training.

- *Voice Agent* is designed to rely solely on speech-based communication without using any virtual or embodiment presence. It utilizes a speech synthesizer module to convert the text into speech (see Figure 5b). gTTS[5] is used as the text-to-speech library. The agent controller manages this library by transferring the text message provided by the interaction engine (via *Interaction State Acts*) and playing it back to the user through a speaker. The voice agent has no appearance or embodiment. A static image representing the agent is shown on the dialogue display of the user while the agent is speaking.



(a) The robotic agent module.

(b) The voice-based agent module.

**Figure 5:** The agent module of REIT is implemented with a robotic agent and a voice-based agent for RoREIT and VoREIT, respectively.

## 3.2. RoREIT and VoREIT

REIT is built upon our previous work (Görer and Aydemir, 2023c), targeting a more flexible architecture that can adapt to different agent configurations, feature sets, and scenarios. Our previous observations align with the common

---

[5]https://gtts.readthedocs.io/en/latest/

"no size fits all" concept in the education literature, which recognizes that students may have different preferences and needs. Likewise, not all educational institutions may have the resources or infrastructure to support and apply a training system with an advanced technological component.

The trainer agent module is designed to be flexible and adaptable to various types of implementations. This allows REIT to be customized based on the specific requirements and use cases. The choice of the agent module may depend on the desired interaction modality and the available resources. By including or excluding different modules and choosing the appropriate agent module, REIT can be tailored to meet specific needs and provide a customized experience for the user. We implement two different versions of REIT; Robotic Requirements Elicitation Interview Trainer RoREIT and Voice-based Requirements Elicitation Interview Trainer VoREIT. Both implementations use the same core components, which are essential for a practical and effective training experience, such as the interaction engine and contextual feedback analyzer. The two systems fundamentally diverge in agent implementation and feedback evaluation aspects, which are influenced by the interaction mode the system offers.

- RoREIT provides an audio-visual interaction with the user. The embodied robotic agent used in this system can communicate with the user through speech and body gestures. Likewise, the user engages with the system via voice and video. The RoREIT's audio-visual interaction enables the delivery of behavioral feedback. The modules required for the behavioral feedback analysis, including stream recorder, facial expression analyzer, and behavioral feedback evaluator, are activated in this system (see the optional components shown in Figure 1). Despite these advantages, the system comes with a higher hardware cost, and the embodied robotic agent can interact with a single student at a time, making it less scalable then the implementation described below.

- VoREIT provides an interaction that is purely audio based. The agent lacks visual representation and has no visual modality in its communication. Video streaming is excluded in both ways of communication between the user and the system. Hence, the behavioral feedback analysis feature is not available in this version of REIT. This design choice facilitates exploring and comprehending users' preferences on two distinct systems with different interaction modalities and feedback utility within the REIT framework. VoREIT provides a more cost-effective and less resource-intensive solution compared to the robotic agent version. However, RoREIT provides more human-like communication and extended feedback evaluation with behavioral analysis.

### 3.3. Interaction Flow

This section outlines the procedure for trainee interviewers utilizing the system. The session is divided into four distinct parts. To begin, the agent greets the interviewer and presents the system. The interview process then takes place, with the agent acting as a stakeholder for the project. Following the interview, the feedback session starts, where the agent acts as a tutor and reinforces feedback on the problematic parts of the interview session. The system concludes the training session by providing an overall analysis, including a review of all interview turns and how much improved during the feedback session, and behavioral performance evaluation. Figure 6 illustrates the flow of interaction, with the interview and feedback processes clearly marked. The interaction steps are as follows:

**T1: Greet the user.** The agent greets the user and exhibits a praising affinity with the target of establishing a positive rapport with the user and creating a friendly and welcoming interaction environment to boost the user's acceptance of the system.

**T2: Introduce the system.** The agent explains the function of the system and defines the user's role as the interviewer and its own role as the interviewee. The agent also provides a brief overview of the scenario for the session.

**T3a: Present the options of the next question to the user.** At this stage, the system presents the user with a screen displaying potential questions for the interviewer to select from, as shown in Figure 7a. The user is prompted to select one of these questions to ask as the next question in the interview process.

**T3b: Start tracking the user facial expressions.** The agent begins to record the user's facial images. These images are later evaluated to provide feedback on the user's behavior during the interview.

**T4: Select option.** The user, who is acting as the interviewer, evaluates the presented options, taking into account the response of the stakeholder (the agent), the interview context, and the direction of the interview. The user then selects one of the options and speaks it out loud. The selected option is highlighted on the dialogue display to confirm the user's choice, as shown in Figure 7b.
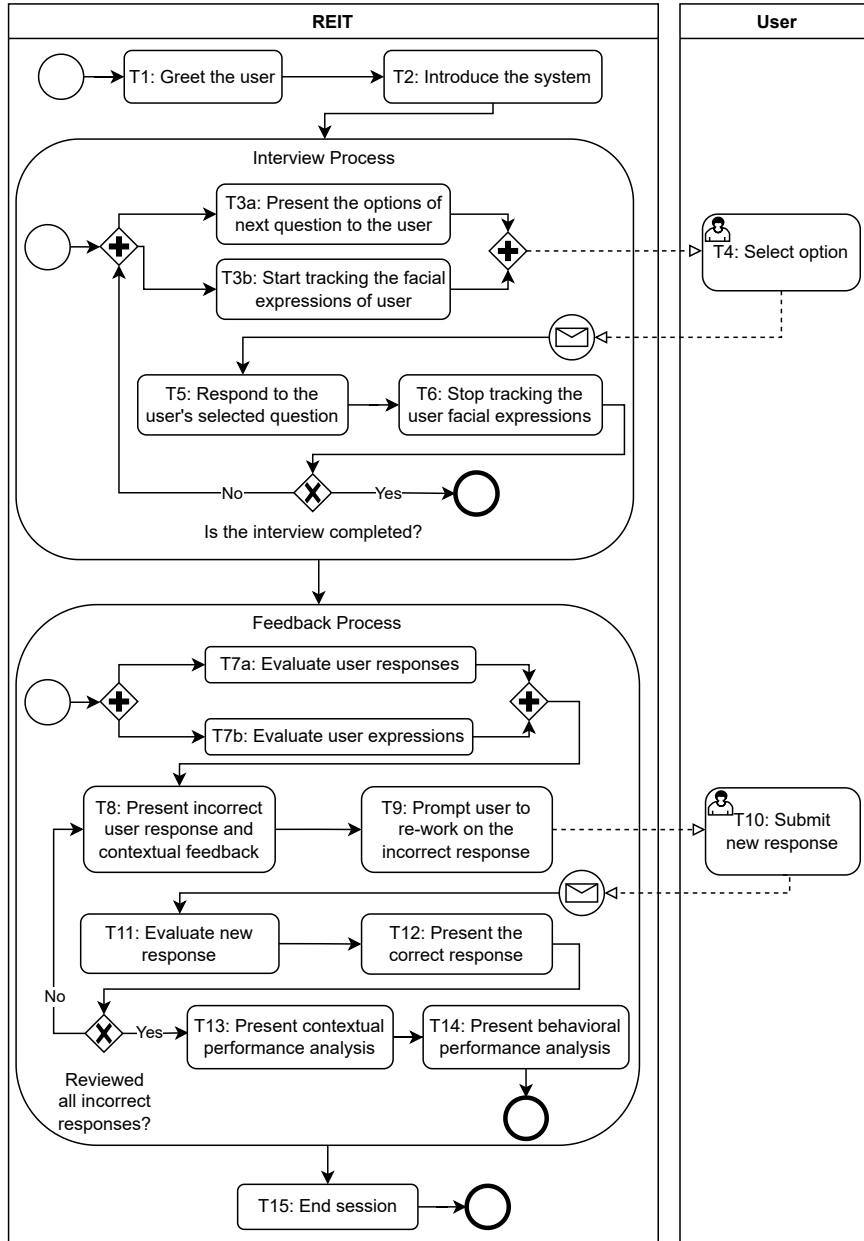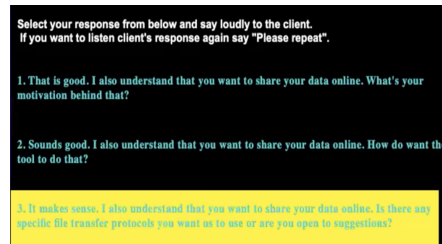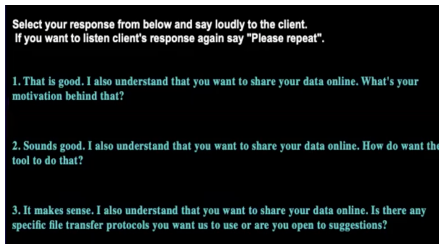
**Figure 6:** The interaction flow between the user and REIT.

**T5: Respond to the user's selected question.** The system recognizes the input from the user and records the user's choice. The agent then responds to the question selected by the user according to the scenario flow.

**T6: Stop tracking the user facial expressions.** The agent stops recording the user's facial images for that turn. If there are more questions in the scenario, this loop continues with the next turn. Otherwise, the interview session is concluded, and the training continues with the feedback session.

**T7a: Evaluate user responses.** The system detects the turns with incorrect responses by comparing the user responses with the correct options in the specified scenario. The user is then notified by the agent that the feedback session has started.

(a) The agent presents the possible questions for the user to select from.

(b) The agent confirms the user's choice by highlighting it.

**Figure 7:** Samples of the interview training system's dialogue displayer states during the interview session.



(a) The agent highlights any incorrect responses given by the user and provides an opportunity for the user to correct them.

(b) The agent notifies the user about the evaluation result of their second attempt on the incorrect response.

**Figure 8:** Samples of the interview training system's dialogue displayer states during the feedback session.

**T7b: Evaluate user facial expressions.** The system analyses the recorded facial images of the user to determine their level of arousal and pleasure during each interview turn. The system then computes the descriptive statistics and displays them, such as Figure 3, to give the user an understanding of their emotional state during the interview. The analysis is presented to the user in the later stages of the feedback session.

**T8: Present incorrect user response and contextual feedback.** The system revisits the incorrect turns that are evaluated in step T7a. For each turn, the agent's previous response to that turn and the available user options are presented. The incorrect response of the user is highlighted in red to remind their previous selection, as shown in Figure 8a. The agent provides the contextual feedback associated with the user's incorrect choice.

**T9: Prompt user to re-work the incorrect response.** Upon providing the feedback, the system gives another chance to the user to identify the correct response.

**T10: Submit new response.** The user re-evaluates the available options considering the interview context and the feedback provided. The user then makes a new choice and speaks it out loud.

**T11: Evaluate new response.** The system compares the user's new choice with the correct option in the predefined scenario and verifies if the new response is correct.

**T12: Present the correct response.** If the new response is correct, the system informs the user, as shown in Figure 8b. If the second trial of the user is also incorrect, the system displays the correct option with green while highlighting the others with red. For the second evaluation, no detailed feedback, like in step T8, is provided to keep the overall session length within limits.

**T13: Present contextual performance analysis.** Upon reexamining all incorrect interview turns, the agent informs the user that the feedback session has ended. The agent provides the user's overall performance on each interview turn by visualizing the accuracy of user choice during interview and feedback sessions. Moreover, the system displays the number of incorrect answers and their mistake categories (as defined by Bano et al. (2019)) both for interview and feedback sessions. A sample visual presentation of the analysis is shown in Figure 2.

**T14: Present behavioral performance analysis.** The system displays the behavioral analysis of the user's facial expressions for each interaction turns, as prepared in step T7b. In the meantime, the agent explains the goal and content of facial expression analysis and how this information might help the user better control their emotions during the interview.

**T15: End Session.** After giving the user a sincere thank you for taking part in the training, the agent requests the user to complete the post-session questionnaires and concludes the session.

## 4. Evaluation

This section presents the design, the execution procedure, and the empirical evaluation results of the user study. Our goal is to examine the comparative evaluations of the two systems: RoREIT with an audio-visual interaction interface involving an embodied physical robotic agent and VoREIT with an audio-only interaction having a virtual voice agent. We designed a user study wherein participants performed consecutive training sessions with both systems. During the sessions, the participants underwent an elicitation interview and received feedback on their performance, allowing them to learn from their mistakes. Half of the participants experimented with RoREIT first and then VoREIT, while the other half experimented with VoREIT first and then RoREIT. We collected users' perceptions of the systems' acceptability and engagement. We also measure the learning gain and interaction experience of the users. The following research questions (**RQs**) are addressed in our study:

**RQ1:** *How do RoREIT and VoREIT influence the learning gain of the participants?*
Both RoREIT and VoREIT systems are designed to deliver comparable interview and contextual feedback sessions, but they have different embodiment forms and interaction interfaces. Previous studies have suggested that robots may be more effective teachers than traditional training tools like web platforms or audiobooks, and that the physically embodied robots compared to their virtual equivalents can lead to measurable learning gains (Han et al., 2005; Leyzberg et al., 2012). These studies, however, did not focus on adult learners in real-world educational setups like ours. This RQ investigates the impact of two systems on participants' learning gains in the context of requirements elicitation interview training. As the measure of learning gain, we use the normalized change in the number of mistakes made by participants between two subsequent interview sessions.

**RQ2:** *How do RoREIT and VoREIT influence the processing speed of the users during the interview session?*
The relationship between the diversity of interaction modalities and cognitive task load is an interesting and complex research topic. Cao et al. (2009) suggest that exposure to diverse interaction modalities can lead to increased task load, as users need to process and adapt to different input and output channels. To investigate whether engagement in an audio-visual interaction with a physical robot might influence participants' task load differently from a voice-only interaction, we quantified participants' response speed during interview turns. In this context, we devise an in-context metric for processing speed, grounded in the principles of cognitive load theory in learning (Kirschner, 2002).

**RQ3:** *Does the participants' processing speed for a question vary based on their performance on the question?*
This RQ is focused on understanding whether there is a relationship between the participants' processing speed for a question and their performance on that question. Specifically, it aims to investigate whether participants' processing speed on the interview turns that they select the incorrect option differs from the turns they successfully respond to. The findings can provide insights into using our context-tailored processing speed measure to predict learners' performance during interviewing and improve the systems accordingly.

**RQ4:** *How do RoREIT and VoREIT influence the perceived acceptance of the underlying system in the dimensions of **RQ4a:** perceived attitudes, **RQ4b:** perceived ease-of-use, **RQ4c:** perceived usefulness*?
An individual's intention to use a technological system is determined by their attitudes toward that technology and their perceived ease of use and usefulness. The technology acceptance model is a framework that helps to understand how and why individuals use technology. This RQ targets to evaluate users' attitudes and perception of ease-of-use and usefulness towards the underlying system used for requirements elicitation interview training. We employed an expanded version of the technology acceptance model questionnaire (Yang and Yoo, 2004) to assess the participants' perceived acceptability of the two systems.

**RQ5:** *How do RoREIT and VoREIT influence the perceived and measured engagement levels during training with the system?*

Engagement is a crucial concept for the tutoring systems used in education because it can impact the effectiveness of the training process (Trowler, 2010; Appleton et al., 2006). A tutoring system should keep the learner connected and engaged in order to maximize the learning outcome. If learners are not engaged in the educational session, they may not be paying enough attention, which can hinder their understanding of the subject. Contrarily, if learners are engaged in the tutoring session, they are more likely to actively participate and pay attention, which can enhance their learning experience and improve their understanding of the learning material. Engagement can also increase the learner's motivation and interest in the subject, leading to greater enjoyment and success in their studies. Hence, we would like to analyze how the proposed systems can succeed in engaging the users during the interview training sessions. To answer this RQ, we examine both perceived and measured engagement. To assess perceived engagement, we gather user responses on a questionnaire about how engaged they feel when using the system. For measured engagement, we extract the arousal levels from the users' speech samples that are collected during the interview sessions.

**RQ6:** *What are the relationships between individual user characteristics (i.e., age, gender, interview experience level, interview anxiety level) and perceived acceptance and engagement of the systems?*

This RQ concerns the relationship between individual aspects of the participants and their perceived acceptance and engagement of the systems. We aim to understand whether certain personal aspects correlate with perceived acceptance and engagement scores. We explored age, gender, interview experience level, and interview anxiety level as personal aspects.

## 4.1. Study Design

To investigate the user experiences on RoREIT and VoREIT, we create a user study employing a combination of between-subject and within-subject design by following the guidelines in Hoffman and Zhao (2020). In successive experiment sessions, participants use both RoREIT and VoREIT, so the system itself is the condition varying within subjects. As the between-subject condition, we created two configurations: *setup A* and *setup B*, where participants are randomly assigned to one of the two setups. In setup A, participants experiment with RoREIT first, followed by VoREIT, whereas in setup B, participants experiment with VoREIT, followed by RoREIT. Every experiment session includes training for requirements elicitation interviews with the given system, followed by completing the post-condition questionnaire. Different scenarios are used in each participant's successive sessions, and scenarios are allocated in random order across the participants. In this way, we ensure that no participant is exposed to the same scenario twice and that the system and scenario pairings happen in equal numbers. The overall study design is shown in Figure 9.



**Figure 9:** The experimental design of the user study.

### 4.1.1. Procedure

Due to the shift towards remote education right after the Covid-19 pandemic restrictions, the experiment is planned to be carried out online using a video conference tool. The experiment announcement, which includes information about the overall study, the informed consent form, technical requirements, and expected duration, is prepared in written and video formats and is shared on the social media channel for the requirements engineering class. Potential participants are encouraged to review the announcement and assess their availability first before expressing their interest in participating in the study. If a candidate meets the eligibility criteria and agrees to the informed consent form, which

covers data protection and permission for video recording during the study, the experimenter reaches out to them to schedule an experiment time slot. Since English is the official language of education and is used by most business professionals in their work activities, including elicitation interviews, the experiment is prepared in English. Similarly, all materials, such as the study introduction, informed consent form, and administered questionnaires, are provided in English.

The target participant population's native language is not English, so they may not speak English fluently. Automated speech recognition systems (ASRs) demonstrate relatively reduced reliability when applied to non-native speakers in comparison to their performance with native speakers, particularly concerning unstructured speech (Radzikowski et al., 2019; Engwall et al., 2022). The structured format of the interview scenario adopted in our study could potentially rule out this issue by engaging ASRs with expected input sets. However, to expedite the prototyping of the systems, we opted for a Wizard of Oz approach. An experimenter assumes the role of the speech-to-text component, transcribing participants' speech into the expected system input. The participants, however, are not aware of the experimenter's presence and are explicitly informed before the study that the experimenter will leave the environment once the study begins and that no help or assistance will be given during the session. This is to prevent any potential participant bias caused by participants' desire to look their best in front of the experimenter, which could cause them to not be sincere and genuine in their behaviors during the experiment.

The below steps are followed during the experimental process:

1) Introduction: The experimenter dials into the video conference at the appointed time. When the participant joins, the experimenter describes the experiment's flow and responds to any questions they may have. Before the experiment begins, the participant is given the pre-experiment questionnaire to complete. The participant is asked to inform the experimenter once they are done with the questionnaire and ready to begin the study.

2) Training: The participant performs the two consecutive sessions using RoREIT and VoREIT in the order of the assigned setup. The flow of a session is detailed below:

   2a) Readiness check: The participant is informed of the physical requirements of the experiment before the study begins.

   – RoREIT condition: The experimenter makes sure the participant's face is visible and their voice can be clearly heard. The researcher then displays the dialogue displayer of RoREIT, which contains the live video of the robot and the dialogue options, on the shared screen (see Figure 10a). After getting the approval of the participant to start, the experimenter initiates RoREIT and notifies the participant that she will leave the environment.

   – VoREIT condition: There is no video interaction in this condition, and the participant conducts an audio-only interview with VoREIT. The experimenter ensures that the participant's camera is off and the participant's voice is audible clearly. The dialogue displayer of VoREIT, which includes the representative image of the text-to-speech agent and the dialogue options (see Figure 10b), is then presented on the shared screen by the experimenter. After receiving the participant's approval to start, the experimenter initiates VoREIT and informs the participant that she will leave the environment.

   2b) Session execution: Even though the participant is told to be left alone during the session, the experimenter stays in the environment to manage the system's speech-to-text component but remains completely invisible to the participant. She provides the selected option to the system after each verbal response of the participant. Through the process outlined in Section 6, the participant conducts the session, which comprises the interview and feedback portions.

   2c) Post-condition questionnaire: The experimenter returns to the call again once the session is over and asks the participant to complete the post-condition questionnaire.

3) Closure: The experiment is concluded once the two consecutive sessions and administered questionnaires are completed. The experimenter thanks the participant for joining and responds to any additional comments or questions from the participant.

To assess the feasibility of the experimental procedure and the anticipated length of the overall experiment, we conducted two pilot trials with our colleagues. The duration of the experiment is estimated to be 60 minutes, divided as follows: 5 minutes for the introduction and completion of the pre-experiment questionnaire; 50 minutes for the two
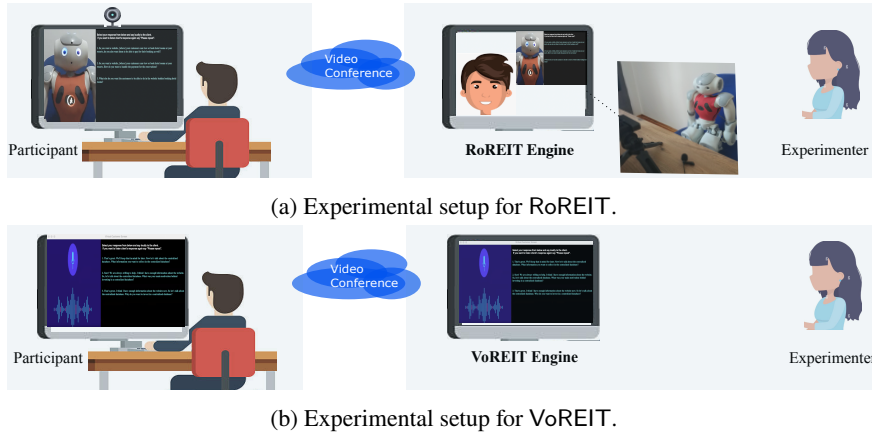
(a) Experimental setup for RoREIT.



(b) Experimental setup for VoREIT.

**Figure 10:** Experimental setup for the user study.

sessions, including completion of the post-condition questionnaire; and 5 minutes for the conclusion. We anticipate some variances in the overall experiment lengths, though, as the length of the feedback phase will depend on how many mistakes the participant makes.

The topological layout of the experimental setup of RoREIT is shown in Figure 10a. The study is placed in a room with a camera, a lapel microphone, the Nao robot sitting on a chair, and the laptop that runs the system components. For VoREIT setup, only the laptop to execute VoREIT system components including the virtual voice agent, is required, as illustrated in Figure 10b. For both systems, the experimenter is in the experiment room to observe the study for speech-to-text wizarding, but she is not visible to the participant. To reduce external distractions, participants are asked to take part in the experiment from a computer with a stable network connection in a well-lit, quiet, and comfortable setting. They are also required to have their camera turned on in RoREIT setup to enable audio-visual interaction.

### 4.1.2. Materials

The materials used in the experiment are available in our shared repository (Görer and Aydemir, 2023a).

***Scenario.*** We used two different scenarios in the experiments to minimize the learning effect over the course of the subsequent sessions. The scenario-system pairing is carried out in a randomized manner to alleviate any potential impacts of the scenario on the system. As one of the scenarios, we used the Cool Ski Resort scenario provided in Ferrari et al. (2020). A real-world use case, designing a website and its associated technological system for a resort company, is employed as the foundational scenario context. It is intended to elicit requirements for the website's design, social media integration, and database usage. We created another scenario, Cool Research Institute, by utilizing a comparable layout of the Cool Ski Resort scenario. In the second scenario, the institution's researcher, who is the stakeholder, wants to improve the processes for their research project. The goal is to create a new publication management system while also improving the existing data processing tool. The scenario aims to gather requirements for both the new publication management system and the enhancement of the existing data processing tool. This will include the addition of new features such as better data visualization and the ability to share data online. We carefully adjusted the dialogue portions for requirements engineer response options to keep the induced mistake counts and types closely similar between the two scenarios. Table 1 shows the total number of induced mistakes for each mistake class as described in Bano et al. (2019), which are counted over the requirements engineer's all possible responses in the given scenario. For both of the scenarios, the number of interaction turns for each interview party is set to a range of 15 to 19.

***Pre-experiment questionnaire.*** We collect the participants' demographics, their opinions about the robots, and their perceived interview anxiety.

- *Demographics:* In order to gain a better understanding of the participant population, we gathered information on their general demographics, which included age, gender, current occupation, and years of working experience. In addition to these basic demographic questions, we also prepared questions using 5-point Likert scale to get participants' proficiency and experience levels in conducting requirements elicitation interviews. We asked about their previous training and experience with eliciting software requirements, as well as the number of times they had practiced and conducted real interviews. Despite our target participant group being students of a graduate

**Table 1**
The occurrences of mistakes induced in *Cool Ski Resort* and *Cool Research Institute* scenarios along with their associated mistake class.

| ID | Mistake Class | Mistake Type | Number of Occurances | |
|----|---------------|--------------|---------------------|---|
| | | | *Cool Ski Resort* | *Cool Research Institute* |
| 1 | *Planning* | Lack of preparation | 4 | 6 |
| 2 | | Lack of planning | 3 | 3 |
| 3 | *Question Omission* | Not identifying stakeholders | 1 | 1 |
| 4 | | Not asking about existing system | 6 | 6 |
| 5 | | Asking long question | 3 | 2 |
| 6 | | Asking unnecessary question | 7 | 7 |
| 7 | *Question Formulation* | Asking customer for solution | 15 | 15 |
| 8 | | Asking vague question | 32 | 33 |
| 9 | | Asking technical question | 5 | 3 |
| 10 | *Order of interview* | Incorrect ending | 6 | 6 |
| 11 | *Customer interaction* | Influencing customer | 9 | 10 |
| 12 | | No rapport with customer | 16 | 18 |
| 13 | *Communication skills* | Unnatural dialogue style | 11 | 13 |

level requirements engineering course, we wanted to gather the self-evaluation of their proficiency in conducting requirements elicitation interviews. We asked them to rate their conceptual understanding of the fundamentals of requirements elicitation interviews, as well as their confidence level for practicing the interviews.

- *Negative Attitudes towards Robots:* We measured the participant's attitudes towards robots. If participants have a bias towards robots, this could impact the experiment result validity. The "Negative Attitude toward Robot Scale" (NARS) measures people's anxiety towards robots. It is designed to gauge people's attitudes and beliefs about robots, including their level of confidence and trust in robots, their understanding of the capabilities of robots, how at ease they are interacting with robots, and the extent to which they believe robots will have a positive or negative impact on society. We used the NARS version provided in Syrdal et al. (2009).

- *Interview Anxiety Questionnaire:* When people feel anxious during an interview, it can impact how they perceive and experience the interview process, including their interactions with the interviewer. McCarthy and Goffin (2004) developed a multidimensional measure of interview anxiety that provides a predictive anxiety level of a person, especially for selection interviews (e.g., job interviews). It typically includes a series of statements that assess the individual's level of nervousness or fear associated with the interview process. These involve the topics such as how confident they feel about their qualifications, how comfortable they are with being evaluated by others, and how much they anticipate experiencing anxiety during an interview. The measure provides an assessment of five interview anxiety dimensions: Communication, Appearance, Social, Performance, and Behavioral. We created a subset of the measure by reducing it to the items that are more relevant to the requirements elicitation interviews. Additionally, we changed the wording slightly to make the items more suitable for requirements elicitation interviews rather than the selection interviews. The adapted version of the measure is presented in Table 2.

***Post-condition questionnaire.*** Through a series of carefully designed and adapted questionnaires, we solicit feedback from participants on their interactions and experiences with the system in order to answer the research questions of the study.

- *Questionnaire on System Design:* We gathered information on the participants' perceived engagement level for the interview system. The questionnaire is designed using a 5-point Likert scale.

- *Technology Acceptance Model:* The technology acceptance model (TAM) was presented by Davis (1989) to reveal two key attitudes that drive the adoption of information technology to accomplish a task: perceived usefulness and ease of use. According to TAM, an individual's attitude towards technology is influenced by their beliefs about its usefulness and ease of use. TAM can help to understand user opinions and predict their intention to use it. This information can be useful for technology developers, as it can help them to improve the

**Table 2**

The interview anxiety questionnaire adapted from McCarthy and Goffin (2004). Each item is rated on a 5-point Likert scale (1 - Strongly disagree to 5 - Strongly agree).

| ID | Anxiety Scale | Item |
|----|---------------|------|
| 1 | Communication | I feel that my verbal communication skills are strong. |
| 2 | Social | While taking an interview, I become concerned that the interviewer will perceive me as socially awkward. |
| 3 | Performance | I am overwhelmed by thoughts of doing poorly when I am in an interview situation. |
| 4 | Appearance | Before an interview, I am so nervous that I spend an excessive amount of time on my appearance. |
| 5 | Communication | During the interviews, I often can not think of a thing to say. |
| 6 | Performance | During an interview, I worry about what will happen if the interviewer is not satisfied with my performance. |
| 7 | Social | I get afraid about what kind of personal impression I am making on interviews. |
| 8 | Appearance | I often feel uneasy about my appearance when I am being interviewed. |
| 9 | Social | I become very uptight about having to interact socially with an interviewer. |
| 10 | Social | I worry about whether the interviewer will like me as a person. |
| 11 | Performance | In interviews, I get nervous about whether my performance is good enough. |
| 12 | Communication | I get anxious in the interviews that I am unable to express my thoughts clearly. |
| 13 | Social | During an interview, I worry that my actions will not be considered socially appropriate. |
| 14 | Performance | During an interview, I am so troubled by thoughts of failing that my performance is reduced. |
| 15 | Communication | During the interviews, I find it hard to understand what the interviewer is asking me. |

technology and make it more appealing to users. Yang and Yoo (2004) refined the original TAM by also taking into account the affective and cognitive aspects of attitude and provided an expanded version. In our research, we used their version, which includes the following sections:

a) *Attitudes towards using* a technology can influence an individual's intention to use it. If an individual has a strong positive attitude towards technology, they will be more likely to use it, even if they perceive it as difficult to use or useless at all. On the other hand, if they have a weak positive attitude towards the technology, they may be more hesitant to use it, even if they perceive it as easy to use or useful.

b) *Perceived ease-of-use* refers to an individual's perception of how easy it is to use the technology. It is one of the key elements affecting a person's decision to use technology. The evaluation of ease-of-use may consider both physical and mental effort and an effort to learn how to use the technology.

c) *Perceived usefulness* is defined as the degree to which a person believes that using a particular system would enhance their job performance. It relates to how much the system can improve a person's task efficacy and how valuable the system is in connection to the content or goal of the task.

- *General remarks:* We asked participants about what they liked and disliked about the system and their suggestions for further improvements as open-ended text. We discuss their responses in Section 5.

### 4.1.3. Dependent Variables

The dependent variables arising from the RQs are *learning gain* (RQ1), *processing speed* (RQ2), *turn-specific processing speed* (RQ3) *perceived acceptance* (RQ4 and RQ6), *perceived engagement* (RQ5 and RQ6), and *measured engagement* (RQ5). Their formal definitions are given below.

*Learning gain* is a measure of the improvement in learning that occurs as a result of a particular instruction or intervention. It can be used to assess the effectiveness of educational programs and materials, and has been applied in a variety of settings, including K-12 education, higher education, and workplace training (Merchant et al., 2014; Roohr et al., 2017). Learning gains have been quantified in a variety of methods. Most frequently, this has been accomplished by tracking changes in the test results of the students during the course of instruction or training sessions. Hake (1998) introduced normalized gain to advocate a consistent analysis over diverse student populations with widely varying initial knowledge states. It is calculated by dividing the amount students learned by the maximum amount they could

have learned, as shown below:

$$\langle g \rangle = \frac{\langle post \rangle - \langle pre \rangle}{100 - \langle pre \rangle} \tag{1}$$

We adapted this measure to our study by considering the number of errors made by participants in a session as the success indicator. The difference between the number of mistakes in two consecutive sessions is used as how much each participant learned in the first session. In Equation 1, 100 is specified as the optimal expected score for the $\langle post \rangle$ evaluation, based on a scale of 0 to 100. For our context, we defined the optimal expected score for $\langle post \rangle$ evaluation as 0, as participants should aim to make no mistakes during the second session to maximize their learning gain.

More formally, let $P^A$ and $P^B$ denote the sets of participants in groups $A$ and $B$, respectively. For each participant $p \in \{P^A \cup P^B\}$, two consecutive interviews $I_i(p)$ are conducted, where $i \in 1, 2$. In group $A$, the first interview is conducted using RoREIT. In contrast, in group $B$, the first interview is conducted using VoREIT. Hence, the learning gain of a participant $p \in P^A$ associates with RoREIT whereas the learning gain of a participant $p \in P^B$ associates with VoREIT. $M_i(p)$ denotes the number of mistakes made by participant $p$ in the interview $I_i(p)$. The learning gain $G(p)$ of participant $p$ is then calculated as follows:

$$G(p) = \frac{M_2(p) - M_1(p)}{0 - M_1(p)} \tag{2}$$

*Processing speed* is a measure of how fast an individual can process information and respond to the surrounding environment (Salthouse, 1996). It is heavily influenced by the cognitive demands of the task at hand (such as reading, interviewing, or talking), the time limitations associated with it, and the social interactions that the person has to manage while performing the task. We designed an in-context measure of *processing speed PS* to quantify how quickly participants process each turn of the interview (Görer and Aydemir, 2023c). It is calculated by dividing the task load $TL$ by the time required to complete the task $RT$ in an interview turn, as shown in Equation 3 where $TL$ denotes the task load, and $RT$ represents the response time required to complete the task, and $t$ denotes the interview turn index of participant $p$ trained with system $sys$. We investigate the measure across the two systems, RoREIT and VoREIT, which employ setups with different interaction interfaces that hold distinct levels of social complexity. Indeed, utilizing response time alone as a metric is inadequate due to the different scenarios presented by each system, which directly impact the nature of the task involved.

$$PS_t(p, sys) = \frac{TL_t(p, sys)}{RT_t(p, sys)} \tag{3}$$

In each interview turn, the participant is presented with three options and expected to select one of them. We measure the response time $RT$ from when the participant is provided with the options until they respond. We calculate the task load per turn as the effort required in evaluating the options. As the options get longer and similar to each other, the participant is expected to put more effort into evaluating them to pick the correct option. The evaluation effort is estimated as the reading effort factored by the difficulty of an interview turn – which indicates how challenging it is for the participant to select the correct answer from the available options. According to the education literature (Ascalon et al., 2007; Shin et al., 2019), one of the elements of the difficulty index is the similarity of the multiple-choice options, which makes it harder for the examinee to eliminate incorrect options and identify the differences between the options. We calculate the similarity of the options using Universal Sentence Encoder[6] (Cer et al., 2018), which converts the text of each option item to a fixed-length vector representation and computes the cosine similarities between the option item vectors. More details of the calculation of *processing speed* are available in our previous work (Görer and Aydemir, 2023c).

The processing speed of participant $p$ for the session conducted with system $sys$ is then given by averaging $PS_t(p, sys)$ over all the turns $T$ as follows:

$$PS(p, sys) = \frac{1}{T} \sum_{t \in \{1...|T|\}} PS_t(p, sys) \tag{4}$$

---

[6]https://tfhub.dev/google/universal-sentence-encoder/4

*Turn-specific processing speed* $PS^\psi(p)$ calculates a participant's $p$ processing speed $PS$ for the given specific group of interview turns $t^\psi$, where $\psi$ denotes the group of the turns, i.e., correctly or incorrectly responded turns. It is given by averaging $PS_t^\psi(p, sys)$ over all turns, occurred in both RoREIT system $R$ and VoREIT system $V$, belonging to the specified group, as follows:

$$PS^\psi(p) = \frac{1}{\sum_{sys \in \{R,V\}} T^\psi(p, sys)} \sum_{sys \in \{R,V\}} \sum_{t}^{|T^\psi(p,sys)|} PS_t(p, sys) \tag{5}$$

*Perceived acceptance* refers to an individual's subjective belief or perception of the level of acceptance and usefulness of a particular technology. We investigated the perceptions of the participant $p$ for the acceptance of the given system $sys$. Upon completion of the training session with each system, the participants are asked to score the acceptance of the underlying system $sys$ by the technology acceptance model questionnaire described in Section 4.1.2. Using a 5-point Likert scale, the questionnaire measures the overall acceptance in three variables; attitudes $PATT(p, sys)$, perceived ease-of-use $PEU(p, sys)$, and perceived usefulness $PU(p, sys)$. All variables are integer values ranging from 1 to 5, with higher values indicating more positive attitudes, ease-of-use, and usefulness towards the technology, respectively.

*Engagement* refers to a state of active and meaningful involvement in a particular activity or interaction. A variety of cognitive, emotional, and behavioral processes are involved in this psychological and social phenomenon. While there has been a significant amount of research on engagement, there is still no consensus on a standard measure of engagement (Salam et al., 2022). This is because engagement is a complex and multifaceted construct that can manifest differently in various contexts and settings. There are many different approaches to measuring engagement, and the choice of measure can depend on a range of factors, including the purpose of the assessment, the specific context or activity being studied, and the characteristics of the individuals being assessed. Some commonly used measures of engagement include self-report questionnaires, behavioral observations, and physiological measures. In our study, we utilize both a self-report questionnaire and an analysis of voice signals to measure engagement.

- *Perceived engagement* score $PE(p, sys)$ is given by participant $p$ by the post-condition questionnaire prepared for the design evaluation of system $sys$. The variable $PE(p, sys)$ takes values in 1,...,5, where higher values indicate higher engagement.

- *Measured engagement* score $ME(p, sys)$ is calculated as the average arousal level of participant $p$ during the interview conducted by the system $sys$. Arousal is one of the important indicators of engagement (Ferrari et al., 2021), and it refers to the level of physiological and psychological activation or stimulation that an individual experiences in response to a particular stimulus or situation. High arousal is often associated with increased engagement and task adaptation in the context of technology use (Beaudry and Pinsonneault, 2010). Voice is an important modality for emotion transfer, and prosodic features, in particular, have been shown to be strongly correlated with arousal levels (Schirmer and Adolphs, 2017). Wagner et al. (2022) propose a transformer-based deep neural network model which constitutes the state-of-art in speech emotion recognition. By using their pre-trained model, we automatically extract arousal levels from speech segments. The overall arousal level across an interview is then calculated by taking the mean of the arousal levels of all turns.

## 4.2. Study Execution

We introduced the experiment to the students of the "Software Requirements Engineering" course in the Computer Engineering Department of Bogazici University. The department offers a one-year non-thesis M.Sc. degree for industry professionals besides a regular M.Sc. program in software engineering. This course is a selective graduate course offered in both programs. In order to allow the students to gauge their self-interest in participating in the study, we explained the main goal of the study, eligibility requirements, how the study will be conducted, and how long it will take. We did not mention the details regarding the experiment's research questions or the content of the scenarios used in the experiments to avoid any potential bias. Although the lecturer encouraged the students to take part in the study as an extracurricular class activity, participation was entirely up to the students. The participating students were offered to receive a small bonus grade. 27 students accepted to take part in the study. The study was conducted in the second semester of the 2021-2022 academic year.

**Table 3**
Demographics data of the participants provided separately for *Group A* and *Group B* as well as for all participants.

| | | Group A n=13 | Group B n=14 | All n=27 |
|---|---|---|---|---|
| Age(years) | range | 23 - 36 | 24 - 42 | 23-42 |
| | mean (SD) | 26.8 (3.8) | 29.3 (5.1) | 28.1 (4.7) |
| Gender | Female | 6 | 3 | 9 |
| | Male | 7 | 11 | 18 |
| Occupation | SW professsional | 9 | 12 | 21 |
| | Other | 4 | 2 | 6 |
| Years of Work Experience | 0-1 years: | 0 | 1 | 1 |
| | 1-3 years: | 10 | 7 | 18 |
| | 4-6 years: | 1 | 3 | 4 |
| | >6 years: | 2 | 3 | 5 |
| Confidence level in practicing RE interviews (1-lowest, 5-highest) | mean (SD) | 2.84 (0.77) | 3.07 (0.70) | 2.96 (0.74) |
| Level of theoretical knowledge on RE interviews (1-lowest, 5-highest) | mean (SD) | 3.00 (0.87) | 3.21 (0.77) | 3.11 (0.83) |
| # of practiced RE interviews | 0 | 6 | 6 | 12 |
| | 1-3 times | 5 | 6 | 11 |
| | 4-6 times | 1 | 0 | 1 |
| | >6 times | 1 | 2 | 3 |
| # of real RE interviews conducted | 0 | 10 | 6 | 16 |
| | 1-3 times | 2 | 5 | 7 |
| | 4-6 times | 0 | 0 | 0 |
| | >6 times | 1 | 3 | 4 |
| Level of Interview Anxiety (1-lowest, 5-highest) | mean (SD) | 2.46 (0.92) | 1.92 (0.79) | 2.18 (0.90) |
| Level of Negative Attitude Towards Robots (1-lowest, 5-highest) | mean (SD) | 2.30 (0.74) | 2.39 (0.80) | 2.35 (0.78) |

The participants, 9 females and 18 males ranging from 23 to 42 years of age (*mean* = 28.1, *SD* = 4.7), were all graduate students of the software engineering master program. We gathered information about their professions and levels of seniority at work. They all had full-time jobs. Out of the 27, 21 were employed in the software industry as developers, testers, analysts, or similar roles. The remaining six were working in fields other than software engineering, such as mechanical engineering and portfolio management. The majority of participants hold entry-level roles in their jobs and have worked for less than a year ($n = 1$), between 1-3 years ($n = 18$), between 4-6 years ($n = 4$), or longer than 6 years ($n = 5$).

We assessed the participants' knowledge of and experience with software requirements elicitation interviews. All of them were receiving training in software requirements engineering through the university course, and three of them had also joined company training programs related to this topic before. 15 of them practiced mock interviews before, one to three times ($n = 11$), four to six times ($n = 1$), and more than six times ($n = 3$). 12 of them had never participated in a mock interview. On the other hand, the number of participants who conducted a real interview is relatively lower. Only 11 of them had an interview with a real stakeholder one to three times ($n = 7$), more than six times ($n = 4$). We also gathered participants' self-evaluation scores on theoretical understanding of requirement elicitation interview techniques (*mean* = 3.11, *SD* = 0.83) and confidence level in practicing an interview (*mean* = 2.96, *SD* = 0.74) over a 5-point Likert scale (1=poor, 5=very good). The interview anxiety level is measured on a scale of 1 to 5, where 1 represents the lowest level of anxiety and 5 represents the highest level of anxiety. The mean anxiety level of the participants was 2.18 (*SD* = 0.90). The participants' negative attitude score towards robots was 2.35 on average (*SD* = 0.78). It is also measured on a scale of 1 to 5, where 1 represents the lowest level of negative attitudes toward robots and 5 represents the highest level of negative attitudes towards robots. None of the participants were native English speakers, but they were all proficient in the language with varying levels. The participants' demographics are given in Table 3.

Out of the 27 participants, 13 are assigned to *Setup A* condition (Group A), and the other 14 are assigned to *Setup B* condition (Group B) in a randomized manner. Using an online scheduling tool, the experiment hours were set in advance based on the participants' preferences. At the appointed time, the participant and experimenter joined the call, and the experiment was carried out using the procedures outlined in Section 4.1.1. The interview session of the experiment lasted 11.81 minutes on average ($SD = 2.71$) for RoREIT condition and 11.22 minutes ($SD = 2.34$) for VoREIT condition. The difference is mostly because of the slightly slower speech rate of the text-to-speech module utilized in RoREIT compared to the one used in VoREIT. The interview session lengths show similar deviations across the participants for both conditions, which is caused by the participant's processing speed and how many turns they have. The participants also spend considerable time in the feedback session with an average of 9.44 minutes ($SD = 2.37$) in RoREIT condition and 6.98 minutes ($SD = 1.93$) in VoREIT condition. The number of mistakes visited and the participants' review time for their second attempts account for most of the differences in feedback session duration among the participants and between conditions. The overall study was completed in 20 days period.

### 4.2.1. Randomization Validation

We checked the validation of randomization to ensure that there is no systematic bias in the way participants are allocated to *Group A* and *Group B*. As the participants' baseline characteristics, we monitored their self-reported NARS scores, interview anxiety scores, and expertise levels in requirements elicitation interviews across the two groups. We believe that the unequal distribution of these elements could influence the results of the related research objectives. We applied Mann-Whithey U test (Nachar et al., 2008) to check if there was a significant difference between the two groups as the data for the three factors are not normally distributed. The participants assigned to the *Group A* and *Group B* did not significantly differ in their NARS scores ($U = 87.0, p\text{-}value = 0.85$) and interview anxiety scores ($U = 121.0, p\text{-}value = 0.13$). Likewise, the two groups' expertise levels in conducting requirements elicitation interviews did not show a significant discrepancy. The average of the self-reported theoretical and practical expertise scores was considered as the expertise level, yet there was no difference ($U = 83.0, p\text{-}value = 0.70$).

To maintain a consistent level of difficulty and interview length, the scenarios employed in the experiment are meticulously designed and randomized across the experimental conditions. We checked if the associated scenarios influence the evaluation of the systems for each dependent variable. A series of statistical tests show that there is not any significant effect of the scenario on the results.

Overall, our analysis reveals that there are no deviations from the intended randomization scheme.

## 4.3. Results

We conduct a series of statistical tests to examine a set of hypotheses derived from our research questions (RQs). The null and alternative hypotheses corresponding to each RQ are presented, along with the specific tests used to analyze them. We first perform Shapiro-Wilk normality tests to determine if the data is are normally distributed. For dependent samples of data that are not normally distributed, we employ the Wilcoxon signed-rank test (Woolson, 2007). For the data that fulfill the characteristics of normal distribution, we use independent $T$-test to compare the means of two independent groups and dependent (paired) $T$-test to compare the means of two measurements taken from the same participants (Cramér, 2016). All hypotheses are tested at a 95% confidence level ($p\text{-}value \leq 0.05$).

For our study, Likert scales are ordinal, with 1 denoting a strong disagreement and 5 denoting a strong agreement. It is not safe to assume that the intervals between the Likert values are the same even though they are ranked in a certain order. Because the mathematical processes required to obtain the mean and standard deviation are inappropriate for ordinal data, the common approach is to use the median as the measure of central tendency (Blaikie, 2003). We use Wilcoxon signed-rank test to evaluate our hypotheses about the variables that are assessed using a Likert scale as it has similar power to the $T$-test even for small sample sizes (de Winter and Dodou, 2010).

The descriptive statistics of the dependent variables of each question are reported in the related research question, including median ($Mdn$) and interquartile range ($IQR$) for the data that are not normally distributed, and mean ($Mean$) and standard deviation ($SD$) for normally distributed data. These statistics provide information on the central tendency and dispersion of the data, and allow for a comparison of the groups being studied. The corresponding hypothesis test results are also presented with $p\text{-}value$ and test statistics.

**RQ1:** *How do RoREIT and VoREIT influence the learning gain of the participants?*
To answer RQ1, we used independent samples of the learning gains of the participants of Group A and Group B. The definition of learning gain measurement is given in Section 4.1.3. In Group A, the participants trained with RoREIT

**Table 4**
The descriptive statistics and hypothesis test results for RQ1.

| | Group A | | Group B | | Independent T-test |
|---|---|---|---|---|---|
| | *Mean* | *SD* | *Mean* | *SD* | p-value (T stat) |
| Learning Gain (G) | 0.35 | 0.28 | 0.08 | 0.39 | **0.05 (2.01)** |

**Table 5**
The descriptive statistics and hypothesis test results for RQ2.

| | RoREIT | | VoREIT | | Wilcoxon signed rank test |
|---|---|---|---|---|---|
| | *Mdn* | *IQR* | *Mdn* | *IQR* | p-value (Z stat) |
| Processing Speed (PS) | 7.11 | 4.36 | 8.06 | 3.05 | 0.86 (181.0) |

first, followed by VoREIT, whereas in Group B, it is the opposite. The participant's learning gain is associated with the system utilized during the initial session, given that the learning impact of the system is measured by the difference in success across consecutive interview sessions. We argue that the learning gains of the participants trained with RoREIT in their first sessions are different from the learning gains of the participants who trained with VoREIT in their first sessions. Formally we have $G_{\text{RoREIT}} = \{G(p_i), i = 1...|P^A|\}$ and $G_{\text{VoREIT}} = \{G(p_i), i = 1...|P^B|\}$. The two-tailed null hypothesis is $H_{10} =$ "the participants' learning gain who are trained with RoREIT is equal to the ones who are trained with VoREIT" (i.e., $\mu_{G_{\text{RoREIT}}} = \mu_{G_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{11} =$ "the participants' learning gain who are trained with RoREIT is not equal to the ones who are trained with VoREIT" (i.e., $\mu_{G_{\text{RoREIT}}} \neq \mu_{G_{\text{VoREIT}}}$). Two-tailed independent T-test reveals significant results ($T = 2.01, p\text{-}value = 0.05$). Hence, we can reject $H_{10}$ in favour of $H_{11}$. The learning gain of the participants trained with RoREIT is higher than those trained with VoREIT. The test result and descriptive statistics are given in Table 4.

**RQ2:** *How do RoREIT and VoREIT influence the processing speed of the users during the interview session?*
We consider dependent samples of the processing speed variable *PS* from the same participant experimented with the two systems. Formally, we have $PS_{\text{RoREIT}} = \{PS(p_i, \text{RoREIT}), i = 1...|P|\}$ and $PS_{\text{VoREIT}} = \{PS(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{20} =$ "the processing speed in RoREIT condition is equal to the one of the VoREIT condition" (i.e., $\mu_{PS_{\text{RoREIT}}} = \mu_{PS_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{21} =$ "the processing speed in RoREIT condition is not equal to the one of the VoREIT condition" (i.e., $\mu_{PS_{\text{RoREIT}}} \neq \mu_{PS_{\text{VoREIT}}}$). Since the processing speed variable violates the normal distribution assumption with Shapiro-Wilk's test result ($W = 0.89, p\text{-}value = 0.01$), we applied Wilcoxon signed-rank test to check whether the processing speed is significantly different in RoREIT and VoREIT conditions ($H_{21}$). Although the processing speed is greater in RoREIT condition, the difference is not significant with $Z = 181.0, p\text{-}value = 0.86$ (as given in Table 5). Hence, we can not reject $H_{20}$.

**RQ3:** *Does the participants' processing speed for a question vary based on their performance on the question?*
To answer this RQ, we compared the processing speed of the participants on turns with no mistake that they responded with a correct option versus the mistaken turns, and see if there is a significant difference between the two. Our goal is to examine whether the processing speed *PS* of the participant is influenced by the accuracy of the answer. We consider dependent samples of the turn-specific processing speed variable $PS^\psi(p)$ for no-mistake ($\psi = NM$) and mistaken ($\psi = M$) responses from each participant. The two-tailed null hypothesis is defined as $H_{30} =$ "The processing speed for the responses with *no-mistake* is equal to the responses with *mistake*." (i.e., $PS^{NM}(p) = PS^M(p)$). The two-tailed alternative hypothesis is $H_{31} =$ "the processing speed for the responses with *no-mistake* is not equal to the ones with *mistake*" (i.e., $PS^{NM}(p) \neq PS^M(p)$). To test the hypothesis, we performed two-tailed Wilcoxon signed-rank test. The difference is significant with ($Z = 325.0, p\text{-}value < 0.001$) and we rejected $H_{30}$. The descriptive statistics and test results are given in Table 6. Participants processed questions they correctly answered more quickly than questions they incorrectly answered.
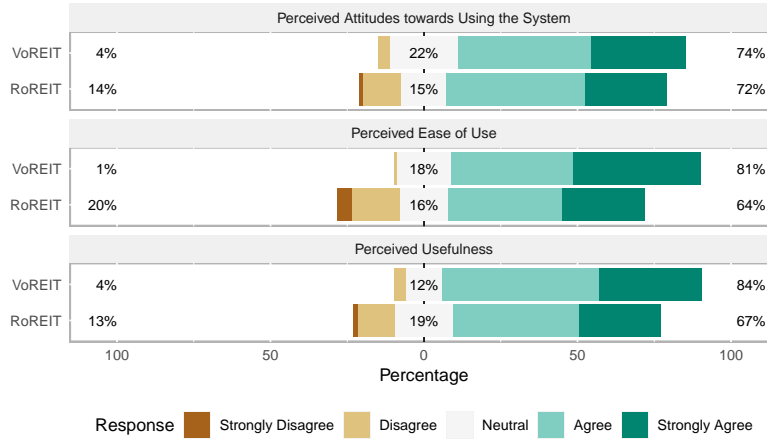
**RQ4:** *How do RoREIT and VoREIT influence the perceived acceptance of the underlying system in the dimensions of* **RQ4a:** *perceived attitudes,* **RQ4b:** *perceived ease-of-use,* **RQ4c:** *perceived usefulness?*
To answer the sub-research questions RQ4a, RQ4b, and RQ4c, we consider dependent samples of the dependent variables *PATT*, *PEU*, and *PU*, respectively, from condition RoREIT and VoREIT. To evaluate the reliability of the

**Table 6**
The descriptive statistics and hypothesis test results for RQ3.

|  | Mistake | | No Mistake | | Wilcoxon signed rank test |
|---|---|---|---|---|---|
|  | $Mdn$ | $IQR$ | $Mdn$ | $IQR$ | p-value (Z stat) |
| Processing Speed (PS) | 7.34 | 2.82 | 8.69 | 3.62 | **<0.001 (325.0)** |



**Figure 11:** The questionnaire results for the technology acceptance model for the conditions of VoREIT and RoREIT in the dimensions of attitudes toward using, ease-of-use, and usefulness.

questions in each of the surveyed dimensions, we employed Cronbach's $\alpha$ analysis, a widely used method for assessing the internal consistency of data. The reliability of the questions in each of the surveyed dimensions was higher than 0.70, indicating high reliability, particularly the questions of the dimensions of perceived attitudes towards using ($\alpha = 0.87$), perceived ease-of-use ($\alpha = 0.89$), and perceived usefulness ($\alpha = 0.87$). The scores of each dimension's questions were added up and averaged for each participant in each experimental condition. The average point of each dimension was then used to run the significance tests and to calculate descriptive statistics, which are given in Table 7.

**RQ4a:** We have $PATT_{\text{RoREIT}} = \{PATT(p_i, \text{RoREIT}), i = 1...|P|\}$ and $PATT_{\text{VoREIT}} = \{PATT(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{4a0}$ ="the perceived attitudes in RoREIT condition is equal to the one of the VoREIT condition" (i.e., $\mu_{PATT_{\text{RoREIT}}} = \mu_{PATT_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{4a1}$ ="the perceived attitudes for RoREIT condition is not equal to the one for the VoREIT condition" (i.e., $\mu_{PATT_{\text{RoREIT}}} \neq \mu_{PATT_{\text{VoREIT}}}$). We applied two-tailed Wilcoxon signed-rank test to test the significance. The difference is not significant ($Z = 59.5, p\text{-}value = 0.41$); hence, we failed to reject $H_{4a0}$. From Figure 11, we see that a high proportion of both groups evaluated the interview experience positively.

**RQ4b:** Formally, we have $PEU_{\text{RoREIT}} = \{PEU(p_i, \text{RoREIT}), i = 1...|P|\}$ and $PEU_{\text{VoREIT}} = \{PEU(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{4b0}$ ="the perceived ease-of-use in RoREIT condition is greater than equal to the one of the VoREIT condition" (i.e., $\mu_{PEU_{\text{RoREIT}}} = \mu_{PEU_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{4b1}$ ="the perceived ease-of-use for RoREIT condition is not equal the one for the VoREIT condition" (i.e., $\mu_{PEU_{\text{RoREIT}}} \neq \mu_{PEU_{\text{VoREIT}}}$). Two-tailed Wilcoxon signed-rank test result reveals a significant difference ($Z = 40.0, p\text{-}value = 0.03$) between the two groups for perceived ease-of-use. As seen from Figure 11, the ratio of responses on the strong agreement of ease-of-use is higher in VoREIT condition compared to RoREIT. The results indicate that VoREIT is perceived to be easier to use compared to RoREIT though both systems are rated highly positive.

**RQ4c:** Formally, we have $PU_{\text{RoREIT}} = \{PU(p_i, \text{RoREIT}), i = 1...|P|\}$ and $PU_{\text{VoREIT}} = \{PU(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{4c0}$ ="the perceived usefulness in RoREIT condition is equal to the one of the VoREIT condition" (i.e., $\mu_{PU_{\text{RoREIT}}} = \mu_{PU_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{4c1}$ ="the perceived usefulness for RoREIT condition is not equal to the one for the VoREIT condition" (i.e., $\mu_{PU_{\text{RoREIT}}} \neq \mu_{PU_{\text{VoREIT}}}$).

**Table 7**
The descriptive statistics and hypothesis test results for RQ4.

| | RoREIT | | VoREIT | | Wilcoxon signed rank test |
| --- | --- | --- | --- | --- | --- |
| | *Mdn* | *IQR* | *Mdn* | *IQR* | p-value (Z stat) |
| Perceived Attitudes (*PATT*) | 4 | 0.75 | 4 | 0.75 | 0.41 (59.5) |
| Perceived Ease of Use (*PEU*) | 4 | 1.5 | 4 | 1 | **0.03 (40.0)** |
| Perceived Usefulness (*PU*) | 4 | 2 | 4 | 1 | 0.13 (13.5) |

**Table 8**
The descriptive statistics and hypothesis test results for RQ5.

| | RoREIT | | VoREIT | | Wilcoxon signed rank test |
| --- | --- | --- | --- | --- | --- |
| | *Mdn* | *IQR* | *Mdn* | *IQR* | p-value (Z stat) |
| Perceived Engagement (*PE*) | 4 | 1 | 4 | 0 | **0.02 (39.0)** |

| | RoREIT | | VoREIT | | Dependent T-test |
| --- | --- | --- | --- | --- | --- |
| | *Mean* | *SD* | *Mean* | *SD* | p-value (T stat) |
| Measured Engagement (*ME*) | 0.46 | 0.09 | 0.45 | 0.09 | 0.65 (0.45) |

Two-tailed Wilcoxon signed-rank test shows the difference is not significant for the perceived usefulness for RoREIT and VoREIT ($Z = 13.5, p\text{-}value = 0.13$). As shown in Figure 11, participants evaluated both systems as highly useful.

**RQ5:** *How do VoREIT and RoREIT influence the perceived and measured engagement levels during training with the system?*
To answer RQ5, we considered both the perceived engagement levels reported by the participants and objectively measured engagement levels from the voice signals of the participants. For the first, we used dependent samples of the perceived engagement scores *PE* provided for both systems. Formally, we have $PE_{\text{RoREIT}} = \{PE(p_i, \text{RoREIT}), i = 1...|P|\}$ and $PE_{\text{VoREIT}} = \{PE(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{510} =$"the perceived engagement in RoREIT condition is equal to the one of the VoREIT condition" (i.e., $\mu_{PE_{\text{RoREIT}}} = \mu_{PE_{\text{VoREIT}}}$). The two-tailed alternative hypothesis is $H_{511} =$"the perceived engagement for RoREIT condition is not equal to the one for the VoREIT condition" (i.e., $\mu_{PE_{\text{RoREIT}}} \neq \mu_{PE_{\text{VoREIT}}}$). The result shows significance ($Z = 39.0, p\text{-}value = 0.04$) so that $H_{510}$ is rejected. The participants evaluated VoREIT as more engaging than RoREIT. As the second hypothesis, we postulate that there is no significant difference in the mean arousal level of the participants' speech across the two systems. Formally, we have $ME_{\text{RoREIT}} = \{ME(p_i, \text{RoREIT}), i = 1...|P|\}$ and $ME_{\text{VoREIT}} = \{ME(p_i, \text{VoREIT}), i = 1...|P|\}$. The two-tailed null hypothesis is $H_{520} =$"the measured engagement level in RoREIT condition is equal to the one of the VoREIT condition" (i.e., $\mu_{ME_{\text{RoREIT}}} = \mu_{ME_{\text{RoREIT}}}$). The two-tailed alternative hypothesis is $H_{521} =$"the measured engagement for RoREIT condition is not equal to the one for the VoREIT condition" (i.e., $\mu_{ME_{\text{RoREIT}}} \neq \mu_{ME_{\text{VoREIT}}}$). The relative $T$-test result shows no significance ($T = 0.45, p\text{-}value = 0.65$) so that $H_{520}$ can not be rejected. The perceived engagement scores are not equated with the measured levels from the voice signals of the participants. The statistics and test results for both dependent variables are provided in Table 8.

**RQ6:** *What are the relationships between individual user characteristics (i.e., age, gender, interview experience level, interview anxiety level) and perceived acceptance and engagement of the systems?*
For each of the user characteristics, we calculated its correlation with the dependent variables, namely perceived attitudes, perceived ease of use, perceived usefulness, and perceived engagement. This analysis is performed for both RoREIT and VoREIT as the dependent variables are collected separately for both of the systems. We use Kendall's rank correlation coefficient $\tau$, which is a statistic used to measure the ordinal association between the two variables. In total, 32 analyses were conducted, and four of them showed statistically significant correlations. The analysis revealed that age has a significant weak positive correlation with perceived attitudes towards using the system for VoREIT ($\tau = 0.35, p\text{-}value = 0.02$). Moreover, the correlation between interview anxiety level and perceived usefulness indicated a significant weak negative correlation for VoREIT ($\tau = -0.35, p\text{-}value = 0.04$). The interview anxiety level was found to have a significant moderate negative correlation with perceived engagement for VoREIT($\tau = -0.47, p\text{-}value = 0.007$). For the RoREIT system, the only significant correlation was observed between

interview experience level and perceived ease of use with a moderate negative correlation ($\tau = -0.41, p\text{-value} = 0.01$). These findings provide valuable insights into the relationship between user characteristics and system perceptions, which can inform the design of future personalized interview training systems. Further discussions are provided in Section 5.

## 5. Discussion

Our evaluation of the effect of RoREIT and VoREIT on the learning gain of the participants suggests that the physically embodied robotic agent of RoREIT resulted in higher learning gains than the disembodied voice agent of VoREIT. This result is consistent with the literature hypothesizing that embodied pedagogical agents have high capabilities to mimic real-world learning conditions and immerse students in learning activities (Grivokostopoulou et al., 2020; Davis et al., 2023). However, existing studies that evaluate embodied trainer agents in empirical research have inconclusive or inconsistent results for this claim (Darwish, 2014). Further research is needed to better understand the underlying processes and factors affecting the learning gain.

Participants demonstrated quicker processing times for questions they answered correctly compared to those they answered incorrectly. Leveraging our in-context processing speed measure, we can pinpoint questions that participants find challenging and provide timely assistance during the interview. This assistance might entail offering additional explanations or context to help learners in arriving at the correct answer.

VoREIT is perceived to be easier to use compared to RoREIT. Several factors may contribute to this. One would consider the voice-only interaction modality of VoREIT to be less challenging than the audio-visual interaction modality of RoREIT, which requires the participants to keep their cameras on. While physical robots may offer a higher sense of social presence, they may also create pressure for users to perform well and interact in a certain way. In contrast, a voice-only agent may feel more anonymous and less intimidating, allowing users to interact more freely. Familiarity with the underlying technology could also be another important factor. People may be unpracticed to interact with physical robots, which can raise their perceived complexity and decrease user comfort. The participants might also think that VoREIT is more accessible and manageable compared to RoREIT, which has limited serving capabilities due to the physical robot and, therefore, is more difficult to use.

Participants have reported a lower level of engagement with RoREIT compared to VoREIT, even though both systems received high engagement ratings. This discrepancy might be influenced by the relatively slower speech rate in RoREIT as slower speech may result in reduced user involvement and interest. Furthermore, the presence effect, as supported by prior research (Li, 2015), is another plausible explanation for this result. They show that users tend to hold more favorable views of physical robots when they are physically present in their environment, as opposed to when they are represented only through digital means, such as video feeds on a screen. Additionally, individual preferences and biases could have played a role in shaping participants' perceptions of engagement with the two systems. Although the measured engagement levels did not reveal a preference for VoREIT over RoREIT, there remains potential for improving RoREIT's perceived engagement. This could be achieved by enhancing the robot's voice for more effective communication or by incorporating the robot's physical presence into the user interaction.

Age has been known to play an important role in shaping one's perspective and adoption of new technologies (Hong et al., 2013). Older adults often exhibit a greater sense of comfort when interacting with familiar technologies, in contrast to advanced technologies with which they have limited exposure. Our empirical data also reveal a positive correlation between the age of the participants and the perceived attitude towards VoREIT. Additionally, our analysis demonstrates a negative correlation between participants' levels of interview experience and their perceived ease of use of RoREIT. Since the participants with more interview experience are more likely to have their own preferences for how interviews should be conducted, their expectations from a robotic training system could be more higher. In cases where the system falls short of meeting these expectations, it can lead to a perception of reduced ease of use, as they may compare the system's performance to their own established standards and practices.

Anxiety can make people self-conscious and focus on how well their performance is perceived by the interviewer rather than on the process itself, which may exacerbate the interview experience for both sides. Previous research on the use of technological tools for job interviews claims that a person's anxiety may be induced by several variables, such as the setting of the interview and the realism of the technology involved (Kwon et al., 2013; Vilar et al., 2020). Our analysis reveals that the participants' levels of anxiety negatively correlate with the perceived usefulness and engagement of VoREIT. One possible explanation for this is that the lack of a physical presence of VoREIT may contribute to feelings of discomfort for some users. When interacting with a disembodied voice, people may be more

likely to feel self-conscious or nervous, particularly if they have higher levels of interview anxiety. The robot's physical embodiment may have contributed to creating a more natural and intuitive interaction environment for the participants, which could have reduced their anxiety levels, hence offering them a better interview experience. However, we did not find any significant correlation between the participants' levels of anxiety and system perceptions for RoREIT. Further research is needed to understand underlying mechanisms and improve interview training systems' design accordingly.

*Qualitative remarks.* We asked the participants their most and least favorite aspects of the systems and their suggestions for improvements as open-ended questions. The participants generally agree that the contextual feedback component common to both systems is beneficial. The design of the dialog options received criticism. Some participants found it constraining their communication during the interview. They would like to be more flexible in the question formulation instead of selecting from the predefined question sets. The participants suggested employing further cases encountered in interviews, such as when the stakeholder deviates from the themes to be covered in the interview. They would like to experience and learn from different challenging cases. Some participants proposed an interactive feedback component to receive personalized feedback based on their level of expertise as an improvement.

The participants enjoyed the emphatic gestures of the robot performed in the feedback session when using RoREIT. One participant highlighted the importance of the eye contact with the stakeholder during the interview, thanks to the physical presence of the robot. The voice quality of the robot was primarily noted as a drawback. Since the participants were non-native English speakers and the interviews are conducted in English, language impediments could have been a potential issue. The Internet connection quality may have also affected the quality of voice communication. Some participants preferred a more human-like robot capable of making facial expressions. Two of them found the robot childish and non-realistic to play the stakeholder role. Nonetheless, it is clear from the participants' remarks that they have high standards for the robotic component, and technical excellence is desired. The behavioral feedback analysis of RoREIT is also frequently noted by the participants as a significant and helpful feature. For VoREIT, mixed views regarding the voice-based agent have been expressed. One participant appreciated how it resembled the voice assistant they use every day. In contrast, another person thought it felt talking on the phone and unenjoyable. Another participant appreciated having no visual interaction during the interview, but some others complained that the agent has no visual appearance and is too artificial.

*Limitations of the system.* Even though REIT is engineered to accommodate any scenario prepared in the expected format, one of its current limitations is the number of available scenarios. We built an additional scenario besides the one provided by Debnath and Spoletini (2020). Having a larger pool of scenarios would facilitate learners to practice the management of a broader range of requirements elicitation challenges. Interview scenarios based on prebuilt conversation graphs restrict users from expressing themselves freely. Having open-ended dialogues can help users expose themselves to diverse interview contexts. However, designing a dialogue system that can identify and address a broader spectrum of user responses poses challenges, particularly in the absence of datasets specifically created for the target task and domain (Ni et al., 2023). In our recent study, we investigate leveraging large language models (LLMs) for requirements elicitation interview dialogue generation through prompt engineering techniques (Görer and Aydemir, 2023b). Our findings show that LLMs face challenges in preserving coherence and logical flow, even when dealing with a small number of interview turns. Nonetheless, these models keep a promise to perform better on task-oriented dialogue generation when pre-trained with in-domain data (Gururangan et al., 2020; Chang et al., 2023).

*Implications on REET.* A considerable amount of resources and effort is required to organize realistic practice requirements elicitation interviews for students. REIT reduces the need for human stakeholders and can be adapted to various agent configurations and scenarios. Educators can specify their own scenarios for the desired domain, requirements, and choices for the students. Depending on the resources available, the agent who plays the stakeholder can be a humanoid robot as in RoREIT, a voice agent as in VoREIT. Other agent forms with tailored functions can also be easily utilized in the architecture. The system can be used for remote (as in our study) or in-person sessions. The participation in our study was voluntary, yet we propose including REIT to the curricula of REET. The system may be used as a practice ground for students to apply their theoretical knowledge. It may also be used as a graded activity where the students are graded based on their performance. Our shared repository includes the code, scenario, and other materials needed to run our implementations and further customize the systems.

## 6. Threats to Validity

In this section, we discuss the potential threats to our system and how we have attempted to mitigate them. The main threats identified, based on Wohlin et al. (2012), are as follows:

*Internal validity.* Maintaining internal validity is essential to ensure that the observed effects are not due to extraneous variables but rather to the investigated treatments. In our study, we utilized a within-subject design, which involves exposing participants to multiple treatments in varying orders across participants. One common internal validity threat for within-subject design is the order effect, where the order in which treatments are administered may affect the participants' responses. To mitigate this, we used counterbalancing, a technique that involves presenting the treatments in different orders across participants. This way, the potential impact of the order of treatments is equally distributed across the participants. Moreover, to mitigate possible maturation effects, the participants experimented with the two conditions sequentially without a break. We ensured they did not receive any external training across the conditions that could impact their performance. We also adjusted the length of the interview scenario and completed the overall session within one hour to minimize possible participant fatigue across the two consecutive sessions.

To minimize the potential social threats to internal validity, we took great care in designing the experiment's introduction for the participants. We provided only high-level context and the study's aims without disclosing details to prevent potential unequalization of experiment conditions, though every participant was treated with both systems. Participation was entirely voluntary, although the instructor encouraged the experiment as an external class exercise and rewarded it with a small bonus for the course grade. Student participants were explicitly informed that their performance during the interview would not affect their course grades. This ensured that the experiment did not generate any undue stress or pressure for the participants and that they could focus on the experiment's objectives without any external pressures.

*External validity.* In our endeavor to ensure the relevance and wider applicability of our study to real-world scenarios, we carefully pursued a high degree of realism in the experimental setup. We employed realistic scenarios during the interview sessions and integrated feedback derived from actual elicitation interviews. We conducted the study within the schedule of a requirements engineering course, and the participants were recruited from this course. They were graduate students and working professionals in various positions, mostly in the software business. Their prior participation in requirements elicitation interviews and years of work experience are also diverse, as shown in Table 3. We argue that our sample population represents the target user group, which is expected to include anyone at various levels of expertise training for requirements elicitation interviews. Still, the system should be evaluated with a broader group of users with various characteristics, including expertise, profession, as well as demographic factors such as age and nationality.

*Construct validity.* Construct validity may be at risk due to participants' discomfort in feeling monitored and evaluated. This can affect the accuracy and reliability of the study's findings since the participants may not perform optimally under this feeling. To counter this social threat known as evaluation apprehension, prior to the session beginning, the participants were explicitly informed by the researcher that they would be alone during the session (see Section 4.1.1). The researcher was in the experiment environment for wizarding speech-to-text functionality but was kept out of sight to prevent the participants from feeling observed or judged. Our within-subject design can raise another threat to construct validity: the interaction of different treatments. When the effects of one treatment or condition interact with another, it becomes difficult to attribute the observed outcomes to a specific treatment. To address this, we used different scenarios in each condition of a participant to minimize the learning effect across the treatments. The scenarios were developed to be highly similar in terms of length and the mistakes induced but had a different context to avoid any potential learning effect across the treatments. We also examined if the scenario itself induced any confounding factor on the results and did not find any, as mentioned in Section 4.2.1.

*Conclusion validity.* To minimize the subjectivity, we used the well-established technology acceptance model survey to get participant opinions for the research questions interested in the users' perceptions. We also developed quantitative measures like processing speed and learning gain based on the current literature to address the related research questions objectively. To ensure that each person received the same treatment during the experiment, we established a standardized experimental technique followed consistently for all participants. The experimental systems were designed to be autonomous, except for wizarded speech-to-text functionality, thereby effectively mitigating potential researcher bias throughout the experiment. For the speech-to-text functionality, though, there is no room for the researcher to introduce any application bias as participant speech is taken as it is and accepted only if it fits one of the available dialogue options. To further improve the experiment's robustness, the technical requirements were communicated to the participants in the introduction of the experiment to eliminate the influence of any random confounding variables. We aim to mitigate any disturbing effects like background noise or distractions from others present in the experiment environment. However, since the experiment was conducted online, we were unable to exert complete control over the participants' experimental environments. Finally, all statistical tests were chosen based on

the distribution and independence checks of the data to ensure that the test assumptions were not violated. By selecting appropriate tests and confirming the suitability of the data for each test, the statistical analyses were better able to provide reliable and meaningful results for the study. We obtain a statistical power, estimated at $1 - \beta = 0.68$ for the chosen alpha level $\alpha = 0.05$ and an assumed effect size of $d = 0.5$ (as informed in Bartlett et al. (2022)) for the repeated measures. An increase in sample size could improve the statistical power of the results. Nevertheless, our experimental outcomes, acquired through the participation of individuals who closely resemble the target user population of the study, yield authentic and relevant findings (Falessi et al., 2018).

## 7. Conclusions and Future Work

This paper introduces REIT, an extensible and configurable requirements elicitation interview training system architecture to support requirements elicitation interviews training. We implement two instances of REIT: *i.* RoREIT with an embodied physical robot, and *ii.* VoREIT with a voice-based virtual agent. We assessed the two systems' advantages and drawbacks in terms of learners' experience and outcomes by conducting a user study with the students of a graduate-level REET course. Our research constitutes pioneering work in the field, incorporating emerging technologies to enhance the training process of requirements elicitation interviews. We share the implementations of RoREIT and VoREIT in our public repository provided in Görer and Aydemir (2023a). We invite the community to further investigate and improve our work and implement other versions of REIT with different agents having other interaction modalities or capabilities.

In our study, the participants showed higher learning gains in RoREIT than in VoREIT, with a significant difference between the two systems. The results indicate that VoREIT was perceived to be easier to use compared to RoREIT. Both systems are rated appreciatively by the participants (i.e., higher than 3 = moderate level), and we do not detect a significant difference in the attitudes and perceived usefulness between RoREIT and VoREIT. However, the participants found VoREIT to be more engaging, although objective measurements of engagement based on arousal levels of participants' speech did not indicate any significant difference between the two systems. The participants responded to both systems at similar speeds, but their responses for the turns with incorrect selection were slower than the correctly replied turns.

As part of future work, we plan to integrate one of the state-of-art speech recognition libraries (Radford et al., 2022) into REIT to replace its human-operated speech-to-text component. Increasing the number of available scenarios for REIT would also provide additional training experience to its users. Our plan is to build a public scenario library where the RE community can contribute with their own scenarios. Considering the proven effect of adaptive systems in education (Ahmad et al., 2017), our system can also offer a more personalized and effective learning experience by offering scenarios customized to the needs of each learner. The complexity level of scenarios and the intensity of feedback can be adapted to the learner's performance, ensuring that each learner receives personalized content appropriate for their level of expertise. In this way, experienced learners would not get bored with simple training content, and inexperienced learners would not get demotivated by overly challenging material.

## Declarations

*Conflicts of Interest.*  The authors have no conflicts of interest to declare that are relevant to the content of this article.

## CRediT authorship contribution statement

**Binnur Görer:** Conceptualization, Method, Software, Validation, Investigation, Data Curation, Writing. **Fatma Başak Aydemir:** Conceptualization, Method, Writing, Supervision.

# References

Ahmad, M.I., Mubin, O., Orlando, J., 2017. A systematic review of adaptivity in human-robot interaction. Multimodal Technologies and Interaction 1, 14.

Alraimi, K.M., Zo, H., Ciganek, A.P., 2015. Understanding the moocs continuance: The role of openness and reputation. Computers & Education 80, 28–38.

Andrews, B., Hejdenberg, J., Wilding, J., 2006. Student anxiety and depression: comparison of questionnaire and interview assessments. Journal of affective disorders 95, 29–34.

Appleton, J.J., Christenson, S.L., Kim, D., Reschly, A.L., 2006. Measuring cognitive and psychological engagement: Validation of the student engagement instrument. Journal of school psychology 44, 427–445.

Ascalon, M.E., Meyers, L.S., Davis, B.W., Smits, N., 2007. Distractor similarity and item-stem structure: Effects on item difficulty. Applied Measurement in Education 20, 153–170.

Bano, M., Zowghi, D., Ferrari, A., Spoletini, P., Donati, B., 2019. Teaching requirements elicitation interviews: an empirical study of learning from mistakes. Requirements Engineering 24, 259–289.

Barros, P., Churamani, N., Sciutti, A., 2020. The facechannel: a fast and furious deep neural network for facial expression recognition. SN Computer Science 1, 1–10.

Bartlett, M.E., Edmunds, C., Belpaeme, T., Thill, S., 2022. Have i got the power? analysing and reporting statistical power in hri. ACM Transactions on Human-Robot Interaction (THRI) 11, 1–16.

Beaudry, A., Pinsonneault, A., 2010. The other side of acceptance: Studying the direct and indirect effects of emotions on information technology use. MIS quarterly , 689–710.

Belpaeme, T., Kennedy, J., Ramachandran, A., Scassellati, B., Tanaka, F., 2018. Social robots for education: A review. Science robotics 3.

Blaikie, N., 2003. Analyzing quantitative data: From description to explanation. Sage.

Cao, Y., Theune, M., Nijholt, A., 2009. Modality effects on cognitive load and performance in high-load information presentation, in: Proceedings of the 14th international conference on Intelligent user interfaces, pp. 335–344.

Caponetto, I., Earp, J., Ott, M., 2014. Gamification and education: A literature review, in: European Conference on Games Based Learning, Academic Conferences International Limited. p. 50.

Cer, D., Yang, Y., Kong, S.y., Hua, N., Limtiaco, N., John, R.S., Constant, N., Guajardo-Cespedes, M., Yuan, S., Tar, C., et al., 2018. Universal sentence encoder. arXiv preprint arXiv:1803.11175 .

Chang, Y., Wang, X., Wang, J., Wu, Y., Zhu, K., Chen, H., Yang, L., Yi, X., Wang, C., Wang, Y., et al., 2023. A survey on evaluation of large language models. arXiv preprint arXiv:2307.03109 .

Ciolacu, M., Tehrani, A.F., Binder, L., Svasta, P.M., 2018. Education 4.0-artificial intelligence assisted higher education: early recognition system with machine learning to support students' success, in: 2018 IEEE 24th International Symposium for Design and Technology in Electronic Packaging(SIITME), IEEE. pp. 23–30.

Cramér, H., 2016. Mathematical Methods of Statistics (PMS-9), Volume 9. Princeton university press.

Darwish, H., 2014. The "persona effect": Shortcomings in the evaluation of pedagogical agents' embodiment, in: 2014 International Conference on Web and Open Access to Learning (ICWOAL), IEEE. pp. 1–5.

Daun, M., Brings, J., Hübscher, C., 2017. How common are controlled experiments with student participants in requirements engineering?: A systematic mapping study on the use and reporting of graduate and undergraduate students in requirements engineering experiments, in: 2017 IEEE 25th International Requirements Engineering Conference Workshops (REW), IEEE. pp. 307–314.

Daun, M., Grubb, A.M., Tenbergen, B., 2021. A survey of instructional approaches in the requirements engineering education literature, in: 2021 IEEE 29th International Requirements Engineering Conference (RE), IEEE. pp. 257–268.

Davis, A., Dieste, O., Hickey, A., Juristo, N., Moreno, A.M., 2006. Effectiveness of requirements elicitation techniques: Empirical results derived from a systematic review, in: 14th IEEE International Requirements Engineering Conference (RE'06), IEEE. pp. 179–188.

Davis, F.D., 1989. Perceived usefulness, perceived ease of use, and user acceptance of information technology. MIS quarterly , 319–340.

Davis, R.O., Park, T., Vincent, J., 2023. A meta-analytic review on embodied pedagogical agent design and testing formats. Journal of Educational Computing Research 61, 30–67.

Debnath, S., Spoletini, P., 2020. Designing a virtual client for requirements elicitation interviews. Requirements Engineering: Foundation for Software Quality. REFSQ 2020 .

Engwall, O., Lopes, J., Cumbal, R., 2022. Is a wizard-of-oz required for robot-led conversation practice in a second language? International Journal of Social Robotics 14, 1067–1085.

Falessi, D., Juristo, N., Wohlin, C., Turhan, B., Münch, J., Jedlitschka, A., Oivo, M., 2018. Empirical software engineering experts on the use of students and professionals in experiments. Empirical Software Engineering 23, 452–489.

Ferrari, A., Huichapa, T., Spoletini, P., Novielli, N., Fucci, D., Girardi, D., 2021. Using voice and biofeedback to predict user engagement during requirements interviews. arXiv preprint arXiv:2104.02410 .

Ferrari, A., Spoletini, P., Bano, M., Zowghi, D., 2020. Sapeer and reversesapeer: teaching requirements elicitation interviews with role-playing and role reversal. Requirements Engineering 25, 417–438.

Garbers, B., Periyasamy, K., 2006. A light weight tool for teaching the development and evaluation of requirements documents, in: 2006 Annual Conference & Exposition, pp. 11–61.

Garcia, I., Pacheco, C., León, A., Calvo-Manzano, J.A., 2019. Experiences of using a game for improving learning in software requirements elicitation. Computer Applications in Engineering Education 27, 249–265.

García, I., Pacheco, C., León, A., Calvo-Manzano, J.A., 2020. A serious game for teaching the fundamentals of iso/iec/ieee 29148 systems and software engineering–lifecycle processes–requirements engineering at undergraduate level. Computer Standards & Interfaces 67, 103377.

Gouaillier, D., Hugel, V., Blazevic, P., Kilner, C., Monceaux, J., Lafourcade, P., Marnier, B., Serre, J., Maisonnier, B., 2008. The nao humanoid: a combination of performance and affordability. CoRR abs/0807.3223 .

Grivokostopoulou, F., Kovas, K., Perikos, I., 2020. The effectiveness of embodied pedagogical agents and their impact on students learning in virtual worlds. Applied Sciences 10, 1739.

Gururangan, S., Marasović, A., Swayamdipta, S., Lo, K., Beltagy, I., Downey, D., Smith, N.A., 2020. Don't stop pretraining: Adapt language models to domains and tasks. arXiv preprint arXiv:2004.10964 .

Görer, B., Aydemir, F.B., 2023a. Emerging Technologies in Requirements Elicitation Interview Training: Robotic and Virtual Tutors— Experimental Material. URL: https://doi.org/10.5281/zenodo.7861906, doi:10.5281/zenodo.7861906.

Görer, B., Aydemir, F.B., 2023b. Generating requirements elicitation interview scripts with large language models, in: To appear in 2023 IEEE 31st International Requirements Engineering Conference Workshops (REW), IEEE. URL: https://aire-ws.github.io/aire23/papers/AIRE_03.pdf.

Görer, B., Aydemir, F.B., 2023c. Roboreit: an interactive robotic tutor with instructive feedback component for requirements elicitation interview training. Journal of Software: Evolution and Process , e2608doi:https://doi.org/10.1002/smr.2608.

Hadar, I., Soffer, P., Kenzi, K., 2014. The role of domain knowledge in requirements elicitation via interviews: an exploratory study. Requirements Engineering 19, 143–159.

Hainey, T., Connolly, T.M., Stansfield, M., Boyle, E.A., 2011. Evaluation of a game to teach requirements collection and analysis in software engineering at tertiary education level. Computers & Education 56, 21–35.

Hake, R.R., 1998. Interactive-engagement versus traditional methods: A six-thousand-student survey of mechanics test data for introductory physics courses. American journal of Physics 66, 64–74.

Han, J., Jo, M., Park, S., Kim, S., 2005. The educational use of home robots for children, in: ROMAN 2005. IEEE International Workshop on Robot and Human Interactive Communication, 2005., IEEE. pp. 378–383.

Hoffman, G., Zhao, X., 2020. A primer for conducting experiments in human–robot interaction. ACM Transactions on Human-Robot Interaction (THRI) 10, 1–31.

Hong, S.J., Lui, C.S.M., Hahn, J., Moon, J.Y., Kim, T.G., 2013. How old are you really? cognitive age in technology acceptance. Decision Support Systems 56, 122–130.

Ibrahim, Z., Soo, M.C., Soo, M.T., Aris, H., 2019. Design and development of a serious game for the teaching of requirements elicitation and analysis, in: 2019 IEEE International Conference on Engineering, Technology and Education (TALE), IEEE. pp. 1–8.

Kakeshita, T., Yamashita, S., 2015. A requirement management education support tool for requirement elicitation process of rebok, in: 2015 3rd International Conference on Applied Computing and Information Technology/2nd International Conference on Computational Science and Intelligence, IEEE. pp. 40–45.

Kirschner, P.A., 2002. Cognitive load theory: Implications of cognitive load theory on the design of learning.

Konlog, R.I., Spoletini, P., 2023. Reit-builder: Customizable training for requirements elicitation interviews, in: REFSQ Co-Located Events. URL: https://ceur-ws.org/Vol-3378/PT-paper5.pdf.

Kwon, J.H., Powell, J., Chalmers, A., 2013. How level of realism influences anxiety in virtual reality environments for a job interview. International journal of human-computer studies 71, 978–987.

Laiq, M., Dieste, O., 2020. Chatbot-based interview simulator: A feasible approach to train novice requirements engineers, in: 2020 10th International Workshop on Requirements Engineering Education and Training (REET), IEEE. pp. 1–8.

van Lamsweerde, A., 2009. Requirements engineering: From system goals to UML models to software. volume 10. Chichester, UK: John Wiley & Sons.

Leoste, J., Jõgi, L., Õun, T., Pastor, L., San Martín López, J., Grauberg, I., 2021. Perceptions about the future of integrating emerging technologies into higher education—the case of robotics with artificial intelligence. Computers 10, 110.

Leyzberg, D., Spaulding, S., Toneva, M., Scassellati, B., 2012. The physical presence of a robot tutor increases cognitive learning gains, in: Proceedings of the annual meeting of the cognitive science society.

Li, J., 2015. The benefit of being physically present: A survey of experimental works comparing copresent robots, telepresent robots and virtual agents. International Journal of Human-Computer Studies 77, 23–37.

Liang, P., De Graaf, O., 2010. Experiences of using role playing andwiki in requirements engineering course projects, in: 2010 5th International Workshop on Requirements Engineering Education and Training, IEEE. pp. 1–6.

Liaw, S.S., 2008. Investigating students' perceived satisfaction, behavioral intention, and effectiveness of e-learning: A case study of the blackboard system. Computers & education 51, 864–873.

Lu, J., Schmidt, M., Lee, M., Huang, R., 2022. Usability research in educational technology: A state-of-the-art systematic review. Educational technology research and development , 1–42.

McCarthy, J., Goffin, R., 2004. Measuring job interview anxiety: Beyond weak knees and sweaty palms. Personnel Psychology 57, 607–637.

Merchant, Z., Goetz, E.T., Cifuentes, L., Keeney-Kennicutt, W., Davis, T.J., 2014. Effectiveness of virtual reality-based instruction on students' learning outcomes in k-12 and higher education: A meta-analysis. Computers & Education 70, 29–40.

Nachar, N., et al., 2008. The mann-whitney u: A test for assessing whether two independent samples come from the same distribution. Tutorials in quantitative Methods for Psychology 4, 13–20.

Nakamura, T., Kai, U., Tachikawa, Y., 2014. Requirements engineering education using expert system and role-play training, in: 2014 IEEE International Conference on Teaching, Assessment and Learning for Engineering (TALE), IEEE. pp. 375–382.

Ni, J., Young, T., Pandelea, V., Xue, F., Cambria, E., 2023. Recent advances in deep learning based dialogue systems: A systematic survey. Artificial intelligence review 56, 3055–3155.

Ochoa, O., Babbit, A., 2019. Incorporating a virtual reality environment in the teaching of analysis of software requirements, in: 2019 IEEE Frontiers in Education Conference (FIE), IEEE. pp. 1–5.

Ogata, S., Matsuura, S., 2012. Training of requirements analysis modeling with uml-based prototype generation tool, in: Proceedings of the 5th India Software Engineering Conference, pp. 105–108.

Paschoal, L.N., de Oliveira, M.M., Chicon, P.M.M., 2018. A chatbot sensitive to student's context to help on software engineering education, in: 2018 XLIV Latin American Computer Conference (CLEI), IEEE. pp. 839–848.

Pelletier, K., McCormack, M., Reeves, J., Robert, J., Arbino, N., Dickson-Deane, C., Guevara, C., Koster, L., Sanchez-Mendiola, M., Bessette, L.S., et al., 2022. 2022 EDUCAUSE Horizon Report Teaching and Learning Edition. Technical Report. EDUC22.

Powell, D.M., Bourdage, J.S., Bonaccio, S., 2021. Shake and fake: The role of interview anxiety in deceptive impression management. Journal of business and psychology 36, 829–840.

Quigley, M., Conley, K., Gerkey, B., Faust, J., Foote, T., Leibs, J., Wheeler, R., Ng, A.Y., et al., 2009. Ros: an open-source robot operating system, in: ICRA workshop on open source software, Kobe, Japan. p. 5.

Radford, A., Kim, J.W., Xu, T., Brockman, G., McLeavey, C., Sutskever, I., 2022. Robust Speech Recognition via Large-Scale Weak Supervision. Technical Report. Technical report, OpenAI, 2022. URL: https://cdn.openai.com/papers/whisper.pdf.

Radzikowski, K., Nowak, R., Wang, L., Yoshie, O., 2019. Dual supervised learning for non-native speech recognition. EURASIP Journal on Audio, Speech, and Music Processing 2019, 1–10.

Roohr, K.C., Liu, H., Liu, O.L., 2017. Investigating student learning gains in college: A longitudinal study. Studies in Higher Education 42, 2284–2300.

Russell, J.A., 1980. A circumplex model of affect. Journal of personality and social psychology 39, 1161.

Salam, H., Celiktutan, O., Gunes, H., Chetouani, M., 2022. Automatic context-driven inference of engagement in hmi: A survey. arXiv preprint arXiv:2209.15370 .

Salthouse, T.A., 1996. The processing-speed theory of adult age differences in cognition. Psychological review 103, 403.

Schirmer, A., Adolphs, R., 2017. Emotion perception from face, voice, and touch: comparisons and convergence. Trends in cognitive sciences 21, 216–228.

Shin, J., Guo, Q., Gierl, M.J., 2019. Multiple-choice item distractor development using topic modeling approaches. Frontiers in psychology 10, 825.

Sitzmann, T., 2011. A meta-analytic examination of the instructional effectiveness of computer-based simulation games. Personnel psychology 64, 489–528.

Syrdal, D.S., Dautenhahn, K., Koay, K.L., Walters, M.L., 2009. The negative attitudes towards robots scale and reactions to robot behaviour in a live human-robot interaction study. Adaptive and emergent behaviour and complex systems .

Timotheou, S., Miliou, O., Dimitriadis, Y., Sobrino, S.V., Giannoutsou, N., Cachia, R., Mones, A.M., Ioannou, A., 2023. Impacts of digital technologies on education and factors influencing schools' digital capacity and transformation: A literature review. Education and information technologies 28, 6695–6726.

Trowler, V., 2010. Student engagement literature review. The higher education academy 11, 1–15.

Vega, K., Fuks, H., Carvalho, G., 2009. Training in requirements by collaboration: Branching stories in second life, in: 2009 Simposio Brasileiro de Sistemas Colaborativos, IEEE. pp. 116–122.

Vilar, E., Noriega, P., Borges, T., Rebelo, F., Ramos, S., 2020. Can an environmental feature influence interview anxiety?, in: International Conference on Human-Computer Interaction, Springer. pp. 351–369.

Wagner, J., Triantafyllopoulos, A., Wierstorf, H., Schmitt, M., Eyben, F., Schuller, B.W., 2022. Dawn of the transformer era in speech emotion recognition: closing the valence gap. arXiv preprint arXiv:2203.07378 .

de Winter, J.F., Dodou, D., 2010. Five-point likert items: t test versus mann-whitney-wilcoxon (addendum added october 2012). Practical Assessment, Research, and Evaluation 15, 11.

Wohlin, C., Runeson, P., Höst, M., Ohlsson, M.C., Regnell, B., Wesslén, A., 2012. Experimentation in software engineering. Springer Science & Business Media.

Woolson, R.F., 2007. Wilcoxon signed-rank test. Wiley encyclopedia of clinical trials , 1–3.

Yang, H.d., Yoo, Y., 2004. It's all about attitude: revisiting the technology acceptance model. Decision support systems 38, 19–31.

Yasin, A., Liu, L., Li, T., Wang, J., Zowghi, D., 2018. Design and preliminary evaluation of a cyber security requirements education game (sreg). Information and Software Technology 95, 179–200.

Zhao, Y., Watterston, J., 2021. The changes we need: Education post covid-19. Journal of Educational Change 22, 3–12.

Zowghi, D., Coulin, C., 2005. Requirements elicitation: A survey of techniques, approaches, and tools, in: Engineering and managing software requirements. Springer, pp. 19–46.