

# Designing PairBuddy—A Conversational Agent for Pair Programming

PETER ROBE and SANDEEP KAUR KUTTAL, University of Tulsa

From automated customer support to virtual assistants, conversational agents have transformed everyday interactions, yet despite phenomenal progress, no agent exists for programming tasks. To understand the design space of such an agent, we prototyped PairBuddy—an interactive pair programming partner—based on research from conversational agents, software engineering, education, human-robot interactions, psychology, and artificial intelligence. We iterated PairBuddy’s design using a series of Wizard-of-Oz studies. Our pilot study of six programmers showed promising results and provided insights toward PairBuddy’s interface design. Our second study of 14 programmers was positively praised across all skill levels. PairBuddy’s active application of soft skills—adaptability, motivation, and social presence—as a navigator increased participants’ confidence and trust, while its technical skills—code contributions, just-in-time feedback, and creativity support—as a driver helped participants realize their own solutions. PairBuddy takes the first step towards an Alexa-like programming partner.

CCS Concepts: • **Human-Centered Computing** → **User studies**;

Additional Key Words and Phrases: Conversational agents, pair programming, user centered design, Wizard of Oz

## ACM Reference format:

Peter Robe and Sandeep Kaur Kuttal. 2022. Designing PairBuddy—A Conversational Agent for Pair Programming. *ACM Trans. Comput.-Hum. Interact.* 29, 4, Article 34 (May 2022), 44 pages.  
<https://doi.org/10.1145/3498326>

## 1 INTRODUCTION

Conversational agents allow humans to use natural language to directly interface with computer agents such as virtual assistants (e.g., Apple’s Siri [244], Google Assistant [245], and Amazon’s Alexa [243]), customer support agents, or individual/social chatbots (e.g., Mitsuku [222], Cleverbot [221], and XiaoIce [224]). Conversational agents mimic human conversations, establish more personal connections, and can even increase accessibility for people with physical disabilities or language barriers. For businesses, conversational agents personalize customer experience while

This material is based upon work supported by the National Science Foundation (CAREER) under award number 2046205. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the view of the NSF. We would like to thank David Magar, Jarow Myers, Sam Gurka, Katherine Kwasny, and Bali Ong for help with the studies; David Piorkowski and Rachel Bellamy for their feedback; and Courtney Spivey for editing.

Authors’ address: P. Robe and S. K. Kuttal, University of Tulsa, Tulsa, 74104 Oklahoma; emails: {pjr144, sandeep-kuttal}@utulsa.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2022 Association for Computing Machinery.

1073-0516/2022/05-ART34 \$15.00

<https://doi.org/10.1145/3498326>

bringing down operational costs. Today, a full two-thirds of the most popular websites use conversational agents to interact with users and address their needs [220]. Despite the phenomenal penetration of conversational agents into domains ranging from business to personal use and entertainment, no conversational agents exist for computer programming tasks.

Extensive research aims to increase programming efficiency for both novice and professional programmers via automated tools, including debugging support [7, 8], intelligent **interactive development environments (IDEs)**, and test case generation [74, 196]. In a parallel research field, **Intelligent Tutoring Systems (ITS)** aim to simulate a human teacher by establishing a tutor-tutee relationship between novice programmers (tutee) and computer agents (tutor) [56, 75, 239]. Within the domain of ITS systems, peer-learning agents directly collaborate with learners, often by frequently switching roles between tutor and tutee [47, 99, 195, 216, 237, 264]. Inspired from these two research domains, we aim to increase programmers' efficiency and facilitate learning; however, we further leverage the anthropomorphic interactions of conversational agents to directly collaborate on programming tasks, facilitate learning-by-doing, engage and motivate programmers, and increase self-efficacy.

In this article, we take the first step toward the creation of "PairBuddy"—an anthropomorphic conversational agent and interactive programming partner. PairBuddy serves to simulate a human during "pair programming," a quintessential and well-established collaboration technique used in education and the industry.

In pair programming, two programmers work collaboratively on the same design, algorithm, code, or test [185, 259, 260]. Programmers switch between the roles of driver (writing code) and navigator (making suggestions). Pair programming provides a variety of benefits, including increased code quality, productivity, creativity, knowledge management, and self-efficacy [27, 44, 52, 61, 66, 121, 159, 160, 185, 201, 206, 255, 259–261, 277]. It even has the potential to reduce gender prejudice by encouraging women to pursue computer science [255]. Pair programming increases programmers' contemporary skills, understanding of fundamental concepts, and intellectual pursuits [44, 159, 201, 206, 255]. However, pair programming has certain limitations including scheduling difficulties, collocating pairs, student resistance to pairing, and the dependency on a partner's programming abilities [84, 101, 181, 258].

### 1.1 Motivational Scenario

To motivate the design of PairBuddy and inform its abilities, we provide the following scenario:

Tiana is a junior CS student who enjoys pair programming, but due to COVID-19, she has to return home. Given the physical and time-zone differences, scheduling pair programming sessions is difficult. Additionally, all of her potential partners either like to work solo, have incompatible cognitive styles, criticize without providing solutions (e.g., "*We definitely cannot finish this*"), are overly competitive, take credit for Tiana's work, or have other problematic habits. As a result, their sessions with Tiana are stressful and unproductive, resulting in the loss of pair programming's benefits. Fortunately, Tiana can use PairBuddy (Figure 1).

Tiana can collaborate with PairBuddy, switching between the roles of driver and navigator to create test cases, write code, and refactor code. Additionally, PairBuddy: (1) is unbiased toward Tiana's gender, ethnicity, and socio-economic status; (2) motivates and encourages her hard work; (3) gives non-judgemental feedback, allowing Tiana to disagree with or ignore PairBuddy without the risk of hurt feelings; (4) provides a sense of security through its active listening and engagement through both voice and text; (5) encourages and instills healthy problem-solving and creative styles; and (6) avoids dominating power dynamics by balancing pair programming roles. For these reasons and more, Tiana prefers PairBuddy as her programming partner to complete difficult assignments comfortably and effectively.



Fig. 1. Tiana interacting with PairBuddy.

## 1.2 Feasibility of PairBuddy

We have already established the feasibility of PairBuddy in our previous research, “Trade-offs for Substituting a Human with an Agent in a Pair Programming Context: The Good, the Bad, and the Ugly” [208]. Using the final iteration of PairBuddy (discussed in Section 5.1), we investigated the trade-offs of substituting an agent in a pair programming context, specifically analyzing what aspects of pair programming were gained, lost, and transformed. We compared and contrasted transcripts, code quality, productivity, and survey results of participants’ self-efficacy and pair programming experience between nine human-human pairs and fourteen human-PairBuddy pairs (using a Wizard of Oz simulation). The results were triangulated with interviews. Upon comparing human-human and human-PairBuddy, we found: (1) no significant differences in code quality, productivity, and self-efficacy before and after the programming task (Table 1); (2) PairBuddy was unable to explain the logic behind its code or participate in idea discussion, while human partners could easily communicate together conceptually; (3) PairBuddy’s instructions were often trusted without question—in contrast with human partners—in one case, a participant responded, “*I’m going to blindly believe you;*” (4) human partners interrupted PairBuddy when they were stuck, were unsure about their problems or wanted clarification, while human-human pairs were more hesitant and asked follow-up questions to understand, clarify, or verify their partners’ decisions; (5) human partners showed the same humility towards PairBuddy as they would with other humans by attributing success to the group through words like “we” and taking personal responsibility for mistakes using “I” or “me;” and (6) humans freely accepted or dismissed PairBuddy’s ideas (as in [135]), yet were more hesitant to dismiss ideas from another human.

To supplement our previous research [208], the present article describes our iterative, user-centered design methodology used to build PairBuddy. We established PairBuddy’s anthropomorphic design space through the integration of novel concepts from an extensive literature review of various domains such as human-computer interaction, software engineering, artificial intelligence, psychology, and education. PairBuddy’s focused design evolved through a series of empirical Wizard of Oz studies to realize a robust pair programming conversational agent. Our methodological approach and design space can be utilized by researchers and practitioners to advance programmers’ interactions in other domains such as ITS and interactive educational platforms.

## 1.3 PairBuddys’ Design Challenges

Since no pair programming conversational agents exist, the challenges of creating one include:

- (1) **An Unknown Design – Interface and Interactions:** PairBuddy’s interface and interaction design must be explored due to the unique properties of pair programming dialogue

Table 1. Performance Metrics from [208] for Human-Human and Human-PairBuddy Studies

	Human-Human Studies		P-Value ( $\alpha < 0.5$ )	Human-PairBuddy Studies	
	Samples	Average		Samples	Average
Code Quality	9	89.33	0.9144	14	98.36
Productivity	9	50.44	0.2191	14	66.00
Self-Efficacy	18	3.83	0.1093	14	3.36

[200]. To support programmer-computer interaction, PairBuddy must imitate an effective programmer by integrating a multitude of technical and soft skills, including diverse problem-solving and creative strategies. Therefore, PairBuddy's design must be informed from multi-disciplinary research.

- (2) **A Specific Domain – Software Development:** PairBuddy needs to be created for the specific domain of software development, which remains relatively unexplored by conversational agent research. Developing software is unique, as it demands the synthesis of requirements; the generation and development of solutions; and the implementation, testing, and refactoring of code. Therefore, pre-existing conversational agents cannot be directly modified for programming, and instead, a new paradigm must emerge to support programmer collaboration with an agent. This shift in approach must be informed from software engineering and education. Specifically, we consider literature on ITS to integrate the programming and educational aspects of pair programming.
- (3) **A Specific Userbase – Programmers:** Programmers are a unique population that have yet to be studied in the realm of conversational agents. From students to professionals, programmers' experiences are diverse; some may use agents to learn programming concepts, while others may develop those same agents. To design PairBuddy for both populations, we must conduct user studies to fully understand the breadth of expectations, preferences, and needs of programmers.

These challenges necessitate a review of existing research on human-computer interactions, human-robotic interactions, software engineering, conversational agents, artificial intelligence, ITS, education, psychology, cognitive science, and management science to inform the design of PairBuddy.

## 1.4 Outline

The remainder of this article is organized as follows: Section 2 compares our research to related work on automated support in programming environments, ITS, and mixed-initiative interfaces; gaps in existing research motivate the necessity of designing a pair programming conversational agent. Section 3 provides background on the Wizard of Oz method. Section 4 details the first iteration of PairBuddy's design space, evaluates the design in a pilot study, and reports lessons learned. Section 5 describes the second iteration of PairBuddy's design, informed from the pilot study, and evaluated in a larger main study. Section 6 discusses the design space and associated challenges for future work. Section 7 summarizes the contributions of our work.

## 2 RELATED RESEARCH

While no research specifically addresses programming conversational agents, a variety of research exists for automated programming support; notably, automated debugging tools embedded in IDEs aid software development tasks and ITS facilitate learning with novice programmers. Furthermore, we contextualize PairBuddy within research on mixed-initiative interfaces.

## 2.1 Automated Support in Programming Environments

Automated debugging support is the most commonly integrated form of programming help within existing IDEs. Additionally, research has explored the integration of intelligent agents directly within IDEs to provide advanced feedback.

**2.1.1 Automated Debugging.** Programmers often leverage tools present in many existing IDEs to detect simple syntax errors in real-time, automate repetitive tasks, and debug code [7, 8]. These tools are popular across almost any programming language as they help minimize simple mistakes and let programmers focus on the bigger picture. Past research has even identified techniques that automatically test a range of function parameters to give automated feedback on code correctness [257].

**2.1.2 Intelligent IDE Assistants.** Intelligent IDE assistants provide feedback to programmers via interfaces embedded within IDEs. For example, PETCHA helps teachers author assignments with test cases and feedback to provide students helpful hints when their code fails [192]. For parallel programming, PAPA utilizes IBM Watson services to provide a dialogue-based chat window integrated within an existing IDE [165]. Specifically, PAPA helps answer common parallel programming questions identified in research and will suggest relevant research papers for other queries.

Rather than replace existing automated tools, we can utilize their functionality to support PairBuddy's interactions with programmers. We aim to advance programming help by leveraging a new form of interaction—conversational agents—which in turn, can help elevate the under-utilization of refactoring tools [171]. PairBuddy's main focus is not the generation of code or feedback itself, but rather, its unique ability to interface this information in a way that promotes the productivity, code quality, and self-efficacy of programmers; particularly by leveraging humans' unique creative abilities.

## 2.2 Intelligent Tutoring Systems

ITS provide immediate and personalized feedback to individual learners in the absence of a human teacher [56, 75, 110, 123, 186, 239]. These computer-assisted learning systems model the cognitive and emotional state of learners to adapt and individualize instruction [42]. ITSs are known to increase student performance well beyond conventional classes [131, 153, 174, 240], and in some cases, outperform human tutoring [230, 231, 239]. They have been developed to teach programming (e.g., C++, Java, Lisp, Prolog), database systems, data structures, SQL queries, and other introductory computer science topics. ITSs teach via dynamic lesson plans—simple questions including multiple choice, true/false or short answer—or worked solutions from instructional resources [56, 140].

**2.2.1 ITS Systems with Dialogue-Based Interactions.** AI-powered tutoring approaches that support dialogue-based interactions function similarly to PairBuddy [141, 274], including examples such as ProPL and HabiPro. ProPL [138] is an agent that interacts with students via natural language to develop pseudocode solutions for a given problem. ProPL follows a hierarchical dialogue tree to guide participants toward the solution by asking thought-provoking questions and resolving students' misconceptions. The agent's dialogue is curated by a human expert (i.e., questions and expected answers) and tailored to a specific programming task. HabiPro [248] simulates a student co-learner in a collaborative learning environment and guides group members to answer programming questions (e.g., identify code mistakes). HabiPro models student behavior and guides sessions via natural language to encourage passive students, avoid off-topic conversation, give ideas, provide examples, and generally promote effective learning.

**2.2.2 Peer Learning Agents.** Rather than replace human tutors, peer-learning agents act as equals in collaborative tasks, often by frequently switching between roles of tutor and tutee [47, 99, 195, 216, 237, 264]. Han et al. [99] correlates the tutor-tutee relationship to that of the driver and navigator roles in pair programming; however, in their work, roles were limited to generating flowcharts as the navigator and implementing them as the driver.

Our research is motivated by ITSs regarding promoting student learning, but we are very different as (1) we use a peer-learning approach rather than strictly defined tutor/tutee roles, (2) PairBuddy interacts with programmers as an anthropomorphic conversational agent, (3) interactions are more comparable to human-human pair programming (e.g., focus on building rapport), (4) we aim to increase the code quality, productivity, and self-efficacy of programmers by utilizing pair programming, and (5) our design targets both novice (student) and professional programmers.

### 2.3 Mixed Initiative Interfaces

Mixed-initiative interfaces are broadly defined as an interface where a computer agent shares task load with a user [63, 179]. Mixed-initiative interfaces help foster creativity when producing artifacts [268] and have been used in video game level design [219], robotics [13, 102], visual analytics [154], and GUI-based task learning [145]. As a mixed-initiative interface, PairBuddy aims to supplement, not replace programmers' creativity by sharing task load and negotiating pair programming roles.

## 3 WIZARD OF OZ METHOD

Due to the limited available research on pair programming conversational agents, we used an iterative approach to understand and design PairBuddy through a series of two Wizard of Oz studies with programmers.

In a Wizard of Oz study, participants interact with an agent whose actions are secretly controlled by a human "wizard." Wizard of Oz is a rapid-prototyping method that examines interfaces that are technically demanding or are yet to be created [93]. It helps develop user-friendly interfaces that promote natural language dialogue, consider the unique qualities of human-agent interaction as distinct from normal human discourse [59], and study user interactions with conversational agents [29, 32, 252]. Wizard of Oz efficiently creates the functionality of a product before it is refined via testing [137, 262]. In the words of human-computer interaction expert Jef Raskin, "*Once the product's task is known, design the interface first; then implement to the interface design.*" Accordingly, we used the Wizard of Oz paradigm to simulate PairBuddy's interface and interactions as we explored the design space of pair programming conversational agents. For the remainder of this article, we often use the term "PairBuddy" to refer to "the wizard's simulation of PairBuddy."

Choosing Wizard of Oz studies helps us to investigate the design space of PairBuddy before starting the implementation, as creating a fully functional conversational agent requires: (1) understanding the domain of pair programming conversational agents, (2) collecting pair programming data used to train an agent's underlying machine learning models, and (3) designing and implementing a custom conversational agent architecture including the application of multiple machine-learning algorithms [85]. Additionally, platforms for developing conversational agents such as IBM Watson Assistant [246], SAP Conversational AI [223], and Oracle Digital Assistant [247] cannot be used due to their focus on simple/basic enterprise problems such as answering questions and solving tasks based on web or enterprise data.

Our user-centered approach involved the exploration and evaluation of PairBuddy's design using two iterations of Wizard of Oz studies with programmers. For the first iteration, we conducted an exploratory pilot study with six university students focused on the initial design of PairBuddy's interface and interactions. For the second iteration, a larger study of fourteen university students



Table 2. Design Decisions Evaluated in the Pilot Study

ID	Design Decision	Description	Example Sources
Interface			
F1	Avatar	Embodied by a dynamic 3D avatar	[25, 28, 43, 151, 229, 270]
F2	Gender	Gender can be toggled	[33, 50, 236]
F3	Voice	Communicate via voice synthesis	[28, 41, 242]
F4	Text Chat	Communicate via shared text chat	
Interaction			
I1	Indirect Driving	Send code via text chat	[10, 19, 263]
I2	Adapted Skill	Balance contributions via frequency	[19, 56, 75, 239, 248]
Soft Skills			
S1	Greeting	Introduce itself	[120]
S2	Motivation	Encourage, recognize, comfort, commend	[18, 45, 64, 78]
Technical Skills			
T1	Write/Feedback Tests	Generate test cases & feedback	[17, 74, 161, 163, 168, 196]
T2	Write Code	Examples from online repositories	[124, 125, 177, 193]
T3	Guidance	Provide direction via user stories	

and professional programmers evaluated a variety of design decisions informed from the pilot study and existing research.

#### 4 PILOT STUDY (ITERATION 1)

The pilot study provided valuable feedback from programmers for the initial design of PairBuddy. The insights gathered from the pilot study served as a first step toward the exploration of the design space of pair programming conversational agents as well as a starting point for future iterations.

##### 4.1 PairBuddy Design (Iteration 1)

PairBuddy's initial iteration primarily served to study its interface and interaction design, so we focused our research review on relevant literature including human-robotic interactions, embodiment, software engineering, management science, ITS, education, cognitive science, psychology, and gender-bias in agent design. Table 2 lists the design decisions made for PairBuddy in the pilot study, and are as follows:

##### Design: Creating the Interface and Interactions

PairBuddy design included anthropomorphic characteristics in its interface and interactions.

**(A) Interface – Embodiment via Avatar, Gender, Voice, and Text:** PairBuddy communicated with participants by means of avatar, voice, and text to enhance human-computer interaction. The avatar was incorporated since avatars make interfaces more human [229] and improve understanding, engagement, and trust in novice programmers [151, 229, 270]. Furthermore, embodiment through avatars can facilitate non-verbal communication in order [25, 43] to help maintain effective pair programming relationships. Agents with avatars are given more personality attributes than those without them [229], but at the cost of heightened expectations [151]. They are particularly important to establishing first impressions [28]. Therefore, we formulated design decision F1: PairBuddy will be embodied by a 3D avatar.

The inevitable gendered attributes spawned from conversational agent embodiment are the target of similar gender-biases present in the real world. Particularly, female conversational agents are more often the target of negative stereotypes, sexual attention, and profanities than male agents

[33]. However, they are also more likely to be forgiven, even if satisfaction ultimately does not change [236]. Additionally, female programmers have voiced concerns about pair programming with male partners [50, 133] due to the expectation of being stereotyped for their gender. To analyze potential gender preferences and bias, we formulated design decision F2: Participants can toggle the agent between the two most common genders. While research has looked at gender bias within conversational agents and pair programming separately, we combine the analysis by evaluating participants' gender preferences for a pair programming conversational agent.

The choice between voice and text for an agent's interface depends on the use case. While text-only interfaces are less intrusive (e.g., website support), research [28, 41] has shown that voice interfaces increase users' trust in conversational agents. Additionally, cognitive science research has demonstrated that simultaneously sharing the same sensory modality for both short-term memory and active-use negatively affects response and accuracy rates under high-load conditions [242]. These findings suggest that using audio responses can help reduce programmers' cognitive load since it does not share the same sensory modality as reading code. Additionally, the absence of a text chat saves additional screen space and reduces context switching between applications. For our design decisions F3 and F4, we explored input modality by supporting both methods of communication, allowing messages to arrive to/from the agent via both voice and text.

**(B) Interactions – Indirect Driving, Timed Feedback, and Adapted Skill:** In effective pair programming, both partners directly contribute to the code as drivers. However, if PairBuddy were to make a mistake and overwrite code, it would be difficult for participants to undo PairBuddy's actions. Since a prominent principle of HCI is that users should remain in control of their work [10, 19, 263], we formulated design decision I1: As a driver, PairBuddy will make indirect contributions through text messages so that participants can reference and modify the code themselves. In this way, the distinction between pair programming roles is less pronounced; control over the roles is not negotiated, rather, PairBuddy's contributions simply encompass both responsibilities: providing code examples and giving feedback.

ITSs such as HabiPro [248] have shown the ability to guide students' behavior by utilizing a learner model [46] to represent and track students' knowledge and progress. In the same way that ITSs individualize learning, we formulated a design decision I2: PairBuddy will adapt its skill level to match each participant such that contributions remain balanced. Pair programming research further informs this decision since programmers prefer when their partner is equally or more competent [164]. In the pilot study, this design decision manifested itself in the frequency of PairBuddy's feedback and code suggestions (refer T2). With this decision, we aim to tailor PairBuddy's interactions with programmers based on their specific needs and abilities.

### **Programmer: Integrating Technical and Soft Skills**

PairBuddy will imitate the characteristics of a competent programmer. A programmer's technical and soft skills are crucial for being effective team members [15, 58, 178, 275] and are considered by managers when they make hiring decisions [150, 157, 169, 209, 218]. Therefore, we integrated both types of skills to mirror a programmer's capabilities.

**(A) Soft Skills – Greeting and Motivation:** While greetings vary between cultures, humans introduce themselves to make their presence known and start conversations. Similarly for agents, Kahn et al. [120] identifies "The Initial Introduction" design pattern where agents use scripted, conventional introductions to recognize and inquire about another. It is an important design choice in human-robotic interaction, as it allows a deepening of relationships while removing initial awkwardness. Based on this researched design pattern, we formulated design decision S1: PairBuddy will introduce itself and greet the participant [19].

Motivation is seen as a driving force that has a substantial impact on a programmer's performance and productivity [18, 45], and comes from either intrinsic or extrinsic sources [45, 64].



Designing extrinsic motivators that synergize with intrinsic motivators requires supporting a person's sense of competence without undermining their self-determination [18]. Fischer et al. [78] found that a higher perceived probability of receiving extrinsic motivation in the form of relational rewards (e.g., praise, recognition, performance feedback) [21] often positively affected creative and innovative outcomes. Therefore, we formulated design decision S2: PairBuddy will motivate using relational rewards in the form of encouragement (e.g., *"We've got this!"*), recognition (e.g., *"I see, good idea!"*), and comforting (e.g., *"That's okay, everyone makes mistakes"*). Additionally, extrinsic motivation is most effective during the stages of the creative process that make the most meaningful contributions to the project [18]. Therefore, PairBuddy will commend success through praise and celebration: *"I knew we could do it!"* or *"We make a great team!"*.

**(B) Technical Skills – Writing Tests/Code and Giving Guidance:** Automated approaches to generating test cases (i.e., a code fragment that specifies inputs and expected results to verify compliance with a requirement [1]) have explored the use of search-based techniques and requirement artifacts. Search-based software testing uses a variety of search algorithms [161, 163, 168] to determine the most efficient path through source code that maximizes code coverage for automatic test case generation [17]. Alternatively, research has demonstrated the feasibility of converting requirement artifacts such as user stories (i.e., a description of a requirement from a user's perspective [2]), acceptance criteria (i.e., the boundaries of a user story [1]), and scenarios (i.e., step-by-step description of a series of events [1]) into test cases [74, 196]. Based on this research, we formulated design decision T1: As a driver, PairBuddy can generate test cases automatically. As a navigator, PairBuddy can give feedback and answer programmers' queries based on the generated solutions. For example, PairBuddy could offer help by asking, *"Would you like me to generate a test case?"*

To further PairBuddy's competency as a driver, we formulated a design decision T2: PairBuddy will provide example code from online repositories (e.g., GitHub [3]), question and answer forums (e.g., Stack Overflow [6]), and package documentation. For example, PairBuddy might send sample code from GitHub into the text chat and ask, *"Is this code example from online useful?"* Past research identifies code-querying algorithms used to search online repositories for semantically similar code [124, 125, 177, 193]. However, these algorithms are not always perfect, so to simulate a realistic implementation, PairBuddy's code recommendations were not an exact match with the task.

To allow PairBuddy to guide participants, we formulated a design decision T3: PairBuddy will use user stories as a basis to track and direct the current objective, encouraging participants to reference them when determining how to proceed. Additionally, PairBuddy might need to ask the participant, *"What user story is that?"* when the current objective is unclear from PairBuddy's perspective. This clarifying dialogue serves a second purpose: participants are forced to refocus on the user stories, and reorient themselves toward the task's goal.

## 4.2 Wizard/PairBuddy Implementation

The wizard's interface and implementation integrated the aforementioned design decisions into the wizard's script and interface for the pilot study, and are detailed as follows:

PairBuddy's interface was implemented within the Eclipse IDE [70] using a modified version of the Saros plugin [72], allowing both voice and text communication. Directly integrating these features allowed participants to seamlessly interact with PairBuddy. The wizard used Saros to monitor the participant's code, but made code contributions via the text chat. PairBuddy was embodied by a 3D avatar via FACSvatar [76]. Using custom networking code, FACSvatar mapped the wizard's facial animations onto PairBuddy's avatar from a remote location, and PairBuddy's voice was generated using Google Text-to-Speech [98]. While participants were free to change the arrangement

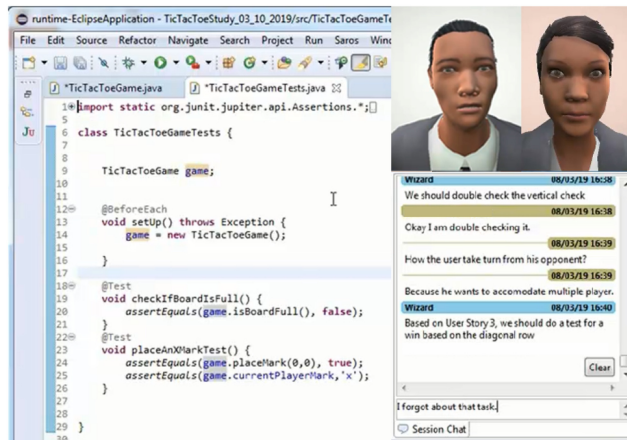


Fig. 2. The pilot study's interface including the Eclipse IDE, text chat, and avatar. Note that each participant only saw one avatar.

of the IDE and avatar windows, none of them did. Figure 2 shows a screenshot of the participant's interface, including the Eclipse IDE and either the man or woman avatar. While participants were informed that they were being recorded, they did not know that the wizard was monitoring their webcam and microphone. As required by IRB, we disclosed the nature of the deception study after the study sessions were completed.

The wizard also simulated the back end design of PairBuddy, including the aforementioned design decisions and common components of conversational agent architecture [85] such as natural language understanding [12, 69, 148, 158, 194, 267], dialogue state tracking [126, 139, 199, 276], dialogue policy [109, 235, 265, 272], and natural language generation [100, 198].

If the participant said, *"This is why I hate coding!"*, a real conversational agent would recognize their words (text-to-speech), classify their intent as "Negative Feedback" (natural language understanding), and log their difficulties (dialogue state tracking). In response, the agent would decide to "Give Motivation" (dialogue policy) and select an appropriate response *"Don't give up, you are almost there!"* (natural language generation). Without an underlying system to power interactions, the wizard adhered to a script of limited dialogue options to simulate a robotic conversational style. Dialogue options were templated according to Shneiderman's guidelines [217] and Neilsen's heuristics [176]. Since communication styles differ by gender [133], we provided gender-inclusive language [77, 167] in our script including non-authoritative suggestions [214] to both engage and motivate programmers.

The wizard exhibited a realistic level of intelligence. If participants asked questions beyond the protocol, the wizard would answer, *"I'm afraid logic isn't my strong suit,"* to simulate the limitations of a potentially automated system [19]. However, to maintain participants' trust and engagement, we designed PairBuddy to give alternative contributions if it could not directly answer queries. Furthermore, the script was designed to vary the responses for the same intent. For example, motivational scripts included, *"We've got this!"* along with five or more similar dialogue templates so participants would not receive repeated phrases from PairBuddy. In general, the script served to simulate the back end of a fully functional conversational agent, including its limitations, through the use of dialogue templates. One researcher simulated PairBuddy as the wizard in the pilot study. To maintain consistency between studies, the wizard was trained in three trial studies prior to any official studies.

Table 3. Demographics of the Pilot Study Participants

P#	Age	Gender	Education	Programming Experience
PG1	24–29	Man	PhD	4+ years
PG2	30–40	Man	PhD	4+ years
PG3	24–29	Man	Masters	4+ years
PU4	19–23	Man	Undergrad	<1 year
PU5	19–23	Woman	Undergrad	<1 year
PU6	19	Man	Undergrad	<1 year

### 4.3 Participants

Seven students majoring in computer science were recruited from our university. However, we only report data from 6 participants (3 undergraduate and 3 graduate) since one participant did not interact with PairBuddy during their study session. Both undergraduate and graduate students were recruited for the study to account for the varying diversity of programming skill. Participants with basic Java programming experience were chosen on a first-come-first-serve basis from a recruitment email. Upon completion, participants were given \$20 in Amazon gift cards. Table 3 shows the demographics of our pilot study participants. These participants are referred to as PU# and PG# for undergraduate and graduate students, respectively. For example, PU4 is the fourth undergraduate participant of the pilot study.

### 4.4 Study Design

Participants completed a consent form and background questionnaire prior to the pilot study. Before starting the task, participants watched video tutorials on the concepts used in the study including the driver and navigator pair programming roles, the think-aloud method, and test-driven development. The think-aloud method encourages participants to vocalize their thoughts and feelings [143, 212]. Test-driven development is a type of extreme programming that prioritizes the creation of test cases before implementing and refactoring code. Test-driven development evaluates participants' knowledge and evokes diverse dialogue since each development stage requires a unique style of thinking.

Participants were asked to use pair programming and test-driven development alongside PairBuddy to complete an implementation of tic-tac-toe: a game where two players take turns marking spaces in a 3x3 board. Based on a list of user stories, acceptance criteria and scenarios, participants were instructed to implement test cases and functionality for turn taking, vertical/horizontal/diagonal wins, a full board test, and a tie game test. As a starting point, code for the board along with three example test cases were provided to participants. Tic-tac-toe was selected for its simplicity since anyone with basic programming experience can understand and implement its requirements without prior knowledge of the domain. Participants wrote Java code in the Eclipse IDE [70] and implemented testing via JUnit [71]. The duration of the task was fixed to 50 minutes to prevent participant fatigue and to ensure that individual study sessions lasted under 90 minutes.

We applied a between-study design for initializing the avatar's gender by starting half of all studies as a man and half as a woman.

A semi-structured interview was conducted using a script, yet individualized questions were used to further explore study-specific events. To avoid social desirability bias [94, 173], we told participants that the study was being used to improve PairBuddy's design and that honest feedback was crucial to evaluate our system.

Table 4. Code Set for PairBuddy's Contributions

Contribution Type	Description	Example
Direction	Guiding participants towards goals	<i>"Write a method... for a horizontal win."</i>
Domain - Help	IDE, language, or domain knowledge	<i>"Looks like... an import error for JUnit."</i>
Method - Add	Provide example code for methods	<i>"Is this code... useful?"</i>
Method - Clarify	Give knowledge about methods	<i>"Get the value... by accessing its coordinate."</i>
Test case - Add	Provide example code for test cases	<i>"I can generate test cases."</i>
Test case - Clarify	Give knowledge about test cases	<i>"To test... we will need to place marks."</i>
Bug - Identify	Identify bugs or mistakes	<i>"Double check the vertical check."</i>
Bug - Fix	Fix bugs or mistakes	<i>"Error is commonly caused by..."</i>
Contribution Source	Description	
PairBuddy Alone	PairBuddy offered contributions	
Human Asked	Participants prompted PairBuddy	
Human Asked - Unanswered	PairBuddy couldn't help programmer	

#### 4.5 Data Analysis

The video, audio, and interviews from our pilot study were transcribed and subsequently analyzed through both quantitative and qualitative measures to evaluate the usability of PairBuddy. We used the Corbin and Strauss variant [232] of Grounded Theory [90] to analyze our transcripts. Specifically, we created a codeset of the types of contributions PairBuddy made by creating, gathering, modifying, and removing contribution types during an iterative, open-coding process [23, 212]. These contribution types allowed us to track participant interaction with PairBuddy over time (Table 3). Using thematic analysis [35], we grouped interview responses into themes that encompass various dimensions of interaction with PairBuddy. For our quantitative analysis, we collected responses to self-efficacy questionnaires and discrete answers to our interview questions. Transcripts were coded by two researchers. Initially, both researchers independently coded the same 20% of the transcripts, and average inter-rater reliability was measured at 86% using the Jaccard index. The remaining transcripts were split and coded separately. Interview transcripts from the pilot study can be found [5].

#### 4.6 Results

The first iteration of PairBuddy allowed us to examine the feasibility of a pair programming conversational agent. The insights and avenues for improvement are as follows:

**(1) Helping Programmers:** PairBuddy helped participants complete the programming task as both a driver and a navigator. Figure 3 details the contributions PairBuddy made in three separate ways: (1) when PairBuddy offered help, (2) when participants asked for help, and (3) when participants asked, but PairBuddy could not provide help. The contributions types include direction, domain-related help, method/test clarification, method/test addition, and bug identification/fixing.

When participants were lost, PairBuddy provided direction through its messages. For example, when PU6 had difficulty writing a method, PairBuddy gave guidance, *"We need to write a method to check for a horizontal win."* Furthermore, PairBuddy provided help for questions about the IDE, JUnit, language, and domain (tic-tac-toe). For example, when PG3 was having difficulties using the testing suite (JUnit), PairBuddy provided help, *"It looks like we may have an import error for JUnit."*

PairBuddy contributed knowledge by clarifying test cases and methods. For example, when PU4 was unsure about how to write a test case, PairBuddy responded, *"To test for horizontal, we will need to place marks at zero zero, one zero, and two zero."* Similarly, PairBuddy helped PU5 with her method, *"If a space is occupied by a '-' then it is considered empty."*

PairBuddy provided sample code for test cases or methods in the text chat. For instance, PairBuddy offered to contribute a test case for PG6, *"I can generate test cases based on the user scenarios. Would you like me to do this?"*, and a method for PG1, *"Does this code help you out any?"*

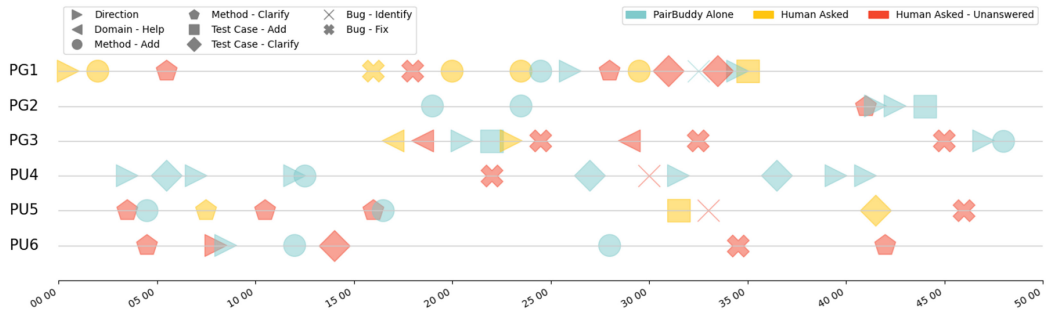


Fig. 3. A timeline illustrating PairBuddy's contributions during the pilot study. Contributions were either offered by PairBuddy (blue) or asked by participants (yellow). In some instances, PairBuddy was unable to provide help (red). Contributions types include direction, domain-related help, method/test case addition, method/test case clarification, and bug identification/fixing.

Finally, PairBuddy provided guidance for finding and fixing bugs. When PG1 made a mistake in the vertical win method, PairBuddy commented, “I think there is a syntax error on line 41.”

PairBuddy failed to answer 23 of the 35 questions (65.7%) asked by participants (red in Figure 3). However, PairBuddy was designed to follow-up with separate contributions if possible (blue in Figure 3). For example, when PU5 asked, “What would the arguments be for the method to find if there is a winner?”, PairBuddy was not designed to answer this difficult type of question, so instead, it made a suggestion from online, “Is this code example from online useful?” However, when PairBuddy failed to provide an example, PU6 complained, “You can't copy and paste it from somewhere else again please?”

**(2) Effect of Programming Experience:** Figure 3 shows the trend that graduate students (PG1–PG3) interacted with PairBuddy during the second half of the study, while undergraduate students (PU4–PU6) interacted from the beginning. One graduate student, PG1, interacted from the start, but PG2 and PG3's first interactions were at 19:00 minutes and 17:30 minutes, respectively. In his interview, PG2 said that he wanted help just-in-time and preferred to work solo on his tasks, “I was thinking, and I want to have time for... quiet and focus.” We conjecture that experienced participants were less trusting initially, but their trust increased overtime. For example, PG3 said, “As it went on, I was like, ‘Hey,... we have the same idea for this.’”

**(3) User Experience:** Participants enjoyed working with PairBuddy and appreciated its company. PG2 found PairBuddy's voice supportive saying, “Just hearing him talk helped me stay on track and stay focused on the task at hand and what needs to be done.” Similarly, PU5 indicated that she often thinks out-loud with a partner in her classes, explaining, “I'll like talk, ‘I don't know what I'm doing,’ and [my friend] will be like, ‘I don't either.’” Additionally, PG3 felt synergy with PairBuddy, commenting, “I think we work together well.”

Participants also enjoyed the motivational aspect of PairBuddy. In fact, all participants responded positively to motivation. For example, when PG1 fixed a bug that he was struggling with, PairBuddy celebrated saying, “Yay! You did it!”, and PG1 smiled responding, “Thanks!” Similarly, PU6 said, “Thank you so much dude. Your motivation is so good. Thank you, thank you, your motivation is getting me through.” However, participants' reactions to positive feedback varied, as PG2 only said, “Okay thanks.”

#### 4.7 Lessons Learned

Feedback from the pilot study revealed many shortcomings of PairBuddy's design and informed modifications to the original design decisions for use in the main study. We found the following limitations of both PairBuddy's design and the study's implementation:



**Avatar Did Not Support Lip-Sync:** Participants rarely looked at the avatar window throughout the study. PU6 put it bluntly, “I forgot he [the avatar] was even there.” While the avatar often went unused, PU4 noted that avatars play a specific role in communication for him, “I’m actually hard of hearing just a little bit, but I take cues, at least for certain words, I take cues from reading lips, so [lip-sync] would actually help me a lot.” Additionally, research finds that lip-synchronization heightens the level of an avatar’s embodiment [91]. However, in the pilot study, PairBuddy did not include a lip-sync feature, making communication less accessible and embodiment weaker.

**Gender Toggle was Not Utilized by Participants:** All six participants never changed the gender of the avatar, which limited our ability to make any explicit determinations on participants’ gender preferences. Additionally, the style of interaction between genders was indiscernible due to the low participant count, causing individual differences to out-shine gender-based differences. To gather relevant gender data, we must use an approach that can discern individual gender preferences.

**Participants’ Bias Toward Text:** While participants were explicitly instructed that they could communicate using either voice or text, participants initially felt more comfortable using text and often forgot that voice was even an option. In fact, 5/6 participants used text rather than voice. However, when PU6 was accustomed to typing messages, he was taken off-guard when PairBuddy responded to his voice, “Wait, can you hear me? That would be kinda cool if you can hear me.” Participants who used text forgot PairBuddy could hear them, often verbalizing their messages before typing them into the chat. For example, PG1 spoke out loud, “How to write the method?” before sending it word-for-word in the text chat. The inclination to vocalize questions suggests that the decision to use text over voice was partially based on preconceptions and habits, rather than utility. While participants were biased toward text, they spent an unnecessary amount of time typing messages. This notion of slow text communication is supported by research on interaction speed [205], which finds that typing is 3x slower than voice recognition (on mobile devices). To avoid confusion between text and voice, PairBuddy’s design should emphasize one medium of communication rather than using both simultaneously.

**Participants Needed More Guidance:** Although we envisioned that PairBuddy would act as both a driver and navigator, design decisions I1 (indirect driving) and T2 (code examples from online) transformed PairBuddy’s role as a driver into a code recommendation system. Since PairBuddy could not directly contribute to the code, its ability to meaningfully guide progress was limited.

**Interjections Become Interruptions:** While participants generally appreciated when PairBuddy interjected with helpful dialogue, they did not like it when interjections became interruptions, especially when they were deep in thought. PG2 explained, “When I’m trying to think about a problem... Most of the time I need a quiet place to think.” Similarly, PG1 complained, “Some of the things were distracting to me and unexpected... That I didn’t like.” While research finds that interruptions greatly reduce progress on writing assignments [80], our participants still wanted an agent that is engaging and socially present. PG3 expressed his preferences, “I would prefer someone who is more engaged,” and PU5 did too, “I don’t feel as crazy or lonely because it’s like ‘Oh, I’m talking to a computer. It’s fine,’ but like when you’re doing it alone it’s like, ‘Ah, there’s no one here.’”

**Lack of Support for Creativity:** As identified in previous research [134], programmers follow the Osborn–Parnes Creative Problem Solving Process when working together. However, with PairBuddy, participants were able to skip major steps along the way. Osborn–Parnes Creativity is frequently used to understand the creative process of an individual [112, 122, 182, 184] and is vital to programmer success [73, 251, 273], especially when solving open-ended problems



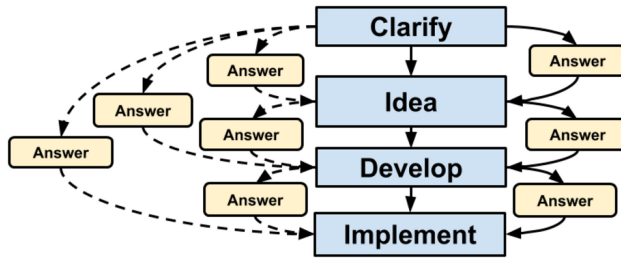


Fig. 4. Creativity stages (blue) progress from Clarify to Implement, but humans skipped the Idea and Develop stages (dashed arrows) when PairBuddy provided answers (code/test cases). Typically, human-human pairs progress through each stage in succession (solid arrows).

[24, 36, 142, 149, 187, 256]. Particularly, it describes how programmers step through a series of creativity stages; Figure 4 illustrates how programmers start at the Clarify stage of creativity, where information, goals, and challenges are identified. In the subsequent Idea stage, a wide variety of potential solutions are generated. Then, in the Develop stage, solutions are evaluated, strengthened, and selected for “best fit.” Finally, in the Implement stage, the resulting solutions are written into code. In human-human pair programming, humans follow the creativity stages in sequence [134] (solid arrows in Figure 4). However, when working with PairBuddy, we observed that participants skipped the Idea and Develop stages by directly copying implementation from PairBuddy’s code examples (dashed arrows in Figure 4). Pair programming has been noted for its ability to teach effective creative strategies [213]; based on feedback from their partners, pair programmers can gain experience by practicing each stage of creativity. Therefore, if participants are to experience the creative benefits of pair programming, PairBuddy must support its partner during every creativity stage: specifically the skipped stages of Idea and Develop.

## 5 MAIN STUDY (ITERATION 2)

Through the insights gained from the pilot study and a more comprehensive research review, we improved the design of PairBuddy and conducted a second “main” Wizard of Oz study.

### 5.1 PairBuddy Design (Iteration 2)

The design of PairBuddy for the main study intends to imitate human pair programming to the highest degree that can be realistically achieved through the capabilities of current research. Table 5 lists the design decisions used in the main study. Those adapted from the pilot study remain white, while new design decisions are highlighted in blue. Design decision modified or added for the main study are as follows:

#### Design: Creating the Interface and Interactions

**(A) Interface – Embodiment via Avatar, Gender, Voice, and Text:** For design decision F1, PairBuddy’s dynamic 3D avatar was modified to include lip-synchronization via the Facerig avatar embodiment software. Furthermore, to explore the trade-offs of using text vs. voice for a pair programming agent, we modified design decisions F3 and F4: communication with PairBuddy will be primarily verbal, while reserving the text chat for resources such as images or links. This decision attempted to model typical human-human remote communication.

**(B) Interactions – Direct Driving, Timed Feedback, Adapted Skill, Typing Speed, and Redirect Suggestions:** To allow PairBuddy to more effectively guide code progress as an active driver, we modified design decision I1: PairBuddy will make direct code contributions through the

Table 5. Design Decisions of PairBuddy's Main Study Implementation

ID	Design Decision	Description	Example Sources
Interface			
F1	Avatar	Embodied by a 3D lip-synced avatar	[25, 28, 43, 151, 229, 270]
F2	Gender	Gender can be toggled	[33, 50, 236]
F3	Voice	Communicate via voice synthesis	[28, 41, 242]
F4	Text Chat	Paste images or links	
Interactions			
I1	Direct Driving	Edit code via IDE	
I2	Adapted Skill	Balance contributions via frequency and size	[19, 56, 75, 239, 248]
I3	Timed Feedback	Feedback at appropriate time	[31, 128, 132, 144, 211]
I4	Typing Speed	Paste small sections of code	
I5	Redirect Suggestions	Participants implement their suggestions	
Soft Skills			
S1	Greeting	Introduce itself	[120]
S2	Motivation	Encourage, recognize, comfort, commend	[18, 45, 64, 78]
S3	I vs. We	Share success and personalize mistakes	[119, 197]
S4	Uncertain/Verification	Show uncertainty via verification of work	[22, 133]
S5	Social Presence	Actively listen rather than interrupt	[53, 80, 104, 175]
Technical Skills			
T1	Write/Feedback Tests	Generate test cases and feedback	[17, 74, 161, 163, 168, 196]
T2	Write/Feedback Code	Generate code and feedback	[60, 124, 125, 177, 193, 274]
T3	Guidance	Provide direction via user stories	
T4	Creativity Support	Prompt divergent and convergent thinking	[115, 116]
T5	Feature Location	Locate code from a description	[147, 156, 210]
T6	Unnecessary Code	Suggest deleting unused code	[228]
T7	Missing Code	Determine where more code is needed	[89]

Decisions modified from the pilot study remain white, while new additions are highlighted in blue.

IDE rather than sending code snippets through the text chat. Furthermore, as an active navigator, PairBuddy will provide more specific feedback. With this new ability, PairBuddy negotiated the exchange of pair programming roles through suggestions such as, “*I can try if you’d like,*” and, “*Actually, I’m a bit stuck. Mind if you take over?*”

To more effectively adapt PairBuddy’s skill-level to that of its programming partners, we modified design decision I2: PairBuddy will tailor the size (rather than only the frequency) of its contributions. For example, when participants struggle with code, PairBuddy will only provide small contributions (e.g., structure of the code or single lines of code) to help them become “unstuck.” Alternatively, PairBuddy will match participants who effortlessly complete the task with larger contributions (e.g., a full test case). For an automated agent, the size of participants’ contributions could be measured using, for example, the number of lines of code, variables, or user stories completed. Additionally, PairBuddy’s intent to balance contributions will regulate its suggestions to switch pair programming roles.

Feedback has substantial impact on learning and achievement [105, 106, 128, 144]. Among its influential properties is the timing of the feedback. Psychological research suggests that feedback on difficult concepts should be delayed, while feedback on simple concepts is more beneficial when immediate [31, 132, 211, 233, 234]. For example, IDEs already provide instant feedback via error highlighting, but when a programmer uses a misguided approach, delayed feedback is preferred to allow them time to evaluate the feasibility of their ideas. Therefore, we formulated design decision I3: PairBuddy will provide appropriately timed feedback [19]. With this decision, we hope to maximize the impact of feedback and eliminate unnecessary interruptions.

Since PairBuddy now directly contributes code through the IDE, it is necessary to consider the speed at which PairBuddy will type. Intentionally limiting an agent’s typing speed could

potentially be unappealing and cause programmers to become impatient, yet pasting large chunks of code could be overwhelming and difficult to comprehend. Therefore, we formulated a design decision I4: PairBuddy will incrementally paste small snippets of code to prevent programmers from becoming impatient or overwhelmed. This design decision was informed from our best reasoning rather than previous research.

In pair programming, the role of the navigator includes providing feedback and suggestions to the driver. Sometimes, the navigator performs “backseat driving” where they instruct the driver directly [117]. However, current technology does not fully support the implementation of arbitrary ideas in this way, so in the pilot study, PairBuddy was designed to immediately admit its limitations saying, *“I’m sorry, I don’t know how to help with this.”* Unfortunately, this response often marked an abrupt end to the conversation. Therefore, to encourage higher engagement, we formulated a design decision I5: PairBuddy will redirect suggestions back to the programmer through dialogue such as *“How would that look like?”*, *“Do you want to try?”*, or *“Can you do that for me?”* However, if its partner insists, only then will PairBuddy admit its limitations.

### **Programmer: Integrating Technical and Soft Skills**

**(A) Soft Skills – Leadership, Uncertainty, Social Presence:** Research finds that a democratic leadership style is most effective for pair programming [133]. A way for PairBuddy to integrate democratic leadership is to practice effective pronoun use. Research has shown that leaders use collective pronouns (e.g., “we” and “us”) to gain influence in a group [119], but use personal pronouns (e.g., “I” and “me”) to “own up” to mistakes [197]. Therefore, we formulated design decision S3: PairBuddy will display leadership by attributing successes to the group while taking ownership for its mistakes. For example, when a test case fails, PairBuddy will take ownership saying, *“I think I made a mistake,”* but if all the test cases pass, PairBuddy might say, *“Great, we did it!”*

Research suggests that conversational agents should pre-emptively use uncertainty to avoid situations where they make large miscommunications [19, 22]. The resulting conversational breakdowns decrease users’ satisfaction, trust, and willingness to continue talking to conversational agents [113, 114, 152]. Additionally, research on pair programming shows that people often ask for verification after each creative stage of development [133]. Therefore, we formulated design guideline S4: PairBuddy will convey uncertainty by asking for verification to prevent conversational breakdown. To potentially support this, the uncertainty of PairBuddy’s dialogue could be based on the confidence of the machine learning algorithms used to generate code and feedback.

In the pilot study, participants voiced complaints about the timing of PairBuddy’s feedback. When they were deep in thought, participants wished that PairBuddy listened rather than interrupted. However, maintaining an active presence as a navigator is key to preserving the balance between roles. To this end, we utilize active listening—a technique to allow a speaker an outlet for self expression [202]—to maintain social presence in a conversation. For conversational agents, social presence is the feeling and perception of interacting with a human being [180], and has been shown to increase the use of conversational agents [175], and can even increase perceived usefulness, trust, and enjoyment [104]. Therefore, we formulated design decision S5: PairBuddy will increase its social presence using active listening through acknowledgements and confirmation questions. In this way, PairBuddy can avoid interruptions by encouraging participants to think their ideas out-loud through dialogue such as, *“What are you thinking?”* and *“What does this code do?”*.

**(B) Technical Skills: Guidance, Creativity Support, Feature Location, Unnecessary Code, and Missing Code:** To encourage higher engagement with PairBuddy relative to the pilot study, we conducted a thorough research review to identify the potential capabilities of automated code and feedback algorithms, as well as techniques to support creative problem solving.

PairBuddy’s power to provide context-specific code and feedback is dependent on the capabilities demonstrated by existing research. Automated code and feedback techniques [60, 274] show

the ability to provide feedback on code using a dataset of past programming solutions. However, this research is either limited by the specificity of the feedback [60] or has only been demonstrated in a simple programming language (iSnap) [274]. To increase the specificity of PairBuddy's contributions, we modified design decisions T1 and T2: PairBuddy will have the limited ability to contribute context-specific code (as the driver) and provide meaningful feedback (as the navigator) based on a dataset of past solutions. For example, if PairBuddy detects that a participant's code is semantically similar yet deviates from a known solution, it might suggest, *"There might be a mistake on line 66."* However, if the participant writes code that deviates from all past solutions within the database, PairBuddy will be unable to provide guidance.

Since PairBuddy attempts to replace a human partner, it will be designed to integrate the creative problem solving stages used by humans. However, in the pilot study, PairBuddy only supported the Clarify stage (via direction from user stories) and the Implement stage (via example implementation from online). Therefore, we formulated design decision B4: PairBuddy will add support for the remaining Idea and Develop stages by using concepts from Idea Garden [115, 116]. Idea Garden suggests the use of probing questions to promote a programmers' use of diverse problem-solving strategies. We integrated two strategies into PairBuddy's script: working backwards from the goal and encouraging divergent thinking before convergent thinking. For example, PairBuddy might ask, *"What are all the possible ways we could do this?"*, *"Why do you think so?"*, or *"What data structure could we use?"* As a navigator, PairBuddy can prompt the participant to move from the Idea to the Develop stage by asking, *"What would that [idea] look like?"* As a driver, PairBuddy will provide the empty structure of the code to help programmers conceptualize the solution. Additionally, PairBuddy can discuss the general structure of the code saying, *"I think we should use a for loop/if statement here,"* or *"Should we use a while loop or a for loop?"*

Both static [156] and dynamic [147] techniques exist to automatically search source code for specific features or descriptions. For the Java programming language, the Eclipse plugin Flat3 [210] searches source code using arbitrary descriptions (e.g., "file saving"). Based on feature location algorithms, we formulated a design decision T5: PairBuddy can identify locations within the code that match a description. For example, if the participant asks, *"Do we return false in the tie game detection function?"*, PairBuddy can identify the *isTied()* function as the context to the question.

Many programming IDEs include refactoring tools that automatically detect uncalled functions or unused variables via data-flow analysis [57]. For the Java programming language, UCDetector [228] is an Eclipse plugin that can identify unused variables, functions, and classes. Based on this functionality, we formulated design decision T6: As a navigator, PairBuddy can detect unnecessary blocks of code and suggest that the programmer consider deleting or modifying them. For example, PairBuddy might suggest, *"I believe we have unnecessary code in the tie test function."*

Research on automated feedback has shown the ability to automatically determine the locations where code is missing based on pre-defined functions in Haskell [89]. Since Haskell is a functional programming language, this research does not necessarily apply to an imperative language like Java. Regardless, the supporting research of previous design decisions [17, 60, 163, 168, 274] would be potentially capable of identifying missing code, so we formulated design decision T7: As a navigator, PairBuddy will direct the driver's attention to locations where more code is needed.

## 5.2 Wizard/PairBuddy Implementation

For the main study, the wizard's interface, implementation, and script integrated the new design decisions for PairBuddy's second iteration, including the following changes from the pilot study:

Before conducting the main Wizard of Oz study, four trial studies provided initial feedback on the task, study design, wizard implementation, and design decisions. During these studies, our

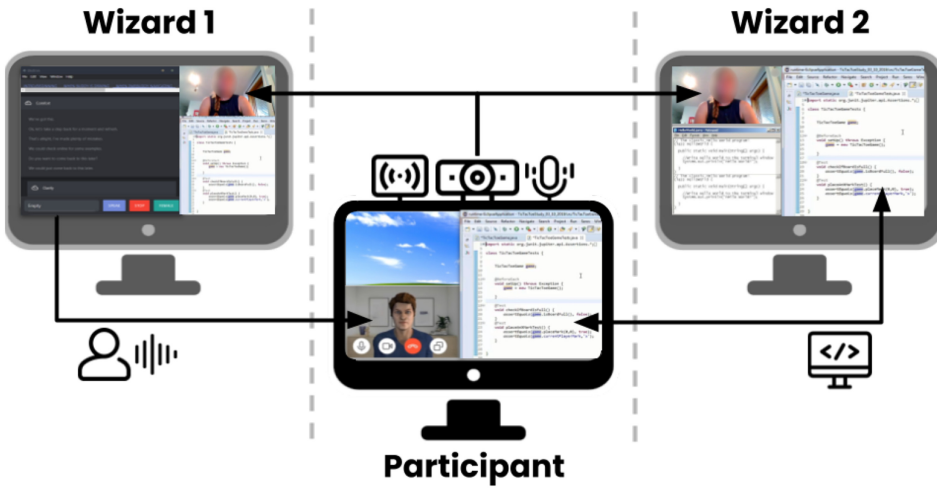


Fig. 5. The main study's wizard implementation. Both wizards (Wizard 1 and Wizard 2) monitored the participant's webcam, microphone, and screen via a communication tool (e.g., Skype, Discord, MS Teams). Wizard 1 used a custom interface to select dialogue templates and send PairBuddy's avatar and voice over the same communication tool. Wizard 2 used a dataset of past solutions to contribute into participants' IDE directly.

single wizard had difficulty balancing all of their new responsibilities as PairBuddy; therefore, our main study leveraged two wizards.

The first wizard (Wizard 1) controlled PairBuddy's avatar and voice interface with participants. To simulate PairBuddy's design decisions, Wizard 1 utilized a basic automated infrastructure where (1) a custom user interface (refer Figure 5) allowed Wizard 1 to quickly select dialogue templates from the script and fill in blanks with contextual information when necessary; (2) using the selected dialogue, PairBuddy's voice was generated via Google Text-to-Speech [98]; (3) the resulting audio facilitated lip-synchronization of PairBuddy's avatar generated using the Facerig embodiment software [227]; and (4) both PairBuddy's avatar and voice were sent to participants via a communication tool (i.e., Skype, Discord, and MS Teams); likewise, both wizards monitored participants' video, audio, and screen directly.

As described in I3: Timed Feedback, the wizards waited an increasing amount of time before providing feedback on incorrect code. As the navigator, Wizard 1 followed a set protocol when providing feedback to participants: (1) If participants explicitly asked, Wizard 1 provided instant feedback; (2) For out-of-turn feedback: (i) If participants made small syntax errors (e.g., wrong number or missing curly braces), Wizard 1 waited 10 seconds after participants finished typing before providing feedback as the navigator to allow participants time to double check their spelling; (ii) For small logical errors (e.g., incorrect conditional statements or variables), Wizard 1 waited 30 seconds before giving feedback since logical errors are more difficult to catch; (iii) When participants made medium to large mistakes (e.g., multiple lines of code that do not match past solutions), it is possible that participants found an alternative solution. However, if they became stuck and stopped discussing ideas, Wizard 1 asked open-ended questions (T4: Creativity Support) before asking participants for the driver role.

The second wizard (Wizard 2) acted as an automated code and test cases generator (design decisions T1 and T2). As the driver, Wizard 2 pasted code snippets into participants' IDE via the Saros remote collaboration tool to correct or continue participants' existing work. To determine code correctness, Wizard 2 simulated a semantic diff (e.g., [136, 204]) between participants' code



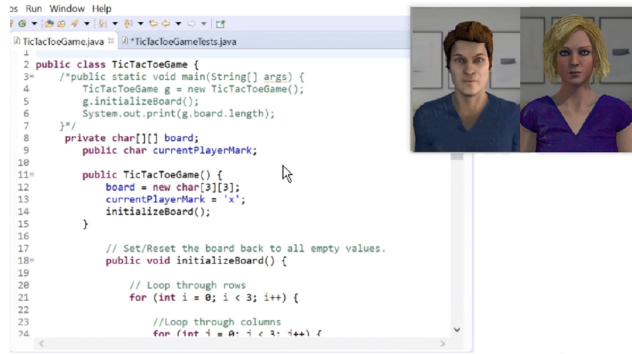


Fig. 6. The main study's interface including the Eclipse IDE and the man/woman avatar. Note, each participant only saw one avatar at a time.

and a database of past solutions collected from both pilot and trial studies. With this knowledge, Wizard 2 followed their own protocol to determine the appropriate time to contribute code as the driver: (1) If participants had difficulty getting started, Wizard 1 offered to drive and Wizard 2 provided some of the structure of a past solution (e.g., function declarations and conditional statements). Using only incremental contributions, Wizard 2 avoided overwhelming participants since PairBuddy could not explain the logic behind its contributions; (2) If progress was slow or participants became stuck, Wizard 2 instead provided small contributions from the most similar past solution to get participants back on track. If participants wrote a large amount of unrecognized code, Wizard 1 offered to comment their code prior to Wizard 2 providing a separate contribution; (3) If participants progressed quickly (i.e., their code matched a past solution or the test cases passed), Wizard 1 offered to drive and Wizard 2 provided a similarly sized contribution. The two wizards maintained close communication throughout the studies, particularly to coordinate code contributions when PairBuddy was the driver.

We assumed PairBuddy's algorithms were competent, so our wizards only made mistakes (intentional or unintentional) either once or twice per session.

Additionally, we chose two avatars (man and woman) based on their professional appearance from the limited selection of Facerig avatars (refer Figure 6).

### 5.3 Participants

Due to the COVID-19 pandemic, study sessions were conducted remotely, and recruitment of participants was done via snowball sampling, social media (Facebook [225], Twitter [226]), and hiring sites (Upwork [83]). Our recruitment approach helped collect participants from universities and industries across the country, in addition to local participants. We gathered 14 participants on a first-come first-serve basis, including 8 students (4 men / 4 women) and 6 professionals (3 men / 3 women). In their background questionnaires (refer Table 6), all participants self-identified in the binary (men and women). We purposefully achieved a gender balance since research has shown a difference in preference and behavior across various genders [39, 48, 87, 160, 166, 215]. All participants were regular programmers with prior experience writing code in Java, although not all had used Java recently. We refer to participants of the main study as either MS# or MP# for students and professionals, respectively. For example, MS4 is the fourth student participant of the main study. To incentivize participation, student participants were given \$20 and professionals were given \$40 in Amazon gift cards. This discrepancy was approved by IRB since recruiting professional developers is more challenging due to their limited availability and high hourly wage [172].



Table 6. Demographics of the Main Study Participants

P#	Age	Gender	Education	Programming Experience
MS1	19–23	Woman	Undergrad	3 years
MS2	19–23	Man	Undergrad	2 years
MS3	19–23	Man	Undergrad	2 years
MS4	19–23	Woman	Undergrad	2 years
MS5	19–23	Woman	Undergrad	3 years
MS6	19–23	Man	Undergrad	3 years
MS7	19–23	Woman	Undergrad	2 years
MS8	19–23	Man	Undergrad	4+ years
MP9	30–40	Woman	Masters	4+ years
MP10	41+	Man	Masters	4+ years
MP11	19–23	Man	Undergrad	3 years
MP12	30–40	Woman	Masters	3+ years
MP13	30–40	Man	Masters	2 years
MP14	19–23	Woman	Undergrad	1 year

## 5.4 Study Design

The study design remained the same as the pilot study except for the following changes:

As previously mentioned, the main study was conducted during the COVID-19 pandemic, so all study sessions were conducted virtually. Prior to this 40 minute study, participants completed a consent form and background-questionnaire. To combat the low level of interaction with PairBuddy in the pilot study, the main study included a tutorial [4] that explained PairBuddy’s interface (i.e., voice communication and direct code editing) and encouraged interaction. Since its learnability should ideally be quick, we did not disclose PairBuddy’s design decisions to evaluate how participants adapt to PairBuddy’s abilities and limitations. Additionally, a pre-study questionnaire was used to establish a baseline for participants’ self-efficacies using a 7 point Likert scale [118] for a maximum score of 63 points. The tutorial and self-efficacy questionnaires can be found [4].

The study was designed to ensure that each participant interacted with both of PairBuddy’s gender embodiments. Rather than evaluating gender preferences directly, we used a within-study design where the avatar’s gender was changed halfway through the study to compare preferences on an individual basis. While we recognize the potential confusion that PairBuddy’s gender change may have caused, it was a necessary trade-off to adequately assess gender-bias. Additionally, no participants reported confusion in their interviews. To counterbalance gender, we evenly distributed PairBuddy’s gender between the men and women participants.

Each study session was followed by an additional self-efficacy questionnaire and interview questions to triangulate our study findings. The same interview questions used in the pilot study were integrated with additional questions related to the usability of the features that were more difficult to evaluate using our think-aloud study.

## 5.5 Data Analysis

The resulting transcripts were analyzed using the same methodology as described in the pilot study (refer Section 4.5), including the contribution type codeset (Table 4). Three researchers independently coded 20% of the transcripts and reached agreement on 93% of the coded data by calculating inter-rater reliability using the Jaccard index. The wizard’s script, study transcripts, and interview transcripts can be found [4].

## 5.6 Limitations

One limitation of the main study is its small sample size of 14 programmers. Although our sample was gender-balanced and included a wide range of skill levels (8 students and 6 professionals), the study's small size did not allow any room for further stratification. On the other hand, our study helped evaluate PairBuddy's usability with a diverse population of programmers. While we only studied one programming language (Java) and one IDE (Eclipse), our focused approach was an appropriate choice to show baseline feasibility. Future studies are needed to evaluate a larger variety of programmers and languages.

Programmers' varied experience with Java, test-driven development, and pair programming may have affected their experience with PairBuddy. For example, PairBuddy provided very few contributions for our most experienced professional, MP10, while providing more contributions to other less-experienced professionals. In the future, more variables should be considered in our analysis.

We acknowledge the variety of other dimensions (e.g., non-binary, race, and ethnicity) to PairBuddy's persona that can influence participants' potential biases as research has identified effects that an avatar's persona has on user behavior [26, 155]. In fact, our brains subconsciously categorize characteristics of identify (e.g., perceived sex and race) in just 200 milliseconds and place people into social categories informed from stereotypes and biases [54]. However, in this study, we focus on PairBuddy's gender due to the well-documented prevalence of gender-bias within computer science [111, 162, 241]. To alleviate these biases in the future, we can include multiple personas and allow participants/users to select their preferred avatar.

Since the wizards manually controlled PairBuddy's actions, there was a noticeable delay between participants' requests and PairBuddy's responses. While our custom wizard interface reduced latency, nonetheless, searching for the appropriate response and filling-in any contextual information produced lag. However, shorter replies like "Okay" were quickly typed to avoid navigating the application's interface.

Our results are based on the simple task of tic-tac-toe, but the game's simplicity may have affected our usability results since one participant mentioned that he might not trust PairBuddy for more complex tasks. While the variation of task complexity needs to be further explored with additional studies, nonetheless, we believe that our task was a good representation of student assignments and that our results confirm the feasibility of PairBuddy in educational settings.

We studied the usability of PairBuddy in a virtual lab setting for only 40 minute sessions. Although this was a good starting point, our study could not provide a long-term perspective regarding the usability and user experience of programmers at this moment. For instance, will interactions with PairBuddy for longer periods of time increase or decrease trust?

Finally, we were forced to conduct a virtual lab study due to COVID-19, which may not mimic traditional lab studies where participants sit in a controlled lab environment. However, we view this as an opportunity to evaluate the usability of PairBuddy in a hybrid lab study consisting of both real-world and controlled lab settings.

## 5.7 Results

The second iteration of PairBuddy served as a realistic portrayal of a pair programming conversational agent. Insights collected are as follows:

**(1) Helping Programmers:** The results from Figure 7 show that participants were most likely to request help from PairBuddy (yellow) during the second half of sessions. We conjecture that

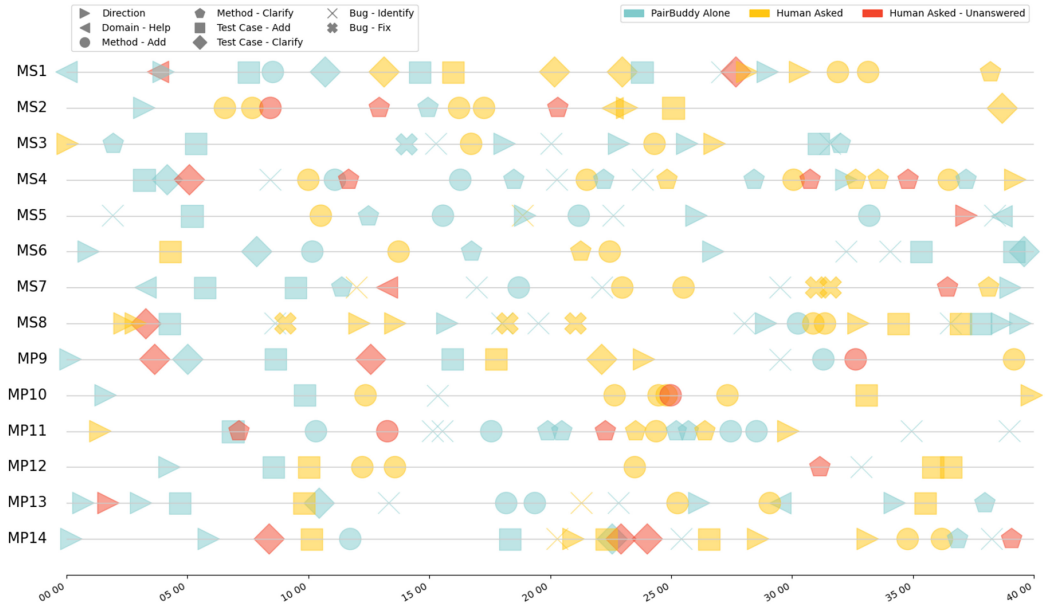


Fig. 7. A timeline illustrating the contributions PairBuddy made in the main study. Contributions were either offered by PairBuddy (blue) or asked by participants (yellow). In some instances, PairBuddy was unable to provide help (red). Contributions include direction, domain-related help, method/test case addition, method/test case clarification, and bug identification/fixing (refer Table 4).

the change in requests overtime is the result of participants' initial distrust in PairBuddy, as noted by MS5's comment, "I think I was a little distrustful at first half of the application, but then after working with it, especially after I saw that it was helping me solve the problem... I trusted it a little more."

Participants' confidence in their coding abilities increased with the usage of PairBuddy. On average, participants' self-efficacy scores increased +3.64 out of a total 63 points (49.71 to 53.07) after their interactions with PairBuddy (Table 7). Only 3/14 participants reported a decrease in self-efficacy. MS6 was by far the largest outlier, with a difference in self efficacy of -18 points: 13 points fewer than any other participant. Otherwise, 10/14 participants saw an increase in their self-efficacy.

**Getting Directions and Domain-Related Help:** Participants received direction from PairBuddy as seen in Figure 7. PairBuddy's guidance helped start participants in the right direction, as PS5 explained, "It definitely helped me get started." Participants received guidance for an unfamiliar language, including MP11 who said, "Getting started was the hardest part for me, like trying to wrap my head around Java again." Some learned new techniques, like PS2 who commented, "How to do test-driven development... the robot helped me make sure that I was writing it correctly, and I liked that."

**Clarify a Method or Test Case:** Participants mentioned that PairBuddy helped clarify task objectives, including MP9 who commented, "I felt like PairBuddy was pretty good at understanding the overall objective of the project." This positive feedback was largely caused by PairBuddy's ability to Clarify the task based on user stories or generated solutions. For example, PairBuddy provided insight when MS6 was trying to finish a method, "I think we are missing code in this function." However, PairBuddy often could not answer Idea-related questions, so when MP14 asked, "Do we

Table 7. Self-efficacy Scores for Main Study Participants from Pre- and Post-study Questionnaires

P#	Pre-Study	Post-Study	Difference
MS1	50	56	+6
MS2	57	61	+4
MS3	50	45	-5
MS4	52	58	+6
MS5	41	56	+15
MS6	46	28	-18
MS7	43	55	+12
MS8	55	51	-4
MP9	48	58	+10
MP10	51	54	+3
MP11	60	61	+1
MP12	44	53	+9
MP13	50	50	0
MP14	49	57	+12
Average	49.71	53.07	+3.64
$\sigma$	5.105	8.102	8.329

Responses to nine questions were scored on a seven point Likert scale for a maximum score of 63 points.

*need to account for the opposing player mark?*”, PairBuddy reflected the question by responding, “*What do you think?*” While such responses continued the conversation, they did not provide the direct, human-like answers that participants desired.

**Adding Methods or Test Cases:** As the driver, PairBuddy wrote both test cases and method functionality. PairBuddy’s contributions helped participants make progress, and almost all participants ended up reusing the code given by PairBuddy. This includes MP11, who successfully adapted PairBuddy’s test case for a vertical win into both horizontal- and diagonal-win test cases. For some participants, PairBuddy’s code offered a new approach. For instance, when MS1 struggled, PairBuddy’s alternative suggestion provided the insight necessary for MS1 to complete her vertical win functionality.

**Identifying or Fixing Bugs:** As both the driver and navigator, PairBuddy helped participants find and fix bugs, leading to increased code quality. As the navigator, PairBuddy would often point to locations of errors and provide hints. For MS1, PairBuddy pointed, “*I think there is a typo on line 107,*” leading her to immediately fix the mistake. On some occasions, PairBuddy had to fix bugs itself when participants struggled to understand the error that PairBuddy was referring to. When MP13 had trouble fixing an identified mistake, PairBuddy took action by asking for permission to drive and fixing the error while completing the remainder of the method.

**When PairBuddy Failed to Answer:** Figure 7 shows the instances where PairBuddy was unable to help participants (red). Either PairBuddy would redirect queries (e.g., “*What do you think?*”) or admit its limitations (e.g., “*Sorry, I’m not good at logic.*”). In the main study, PairBuddy was unable to directly answer user queries 26 out of 114 times (22.8%); a significant decrease from 23 out of 35 times (65.7%) in the pilot study. The increased utility from PairBuddy’s answers (combined with its increased out-of-turn dialogue) may have contributed to participants’ overall higher interaction. PairBuddy’s inability to provide help like a human negatively affected some participants’ trust and confidence in PairBuddy’s abilities. After his queries went unanswered three times, MS2 mentioned

Table 8. Responses to Interview Questions in the Main Study

Interview Question	Yes	Neutral	No
Did you enjoy working with PairBuddy?	12	1	1
Did PairBuddy help you solve the problem?	13	0	1
Did you learn anything from PairBuddy?	6	3	5
Did you trust PairBuddy?	12	2	0
Did you trust PairBuddy's navigation?	12	1	1
Did you trust PairBuddy's navigation over a human?	5	6	3
Did you trust PairBuddy's driving?	8	5	1
Did you trust PairBuddy's driving over a human?	4	4	6
Was PairBuddy's avatar helpful?	6	6	2
Was PairBuddy's voice helpful?	13	1	0
Would more text have been helpful?	2	5	7
Do you prefer a casual over a polite agent?	8	4	2
Was PairBuddy too casual (yes) or too polite (no)?	0	14	0
Do you prefer human-like dialogue over robotic?	5	6	3
Do you prefer positive over neutral feedback?	12	2	0
Do you like negative feedback?	6	7	1

PairBuddy's inaction as a negative aspect in his interviews, *"He couldn't respond to that. At least not in a way that made sense."* Participants expressed concerns about PairBuddy's inability to answer "why" questions or give reasoning behind the code it wrote. After PairBuddy failed to help, MP9 later commented in her interviews, *"I felt like they [PairBuddy] knew what to do. But they weren't always able to communicate to me the why."*

**(2) Usability of PairBuddy:** In their interviews, participants expressed mostly positive experiences with PairBuddy (Table 8):

#### ***Enjoyed/Helped/Learned:***

**Enjoyed:** A vast majority (12/14) of participants enjoyed working with PairBuddy. MP10 was the most enthusiastic commenting, *"It is really, really awesome."* PairBuddy's design surprised many participants, including MS5 who said, *"I thought it's really cool technology. I didn't know stuff like that was out there, really. So it's pretty novel."* The only participant who disliked PairBuddy was MS6, who mentioned, *"Not particularly...There wasn't as much human interaction."*

**Helped:** Most participants (13/14) expressed that PairBuddy helped them solve the task. MS8 thought PairBuddy saved him time saying, *"It did, yea. It saved me some time because a lot of that time would have been spent doing trial and error."* Even MS6, who disliked PairBuddy, indicated that it helped him solve the problem, *"[PairBuddy] helped me understand... how to start approaching things so that I could start going myself."*

**Learned:** However, only 6/14 participants said that they learned from PairBuddy, and 3/14 remained neutral. Some participants mentioned that they learned task-related concepts, including MS6 who commented, *"I think learning the general structure of what pair programming is."* Others learned new creative strategies, including MS3 who said, *"I got it like a new perspective... I guess I still need to start thinking more outside the box."* On the other hand, 5/14 responded they did not learn from PairBuddy. Unlike a human, PairBuddy did not explain the code it wrote; MS3 described that PairBuddy *"left it up to [their] interpretation."* PairBuddy's inability to express its reasoning may have contributed to lower satisfaction with PairBuddy.

**Future Design:** To improve PairBuddy's design going forward, we will explore more ways to impart knowledge through code and test case generation to the extent that current research allows (details in Section 6 (1)).

**Trusted:** Most (12/14) participants trusted PairBuddy overall (Table 8). Participant's trust may have been based on the assumption that, "[PairBuddy] was familiar with what the program should look like," as described by MS7.

**Navigator and Driver Roles:** PairBuddy was trusted as a navigator (12/14) more than as a driver (8/14).

PairBuddy acted as an *active navigator*, using soft skills (leadership, motivation, uncertainty, and social presence) to help increase participants' confidence and trust. 12/14 participants indicated trust and confidence, as confirmed by MS2's comment, "*I thought it [navigation] was it's [PairBuddy's] most useful function.*" Similarly, MS6 mentioned, "*One thing that was helpful is that when they made their navigator suggestions, like all of them seemed reasonable.*" The motivation PairBuddy provided was praised by participants like MS7 who explained that, "*It was cute that it said, 'Good job,' and gave positive affirmations.*" PairBuddy's social presence provided a sense of security toward participants' work, as evidenced by MS2's comment, "*I think that I can keep on coding and not like, worry about checking it three more times because PairBuddy's got my back.*" MP9 treated PairBuddy as someone to talk to, commenting, "*I could say things in the same way that I would to a human.*"

PairBuddy acted as an *active driver*, using technical skills to make contributions, give just-in-time feedback, and support divergent thinking. Participants appreciated PairBuddy's divergent thinking including MS5 who commented, "*I'm gonna approach coding in the future... [by] stepping back and looking for alternate approaches.*" Although participants trusted PairBuddy's code and just-in-time feedback, they wanted even more assistance. This desire was evidenced by MS3's suggestion, "*Maybe have it give more frequent and more pointed advice on the code itself.*" MS7 wished PairBuddy's feedback was more in-line with their own ideas, "*I think my trust would decrease if I noticed that the suggestions it was giving me weren't like really aligned with what I was trying to do.*"

**Preference for PairBuddy vs. Human:** Participants had mixed opinions when asked to compare PairBuddy with a human. Some participants preferred PairBuddy's collaboration style including MS5 who commented, "*I think I was getting feedback from it akin to what I would from humans, and in fact, [it] is a little less intrusive than some humans are when it comes to how I'm writing my code.*" However, many participants wished to work with humans to discuss ideas, since PairBuddy couldn't explain its own. In fact, participants ignored PairBuddy's suggestions a total of 16 times and often chose to pursue their own ideas. MS3 wished PairBuddy provided explanations as a driver, and described that, "*The code that it [PairBuddy] wrote was functional for what it did, but it [PairBuddy] never actually explained what it drove.*" Participants also expected PairBuddy to offer more contextual feedback, as MS6 mentioned, "*If I were to ask a human, they wouldn't just say 'I may have made a mistake,' they would explain what they did and explain why they were thinking they made a mistake or why they thought it wasn't a mistake.*"

**Future Design:** Going forward, we will improve PairBuddy's ability to generate explanations and contribute to discussions as much as current technology will allow (details in Section 6 (1)).

**Embodiment:** The embodiment of PairBuddy's voice was much appreciated, while the presence of PairBuddy's avatar received mixed feelings. Feedback on the presence of the text chat leaned negative (see Table 8).

**Avatar:** Most participants preferred (6/14) or were neutral (6/14) towards the presence of an avatar. Some participants found it natural to interface with an avatar, as described by MS7, "*It*



*feels more natural to talk to it than if there's no avatar.*" Similarly, MS4 felt a closer connection to PairBuddy when personified, saying, *"[I felt connected] more so than if it was just a voice... or just text."* On the other hand, 2/14 participants did not like the avatar. These differing opinions coincide with mixed empirical evidence for the necessity of avatars for agent embodiment [65, 103, 107, 170, 238, 269].

**Voice:** Voice was PairBuddy's most desired feature (13/14 preferred, 0/14 not preferred). The ability to communicate with PairBuddy verbally was very intuitive for participants. MP14 particularly enjoyed being able to think out loud, commenting, *"I liked being able to speak... [so that PairBuddy can] understand what I'm saying while I'm talking out loud rather than having to type out everything."* MS2 found PairBuddy's voice to be less distracting when he was deep in thought saying, *"I feel like the voice is better when you're focused on programming."*

**Text:** The desire for text messages was mixed (2/14 preferred, 7/14 not preferred). Participants who advocated for text suggested it as an addition, rather than a replacement, of voice communication. MS5 suggested to display PairBuddy's messages as text bubbles, explaining, *"I also could also see like a little pop up in the corner maybe being helpful."* However, most participants thought text would be distracting, as MP9 described, *"I wouldn't want to be context switching between the programming and the messages."* These results contradict preferences for text chat in the pilot study, and may be influenced by the lack of text messages in the main study.

**Gender:** Participants' gender preferences were identified by asking whether they preferred the first or second avatar/voice. Many participants' reasoning for their choice was not gender-related; MP13's reason was *"just because I started with it,"* while MS4's was because *"she sounded a little less robotic."* While previous research has identified gender preferences and bias toward conversational agents [33, 50, 236], no useful patterns emerged from our data that evidenced any effect of PairBuddy's gender.

**Future Design:** We will keep the same embodiment design with additional choices for PairBuddy's gender, ethnicity, and accompanying voice tones.

**Tone/Style/Feedback:** Participants preferred a casual tone, had mixed opinions of PairBuddy's dialogue style, and desired non-neutral feedback (Table 8).

**Polite and casual tone:** Most participants preferred if PairBuddy used a casual (8/14) over a polite tone (2/14). Since PairBuddy's tone leaned more casual, all participants were satisfied with PairBuddy's tone (14/14). Even participants who preferred polite dialogue liked PairBuddy's tone the way it was. In his interview, MP10 commented, *"I felt it was excellent. It didn't feel to be like trying to be polite. It didn't feel to be rude as well."* Similarly, MP14 thought PairBuddy's tone achieved a good balance, saying, *"And the tone was like, not very formal, which I liked, but I wouldn't want it more casual than it is now."*

**Human vs. robotic style:** Participants had mixed opinions on whether PairBuddy should be human-like or robotic. 5/14 participants favored human-like dialogue, while 3/14 wished PairBuddy's dialogue was more robotic. MS3 idealized a more human PairBuddy explaining, *"Ultimately, I think having it speak as human-like as possible is the goal."* MP10 described that he enjoyed the social aspect of PairBuddy's human-like dialogue, *"I enjoy... my buddies and working with them... Closer to that, I think more people would like [PairBuddy] because that is how they work in real world."* Participants who preferred robotic dialogue mentioned that it would be strange or disturbing if PairBuddy was more human-like, including MS1, who commented, *"It would be weirder if it [PairBuddy] was more like a human; You kind of expect some level of not human dialogue."* Similarly, MS5 explained her feelings toward PairBuddy's style, *"I think there's something a little creepy about something that's like really close to being human, but not quite."*

**Positive/Neutral/Negative Feedback:** All participants preferred (12/14) or were neutral (2/14) toward PairBuddy's positive feedback. MS4 voiced her opinion, saying, *"It [PairBuddy] wasn't rude about it. He was very like, 'I don't know if that's gonna work,' rather than like, 'Wow, why'd you do that?' I think it made it... a little more trustworthy."* MP10 emphasized the need for exclusively positive feedback, *"The way [humans] think and the way they solve problems requires a lot of motivation and requires a lot of positive praise, even if they are doing things wrong."*

Participants interpreted negative feedback as identifying mistakes in the code, so most of them liked (6/14) or were neutral (7/14) towards such responses. This misunderstanding was evidenced by MS8's comment, *"If there are negative parts to what I'm doing, I'd like to know."*

**Future Design:** We will provide personalization for PairBuddy's script across tone, style, and feedback (details in Section 6 (6)). To avoid the "uncanny valley" effect where realism becomes unsettling [86], PairBuddy's movement and stylization will be further considered. Finally, PairBuddy will continue to give encouragement while still providing informative feedback.

**(3) Pair Jelling with PairBuddy as a Programming Partner:** Participants' interactions with PairBuddy were fewer in the first half of sessions, but increased by the second half. Infrequent interactions may have been caused by participants' unfamiliarity with PairBuddy as evidenced by MS7's interview, *"It was a little bit awkward initially to figure out what things I could say that it would actually respond to... But I figured that out pretty quickly once we started working more."* Adjusting to PairBuddy in this way is similar to the pair "jelling" period of human-human pair programming, where programmers take time to get accustomed to each other's personality, style, and abilities [117]. However, the jelling period with PairBuddy often ended when participants understood its utility as MS5 explained, *"Seeing consistent results in regards to... positively affecting my performance would increase my trust in it [PairBuddy]."* Research corroborates this claim, as potential productivity has shown to be a significant motivator for the utilization of conversation agents [34].

**(4) Unrealistic Assumptions Regarding PairBuddy's Capabilities:** Our participants assumed that PairBuddy would know all the answers as described by PS1, *"I assumed it [PairBuddy] knew the answer."* Similarly, PS11 commented, *"I went in assuming that... you guys had already tested it quite a bit, and so the robot was really familiar with what the program should look like."* This assumption of PairBuddy's knowledge is untrue since design decisions T1 and T2 provide PairBuddy with only limited abilities since automatic code and test case generation are topics of ongoing research. In the future, the inaccurate expectations of PairBuddy's abilities can be addressed through further clarifications of PairBuddy's limitations in its script, *"Sorry, I have limited capabilities. Researchers are working to make me smart enough to answer all your questions. Follow the link I sent in the text chat to read the current research on automated code and test case generation."* However, PS3 understood PairBuddy's limitations, especially for more complicated tasks that require more discussions, *"I guess it depends on the complexity and the nuance of the situation and the [user] stories that I need to take care of, so for simpler problems or more straightforward tasks, I'd say that prefer PairBuddy. But for more complex scenarios, I would rather have a human that can provide more... pointed and specific feedback and advice."*

**(5) Helpful as a Non-Judgemental Partner and When Working Solo:** PairBuddy helped participants by acting as a non-judgemental partner. PS5 explained, *"I think the big thing is sometimes as a programmer, it's embarrassing when you make mistakes, so you're stuck on something around your colleagues, but PairBuddy I don't think judges me."* PS2 described the lack of pressure to perform, *"I feel like whenever I program alone, there's less pressure to like write something that's correct the first time... So I feel like in the end, it [PairBuddy] helped me write better code."*

Table 9. Implementation Details of the Wizard in the Main Study and PairBuddy in the Future

ID: Design Decision	Wizard	PairBuddy Implementation using
<b>Interface</b>		
F1: Avatar	Facerig software	Same or alternative software
F2: Gender	Facerig's genders and Google Text-to-Speech voices	Same, consider race/accent and use GenderMag [40, 87, 108, 250]
F3: Voice	Google Text-to-Speech	Same
F4: Text Chat	Skype, MS Teams, Discord, etc.	Same or consider IDE integration
<b>Interactions</b>		
I1: Direct Driving	Saros collaboration plugin for Eclipse IDE	Plugins for other IDEs and languages
I2: Adapted Skill	Intuitive contribution detection	Discrete progress detection (e.g., # lines of code)
I3: Timed Feedback	Rule based (refer Section 5.2)	Same or reinforcement learning
I4: Typing Speed	Copy/paste small sections	Same
I5: Redirect Suggestions	Follow-up question templates	Same or reinforcement learning
<b>Soft Skills</b>		
S1: Greeting	Dialogue templates	Same
S2: Motivation	Dialogue templates	Same and emotion detection algorithm
S3: I vs. We	Dialogue template wording	Same
S4: Uncertain/Verification	Certainty based on correctness of contribution	Certainty based on confidence of algorithms powering code, test cases, and feedback
S5: Social Presence	Active listening via wizard	Active listening via ML monitoring microphone
<b>Technical Skills</b>		
T1: Write/Feedback Tests	Diff from past solution dataset	Same and generate tests from comments or code coverage. Further research needed.
T2: Write/Feedback Code	Diff from past solution dataset	Same and language models such as GPT-3. Further research needed.
T3: Guidance	Based on passing test cases or requirement completion	Same, but further research for matching test cases/code to requirements
T4: Creativity Support	Idea Garden prompts or provide code structure	Same, limited support for Idea and Develop creativity stages [134, 200]
T5: Feature Location	Intuition	Static and dynamic techniques. Further research needed.
T6: Unnecessary Code	Diff from past solution dataset	Dead code via data flow diagrams
T7: Missing Code	Diff from past solution dataset	As supported by T1-T6 design decisions

The difficulty of implementing each design decision includes: Difficult (peach), Medium (gray), and Easy (white).

**(6) High Rate of Role Exchange:** Both iterations of PairBuddy shared a similar rate of role exchange. When normalized over 40 minute sessions, role exchange occurred on average 6.9 and 9.0 times over pilot and main studies, respectively. This difference is not unexpected; PairBuddy interacted with participants far more in the main study than the pilot study (refer Figure 3 vs. Figure 7). However, in a comparable human-human pair programming study [208], role exchange averaged 6.0 over 40 minute studies (3.0 fewer than the main study). Likely, this is due to humans' superior communication skills; human programmers can discuss ideas directly, while PairBuddy relied on code contributions as the driver to communicate ideas.

## 6 DISCUSSIONS AND FUTURE WORK

Our results indicate that PairBuddy was an effective pair programming partner and was enjoyed by our study participants. Table 9 compares the wizard's implementation of PairBuddy to the potential future implementation of PairBuddy. A majority of our design decisions for the wizard are supported by state-of-the-art research; however, the following describes PairBuddy's future design space and the associated challenges:

**(1) Supporting Method-level Code and Test Case Generation:** Designing a conversational agent for the programming domain is challenging. Unlike other conversational agents, PairBuddy

must uniquely support all stages of software development (i.e., clarifying requirements, discussing ideas, designing solutions, and implementing code). Support for each phase requires different knowledge and a different approach.

Moving forward, we need to research automated test case, code, and explainable feedback generation for PairBuddy. Currently, language models such as GPT-3 [37] can be trained from GitHub [37, 95–97] to generate code, while automated test cases can be generated by tools such as Randoop [183] and EvoSuite [81, 82, 203]. However, these tools must be adapted in order for PairBuddy to explain the decisions made and give appropriate feedback. To directly answer the “why” questions that programmers ask, Ko et al. [129, 130] created the automated tool Whyline for both Alice [55] and Java, but similar research must be conducted in other programming languages and domains to allow agents to more directly support all stages of software development.

**(2) Supporting Diverse Problem Solving Strategies:** An observed advantage of PairBuddy between the pilot and main studies was the increased consideration of humans’ unique problem solving strategies, as evidenced by MS5’s comment, “*I’m gonna approach coding in the future... [by] stepping back and looking for alternate approaches.*” This suggests that integrating problem-solving strategies in intelligent programming environments can help instill effective programming practices, especially for students. In the future, programming environments should integrate problem-solving strategies such as working backwards, divergent thinking, divide and conquer, analogy, generalization, and “sleep on it” [142, 187, 256] enable programmers to make progress on their tasks. Previous research from Idea Garden [115, 116] shows that these problem-solving strategies can be implemented by presenting suggestions via language-independent templates, which are informed by language-dependent information about user tasks and progress.

**(3) Supporting Learning Preferences: Solo or Collaboration:** Research has shown that collaborative learning is more effective than traditional methods (such as lectures) since collaboration allows students to build their own mental models based on the discussion and knowledge transfer that occurs during the problem-solving process [16]. But still, many of our participants preferred to work solo. In the pilot study, some participants only used PairBuddy as a resource for code examples, while in the main study, a few preferred that PairBuddy only observe and identify errors. Still, others may prefer full collaboration to enjoy the benefits of pair programming. To support the diversity of learning preferences among programmers (solo vs. collaborative), we will provide different options for PairBuddy’s behavior including “Navigator Only,” “Driver Only,” or “Full Pair Programmer” modes.

PairBuddy could even support a “Human-Human Mode” for those who prefer human partners. Utilizing many of our current decisions, PairBuddy can act as a group facilitator agent by listening to human conversation and providing feedback [14, 38, 67, 92, 146, 188, 189, 200, 207, 249, 253]. For pair programming, group facilitator agents could help equalize member participation, clarify processes, resolve conflict, encourage idea development [189], or teach various social skill (e.g., active listening [9]).

**(4) Implementing PairBuddy as a Task-Oriented vs. Non-Task-Oriented Agent:** PairBuddy’s design needs to fall between the two common types of conversational agents: task-oriented and non-task-oriented. Task-oriented agents perform a variety of tasks and services for users (e.g., Amazon’s Alexa [243]) and provide concise, direct answers to user queries based on knowledge sources (e.g., Bing QA [254]). Non-task-oriented agents facilitate natural interaction between humans and electronic devices (e.g., Meena [11]). For pair programming, a task-oriented agent can assist user queries, while a non-task-oriented agent can engage in rapport making and off-topic dialogue. Research [51] has even identified methods to integrate both task- and

non-task-oriented agents together. In our human-agent study, we found that 89.28% (3967/4443) of programmers' utterances were task-oriented. The abundance of task-oriented dialogue suggests that as a first step, it may be advantageous to implement PairBuddy as a task-oriented conversational agent that utilizes a pipeline architecture for natural language understanding, dialogue state tracking, dialogue policy, and natural language generation [85]. Regardless, we must collect a large enough dataset of pair programming dialogue to train machine learning models to generate automated responses to programmers' queries.

Machine learning algorithms depend upon the quality of the data collected. Specifically, it is difficult to gather unbiased (e.g., gender, culture) and diverse (e.g., different programming languages and tasks) data. For pair programming tasks, the stylistic differences between human-human and human-agent dialogue have implications toward the choice of training dataset [200]. Furthermore, if the quantity or quality of training data is insufficient, PairBuddy may make inaccurate decisions that would likely frustrate programmers and negatively impact user experience.

As a first step, Robe and Kuttal [200] investigated the feasibility of natural language understanding algorithms on human-human pair programming conversations. They achieved a 57.9% classification accuracy when using shallow machine learning algorithms, motivating the exploration of state-of-the-art deep-learning language models (e.g., BERT [69], GPT-3 [37], and XLNet [267]).

**(5) Designing for Varied Expertise of Programmers:** Our main study explored the similarities and differences between expert (professional) and novice (student) programmers. However, self-efficacy scores and interview questions revealed only marginal differences. For instance, students considered PairBuddy's feelings (5/8) much more than professionals (1/6), suggesting that professionals used PairBuddy more as a tool rather than a partner. However, participants' familiarity in the applicable domain (Java) may have mattered more, as two professionals (MP13 and MP14) had less experience than many students. Particularly, our most experienced professional (MP10) uniquely ignored PairBuddy's requests, opposed negative feedback, and excelled at the task. In regard to trusting PairBuddy, there was a slight discrepancy: 3/6 professionals vs. 5/8 students. While our interviews were not conclusive, individuals with particularly low or high experience with a domain tend to trust computer answers more [49, 266]. Hence, future work must investigate the range of programming experience within both educational and professional settings. Initially, research should focus on educational applications since classroom assignments can be designed for a specific task or programming language.

**(6) Manifesting Anthropomorphic Features of a Programmer:** Although we strive to integrate anthropomorphic features into PairBuddy, we received mixed reactions from our participants. While PairBuddy's voice felt natural, its avatar was less impactful, and many participants even minimized the avatar window. Preferences toward a more robotic or a more human-like PairBuddy were very mixed as well. For the past 20 years, researchers have argued in favor or against including anthropomorphic features in intelligent agents [88, 127]. Therefore, it is still an open-ended question whether PairBuddy should use features such as embodiment or emotional intelligence. In future studies, we will investigate whether anthropomorphic features are appropriate for PairBuddy.

Preferences toward PairBuddy's conversational style varied across many dimensions in our interviews (e.g., human-like vs. robotic, polite vs. casual). Many of these characteristics mirror the high-consideration/high-involvement paradigm [190, 191]. High-consideration is characterized by slower turn-taking, non-imposing speech, longer pauses between turns, and avoiding interruptions (i.e., polite, robotic). Conversely, high-involvement entails a faster rate of speech, more turn-taking, few inter-turn pauses, and frequent initiations of simultaneous speech (i.e., casual, human).



Table 10. How PairBuddy’s Design Decisions Influenced its Conversational Style

	High-Consideration	High-Involvement
Driver	S4: Uncertain/Verification	I5: Redirect Suggestions
Both	F4: Text Chat	F3: Voice
	I3: Timed Feedback	S1: Greeting
Navigator		S2: Motivation
		S5: Social Presence

In the pilot study, many of PairBuddy’s design decisions were high-consideration. However, based on the lessons learned, PairBuddy’s conversational style largely became high-involvement in the main study to more accurately emulate human pair programming behavior. Table 10 lists how PairBuddy’s design decisions influenced its conversational style in different roles. Participants’ opinions toward PairBuddy’s conversational style often contradicted one another, hindering the design of a universally accessible PairBuddy. To accommodate individual differences, PairBuddy’s future design should be malleable, allowing users to tune PairBuddy’s parameters to their liking [20]. One possible solution is to make PairBuddy’s script interchangeable. For example, a high-consideration version of PairBuddy’s script might say, “*I’m 80% certain that there is an error on line 45,*” while a high-involvement version would say, “*I think we might have made a mistake on line 45.*” However, interchangeable scripts only go so far, so in the future, we will explore additional avenues to provide a personalized PairBuddy experience.

Furthermore, dialogues are often multi-modal and involve both verbal and non-verbal inputs [30, 62, 68, 271]. Detecting these communication cues from users may help agents build rapport and increase accuracy, engagement, and empathy with their human partners. The lack of non-verbal cues may limit an agent’s human-like feel and affect its usability. Therefore, in the future, each dialogue template of PairBuddy’s script will be accompanied by non-verbal meta-data (e.g., avatar facial expressions, UI events) to facilitate multi-modal communication and build rapport with its partners. Additionally, non-verbal cues can help determine a programmer’s interruptibility. While supported in social settings [79], further work must adapt interruptibility research for programmers and further synchronize its interruptions with the timing of PairBuddy’s feedback (design decision I3).

## 7 CONCLUSION

In this research, we explore the uncharted territory of interactive pair programming conversational agents through the user-centered prototyping of our agent—PairBuddy. This work makes several contributions, including ones that generalize beyond pair programming:

- Anthropomorphic design space of PairBuddy arose from integrating diverse interface and interaction mechanisms (embodiment, dialogue styles, and agent actions) and programmer characteristics (technical skills and soft skills). These design decisions stem from an integration of novel concepts from our extensive literature review of various domains such as human-computer interaction, software engineering, artificial intelligence, psychology, and education. These design decisions can be utilized for advancing programmers’ interactions in ITS and interactive educational platforms.
- Our focused design, evaluation, and refinement cycles through the use of two Wizard of Oz studies incrementally evolved PairBuddy’s functionality to realize a robust conversational agent for pair programming. This methodological approach can drive the design for



programming conversational agents in other domains of programming including educating children, end-user programmers, and people with disabilities.

- The wizard’s script and study materials for both pilot [5] and main studies [4] are available online for reproducibility by researchers and practitioners.

Our study results showed programmers’ positive attitudes towards using PairBuddy. The results confirm the feasibility of PairBuddy as a programming partner that can significantly advance programmer-computer interactions. PairBuddy has significant potential to change how programming is learned and how programming is done. In the words of one of our participants with 20+ years of experience, “*I think what I learned out of this [study] is that [PairBuddy has] a lot of potential and... it excites me a lot where technology is going.*”

## AUTHOR STATEMENT

This work is not related to any prior or concurrent publications, and its contributions stand on its own.

## REFERENCES

- [1] 2017. ISO/IEC/IEEE international standard - systems and software engineering—vocabulary. *ISO/IEC/IEEE 24765:2017(E)* (2017), 1–541. Retrieved on 17 March, 2022 from <https://doi.org/10.1109/IEEESTD.2017.8016712>
- [2] 2018. ISO/IEC/IEEE international standard - systems and software engineering – developing information for users in an agile environment. *ISO/IEC/IEEE 26515:2018(E)* (2018), 1–32. Retrieved on 17 March, 2022 from <https://doi.org/10.1109/IEEESTD.2018.8584455>
- [3] 2020. GitHub. Retrieved from <http://github.com>.
- [4] 2020. Main Study Supporting Material. Retrieved on 17 March, 2022 from [https://drive.google.com/drive/folders/179OnUEVqPHQjW\\_K38ynUopyL4MDKIG-u?usp=sharing](https://drive.google.com/drive/folders/179OnUEVqPHQjW_K38ynUopyL4MDKIG-u?usp=sharing).
- [5] 2020. Pilot Study Supporting Material. Retrieved on 17 March, 2022 from [https://drive.google.com/drive/folders/1WgINmA\\_iz3iONfpaFNB70oCN25mJ76k?usp=sharing](https://drive.google.com/drive/folders/1WgINmA_iz3iONfpaFNB70oCN25mJ76k?usp=sharing).
- [6] 2020. StackOverflow. Retrieved on 17 March, 2022 from <http://stackoverflow.com>.
- [7] 2021. JetBrains. Retrieved on 17 March, 2022 from <https://www.jetbrains.com/>.
- [8] 2021. Visual Studio. Retrieved on 17 March, 2022 from <https://visualstudio.microsoft.com/>.
- [9] Tahir Abbas, Vassilis-Javed Khan, Ujwal Gadiraju, and Panos Markopoulos. 2020. Trainbot: A conversational interface to train crowd workers for delivering on-demand therapy. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 8, 1 (Oct. 2020), 3–12. Retrieved from <https://ojs.aaai.org/index.php/HCOMP/article/view/7458>.
- [10] Gregory D. Abowd and Alan J. Dix. 1992. Giving undo attention. *Interacting with Computers* 4, 3 (12 1992), 317–342. DOI : [https://doi.org/10.1016/0953-5438\(92\)90021-7](https://doi.org/10.1016/0953-5438(92)90021-7) arXiv:<https://academic.oup.com/iwc/article-pdf/4/3/317/2175174/iwc4-0317.pdf>.
- [11] Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. Towards a human-like open-domain chatbot. arXiv:2001.09977. Retrieved from <https://arxiv.org/abs/2001.09977>.
- [12] Ali Ahmadvand, Jason Ingyu Choi, and Eugene Agichtein. 2019. Contextual dialogue act classification for open-domain conversational agents. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1273–1276.
- [13] M. Ai-Chang, J. Bresina, L. Charest, A. Chase, J. C.-J. Hsu, A. Jonsson, B. Kanefsky, P. Morris, Kanna Rajan, J. Yglesias, B. G. Chafin, W. C. Dias, and P. F. Maldague. 2004. MAPGEN: Mixed-initiative planning and scheduling for the Mars Exploration Rover mission. *IEEE Intelligent Systems* 19, 1 (2004), 8–12. DOI : <https://doi.org/10.1109/MIS.2004.1265878>
- [14] Milam Aiken, Mahesh Vanjani, and James Krosch. 1995. Group decision support systems. *Review of Business* 16, 3 (2020/3/2/ 1995), 38+.
- [15] B. Al-Ani and D. Redmiles. 2009. In strangers we trust? Findings of an empirical study of distributed teams. In *Proceedings of the 2009 4th IEEE International Conference on Global Software Engineering*. 121–130.
- [16] Maryam Alavi. 1994. Computer-mediated collaborative learning: An empirical evaluation. *MIS Quarterly* 18, 2 (1994), 159–174.
- [17] S. Ali, L. C. Briand, H. Hemmati, and R. K. Panesar-Walawege. 2010. A systematic review of the application and empirical investigation of search-based test case generation. *IEEE Transactions on Software Engineering* 36, 6 (2010), 742–762.

- [18] Teresa M. Amabile and Michael G. Pratt. 2016. The dynamic componential model of creativity and innovation in organizations: Making progress, making meaning. *Research in Organizational Behavior* 36 (2016), 157–183. DOI : <https://doi.org/10.1016/j.riob.2016.10.001>
- [19] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N. Bennett, Kori Inkpen, Jaime Teevan, Ruth Kikin-Gil, and Eric Horvitz. 2019. Guidelines for human-AI interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, Article 3, 13 pages. DOI : <https://doi.org/10.1145/3290605.3300233>
- [20] Ofer Arazy, Oded Nov, and Nanda Kumar. 2015. Personalization: UI personalization, theoretical grounding in HCI and design research. *AIIS Transactions on Human-Computer Interaction* 7, 2 (2015), 43–69.
- [21] Michael Armstrong. 2012. *Armstrong's Handbook of Reward Management Practice: Improving Performance Through Reward* (12 ed.). Kogan Page Publishers.
- [22] Zahra Ashktorab, Mohit Jain, Q. Vera Liao, and Justin D. Weisz. 2019. Resilient chatbots: Repair strategy preferences for conversational breakdowns. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, Article 254, 12 pages. DOI : <https://doi.org/10.1145/3290605.3300484>
- [23] D. C. Hoaglin, B. A. Kitchenham, S. L. Pfleeger, and J. Rosenberg. 2002. "Preliminary guidelines for empirical research in software engineering". In *Proceedings of the IEEE Transactions on Software Engineering*, Vol. 28. 721–734.
- [24] Claudio Barra and Broderick Crawford. 2007. Fostering creativity thinking in agile software development. In *Proceedings of the Symposium of the Austrian HCI and Usability Engineering Group*, Vol. 4799. 415–426. DOI : [https://doi.org/10.1007/978-3-540-76805-0\\_37](https://doi.org/10.1007/978-3-540-76805-0_37)
- [25] Amy L. Baylor and Soyoung Kim. 2009. Designing nonverbal communication for pedagogical agents: When less is more. *Computers in Human Behavior* 25, 2 (2009), 450–457. DOI : <https://doi.org/10.1016/j.chb.2008.10.008>
- [26] Tara Behrend, Steven Toaddy, Lori Foster Thompson, and David J. Sharek. 2012. The effects of avatar appearance on interviewer ratings in virtual employment interviews. *Computers in Human Behavior* 28, 6 (2012), 2128–2133. DOI : <https://doi.org/10.1016/j.chb.2012.06.017>
- [27] A. Belshee. 2005. Promiscuous pairing and beginner's mind: Embrace inexperience. In *Proceedings of the Agile Development Conference*. 125–131. DOI : <https://doi.org/10.1109/ADC.2005.37>
- [28] Gary Bente, Sabine Rüggenberg, Nicole C. Krämer, and Felix Eschenburg. 2008. Avatar-mediated networking: Increasing social presence and interpersonal trust in net-based collaborations. *Human Communication Research* 34, 2 (April 2008), 287–318. DOI : <https://doi.org/10.1111/j.1468-2958.2008.00322.x> arXiv:<https://academic.oup.com/hcr/article-pdf/34/2/287/22325251/jhumcom0287.pdf>
- [29] Timothy Bickmore and Justine Cassell. 2001. Relational agents: A model and implementation of building user trust. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 396–403.
- [30] Dan Bohus, Chit W. Saw, and Eric Horvitz. 2014. Directions robot: In-the-wild experiences and lessons learned. In *Proceedings of the 2014 International Conference on Autonomous Agents and Multi-Agent Systems*. 637–644.
- [31] Yvonne Brackbill, William E. Boblitt, Douglas Davlin, and John E. Wagner. 1963. Amplitude of response and the delay-retention effect. *Journal of Experimental Psychology* 66, 1 (1963), 57.
- [32] Jay Bradley, David Benyon, Oli Mival, and Nick Webb. 2010. Wizard of Oz experiments and companion dialogues. In *Proceedings of the 24th BCS Interaction Specialist Group Conference*. British Computer Society, 117–123.
- [33] Sheryl Brahnam and Antonella De Angeli. 2012. Gender affordances of conversational agents. *Interacting with Computers* 24, 3 (April 2012), 139–153. DOI : <https://doi.org/10.1016/j.intcom.2012.05.001> arXiv:<https://academic.oup.com/iwc/article-pdf/24/3/139/2027399/iwc24-0139.pdf>
- [34] Petter Bae Brandtzaeg and Asbjørn Følstad. 2017. Why people use chatbots. In *Proceedings of the International Conference on Internet Science*. Springer, 377–392.
- [35] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative Research in Psychology* 3, 2 (Jan. 2006), 77–101. DOI : <https://doi.org/10.1191/1478088706qp0630a>
- [36] Tim Brown. 2009. *Change by Design: How Design Thinking Transforms Organizations and Inspires Innovation*. Harper-Business.
- [37] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.
- [38] Tung X. Bui. 1987. *Co-op: A Group Decision Support System for Cooperative Multiple Criteria Group Decision Making*. Springer, Berlin.
- [39] Margaret Burnett, Anicia Peters, Charles Hill, and Noha Elarief. 2016. Finding gender-inclusiveness software issues with GenderMag: A field investigation. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*. ACM, 2586–2598.
- [40] Margaret Burnett, Simone Stumpf, Jamie Macbeth, Stephann Makri, Laura Beckwith, Irwin Kwan, Anicia Peters, and William Jernigan. 2016. GenderMag: A method for evaluating software's gender inclusiveness. *Interacting with Computers* 28, 6 (Oct. 2016), 760–787.

- [41] Ramón Burri. 2018. *Improving User Trust Towards Conversational Chatbot Interfaces with Voice Output*. Master's Thesis. KTH, School of Electrical Engineering and Computer Science (EECS).
- [42] ROSE Carolyn. 2007. Tools for authoring a dialogue agent that participates in learning studies. *Artificial Intelligence in Education: Building Technology Rich Learning Contexts That Work* 158 (2007), 43.
- [43] Justine Cassell, Yukiko I. Nakano, Timothy W. Bickmore, Candace L. Sidner, and Charles Rich. 2001. Non-verbal cues for discourse structure. In *Proceedings of the 39th Annual Meeting of the Association for Computational Linguistics*. 114–123.
- [44] Mehmet Celepkolu and Kristy Elizabeth Boyer. 2018. Thematic analysis of students' reflections on pair programming in CS1. In *Proceedings of the 49th ACM Technical Symposium on Computer Science Education*. ACM, New York, NY, 771–776. DOI: <https://doi.org/10.1145/3159450.3159516>
- [45] Christopher P. Cerasoli, Jessica M. Nicklin, and Michael T. Ford. 2014. Intrinsic motivation and extrinsic incentives jointly predict performance: A 40-year meta-analysis. *Psychological Bulletin* 140, 4 (2014), 980.
- [46] Hyun Jin Cha, Yong Se Kim, Seon Hee Park, Tae Bok Yoon, Young Mo Jung, and Jee-Hyong Lee. 2006. Learning styles diagnosis based on user interface behaviors for the customization of learning interfaces in an intelligent tutoring system. In *Proceedings of the International Conference on Intelligent Tutoring Systems*. Springer, 513–524.
- [47] Tak-Wai Chan. 1996. Learning companion systems, social learning systems, and the global social learning club. *Journal of Artificial Intelligence in Education* 7, 2 (1996), 125. Retrieved from <https://www.learntechlib.org/p/82394>.
- [48] Gary Charness and Uri Gneezy. 2012. Strong evidence for gender differences in risk taking. *Journal of Economic Behavior & Organization* 83, 1 (2012), 50–58.
- [49] J. Y. C. Chen and M. J. Barnes. 2014. Human-agent teaming for multirobot control: A review of human factors issues. *IEEE Transactions on Human-Machine Systems* 44, 1 (2014), 13–29.
- [50] K. S. Choi. 2013. Evaluating gender significance within a pair programming context. In *Proceedings of the 2013 46th Hawaii International Conference on System Sciences*. 4817–4825.
- [51] Tai-Liang Chou and Yu-Ling Hsueh. 2019. A task-oriented chatbot based on LSTM and reinforcement learning. In *Proceedings of the 2019 3rd International Conference on Natural Language Processing and Information Retrieval*. ACM, New York, NY, 87–91. DOI: <https://doi.org/10.1145/3342827.3342844>
- [52] Alistair Cockburn and Laurie Williams. 2001. *Extreme Programming Examined*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, Chapter The Costs and Benefits of Pair Programming, 223–243. Retrieved from <http://dl.acm.org/citation.cfm?id=377517.377531>.
- [53] Michelle Cohn, Chun-Yen Chen, and Zhou Yu. 2019. A large-scale user study of an alexa prize chatbot: Effect of TTS dynamism on perceived quality of social dialog. In *Proceedings of the 20th Annual SIGdial Meeting on Discourse and Dialogue*. 293–306.
- [54] Juan Manuel Contreras, Mahzarin R. Banaji, and Jason P. Mitchell. 2013. Multivoxel patterns in fusiform face area differentiate faces by sex and race. *PLoS One* 8, 7 (2013), e69684.
- [55] Stephen Cooper, Wanda Dann, and Randy Pausch. 2000. Alice: A 3-D tool for introductory programming concepts. *Journal of Computing Sciences in Colleges* 15, 5 (April 2000), 107–116.
- [56] Tyne Crow, Andrew Luxton-Reilly, and Burkhard Wünsche. 2018. Intelligent tutoring systems for programming education: A systematic review. In *Proceedings of the 20th Australasian Computing Education Conference*. 53–62. DOI: <https://doi.org/10.1145/3160489.3160492>
- [57] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. 1991. Efficiently computing static single assignment form and the control dependence graph. *ACM Transactions on Programming Languages and Systems* 13, 4 (Oct. 1991), 451–490. DOI: <https://doi.org/10.1145/115372.115320>
- [58] Fabio Q. B. da Silva, Catarina Costa, A. Cesar C. Franca, and Rafael Prikladinicki. 2010. Challenges and solutions in distributed software development project management: A systematic literature review. In *Proceedings of the 2010 5th IEEE International Conference on Global Software Engineering*. IEEE, 87–96.
- [59] Nils Dahlbäck, Arne Jönsson, and Lars Ahrenberg. 1993. Wizard of Oz studies: Why and how. In *Proceedings of the 1st International Conference on Intelligent User Interfaces*. 193–200.
- [60] M. Day, M. R. Penumala, and J. Gonzalez-Sanchez. 2019. Annete: An intelligent tutoring companion embedded into the eclipse IDE. In *Proceedings of the 2019 IEEE 1st International Conference on Cognitive Machine Intelligence*. 71–80.
- [61] Claudio León de la Barra and Broderick Crawford. 2007. Fostering creativity thinking in agile software development. In *Proceedings of the HCI and Usability for Medicine and Health Care*. Andreas Holzinger (Ed.). Springer Berlin, 415–426.
- [62] Harm De Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron Courville. 2017. Guess-what?! visual object discovery through multi-modal dialogue. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5503–5512.
- [63] Hans Dechert and Manfred Raupach. 1987. Conversational style. *Psycholinguistic Models of Production* (1987), 251–267.

- [64] Edward L. Deci, Anja H. Olafsen, and Richard M. Ryan. 2017. Self-determination theory in work organizations: The state of a science. *Annual Review of Organizational Psychology and Organizational Behavior* 4, 1 (2017), 19–43.
- [65] Doris M. Dehn and Susanne van Mulken. 2000. The impact of animated interface agents: A review of empirical research. *International Journal of Human-Computer Studies* 52, 1 (Jan. 2000), 1–22. DOI : <https://doi.org/10.1006/ijhc.1999.0325>
- [66] Tom DeMarco and Timothy Lister. 1987. *Peopleware: Productive Projects and Teams*. Dorset House Publishing Co., Inc., New York, NY.
- [67] Gerardine DeSanctis and R. Brent Gallupe. 1987. A foundation for the study of group decision support systems. *Management Science* 33, 5 (May 1987), 589–609.
- [68] David DeVault, Ron Artstein, Grace Benn, Teresa Dey, Ed Fast, Alesia Gainer, Kallirroi Georgila, Jon Gratch, Arno Hartholt, Margaux Lhommet, Gale Lucas, Stacy Marsella, Fabrizio Morbini, Angela Nazarian, Stefan Scherer, Giota Stratou, Apar Suri, David Traum, Rachel Wood, Yuyu Xu, Albert Rizzo, and Louis-Philippe Morency. 2014. SimSensei kiosk: A virtual human interviewer for healthcare decision support. In *Proceedings of the 2014 International Conference on Autonomous Agents and Multi-agent Systems*. 1061–1068.
- [69] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. Retrieved from <https://arxiv.org/abs/1810.04805>.
- [70] Eclipse 2019. Eclipse IDE. Retrieved on 17 March, 2022 from <https://www.eclipse.org/>.
- [71] Eclipse 2020. Junit. Retrieved on 17 March, 2022 from <https://junit.org/junit5/>.
- [72] Stephan Salinger, Christopher Oezbek, Karl Beecher, and Julia Schenk. 2010. *Saros: An Eclipse Plug-in for Distributed Party Programming (CHASE'10)*. Association for Computing Machinery, New York, NY, USA, 48–55. <https://doi.org/10.1145/1833310.1833319>
- [73] Berland Edelman and Inc. 2010. *Creativity and Education: Why it Matters*. Retrieved September 18th, 2019 from [http://www.adobe.com/aboutadobe/pressroom/pdfs/Adobe\\_Creativity\\_and\\_Education\\_Why\\_It\\_Matters\\_study.pdf](http://www.adobe.com/aboutadobe/pressroom/pdfs/Adobe_Creativity_and_Education_Why_It_Matters_study.pdf).
- [74] R. Elghondakly, S. Moussa, and N. Badr. 2015. Waterfall and agile requirements-based model for automated test cases generation. In *Proceedings of the 2015 IEEE 7th International Conference on Intelligent Computing and Information Systems*. 607–612.
- [75] Martha Evens and Joel Michael. 2006. *One-on-one Tutoring By Humans and Computers*. Psychology Press.
- [76] Stef van der Struijk, Hung-Hsuan Huang, Maryam Sadat Mirzaei, and Toyoaki Nishida. 2018. FACSvatar: An Open Source Modular Framework for Real-Time FACS Based Facial Animation. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents (Sydney, NSW, Australia) (IVA'18)*. Association for Computing Machinery, New York, NY, USA, 159–164. <https://doi.org/10.1145/3267851.3267918>
- [77] Raoul Festante. 2007. *An Introduction to the Theory of Gender-neutral Language*. BoD–Books on Demand.
- [78] Carmen Fischer, Charlotte P. Malycha, and Ernestine Schafmann. 2019. The influence of intrinsic motivation and synergistic extrinsic motivators on creativity and innovation. *Frontiers in Psychology* 10 (2019), 137. DOI : <https://doi.org/10.3389/fpsyg.2019.00137>
- [79] James Fogarty, Scott E. Hudson, Christopher G. Atkeson, Daniel Avrahami, Jodi Forlizzi, Sara Kiesler, Johnny C. Lee, and Jie Yang. 2005. Predicting human interruptibility with sensors. *ACM Transactions on Computer-Human Interaction* 12, 1 (March 2005), 119–146. DOI : <https://doi.org/10.1145/1057237.1057243>
- [80] Cyrus K. Foroughi, Nicole E. Werner, Erik T. Nelson, and Deborah A. Boehm-Davis. 2014. Do interruptions affect quality of work? *Human Factors* 56, 7 (2014), 1262–1271. DOI : <https://doi.org/10.1177/0018720814531786> PMID: 25490806.
- [81] Gordon Fraser and Andrea Arcuri. 2012. Whole test suite generation. *IEEE Transactions on Software Engineering* 39, 2 (2012), 276–291.
- [82] Gordon Fraser and Andrea Arcuri. 2013. Evosuite: On the challenges of test case generation in the real world. In *Proceedings of the 2013 IEEE 6th International Conference on Software Testing, Verification and Validation*. IEEE, 362–369.
- [83] Freelancing Platform [n. d.]. Upwork Inc.
- [84] Hans Gallis and Erik Arisholm. 2002. A transition from partner programming to pair programming-an industrial case study. In *Proceedings of the Workshop: "Pair Programming Installed" at Object-Oriented Programming, Systems, Languages and Applications*.
- [85] Jianfeng Gao, Michel Galley, and Lihong Li. 2018. Neural approaches to conversational ai. In *Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. 1371–1374.
- [86] Tom Geller. 2008. Overcoming the uncanny valley. *IEEE Computer Graphics and Applications* 28, 4 (2008), 11–17.
- [87] GenderMag 2019. GenderMag. Retrieved on 17 March, 2022 from <http://gendermag.org/>.
- [88] Stella George. 2019. From sex and therapy bots to virtual assistants and tutors: How emotional should artificially intelligent agents be? In *Proceedings of the 1st International Conference on Conversational User Interfaces*. ACM, New York, NY, Article 19, 3 pages. DOI : <https://doi.org/10.1145/3342775.3342807>

- [89] Alex Gerdes, Bastiaan Heeren, Johan Jeuring, and L. Thomas van Binsbergen. 2016. Ask-Elle: An adaptable programming tutor for Haskell giving automated feedback. *International Journal of Artificial Intelligence in Education* 27, 1 (Feb. 2016), 65–100. DOI: <https://doi.org/10.1007/s40593-015-0080-x>
- [90] Barney G. Glaser and Anselm L. Strauss. 1967. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. Aldine de Gruyter, New York, NY.
- [91] M. Gonzalez-Franco, A. Steed, S. Hoogendyk, and E. Ofek. 2020. Using facial animation to increase the enface-ment illusion and avatar self-identification. *IEEE Transactions on Visualization and Computer Graphics* 26, 5 (2020), 2023–2029.
- [92] Paul Gray. 1987. Group decision support systems. *Decision Support Systems* 3, 3 (Sept. 1987), 233–242.
- [93] Paul Green and Lisa Wei-Haas. 1985. *The Wizard of Oz: A Tool for Rapid Development of User Interfaces*. Technical Report UU-CS-2015-019. Utrecht University, Utrecht, The Netherlands.
- [94] Pamela Grimm. 2010. Social Desirability Bias. John Wiley Sons, Ltd. <https://doi.org/10.1002/9781444316568.wiem02057> arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781444316568.wiem02057>
- [95] GTP-3 2020. Tweet of Code Generated from GPT-3. Retrieved on 17 March, 2022 from <https://twitter.com/sharifshameem/status/1282676454690451457>.
- [96] GTP-3 2020. YouTube Video Showing Generation of Code by GPT3. Retrieved on 17 March, 2022 from <https://www.youtube.com/watch?v=utuz7wBGjKM>.
- [97] GTP-3 2020. YouTube Video Showing Generation of Code by GPT3. Retrieved 24 June, 2020 from GPT3:<https://www.youtube.com/watch?v=y5-wzgIySb4>.
- [98] GTTS 2019. Google Text-to-speech Python Library. Retrieved on 17 March, 2022 from <https://github.com/pndurette/gTTS>.
- [99] Keun-Woo Han, EunKyoung Lee, and Youngjun Lee. 2010. The impact of a peer-learning agent based on pair programming in a programming course. *IEEE Transactions on Education* 53, 2 (June 2010), 318–327. DOI: <https://doi.org/10.1109/TE.2009.2019121>
- [100] Qinghong Han, Yuxian Meng, Fei Wu, and Jiwei Li. 2020. Non-autoregressive neural dialogue generation. Retrieved from <https://arxiv.org/abs/2002.04250>.
- [101] Brian F. Hanks. 2004. Distributed pair programming: An empirical study. In *Proceedings of the Extreme Programming and Agile Methods - XP/Agile Universe 2004*. Carmen Zannier, Hakan Erdogmus, and Lowell Lindstrom (Eds.). Springer Berlin, 81–91.
- [102] Benjamin Hardin and Michael A. Goodrich. 2009. On using mixed-initiative control: A perspective for managing large-scale robotic teams. In *Proceedings of the 2009 4th ACM/IEEE International Conference on Human-Robot Interaction*. ACM, New York, NY, 165–172. DOI: <https://doi.org/10.1145/1514095.1514126>
- [103] Dai Hasegawa, Justine Cassell, and Kenji Araki. 2010. The role of embodiment and perspective in direction-giving systems. In *Proceedings of the 2010 AAAI Fall Symposium Series*.
- [104] Khaled Hassanein and Milena Head. 2007. Manipulating perceived social presence through the web interface and its impact on attitude towards online shopping. *International Journal of Human-Computer Studies* 65, 8 (2007), 689–708. DOI: <https://doi.org/10.1016/j.ijhcs.2006.11.018>
- [105] John Hattie. 1999. Influences on student learning. *Inaugural Lecture Given on August 2, 1999* (1999), 21.
- [106] John Hattie and Helen Timperley. 2007. The power of feedback. *Review of Educational Research* 77, 1 (2007), 81–112.
- [107] Renate Häußelschmid, Max von Bülow, Bastian Pflöging, and Andreas Butz. 2017. SupportingTrust in autonomous driving. In *Proceedings of the 22nd International Conference on Intelligent User Interfaces*. ACM, New York, NY, 319–329. DOI: <https://doi.org/10.1145/3025171.3025198>
- [108] Charles Hill. 2017. The sum of its parts: Investigating the component pieces of GenderMag. Retrieved on 17 March, 2022 from <http://hdl.handle.net/1957/61887>.
- [109] Xinting Huang, Jianzhong Qi, Yu Sun, and Rui Zhang. 2020. Semi-Supervised Dialogue Policy Learning via Stochastic Reward Estimation. In *ACL*. 660–670. <https://doi.org/10.18653/v1/2020.acl-main.62>
- [110] Niklas Humble and Peter Mozelius. 2019. Teacher-supported AI or AI-supported teachers? In *European Conference on the Impact of Artificial Intelligence and Robotics 2019 (ECLAIR'19)*, Oxford, UK, Vol. 1. Academic Conferences and Publishing International Limited, 157–164.
- [111] Lilly Irani. 2004. Understanding gender and confidence in CS course culture. *ACM SIGCSE Bulletin* 36, 1 (2004), 195–199. ACM, New York, NY. DOI: <https://doi.org/10.1145/971300.971371>
- [112] Scott G. Isaksen and Donald J. Treffinger. 2004. Celebrating 50 years of reflective practice: Versions of creative problem solving. *The Journal of Creative Behavior* 38, 2 (June 2004), 75–101.
- [113] Mohit Jain, Pratyush Kumar, Ishita Bhansali, Q. Vera Liao, Khai Truong, and Shwetak Patel. 2018. FarmChat: A conversational agent to answer farmer queries. *Proceedings of the ACM on Interactive Mobile Wearable Ubiquitous Technologies* 2, 4 (Dec. 2018), Article 170, 22 pages. DOI: <https://doi.org/10.1145/3287048>
- [114] Mohit Jain, Pratyush Kumar, Ramachandra Kota, and Shwetak N. Patel. 2018. Evaluating and informing the design of chatbots. In *Proceedings of the 2018 Designing Interactive Systems Conference*. 895–906.



- [115] William Jernigan, Amber Horvath, Michael Lee, Margaret Burnett, Taylor Cuiity, Sandeep Kuttal, Anicia Peters, Irwin Kwan, Faezeh Bahmani, Andrew Ko, and Christopher Mendez. 2017. General principles for a generalized idea garden. *Journal of Visual Languages & Computing* 39, C (May 2017), 51–65. DOI : <https://doi.org/10.1016/j.jvlc.2017.04.005>
- [116] Will Jernigan, Amber Horvath, Michael Lee, Margaret Burnett, Cuiity Taylor, Sandeep Kuttal, Anicia Peters, Irwin Kwan, Faezeh Bahmani, and Andrew Ko. 2015. A principled evaluation for a principled idea garden. Retrieved on 17 March, 2022 from <https://doi.org/10.1109/VLHCC.2015.7357222>
- [117] Danielle Jones and Scott Fleming. 2013. What use is a backseat driver? A qualitative investigation of pair programming. In *Proceedings of IEEE Symposium on Visual Languages and Human-Centric Computing*, 103–110. DOI : <https://doi.org/10.1109/VLHCC.2013.6645252>
- [118] Ankur Joshi, Saket Kale, Satish Chandel, and D. Kumar Pal. 2015. Likert scale: Explored and explained. *Current Journal of Applied Science and Technology* 7, 4 (2015), 396–403.
- [119] Ewa Kacewicz, James W. Pennebaker, Matthew Davis, Moongee Jeon, and Arthur C. Graesser. 2014. Pronoun use reflects standings in social hierarchies. *Journal of Language and Social Psychology* 33, 2 (2014), 125–143. DOI : <https://doi.org/10.1177/0261927X13502654> arXiv:<https://doi.org/10.1177/0261927X13502654>
- [120] Peter H. Kahn, Nathan G. Freier, Takayuki Kanda, Hiroshi Ishiguro, Jolina H. Ruckert, Rachel L. Severson, and Shaun K. Kane. 2008. Design patterns for sociality in human-robot interaction. In *Proceedings of the 3rd ACM/IEEE International Conference on Human Robot Interaction* ACM, New York, NY, 97–104. DOI : <https://doi.org/10.1145/1349822.1349836>
- [121] Neha Katira, Laurie Williams, Laurie Williams, Eric Wiebe, Carol Miller, Suzanne Balik, and Ed Gehringer. 2004. On understanding compatibility of student pair programmers. *SIGCSE Bull.* 36, 1 (March 2004), 7–11. DOI : <https://doi.org/10.1145/1028174.971307>
- [122] R. K. Kavitha and M. S. Irfan Ahmed. 2013. Knowledge sharing through pair programming in learning environments: An empirical study. *Education and Information Technologies* 20, 2 (Oct. 2013), 319–333.
- [123] Greg P. Kearsley. 1987. *Artificial Intelligence and Instruction: Applications and Methods*. Addison-Wesley Longman Publishing Co., Inc.
- [124] Iman Keivanloo, Juergen Rilling, and Ying Zou. 2014. Spotting working code examples. In *Proceedings of the 36th International Conference on Software Engineering*. ACM, New York, NY, 664–675. DOI : <https://doi.org/10.1145/2568225.2568292>
- [125] Kisub Kim, Dongsun Kim, Tegawendé F. Bissyandé, Eunjong Choi, Li Li, Jacques Klein, and Yves Le Traon. 2018. FaCoY: A code-to-code search engine. In *Proceedings of the 40th International Conference on Software Engineering*. ACM, New York, NY, 946–957. DOI : <https://doi.org/10.1145/3180155.3180187>
- [126] Sungdong Kim, Sohee Yang, Gyuwan Kim, and Sang-Woo Lee. 2020. Efficient Dialogue State Tracking by Selectively Overwriting Memory. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 567–582. <https://doi.org/10.18653/v1/2020.acl-main.53>
- [127] Lorenz Cuno Klopfenstein, Saverio Delpriori, Silvia Malatini, and Alessandro Bogliolo. 2017. The rise of bots: A survey of conversational interfaces, patterns, and paradigms. In *Proceedings of the 2017 Conference on Designing Interactive Systems*. ACM, New York, NY, 555–565. DOI : <https://doi.org/10.1145/3064663.3064672>
- [128] Avraham N. Kluger and Angelo DeNisi. 1996. The effects of feedback interventions on performance: A historical review, a meta-analysis, and a preliminary feedback intervention theory. *Psychological Bulletin* 119, 2 (1996), 254.
- [129] Andrew J. Ko and Brad A. Myers. 2004. Designing the whyline: A debugging interface for asking questions about program behavior. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, 151–158. DOI : <https://doi.org/10.1145/985692.985712>
- [130] Andrew J. Ko and Brad A. Myers. 2009. Finding causes of program output with the Java whyline. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, 1569–1578. DOI : <https://doi.org/10.1145/1518701.1518942>
- [131] James A. Kulik and J. D. Fletcher. 2016. Effectiveness of intelligent tutoring systems: A meta-analytic review. *Review of Educational Research* 86, 1 (2016), 42–78. DOI : <https://doi.org/10.3102/0034654315581420>
- [132] James A. Kulik and Chen-Lin C. Kulik. 1988. Timing of feedback and verbal learning. *Review of Educational Research* 58, 1 (1988), 79–97. DOI : <https://doi.org/10.3102/00346543058001079> arXiv:<https://doi.org/10.3102/00346543058001079>
- [133] Sandeep Kuttal, Kevin Gerstner, and Alexandra Bejarano. 2019. Remote pair-programming in online CS education: Investigating through a gender lens. In *Proceedings of the IEEE Symposium on Visual Languages & Human-Centric Computing*.
- [134] S. K. Kuttal, J. Myers, S. Gurka, D. Magar, D. Piorkowski, and R. Bellamy. 2020. Towards designing conversational agents for pair programming: Accounting for creativity strategies and conversational styles. In *Proceedings of the 2020 IEEE Symposium on Visual Languages and Human-Centric Computing*. 1–11.
- [135] Danielle L. Jones and Scott D. Fleming. 2013. What use is a backseat driver? A qualitative investigation of pair programming. In *Proceedings of IEEE Symposium on Visual Languages and Human-Centric Computing*. 103–110.



- [136] Shuvendu K. Lahiri, Chris Hawblitzel, Ming Kawaguchi, and Henrique Rebêlo. 2012. Symdiff: A language-agnostic semantic diff tool for imperative programs. In *Proceedings of the International Conference on Computer Aided Verification*. Springer, 712–717.
- [137] Thomas K. Landauer. 1987. Psychology as a mother of invention. *ACM SIGCHI Bulletin* 18, 4 (1987), 333–335.
- [138] H. Lane and Kurt Vanlehn. 2005. Teaching the tacit knowledge of programming to novices with natural language tutoring. *Computer Science Education* 15, 3 (Sep. 2005), 183–201. DOI: <https://doi.org/10.1080/08993400500224286>
- [139] Hung Le, Richard Socher, and Steven C. H. Hoi. 2020. Non-autoregressive dialog state tracking. In *Proceedings of the International Conference on Learning Representations*. Retrieved from [https://openreview.net/forum?id=H1e\\_cC4twS](https://openreview.net/forum?id=H1e_cC4twS).
- [140] Nguyen-Thinh Le. 2016. A classification of adaptive feedback in educational systems for programming. *Systems* 4, 2 (May 2016), 22. DOI: <https://doi.org/10.3390/systems4020022>
- [141] Nguyen-Thinh Le, Sven Strickroth, Sebastian Gross, and Niels Pinkwart. 2013. A review of AI-supported tutoring approaches for learning programming. *Advanced Computational Methods for Knowledge Engineering* 479 (Jan. 2013), 267–279. DOI: [https://doi.org/10.1007/978-3-319-00293-4\\_20](https://doi.org/10.1007/978-3-319-00293-4_20)
- [142] Marvin Levine. 1988. *Effective Problem Solving*. Prentice Hall.
- [143] Clayton Lewis. 1982. *Using the “thinking-aloud” Method in Cognitive Interface Design*. IBM TJ. Watson Research Center, Yorktown Heights, NY.
- [144] Shaofeng Li. 2010. The effectiveness of corrective feedback in SLA: A meta-analysis. *Language Learning* 60, 2 (2010), 309–365.
- [145] Toby Jia-Jun Li, Jingya Chen, Tom Mitchell, and Brad Myers. 2020. Towards Effective Human-AI Collaboration in GUI-Based Interactive Task Learning Agents. *CHI 2020 Workshop on Artificial Intelligence for HCI: A Modern Approach (AI4HCI)*. <https://doi.org/arXiv:2003.02622>
- [146] Moez Limayem, Probir Banerjee, and Louis Ma. 2006. Impact of GDSS: Opening the black box. *Decision Support Systems* 42, 2 (Nov. 2006), 945–957.
- [147] Dapeng Liu, Andrian Marcus, Denys Poshyvanyk, and Vaclav Rajlich. 2007. Feature location via information retrieval based filtering of a single scenario execution trace. In *Proceedings of the ACM/IEEE International Conference on Automated Software Engineering*, 234–243. DOI: <https://doi.org/10.1145/1321631.1321667>
- [148] Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4487–4496.
- [149] Zhiqiang Liu and Dieter J. Schonwetter. 2004. Teaching creativity in engineering. *International Journal of Engineering Education* 20, 5 (2004), 801–808.
- [150] Ju Long. 2009. Open source software development experiences on the students’ resumes: Do they count?-insights from the employers’ perspectives. *Journal of Information Technology Education: Research* 8, 1 (2009), 229–242.
- [151] Irene Lopatovska and Harriet Williams. 2018. Personification of the amazon alexa: BFF or a mindless companion. In *Proceedings of the 2018 Conference on Human Information Interaction & Retrieval*. ACM, New York, NY, 265–268. DOI: <https://doi.org/10.1145/3176349.3176868>
- [152] Ewa Luger and Abigail Sellen. 2016. “Like having a really bad PA” the gulf between user expectation and experience of conversational agents. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, 5286–5297.
- [153] Wenting Ma, Olusola O. Adesope, John C. Nesbit, and Qing Liu. 2014. Intelligent tutoring systems and learning outcomes: A meta-analysis. *Journal of Educational Psychology* 106, 4 (2014), 901.
- [154] Stephen Makonin, Daniel McVeigh, Wolfgang Stuerzlinger, Khoa Tran, and Fred Popowich. 2016. Mixed-initiative for big data: The intersection of human + visual analytics + prediction. In *Proceedings of the 2016 49th Hawaii International Conference on System Sciences*, 1427–1436. DOI: <https://doi.org/10.1109/HICSS.2016.181>
- [155] Divine Maloney. 2018. Mitigating negative effects of immersive virtual avatars on racial bias. In *Proceedings of the 2018 Annual Symposium on Computer-Human Interaction in Play Companion Extended Abstracts*, 39–43.
- [156] A. Marcus, V. Rajlich, J. Buchta, M. Petrenko, and A. Sergeyev. 2005. Static techniques for concept location in object-oriented code. In *Proceedings of the 13th International Workshop on Program Comprehension*, 33–42.
- [157] Jennifer Marlow and Laura Dabbish. 2013. Activity traces and signals in software developer recruitment and hiring. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work*. ACM, 145–156.
- [158] Akane Matsushima, Natsuki Oka, Chie Fukada, and Kazuaki Tanaka. 2019. Understanding dialogue acts by bayesian inference and reinforcement learning. In *Proceedings of the 7th International Conference on Human-Agent Interaction*. ACM, New York, NY, 262–264. DOI: <https://doi.org/10.1145/3349537.3352786>
- [159] Charlie McDowell, Linda Werner, Heather Bullock, and Julian Fernald. 2002. The effects of pair-programming on performance in an introductory programming course. In *Proceedings of the 33rd SIGCSE Technical Symposium on Computer Science Education*. ACM, New York, NY, 38–42. DOI: <https://doi.org/10.1145/563340.563353>
- [160] Charlie McDowell, Linda Werner, Heather E. Bullock, and Julian Fernald. 2003. The impact of pair programming on student performance, perception and persistence. In *Proceedings of the 25th International Conference on Software*

- Engineering. IEEE Computer Society, Washington, DC, 602–607. Retrieved from <http://dl.acm.org/citation.cfm?id=776816.776899>.
- [161] Phil McMinn. 2004. Search-based software test data generation: A survey: Research Articles. *Software Testing, Verification and Reliability* 14, 2 (June 2004), 105–156. DOI : <https://doi.org/10.1002/stvr.294>
  - [162] Paola Medel and Vahab Pournaghshband. 2017. Eliminating gender bias in computer science education materials. In *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education*. ACM, New York, NY, 411–416. DOI : <https://doi.org/10.1145/3017680.3017794>
  - [163] Meiliana, Irwandhi Septian, Ricky Setiawan Alianto, Daniel, and Ford Lumban Gaol. 2017. Automated test case generation from UML activity diagram and sequence diagram using depth first search algorithm. *Procedia Computer Science* 116 (2017), 629–637. DOI : <https://doi.org/10.1016/j.procs.2017.10.029>
  - [164] Grigori Melnik and Frank Maurer. 2002. Perceptions of agile practices: A student survey. In *Proceedings of the Conference on Extreme Programming and Agile Methods*. Springer, 241–250.
  - [165] Suejb Memeti and Sabri Pllana. 2018. PAPA: A parallel programming assistant powered by IBM Watson cognitive computing technology. *Journal of Computational Science* 26 (2018), 275–284. <https://doi.org/10.1016/j.jocs.2018.01.001>
  - [166] Christopher Mendez, Hema Susmita Padala, Zoe Steine-Hanson, Claudia Hilderbrand, Amber Horvath, Charles Hill, Logan Simpson, Nupoor Patil, Anita Sarma, and Margaret Burnett. 2018. Open source barriers to entry, revisited: A sociotechnical perspective. In *Proceedings of the 40th International Conference on Software Engineering*. ACM, 1004–1015.
  - [167] Casey Miller and Kate Swift. 2001. *The Handbook of Nonsexist Writing*. iUniverse.
  - [168] Matheus Monteiro Mariano, Érica F. Souza, André T. Endo, and Nandamudi L. Vijaykumar. 2019. Analyzing graph-based algorithms employed to generate test cases from finite state machines. In *2019 IEEE Latin American Test Symposium (LATS)*. 1–6. <https://doi.org/10.1109/LATW.2019.8704603>
  - [169] Dana Movshovitz-Attias, Yair Movshovitz-Attias, Peter Steenkiste, and Christos Faloutsos. 2013. Analysis of the reputation system and user contributions on a question answering website: Stackoverflow. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. ACM, 886–893.
  - [170] Susanne van Mulken, Elisabeth André, and Jochen Müller. 1999. An empirical study on the trustworthiness of life-like interface agents. In *Proceedings of the HCI International'99 8th International Conference on Human-Computer Interaction on Human-Computer Interaction: Communication, Cooperation, and Application Design-Volume 2 - Volume 2*. L. Erlbaum Associates Inc., Hillsdale, NJ, 152–156. Retrieved from <http://dl.acm.org/citation.cfm?id=647944.741893>.
  - [171] Emerson Murphy-Hill and Andrew P. Black. 2007. Why don't people use refactoring tools? In *Proceedings of the 1st Workshop on Refactoring Tools*. 61–62.
  - [172] Alena Naiakshina, Anastasia Danilova, Eva Gerlitz, and Matthew Smith. 2020. On conducting security developer studies with CS students: Examining a password-storage study with CS students, freelancers, and company developers. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, 1–13.
  - [173] Anton J. Nederhof. 1985. Methods of coping with social desirability bias: A review. *European Journal of Social Psychology* 15, 3 (1985), 263–280.
  - [174] J. C. Nesbit, O. O. Adesope, Q. Liu, and W. Ma. 2014. How effective are intelligent tutoring systems in computer science education? In *Proceedings of the 2014 IEEE 14th International Conference on Advanced Learning Technologies*. 99–103.
  - [175] Magdalene Ng, Kovila PL Coopamootoo, Ehsan Toreini, Mhairi Aitken, Karen Elliot, and Aad van Moorsel. 2020. Simulating the effects of social presence on trust, privacy concerns & usage intentions in automated bots for finance. In *2020 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*. IEEE, 190–199.
  - [176] Jakob Nielsen and Rolf Molich. 1990. Heuristic evaluation of user interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 249–256.
  - [177] Haoran Niu, Iman Keivanloo, and Ying Zou. 2017. Learning to rank code examples for code search engines. *Empirical Software Engineering* 22, 1 (2017), 259–291.
  - [178] John Noll, Sarah Beecham, and Ita Richardson. 2010. Global software development and collaboration: Barriers and solutions. *ACM Inroads* 1, 3 (2010), 66–78.
  - [179] David Novick and Stephen Sutton. 1997. What is mixed-initiative interaction? In *Proceedings of the AAAI Spring Symposium on Computational Models for Mixed Initiative Interaction*.
  - [180] Catherine S. Oh, Jeremy N. Bailenson, and Gregory F. Welch. 2018. A systematic review of social presence: Definition, antecedents, and implications. *Frontiers in Robotics and AI* 5 (2018), 114. DOI : <https://doi.org/10.3389/frobt.2018.00114>
  - [181] Andy Oram and Greg Wilson. 2010. *Making Software: What Really Works, and Why We Believe It* (1st ed.). O'Reilly Media, Inc.
  - [182] A. F. Osborn. 1957. *Applied Imagination: Principles and Procedures of Creative Thinking*. Charles Scribner's Sons.

- [183] Carlos Pacheco and Michael D. Ernst. 2007. Randoop: Feedback-directed random testing for Java. In *Proceedings of the Companion to the 22nd ACM SIGPLAN Conference on Object-oriented Programming Systems and Applications Companion*. 815–816.
- [184] M. Page-Jones. 1988. *The Practical Guide to Structured Systems Design*. Prentice Hall.
- [185] David Walsh Palmieri. 2002. Knowledge management through pair programming. Retrieved on 17 March, 2022 from <http://www.lib.ncsu.edu/resolver/1840.16/1429>.
- [186] Nelishia Pillay. 2003. Developing intelligent programming tutors for novice programmers. *SIGCSE Bull.* 35, 2 (June 2003), 78–82. DOI: <https://doi.org/10.1145/782941.782986>
- [187] George Polya. 2004. *How to Solve It: A New Aspect of Mathematical Method*. Vol. 85. Princeton university press.
- [188] Marshall Scott Poole, Michael Holmes, Richard Watson, and Gerardine DeSanctis. 1993. Group decision support systems and group communication: A comparison of decision making in computer-supported and nonsupported groups. *Communication Research* 20, 2 (1993), 176–213.
- [189] Marshall Scott Poole, Michael Homes, and Gerardine DeSanctis. 1988. Conflict management and group decision support systems. In *Proceedings of the 1988 ACM Conference on Computer-supported Cooperative Work - CSCW'88*. ACM Press.
- [190] Sihang Qiu, Ujwal Gadiraju, and Alessandro Bozzon. 2020. Estimating conversational styles in conversational micro-task crowdsourcing. *Proceedings of the ACM Human-Computer Interaction* 4, CSCW1 (May 2020), Article 032, 23 pages. DOI: <https://doi.org/10.1145/3392837>
- [191] Sihang Qiu, Ujwal Gadiraju, and Alessandro Bozzon. 2020. Improving worker engagement through conversational microtask crowdsourcing. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* ACM, New York, NY, 1–12. DOI: <https://doi.org/10.1145/3313831.3376403>
- [192] Ricardo Alexandre Peixoto Queirós and José Paulo Leal. 2012. PETCHA: A programming exercises teaching assistant. In *Proceedings of the 17th ACM Annual Conference on Innovation and Technology in Computer Science Education*. 192–197.
- [193] M. Raghothaman, Y. Wei, and Y. Hamadi. 2016. SWIM: Synthesizing what I mean - code search and idiomatic snippet synthesis. In *Proceedings of the 2016 IEEE/ACM 38th International Conference on Software Engineering*. 357–367.
- [194] Vipul Raheja and Joel Tetreault. 2019. Dialogue act classification with context-aware self-attention. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 3727–3733.
- [195] Jorge Ramirez Uresti and Benedict Du Boulay. 2004. Expertise, motivation and teaching in learning companion systems. *International Journal of Artificial Intelligence in Education* 14, 2 (Jan. 2004), 193–231.
- [196] Prerana Pradeepkumar Rane. 2017. *Automatic Generation of Test Cases for Agile Using Natural Language Processing*. Ph. D. Dissertation. Virginia Tech.
- [197] James L. Reinertsen. 2000. Let's talk about error. *BMJ* 320, 7237 (2000), 730. DOI: <https://doi.org/10.1136/bmj.320.7237.730> arXiv:<https://www.bmj.com/content/320/7237/730.full.pdf>.
- [198] Ehud Reiter and Robert Dale. 2000. *Building Natural Language Generation Systems*. Cambridge university press.
- [199] Liliang Ren, Jianmo Ni, and Julian McAuley. 2019. Scalable and Accurate Dialogue State Tracking via Hierarchical Sequence Generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 1876–1885. <https://doi.org/10.18653/v1/D19-1196>
- [200] P. Robe, S. Kaur Kuttal, Y. Zhang, and R. Bellamy. 2020. Can machine learning facilitate remote pair programming? Challenges, insights implications. In *Proceedings of the 2020 IEEE Symposium on Visual Languages and Human-Centric Computing*. 1–11.
- [201] Fernando J. Rodríguez, Kimberly Michelle Price, and Kristy Elizabeth Boyer. 2017. Exploring the pair programming process: Characteristics of effective collaboration. In *Proceedings of the 2017 ACM SIGCSE Technical Symposium on Computer Science Education*. ACM, New York, NY, 507–512. DOI: <https://doi.org/10.1145/3017680.3017748>
- [202] Carl Ransom Rogers and Richard Evans Farson. 1957. *Active Listening*. Industrial Relations Center of the University of Chicago Chicago, IL.
- [203] José Miguel Rojas, José Campos, Mattia Vivanti, Gordon Fraser, and Andrea Arcuri. 2015. Combining multiple coverage criteria in search-based unit test generation. In *Proceedings of the International Symposium on Search Based Software Engineering*. Springer, 93–108.
- [204] Reudismam Rolim, Gustavo Soares, Loris D'Antoni, Oleksandr Polozov, Sumit Gulwani, Rohit Gheyi, Ryo Suzuki, and Björn Hartmann. 2017. Learning syntactic program transformations from examples. In *Proceedings of the 2017 IEEE/ACM 39th International Conference on Software Engineering*. 404–415. DOI: <https://doi.org/10.1109/ICSE.2017.44>
- [205] Sherry Ruan, Jacob O. Wobbrock, Kenny Liou, Andrew Ng, and James A. Landay. 2018. Comparing speech and keyboard text entry for short messages in two languages on touchscreen phones. *Proceedings of the ACM on Interactive Mobile Wearable Ubiquitous Technologies* 1, 4 (Jan. 2018), Article 159, 23 pages. DOI: <https://doi.org/10.1145/3161187>

- [206] Omar Ruvalcaba, Linda Werner, and Jill Denner. 2016. Observations of pair programming: Variations in collaboration across demographic groups. In *Proceedings of the 47th ACM Technical Symposium on Computing Science Education*. ACM, New York, NY, 90–95. DOI : <https://doi.org/10.1145/2839509.2844558>
- [207] V. Sambamurthy and Wynne W. Chin. 1994. The effects of group attitudes toward alternative GDSS designs on the decision-making performance of computer-supported groups\*. *Decision Sciences* 25, 2 (1994), 215–241.
- [208] Kate Kwasny, Peter Robe, Sandeep Kaur Kuttal, and Bali Ong. 2021. Trade-offs for substituting a human with an agent in a pair programming context: The good, the bad, and the ugly. (2021). Retrieved from [https://drive.google.com/drive/folders/1kIkt\\_d-q5F\\_DMsshhtv\\_Mbrm2YhtDdW?usp=sharing](https://drive.google.com/drive/folders/1kIkt_d-q5F_DMsshhtv_Mbrm2YhtDdW?usp=sharing).
- [209] Anita Sarma, Xiaofan Chen, Sandeep Kuttal, Laura Dabbish, and Zhendong Wang. 2016. Hiring in the global stage: Profiles of online contributions. In *Proceedings of the 2016 IEEE 11th International Conference on Global Software Engineering*. IEEE, 1–10.
- [210] T. Savage, M. Reville, and D. Poshyanyk. 2010. FLAT3: Feature location and textual tracing tool. In *Proceedings of the 2010 ACM/IEEE 32nd International Conference on Software Engineering*, Vol. 2. 255–258.
- [211] Marvin L. Schroth and Elissa Lund. 1993. Role of delay of feedback on subsequent pattern recognition transfer tasks. *Contemporary Educational Psychology* 18, 1 (1993), 15–22.
- [212] Carolyn B. Seaman. 1999. Qualitative methods in empirical studies of software engineering. *IEEE Transactions on Software Engineering* 25, 4 (1999), 557–572.
- [213] Young-Ho Seo and Jong-Hoon Kim. 2016. Analyzing the effects of coding education through pair programming for the computational thinking and creativity of elementary school students. *Indian Journal of Science and Technology* 9, 46 (Dec. 2016). DOI : <https://doi.org/10.17485/ijst/2016/v9i46/107837>
- [214] Michael Seymour, Kai Riemer, and Judy Kay. 2017. Interactive realistic digital avatars-revisiting the uncanny valley. (2017). In *Proceedings of the Hawaii International Conference on System Sciences*.
- [215] Arun Shekhar and Nicola Marsden. 2018. Cognitive walkthrough of a learning management system with gendered personas. In *Proceedings of the 4th Conference on Gender & IT*. ACM, 191–198.
- [216] Leonid Sheremetov and Adolfo Guzmán Arenas. 2002. EVA: An interactive web-based collaborative learning environment. *Computers & Education* 39, 2 (2002), 161–182. DOI : [https://doi.org/10.1016/S0360-1315\(02\)00030-1](https://doi.org/10.1016/S0360-1315(02)00030-1)
- [217] Ben Shneiderman. 1982. Designing computer system messages. *Communication of the ACM* 25, 9 (1982), 610–611.
- [218] Leif Singer, Fernando Figueira Filho, Brendan Cleary, Christoph Treude, Margaret-Anne Storey, and Kurt Schneider. 2013. Mutual assessment in the social programmer ecosystem: An empirical investigation of developer profile aggregators. In *Proceedings of the 2013 Conference on Computer Supported Cooperative Work*. ACM, 103–116.
- [219] Gillian Smith, Jim Whitehead, and Michael Mateas. 2010. Tanagra: A mixed-initiative level design tool. In *Proceedings of the 5th International Conference on the Foundations of Digital Games*. ACM, New York, NY, USA, 209–216. DOI : <https://doi.org/10.1145/1822348.1822376>
- [220] Social Bot 2020. Chatbot Statistics. Retrieved 17 March, 2022 from <https://www.smallbizgenius.net/by-the-numbers/chatbot-statistics/#gref>.
- [221] Social Bot [n. d.]. Cleverbot. Retrieved 17 March, 2022 from <https://www.cleverbot.com/>.
- [222] Social Bot [n. d.]. Mitsuku. Retrieved 15 June, 2021 from <https://www.pandorabots.com/mitsuku/>.
- [223] Social Bot [n. d.]. SAP Conversational AI. Retrieved 17 March, 2022 from <https://www.sap.com/products/conversational-ai.html>.
- [224] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. 2020. The Design and Implementation of XiaoIce, an Empathetic Social Chatbot. *Computational Linguistics* 46, 1 (03 2020), 53–93. [https://doi.org/10.1162/coli\\_a\\_00368](https://doi.org/10.1162/coli_a_00368) arXiv:[https://direct.mit.edu/coli/article-pdf/46/1/53/1847834/coli\\_a\\_00368.pdf](https://direct.mit.edu/coli/article-pdf/46/1/53/1847834/coli_a_00368.pdf)
- [225] Social Media Website [n. d.]. Facebook. Retrieved 17 March, 2022 from [www.facebook.com](http://www.facebook.com).
- [226] Social Media Website [n. d.]. Twitter. Retrieved 17 March, 2022 from [www.twitter.com](http://www.twitter.com).
- [227] Software Application [n. d.]. Facerig. Retrieved 15 June, 2021 from <https://facerig.com/>.
- [228] Jörg Spieler. 2021. *UCDetector*. Retrieved from <http://www.ucdetector.org/>.
- [229] Lee Sproull, Mani Subramani, Sara Kiesler, Janet H. Walker, and Keith Waters. 1996. When the interface is a face. *Human-Computer Interaction* 11, 2 (June 1996), 97–124. DOI : [https://doi.org/10.1207/s15327051hci1102\\_1](https://doi.org/10.1207/s15327051hci1102_1)
- [230] Saiying Steenbergen-Hu and Harris Cooper. 2013. A meta-analysis of the effectiveness of intelligent tutoring systems on K–12 students’ mathematical learning. *Journal of Educational Psychology* 105, 4 (2013), 970.
- [231] Saiying Steenbergen-Hu and Harris Cooper. 2014. A meta-analysis of the effectiveness of intelligent tutoring systems on college students’ academic learning. *Journal of Educational Psychology* 106, 2 (2014), 331.
- [232] Anselm L. Strauss and Juliet M. Corbin. 1998. *Basics of Qualitative Research: Techniques and Procedures for Developing Grounded Theory*. Sage Publications, Thousand Oaks, Calif. XIII, 312 s pages.
- [233] Persis T. Sturges. 1972. Information delay and retention: Effect of information in feedback and tests. *Journal of Educational Psychology* 63, 1 (1972), 32.



- [234] Linda K. Swindell and Walter F. Walls. 1993. Response confidence and the delay retention effect. *Contemporary Educational Psychology* 18, 3 (1993), 363–375.
- [235] Ryuichi Takanobu, Runze Liang, and Minlie Huang. 2020. Multi-Agent Task-Oriented Dialog Policy Learning with Role-Aware Reward Decomposition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 625–638. <https://doi.org/10.18653/v1/2020.acl-main.59>
- [236] Diana-Cezara Toader, Grația Boca, Rita Toader, Mara Măcelaru, Cezar Toader, Diana Ighian, and Adrian T. Rădulescu. 2019. The effect of social presence and chatbot errors on trust. *Sustainability* 12, 1 (Dec 2019), 256. DOI: <https://doi.org/10.3390/su12010256>
- [237] Jorge A. Ramirez Uresti. 2000. Should I teach my computer peer? Some issues in teaching a learning companion. In *Proceedings of the Intelligent Tutoring Systems*. Gilles Gauthier, Claude Frasson, and Kurt VanLehn (Eds.). Springer Berlin, 103–112.
- [238] Susanne van Mulken, Elisabeth André, and Jochen Müller. 1998. The persona effect: How substantial is it? In *Proceedings of the People and Computers XIII*. Hilary Johnson, Lawrence Nigay, and Christopher Roast (Eds.). Springer, 53–66.
- [239] Kurt VanLehn. 2011. The relative effectiveness of human tutoring, intelligent tutoring systems, and other tutoring systems. *Educational Psychologist* 46, 4 (2011), 197–221. DOI: <https://doi.org/10.1080/00461520.2011.611369>
- [240] Kurt VanLehn, Arthur C. Graesser, G. Tanner Jackson, Pamela Jordan, Andrew Olney, and Carolyn P. Rosé. 2007. When are tutorial dialogues more effective than reading? *Cognitive Science* 31, 1 (2007), 3–62. DOI: <https://doi.org/10.1080/03640210709336984>
- [241] Roli Varma. 2010. Why so few women enroll in computing? Gender and ethnic differences in students' perception. *Computer Science Education* 20, 4 (Dec. 2010), 301–316. DOI: <https://doi.org/10.1080/08993408.2010.527697>
- [242] Nicolas Vermeulen, Olivier Corneille, and Paula M. Niedenthal. 2008. Sensory load incurs conceptual processing costs. *Cognition* 109, 2 (2008), 287–294. DOI: <https://doi.org/10.1016/j.cognition.2008.09.004>
- [243] Virtual Assistant [n. d.]. Amazon Alexa. Retrieved 17 March, 2022 from <https://developer.amazon.com/en-US/alexa>.
- [244] Virtual Assistant [n. d.]. Apple Siri. Retrieved 17 March, 2022 from <https://www.apple.com/siri/>.
- [245] Virtual Assistant [n. d.]. Google Assistant. Retrieved 17 March, 2022 from <https://assistant.google.com/>.
- [246] Virtual Assistant [n. d.]. IBM Watson Assistant. Retrieved 17 March, 2022 from <https://cloud.ibm.com/apidocs/assistant/assistant-v2>.
- [247] Virtual Assistant 2022. Oracle Digital Assistant. Retrieved 17 March, 2022 from <https://www.oracle.com/solutions/chatbots/>.
- [248] Aurora Vizcaíno. 2005. A Simulated Student Can Improve Collaborative Learning. *International Journal of Artificial Intelligence in Education* 15, 1 (Jan. 2005), 3–40.
- [249] Doug Vogel and Jay Nunamaker. 1990. Group decision support system impact: Multi-methodological exploration. *Information & Management* 18, 1 (1990), 15–28.
- [250] Mihaela Vorvoreanu, Lingyi Zhang, Yun-Han Huang, Claudia Hilderbrand, Zoe Steine-Hanson, and Margaret Burnett. 2019. From gender biases to gender-inclusive design: An empirical investigation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, Article 53, 14 pages. DOI: <https://doi.org/10.1145/3290605.3300283>
- [251] Tony Wagner and Robert A. Compton. 2012. *Creating Innovators: The Making of Young People Who Will Change the World*. Simon and Schuster.
- [252] Pierre Wargnier, Giovanni Carletti, Yann Laurent-Corniquet, Samuel Benveniste, Pierre Jouvelot, and Anne-Sophie Rigaud. 2016. Field evaluation with cognitively-impaired older adults of attention management in the embodied conversational agent louise. In *Proceedings of the 2016 IEEE International Conference on Serious Games and Applications for Health*. IEEE, 1–8.
- [253] Richard T. Watson, Gerardine DeSanctis, and Marshall Scott Poole. 1988. Using a GDSS to facilitate group consensus: Some intended and unintended consequences. *MIS Quarterly* 12, 3 (1988), 463–478.
- [254] Website [n. d.]. Bing QA. Retrieved 17 March, 2022 from <https://www.microsoft.com/en-us/research/project/open-domain-question-answering/>.
- [255] Linda L. Werner, Brian Hanks, and Charlie McDowell. 2004. Pair-programming helps female computer science students. *Journal on Educational Resources in Computing* 4, 1 (March 2004), Article 4, 4. DOI: <https://doi.org/10.1145/1060071.1060075>
- [256] Wayne A. Wickelgren. 1974. *How to Solve Problems: Elements of a Theory of Problems and Problem Solving*. WH Freeman San Francisco.
- [257] Alex C. Williams, Harmanpreet Kaur, Shamsi Iqbal, Ryen W. White, Jaime Teevan, and Adam Fourney. 2019. Mercury: Empowering programmers' mobile work practices with microp productivity. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology*. ACM, New York, NY, 81–94. DOI: <https://doi.org/10.1145/3332165.3347932>

- [258] Laurie Williams and Bob Kessler. 2000. The effects of “pair-pressure” and “pair-learning” on software engineering education. In *Proceedings of the 13th Conference on Software Engineering Education & Training*. IEEE Computer Society, Washington, DC, 59–. Retrieved from <http://dl.acm.org/citation.cfm?id=794188.794326>.
- [259] Laurie Williams and Robert Kessler. 2002. *Pair Programming Illuminated*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA.
- [260] L. Williams, C. McDowell, N. Nagappan, J. Fernald, and L. Werner. 2003. Building pair programming knowledge through a family of experiments. In *Proceedings of the 2003 International Symposium on Empirical Software Engineering*, 2003. 143–152. DOI: <https://doi.org/10.1109/ISESE.2003.1237973>
- [261] Laurie A. Williams, Eric N. Wiebe, Kai Yang, Miriam Ferzli, and Carol Miller. 2002. In support of pair programming in the introductory computer science course. *Computer Science Education* 12, 3 (2002), 197–212.
- [262] James Wilson and Daniel Rosenberg. 1988. Rapid prototyping for user interface design. In *Proceedings of the Handbook of Human-computer Interaction*. Elsevier, 859–875.
- [263] Ian H. Witten, Craig G. Nevill-Manning, and D. L. Maulsby. 1996. Interacting with learning agents: Implications for ml from hci. In *Proceedings of the Workshop on Machine Learning meets Human-Computer Interaction, ML*, Vol. 96. 51–58.
- [264] Beverly Park Woolf. 2008. *Building Intelligent Interactive Tutors: Student-Centered Strategies for Revolutionizing e-Learning*. Morgan Kaufmann Publishers Inc., San Francisco, CA.
- [265] Qingyang Wu, Yichi Zhang, Yu Li, and Zhou Yu. 2021. Alternating Recurrent Dialog Model with Large-scale Pre-trained Language Models. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*. Association for Computational Linguistics, Online, 1292–1301. <https://doi.org/10.18653/v1/2021.eacl-main.110>
- [266] Yvonne Wærn and Robert Ramberg. 1996. People’s perception of human and computer advice. *Computers in Human Behavior* 12, 1 (1996), 17–27. DOI: [https://doi.org/10.1016/0747-5632\(95\)00016-X](https://doi.org/10.1016/0747-5632(95)00016-X)
- [267] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R. Salakhutdinov, and Quoc V. Le. 2019. XLNet: Generalized autoregressive pretraining for language understanding. In *Proceedings of the Advances in Neural Information Processing Systems 32*. H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Eds.). Curran Associates, Inc., 5753–5763. Retrieved from <http://papers.nips.cc/paper/8812-xlnet-generalized-autoregressive-pretraining-for-language-understanding.pdf>.
- [268] Georgios N. Yannakakis, Antonios Liapis, and Constantine Alexopoulos. 2014. *Mixed-Initiative Co-Creativity*. Foundations of Digital Games.
- [269] Nick Yee, Jeremy N. Bailenson, and Kathryn Rickertsen. 2007. A meta-analysis of the impact of the inclusion and realism of human-like faces on user experiences in interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, New York, NY, 1–10. DOI: <https://doi.org/10.1145/1240624.1240626>
- [270] Mohan Zalake, Julia Woodward, Amanpreet Kapoor, and Benjamin Lok. 2018. Assessing the Impact of virtual human’s appearance on users’ trust levels. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*. ACM, New York, NY, 329–330. DOI: <https://doi.org/10.1145/3267851.3267863>
- [271] Jiaping Zhang, Tiancheng Zhao, and Zhou Yu. 2018. Multimodal Hierarchical Reinforcement Learning Policy for Task-Oriented Visual Dialog. In *Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*. Association for Computational Linguistics, Melbourne, Australia, 140–150. <https://doi.org/10.18653/v1/W18-5015>
- [272] Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. Task-oriented dialog systems that consider multiple appropriate responses under the same context. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 9604–9611.
- [273] Yong Zhao. 2012. *World Class Learners: Educating Creative and Entrepreneurial Students*. Corwin Press.
- [274] R. Zhi, S. Marwan, Y. Dong, N. Lytle, T. W. Price, and T. Barnes. 2019. Toward data-driven example feedback for novice programming. In *Proceedings of the 12th International Conference on Educational Data Mining*.
- [275] C. Zhou, S. K. Kuttal, and I. Ahmed. 2018. What makes a good developer? An empirical study of developers’ technical and social competencies. In *Proceedings of the 2018 IEEE Symposium on Visual Languages and Human-Centric Computing*. 319–321.
- [276] Su Zhu, Jieyu Li, Lu Chen, and Kai Yu. 2020. Efficient context and schema fusion networks for multi-domain dialogue state tracking. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 766–781. <https://doi.org/10.18653/v1/2020.findings-emnlp.68>
- [277] Franz Zieris and Lutz Prechelt. 2014. On knowledge transfer skill in pair programming. In *Proceedings of the 8th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*. ACM, New York, NY, Article 11, 10 pages. DOI: <https://doi.org/10.1145/2652524.2652529>

Received November 2020; revised July 2021; accepted November 2021