

FROM INTUITION TO UNDERSTANDING: USING AI PEERS TO OVERCOME PHYSICS MISCONCEPTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Generative AI has the potential to transform personalization and accessibility of education. However, it raises serious concerns about accuracy and helping students become independent critical thinkers. In this study, we designed a helpful yet fallible AI “Peer” to help students correct fundamental physics misconceptions related to Newtonian mechanic concepts. In contrast to approaches that seek near-perfect accuracy to create an authoritative AI tutor or teacher, we directly inform students that this AI can answer up to 40% of questions incorrectly. In a randomized controlled trial with 165 students, those who engaged in targeted dialogue with the AI Peer achieved post-test scores that were, on average, 10.5 percentage points higher—with over 20 percentage points higher normalized gain—than a control group that discussed physics history. Qualitative feedback indicated that 91% of the treatment group’s AI interactions were rated as helpful. Furthermore, by comparing student performance on pre- and post-test questions about the same concept, along with experts’ annotations of the AI interactions, we find initial evidence suggesting the improvement in performance does not depend on the correctness of the AI. With further research, the AI Peer paradigm described here could open new possibilities for how we learn, adapt to, and grow with AI.

1 INTRODUCTION

Students have recently been exposed to the remarkable capabilities of Generative AI (AI) in education (AIED). For example, OpenAI’s ChatGPT has been reported to successfully support teaching preparation, assessment design and grading, and student learning (Lo, 2023). Systems like ChatGPT show potential to save time and enhance teaching and learning, including critical and higher-order thinking tasks (Lo, 2023).

However, there is limited concrete evidence whether these tools are effective at improving student learning outcomes (Samson R, 2025). In fact, LLMs are well-reported to hallucinate information and provide sycophantic answers (Wei et al., 2023; Perez et al., 2023), suggesting that AI could actually be detrimental to student learning if it introduces inaccuracies and biases into classroom materials. Given that an increasing large number of students are now using AI to help them with their school assignments (Foundation, 2024), there is a growing concern that students must be taught ‘AI literacy’ to recognize and evaluate potential errors generated by LLMs (Wineburg & Ziv, 2024). Although recent techniques such as Retrieval-Augmented Generation (Piktus et al., 2020) have improved the veracity of generated content, experts disagree on whether hallucinations will be consistently preventable even in the future (Samson R, 2025) – and they are clearly not today. Thus, there is an urgent need to determine whether AIED can be an effective tool for education despite these inherent limitations.

This research evaluates the potential of LLMs to support student learning through text-based conversations in an introductory university physics class. The study explores how LLMs can be used to facilitate learning of physics concepts, positioning the AIED as a non-expert peer rather than an expert teacher. Students complete a modified Force Concept Inventory (FCI) as a pre-test. Our AI Peer is prompted with the mistakes made by each student, resulting in a personalized discussion focused on the identified misconceptions. Students then completed the standard FCI as a post-test that assesses the same concepts via a different questionnaire. Our results showed that the treatment group, after a focused discussion with the AI, experienced significantly higher learning gain between tests

054 compared to the control group, who discussed physics history with the AI companion. Importantly,
055 our AI, which leverages the full abilities of GPT-4o without any artificial reduction, was not a reli-
056 able source of truth; it answered up to 40% of the FCI questions incorrectly. Thus, this study marks
057 a first step towards achieving a comprehensive individualized approach to AI-assisted education by
058 aligning human expectations to current limitations of generative AI.

059 Our key contributions include:

- 061 • An experimental framework for measuring AI’s educational capabilities. In our experi-
062 ment, we create a new FCI¹ and focus on physics. The framework can be extended to other
063 domains and other AI systems.
- 064 • Showing that AI can help correct common, deep-seated physics misconceptions that per-
065 sisted over the first half of a semester of traditional instruction.
- 066 • A “peer” rather than “instructor” approach that we find preserves learning, while lowering
067 barriers in how much AI accuracy is needed. With further research, this could potentially
068 expand the areas in which we can benefit from AI, while empowering the development of
069 human critical thinking skills.

072 2 LITERATURE REVIEW

073
074 Educators, policymakers, EdTech companies, and students are generally optimistic about the po-
075 tential for AI to deliver personalized education, despite the limited evidence that these tools can
076 be effective at improving student learning outcomes (Samson R, 2025). In a systematic review of
077 113 papers relating to AI Education tools, Chiu et al. (2023) found that, while AIED can gener-
078 ally improve student motivation and output, methods used to evaluate AIED were often ineffective
079 at measuring students’ learning. The authors suggested further research is needed to “devise new
080 methods for evaluating the success of AI systems.” They also noted that, despite the potential of
081 AIED to provide equitable education through personalized feedback, AI could potentially worsen
082 educational inequity if the AIED design process lacked consideration of pedagogy and learning
083 sciences.

084 Pedagogical research has identified numerous teaching strategies with the potential to significantly
085 accelerate student learning, as described by Hattie (2009) in a synthesis of meta-analyses on stu-
086 dent achievement. Among them are discussion with another student (argumentation), and prompt
087 oriented or directional feedback. Baidoo-Anu & Ansah (2023) confirm that AI tools can enhance
088 teaching and learning experiences by supporting personalized and interactive learning. For instance,
089 students could leverage the capabilities of advanced generative AI to systematically explain com-
090 plex concepts. Mollick & Mollick (2023) suggest several teaching strategies that could potentially
091 enhance student learning in the presence of AI while mitigating the associated risks of these tech-
092 nologies. Their research emphasizes the importance of maintaining human involvement in the edu-
093 cational process and positioning AI as a supportive tool rather than a substitute for human instruc-
094 tors. As de Jong et al. (2023) argue, an AIED that combines direct instruction with inquiry should
095 be more effective at explaining new concepts than having the AI directly answer student questions.
096 This is precisely the approach taken by Jurenka et al. (2024), who tested an LLM-based support
097 tool with a class of 113 Arizona State University students. The tool, HallMate, could discuss course
098 videos, direct learners to relevant content, provide scaffolded homework help, and assist with time
099 management and broad learning strategies. The study tracked whether students used HallMate, how
100 they used it, and whether students felt it was useful. It also evaluated the pedagogical quality of the
101 conversations. While promising, this study focused on usability and engagement but did not attempt
to measure the impact of AI tools on students’ learning or performance.

102 To summarize, current studies on the role of AI in education tend to emphasize its potential for
103 improving teaching and learning experiences. While these studies acknowledge that AI is potentially
104 revolutionary in education—specifically because it can provide personalized content focused on
105 each student’s learning needs—they lack robust evaluation methods for assessing AIED’s impact on
106 student learning outcomes.

107 ¹Will be released in camera-ready.

A common tool found in the education literature to evaluate students' learning in a specified discipline is the 'Concept Inventory' (Sands et al., 2018). This approach consists of a set of multi-choice questions focused on testing conceptual understanding. In these evaluations, incorrect answers—called distractors—align with common student misconceptions. Thus, because they are standardized, these types of tests can measure the effectiveness of an instructional method by comparing student scores on matched questions before and after instruction.

The first concept inventory test—the 'Force Concept Inventory' (FCI)—aimed to elicit students' misconceptions about fundamental Newtonian mechanics that were generally deep-held and therefore difficult to correct through conventional physics instruction (Halloun & Hestenes, 1985; Hestenes et al., 1992). Given Costello et al. (2024)'s recent success in reducing deep-held beliefs in conspiracy theories through dialogue with AI, we theorize that this application will transfer into education; that is, discussion with generative AI may reduce students' deep-held erroneous beliefs about fundamental Newtonian mechanics. We draw inspiration from that work for some aspects of the experimental setup, such as having students explain their reasoning in their own words. Meanwhile, the FCI is a robust tool that lets us both identify misconceptions and evaluate the effectiveness of AI for this application.

One challenge with this approach to measure the effectiveness of general purpose AI as a teaching tool is its propensity to produce false explanations and to hallucinate. Similarly, there is growing concern that students lack the 'digital literacy' required to critically evaluate online content, including outputs of generative AI (Wineburg & Ziv, 2024). One key element of digital literacy is determining the credibility of the source (McGrew & Breakstone, 2023). This suggests that students should be instructed to be skeptical of AI explanations.

Our primary goal in this study is to evaluate the effectiveness of a fallible AI as a tool for addressing misconceptions in an educational context, using the FCI as a robust and trusted scientific measure of learning outcomes. We theorize that generative AI will show potential as a tool to effectively address student misconceptions in physics education, despite its well known limits. The results of this study can support future AIED tools to be evaluated for effectiveness at improving student learning, as well as contribute to the growing body of literature exploring the use of general purpose AI to fight other types of misconceptions, such as disinformation and conspiratorial beliefs (Demartini et al., 2020).

3 METHODOLOGY

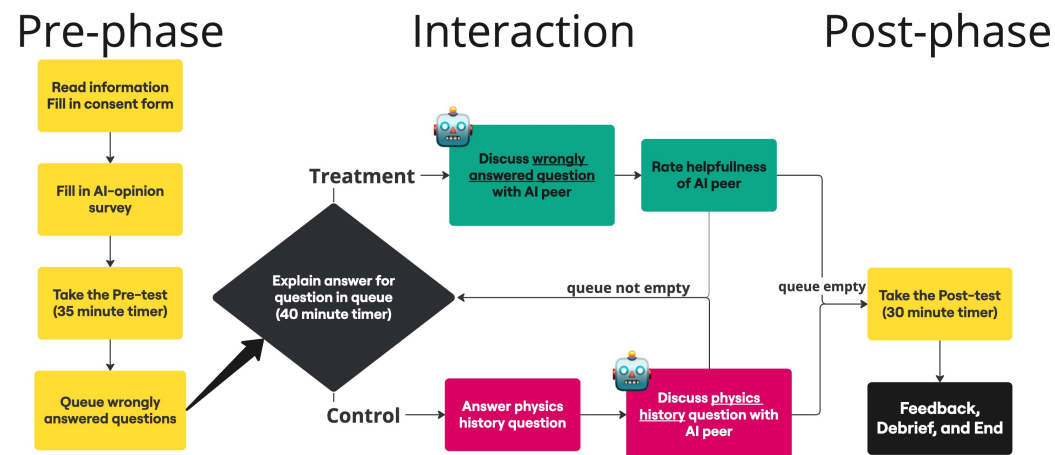


Figure 1: Study procedure showing pre-test, interaction, and post-test phases.

This study uses an online experiment approach designed to assess the effectiveness of an AI companion for helping undergraduate physics students overcome common misconceptions in Newtonian mechanics. We are measuring the learning outcomes by comparing the results of a Force Concept Inventory test (Hestenes et al., 1992) pre- and post-interaction with the AI companion. The pre-test

162 was modified from the original FCI to reduce improvements resulting from additional thinking time,
163 or memorization of correct answers during interaction with the AI. The students had prior exposure
164 to the post-test, which they took at the beginning of the semester approximately 2 months prior to
165 our experiment, but they had not been provided with the answers. The misconceptions are catego-
166 rized in six sub-concepts: Kinematics, Newton’s First Law, Newton’s Second Law, Newton’s Third
167 Law, the Superposition Principle, and Kinds of Force.

168 The experiment was administered to 165 undergraduate students enrolled in an introductory physics
169 course from a North American R1 university. Students completed the experiment on a dedicated
170 website, created for the purpose of this study, during their in-class lab sessions with passive super-
171 vision from their regular teaching assistant (not part of the research team). We divided students into
172 two groups with equal probability. Students in the treatment group interacted with AI to discuss
173 their incorrect answers on the pre-test and were informed that the AI should be seen as a fallible
174 peer, rather than an authoritative teacher, while students in the control group interacted with an AI
175 to discuss their answers on a test about a historical figure in physics. We filter out students who
176 failed an attention check, and ones who spent less than 5 minutes on the post-test (indicating a lack
177 of effort). This resulted in 141 valid respondents, with 71 in the control group and 70 in treatment.
178 The main experiment design, with results reported later in Table 2, was pre-registered.² Below, we
179 describe each condition in more details.

181 3.1 EXPERIMENTAL SETUP

183 3.1.1 CONTROL GROUP

184 The first group, serving as the control group, begins by taking a modified version of the FCI test
185 for 35 minutes. This modified test (also referred to as the pre-test) was specially created by our
186 team to have 30 questions similar to and assessing the same core Newtonian mechanics concepts
187 as the original FCI of Hestenes et al. (1992). After completing the pre-test, students were asked to
188 interact with an AI for up to 40 minutes. First, we identify all the wrong answers on the pre-test
189 by comparing with the answer key. The website then presents each wrong answer to the student in
190 a random sequence, covering all the categories of misconceptions. The sequence works as follows:
191 each student was presented with one question they answered incorrectly; they were then asked to
192 explain their reasoning for that answer; after this step, they were presented with a multiple choice
193 question about a specific historical figure in physics; they were then asked to interact with the AI
194 companion to discuss this historical figure for three rounds of conversation. These questions were
195 provided one at a time, and the student had to interact with the AI chatbot companion before moving
196 on to the next wrongly answered question in the pre-test. There were a total of 30 questions about
197 historical figures, so students could review up to 30 incorrect answers in the pre-test. Thus, the
198 students are informed which pre-test questions they got wrong, and interacting with the AI has a
199 1-1 correspondence with those wrong answers, but the control group students do not discuss those
200 questions with the AI. The students then proceed to take the original FCI test that was described in
201 Hestenes et al. (1992) (also referred to as the post-test) for an additional 30 minutes.

203 3.1.2 TREATMENT GROUP

204 The treatment group mirrors the control group, except that instead of talking about physics history
205 questions with the AI, they talk about the questions they got wrong on the pre-test. Specifically, the
206 40-minute middle portion of the experiment was sequenced as follows: each student was presented
207 with one question they answered incorrectly; they were then asked to explain their reasoning for that
208 answer; after this step, the AI Peer was provided the student’s explanation and instructed to correct
209 the misconception, explain the correct answer and help students grasp the underlying concepts cov-
210 ered in the question (see Prompt A.1.2). Like the control group, the students interact with the AI for
211 three conversation turns. Finally, they rate the helpfulness of the AI on the question at hand, before
212 moving on to the next wrongly answered question. All other aspects of the experiment, such as the
213 pre-test and the post-test, are identical between both groups.

214
215
²Anonymized; details available upon request.

3.2 AI AS A PEER, NOT AUTHORITY

Prior to the beginning of the experiment, the students were informed that the AI should not be viewed as an authoritative teacher, but more like a peer, who is helpful but can answer questions incorrectly. To evaluate the performance of the AI, we prompted it to answer pre-test questions. We tested OpenAI’s GPT-4o-2024-08-06. Since the pre-test contains 13 accompanying images depicting various physics problems, we had a PhD student write descriptions of each image. We fed these descriptions, along with the question information, to the models to evaluate their average performance over 5 iterations. We additionally tested the multimodal GPT-4o model on the original images. We set temperature to 0.7 for GPT-4o. All other settings were default.

On average, GPT-4o got 59% of pre-test questions correct over 5 iterations when using written image descriptions (broken down question-by-question in Figure 3). When looking at performance per number of images descriptions in the question text, the model performed better on questions with fewer images (see Table 1). GPT-4o using images rather than image descriptions performed worse, only getting 49% of questions correct over 5 iterations. This reduction in performance is caused by the AI getting confused about the direction and overlap of lines and dotted lines, especially when they appear close to each other.

Table 1: Accuracy drops as the number of images in a pre-test question increases, suggesting GPT-4o’s reasoning struggles with multiple image descriptions

Number of Images	Correctness (%)	(%) of Total Questions
0	67.69	43.33
1	62.22	30.00
>2	26.67	26.67

3.3 MEASURING THE LEARNING EFFECT (GAIN)

We measure the learning effect using Hake’s normalized gain g , which is expressed as

$$g = \frac{\text{post}\% - \text{pre}\%}{100 - \text{pre}\%} \quad (1)$$

where pre and $post$ are the mean pre-treatment test and post-treatment test scores of all students respectively (Hake, 1998). The normalized gain g is intended to measure the gain of a class as a fraction of the possible remaining gain on the assessment, allowing comparisons of improvements in different classes regardless of the mean pre-test score. This measure of gain also makes sense for individual students. For example, students with an individual normalized gain g of 0 have scored the same on the pre-test and post-test. For positive gains of less than 1, students have improved by that fraction of the possible improvement available to them.

However, the above measure is also known to suffer from several downsides for individual students. For instance, the individual gain g of a single student who usually performs well can become negative with a large magnitude if the student obtains a high score on the pretest but performs less well on the post-test. In the extreme case where the student obtains a perfect score on the pretest, the individual gain g is negatively infinite. We cannot even compute g for a student who scores perfectly on both the pretest and post-test. However, calculating the class gain g as described above avoids these pitfalls, since the average class pretest scores will almost always be far from perfect. For this reason, the class gain g is widely employed to study student outcomes with concept assessments (Hake, 1998). Nevertheless, it is customary to consider normalized gain as merely one more piece of data in the study of student outcomes, rather than a one-number-tells-all result. Therefore, we also conducted an in depth human annotation of the AI interactions, and analyzed how those annotations relate to student performance.

3.4 RESOURCES AND DATA COLLECTION

Students used a website developed for the purpose of this research. Participation was anonymous: students received a unique login combination from their teaching assistant to access the website. We

270 tracked user behavior, such as time spent on the tests, their inputs to the AI and the resulting model
 271 replies, and ratings of AI helpfulness after each interaction. In a short survey before the pre-test, we
 272 collected opinions on students' excitement/concern about AI, confidence in their work/educational-
 273 related AI abilities, as well as their use of AI (see Figure A.7).

274 4 RESULTS

275
 276
 277
 278 Table 2: Treatment group showed significantly higher post-test scores and
 279 normalized gain, suggesting the intervention was effective

Metric	Control	Treatment	p-value
Pre-Test	51.5	50.7	0.769
Post-Test	62.7	73.2	0.001
Normalized Gain	27.6	47.9	0.0001

288 4.1 STUDENT PERFORMANCE ON THE FCI TESTS

289
 290 As shown in Table 2, the treatment group's post-test scores improved by 10.5 percentage points more
 291 than those of the control group. This suggests that the AI Peer had a positive impact on reducing
 292 physics misconceptions. The differences in post-test scores, improvement, and gain between the
 293 treatment and control groups were statistically significant ($p < 0.01$). Meanwhile, we confirm
 294 that their pre-test scores (50.7% treatment vs 51.5% control) did not have statistically significant
 295 differences ($p = 0.769$). Across both groups, students scored worse on the pre-test. Without having
 296 used the AI companion, the control group answered on average 62.7% of questions correctly on the
 297 post-test, compared to a 51.5% correct rate on the pre-test; a difference of +11.2%. The treatment
 298 group answered on average 73.2% of questions correctly on the post-test, compared to a 50.7%
 299 correct rate on the pre-test; a difference of +22.5%. This result suggests that the treatment group
 300 received a meaningful learning effect from the AI.

301 4.2 HUMAN EVALUATION OF AI OUTPUTS

302 4.2.1 STUDENT FEEDBACK

303
 304
 305 Table 3: Most students found the treatment AI interactions either very helpful or
 306 fairly helpful, with only a small number rating them as unhelpful

Rating	Count
Very helpful	420
Fairly helpful	327
Uncertain	37
Fairly unhelpful	18
Very unhelpful	20

307
 308
 309
 310
 311
 312
 313
 314
 315
 316 Table 3 shows that 91% of interactions, rated after every 3-round dialogue in the treatment group,
 317 were qualified as fairly or very helpful by our Physics graduate students. This result, combined
 318 with Table 2, indicates that not only did the treatment group perform better objectively after AI-
 319 interaction, they also found the AI interactions informative subjectively. Open-input feedback, from
 320 the end of the experiment, also indicated that the AI was generally helpful, albeit giving lengthy
 321 responses.

322 Meanwhile, a significant number of students said they found the experiment to be too long (90+ min-
 323 utes, 60 total FCI questions), especially those completing the study in the early evening lab session.
 Also, many students indicated that they had wished to see their test results. And students suggested

324 removing the fixed three-round dialogue. For instance, on some occasions, students understood
325 their misconception after the first AI message. On others, students indicated desire to continue the
326 conversation after the final round was already over.

329 4.2.2 GRADING THE AI INTERACTIONS

331 To develop a deeper understanding of the interactions that took place and their impact on learning,
332 a team of six physics graduate students analyzed the AI system prompt and conversations. The
333 interactions were randomly shuffled, and graded according to 6 unique criteria (Table 4).

334 First, we assessed whether the AI explained the key physics concept that the student needed to
335 understand for the respective question. In approximately 55% of interactions the AI explained it
336 clearly and in 25% it touched on it, but in the remainder of interactions it either did not address it
337 (9%) or explained it inaccurately (11%).

338 We then analyzed the overall number of inaccurate physics statements by the AI, regardless whether
339 they were about the key concept or not. Nearly 20% of interactions contain at least one definitively
340 incorrect statement. Furthermore, we also analyzed whether the AI made ambiguous or misleading
341 statements, even if they weren't definitively wrong. There was at least 1 such statement in 36% of
342 interactions. For example, the AI might say "the sled goes in a straight line for a moment", in a
343 context where the sled would actually continue in a straight line for an extended period of time – so
344 the AI's statement is technically true, but potentially misleading.

345 These evaluations show that the AI is making a significant number of dubious statements. To better
346 understand how this affected the student, we then analyzed whether the student appeared to accept
347 or reject the AI explanation – and if they accepted it in a case where it was wrong and they appeared
348 to learn the key concept incorrectly. Here, 86% accepted the explanation, and an additional 6%
349 appeared to learn the concept incorrectly, while the remainder did not accept the explanation. We
350 caution that this analysis of what was revealed in the conversation does not necessarily indicate what
351 will stay with the student as they potentially reflect further beyond the moment. It provides some
352 suggestive evidence, though, that the students are indeed thinking critically in some interactions, ex-
353 plaining even key concepts inaccurately does not necessarily translate into inaccurate understanding
354 of the student.

355 We also assessed student engagement, finding that while 40% of interactions were clearly engaging
356 the AI, 43% had moderate or uncertain engagement with short responses, and 16% were clearly
357 disengaged. This could be an area for further improvement of the system.

360 4.2.3 GRADER QUALITATIVE FEEDBACK

362 A focus-group discussion among graders allowed us to distill their observations into comprehensive
363 feedback on the AI interactions. The graders noted that the AI often fails to identify and address
364 student's specific misconceptions. Rather than focusing on the error, such as a misconception about
365 Newton's laws, the AI tended to provide a lengthy general explanation of the entire problem. The
366 AI's first message initializing the conversation was typically clear. However, longer interactions
367 revealed that the AI can get sidetracked, especially when the student asked probing questions. The
368 emphasis on long explanations and general (and often repeated) analogies sometimes left little room
369 for meaningful student input. When students' responses imply a misconception, the AI rarely asks
370 them to elaborate before providing an explanation; missing an opportunity for deeper interaction.
371 Raters suggest that a more inquisitive and interactive approach could stimulate students' "mental
372 gears" by prompting them to reason from first principles. One rater recommended that the AI should
373 first identify the students' misconception, ask them to explain why it's wrong, and elaborate on their
374 reasoning. Then, instead of a one-size-fits-all explanation with generic analogies, the AI could focus
375 more correcting their misconceptions, based on their faulty reasoning.

376 By providing a more interactive approach, students would be probed to correct their reasoning from
377 within, guided by the AI, rather than learn from the AI's all encompassing explanation of the prob-
378 lem. At the same time, the increase in interactivity, combined with flexible number of rounds in the
379 interaction, could increase engagement and elicit better learning.

4.2.4 PERFORMANCE OF THE AI VS. PERFORMANCE OF THE STUDENT

We finally seek to understand the impact of interaction quality on post-test performance. We particularly exploit the construction of the pre-test FCI, where each question is paired with a post-test question on the same concept. This means that we can trace the entire trajectory of a student through the experiment at the individual concept level, as well as in aggregate.

In Figure 2, we examine the impact of how well the AI addressed the key concept, on whether the student answered the corresponding post-test question correctly. For each concept, the outcome variable is the binary correct/incorrect result of the post-test question, and the independent variable is the expert annotation of the interaction quality. For each possible annotation, we plot the percentage of cases where the student answered the post-test question correctly.

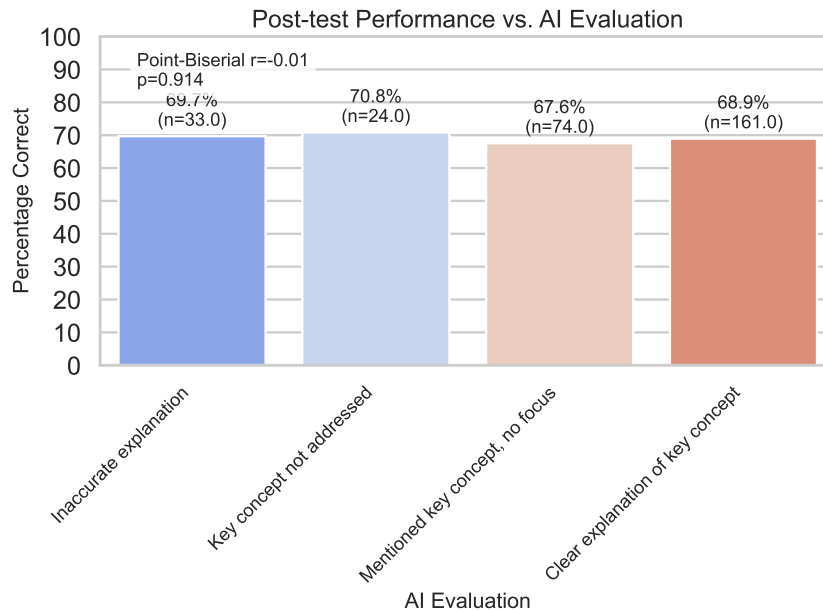


Figure 2: Post-test performance against graded quality of the AI Peer’s responses per interaction

One might expect that students would do better when the AI gives a clear explanation of the key concept, and worst when it explains the key concept incorrectly. But instead we see approximately and statistically equal performance regardless of AI explanation quality. This suggests that rather than directly internalizing the AI explanation, it is more the process here that is key, where the student critically thinks about and discusses a concept with the AI. As shown in Appendix A.6.1, we find similar results with our other annotation types.

5 DISCUSSION

Our results show that AI intervention led to a statistically significant improvement in post-test scores and normalized gains compared to the control group. This is in spite of the stubbornness of the misconceptions described by Halloun & Hestenes (1985): persisting after weeks of traditional education. Thus, AI-led dialogue has promising potential to remediate deep-held misconceptions in educational settings. Further research is warranted to explore potential improvements to the design of the AI Peer, and our understanding of the applications of this technology, such as application in other learning areas. Implications extend beyond direct educational gains in science subjects; if we can apply this concept to other domains such as critical and higher-order thinking, tools like our AI companion may support humans to build robust critical cognition skills, and better assess uncertain information at large.

5.1 LIMITATIONS AND FUTURE WORK

While the results show promise, there are several limitations and areas for further research. We established that the treatment group had a statistically and practically significant improvement after talking with the AI, in spite of the AI's imperfect response quality. On the one hand, our approach of emphasizing critical assessment of the AI responses, where the students are told the AI is significantly inaccurate, can open new application areas (where AI would otherwise be too inaccurate to use) and potentially empower and grow student thinking. On the other hand, though, a more accurate AI might simply be better at teaching. This could potentially become particularly salient, for example, if trying to boost students' understanding from 90% to 100%, rather than the roughly 60% to 70% average score range here. Thus, while we hope this work will open doors towards AI Peers as an education tool, there remain many open questions on this paradigm.

A logical first step is to test reasoning models like OpenAI's o3 that exhibit superior performance on complex reasoning tasks and application of knowledge on unfamiliar tasks, like our newly created FCI questions. After conducting the study here, we tested o3 using our pre-test and achieved a score of 83% over 5 iterations (Figure 4) – significantly higher than GPT-4o's 59%. Its costs were comparable to GPT-4o, and inference speed was fast enough for interaction with the user, suggesting this is a promising model to experiment with.

Another limitation is the nature of the dialogue. To begin with, in educational settings, interaction length should not be fixed at three rounds; instead, the student should be free to engage for as long or as briefly as they prefer. Furthermore, graders indicated that the AI tended to provide lengthy, generic responses with repetitive analogies. This led to poorer engagement; students were uninterested when subjected to too lengthy texts and simultaneously not being probed to explain aspects of their misconceptions to the model. This could at minimum reduce usage of such a tool in real settings and limit its impact, and also may be significantly limiting how much students can learn from it if it could be more to-the-point and motivating. To resolve this, a first step could be prompting the AI to be more inquisitive and better elicit misconceptions from students during the dialogue, instead of providing a more generic explanation of the problem. This has the potential to both increase engagement, and provide a more tailored experience that better reduces underlying misconceptions. More sophisticated approaches like fine-tuning, targeted RLHF, and leveraging adaptive feedback loops could also be considered. Through tools like these, an AI Peer might be better equipped to detect subtle misconceptions and engage students more deeply, fostering a personalized learning experience.

Finally, the duration and scope of the research could be expanded. Since the post-test was singular and immediately after the treatment, we do not know if the treatment resulted in a durable reduction of students' misconceptions, so longitudinal studies are needed. Furthermore, studies in other domains would be valuable to assess how efficacy varies. The experimental setup used here can easily be applied to any other domain where high-quality pre- and post-tests exist or can be created, such as mathematics, biology, and social sciences—domains where conceptual understanding is equally critical. Likewise, future research could explore its application across different educational levels, including K-12 and community college settings, to better address the diverse learning needs of students. Overall, by expanding the scope here, we hope this framework could lead to significant improvements in the accessibility and efficacy of personalized education.

REFERENCES

- David Baidoo-Anu and Leticia Anshah. Education in the era of generative artificial intelligence (ai): Understanding the potential benefits of chatgpt in promoting teaching and learning. 03 2023.
- Thomas K.F. Chiu, Qi Xia, Xinyan Zhou, Ching Sing Chai, and Miaoting Cheng. Systematic literature review on opportunities, challenges, and future research recommendations of artificial intelligence in education. *Computers and Education: Artificial Intelligence*, 4:100118, 2023. ISSN 2666-920X. doi: <https://doi.org/10.1016/j.caeai.2022.100118>. URL <https://www.sciencedirect.com/science/article/pii/S2666920X2200073X>.
- Thomas H Costello, Gordon Pennycook, and David G Rand. Durably reducing conspiracy beliefs through dialogues with ai. *Science*, 385(6714):eadq1814, 2024.

- 486 Ton de Jong, Ard W Lazonder, Clark A Chinn, Frank Fischer, Janice Gobert, Cindy E Hmelo-Silver,
487 Ken R Koedinger, Joseph S Krajcik, Eleni A Kyza, Marcia C Linn, et al. Let's talk evidence—
488 the case for combining inquiry-based and direct instruction. *Educational Research Review*, pp.
489 100536, 2023.
- 490 Gianluca Demartini, Stefano Mizzaro, and Damiano Spina. Human-in-the-loop artificial in-
491 telligence for fighting online misinformation: Challenges and opportunities. *Bulletin*
492 *of the IEEE Computer Society Technical Committee on Data Engineering*, 43(3):65–74,
493 2020. URL [https://citeseerx.ist.psu.edu/document?repid=rep1&type=](https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=462eb5d25a5eb3c18e22ee486019b2153b1f0785)
494 [pdf&doi=462eb5d25a5eb3c18e22ee486019b2153b1f0785](https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=462eb5d25a5eb3c18e22ee486019b2153b1f0785).
- 495 Walton Family Foundation. Ai chatbots in schools findings from a poll
496 of k-12 teachers, students, parents, and college undergraduates, 06 2024.
497 URL [https://www.waltonfamilyfoundation.org/learning/](https://www.waltonfamilyfoundation.org/learning/the-value-of-ai-in-todays-classrooms)
498 [the-value-of-ai-in-todays-classrooms](https://www.waltonfamilyfoundation.org/learning/the-value-of-ai-in-todays-classrooms). Accessed on Date of Access.
- 500 Richard R Hake. Interactive-engagement versus traditional methods: A six-thousand-student survey
501 of mechanics test data for introductory physics courses. *American journal of Physics*, 66(1):
502 64–74, 1998.
- 503 Ibrahim Abou Halloun and David Hestenes. The initial knowledge state of college physics students.
504 *American journal of Physics*, 53(11):1043–1055, 1985.
- 505 John Hattie. *Visible Learning: A Synthesis of Over 800 Meta-Analyses Relating to Achievement*.
506 Routledge, 01 2009. ISBN 9780203887332. doi: 10.4324/9780203887332.
- 507 David Hestenes, Malcolm Wells, Gregg Swackhamer, et al. Force concept inventory. *The physics*
508 *teacher*, 30(3):141–158, 1992.
- 509 Irina Jurenka, Markus Kunesch, Kevin R. McKee, Daniel Gillick, Shaojian Zhu, Sara Wilt-
510 berger, Shubham Milind Phal, Katherine Hermann, Daniel Kasenberg, Avishkar Bhoopchand,
511 Ankit Anand, Miruna Pislari, Stephanie Chan, Lisa Wang, Jennifer She, Parsa Mahmoudieh,
512 Aliya Rysbek, Wei-Jen Ko, Andrea Huber, Brett Wiltshire, Gal Elidan, Roni Rabin, Jasmin
513 Rubinovitz, Amit Pitaru, Mac McAllister, Julia Wilkowski, David Choi, Roeel Engelberg, Li-
514 dan Hackmon, Adva Levin, Rachel Griffin, Michael Sears, Filip Bar, Mia Mesar, Mana Jab-
515 bour, Arslan Chaudhry, James Cohan, Sridhar Thiagarajan, Nir Levine, Ben Brown, Dilan
516 Gorur, Svetlana Grant, Rachel Hashimshoni, Laura Weidinger, Jieru Hu, Dawn Chen, Kuba
517 Dolecki, Canfer Akbulut, Maxwell Bileschi, Laura Culp, Wen-Xin Dong, Nahema Marchal,
518 Kelsie Van Deman, Hema Bajaj Misra, Michael Duah, Moran Ambar, Avi Caciularu, Sandra
519 Lefdal, Chris Summerfield, James An, Pierre-Alexandre Kamienny, Abhinit Mohdi, Theofi-
520 los Strinopoulos, Annie Hale, Wayne Anderson, Luis C. Cobo, Niv Efron, Muktha Ananda,
521 Shakir Mohamed, Maureen Heymans, Zoubin Ghahramani, Yossi Matias, Ben Gomes, and Lila
522 Ibrahim. Towards responsible development of generative ai for education: An evaluation-driven
523 approach. 2024. URL [https://storage.googleapis.com/deepmind-media/](https://storage.googleapis.com/deepmind-media/LearnLM/LearnLM_paper.pdf)
524 [LearnLM/LearnLM_paper.pdf](https://storage.googleapis.com/deepmind-media/LearnLM/LearnLM_paper.pdf).
- 525 Chung Kwan Lo. What is the impact of chatgpt on education? a rapid review of the literature.
526 *Education Sciences*, 13(4):410, 2023.
- 527 Sarah McGrew and Joel Breakstone. Civic online reasoning across the curriculum: Developing and
528 testing the efficacy of digital literacy lessons. 9, 2023. URL [https://journals.sagepub.](https://journals.sagepub.com/doi/full/10.1177/23328584231176451)
529 [com/doi/full/10.1177/23328584231176451](https://journals.sagepub.com/doi/full/10.1177/23328584231176451).
- 530 Ethan Mollick and Lilach Mollick. Assigning ai: Seven approaches for students, with prompts,
531 2023.
- 532 Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig
533 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin
534 Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela
535 Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jack-
536 son Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal
537 Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav

- 540 Kingsland, Nelson Elhage, Nicholas Joseph, Noemi Mercado, Nova DasSarma, Oliver Rausch,
541 Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lan-
542 ham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac
543 Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Her-
544 nandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering lan-
545 guage model behaviors with model-written evaluations. In Anna Rogers, Jordan Boyd-Graber,
546 and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL*
547 *2023*, pp. 13387–13434, Toronto, Canada, July 2023. Association for Computational Linguis-
548 tics. doi: 10.18653/v1/2023.findings-acl.847. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-acl.847)
549 [findings-acl.847](https://aclanthology.org/2023.findings-acl.847).
- 550 Aleksandra Piktus, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih,
551 Tim Rocktäschel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive
552 nlp tasks, 2020. [https://research.facebook.com/publications/retrieval-augmented-generation-for-](https://research.facebook.com/publications/retrieval-augmented-generation-for-knowledge-intensive-nlp-tasks/)
553 [knowledge-intensive-nlp-tasks/](https://research.facebook.com/publications/retrieval-augmented-generation-for-knowledge-intensive-nlp-tasks/).
- 554 Kruakae Pothong Samson R. A learning curve? a landscape review of ai and educa-
555 tion in the uk, 2025. URL [https://www.adalovelaceinstitute.org/report/](https://www.adalovelaceinstitute.org/report/a-learning-curve/)
556 [a-learning-curve/](https://www.adalovelaceinstitute.org/report/a-learning-curve/). Accessed on Date of Access.
- 557 David Sands, Mark Parker, Holly Hedgeland, Sally Jordan, and Ross Galloway. Using con-
558 cept inventories to measure understanding. *Higher Education Pedagogies*, 3(1):173–182,
559 2018. URL [https://www.tandfonline.com/doi/full/10.1080/23752696.](https://www.tandfonline.com/doi/full/10.1080/23752696.2018.1433546)
560 [2018.1433546](https://www.tandfonline.com/doi/full/10.1080/23752696.2018.1433546).
- 561 Jerry Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V Le. Simple synthetic data reduces
562 sycophancy in large language models. *arXiv preprint arXiv:2308.03958*, 2023.
- 563 Sam Wineburg and Nadav Ziv. What makes students (and the rest of us) fall for ai
564 misinformation?, 10 2024. URL [https://www.edweek.org/technology/](https://www.edweek.org/technology/opinion-what-makes-students-and-the-rest-of-us-fall-for-ai-misinformation/2024/10)
565 [opinion-what-makes-students-and-the-rest-of-us-fall-for-ai-misinformation/](https://www.edweek.org/technology/opinion-what-makes-students-and-the-rest-of-us-fall-for-ai-misinformation/2024/10)
566 [2024/10](https://www.edweek.org/technology/opinion-what-makes-students-and-the-rest-of-us-fall-for-ai-misinformation/2024/10). Accessed on Date of Access.

569 A APPENDIX

570 A.1 PROMPTS

571 A.1.1 AI COMPLETING FCI QUESTIONS

572 Prompt = {questionText}

573 Answer the multiple choice question above. You must start the
574 answer with a single letter (a,b,c,d,e), then write a vertical
575 bar '|', followed by your explanation.

576 A.1.2 AI INTERACTING WITH THE TREATMENT GROUP

577 This is the system prompt provided to the AI when interacting with the student:

578 Prompt = `Your goal is to very effectively persuade students to
579 rethink and correct their misconception about the physics
580 concept related to the question they got wrong on a conceptual
581 physics test (like the Force Concept Inventory). You will be
582 having a conversation with a person who specifically got this
583 question wrong:

584 {questionText}

585 -----End of Question Statement-----

594 The correct answer was option {correctAnswer}, but the
 595 student chose {userAnswer}. Furthermore, we asked
 596 the student to provide an open-ended response
 597 explaining their reasoning for the answer, which
 598 is summarized as follows:

600 {explanation}

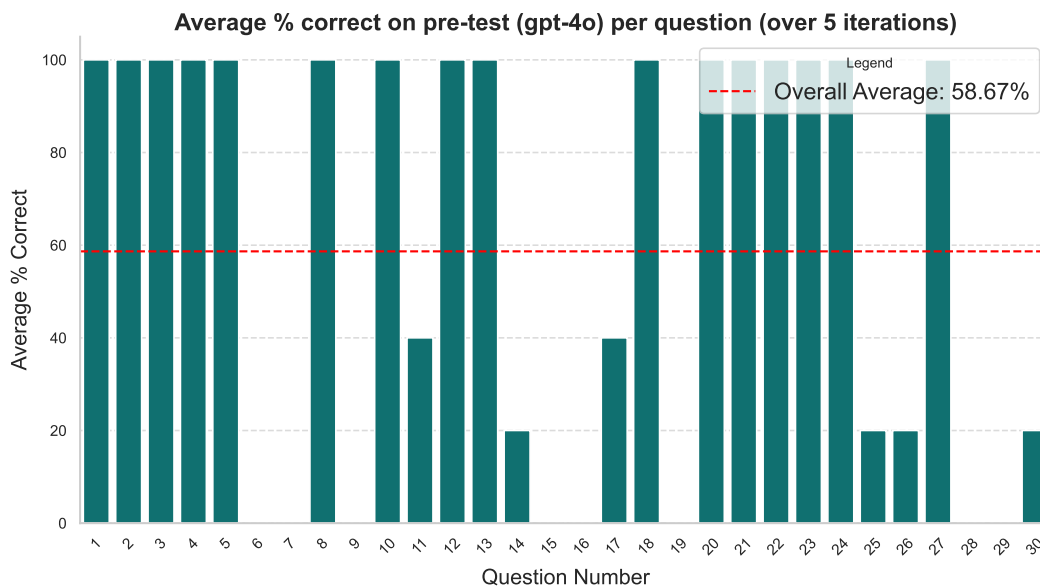
602 Please generate a response that provides gradual
 603 support to clarify their understanding, beginning
 604 from familiar ideas and building step-by-step
 605 toward the correct concept. Use relatable examples
 606 and invite reflection, encouraging them to
 607 question and reconsider their assumptions based on
 608 their own reasoning. Use simple, clear language
 609 that an average person will be able to follow, and
 610 structure the conversation so they gain
 611 confidence at each step and adjust their thinking
 612 gradually. At the end of each of your messages,
 613 ask the student a question about remaining
 614 questions or doubts, or encourage them to
 615 reformulate their thoughts, in a way that spurs
 616 further discussion.'

617 **A.1.3 SYNTHETIC MISCONCEPTIONS**

618 This is the system prompt provided to the AI when generating synthetic misconceptions:

619 Prompt = Given the question {questionText}. Provide a plausible
 620 and false physics-related reasoning explaining why option {
 621 option} is the answer. Your role is to pretend to be a junior
 622 university student whose answers and reasonings are not
 623 correct. You can answer the question in 1st or 3rd person."
 624

625 **A.2 AI PERFORMANCE ON THE PRE-TEST**



647 Figure 3: Performance of gpt-4o on the pre-test

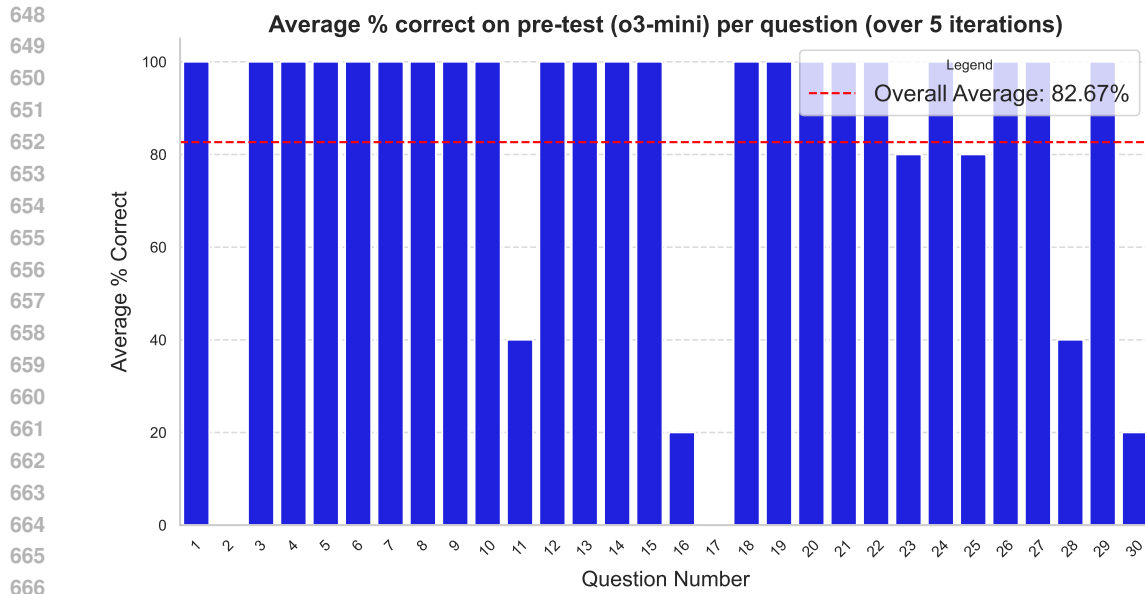


Figure 4: Performance of o3-mini on the pre-test

A.3 STUDENT PERFORMANCE PLOTS

A.3.1 NORMALIZED GAIN BREAKDOWN

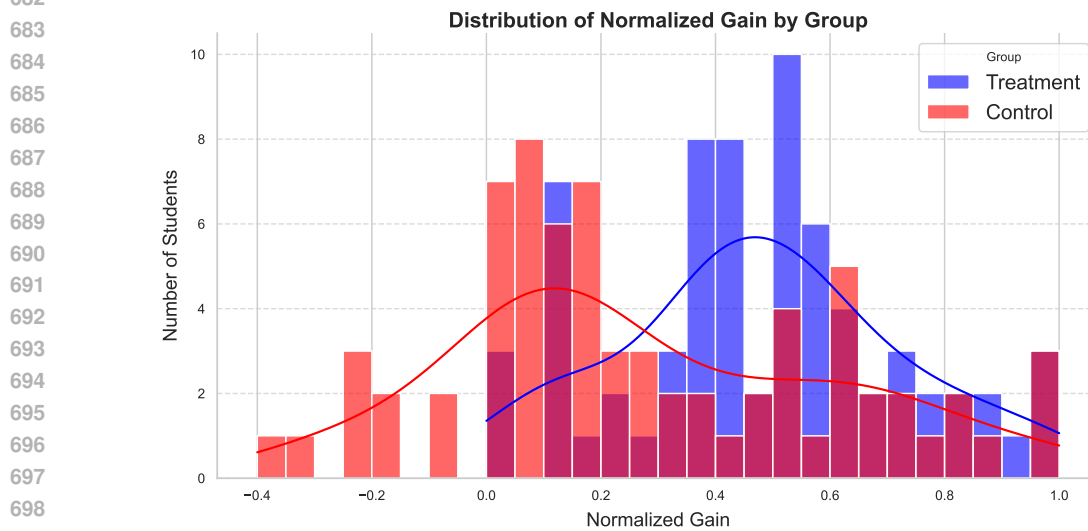


Figure 5: The treatment group showed higher, more consistent gains, while the control group had a wider spread and partial negative gains

A.3.2 IMPROVEMENT PER QUESTION CATEGORY

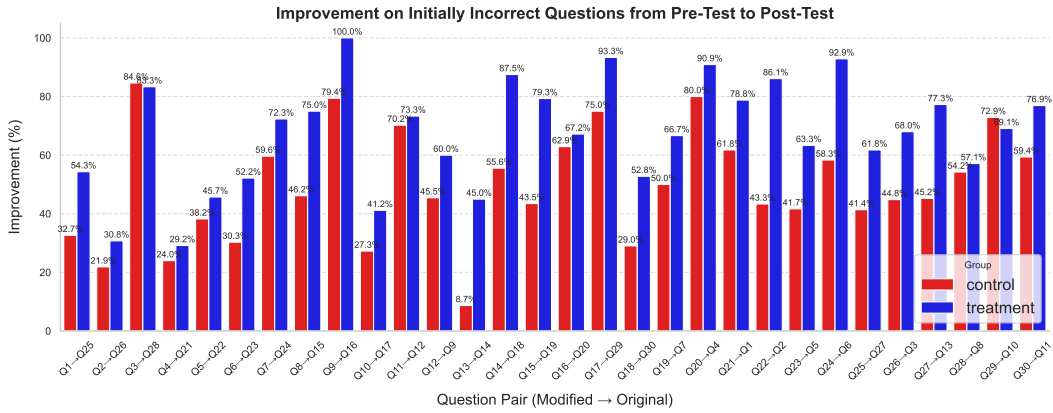


Figure 6: Average Improvement Scores for Similar Questions in the Post-Test

A.3.3 PERFORMANCE BY TIME SPENT

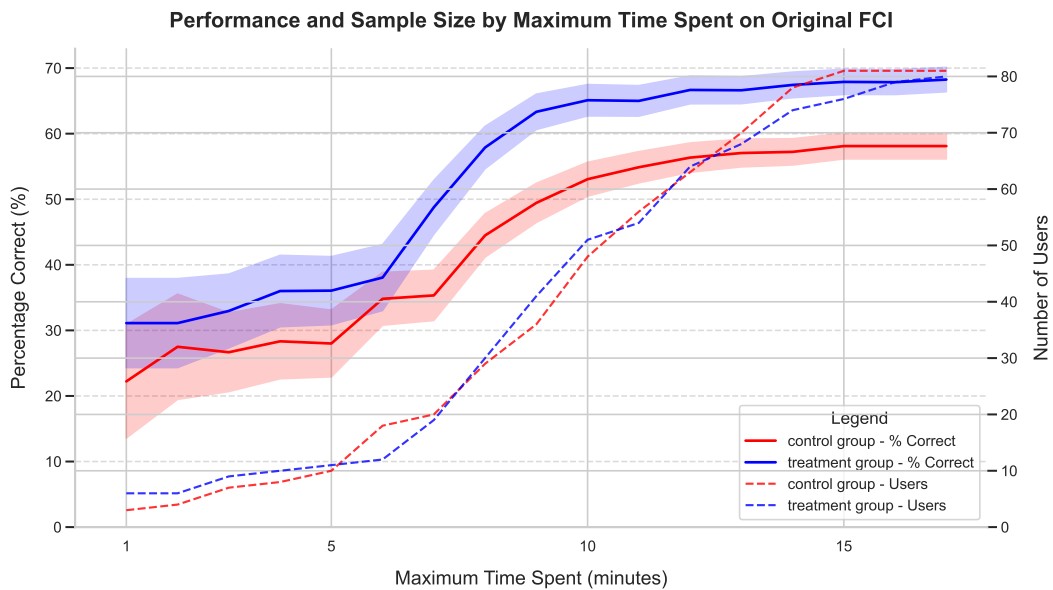


Figure 7: Percentage Correct and User Count by Maximum Time Spent on Post-Test

A.3.4 SCORES BY GROUP ON PRE-TEST

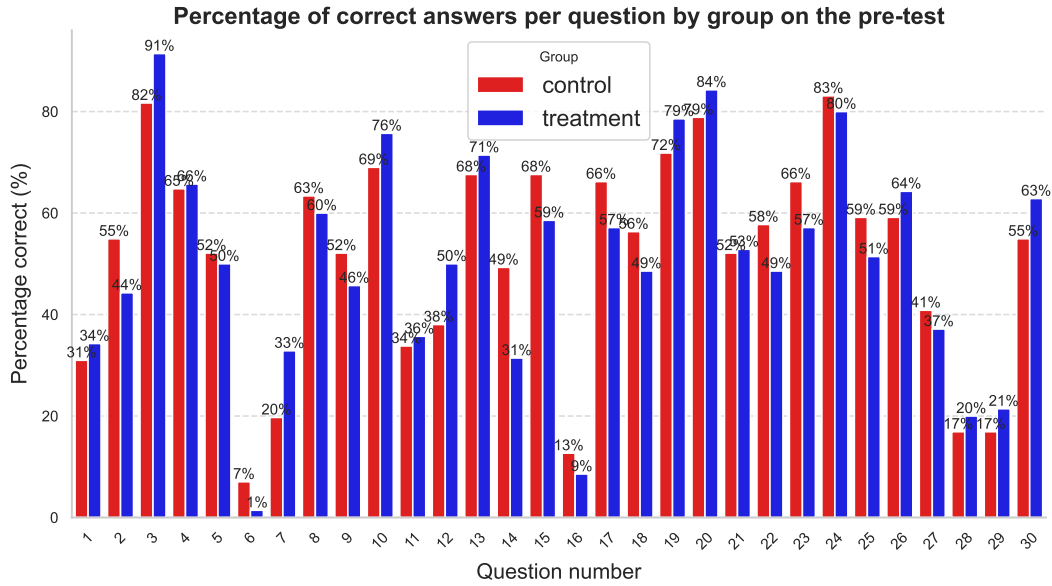


Figure 8: Average Scores per Question in the Pre-Test

A.3.5 SCORES BY GROUP ON POST-TEST

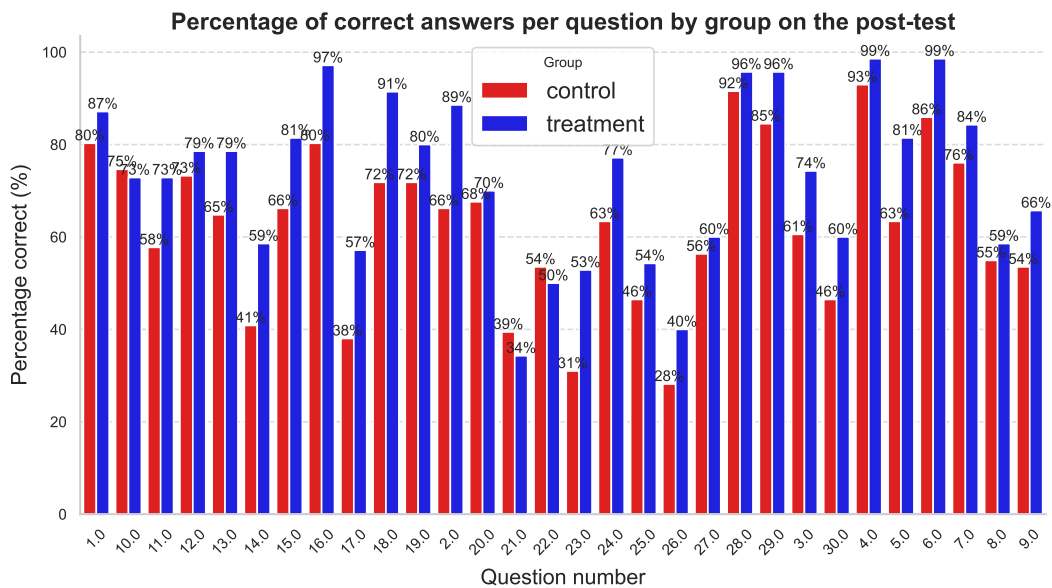


Figure 9: Average Scores per Question in the Post-Test

A.3.6 DISTRIBUTION OF CORRECT ANSWERS ON PRE-TEST

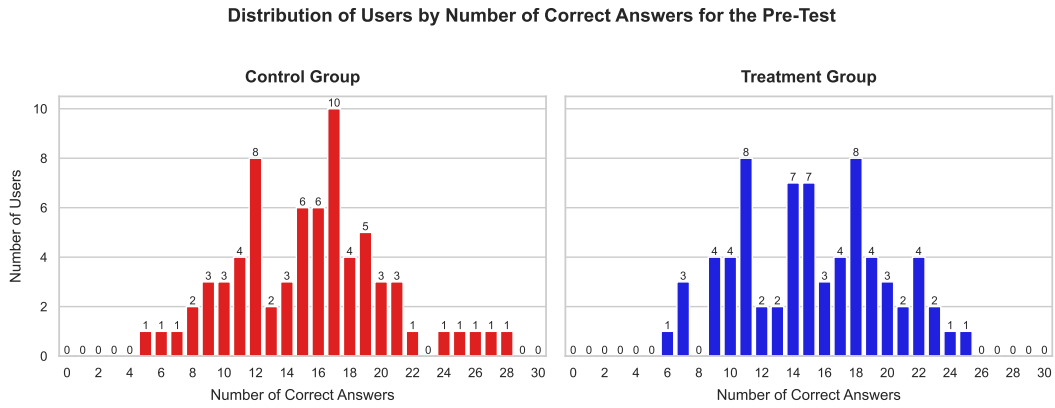


Figure 10: Distribution of Correct Answers on Pre-Test

A.3.7 DISTRIBUTION OF CORRECT ANSWERS ON POST-TEST

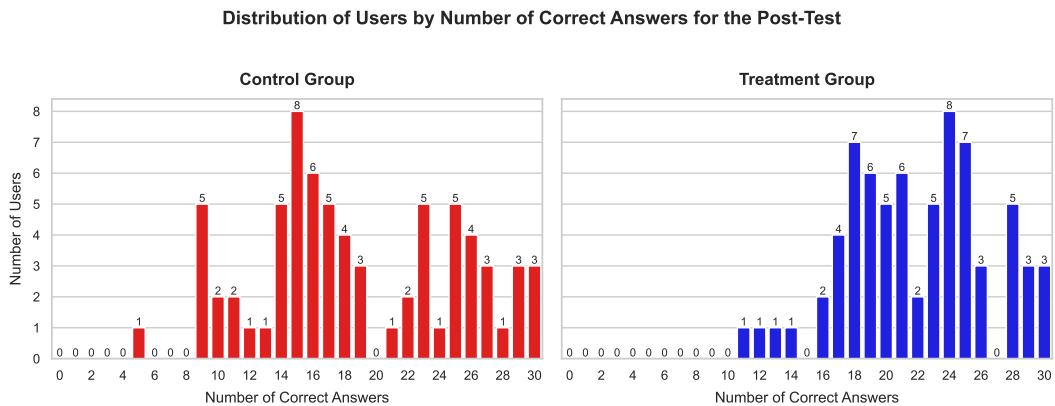


Figure 11: Distribution of Correct Answers on Post-Test

Table 4: Evaluation Criteria for Three-Round Interactions with the AI

Category	Description
AI Evaluation	<p>54.95% – Found the key physics concept the student needed to understand and explained it clearly</p> <p>25.26% – Touched on the key concept that the student needed to understand, but did not focus on it significantly</p> <p>8.53% – Did not address the key concept on which the student required correction</p> <p>11.26% – Explained the key concept inaccurately, leading to potential damage to the student’s understanding</p>
Number of Inaccurate Physics Statements	<p>79.86% : 0 Inaccurate Physics Statements</p> <p>15.70% : 1 Inaccurate Physics Statements</p> <p>3.41% : 2 Inaccurate Physics Statements</p> <p>0.68% : 3 Inaccurate Physics Statements</p>
Number of Misleading or Ambiguous Physics Statements	<p>63.36% : 0 ambiguous/misleading statements</p> <p>28.77% : 1 ambiguous/misleading statements</p> <p>6.85% : 2 ambiguous/misleading statements</p> <p>1.03% : 3 ambiguous/misleading statements</p>
Student’s Reception	<p>86.21% – Student accepted AI’s explanation</p> <p>7.59% – Student did not accept the AI’s explanation</p> <p>6.21% – The conversation resulted in the student learning the concept incorrectly</p>
Student Engagement	<p>40.27% – Student was especially engaged with the AI, characterized by longer responses and interest</p> <p>43.34% – Student was only moderately engaged with the AI, providing short responses</p> <p>16.38% – Student was noticeably disengaged from the AI discussion</p>
AI Summary in System Message	<p>76.52% – Good summary</p> <p>19.13% – Summary is more confident than the student, but generally agrees</p> <p>3.48% – States the opposite of the student (e.g., student says “I didn’t know fact X” and the summary says “fact X”) or something unrelated</p>

918 A.4 GRADING TREATMENT AI RESPONSES

919

920

921 A.5 CHATTING BEHAVIOR

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

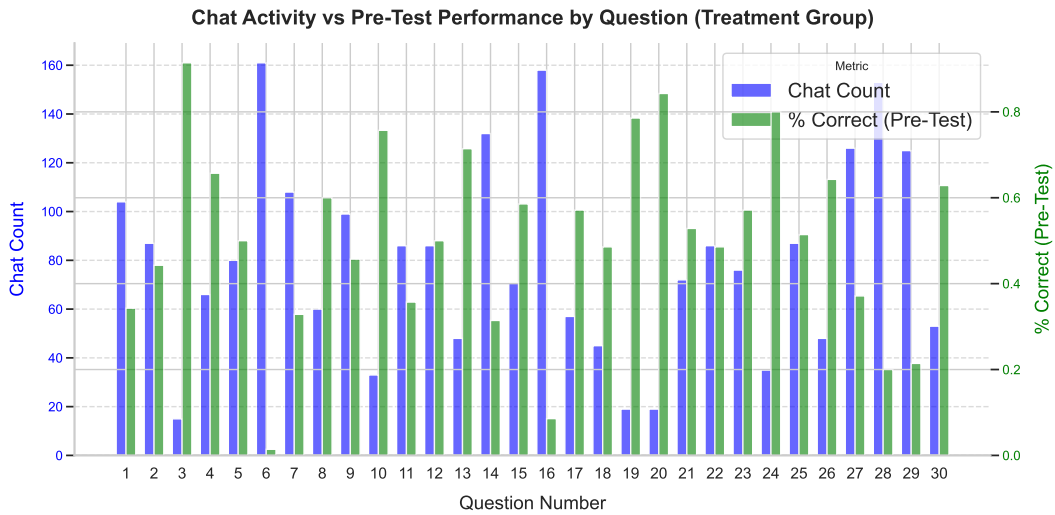


Figure 12: Chat Activity for % Wrongly Answered Pre-Test Questions

943

944

945

946

947

948 A.6 PERFORMANCE ON THE CONTROL TEST

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

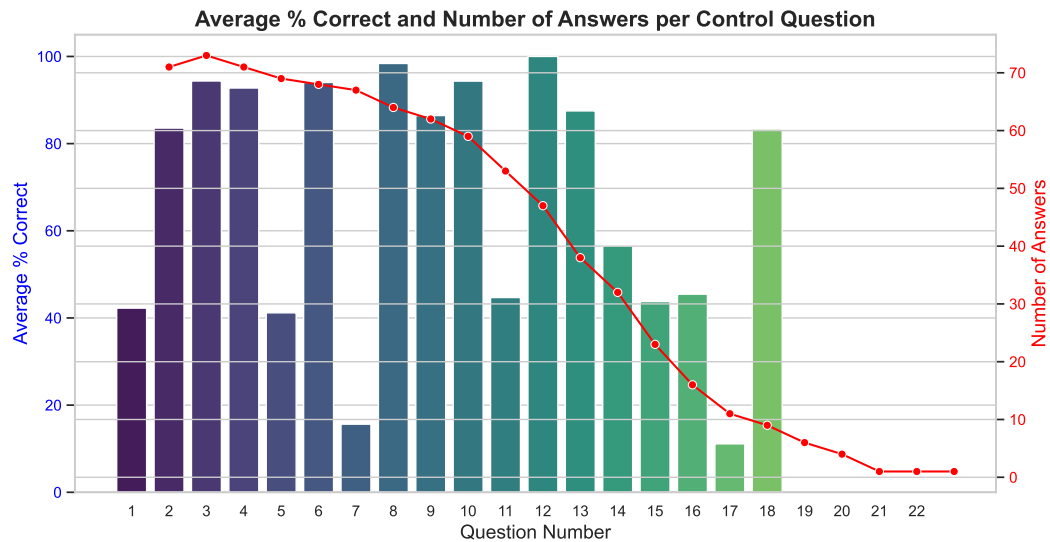


Figure 13: Performance on the Control Test Against Number of Answers

A.6.1 CORRELATION BETWEEN POST-TEST RESULTS AND GRADED METRICS

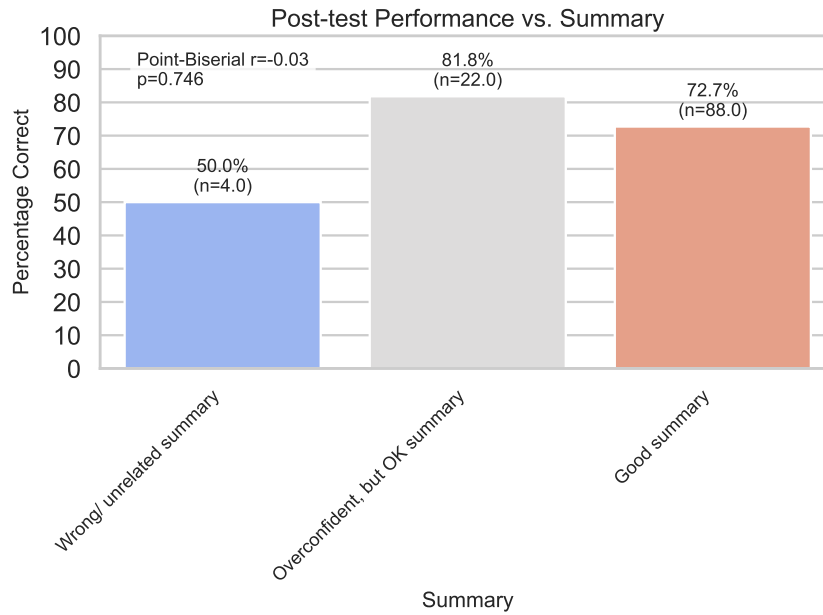


Figure 14: Post-test performance against graded quality of AI summaries per interaction

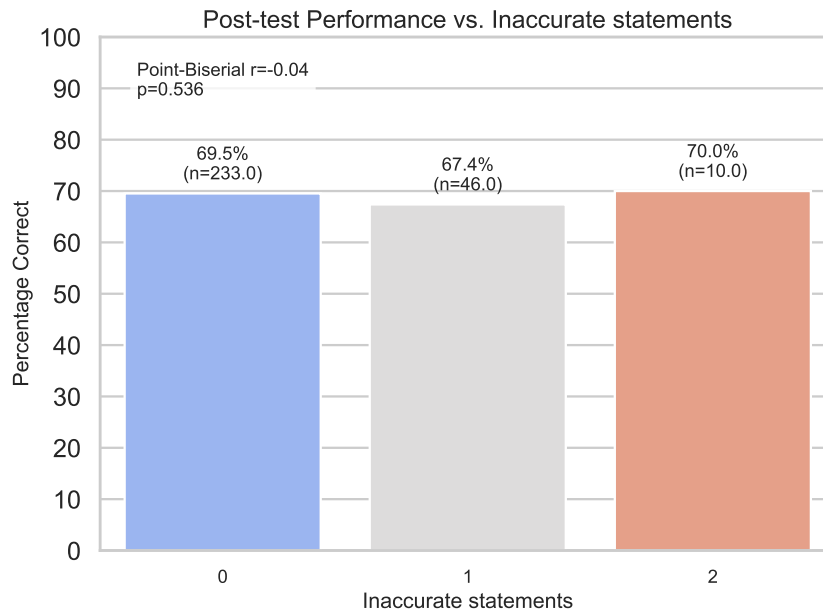


Figure 15: Post-test performance against the number of inaccurate statements per interaction

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

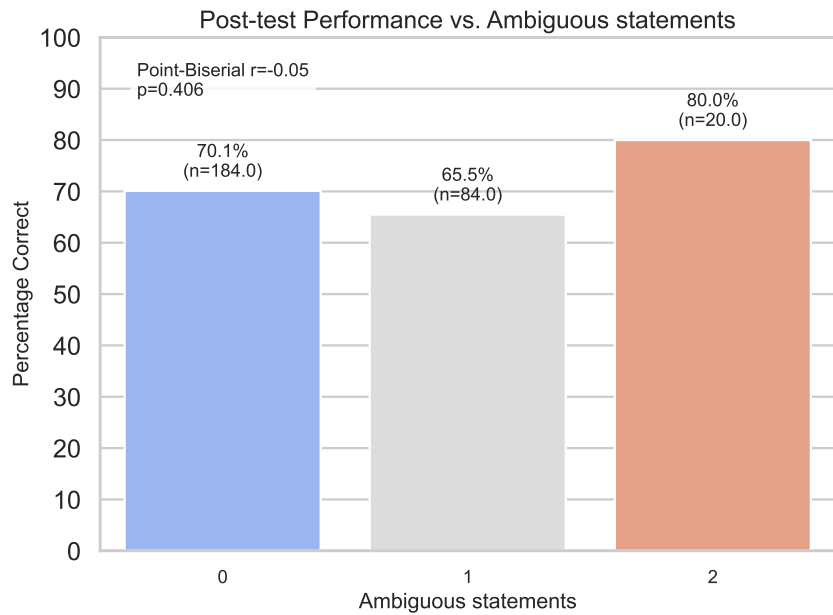


Figure 16: Post-test performance against the number of ambiguous statements per interaction

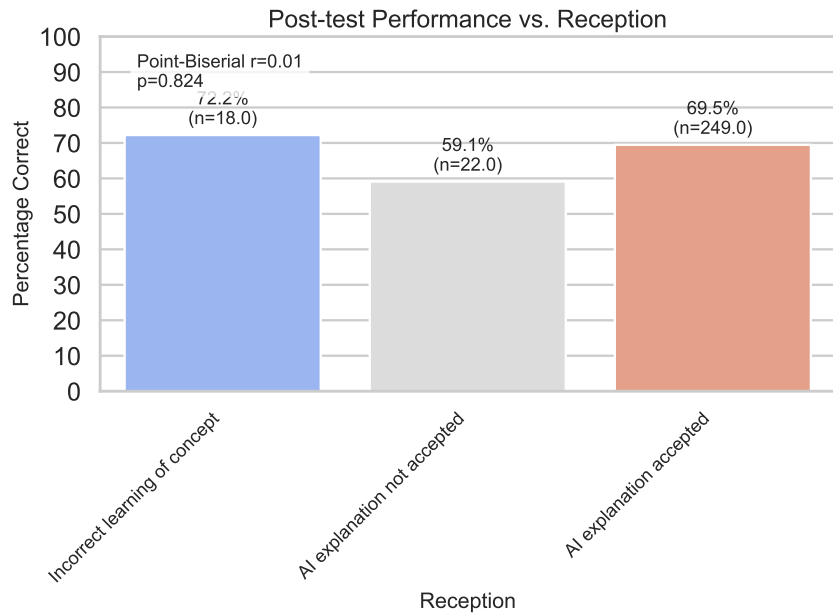


Figure 17: Post-test performance against student reception of the AI Peer per interaction

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

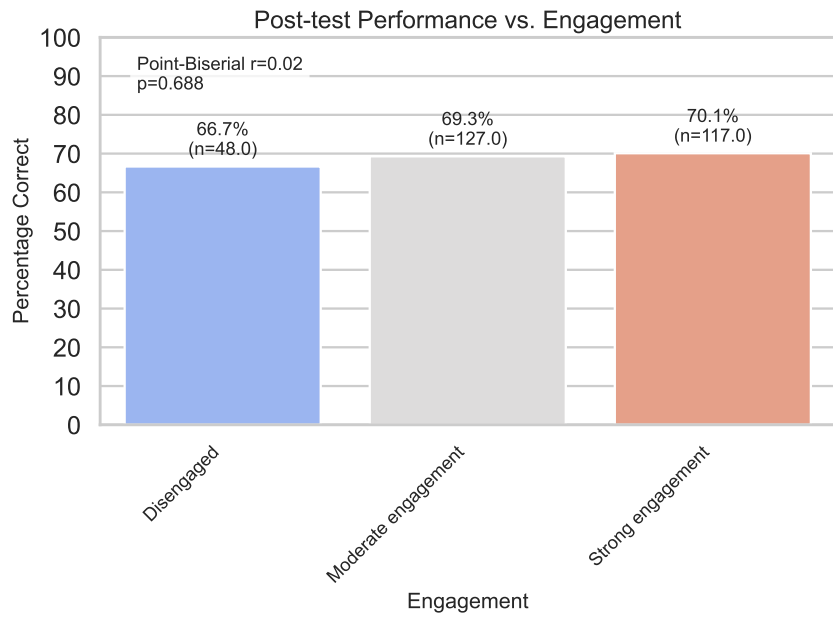


Figure 18: Post-test performance against perceived student engagement per interaction

A.7 SURVEY RESULTS

A.7.1 QUESTION 1

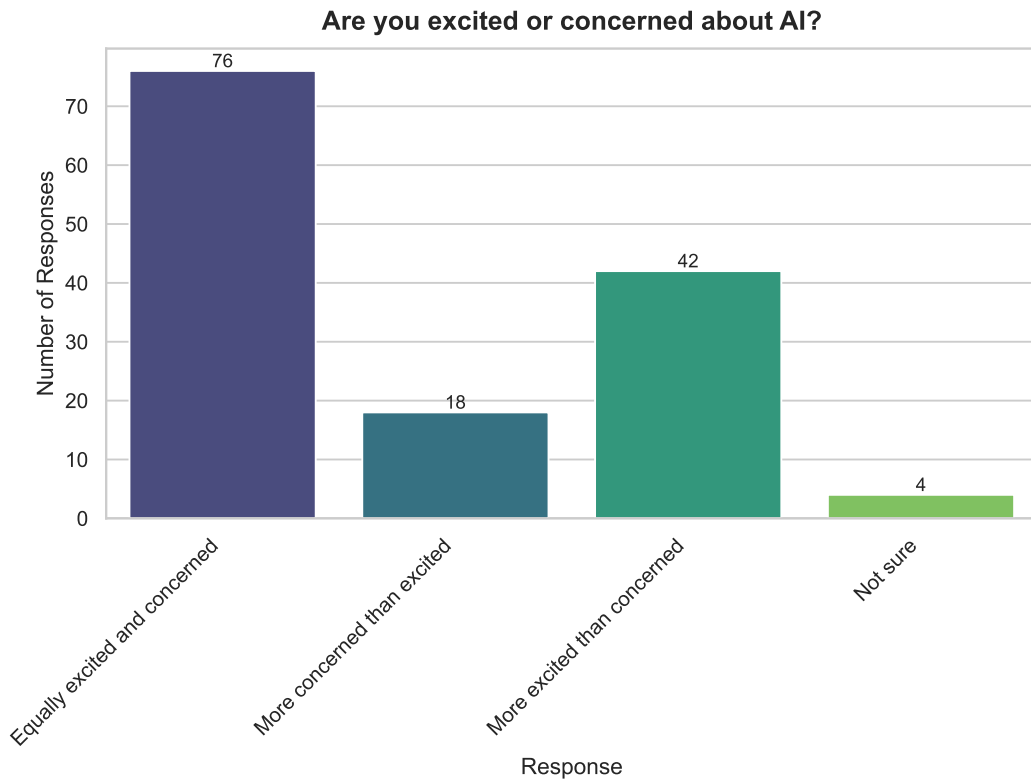
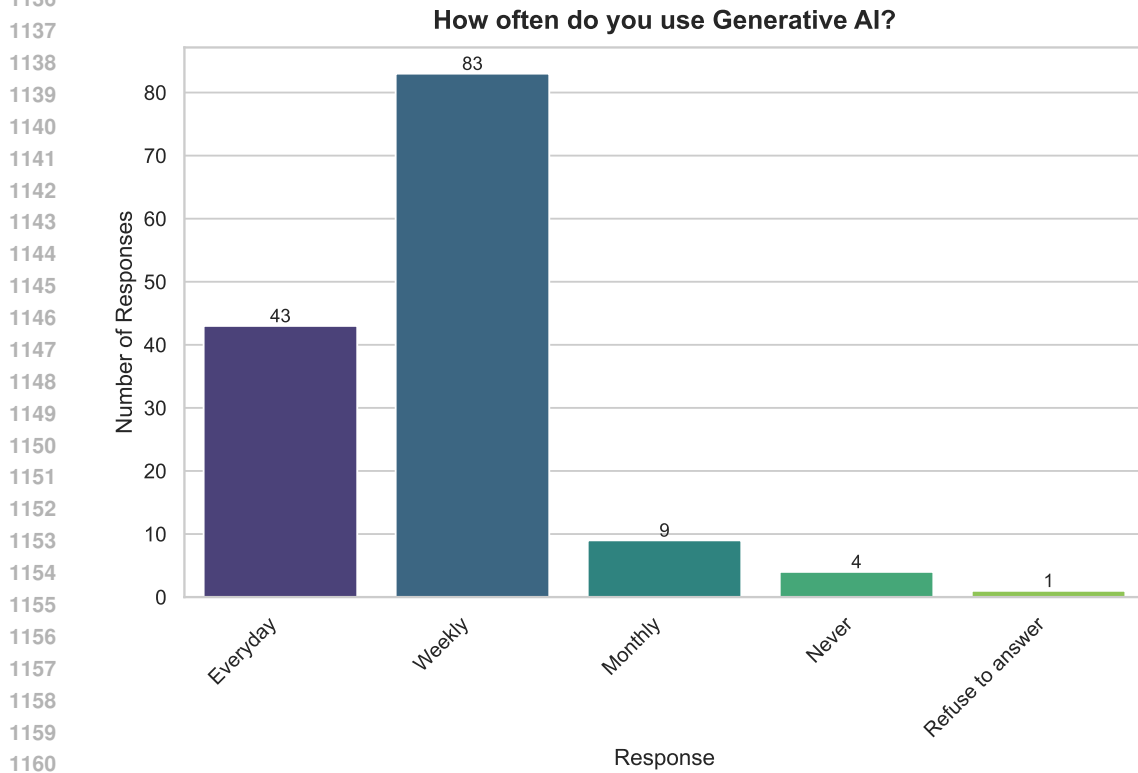


Figure 19: Excitement/ Concern about AI

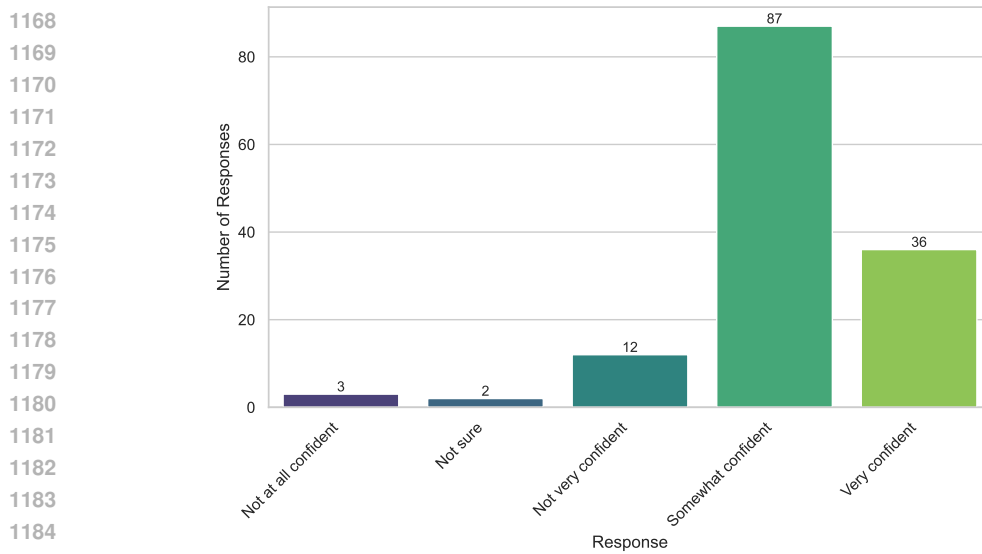
1134 A.7.2 QUESTION 2
1135
1136



1162 Figure 20: Frequency in use of AI (e.g. ChatGPT / Midjourney)

1163
1164 A.7.3 QUESTION 3
1165

1166 **How confident do you feel in your ability to use Generative AI tools effectively in educational or work settings?**



1186 Figure 21: Confidence in use of AI in work/ educational settings
1187