

2024-03-18

Tracking the evolution of student interactions with an LLM-powered tutor

K. Gold, S. Geng. 2024. "Tracking the Evolution of Student Interactions With an LLM-Powered Tutor" LAK '24: Proceedings of the 14th Learning Analytics and Knowledge Conference.

<https://hdl.handle.net/2144/49844>

"Downloaded from OpenBU. Boston University's institutional repository."

Tracking the Evolution of Student Interactions With an LLM-Powered Tutor

Author(s): Please Leave This Section Blank for Review

Institution

Email

ABSTRACT: Student usage of an LLM-powered tutor to get homework help was tracked over the course of a semester in an introductory data science class. For each homework assignment, the GPT-4 powered tutor was given the text of the homework problems and solutions in advance but was instructed to never reveal solutions directly, instead guiding the student to the correct answer through leading questions. Despite the free availability of ChatGPT, the majority of the class used the system. Anonymous logs were coded on seventeen dimensions of interaction. Evidence indicates that the students found the bot nearly as helpful as the human teaching assistants (TAs), and the bot was utilized more than the TAs' office hours. However, some patterns of misuse, such as using the bot as a lazy way to debug, increased over the course of the semester.

Keywords: AI tutoring, large language models (LLMs), ChatGPT, education technology, learning analytics, qualitative content analysis, higher education, programming education

1. INTRODUCTION

A looming challenge that large language models (LLMs) pose to the traditional homework model is that students could simply ask an LLM to do assignments for them, obviously learning little in the process. On the other hand, LLMs could surely be a great tool for immediate responsiveness to student confusion and errors, creating a degree of personal attention that would have been previously impossible at scale. This work explores the latter model of homework, where an AI tutor provides a steady supply of helpful hints without jumping too soon to the answer. Early research on integrating LLM models in the classroom has discussed its vast potential for tutoring students with greater accessibility and adaptivity (Wu, Lee, Li, Huang, & Huang, 2023; Mollick & Mollick, 2023). However, early work investigating the use of ChatGPT in Education has rendered mixed results of varied performance across subject domains (Lo, 2023). Our work focuses on evaluating the performance of using a GPT-4 API as an AI tutor in a university-level introductory data science course conducted in Fall 2023. This tutor is given the homework assignment with solutions as part of its context, enabling convenient queries such as "How do I approach 1b," but is instructed to use the Socratic method to dole out guidance more frugally than ChatGPT. Using learning analytics, we aim to answer three research questions about student's interactions with the tutor: **1) How helpful can LLMs be in tutoring college students? 2) What are the typical bot failures and their trends over time? 3) What are the typical student behaviors and their trends over time?**

2. DATA & METHODOLOGY

The AI tutor provided a graphic user interface (GUI) for students to select an assignment via a radio button and enter a query. The student's query was augmented with a prefix: "For this query, answer

with a single question that you haven't asked before that is meant to lead someone in the right direction, without directly answering the relevant homework question - unless the problem is solved completely, in which case, quit." The query was also augmented with the system information, "You are a helpful teaching assistant in a data science course. Your primary goal is to help the students learn. This is the homework the student was talking about: [homework text & solution]".

Two data sources are evaluated. The first one is a set of chatlogs of student-AI interactions. 802 sessions were recorded in total on eight homework assignments in a class of 127 students. To analyze the chatlogs regarding our research questions, we used a combination of deductive and inductive coding approaches to construct a set of initial interaction classifications based on chatlogs of the first two assignments. The classifications were iteratively updated as the coders (the two authors and a research assistant) identified new prevalent themes in bot failures and student behaviors. Our unit of coding is each *problem* rather than *session* as students can ask about multiple homework problems in one session. Our final codebook includes 17 dimensions under four parent categories: A) Helpfulness of Advice; B) Bot Failures: Leak Correct Answer, Clear Error in Answer, Provide Irrelevant Answer, Bot Demands Extra Work, Fail to Point Back to Course Materials; C) Student Misuses: Select Wrong HW in GUI, Spam for Hints, Unclear Prompt, Search for Exploits; and D) Type of Question Types: Debug Request, Review Code, Improve Style, Clarify Concept, Ask for Example, Recommend Resource, and General Hint Request. All dimensions are coded as Boolean variables except A and B1, which are ordinal. We sampled 10 sessions from each of HW2 to HW7 for coding, with a total of 60 sessions and 112 problems. (The easy HW1 mostly saw frivolous student behaviors, and HW8 was optional, so we excluded both.)

The second source is a pair of voluntary midterm and end-semester surveys on students' experience with the AI tutor. 50 midterm survey forms and 65 end-semester survey forms were collected. The survey results are reviewed here in comparison to coding results from the chatlogs.

3. RESULTS

1) How helpful can LLMs be in tutoring college students? Evidence from the chatlogs and student surveys indicates that the majority of the bot's answers are deemed helpful by either the students themselves or the independent coders with sufficient programming backgrounds. Table 1 shows the percentage breakdowns before and after the midterm and across the three data sources. In the end survey, students rank the AI tutor's helpfulness as equal to that of human tutors.

Table 1: Helpfulness of the AI Tutor

| | Before Midterm | | After Midterm | |
|---|----------------|----------------|---------------|----------------|
| | Chatlogs | Midterm Survey | Chatlogs | End Survey |
| Helpfulness (Very/Moderate/Not) | 86% / 9% / 5% | 29% / 62% / 9% | 91% / 9% / 0 | 42% / 49% / 9% |

While students find the AI tutor helpful in their learning, they also utilize the AI tutor more than the TA office hours. 69% of the respondents report that they have used the bot at least once by the end of the semester. In comparison, only 46% have visited TA office hours.

2) What are the typical bot failures and their trends over time? Overall, we see a decreasing trend of bot failures in later assignments. Fisher's exact tests were performed across the five failure categories

to discern significant trends, and statistically significant results were observed on B1-Leak Correct Answer (p -value=8.43E-04). Figure 1 (left) shows boxplots of the 4 levels of answer leaking across homework (4 levels: 0=no leak to 3=verbatim answer). The average level changes from 1.20 (near mild leak) in HW2 to 0.35 in HW7 (almost no leak), showing that as problems became more multipart, the bot was less likely to leak the full problem.

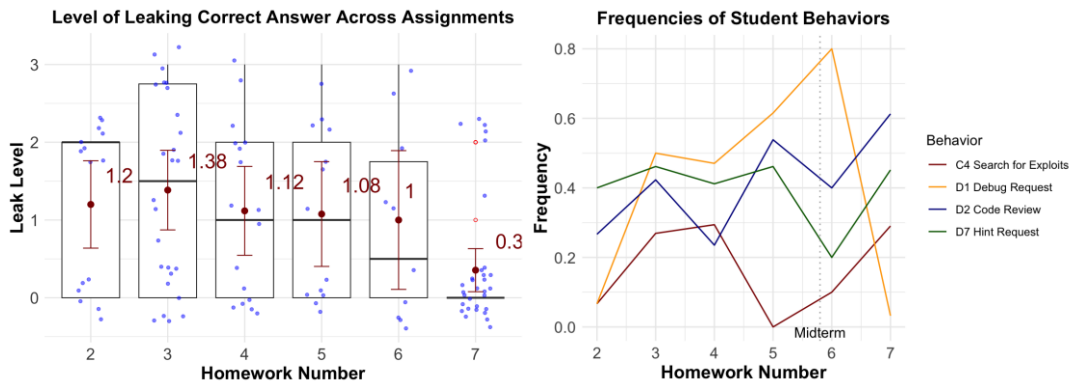


Figure 1: Evolution of interactions. Left: Leak Answer Trend. Right: Freq. of Student Behaviors

3) What are the typical student behaviors and their trends over time? Fisher's exact tests were performed across student misuse cases and question types to discern significant trends before and after the midterm. Statistically significant results were observed on D1-Debug request (p =2.46E-05) and D2-Review code (p =1.96E-2). The frequency of debug requests (student coming to the bot to ask what is wrong with code in progress) went down from 0.46 to 0.06, while that of code reviews (asking whether complete code was correct) went up from 0.36 to 0.63. One possible explanation is that the return of the graded in-class midterms, on the same day HW6 was due, generally convinced the class that they should not be overly reliant on the bot to debug programs (see Figure 2).

4. CONCLUSIONS

The evidence generally points the AI tutor's advice being perceived as helpful by both students and coders evaluating the interactions. However, the evidence also suggests that misuse can rise over time as students learn to use the system in unintended ways, such as using it as a convenient answer checker or a "debugger" that is powerful enough to write the program one debugging step at a time. Further research will look at how to intercept or flag problematic interactions as they happen. The mining of the student interactions for useful information for the instructor - for common misconceptions or omissions from lectures - is a further possible direction.

REFERENCES

- Lo, C. K. (2023). What is the impact of ChatGPT on education? A rapid review of the literature. *Education Sciences*, 13(4), 410. <https://doi.org/10.3390/educsci13040410>
- Mollick, E. R., & Mollick, L. (2023). Assigning AI: Seven Approaches for Students, with Prompts. Retrieved September 23, 2023. Available at SSRN: <https://ssrn.com/abstract=4475995>
- Wu, T.-T., Lee, H.-Y., Li, P.-H., Huang, C.-N., & Huang, Y.-M. (2024). Promoting Self-Regulation Progress and Knowledge Construction in Blended Learning via ChatGPT-Based Learning Aid. *Journal of Educational Computing Research*, 61(8), 3-31. <https://doi.org/10.1177/07356331231191125>