# Can GPT-4 Replicate Empirical Software Engineering Research?

JENNY T. LIANG, Carnegie Mellon University, USA
CARMEN BADEA, Microsoft Research, USA
CHRISTIAN BIRD, Microsoft Research, USA
ROBERT DELINE, Microsoft Research, USA
DENAE FORD, Microsoft Research, USA
NICOLE FORSGREN, Microsoft Research, USA
THOMAS ZIMMERMANN, Microsoft Research, USA

Empirical software engineering research on production systems has brought forth a better understanding of the software engineering process for practitioners and researchers alike. However, only a small subset of production systems is studied, limiting the impact of this research. While software engineering practitioners benefit from replicating research on their own data, this poses its own set of challenges, since performing replications requires a deep understanding of research methodologies and subtle nuances in software engineering data. Given that large language models (LLMs), such as GPT-4, show promise in tackling both software engineering- and science-related tasks, these models could help democratize empirical software engineering research.

In this paper, we examine LLMs' abilities to perform replications of empirical software engineering research on new data. We specifically study their ability to surface assumptions made in empirical software engineering research methodologies, as well as their ability to plan and generate code for analysis pipelines on seven empirical software engineering papers. We perform a user study with 14 participants with software engineering research expertise, who evaluate GPT-4-generated assumptions and analysis plans (i.e., a list of module specifications) from the papers. We find that GPT-4 is able to surface correct assumptions, but struggle to generate ones that reflect common knowledge about software engineering data. In a manual analysis of the generated code, we find that the GPT-4-generated code contains the correct high-level logic, given a subset of the methodology. However, the code contains many small implementation-level errors, reflecting a lack of software engineering knowledge. Our findings have implications for leveraging LLMs for software engineering research as well as practitioner data scientists in software teams.

## 1 INTRODUCTION

Empirical software engineering research on production systems has brought forth a better understanding of the software engineering process for practitioners and researchers alike. For example, empirical studies performed on Microsoft code bases have studied how bugs get fixed [25] and the effects of distributed teams on code quality [10], while studies performed at Bloomberg [38], Meta [17], and Google [52] have revealed key insights on implementing automated program repair systems and static analysis tools in practice.

Yet, only a small subset of production systems are rigorously studied. This limits the impact of empirical software engineering research, as software engineering practitioners are rarely able to reap the benefits from running analyses on their own data, since practitioners often lack the expertise and time to replicate empirical research. Additionally, barriers to the adoption of research include questionable assumptions and limited generalizability [44]. In recent years, data scientists have played an increasingly important role in software teams [36, 37] to help software engineers

understand more about topics such as productivity and code quality in their own teams and projects. Further, software engineers have reported an interest in obtaining answers to questions related to software development, spanning topics such as bug measurements, development practices, and testing practices [9, 32]. Thus, replicating empirical software engineering research could be a potential avenue for software teams to gain insights from their own software artifacts and data.

However, performing replications poses its own set of challenges, as it requires a deep understanding of research methodologies and subtle nuances in software engineering data [36]. While prior has work begun addressing this issue by creating domain-specific languages [33, 34] or programming environments [16] to help automate statistical analyses, these approaches do not directly address study replication, especially for software engineering contexts. Given that large language models (LLMs), such as GPT-4 [47], show promise in tackling software engineering- [31, 68] and science-related [59] tasks, these models may help replicate software engineering research studies. Studying an LLM's ability to replicate empirical software engineering research has the potential to help broaden the impact and democratize this research. This could allow developers to learn insights about their code bases and work habits, potentially helping to increase developer productivity.

In this paper, we examine GPT-4's abilities to perform replications on empirical software engineering research papers. Following the definition from the SIGSOFT Empirical Standards [2], *replication* is applying the same research methodologies from a given research paper on a different set of data. For a study to be replicated accurately, one must not only understand and recreate the conditions of the original study, but also develop knowledge of the research methodology, data, and code [37]. Therefore, being aware of the assumptions made in the original research and how to implement the methodology is essential. If assumptions of an empirical study are not met, the validity of the results can be compromised [12, 49]. Thus, we ask the following research questions:

**RQ1** Can GPT-4 identify assumptions from research methodology?

**RQ2** Can GPT-4 generate an analysis pipeline to replicate research methodology?

To answer these questions, we evaluated GPT-4's ability to generate assumptions, analysis plans (i.e, a list of module specifications), and code on seven empirical software engineering papers. We ran a user study with 14 software engineering researchers, who evaluated GPT-4-generated assumptions and analysis plans. We then perform a manual analysis of the GPT-4-generated code. We find that GPT-4 surfaces mostly correct assumptions, but can struggle to generate ones that reflect common but implicit knowledge about software engineering data. Further, we observe that GPT-4 can generate analysis plans that correctly outline the modules needed to replicate the study, but are limited by the quality and detail of the methodology as written the original research paper. Finally, we find that the GPT-4-generated code contains the correct high-level logic, given a subset of the methodology. However, the code contains many small implementation-level errors, such as looking at the incorrect tables in a database. Our findings have implications for leveraging LLMs for software engineering research, such as teaching GPT-4 software engineering domain knowledge.

## 2 RELATED WORK

Below, we discuss prior research that leverages LLMs for software engineering (Section 2.1) and science (Section 2.2). Since the field is developing quickly, the discussion below offers a snapshot of the field as of September 2023.

### 2.1 Language Models for Software Engineering

Language models have been applied to many tasks in software engineering. The most prominent task is code generation, where models such as Codex [13] have shown strong performance [65] in providing code suggestions to developers, given natural language or code. With the emergence of

publicly accessible LLMs such as GPT-4 [47], LLMs have been applied to a wide variety of software engineering tasks. Zheng et al. [68] performed a survey of 123 papers and identified seven software engineering tasks LLMs have been applied to: code generation, code summarization (i.e., generating comments for code), code translation (i.e., converting code from one programming language to another), vulnerability detection (i.e., identifying and fixing defects in programs), code evaluation and testing, code management (i.e., code maintenance activities such as version control), and Q&A interaction (i.e., using Q&A platforms such as StackOverflow). In a separate literature review of 229 research papers on language models for software engineering, Hou et al. [31] found that LLMs had been used on a large variety of software engineering datasets, which included data such as source code, bugs, patches, code changes, test suites, Stack Overflow, API documentation, code comments, and project issues. They found the papers spanned themes such as software development (e.g., API recommendation [62]), software maintenance (e.g., merge conflict repair [67]), software quality assurance (e.g., flaky test prediction [19]), requirements engineering (e.g., requirements classification [30]), and software design (e.g., software specification synthesis [45]).

Recently, OpenAI released the Code Interpreter feature [1] for ChatGPT that allows ChatGPT to write and execute Python code, which could potentially be used to generate analyses. Our work extends our understanding for tools that generate code, such as Code Interpreter, by observing how LLMs generate code for such analyses.

Thus, based on this literature, LLMs can handle a wide variety of software engineering tasks and data. However, these approaches can require fine-tuning models or specialized approaches. Our study extends from this literature by examining whether pre-trained LLMs like GPT-4 reflect this software engineering domain expertise off-the-shelf without any additional training.

## 2.2 Language Models for Science

Additional work has investigated the using language models for science. Prior work has shown that LLMs can reflect factual knowledge in several scientific domains. Language models such as MultiVerS [59] can validate claims against scientific literature in the domains of COVID-19, public health, and climate change. However, Auer et al. [4] found that ChatGPT struggled to answer challenging questions derived from research papers across many topics, including computer science, engineering, chemistry, geology, immunology, genetics, economics, and urban studies.

Similar studies have been performed in computer science. In natural language processing (NLP), Gao et al. [23] investigated whether LLMs could generate a survey of knowledge related to NLP concepts, such as A* search. In an evaluation with NLP experts, the authors found that GPT-4 could generate reasonable explanations of these concepts but could sometimes generate factually incorrect knowledge. Researchers have also applied LLMs for scientific research in human-computer interaction (HCI). Wu et al. [63] replicated seminal crowdsourcing papers in HCI using LLMs, while other work has found that GPT-3 could generate synthetic HCI data for both open- [28] and closed-ended [55] questions in interviews and surveys. Lastly, Xiao et al. [64] found that GPT-3 could be applied to perform deductive qualitative coding on datasets.

Our work builds upon the language models for software engineering literature by studying the extent to which LLMs can replicate research papers via analyzing methodology and writing code pipelines to repeat the analysis, rather than returning factual knowledge or generating research data. Further, compared to prior work, our study specifically focuses on quantitative empirical research methods, rather than qualitative ones. Finally, we extend our understanding of LLMs' scientific knowledge by studying its performance in the domain of software engineering research.

Table 1. Summary of the papers selected for GPT-4 to generate assumptions, analysis plans, and code. We report each paper's venue, number of citations in the ACM Digital Library, and a brief description of the paper's analysis. We also report number of assumptions and modules generated by GPT-4.

| Paper | Venue | Citations | Analysis Description | # of GPT-4 Generated | |
|---|---|---|---|---|---|
| | | | | Assumps. | Modules |
| [35] | ICSE'11 | 94 | Analysis of API-level refactorings and bug fixes in large open-source projects. The authors identify bug fix revisions and bug-introducing changes and use Change Distiller [20] to compute syntactic program differences. | 18 | 4 |
| [54] | ICSE'16 | 85 | Analysis of performance-related issues of JavaScript projects on GitHub. The authors identify performance-related issues and filter the issues to a set of changes that result in statistically significant performance improvements. | 17 | 4 |
| [18] | MSR'11 | 111 | Analysis of bug-fixing commits. The authors find all bug-fixing commits, identify lines that changed in each bug-fixing commit, and find the commit responsible for the previous version for each of the changed lines. | 15 | 4 |
| [27] | MSR'14 | 168 | Analysis of sentiments expressed in commit comments on GitHub. The authors use sentistrength to analyze commits of popular GitHub repositories. | 14 | 2 |
| [48] | MSR'14 | 103 | Analysis of the sentiment in security-related comments on commits and pull requests. The authors filter for security-related comments through keyword search and then perform sentiment analysis using nltk. | 13 | 3 |
| [22] | ESEC/FSE'22 | 1 | Analysis of review comments in pull requests on popular Java projects. The authors use a Hurdle model to understand the impact of file position during code review on number of comments it receives. | 16 | 3 |
| [56] | ICSE'22 | 5 | Analysis of commit messages on large GitHub projects. The authors then train two classifiers to classify commit messages with "what" and "why" information with various types of machine learning models. | 18 | 3 |

## 3 METHODOLOGY

To answer the research questions, we selected seven empirical software engineering papers (Section 3.1). We leveraged LLMs to automate the tasks that data scientists would perform in today's practice to replicate an empirical study in their own context: analyze the assumptions of the methodology, plan the analysis pipeline, and implement the code. Therefore, we prompted GPT-4 to generate assumptions, analysis plans, and code for each one of the papers (Section 3.2). To evaluate the model's outputs, we performed a user study (Section 3.3) with 14 participants with software engineering research expertise to evaluate the assumptions and analysis plans generated by the model. We then performed a manual evaluation to evaluate the generated code (Section 3.4). Finally, we performed quantitative and qualitative analysis on the collected data (Section 3.5). Materials used in this study, such as the protocols, GPT-4 generated data, and exact prompts, are available in the supplemental materials [5].

## 3.1 Paper Selection

We used the following process to select empirical software engineering papers to prompt GPT-4 to generate assumptions, analysis plans, and code. This process yielded seven research papers (see Table 1). The selection criteria for the research papers were:

- **Has a quantitative empirical analysis on software engineering data**. We focused on analyses that could be replicated through code rather than through manual means and could be derived from Git and GitHub data.
- **Is after 2010**. GitHub was created in 2010. Since the generated code relies on a Git and GitHub database, we identified research papers that utilized a similar type of data.
- **Has an approximately 1-page methodology section**. We ensured the methodology was short enough in length for participants to read through and evaluate three of them in a 1-hour long user study.

To obtain a diverse set of empirical software engineering papers to evaluate on GPT-4, the first author applied the selection criteria to three different sets of papers:

(1) **Empirical papers in software engineering venues**. Software engineering conferences contain numerous empirical studies on software engineering data. Thus, we searched on the ACM Digital Library for papers with "empirical" in the *title*, had the term "software engineering" in the *publication venue*, and whose *content type* was "Research Article". We then sorted by citation and selected the first 20 results to limit the search results. We identified 2 papers from this set [i.e., 35, 54].

(2) **Papers from the International Conference on Mining Software Repositories**. The International Conference on Mining Software Repositories (MSR) is a popular venue for publishing empirical software engineering research. Therefore, we searched on the ACM Digital Library for papers whose *publication venue* was MSR and whose *content type* was "Research Article". We then sorted by citation and selected the first 20 results to limit the search results. We identified 3 papers from this set [i.e., 18, 27, 48].

(3) **Papers published after 2021**. Since GPT-4 was pre-trained during September 2021 [47], we sampled papers that were created after the pre-training period to reduce any potential biases that could be influenced by pre-training. In particular, we applied the selection criteria to ACM Distinguished Paper Awards from top conferences with empirical software engineering studies: ESEC/FSE'22, ICSE'22, ICSE'23, MSR'22, and MSR'23. We identified 2 papers from this set [i.e., 22, 56].

## 3.2 Prompting GPT-4

We present an overview of the prompting strategy applied to GPT-4 in Figure 1. To answer the research questions, we prompted GPT-4 to generate assumptions (**RQ1**). We also prompted GPT-4 to generate analysis plans (i.e., a list of module specifications) as well as code to implement the analysis plan (**RQ2**). We separated the analysis pipeline into two steps—analysis plans and code—to distinguish the higher-level abstraction of creating modules from the lower-level details of writing implementations. This is because code implementation is also influenced by higher levels software design, like modules [41]; thus, studying both abstractions and implementations could reveal a more holistic understanding of GPT-4's ability to generate analysis pipelines.

In each of the prompts, we first provided GPT-4 with one to two example input-output pairs for few-shot learning of the task and output format. Next, we provided text-based instructions for GPT-4 (see Figure 1). The instructions included an explanation of GPT-4's *role* as a software engineering researcher; *skills* it has, such as in software engineering replication or in MySQL; and additional *information* about the given task, such as a description of the inputs. Finally, the prompt
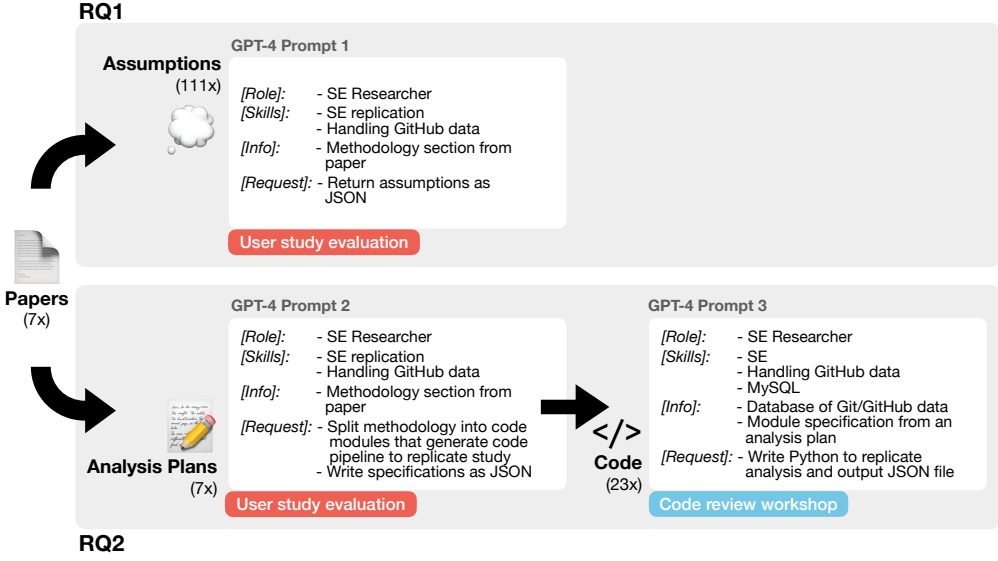
Fig. 1. An overview of the prompts used in the study. To answer RQ1, we used Prompt 1 to generate assumptions for each paper and evaluated them in a user study. To answer RQ2, we used Prompt 2 to create an analysis plan and evaluated it in a user study. We also used Prompt 3 to create the code modules of the analysis plan and evaluated it in a code workshop. The prompts in the figure are only summaries of the actual; for the complete prompts, see the supplemental materials [5].
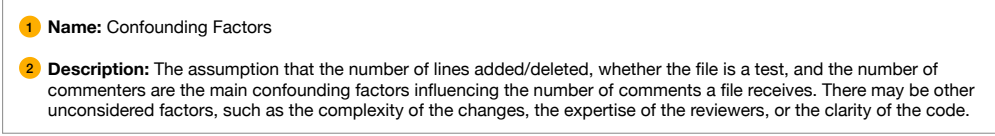


Fig. 2. Example GPT-4 generated assumption, given the methodology from Fregnan et al. [22]. Assumptions contain a name ① and a description ②.

included a *request* to extract, gather, or derive the specific data, following the format provided in the example input-output pairs. We ran the outputs on GPT-4 on the OpenAI Python API in September 2023 with the default model parameters, except for temperature, which was set to 0 to have deterministic outputs. The exact prompts and the generated assumptions, analysis plans, and code are available in the supplemental materials [5].

*3.2.1  Data Format.* Below, we further elaborate on of each of the generation types. We report the number of each data type generated by GPT-4 using the multiplication symbol (×).

*Assumptions (*111×*).* The GPT-4 generated assumptions contain two pieces of information: a *name* and a *description* (see Figure 2). We prompted GPT-4 to generate assumptions that underlie the given research methodology. Next, we prompted GPT-4 to generate assumptions about applying the methodology to a different dataset.

*Analysis Plan (*7×*).* Given a research paper methodology, GPT-4 generates an analysis plan that is represented as a list of code modules. It generates a code module has a *title*, *input* (i.e., which may include one or more outputs from other modules in the analysis plan), *output*, *description*, and

(1) **Name:** Data Extraction

(2) **Description:** This module is responsible for extracting the required data from the dataset provided by Gousios. The data includes commit comments from 90 of the top-starred software projects for the top programming languages on Github.

(3) **Inputs:** None

(4) **Outputs:** Extracted commit comments

(5) **Methodology text:** We extracted our analysis data from the dataset provided by Gousios [2], this dataset includes 90 of the 10-top starred software projects for the top programming languages on Github. We analyzed a total 60425 commit comments.

Fig. 3. Example GPT-4 generated module, given the methodology from Guzman et al. [27]. Modules contain a name (1), description (2), inputs (3), a description of outputs (4), and corresponding methodology text (5).

```
1  import json
2  import nltk
3  from nltk.sentiment import SentimentIntensityAnalyzer
4  from nltk import download
5  download('vader_lexicon')
6
7  # Load the comments
8  with open('./outputs/identification_of_security-related_comments_and_discussions.json', 'r') as file:
10     comments = json.load(file)
11
12 # Initialize the sentiment analyzer
13 sia = SentimentIntensityAnalyzer()
14
15 # Perform sentiment analysis on each comment
16 sentiment_analysis_results = []
17 for comment in comments:
18     sentiment = sia.polarity_scores(comment)
19     sentiment_analysis_results.append({
20         'comment': comment,
21         'sentiment': sentiment
22     })
23
24 # Save the sentiment analysis results to a JSON file
25 with open('./outputs/sentiment_analysis.json', 'w') as file:
26     json.dump(sentiment_analysis_results, file, indent=4)
```

Fig. 4. Example GPT-4 generated code, given the module specification from the analysis plan for Pletea et al. [48]. Generated code reads data from JSON file or database (1), runs additional logic, then outputs the result as a JSON object (2).

*corresponding methodology text* (see Figure 3). We prompted GPT-4 to divide the methodology text into a set of code modules and generate specifications following the above metadata.

*Code (23×).* Given a module specification, GPT-4 generates a piece of Python code that implements the specification (see Figure 4). We chose to instruct GPT to generate code that would output data as a JSON object into a file, so that it may be reused by other modules by reading in the file. However, the code may also make queries to an existing database filled Git and GitHub data following predefined schema; the schema is available in the supplemental materials [5]. To assist with code generation of downstream modules, in addition to prompting GPT-4 to generate code, we also instruct it to return an example JSON object and a natural language description of the JSON object. The example object and description from all previous modules is provided in the prompt for downstream modules in case they depend on it. This way modules that rely on earlier modules in the pipeline are able to update the *input* in their specifications to match this data representation.

## 3.3    User Study

We performed a user study with software engineering experts to validate the GPT-4-generated assumptions (**RQ1**) and analysis plans (**RQ2**) for each paper. The user study was reviewed and approved by our institution's Institutional Review Board. The survey and interview instruments used in the user study protocol (Section 3.3.2) are included in the supplemental material [5].

*3.3.1    Participants.* For the user study, we relied on snowball sampling to recruit participants with software engineering research experience at a large company that uses technology to deliver its services (see Table 2). Our inclusion criteria were people who obtained a Ph.D. studying software engineering-related topics and people who were software engineering data scientists.

*Recruitment.* We compiled a list of potential participants in the co-authors' network who met the inclusion criteria. We then sent them invitations to participate in the study. Potential participants suggested other individuals who met the inclusion criteria, who were also sent study invitations. In total, 24 invitations were sent, with 14 participants participating in the user study. We note that having 14 participants allowed each paper to have three to four participants review the GPT-4 generated assumptions and analysis plans, reflecting common practices in academic peer-review. Additionally, while performing open-coding on the interview data (Section 3.5), no new codes were added after 8 interviews, indicating that 14 participants was sufficient to achieve saturation.

*Demographics.* Participants (7 women, 7 men) were largely located in the United States and Canada ($N = 12$), with a few participants also located in India ($N = 2$). Overall, participants were experienced in software engineering-related research. Participants reported having between 2 to 14 years of software engineering research experience ($\mu = 7.1$ years) and 5 to 20 years of programming experience ($\mu = 8.4$ years). Additionally, participants reported publishing between 0 to 50 publications in top-tier software engineering venues ($\mu = 8.2$ publications) and serving as a reviewer 0 to 20 times at such venues ($\mu = 8.2$ times).

Participants also reported being familiar with LLMs, with 86% of participants using these models at least on a weekly basis. Participants also reported an ability to analyze and evaluate AI applications based on a validated AI literacy instrument [60]. All participants reported being able to choose a proper solution when presented with multiple solutions from AI agents. Additionally, 93% of participants reported being able to evaluate the capabilities and limitations of AI applications after using them. Finally, 71% reported being able to select an appropriate AI for a particular task.

*3.3.2    Protocol.* The user study consisted of a survey to gather participants' background information and an interview to collect participants' evaluations of the GPT-4-generated assumptions and analysis plans. Participants were compensated with a $50 gift card.

*Survey.* We designed a brief 10-minute Microsoft Forms survey to collect participants' demographics, programming background, research background, and familiarity with AI systems. Example topics included: how long participants had performed software engineering research, how many software engineering venues they had reviewed, how many years of programming experience they had, and how often they used LLMs.

We collected information about participants' gender following best practices from the HCI Guidelines for Gender Equity and Inclusivity [53]. To collect information on participants' AI literacy, we used the validated instrument from Wang et al. [60]. Because participants would be evaluating GPT-4 outputs in the interview, we used questions related to the evaluation construct to understand the degree to which they could evaluate the strengths and weaknesses of AI models.

Table 2. An overview of the participants in the user study. We report each participant's number of publications in top-tier software engineering venues, number of times served as a reviewer in software engineering venues, and years of experience doing software engineering research as well as programming. We also report the frequency of large language model usage, location, gender, and job position.

| ID | # of | | # of Years | | LLM Usage | Location | Gender | Position |
| | Pubs. | Revs. | Research | Coding | | | | |
|---|---|---|---|---|---|---|---|---|
| P1 | 15 | 20 | 7 | 11 | Daily | U.S. | Woman | Post Doc Researcher |
| P2 | 6 | 2 | 5 | 10 | Daily | U.S. | Woman | Ph.D. Student |
| P3 | 3 | 15 | 8 | 15 | Daily | U.S. | Man | Security Program Manager |
| P4 | 4 | 1 | 5 | 10 | Daily | U.S. | Man | Senior Scientist |
| P5 | 10+ | 10+ | 5 | 10+ | Weekly | U.S. | Man | Principal Research Manager |
| P6 | 10 | 20 | 14 | 20 | Daily | U.S. | Man | Principal Research Manager |
| P7 | 3 | 10+ | 10 | 15 | Daily | U.S. | Man | Senior Technical Program Manager |
| P8 | 3 | 5 | 7 | 12 | Weekly | U.S. | Woman | Software Engineer II |
| P9 | 0 | 0 | 2.5 | 8 | Weekly | U.S. | Woman | Data Scientist |
| P10 | 2 | 0 | 1.5 | 5 | Daily | India | Woman | Senior Researcher |
| P11 | 50 | 20 | 15 | 20 | Yearly | U.S. | Man | Principal Engineering Manager |
| P12 | 2 | 0 | 2 | 12 | Weekly | India | Man | Senior Researcher |
| P13 | 4 | 2 | 5 | 15 | Weekly | U.S. | Woman | Principal Applied Research Scientist |
| P14 | 6 | 10 | 13 | – | Yearly | Canada | Woman | Staff Researcher |

*Interview.* The first author also conducted 60-minute semi-structured interviews with participants, where participants evaluated the GPT-4-generated assumptions and analysis plans for two research papers. To reduce participant fatigue during the interview, the interview did not include an evaluation of the generated code; instead, the authors performed a manual analysis (see Section 3.4). Interviews were recorded and transcribed. Recordings were deleted after transcription. Additionally, the papers were paired by length so the interview would stay under the allotted time.

Participants were provided with a document containing instructions with all relevant materials for them to complete the study. In the document, participants were instructed to act as a software engineering research consultant applying the assigned research papers' methodology to company data. Their task was to evaluate outputs from an LLM tool that would assist them in the replication. Next, the instructions included grading rubrics[1] to standardize evaluation of the assumptions and generated analysis plans. Finally, the instructions included space for participants to grade the assumptions and analysis plans according to the aforementioned rubrics.

During the interview, participants were presented with the instructions document. Participants then were debriefed on the task instructions. To evaluate the papers, they first read the abstract and methodology of the paper. After being provided the rubrics for the assumptions, the participants silently graded the assumptions by recording their scores in the instructions document. To reduce memory biases, participants could refer to the methodology while grading. To keep the interview within the allotted 60 minutes (20 minutes per paper), participants could skip grading some assumptions if needed. Next, the participant then discussed with the interviewer their impressions on the assumptions and on how to improve the output. The same process was repeated for the analysis plan. After the participant evaluated the first paper, the protocol was repeated for the second paper.

*Data.* We now describe the data and collected and analyzed through the user study. Participants rated assumptions and analysis plans on a 3-point Likert scale for correctness (i.e., not correct, partially correct, correct) and a 5-point scale for all other constructs. Assumptions were evaluated based on ***correctness***, ***relevance*** (i.e., whether the assumption was necessary to consider in the

---

[1]The rubrics are available in the supplementary material.

replication), and ***insightfulness*** (i.e., whether the assumption reflected a deep understanding of software engineering research methodology or data).

Meanwhile, analysis plans were graded both in their entirety as well as at the individual module level. Analysis plans in their entirety were evaluated based on ***correctness***. At the module-level, the plans were evaluated based on ***correctness*** and ***descriptiveness*** (i.e., whether the instructions provided were descriptive enough for another person to replicate the paper).

*Piloting.* Following best practices in human subjects experiments in software engineering [40], we piloted both the survey and interview with two software engineering researchers to clarify wording. Between each pilot, the survey and interview were updated based on the participants' feedback. Pilot participants' data were not included in the final analysis.

## 3.4   Manual Code Review

Three authors also performed a manual review of the 23 GPT-4-generated code snippets (**RQ2**). The three authors convened in a series of code review workshops to review each of the GPT-4-generated code modules for all the papers. For each code module, the authors reviewed the module specification corresponding to the generated code. The authors then read through the generated code. Together, the authors identified any instances of incorrect and correct code logic. Each incorrect and correct logic instance was noted down only upon consensus following discussion. This process generated a dataset of examples of correct and incorrect logic from the code generated by GPT-4. In total, the analysis with the authors took 7 hours (i.e., 21 person-hours).

Additionally, to understand the correctness of the generated code in terms of execution, the first author manually ran each program. For each file, the missing dependencies were installed with `pip`, if they existed. Finally, the code was run with the Python interpreter and the result of the execution (i.e., pass, fail) was recorded.

## 3.5   Analysis

To analyze the data, we used both quantitative and qualitative analysis techniques.

*Quantitative Analysis.* For the quantitative analysis of survey data, we followed the best practices outlined by Kitchenham and Pfleeger [39] by reporting how often participants agreed or strongly agreed with statements, as well as rated outputs as partially correct or fully correct, relevant or very relevant, insightful or very insightful, and descriptive or very descriptive. For the quantitative analysis of the participant-graded data, we calculated descriptive statistics across all the ratings from all seven papers for the assumptions, analysis plans, and module grading constructs.  Finally, we report the percentage of files that were executed as-is without any errors by Python.

*Qualitative Analysis.* While performing qualitative analysis, we followed best practices from Hammer and Berland [29] by interpreting codes as tabulated claims about the data to be investigated in future research. To qualitatively analyze the interview data, the first author performed open-coding. Only the first author performed open-coding since she also conducted the interviews, and therefore had the best understanding of participants' statements within the study team. To open-code the interview data, the first author first read through the interview transcripts. She identified statements that described strengths and limitations of the GPT-4-generated assumptions and analysis plans and inductively generated codes. These statements were assigned with one or more codes, where each code had a name and a description. After open-coding, the first author performed axial coding on the resulting set of codes to extract out broader themes about GPT-4's performance on replicating empirical software engineering papers. The same process was followed for the dataset of examples of incorrect and correct logic from the code generated by GPT-4.
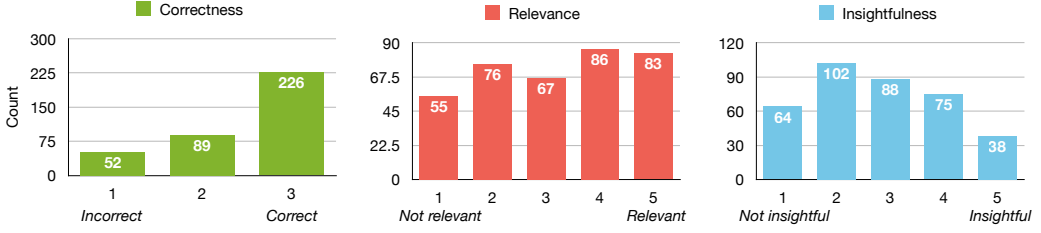
Fig. 5. The distribution of participants' scoring of the GPT-4-generated assumptions by **correctness** (left), **relevance** (middle), and **insightfulness** (right) for all papers.

## 4  RESULTS

In this section, we report our results to the two research questions. Overall, we find that GPT-4 has some understanding of software engineering research methodology and data. However, we also observe that GPT-4 exhibits many limitations to their knowledge, such as lacking basic knowledge of software engineering data. We elaborate further on our findings below.

### 4.1  Can GPT-4 identify assumptions from research methodology? (RQ1)

We report our findings on participants' ratings and impressions of the assumptions generated by GPT-4 on the seven research papers.

*Quantitative Results.* Figure 5 shows the distributions of participants' ratings on the assumptions by **correctness**, **relevance**, and **insightfulness** across all the papers. Overall, we observe that participants rated the assumptions ($N = 367$ scores) as being high in **correctness** ($\mu = 2.5$ out of 3), as 86% of assumptions were graded as partially correct or fully correct. The assumptions were also rated moderately high in terms of **relevance** ($\mu = 3.2$ out of 5), as 46% of assumptions were rated being relevant to the replication. Finally, we observe that participants rated the assumptions lowly in terms of **insightfulness** ($\mu = 2.8$ out of 5), as only 31% of assumptions were rated as insightful.

*Qualitative Results.* We elicited 12 codes related to GPT-4's capabilities and limitations for generating assumptions. Across the codes, three main themes emerged: *reasons for positive ratings* (👍), *reasons for negative ratings* (👎), and *ways of improving outputs* (💡). Below, we report the number of participants who mentioned this code using the multiplication symbol (×).

👍 **Correct** *(13×).* A majority of participants found the assumptions to be correct, as they *"seemed to match the assumptions that were swimming in my head"* (P14). Additionally, *"[assumptions] that weren't [fully correct]...were usually partial correctness"* (P2). As a result, some participants were surprised at GPT-4's capability: *"if an LLM generated these, I'm very impressed"* (P1).

"*Some of the insights were similar to researchers reading the paper could generate.*" (P6)

👍 **Comprehensive** *(12×).* Many participants also found the set of assumptions to be fairly comprehensive in terms of the breadth of topics covered. When asked to generate any assumptions that GPT-4 had missed, 7 participants were not able to suggest any assumptions. Even while taking notes of assumptions while reading the paper, participants still described the set of assumptions to be comprehensive:

"*Based on the notes I took, the assumptions are pretty comprehensive.*" (P4)

↻ *Insightful* (8×). Less frequently, participants also mentioned the assumptions to be insightful, as GPT-4 *"[brought] up some great points about the paper and the assumptions made"* (P1). Others were impressed that the assumptions were generated *"not specifically mentioned in the paper"* (P12).

> "*The coolest one by far...was the Confounding Factors one... That was not explicitly caught out in the paper. I even had to go back and look.*" (P3)

↻ *Relevant* (6×). Additionally, some participants felt the assumptions were relevant to consider for replication. Because of this, participants said the set of assumptions could be helpful to *"someone [who] is not a researcher"* (P6).

> "*And relevance-wise, it is able to extract all the relevant assumptions for this.*" (P5)

↻ *Lack of software engineering knowledge* (9×). The most frequent reason for negative ratings was because GPT-4 lacked knowledge of software engineering data or technologies that participants felt were obvious, as it made assumptions that were *"not representative of the real world"* (P1). This often lowered the **relevance** and **insightfulness** scores. Participants noted that GPT-4 surfaced assumptions about the availability of certain software engineering artifacts, such as *"logs"* (P12) and *"code commits"* (P3). P13 noted, *"What kind of pull request wouldn't show you...[the lines of code changed]?...That's like saying they're assuming a computer is an electronic with binary values."* Participants also noted GPT-4's lack of knowledge in technologies:

> "*The assumption that the GitHub project is compatible with...V8 Spider Monkey. I mean, these are super hugely popular JavaScript engines...so it's almost like some domain knowledge that's not explicit in the paper that the model is missing here.*" (P4)

↻ *Not correct* (7×). Participants also frequently noted that some of the assumptions generated by GPT-4 were *"incorrect"* (P5), such as not being an assumption made in the paper. Some assumptions were not assumptions *"because [the authors] tested their data and they were like, 'This is why we're applying [a methodology], because we observed this'"* (P2). Others were *"facts"* (P8) rather than assumptions. Participants also noted instances where GPT-4 would *"couch its answers... It always gives a second opinion of 'Well, I'm going to tell you something, but I'm also going to tell you it might not be that thing.'"* (P4)

> "*The ones that are off were where it wasn't an assumption. It was like something they actually said they were doing.*" (P13)

↻ *Missing assumptions* (4×). Some participants identified valid assumptions from the paper methodology that GPT-4 did not generate. Identifying additional assumptions usually required additional effort from the participant by taking physical notes on assumptions while reading the paper methodology. Participants mentioned that the missed assumptions were usually minor:

> "*One of them I don't remember seeing is the time window that commits in a 30-minute window are merged and treated as one commit... I don't think that's a super important one, but I did write it down.*" (P7)

↻ *Not relevant* (3×). Participants also noted some of the assumptions were not very relevant for the replication, limiting how *"useful"* (P1) the assumption was. In one case, one participant noted it was because GPT-4 did not reflect a proper understanding of what a mining challenge was: *"I think [considering data completeness] is irrelevant. The data set is...the same for everyone, so for all intents and purposes, for the challenge, it is considered the complete data set."* (P14)

> "*I don't think all of them...are very relevant to the work itself.*" (P1)

⛒ ***Not insightful*** *(2×).* A few participants noted that some of the assumptions generated by GPT-4 were not insightful, as they were *"generic or...simple"* (P6). Another time, one participant noted the assumptions were what was expected, given the domain:

> "*They were just par for the course of any...sentiment analysis that you do. So they were less useful and more obvious.*" (P7)

♡ ***Reducing repetition*** *(9×).* Participants often noted the set of assumptions were repetitive. Some participants identified word-for-word duplicated assumptions generated by GPT-4: *"There's a couple weird duplicates, which I don't know if it's a data copying error or if the model actually repeated itself"* (P4). However, assumptions were also repetitive based on topic that *"that syntactically might be slightly different, but semantically if you look at them, it's...the same."* (P12)

> "*Some of the security-related labeling came up in several different forms. This is some-thing that clued me in...that it probably was [LLM-generated] because that's something that LLMs tend to do.*" (P3)

♡ ***Explaining how to work around assumptions*** *(8×).* Participants noted that GPT-4 often generated text about when an assumption may not hold *"that were counter to what was proposed in the paper"* (P9), after presenting the main assumption. However, participants noted that no suggestions were provided on how to address the assumption if it did not hold:

> "*It would be helpful if there was a suggestion to the replicator on how to overcome or handle that assumption.*" (P8)

♡ ***Providing sources*** *(6×).* Participants also mentioned that they would like a link to the source of the assumption, such as by *"[giving] line numbers"* (P13). This could help explain *"why these assumptions are important or have been made"* (P5) or build confidence in GPT-4's responses:

> "*If there's an AI, it should have a confidence interval telling me, 'I'm 99% confident this assumption is correct...or is explicitly there.'*" (P11)

> ❯ **Key findings:** 86% of assumptions ratings considered the assumptions to be partially correct or fully correct, 46% were considered to be relevant, and 31% to be insightful. Participants also noted that the assumptions were ***correct*** and ***comprehensive***, but had a ***lack of software engineering knowledge***.

## 4.2 Can GPT-4 generate an analysis pipeline to replicate research methodology? (RQ2)

We report our findings on participants' ratings and impressions of the analysis plans generated by GPT-4 on the seven research papers (Section 4.2.1). We also report the results from the manual evaluation performed by three authors on the GPT-4-generated code (Section 4.2.2).

*4.2.1 Analysis Plan.* We report our findings on participants' ratings and impressions of the analysis plans and modules generated by GPT-4 on the seven research papers.

*Quantitative Results.* Figure 6 shows the distributions of participants' ratings on the analysis plans by ***correctness***, as well as individual modules by ***correctness*** and ***descriptiveness*** across all the papers. We observe that participants rated the analysis plans ($N = 27$ scores) as being moderate in terms of ***correctness*** ($\mu = 2.1$ out of 3), as 89% of analysis plans were graded as partially or fully correct. Meanwhile, the individual modules ($N = 87$ scores) performed comparatively in terms of ***correctness*** ($\mu = 2.4$ out of 3), with 72% of modules rated as being partially or fully correct. However, the distribution for the analysis plans is more skewed towards being partially correct compared to
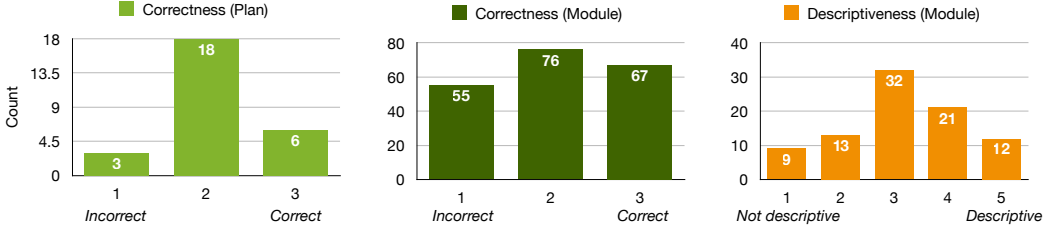
Fig. 6. The distribution of participants' scoring of the GPT-4-generated analysis plans by ***correctness*** (left), as well as individual modules by ***correctness*** (middle) and ***descriptiveness*** (right), for all papers.

the distribution of the modules. Finally, we observe that participants rated the assumptions lowly in terms of ***descriptiveness*** ($\mu$ = 3.2 out of 5), as only 38% of assumptions were rated as descriptive.

*Qualitative Results.* We identified 6 codes related to GPT-4's capabilities and limitations for generating analysis plans. Similar to the assumptions, three themes emerged from these codes: *reasons for positive ratings* (👍), *reasons for negative ratings* (👎), and *ways of improving outputs* (💡).

👍 ***High-level structure*** *(11×).* Participants noted that the analysis plans were correct in their *"high level steps"* (P1), which was the main reason for positive ***correctness*** scores: *"[The modules] broke up how I thought it would."* (P3) Participants mentioned this could be a *"useful [starting point] for somebody else if you want to replicate the analysis"* (P10).

> "*I thought it was able to chunk [the methodology] well into the different pieces, which were outlined in the paper.*" (P13)

👎 ***Not descriptive*** *(13×).* Participants frequently noted that the description and methodology text corresponding to the analysis plan were not descriptive. Some participants felt the *"detail was superfluous...which was a distraction to actually what's being done"* (P9). Further, the descriptions were written largely for software engineering experts: *"I would know what I would do in general...but I'm from the area"* (P11). This is because smaller details were often not *"super clear in the methodology"* (P4) and were not elaborated upon: *"Everything is really close to being almost exactly as I'd want it. But everything is missing, just like a little bit crisper description"* (P4).

> "*[The analysis plan] contains methodology text...but I would have wanted it to be a little bit more explicit.*" (P13)

👎 ***Not correct*** *(9×).* Participants noted some of the analysis plans contained details that were incorrect. One participant noted incorrect module ordering: *"The last step makes no sense in this order. Why would I adjust the time zone at the very end instead of before?"* (P11) Another participant noted that some of the modules created were useless, as they did not consider the replication context: *"The dataset loading piece is going to be the primary difference for replicating the study... I would completely throw away the data set loading thing because I'm going to load a different data set."* (P3) Other participants noted that the inputs of the module specification could also be incorrect:

> "*[The module] says for the refactoring revision identification, the input is two program versions. But that's not necessarily true because we need the entire change history and then we look at a pair of revisions [and] see if there was a refactoring revision.*" (P5)

💡 ***Providing additional context*** *(6×).* Participants noted one way to improve the analysis plan and module outputs was by providing context on related research artifacts, such as a *"Git repository*

*of the replication package"* (P2) and *"paper references"* (P12) or context about the replication context, such as information about the *"new dataset"* (P3).

> "*So [let's say] I'm not even using GitHub. I'm using SVN. Can I tell the model I got this project and guide me to the next step in a...way that adapts to the user's scenario?*" (P6)

⚲ **Improving modularization** *(3×)*. A few participants said that the modularization could also be improved. In particular, participants noted that the responsibilities divided between modules were often uneven and could be divided into additional modules: *"I think the Automatic Identification module is very big. I would have chunked it up into two smaller steps."* (P1)

> "*It's funny because all the work is actually in one of the modules.*" (P7)

⚲ **Providing sources** *(3×)*. Some participants wanted a link to the source of the modules' scope, such as through specific *"reference[s] to [the] methodology section that it references"* (P9).

> "*Maybe giving some references would be better. Yeah, to the methodology text.*" (P10)

*4.2.2 Code.* We report our findings from the manual analysis of the GPT-4-generated code from the module specifications. We found that 7 of the 23 generated code modules (30%) were executable without any modifications with the Python interpreter.

Below, we describe the 11 codes based on correct behaviors and incorrect behaviors within the code generated by the model. Across the codes, we identified two themes: *data* (🍮) and *logic* (⚙). We report the number of Python files that contained this code using the multiplication symbol (×).

*Correct behaviors.* We elicited 3 types of correct behaviors within the code generated by GPT-4.

⚙ **Correct API usage** *(18×)*. As found in prior work [42], the GPT-4-generated code was most often successful in generating boilerplate code for APIs. We noted that GPT-4 generated API calls correctly for many Python APIs, such as `difflib`, `json`, `itertools`, `mysql`, `nltk`, `re`, `scikit-learn`, `sentistrength`, and `subprocess`.

⚙ **High-level structure** *(15×)*. The overall logic of the code modules often followed the correct sequence of steps that was described in the methodology text and description. This aligns with previous work [8, 42], which has noted that code generation models can help developers scaffold solutions to open-ended problems. For example, the second code module in the Pletea et al. [48] analysis plan described steps to instantiate a list of keywords, Porter stem the keywords, and filter GitHub comments based on the Porter stemmed keywords. GPT-4 generated code corresponding to each one of those steps that contained generally correct logic to accomplish the step.

🍮 **Correct data source** *(2×)*. We noticed GPT-4 selecting the correct data source (i.e., a SQL table or loading in from a pre-defined JSON file) infrequently. For instance, the first module in the Pletea et al. [48] analysis plan provided instructions to *"[extract] the relevant tables containing comments on commits and pull requests"*, which GPT-4 successfully did.

*Incorrect behaviors.* We identified 8 codes for the incorrect behaviors in the generated code.

🍮 **Incorrect data source** *(15×)*. GPT-4 also frequently generated code which did not identify the correct data sources (i.e., a SQL table or a pre-defined JSON file) to load from. This aligns with prior work, which has found that selecting from complex database schemas is challenging for language models [66]. In the second module of the Selakovic and Pradel [54] paper, the generated code incorrectly tried to identify tests in the issue body. Additionally, there sometimes were slight errors in how the data was handled or queried—in the first module of the Fregnan et al. [22] analysis plan, only open pull requests are considered, despite it not being specified in the text. We also noticed that more complex SQL queries tended to contain incorrect logic.

⚙ **Missing methodology steps** *(13×)*. GPT-4 often missed generating code for specific steps detailed in the methodology text. This commonly happened for methodology that was lengthy or was very detailed. For instance, in the first module of the Tian et al. [56] analysis plan described logic to filter out non-English commits, which was not implemented by GPT-4. Other times, the methodology was not implemented by GPT-4 but was instead left for a human programmer to complete in a code comment. For instance, the second module in the Selakovic and Pradel [54] analysis plan detailed instructions to *"measure the execution times [of performance fixes] and keep only issues where the optimization leads to a statistically significant performance improvement."* In this case, GPT-4 generated a code comment: `# Check if the issue leads to a statistically significant performance improvement.`

⚙ **Incorrect logic** *(12×)*. We also observed errors in how GPT-4 generated code to implement the methodology. Errors included small logic errors, such as instantiating global variables inside loops or using the wrong type of data structure. Other times, there were errors in the general approach to implement a part of the methodology. For instance, in the second module of the Eyolfson et al. [18] analysis plan, the code identified lines that a commit changed by looking at brand new lines added in a commit rather than examining existing lines that were removed.

⚙ **Guessed implementation** *(7×)*. We also observed errors when there was as lack of detail in the methodology, where GPT-4 "guessed" an implementation for the behavior. For example, the first module for Fregnan et al. [22] analysis plan calls for determining *"whether the file is a test"*, to which GPT-4 implements by checking if the term "test" exists in the file name. While this implementation is technically correct, there could be more sophisticated ways of identifying tests, such as analyzing file source code.

🗄 **Hallucinated data** *(5×)*. We noticed a few errors due to GPT-4 hallucinating different data sources. In some cases, it would hallucinate data files that did not exist, such as a JSON file containing time zones for the fifth module in the analysis plan of Eyolfson et al. [18].

⚙ **Incorrect API usage** *(2×)*. We also noticed errors in GPT-4 using existing APIs incorrectly based on the documentation. For the second module of the Guzman et al. [27] analysis plan, we noted that the usage of the `sentistrength` library was incorrect, as it relied on a parameter that returned the incorrect type of data for the analysis.

⚙ **Hallucinated APIs** *(2×)*. GPT-4 also made errors by hallucinating APIs that did not exist. In the first module of the Kim et al. [35] analysis plan, the specification calls for using Change Distiller [20]. GPT-4 invented an API for Change Distiller, calling hallucinated methods such as `compute_differences()`.

---

❯ **Key findings:** 89% of analysis plan ratings graded the plans to be partially correct or fully correct. Participants also noted that the analysis plans provided a ***high-level structure***, but were ***not descriptive***. Finally, the generated code was most successful in providing ***high-level structure*** and ***API usage*** examples, but often had ***missing methodology steps*** and pulled from an ***incorrect data source***.

---

## 5  DISCUSSION & FUTURE WORK

In this section, we discuss our results with respect to prior work. This generates implications for future work, which we describe below.

## 5.1  GPT-4: The New Software Engineering Research Assistant?

Based on our results, we find that through pre-training alone, LLMs have some expertise in software engineering research to perform study replications, as many of the derived assumptions and analysis plans received high **correctness** scores. Further, the analysis plans and code seemed generally correct in implementing the research methodology in its **high-level structure**. However, based on our findings, GPT-4's domain expertise does not match that of a software engineering research expert. GPT-4 exhibited major gaps in its knowledge of software engineering research and data, which prevents it from performing replications autonomously. For example, assumptions were rated low in terms of **insightfulness** since GPT-4 generated assumptions that indicate a **lack of software engineering knowledge**. Further, generated code often used the **incorrect data source**, provided **guessed implementations**, and **hallucinated APIs** to assist with the analysis of the data. These results align with prior work [23], which found that LLMs reflected some domain knowledge in technical fields like NLP, but made errors such as providing incorrect information.

To enable autonomous replication of empirical software engineering studies, future work could investigate methods for teaching LLMs like GPT-4 software engineering expertise. Theories of developer expertise has noted the importance of task-specific [7] or domain expertise [43] in software development. While GPT-4 has some demonstrated knowledge of programming, it is unable to apply software engineering domain expertise in the code that it generates, such as knowing what database table to query from.

Furthermore, teaching LLMs basic knowledge about software engineering, for example that pull requests usually having access to changed lines, could be another avenue for increasing these models' general knowledge about the field and allow them to reason about the domain better. This could be done through exposing LLMs to more software engineering knowledge by fine-tuning code LLMs like Code LLaMA [51], pre-training specialized language models [26], or providing additional context about a user's replication context (e.g., the data available) or the study's materials (e.g., datasets, code scripts, or external references) while prompting.

However, in teaching domain knowledge to LLMs, it is important that LLMs reflect some level of understanding of the domain rather than performing the task by overfitting to spurious features. For example, we noticed that GPT-4 would create modules based on the section headers of a research paper, rather than based on an understanding of the methodology itself. Therefore, future work could investigate ways to design modules based on a semantic understanding of the text.

## 5.2  GPT-4 Can Only Be as Good as the Research Paper Methodology

In the analysis generated analysis plans and code, we noticed that the code was often incorrect due to the specification being **not descriptive**. This was caused by a lack of detail in the original paper, which left certain methodological details (e.g., the procedure used to identify tests on pull requests) ambiguous. One way GPT-4's generated code could have been improved is if the paper methodology itself were written in a more detailed and systematic way by the authors.

Overall, our exploration in replicating empirical software engineering papers with GPT-4 sheds light on one factor of the replication crisis in empirical studies in computer science [15]—the lack of detail in the methodology section. Participants noted that these sections were not descriptive enough for individuals unfamiliar with software engineering research, such as software engineers:

> "[If] I hand [the methodology] off to a [software engineering] intern...and be like 'Hey write this [in code]!', I would have wanted to be a little bit more explicit." (P13)

GPT-4 could help alleviate this issue by signaling when a methodology section is sufficiently detailed, as the degree of correctness of the generated code could be used as a sign of how well-written the methodology is. If the methodology is too vague, GPT-4 could potentially also be used

to add in fine-grained, step-by-step implementation details to help increase the accuracy of the generated code, which could be another avenue for future research.

### 5.3 LLM-powered Experiences for Empirical Research: Opportunities & Challenges

Our findings have implications for future LLM-powered experiences for researchers. Given that GPT-4 can get the correct **high-level structure** of the analysis plan and code, GPT-4 could be a useful tool for writing analysis pipelines to replicate empirical research. However, given GPT-4's tendency to propagate errors, such as **hallucinating data** and **hallucinating APIs**, they currently are better suited for scaffolding out analyses in code. However, these experiences should provide gaps for users to use their own expertise and write implementations for the parts of the generated code the model has low confidence on. This is because verifying code is a time-consuming activity for users interacting with code generation tools [46].

Another avenue to apply LLMs like GPT-4 in empirical research is during peer review. In our study, GPT-4 was able to provide assumptions that were high in **correctness** scores. Further, participants reported that reading GPT-4-generated assumptions made them more *"critical"* (P9) of the methodology, which could be *"valuable as a reviewer"* (P8). However, these experiences should ensure that the assumptions are not **repetitive** to reduce the cognitive overhead of users.

Finally, our qualitative results underscore the importance of implementing mechanisms to build trust with users while using AI to design analyses, as found in prior work [24]. Participants wanted GPT-4 to **provide sources** for where assumptions and analysis plans had originated from to place more confidence in the outputs. Therefore, applying directly tying each assumption, analysis plan, and line of code to specific sections in the methodology text could be beneficial to users.

## 6 THREATS TO VALIDITY

*External Validity.* Empirical studies have difficulty generalizing [21]. In our study, the assumptions, analysis plans, and code are generated by GPT-4 [47] during August 2023, rather than other LLMs such as PaLM [14], PaLM 2 [3], LLaMA [57], and LLaMA 2 [58]. GPT-4 has been shown to outperform these other models on diverse tasks [58] which suggests that these models (at least "off-the-shelf" with no fine-tuning) would perform no better. Nonetheless, it is unclear the extent to which our results may generalize to other LLMs or future versions of GPT models that may behave differently.

Additionally, we evaluated GPT-4 on seven empirical software engineering papers, which are not representative of all the methodologies or topics performed in the field. These papers are also not representative of how all research papers are written or are structured due to the paper selection criteria. Thus, the study results may not generalize to all studies in empirical software engineering.

*Internal Validity.* How a prompt is crafted can affect the model's performance on a task [6, 11, 50, 61]. While we eliminated the non-determinism of model outputs by setting the temperature to 0, the results from this study could be influenced by how prompts were constructed, as a more effective prompt could elicit different outputs thus alter the ratings and reactions from study participants.

Additionally, user study participants could have misunderstood the wording of the survey and interview questions. To reduce this threat, we piloted the survey and interview with two software engineering researchers to clarify wording of the protocols.

*Construct Validity.* User study participants graded the assumptions and analysis plans along subjective criteria, such as **insightfulness**, **relevance**, **descriptiveness**. Thus, participants' ratings may be inconsistent or may not accurately reflect these constructs. To reduce this threat, we developed grading rubrics and piloted them with two software engineering researchers.

## 7 CONCLUSION

In this study, we investigated GPT-4's ability to replicate empirical software engineering papers. We prompted GPT-4 to generate assumptions, analysis plans (i.e., a sequence of module specifications), and code, given a research paper's methodology for seven empirical software engineering papers.

To evaluate the assumptions and analysis plans, we ran a user study with 14 software engineering researchers and data scientists. We also manually reviewed the code generated by GPT-4 with three coauthors. Overall, we find that GPT-4 is able to generate assumptions that are correct, but lack domain knowledge in software engineering. Further, we find that the code generated by GPT-4 is correct in its high-level structure, but can contain errors in its lower-level implementation. This is also reflected in the analysis plans generated by GPT-4, which are correct in its high-level structure, but lack detail in its description. Our findings have implications on leveraging LLMs for software engineering research, such as teaching LLMs like GPT-4 more domain knowledge in software engineering in both natural language and code. We have made the study's prompts; GPT-4-generated assumptions, analysis plans, and code; and survey and interview protocols available in the supplemental material [5].

## DATA AVAILABILITY

Our supplemental materials are available on HotCRP and a private link on Figshare [5]. The data includes the list of assumptions, analysis plans, and code generated by GPT-4 on the seven papers that were evaluated in this study; participants' ratings of the assumptions and analysis plans; example prompts used to generate the assumptions, analysis plans, and code; the schema of the predefined SQL database; and the survey and interview protocols from the user study. Upon the acceptance of the paper, we will make the supplemental materials publicly available.

## ACKNOWLEDGMENTS

## REFERENCES

[1] 2023. ChatGPT Plugins. Retrieved September 25, 2023 from https://openai.com/blog/chatgpt-plugins#code-interpreter.

[2] 2023. Standards | Empirical Standards. Retrieved September 25, 2023 from https://sigsoft.org/EmpiricalStandards/docs/?standard=Replication.

[3] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. 2023. PaLM 2 Technical Report. *arXiv preprint arXiv:2305.10403* (2023).

[4] Sören Auer, Dante AC Barone, Cassiano Bartz, Eduardo G Cortes, Mohamad Yaser Jaradeh, Oliver Karras, Manolis Koubarakis, Dmitry Mouromtsev, Dmitrii Pliukhin, Daniil Radyush, et al. 2023. The SciQA Scientific Question Answering Benchmark for Scholarly Knowledge. *Scientific Reports* 13, 1 (2023), 7240. https://doi.org/10.1038/s41598-023-33607-z

[5] Anonymous Author(s). 2024. Supplemental Materials to "Can Large Language Models Replicate Empirical Software Engineering Research". The materials are uploaded as Supplementary Material to the HotCrp submission system and Figshare at the following url: https://figshare.com/s/603602d213e22cdcc491.

[6] Stephen Bach, Victor Sanh, Zheng Xin Yong, Albert Webson, Colin Raffel, Nihal V. Nayak, Abheesht Sharma, Taewoon Kim, M Saiful Bari, Thibault Fevry, Zaid Alyafeai, Manan Dey, Andrea Santilli, Zhiqing Sun, Srulik Ben-david, Canwen Xu, Gunjan Chhablani, Han Wang, Jason Fries, Maged Al-shaibani, Shanya Sharma, Urmish Thakker, Khalid Almubarak, Xiangru Tang, Dragomir Radev, Mike Tian-jian Jiang, and Alexander Rush. 2022. PromptSource: An

Integrated Development Environment and Repository for Natural Language Prompts. In *Annual Meeting of the Association for Computational Linguistics (ACL): System Demonstrations*. Association for Computational Linguistics, 93–104. https://doi.org/10.18653/v1/2022.acl-demo.9

[7]  Sebastian Baltes and Stephan Diehl. 2018. Towards a Theory of Software Development Expertise. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 187–200. https://doi.org/10.1145/3236024.3236061

[8]  Shraddha Barke, Michael B James, and Nadia Polikarpova. 2023. Grounded Copilot: How Programmers Interact with Code-generating Models. *Proceedings of the ACM on Programming Languages* 7, OOPSLA1 (2023), 85–111. https://doi.org/10.1145/3586030

[9]  Andrew Begel and Thomas Zimmermann. 2014. Analyze This! 145 Questions for Data Scientists in Software Engineering. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 12–23. https://doi.org/10.1145/2568225.2568233

[10] Christian Bird, Nachiappan Nagappan, Premkumar Devanbu, Harald Gall, and Brendan Murphy. 2009. Does Distributed Development Affect Software Quality? An Empirical Case Study of Windows Vista. *Commun. ACM* 52, 8 (2009), 85–93. https://doi.org/10.1109/ICSE.2009.5070550

[11] 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 (NeurIPS)* 33 (2020), 1877–1901.

[12] J. Carver, J. VanVoorhis, and V. Basili. 2004. Understanding the Impact of Assumptions on Experimental Validity. In *International Symposium on Empirical Software Engineering (ISESE)*. 251–260. https://doi.org/10.1109/ISESE.2004.1334912

[13] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating Large Language Models Trained on Code. *arXiv preprint arXiv:2107.03374* (2021).

[14] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. PaLM: Scaling Language Modeling with Pathways. *arXiv preprint arXiv:2204.02311* (2022).

[15] Andy Cockburn, Pierre Dragicevic, Lonni Besançon, and Carl Gutwin. 2020. Threats of a Replication Crisis in Empirical Computer Science. *Commun. ACM* 63, 8 (2020), 70–79. https://doi.org/10.1145/3360311

[16] Robert A DeLine. 2021. Glinda: Supporting Data Science with Live Programming, GUIs and a Domain-specific Language. In *ACM CHI Conference on Human Factors in Computing Systems*. 1–11. https://doi.org/10.1145/3411764.3445267

[17] Dino Distefano, Manuel Fähndrich, Francesco Logozzo, and Peter W O'Hearn. 2019. Scaling Static Analyses at Facebook. *Commun. ACM* 62, 8 (2019), 62–70. https://doi.org/10.1145/3338112

[18] Jon Eyolfson, Lin Tan, and Patrick Lam. 2011. Do Time of Day and Developer Experience Affect Commit Bugginess?. In *Working Conference on Mining Software Repositories (MSR)*. 153–162. https://doi.org/10.1145/1985441.1985464

[19] Sakina Fatima, Taher A Ghaleb, and Lionel Briand. 2022. Flakify: A Black-box, Language Model-based Predictor for Flaky Tests. *IEEE Transactions on Software Engineering* (2022). https://doi.org/10.1109/TSE.2022.3201209

[20] Beat Fluri, Michael Wursch, Martin PInzger, and Harald Gall. 2007. Change Distilling: Tree Differencing for Fine-grained Source Code Change Extraction. *IEEE Transactions on Software Engineering* 33, 11 (2007), 725–743. https://doi.org/10.1109/TSE.2007.70731

[21] Bent Flyvbjerg. 2006. Five misunderstandings about case-study research. *Qualitative Inquiry* 12, 2 (2006), 219–245. https://doi.org/10.1177/1077800405284363

[22] Enrico Fregnan, Larissa Braz, Marco D'Ambros, Gül Çalıklı, and Alberto Bacchelli. 2022. First Come First Served: The Impact of File Position on Code Review. In *ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 483–494. https://doi.org/10.1145/3540250.3549177

[23] Fan Gao, Hang Jiang, Moritz Blum, Jinghui Lu, Yuang Jiang, and Irene Li. 2023. Large Language Models on Wikipedia-Style Survey Generation: an Evaluation in NLP Concepts. *arXiv preprint arXiv:2308.10410* (2023).

[24] Ken Gu, Madeleine Grunde-McLaughlin, Andrew M McNutt, Jeffrey Heer, and Tim Althoff. 2023. How Do Data Analysts Respond to AI Assistance? A Wizard-of-Oz Study. *arXiv preprint arXiv:2309.10108* (2023).

[25] Philip J Guo, Thomas Zimmermann, Nachiappan Nagappan, and Brendan Murphy. 2010. Characterizing and Predicting which Bugs Get Fixed: An Empirical Study of Microsoft Windows. In *ACM/IEEE International Conference on Software Engineering (ICSE)*. 495–504. https://doi.org/10.1145/1806799.1806871

[26] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *Annual Meeting of the Association for Computational Linguistics (ACL)*. 8342–8360. https://doi.org/10.18653/v1/2020.acl-main.740

[27] Emitza Guzman, David Azócar, and Yang Li. 2014. Sentiment Analysis of Commit Comments in GitHub: An Empirical Study. In *Working Conference on Mining Software Repositories (MSR)*. 352–355. https://doi.org/10.1145/2597073.2597118

[28] Perttu Hämäläinen, Mikke Tavast, and Anton Kunnari. 2023. Evaluating Large Language Models in Generating Synthetic HCI Research Data: A Case Study. In *ACM CHI Conference on Human Factors in Computing Systems (CHI)*. 1–19. https://doi.org/10.1145/3544548.3580688

[29] David Hammer and Leema K Berland. 2014. Confusing Claims for Data: A Critique of Common Practices for Presenting Qualitative Research on Learning. *Journal of the Learning Sciences* 23, 1 (2014), 37–46. https://doi.org/10.1080/10508406.2013.802652

[30] Tobias Hey, Jan Keim, Anne Koziolek, and Walter F Tichy. 2020. Norbert: Transfer Learning for Requirements Classification. In *IEEE International Requirements Engineering Conference (RE)*. IEEE, 169–179. https://doi.org/10.1109/RE48521.2020.00028

[31] Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong Wang, Li Li, Xiapu Luo, David Lo, John Grundy, and Haoyu Wang. 2023. Large Language Models for Software Engineering: A Systematic Literature Review. *arXiv preprint arXiv:2308.10620* (2023).

[32] Hennie Huijgens, Ayushi Rastogi, Ernst Mulders, Georgios Gousios, and Arie van Deursen. 2020. Questions for Data Scientists in Software Engineering: A Replication. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 568–579. https://doi.org/10.1145/3368089.3409717

[33] Eunice Jun, Maureen Daum, Jared Roesch, Sarah Chasins, Emery Berger, Rene Just, and Katharina Reinecke. 2019. Tea: A High-level Language and Runtime System for Automating Statistical Analysis. In *ACM Symposium on User Interface Software and Technology (UIST)*. 591–603. https://doi.org/10.1145/3332165.3347940

[34] Eunice Jun, Audrey Seo, Jeffrey Heer, and René Just. 2022. Tisane: Authoring Statistical Models via Formal Reasoning from Conceptual and Data Relationships. In *ACM CHI Conference on Human Factors in Computing Systems (CHI)*. 1–16. https://doi.org/10.1145/3491102.3501888

[35] Miryung Kim, Dongxiang Cai, and Sunghun Kim. 2011. An Empirical Investigation into the Role of API-level Refactorings during Software Evolution. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 151–160. https://doi.org/10.1145/1985793.1985815

[36] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2016. The Emerging Role of Data Scientists on Software Development Teams. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 96–107. https://doi.org/10.1109/TSE.2017.2754374

[37] Miryung Kim, Thomas Zimmermann, Robert DeLine, and Andrew Begel. 2017. Data Scientists in Software Teams: State of the Art and Challenges. *IEEE Transactions on Software Engineering* 44, 11 (2017), 1024–1038. https://doi.org/10.1109/TSE.2017.2754374

[38] Serkan Kirbas, Etienne Windels, Olayori McBello, Kevin Kells, Matthew Pagano, Rafal Szalanski, Vesna Nowack, Emily Rowan Winter, Steve Counsell, David Bowes, et al. 2021. On the Introduction of Automatic Program Repair in Bloomberg. *IEEE Software* 38, 4 (2021), 43–51. https://doi.org/10.1109/MS.2021.3071086

[39] Barbara A. Kitchenham and Shari Lawrence Pfleeger. 2008. Personal Opinion Surveys. In *Guide to Advanced Empirical Software Engineering*, Forrest Shull, Janice Singer, and Dag I. K. Sjøberg (Eds.). Springer, 63–92. https://doi.org/10.1007/978-1-84800-044-5_3

[40] Amy J Ko, Thomas D LaToza, and Margaret M Burnett. 2015. A Practical Guide to Controlled Experiments of Software Engineering Tools with Human Participants. *Empirical Software Engineering* 20, 1 (2015), 110–141. https://doi.org/10.1007/s10664-013-9279-3

[41] Jenny T Liang, Maryam Arab, Minhyuk Ko, Amy J Ko, and Thomas D LaToza. 2023. A Qualitative Study on the Implementation Design Decisions of Developers. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 435–447. https://doi.org/10.1109/ICSE48619.2023.00047

[42] Jenny T Liang, Chenyang Yang, and Brad A Myers. 2023. Understanding the Usability of AI Programming Assistants. *arXiv preprint arXiv:2303.17125* (2023).

[43] Jenny T Liang, Thomas Zimmermann, and Denae Ford. 2022. Understanding Skills for OSS Communities on GitHub. In *ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 170–182. https://doi.org/10.1145/3540250.3549082

[44] David Lo, Nachiappan Nagappan, and Thomas Zimmermann. 2015. How Practitioners Perceive the Relevance of Software Engineering Research. In *ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE)*. 415–425. https://doi.org/10.1145/2786805.2786809

[45] Shantanu Mandal, Adhrik Chethan, Vahid Janfaza, SM Mahmud, Todd A Anderson, Javier Turek, Jesmin Jahan Tithi, and Abdullah Muzahid. 2023. Large Language Models Based Automatic Synthesis of Software Specifications. *arXiv preprint arXiv:2304.09181* (2023).

[46] Hussein Mozannar, Gagan Bansal, Adam Fourney, and Eric Horvitz. 2022. Reading between the Lines: Modeling User Behavior and Costs in AI-assisted Programming. *arXiv preprint arXiv:2210.14306* (2022).

[47] OpenAI. 2023. GPT-4 Technical Report. *ArXiv* abs/2303.08774 (2023).

[48] Daniel Pletea, Bogdan Vasilescu, and Alexander Serebrenik. 2014. Security and Emotion: Sentiment Analysis of Security Discussions on GitHub. In *Working Conference on Mining Software Repositories (MSR)*. 348–351. https://doi.org/10.1145/2597073.2597117

[49] Lutz Prechelt. 2021. On Implicit Assumptions Underlying Software Engineering Research. In *International Conference on Evaluation and Assessment in Software Engineering (EASE)*. 336–339. https://doi.org/10.1145/3463274.3463356

[50] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shotTtext-to-image Generation. In *International Conference on Machine Learning (ICML)*. PMLR, 8821–8831.

[51] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code LLaMA: Open Foundation Models for Code. *arXiv preprint arXiv:2308.12950* (2023).

[52] Caitlin Sadowski, Edward Aftandilian, Alex Eagle, Liam Miller-Cushon, and Ciera Jaspan. 2018. Lessons from Building Static Analysis Tools at Google. *Commun. ACM* 61, 4 (2018), 58–66. https://doi.org/10.1145/3188720

[53] Morgan Klaus Scheuerman, Katta Spiel, Oliver L Haimson, Foad Hamidi, and Stacy M Branham. 2020. HCI Guidelines for Gender Equity and Inclusivity. In *UMBC Faculty Collection*.

[54] Marija Selakovic and Michael Pradel. 2016. Performance Issues and Optimizations in JavaScript: An Empirical Study. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 61–72. https://doi.org/10.1145/2884781.2884829

[55] Mikke Tavast, Anton Kunnari, and Perttu Hämäläinen. 2022. Language Models Can Generate Human-like Self-reports of Emotion. In *ACM International Conference on Intelligent User Interfaces (IUI)*. 69–72. https://doi.org/10.1145/3490100.3516464

[56] Yingchen Tian, Yuxia Zhang, Klaas-Jan Stol, Lin Jiang, and Hui Liu. 2022. What Makes a Good Commit Message?. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 2389–2401. https://doi.org/10.1145/3510003.3510205

[57] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. LLaMA: Open and Efficient Foundation Language Models. *arXiv preprint arXiv:2302.13971* (2023).

[58] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. LLaMA 2: Open Foundation and Fine-tuned Chat Models. *arXiv preprint arXiv:2307.09288* (2023).

[59] David Wadden, Kyle Lo, Lucy Lu Wang, Arman Cohan, Iz Beltagy, and Hannaneh Hajishirzi. 2022. MultiVerS: Improving Scientific Claim Verification with Weak Supervision and Full-document Context. In *Findings of the Association for Computational Linguistics: NAACL*. 61–76. https://doi.org/10.18653/v1/2022.findings-naacl.6

[60] Bingcheng Wang, Pei-Luen Patrick Rau, and Tianyi Yuan. 2022. Measuring User Competence in Using Artificial Intelligence: Validity and Reliability of Artificial Intelligence Literacy Scale. *Behaviour & Information Technology* (2022), 1–14. https://doi.org/10.1080/0144929X.2022.2072768

[61] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought Prompting Elicits Reasoning in Large Language Models. *Advances in Neural Information Processing Systems* 35 (2022), 24824–24837.

[62] Moshi Wei, Nima Shiri Harzevili, Yuchao Huang, Junjie Wang, and Song Wang. 2022. Clear: Contrastive Learning for API Recommendation. In *IEEE/ACM International Conference on Software Engineering (ICSE)*. 376–387. https://doi.org/10.1145/3510003.3510159

[63] Tongshuang Wu, Haiyi Zhu, Maya Albayrak, Alexis Axon, Amanda Bertsch, Wenxing Deng, Ziqi Ding, Bill Guo, Sireesh Gururaja, Tzu-Sheng Kuo, et al. 2023. LLMs as Workers in Human-Computational Algorithms? Replicating Crowdsourcing Pipelines with LLMs. *arXiv preprint arXiv:2307.10168* (2023).

[64] Ziang Xiao, Xingdi Yuan, Q Vera Liao, Rania Abdelghani, and Pierre-Yves Oudeyer. 2023. Supporting Qualitative Analysis with Large Language Models: Combining Codebook with GPT-3 for Deductive Coding. In *ACM International Conference on Intelligent User Interfaces (IUI)*. 75–78. https://doi.org/10.1145/3581754.3584136

[65] Frank F Xu, Uri Alon, Graham Neubig, and Vincent Josua Hellendoorn. 2022. A Systematic Evaluation of Large Language Models of Code. In *ACM SIGPLAN International Symposium on Machine Programming (MAPS)*. 1–10. https://doi.org/10.1145/3520312.3534862

[66] Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2018. Spider: A Large-Scale Human-Labeled Dataset for Complex and Cross-Domain Semantic Parsing and Text-to-SQL Task. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 3911–3921. https://doi.org/10.18653/v1/D18-1425

[67] Jialu Zhang, Todd Mytkowicz, Mike Kaufman, Ruzica Piskac, and Shuvendu K Lahiri. 2022. Using Pre-trained Language Models to Resolve Textual and Semantic Merge Conflicts (Experience Paper). In *ACM International Symposium on Software Testing and Analysis (ISSTA)*. 77–88. https://doi.org/10.1145/3533767.3534396

[68] Zibin Zheng, Kaiwen Ning, Jiachi Chen, Yanlin Wang, Wenqing Chen, Lianghong Guo, and Weicheng Wang. 2023. Towards an Understanding of Large Language Models in Software Engineering Tasks. *arXiv preprint arXiv:2308.11396*

(2023).