# Understanding Users' Dissatisfaction with ChatGPT Responses: Types, Resolving Tactics, and the Effect of Knowledge Level

Yoonsu Kim
yoonsu16@kaist.ac.kr
Graduate School of AI, KAIST
Daejeon, Republic of Korea

Jueon Lee
audreylee@snu.ac.kr
College of Liberal Studies, SNU
Seoul, Republic of Korea

Seoyoung Kim
youthskim@kaist.ac.kr
School of Computing, KAIST
Daejeon, Republic of Korea

Jaehyuk Park
jp@kdischool.ac.kr
School of Public Policy and
Management, KDI
Sejong, Republic of Korea

Juho Kim
juhokim@kaist.ac.kr
School of Computing, KAIST
Daejeon, Republic of Korea

## ABSTRACT

Large language models (LLMs) with chat-based capabilities, such as ChatGPT, are widely used in various workflows. However, due to a limited understanding of these large-scale models, users struggle to use this technology and experience different kinds of dissatisfaction. Researchers have introduced several methods, such as prompt engineering, to improve model responses. However, they focus on enhancing the model's performance in specific tasks, and little has been investigated on how to deal with the user dissatisfaction resulting from the model's responses. Therefore, with ChatGPT as the case study, we examine users' dissatisfaction along with their strategies to address the dissatisfaction. After organizing users' dissatisfaction with LLM into seven categories based on a literature review, we collected 511 instances of dissatisfactory ChatGPT responses from 107 users and their detailed recollections of dissatisfactory experiences, which we released as a publicly accessible dataset. Our analysis reveals that users most frequently experience dissatisfaction when ChatGPT fails to grasp their intentions, while they rate the severity of dissatisfaction related to accuracy the highest. We also identified four tactics users employ to address their dissatisfaction and their effectiveness. We found that users often do not use any tactics to address their dissatisfaction, and even when using tactics, 72% of dissatisfaction remained unresolved. Moreover, we found that users with low knowledge of LLMs tend to face more dissatisfaction on accuracy while they often put minimal effort in addressing dissatisfaction. Based on these findings, we propose design implications for minimizing user dissatisfaction and enhancing the usability of chat-based LLM.

## CCS CONCEPTS

• **Human-centered computing** → **Empirical studies in HCI**.

## KEYWORDS

Large Language Models, Chat-based LLM, ChatGPT, User-side dissatisfaction, Resolving tactics, Knowledge-level, datasets

## 1 INTRODUCTION

Large Language Models (LLM) have exhibited remarkable performance across various tasks (e.g., language generation [62] and reasoning [36]), and they have become more accessible with integration into chat interfaces and instruction tuning [63], such as ChatGPT [1]. As a result, many people are increasingly incorporating this technology into their workflows across various domains such as education [43, 74], healthcare [47, 56, 74], and law [13, 60].

When using a chat-based LLM, natural language prompts play a crucial role because they are the primary medium for interaction between the user and the model [20, 94, 100]. Accordingly, prompt engineering—aimed at enhancing the quality of model responses to get desired responses from the model—has been a popular stream of research. As various people use LLMs in their workflows, researchers and practitioners have published various guidelines, tools, books, and even online courses for prompt engineering, not only for developers but also for laypeople [3, 72, 88, 100].

However, despite the proliferation of these resources, end-users often encounter dissatisfaction during conversations with LLMs. When end-users have limited knowledge about LLMs, they may have incorrect expectations about the model's behavior, which can further contribute to their dissatisfaction. This dissatisfaction may arise from various known limitations of LLMs, including hallucination [9, 39, 51], inconsistency [24, 38, 51], unfavorable tone and format [8, 71, 90], and lack of transparency [17, 79]. In addition, such dissatisfaction can become more critical when end-users utilize LLMs for practical purposes.
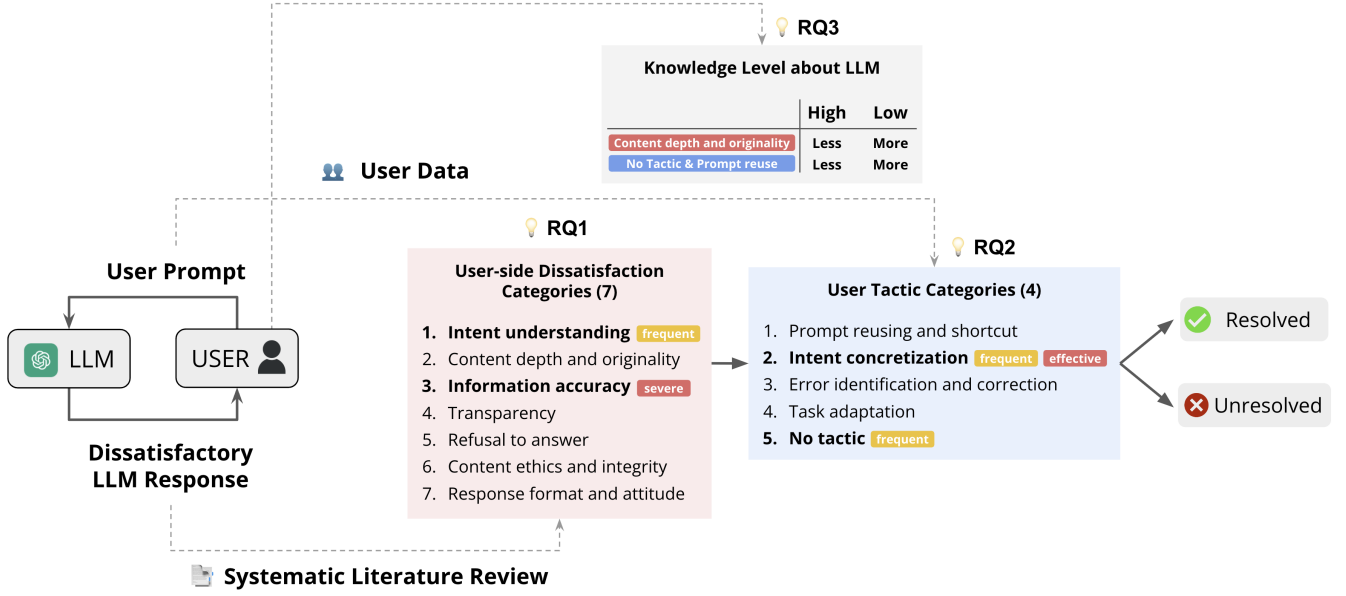
---

[1] https://chat.openai.com/

**Figure 1: Overview of our research questions and findings.**

Little previous research, however, has investigated users' dissatisfaction during conversations with LLMs. In particular, existing prompt engineering techniques mainly focus on enhancing the model's performance in specific tasks, and little has been investigated on how users should respond to dissatisfactions they face from LLMs' responses during the conversation. Therefore, in our research, with ChatGPT as the case study, we aim to understand the dissatisfaction experienced by the users during the conversations. We focus on situations where users seek practical assistance from ChatGPT within their workflows (e.g., translation, email writing, and programming) rather than situations where users intentionally provoke dissatisfactory responses from ChatGPT and test its boundaries and limitations. Specifically, we explore the types of dissatisfaction users experience during the conversation, how serious each type of dissatisfaction is, and how users address dissatisfaction in the subsequent prompts. Furthermore, building upon prior research that demonstrated how users' experiences with technological failure depend on their knowledge of that technology in the context of conversational agent [53], we investigate how dissatisfaction and user responses vary based on the user's knowledge level of LLMs.

At first, we conducted a systematic literature review of papers dealing with limitations and challenges associated with LLMs and identified seven user-side dissatisfaction categories stemming from LLM responses (Table 1). Then, using ChatGPT as a case study, we collected how much users confront these seven dissatisfaction categories and how they respond to them during actual conversations through our data collection system (Figure 2). As a result, we collected 307 ChatGPT conversation logs from 107 respondents, which contained 511 user-side dissatisfactions on ChatGPT responses. Through a quantitative analysis, we found that users most frequently experienced dissatisfaction in terms of ChatGPT's poor

understanding of users' intent, while users felt the most severe dissatisfaction related to inaccuracies in information. We also conducted a qualitative analysis of users' behavior to address the dissatisfaction at subsequent prompts, which resulted in four tactic categories (Table 3): 'prompt reusing', 'intent concretization', 'error identification and correction', 'task adaptation', and 'no tactic'. Moreover, we analyzed differences in dissatisfaction and tactics across the users' knowledge levels on LLMs and confirmed that low-knowledge users more frequently experienced dissatisfaction regarding ChatGPT's responses being too general and lacking originality. We also observed that low-knowledge users often resorted to 'no tactic' or 'prompt reusing', which involved minimal efforts in prompt crafting when they experienced dissatisfaction.

Based on our findings, we suggest design implications to improve the usability of LLMs for users, leveraging the occurrence of dissatisfaction and corresponding tactics during the conversation. We also suggest that the responses of LLMs could be more tailored to the user's knowledge level. Furthermore, we release the actual user data we collected as a publicly available dataset [2] to aid relevant research. The contributions of our research are as follows:

- Categorization and analysis of user-side dissatisfaction and corresponding tactics at the conversational turn level.
- Investigation of how dissatisfaction and tactics appear differently depending on users' knowledge level of LLMs.
- A dataset containing specific user interactions and their experiences of dissatisfaction in actual conversations with ChatGPT, thereby offering resources for further research on user-centric LLMs.

---

[2]https://chatgpt-analysis.kixlab.org

## 2 RELATED WORK

We review related work in (1) limitations and user challenges in LLMs and (2) user's strategies to overcome those challenges in Language Models.

### 2.1 Limitations and User Challenges in LLMs

A rich body of previous work has addressed various limitations associated with language models, including hallucination [9, 39, 51], inconsistency in reasoning [24, 38, 51], and numerical computation [67, 93]. Zhao et al. [97] reviewed major challenges in recent large language models in terms of three basic types of ability of LLMs: language generation, knowledge utilization, and complex reasoning. Borji [14] organized ChatGPT's failures into eleven distinct categories, including reasoning, factual errors, math, coding, and bias.

However, how users actually experience may be different from LLM's failures. Thus, several studies investigated challenges that can be experienced from the user's side [11, 14, 74]. Behrooz [11] points out the core challenges of research chatbots like OpenAI's ChatGPT, Meta AI's BlenderBot, and Google's LaMDA, especially related to user perceptions. These challenges encompass the lack of conversational context [16, 82], the speaker perception void [33], and the lack of expectation baseline [77].

While a stream of research has explored the limitations of language models and the user challenges when interacting with them, there is a lack of comprehensive categorization of the user-side dissatisfaction and how often and seriously users experience each dissatisfaction in the context of users' actual conversation situations. Understanding the user-side dissatisfactions arising from practical usage can provide insights into building LLMs with better usability. To this end, our paper investigates how users experience dissatisfaction and the severity of these dissatisfactions by analyzing users' conversation logs with LLMs.

### 2.2 User's Strategies to Overcome Challenges in Language Models

To improve the usability of language models, it is important to understand users' current practices to overcome the challenges they face. For this, previous research has delved into how users react and overcome challenges encountered while interacting with various language models. Porcheron et al. [65] and Luger et al. [53] examined how users interact with a conversational agent in voice user interfaces (VUI). Specifically, Myers et al. [58] identified ten main categories of tactics users employ to overcome challenges encountered in VUI, and discovered patterns of tactics. Although LLMs and VUIs share the same characteristic in that users communicate with AI agents via natural language, how users overcome challenges may differ as LLMs use text prompting, which may allow more careful prompting strategies compared to VUIs.

Accordingly, prompt engineering techniques have been extensively studied to address challenges in LLMs [54, 72, 85, 88, 100]. For instance, Chain-of-Thought Prompting (CoT) is renowned for improving LLM's reasoning performance by integrating intermediate reasoning steps into prompts [85]. Building upon the effectiveness of CoT, researchers have explored variants like Zero-shot CoT [46], Auto-CoT [96], and Self-Consistency (CoT-SC) [84] and showed that those methods can mitigate LLM's deficiency in reasoning. Specifically, CoT-SC is also known for mitigating LLM's inconsistency issue. Madaan et al. [55] also showed that transforming a certain task into a code generation task can be effective in addressing reasoning and inconsistency issues in LLMs.

However, previous studies investigate how to enhance the model's performance in specific tasks (e.g., reasoning), and they rarely addressing how to handle the dissatisfaction experienced by users during conversations with LLM, stemming from the responses they receive. Therefore, in this paper, we investigate users' behaviors when they encounter dissatisfaction from their actual conversations with LLM. Through this, we analyze users' tactics to address their dissatisfaction and their effectiveness. This will provide insights into how LLM and its interface can be further developed to aid users when they encounter dissatisfaction in the middle of the conversation.

## 3 SYSTEMATIC LITERATURE REVIEW: CATEGORIZING USER-SIDE DISSATISFACTION

To understand and categorize the dissatisfaction points that users encounter when using LLMs for practical purposes, we conducted a systematic literature review to investigate the challenges, limitations, and failures identified in previous research within the LLM context. We focused on user-side dissatisfaction experiences directly arising from LLM responses. For this purpose, we scrutinized a total of 59 papers and conducted qualitative coding, which resulted in 19 codes representing user-side dissatisfaction points from LLM responses. These points were subsequently categorized into seven themes (Table 1). The seven themes were provided as multiple-choice items in our data collection, allowing users to select the dissatisfaction points they have experienced from LLM responses.

### 3.1 Search Keywords

We first conducted an extensive search on Google Scholar [3], ACM Digital Library [4], and arXiv [5] using the combination of "Large Language Models(LLMs)," and "ChatGPT," with "Challenges," "Limitations," and "Difficulties" as search keywords. The reason we specifically included ChatGPT as a search keyword is because ChatGPT has been one of the most extensively used LLMs and has been widely adopted across a variety of domains, such as the medical domain and education. Considering the temporal progress in LLM technologies, we restricted the search period to after 2021. To not exclude papers that might be relevant but do not explicitly contain our search keywords, we extended our search by traversing the citation graph of the initial set of papers. We explored the papers that are either cited by or cite the papers within our initial set and gathered any papers that discuss user-side dissatisfaction, challenges, or difficulties with the use of LLMs, as well as instances of LLM failures.

---

[3]https://scholar.google.com
[4]https://dl.acm.org
[5]https://arxiv.org

## 3.2 Exclusion Criteria and Filtering

As the result of the search process (Section 3.1), we collected 1,249 papers. After removing duplicates, we had 866 papers. To focus on user-side dissatisfaction with LLM responses, we set the 4 exclusion criteria and filtered papers based on them.

**EC1.** We excluded papers that used the terms "limitation," "challenge," or "difficulty" in a general sense, not specifically about LLMs.

**EC2.** We excluded papers that focused solely on the technical challenges or limitations of LLMs.

**EC3.** We excluded papers that discussed potential risks of LLM usage, such as the overreliance of students on LLMs for learning [41, 69] or the potential for privacy issues [31].

**EC4.** We excluded papers that discuss the difficulty of tuning or maintaining LLMs that are not directly related to LLMs responses.

We filtered papers following these criteria, resulting in 59 papers. This allows us to include papers that discuss the practical application of LLMs in specific domains or workflows intended to enhance productivity, which resulted in diverse fields such as education, healthcare, and research.

## 3.3 Analysis Procedure

To analyze and categorize the user-side dissatisfaction from LLM responses, our initial step involved reading 59 papers and compiling a comprehensive list related to user-side dissatisfaction, challenges, or difficulties with the use of LLMs, as well as instances of LLM failures. Two authors then independently conducted open coding on the compiled list. Our primary focus was on identifying aspects of user-side dissatisfaction that emanated from interactions with LLM responses. Following the individual open coding phase, the two authors engaged in collaborative and iterative discussions. These discussions were instrumental in consolidating and refining the initially identified codes. The authors worked together to ensure that the codes accurately captured the nuances of user dissatisfaction associated with LLM responses. Subsequently, to establish relationships among these codes, all authors participated in axial coding [78]. This involved a series of successive discussions aimed at clustering the individual codes into broader, more abstract categories. The goal was to identify common threads and overarching themes that emerged from the data. The axial coding process culminated in the consolidation of the identified aspects of user-side dissatisfaction into seven main themes. (Table 1) These themes encapsulated the various dimensions of user dissatisfaction when interacting with LLM responses. The dissatisfaction themes were later used when collecting data from users, which is explained in detail in Section 4.

## 3.4 Result: Categorizing User-side Dissatisfaction

We categorized the various aspects of user dissatisfaction arising from LLM responses into 19 distinct codes, further organized into seven overarching themes. The detailed information is denoted in Table 1. All paper lists are in the Appendix A.1.

**Theme 1. Intent Understanding ($D_{intent}$)** This theme encompasses issues related to LLM's failure to correctly interpret or reflect the user's intent, instructions, or context. Three codes (C1, C2, C17) fall into this theme. LLM outputs often fail to align with the users' needs and expectations [42]. ChatGPT has been found to suggest unnecessary out-of-context actions in medical use [70], and to use the wrong tone or be excessively literal due to its low understanding of non-literal language such as sarcasm [71].

**Theme 2. Content Depth and Originality ($D_{depth}$)** Users experienced this type of dissatisfaction when they expected more in-depth and creative answers catered to their specific needs, but LLM gave responses that were perceived as overly general, lacking originality, or requiring more diversity. ChatGPT rarely diverges from the topic, generating less diverse content than humans [14, 30]. Concerns rise on unvarying and repetitive ChatGPT outputs which are results of generation based on past data [45]. ChatGPT showed weaknesses in providing practical examples in academic writing [47].

**Theme 3. Information Accuracy ($D_{acc}$)** Dissatisfactions related to false, outdated, or inaccurate information in responses fall under this theme. In addition, inconsistencies within one response or in conversation beyond one answer also belong to this theme. Users were dissatisfied when LLMs provided incorrect or conflicting information, eroding trust in the system's reliability. ChatGPT is incompetent in correctly calculating large numbers [8], and bases its answers on training data up to a certain point in the past - September 2021 is the cutoff in the latest released version of ChatGPT-therefore generating outdated and wrong information when facts change over time [92]. Language models are known to show inconsistency in their claims and explanations [6]. ChatGPT has limited reasoning capabilities, including inductive, spatial, and mathematical reasoning [9, 95]. Hallucination, the generation of absurd output that contradicts the source or cannot be verified from it, is a threat in real-world applications since the wrong output can cause harm when people trust the outcome of LLMs without further inspection [39]. Sycophancy, a behavior where LLMs contradict their original output in order to agree with human input, is also a reason for concern about inaccurate and trustworthy generation [64].

**Theme 4. Transparency ($D_{trans}$)** Users experiencing difficulties in understanding the underlying reasoning or criteria behind LLM responses led to dissatisfaction related to transparency. Users desired more transparency in how the language model generated its answers, especially when complex or critical information was involved. The 'black box' nature of LLMs makes it difficult for users to interpret the reasons behind their outputs.

**Theme 5. Refusal to Answer ($D_{refuse}$)** Responses where LLMs avoided providing answers, often using phrases like "As a language model, I am not capable..." or similar, were categorized under this theme. Users were frustrated when the system declined to provide information or guidance. ChatGPT may refrain from giving its direct opinion [14], and refuse to verify if a claim can be considered misinformation when the claim is closely related to social issues [9]. Refusing is also found in questions regarding information in a time point outside ChatGPT's training data cutoff [30].

**Theme 6. Content ethics and integrity ($D_{ethic}$)** This theme represents the presence of unlawful, unethical, harmful, or biased content

| Category (7) | Description | Code (19) | Example |
|---|---|---|---|
| Intent Understanding ($D_{intent}$) | This response does not correctly reflect the user's intent, instruction, or context. | C1. Response does not meet users' intent or instruction. | [42] |
| | | C2. Response is not aligned with the user's context. | [70] |
| | | C17. The tone or communication style is disappointing. | [71] |
| Content Depth and Originality ($D_{depth}$) | This response is overly general, lacks originality, or needs more diversity. | C3. Response is too general. | [14] |
| | | C4. Response lacks originality. | [45] |
| | | C5. Response lacks information. | [47] |
| Information Accuracy ($D_{acc}$) | This response contains false/inaccurate information or inconsistency. | C6. The response contains incorrect information. | [8] |
| | | C7. Response is based on training data cut off at a certain date, and has limited access to newly created data. | [92] |
| | | C8. Response is inconsistent. | [6] |
| | | C9. ChatGPT struggles with reasoning. | [95] |
| | | C10. (Hallucination) ChatGPT fabricates contents that conflict with the source content or cannot be verified from existing sources. | [39] |
| | | C19. (Sycophancy) ChatGPT excessively conforms to the user. | [64] |
| Transparency ($D_{trans}$) | It is difficult to understand the underlying reasoning or criteria of this response. | C11. It's difficult to understand the reasons, criteria, logic, and evidence behind the responses. | [79] |
| Refusal to Answer ($D_{refuse}$) | ChatGPT avoids answering by saying something similar to "As a language model, I am not capable …" | C12. ChatGPT avoids giving its own opinion by saying something similar to "As a language model, I am not capable …" | [14] |
| | | C13. ChatGPT avoids talking about difficult or controversial issues by saying something similar to "As a language model, I am not capable ..." | [9] |
| | | C7. Response is based on training data cut off at a certain date, and has limited access to newly created data. | [30] |
| Content Ethics and Integrity ($D_{ethic}$) | This response contains unlawful, unethical, harmful, or biased content. | C14. Response contains unlawful content | [51] |
| | | C15. Response contains unethical, harmful content. | [86] |
| | | C16. Response contains biased content. | [18] |
| Response Format and Attitude ($D_{format}$) | The format of this response — including but not limited to tone, length, structure, and attitude — is disappointing. | C17. The tone or communication style is disappointing. | [30] |
| | | C18. Response is overly detailed or too long | [90] |
| | | C19. (Sycophancy) ChatGPT excessively conforms to the user. | [9] |

Table 1: 7 category and corresponding 19 codes of user-side dissatisfaction from LLM Responses.

in LLM responses. Illegal and dangerous information was found to be accessible through LLMs [90], as well as stereotypes, discriminatory views, and performance disparity in certain groups [86]. The risk of LLMs not only generating but potentially magnifying existing social biases is a matter of concern as well [18].

**Theme 7. Response Format and Attitude ($D_{format}$)** Dissatisfaction with the format of responses, including tone, length, structure, and overall attitude, was captured within this theme. This dissatisfaction can arise when users have expectations regarding the manner in which responses were delivered and the tone used by

the LLM. ChatGPT's choice of words and formal, dry tone [30], as well as extensive and detailed responses [90] are quite different from human-generated text, which was colloquial and shorter.

These seven themes collectively offer a structured framework for understanding the multifaceted nature of user dissatisfaction with ChatGPT responses. Our survey utilized these themes as a basis for systematically investigating and quantifying user dissatisfaction.

# 4 DATA COLLECTION

Based on the categorization of user-side dissatisfaction from LLM responses, we collected the actual user's ChatGPT conversation log data with the dissatisfaction through a data collection system we designed and implemented. Our system targeted individuals who have utilized ChatGPT for practical purposes such as increasing productivity or efficiency in work, study, or hobbies. This process aims to address the following three research questions:

**RQ1.** What and How much dissatisfaction do users experience from LLM-generated responses?

**RQ2.** How do users address these dissatisfactions in their subsequent prompts during the conversation with LLM?

**RQ3.** How do user dissatisfaction and tactics vary depending on users' knowledge level regarding LLMs?

## 4.1 Data Collection System Design

To collect users' ChatGPT conversation log data in the wild, we designed and implemented a data collection system that includes the following four stages.

**Stage 1. Answering a General Questionnaire** In the first stage, we collected demographic information of participants such as gender, age, occupation, and overall experiences with ChatGPT (e.g., the frequency and period of using ChatGPT in their workflow). We also asked about the participants' knowledge level regarding Large Language Models (LLM) ("Regarding the mechanisms of Large-language models such as ChatGPT, how much do you agree with the following statement?"). All questions in this stage were measured through a 7-point Likert scale.

**Stage 2. Looking Through ChatGPT Chat History** In stage 2, participants were instructed to review their ChatGPT conversation history that had happened within 30 days. While reviewing, we asked the participants to find a conversation in which they experienced dissatisfaction with ChatGPT responses. To facilitate participants to think of various cases of dissatisfaction, we provided the descriptions of dissatisfaction categories derived from our systematic literature review as examples.

**Stage 3. Submitting Dissatisfactory Conversations** Based on their reflections regarding dissatisfaction in stage 2, we requested the participants to share a ChatGPT conversation link [6] within the past 30 days in which they experienced at least one dissatisfactory response. The participants can input the link into our system. To collect the conversation data with the details of the context, we also asked them to provide information about the purpose of the conversation, the reasons for using ChatGPT in that context, and the version of ChatGPT they used in this conversation, like GPT-3.5. Lastly, we asked the participants how much they remembered the conversation.

**Stage 4. Answering Questions About Dissatisfactory Responses** The participant's shared link was processed by transforming ChatGPT responses and user prompts to be presented as selectable components in the system (Fig 2-a). The system also allowed participants to provide specific experiences of dissatisfactory responses by selecting each response (Fig 2-b~f). For each selected response, participants were asked to (1) rate the overall level of dissatisfaction on a scale of 1 to 10 (1: a little dissatisfied, 10: extremely dissatisfied) (Fig 2-b), (2) choose one or more dissatisfaction categories from the given seven categories, or optionally describe a custom dissatisfaction point for dissatisfaction (Fig 2-c), (3) rate the level of dissatisfaction for each selected category on a scale of 1 to 10 (1: a little dissatisfied, 10: extremely dissatisfied) (Fig 2-c), (4) provide a detailed free-form explanation for their dissatisfaction (Fig 2-d), (5) select a prompt among the subsequent conversations in which they tried to resolve the dissatisfaction (Fig 2-e), (6) describe their tactic to address the dissatisfaction in the prompt (Fig 2-f), (7) rate the effectiveness of their tactic on a scale of 1 to 10 (1: not effective, 10: highly effective) (Fig 2-f), (8) provide a written explanation of the reasons for their effectiveness rating (Fig 2-f). In cases where there was no subsequent prompt or the conversation ended after dissatisfaction, participants were asked to provide written reasons instead of responding to (5)-(8).

## 4.2 Collected Data

*4.2.1 Participants and Collected Data.* We distributed the data collection system to people over the age of 18 globally through the Prolific platform [7]. Participants who provided at least two ChatGPT conversation links and evaluated at least one dissatisfactory response for each link received a compensation of £6. For each additional dissatisfactory response submitted from a single conversation link, participants received an additional £0.75 per response. For each additional conversation link provided beyond the initial two, participants received an additional £1.5 per link. We limited the number of maximum conversation links that can be submitted to five for each participant to prevent one participant from providing lots of conversation links. In total, we collected 307 ChatGPT conversation links, 511 dissatisfactory ChatGPT responses, and 615 user responses regarding those dissatisfactions from 107 individuals. Each user submitted an average of 2.87 links (std=1.21), 4.78 dissatisfactory ChatGPT responses (std=5.61), and 5.75 responses regarding those dissatisfactions (std=6.62). This study was approved by our institution's IRB, and we received consent from participants for the release of datasets.

*4.2.2 Data Filtering and Pre-processing.* To ensure the quality and reliability of the data collected from our system, two authors reviewed all the data together according to the following criteria and conducted filtering or pre-processing where necessary.

**Filtering Process** The data was filtered out at three levels: (1) user, (2) conversation, and (3) dissatisfactory responses.

*1. User-Level Filtering* We identified that one participant provided altogether contradictory responses, which contradicted the dissatisfactory response and the effectiveness of the prompt in resolving the dissatisfaction. Consequently, all data from this user were excluded.

*2. Conversation-Level Filtering* The conversation-level filtering was conducted based on the following four criteria, and a total of 20 conversations were filtered out. The detailed reason for each criteria is in the Appendix (Sec A.2).

(1) Conversation older than 30 days.

(2) Conversation with a memory level of 3 or lower.

(3) Conversation for fun or testing purposes.

---

[6] https://help.openai.com/en/articles/7943611-create-a-shared-link

[7] https://www.prolific.co/

Figure 2: Screenshot of the data collection system.

(4) Conversation from versions other than GPT-3.5.

**3. *Response-Level Filtering*** Response-level filtering was conducted based on the following four criteria, leading to the exclusion of a total of 16 dissatisfactory ChatGPT responses.

(1) Dissatisfaction due to ChatGPT's error messages
(2) Unconvincing dissatisfaction
(3) Mismatch between score and reason
(4) No Correlation between selected dissatisfactions and subsequent prompts for resolving that dissatisfaction

Detailed reasons and examples of each filtering case can be found in the Appendix and supplementary material. Please note that when filtering at the response level, all associated subsequent prompts and tactic data related to that response were also filtered. When filtering at the conversation level, all data related to the ChatGPT dissatisfactory responses and user prompts within that conversation were also filtered out. When filtering at the user level, all data provided by that user were excluded.

**Pre-processing Process** The data pre-processing process primarily involved the reassignment of dissatisfaction categories. This

step was undertaken to deal with cases where participants incorrectly selected dissatisfaction categories or opted for the 'other' option when evaluating the dissatisfaction category. Two authors examined all the data and carried out reassignment according to the following two criteria, proceeding only when a consensus was reached. Detailed examples of each case where reassignment occurred can be found in the supplemental.

Criterion 1: Reassigning 'other' to a specific category. For the 'other' option, when we found that there was a more suitable match with another category that was not selected based on the dissatisfaction reason, the 'other' score was reallocated to the corresponding category. As a result of this criterion, four entries were reassigned to the $D_{intent}$ category, two to $D_{depth}$, three to $D_{acc}$, and five to $D_{format}$.

Criterion 2: Reassigning an incorrectly selected category to another. If a participant had only checked one dissatisfaction category, and upon reviewing the dissatisfaction reason and conversation, it was evident that the selected category was not appropriate but another category was a better fit, the score was reassigned to the more suitable category. Using this criterion, three entries were reallocated from $D_{intent}$ to $D_{format}$, three from $D_{intent}$ to $D_{acc}$, two from $D_{acc}$ to $D_{intent}$, one from $D_{acc}$ to $D_{depth}$, and two from $D_{depth}$ to $D_{format}$.

*4.2.3 Dataset.* After filtering and pre-processing, we built a dataset on end-users' dissatisfaction with ChatGPT and their responses. The dataset is hierarchically organized, comprising the following components:

(1) User (N=94)
(2) ChatGPT conversation links and logs (N=249)
(3) User's recollected experience data on dissatisfactory ChatGPT responses (N=377)
(4) User's strategies to respond to the dissatisfactory response (N=459)

Here, the user's strategies were qualitatively analyzed, resulting in the creation of 13 tactic codes categorized into four themes. More detail of this is in Sec 5.2. Each data is also labeled as corresponding tactic codes by the authors. With this dataset, we conducted a quantitative and qualitative analysis to answer our research questions. We provide this dataset to facilitate future research about user experiences on chat-based LLMs. In releasing the dataset, we took careful consideration by masking all sensitive information related to their privacy and personal information. A more detailed description about the dataset can be accessed through our project website [8].

## 5 DATA ANALYSIS AND RESULTS

In this section, we present the analysis method and results that answer our research questions based on the constructed dataset. Firstly, we present the analysis of the types of dissatisfaction users face in LLM responses (RQ1). Next, we present how users respond to dissatisfaction through qualitative analysis (Table 3) and analyze the effectiveness of the tactics users use (RQ2). Finally, we present how users' knowledge level regarding LLM influences their experiences

---

[8]https://chatgpt-analysis.kixlab.org

of dissatisfaction and their behaviors when they face dissatisfaction (RQ3).

## 5.1 RQ1. Analysis of how users experience dissatisfaction

*5.1.1 Dissatisfaction Category Analysis.* We analyzed the count, distribution, and dissatisfaction score of the seven categories of dissatisfaction organized through a systematic literature review in Section 3, and the results are described in Table 2. In terms of the count of each category, $D_{intent}$ accounted for the largest proportion (32.18%), while $D_{trans}$, $D_{refuse}$, and $D_{ethic}$ constituted significantly smaller proportions compared to the other categories. To investigate the severity degree of user dissatisfaction in each category, we conducted Kruskal-Wallis test and confirmed significant differences between categories ($\chi^2$ = 17.6, p-value < 0.01, df = 6). In particular, we found that $D_{acc}$'s dissatisfaction score was the highest, and its score was statistically significantly higher than $D_{depth}$ through Dwass-Steel-Critchlow-Fligner(DSCF) pairwise comparison (p-value=0.008). This means that users are statistically significantly more dissatisfied with dissatisfaction due to $D_{acc}$ than $D_{depth}$.

Considering that each user provided multiple dissatisfactory responses, we also conducted a user-level analysis, accounting for potential correlations among the data submitted by the same user. To achieve this, we normalized each dissatisfaction category data by dividing them by the number of dissatisfactory responses each user submitted. This method allowed us to express each data point as the frequency of how often each user experienced dissatisfaction in a certain category. The analysis results are presented in Table 2 in the "User-level" analysis column. The mean frequency value of $D_{intent}$ was 0.47, indicating that if a user has experienced 100 dissatisfactory ChatGPT responses, on average, 47 of them fall into the $D_{intent}$ category. Furthermore, the Kruskal-Wallis test result shows statistically significant differences in user-level frequency values between each category ($\chi^2$ = 9.93, p-value < 0.01, df = 6). In the user-level analysis, we can see a similar tendency to the response-level analysis, users experience $D_{intent}$ the most frequently. Following this, the second most frequently encountered dissatisfaction is $D_{depth}$. However, the standard deviation of $D_{depth}$ is 0.35, which is much higher than other categories, indicating that the frequency of experiencing $D_{depth}$ varies significantly from user to user.

*5.1.2 Co-occurrence Analysis.* In a single dissatisfactory response, multiple dissatisfaction categories can co-occur. For example, a user may simultaneously experience dissatisfaction with the lack of originality ($D_{depth}$) and the length ($D_{format}$) of ChatGPT's response at the same time. Therefore, we analyzed co-occurrence patterns to investigate the correlations between each category of dissatisfaction. Results are presented in Fig 3 and the value at (*i, j*) in this matrix represents the frequency of when the *i*-th row was selected as a source of dissatisfaction, the *j*-th column was also selected together. The result shows that $D_{intent}$ frequently appears concurrently with all other categories. Also, while $D_{trans}$ and $D_{ethic}$ have relatively low counts, they co-occur with $D_{intent}$ more than half the times in each occurrence.

| Dissatisfaction Category | Response-level analysis | | User-level analysis |
|---|---|---|---|
| | Count: N (%) | Dissatisfaction Score: mean (std)* | Frequency: mean (std)* |
| $D_{intent}$ | **168 (32.18%)** | 5.56 (2.94) | **0.47 (0.03)** |
| $D_{depth}$ | 107 (20.50%) | **5.09 (2.69) *** | 0.33 (0.35) |
| $D_{acc}$ | 83 (15.90%) | **6.52 (2.76) *** | 0.20 (0.03) |
| $D_{trans}$ | 27 (5.17%) | 4.81 (3.13) | 0.08 (0.02) |
| $D_{refuse}$ | 27 (5.17%) | 6.37 (2.68) | 0.09 (0.02) |
| $D_{ethic}$ | 4 (0.77%) | 6.25 (3.20) | 0.01 (0.01) |
| $D_{format}$ | 106 (20.31%) | 6.14 (3.04) | 0.27 (0.03) |

Table 2: Analysis results on the count, dissatisfaction score, and user-level frequency for the dissatisfaction category (* p-value < 0.01)



Figure 3: Normalized Co-occurrence matrix of dissatisfaction category. The value at $(i, j)$ in this matrix represents the frequency of when the $i$th row was selected as a dissatisfaction point, the $j$th column was also selected as a dissatisfaction.

## 5.2 RQ2. Analysis of how users respond to dissatisfaction

*5.2.1 Categorizing Tactics for Resolving Dissatisfaction.* Through qualitative analysis, we categorized users' tactics to understand and analyze how users address their dissatisfaction from ChatGPT's response through subsequent prompts. Two authors independently conducted open coding by reviewing ChatGPT conversation log data, user-side dissatisfactions on ChatGPT responses, employed tactics in subsequent prompts, and user-reported effectiveness and the reasons for these tactics. After completing the open coding, the two authors engaged in an iterative process of code consolidation. To precisely capture and categorize the subtleties of user tactics, both authors iterated all data together, making a code set through discussion. We proceeded with these processes until the authors met a common ground. After two times of iterations, we identified the user's tactic with 13 codes as presented in Table 3. To establish relationships between these codes and identify overarching themes, axial coding [78] was performed. Through this coding process, we identified four main themes of the user's tactics, as presented in Table 3.

**Tactic Category 1. Prompt Reusing and Shortcut** This category of tactic represents users either reusing prompts or employing a single word to request similar or diverse responses, often requiring minimal effort in crafting the prompt. This category comprises three tactics. First, users just reuse the exact same prompt as the previous one or paraphrase it slightly (T1). Second, users use a single word like 'more' or 'another' as a shortcut to get either similar responses from the previous turn or a wider range of responses from ChatGPT (T2). Last, users retry by adding emphasis through formatting, such as using all capital letters or using double quotation marks (T3).

**Tactic Category 2. Intent Concretization** This category encompasses four tactics of users trying to concretize their intent and context to get a more appropriate response. Users further specify their needs by providing more detailed or direct instructions (T4), giving additional context or explanation (T5). For example, if users ask ChatGPT to recommend a dinner menu and they doesn't like ChatGPT's answer, they can further specify their needs by saying, "Recommend a **healthy** dinner menu using **tomatoes**" (T4), or explain their context by saying, "I'm going to invite a guest to my house for my dinner" (T5). And users concretize their intent by adding specific conditions related to the format such as "make it

| Category (4) | Tactic Code (13) |
|---|---|
| Prompt Reusing and Shortcut ($T_{repeat}$) | T1: Re-using an identical prompt or slightly paraphrasing it |
| | T2: Using the specific word (e.g., more, another) that implies requesting different or more outputs for the same task as the previous prompt |
| | T3: Re-using an identical prompt but adding emphasis through formatting (e.g., using all capital letters, using double quotation marks) |
| Intent Concretization ($T_{specify}$) | T4: Specifying user intent by providing detailed or direct instructions |
| | T5: Specifying user intent by providing additional context or explanation |
| | T6: Adding format-specific conditions (e.g., make it shorter, provide in list format) |
| | T7: Adding tone-specific conditions (e.g., make it casual) |
| Error Identification and Correction ($T_{error}$) | T8: Pointing out errors or mistakes |
| | T9: Providing the correct answer or hints |
| | T10: Asking clarification questions |
| Task Adaptation ($T_{adapt}$) | T11: Adapting by shifting to another topic or task that is different from the original intent. |
| | T12: Breaking down the original task into smaller subtasks |
| | T13: Asking follow-up questions deviating from the original task |
| **No Tactic** | No further prompting to address the dissatisfaction and even terminating the conversation due to dissatisfaction |

**Table 3: User tactic category**

shorter" (T6), and adding specific conditions related to the tone, such as "make it casual" (T7).

**Tactic Category 3. Error Identification and Correction** This category mainly contains tactics when there are some errors in the ChatGPT's response, and the users point out or correct them. Users simply say "It was wrong." or point out the part that is wrong (T8), give the correct answer or hints of the correct answer (T9), and ask a clarification question to confirm the error or doubtful aspects such as by asking "Can you confirm that ... ?" or "Are you sure ...?"(T10).

**Tactic Category 4. Task Adaptation** This category represents the user adjusting to another task instead of the original task where the user felt dissatisfied. Users adapt their task by altering their initial task to a different one (T11). For instance, if users initially ask for the latest information and ChatGPT says it can only answer up to 2021 information, then they can slightly adjust their original task and ask for 2021 information rather than the latest information. Users also adjust their original task by dividing it into smaller and more manageable subtasks (T12). For example, when users ask ChatGPT for a complex math problem, they can ask them in intermediate steps. Finally, Users ask follow-up questions deviating from the original task, such as asking follow-up questions about parts that lack details or are unfamiliar to them in ChatGPT's responses. (T13).

*5.2.2 Tactic Category Analysis.* After creating the tactic categories, we categorized users' prompts into four tactic categories or **No Tactic**. **No Tactic** indicates no further prompting to address the

dissatisfaction and even terminating the conversation due to dissatisfaction. Here, note that a single user prompt can encompass multiple tactic categories if the prompt contains multiple requests. We conducted response-level analysis for the count, distribution, and effectiveness of each tactic as well as user-level analysis for frequency (Table 4). Notably, we observed that $T_{specify}$ stands out as the dominant category, and it accounts for over half of the distribution (58.6%) among the four tactic categories without **No Tactic**. In addition, we analyzed the effectiveness of each tactic based on users' rating of the effectiveness score between 1 and 10. We conducted a Kruskal-Wallis test and confirmed that there are statistically significant differences between the effectiveness scores of each tactic ($\chi^2$ = 23.1, p-value < 0.01, df = 4). Specifically, we found that $T_{specify}$, a tactic for users to further specify their own intents, is most effective with a mean score 0f 6.04, highest of all categories.

*5.2.3 Dissatisfaction Category and Corresponding Tactics: Whether the dissatisfaction was solved.* We investigated how users applied different tactics to address each dissatisfaction category and whether these tactics resolved the dissatisfaction. Firstly, we analyzed the distribution of tactics used for each dissatisfaction category (Fig. 4(a)), and drew a Sankey diagram to visualize the overall flow of tactics used by each dissatisfaction category (Fig. 4(b)). We observed that $T_{specify}$ is the dominant tactic across various dissatisfaction categories. However, when users encounter dissatisfaction related to the accuracy of information ($D_{acc}$), they tend to employ $T_{error}$ rather than $T_{specify}$. Lastly, in cases of $D_{trans}$, $D_{refuse}$, and $D_{ethic}$, users often resort to **No Tactic**, ending up the conversation. The

| Tactic Category | Tactic Code | Response-level analysis | | | | User-level analysis | |
|---|---|---|---|---|---|---|---|
| | | Count: N (%) | | Effectiveness Score: mean (std) | | Frequency: mean (std) | |
| | | Category | Code | Category* | Code | Category | Code |
| $T_{repeat}$ | T1 | 45 (9.4%) | 29 (5.8%) | 4.04 (3.16) | 4.45 (3.15) | 0.09 (0.20) | 0.07 (0.18) |
| | T2 | | 18 (3.6%) | | 3.06 (3.06) | | 0.02 (0.09) |
| | T3 | | 2 (0.4%) | | 1.00 (0.00) | | 0.00 (0.04) |
| $T_{specify}$ | T4 | **183 (38.4%)** | 122 (24.4%) | **6.04 (3.44)** | 6.25 (3.53) | **0.33 (0.34)** | 0.22 (0.28) |
| | T5 | | 26 (5.2%) | | 5.35 (3.33) | | 0.06 (0.15) |
| | T6 | | 40 (8.0%) | | 6.45 (3.16) | | 0.08 (0.17) |
| | T7 | | 11 (2.2%) | | 4.73 (3.04) | | 0.02 (0.10) |
| $T_{error}$ | T8 | 73 (15.3%) | 53 (10.6%) | 4.19 (2.95) | 4.26 (2.99) | 0.10 (0.22) | 0.06 (0.16) |
| | T9 | | 13 (2.6%) | | 4.62 (2.66) | | 0.02 (0.09) |
| | T10 | | 10 (2.0%) | | 3.80 (3.16) | | 0.03 (0.10) |
| $T_{adapt}$ | T11 | 12 (2.5%) | 7 (1.4%) | 5.17 (3.04) | 4.57 (3.21) | 0.04 (0.11) | 0.03 (0.10) |
| | T12 | | 2 (0.4%) | | 8.00 (0.00) | | 0.00 (0.03) |
| | T13 | | 3 (0.6%) | | 4.67 (3.22) | | 0.00 (0.04) |
| **No Tactic** | | 164 (34.4%) | 164 (32.8%) | - | - | 0.47 (0.38) | 0.47 (0.36) |

**Table 4: Analysis results on the count, effectiveness score, and user-level frequency for the tactic category (* p-value < 0.01)**



**Figure 4: (a) Distribution of tactic categories by dissatisfaction category. (b) Sankey diagram to visualize how users respond among four tactic categories or No Tactic after experiencing each of the dissatisfaction categories. Note that the count in the Sankey diagram can be greater than the count of response-level analysis in Table 2 and 4. This is because one response can include multiple dissatisfaction categories and multiple tactic categories, and they were counted multiple times to draw a Sankey diagram.**

proportion and visualization of whether or not dissatisfaction has been resolved by each tactic can be seen in Fig. 5(a). Fig. 5(a) illustrates that the users managed to resolve their dissatisfaction by 58 % by utilizing tactics. Notably, $T_{specify}$ was an effective way of resolving dissatisfaction in many cases (67%), while with other tactics, there were more cases where dissatisfaction remained unsolved. Fig. 5(b) shows which tactics users use for each dissatisfaction category and how this eventually leads to resolve the dissatisfaction.

Through this analysis, we can observe the overall flow of how users, while conversing with ChatGPT, experience various dissatisfactions in what proportion, how they respond to them using different tactics, and how this leads to the resolution of these dissatisfactions. When users encounter dissatisfaction, approximately 34% opt for **No Tactic** while 66% employ tactics. However, it can be seen that approximately 58% of dissatisfactions are resolved through tactics.

Figure 5: (a) A Sankey diagram that visualizes whether users resolved their dissatisfaction using each of the tactic categories. (b) The overall visualization of how users respond among the four tactic categories after experiencing each of the dissatisfaction categories and finally whether that dissatisfaction was solved or not.



Figure 6: Distribution of participants' knowledge level regarding LLM on a 7-point scale (1: very low, 7: very high). None of the participants reported a knowledge level of 1.

In the end, users manage to resolve only 28% of their dissatisfactions using tactics, leaving 72% of dissatisfactions unresolved.

## 5.3 RQ3. Analysis of how dissatisfaction and tactics vary based on the user's knowledge level of LLMs

We analyzed how users' experience of dissatisfaction and their tactics differ depending on their knowledge levels regarding LLMs. First, we examined the distribution of users' knowledge levels regarding LLMs in our dataset, as depicted in Fig. 6. We collected the knowledge level data about LLMs on a 7-point scale, where 1 indicates very low kn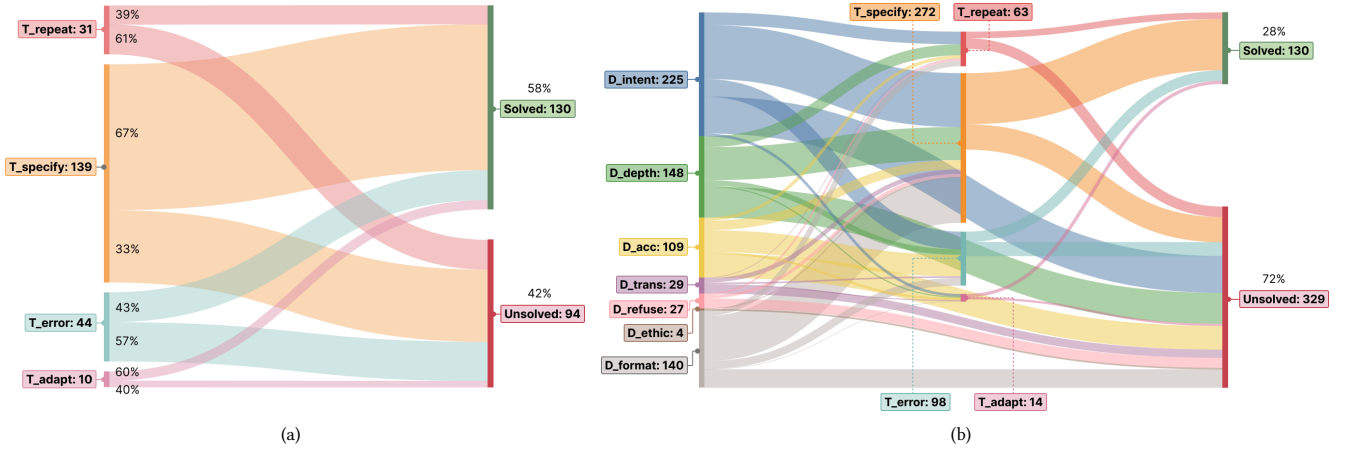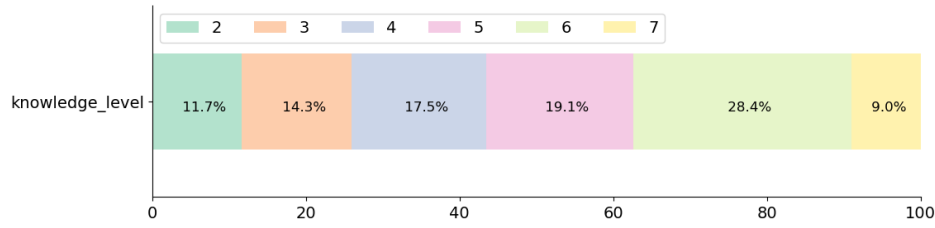owledge, and 7 indicates very high knowledge. We divided the groups into "low knowledge level" (those with a knowledge level 1-3) and " high knowledge level" (those with a knowledge level 5-7), as four lies in the middle of the 7-point scale.

To investigate whether there is a difference in the distribution of dissatisfaction categories between the two groups, we conducted a chi-square test for the dissatisfaction categories of each group and found that there were statistically significant differences in the distribution of dissatisfaction categories by different knowledge groups ($\chi^2$ = 17.7, p-value < 0.01). Specifically, we observed that the low-knowledge group experiences $D_{depth}$ ((count: 26.97%,

user-level frequency: 0.38)) and $D_{refuse}$ (count: 8.55%, user-level frequency: 0.14) more frequently, while the high-knowledge group experiences $D_{acc}$ (count: 17.38%, user-level frequency: 0.24) and $D_{format}$ (count: 24.82%, user-level frequency: 0.28) more frequently. On the other hand, we conducted a Mann-Whitney U test to investigate the differences in dissatisfaction scores between knowledge groups, but there were no significant differences.

Similarly, we conducted a chi-square test for tactic categories and found significant differences in the count of tactic categories among the two groups ($\chi^2$ = 21.6, p-value < 0.01). In particular, **No Tactic** was more prevalent in the low-knowledge group. Additionally, $T_{repeat}$, which involves minimal prompt engineering, was more commonly used in the low-knowledge group, while $T_{error}$, aimed at pointing out and rectifying errors in ChatGPT responses, was more prevalent in the high-knowledge group. Furthermore, to compare and analyze the effectiveness of the tactics used in each knowledge group, we performed a Mann-Whitney U test on the effectiveness scores of tactic categories, which were collected from users. Through this test, we found that the effectiveness of the $T_{repeat}$ was statistically higher in the high-knowledge group (p-value < 0.01, effect size= 0.5789). Fig. 7(a) and 7(b) present Sankey diagrams that illustrate how users in the low-knowledge and high-knowledge

| Dissatisfaction Category | Response-level analysis | | | | User-level analysis | |
|---|---|---|---|---|---|---|
| | Count: N (%) * | | Dissatisfaction Score: mean(std) | | Frequency: mean (std) | |
| | high | low | high | low | high | low |
| $D_{intent}$ | 89 (31.56%) | 45 (29.61%) | 5.91 (2.85) | 5.18 (3.08) | 0.43 (0.30) | 0.49 (0.39) |
| $D_{depth}$ | **50 (17.73%)** | **41 (26.97%)** | 5.02 (2.70) | 5.22 (2.72) | **0.30 (0.31)** | **0.38 (0.38)** |
| $D_{acc}$ | 49 (17.38%) | 18 (11.84%) | 6.73 (2.85) | 6.5 (2.62) | **0.24 (0.29)** | **0.14 (0.21)** |
| $D_{trans}$ | 12 (4.26%) | 9 (5.92%) | 5.25 (3.33) | 3.67 (3.00) | 0.07 (0.16) | 0.10 ( 0.23) |
| $D_{refuse}$ | **11 (3.90%)** | **13 (8.55%)** | 6.82 (2.79) | 6.92 (2.02) | **0.07 (0.16)** | **0.14 (0.26)** |
| $D_{ethic}$ | 1 (0.35%) | 3 (1.97%) | 3 (-) | 7.33 (2.89) | 0.01 (0.07) | 0.03 (0.08) |
| $D_{format}$ | 70 (24.82%) | 23 (15.13%) | 6.66 (2.86) | 5.7 (3.36) | **0.28 (0.32)** | **0.25 (0.37)** |

Table 5: Dissatisfaction category for knowledge level high and low group (* p-value < 0.01)

| Tactic Category | Tactic Code | Response-level analysis | | | | | | | | User-level analysis | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Count: N (%) | | | | Effectiveness Score: mean (std) | | | | Frequency: mean (std) | | | |
| | | Category* | | Code | | Category | | Code | | Category | | Code | |
| | | high | low | high | low | high | low | high | low | high | low | high | low |
| $T_{repeat}$ | T1 | 16 (6.11%) | **19 (14.5%)** | 12 (4.4%) | 9 (6.5%) | **5.06 (3.00)\*** | **2.37 (2.27)\*** | 4.75 (2.96) | 3.00 (3.00) | 0.08 (0.17) | 0.11 (0.25) | 0.06 (0.14) | 0.08 (0.22) |
| | T2 | | | 4 (1.5%) | 12 (8.7%) | | | 6.00 (3.37) | 1.67 (1.15) | | | 0.02 (0.08) | 0.04 (0.13) |
| | T3 | | | 1 (0.4%) | 0 (0.0%) | | | 1 (-) | - (-) | | | 0.00 (0.02) | 0.23 (0.35) |
| $T_{specify}$ | T4 | 111 (42.37%) | 52 (39.7%) | 84 (30.8%) | 24 (17.4%) | 5.88 (3.56) | 6 (3.33) | 5.77 (3.71) | 7.17 (2.78) | 0.34 (0.31) | 0.39 (0.41) | 0.23 (0.25) | - |
| | T5 | | | 13 (4.8%) | 12 (8.7%) | | | 5.00 (3.03) | 5.42 (3.73) | | | 0.05 (0.11) | 0.10 (0.22) |
| | T6 | | | 19 (7.0%) | 13 (9.4%) | | | 7.00 (2.83) | 6.00 (3.70) | | | 0.09 (0.20) | 0.07 (0.14) |
| | T7 | | | 4 (1.5%) | 7 (5.1%) | | | 6.00 (4.08) | 4.00 (2.31) | | | 0.01 (0.05) | 0.05 (0.18) |
| $T_{error}$ | T8 | **49 (18.70%)** | 8 (6.1%) | 37 (13.6%) | 5 (3.6%) | 3.53 (2.60) | 5.75 (3.06) | 3.81 (4.08) | 5.2 (3.83) | 0.12 (0.24) | 0.08 (0.18) | 0.07 (0.18) | 0.05 (0.12) |
| | T9 | | | 7 (2.6%) | 2.00 (1.4%) | | | 3.57 (1.90) | 7.00 (1.41) | | | 0.03 (0.09) | 0.02 (0.10) |
| | T10 | | | 6 (2.2%) | 2 (1.4%) | | | 2 (2.45) | 7.00 (1.41) | | | 0.03 (0.12) | 0.03 (0.12) |
| $T_{adapt}$ | T11 | 5 (1.91%) | 1 (0.8%) | 4 (1.5%) | 1 (0.7%) | 4.40 (3.29) | 1 (-) | 5.25 (3.10) | 1 (-) | 0.03 (0.11) | 0.01 (0.07) | 0.03 (0.11) | 0.01 (0.07) |
| | T12 | | | 0 (0.0%) | 0 (0.0%) | | | - | - | | | - | - |
| | T13 | | | 1 (0.4%) | 0 (0.0%) | | | 1 (-) | - | | | 0.00 (0.01) | - |
| **No Tactic** | | 81 (30.92%) | **51 (38.9%)** | 81 (29.7%) | 51 (37.0%) | - | - | - | - | 0.47 (0.37) | 0.44 (0.39) | 0.47 (0.37) | 0.44 (0.39) |

Table 6: Tactic category and code for knowledge level high and low group (* p-value <0.01)

groups experience dissatisfaction categories from ChatGPT's responses, respond to the dissatisfactions with each tactic category at user prompts, and whether these tactics ultimately resolve their dissatisfactions or not. Through this, we can see that the rate of resolving dissatisfaction in the high-knowledge group (29%) is higher than low-knowledge group (23.5%).

## 6  DISCUSSION

In this section, we first discuss the interpretation of our results and their implications. Second, we suggest design implications for building LLMs with better usability based on our study results. Lastly, we discuss the limitations of our study and future work.

### 6.1  Interpretation of results

Building upon the analysis of user-side dissatisfaction and corresponding user tactics during the conversation, we discuss the most

prevalent, severe, and unaddressed categories of dissatisfaction and their implications. We also discuss the differences in dissatisfaction and corresponding tactics across users with different knowledge levels about LLM.

*6.1.1  The Most Prevalent Dissatisfaction and Tactics.* Our results suggest that despite the advances in LLMs to align with the user intent, there still exists much room for improvement from the users' perspective. With recent advancements in LLMs and the introduction of techniques to align LLMs with user intents, such as Reinforcement Learning from Human Feedback (RLHF), LLMs are now known to better align with human intent than before [19, 63, 101]. However, we found that $D_{intent}$, the dissatisfaction in terms of understanding users' intent, is the most prevalent (Table 2) and frequently co-occurring with other dissatisfaction categories (Figure 3). We also discovered that users frequently use $T_{specify}$ that further
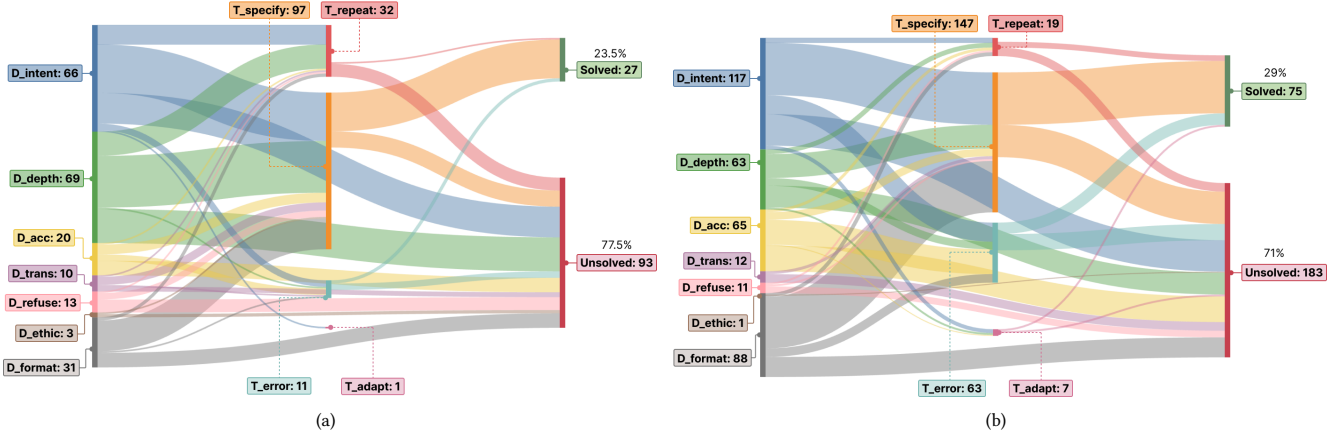
**Figure 7: Sankey diagrams by users' knowledge level of LLMs that visualize how users respond among four tactic categories after experiencing each of the dissatisfaction categories and finally whether that dissatisfaction was solved or not. (a): Low-knowledge group's Sankey diagram (b): High-knowledge group's Sankey diagram.**

specify their intent to address the dissatisfaction. Moreover, users rated $T_{specify}$ as the most effective among tactic categories, but there are still many cases (about 42%) where dissatisfaction was not resolved despite using this tactic. This may be because users have difficulty clearly representing their intent. Previous work on web search and information retrieval has also noticed this problem [35], and there exist several methods to better support users to specify their intent in these domains, such as context-sensitive query auto-completion [10] and context-based term suggestions [73]. Similarly, in the context of LLMs, further research can investigate methods to help users specify their intents based on their context.

*6.1.2 The Most Severe or Unaddressed Dissatisfaction.* Self-reported scores on the level of dissatisfaction show that users perceived the dissatisfaction of $D_{acc}$ to be the most severe (Table 5). This shows that users feel a high level of dissatisfaction with limitations of LLMs related to information accuracy such as hallucination [9, 39, 51, 91], inconsistency or incorrectness in the responses [24, 37, 38, 51], and the inability of ChatGPT to provide up-to-date information [6, 97]. Furthermore, our findings show that users tend to respond to this dissatisfaction primarily by pointing out LLM's faults or correcting them ($T_{error}$), but more than half of them (57%) nevertheless fail to resolve this dissatisfaction.

We also found that when users encountered dissatisfaction when their prompts were refused to answer ($D_{refuse}$), when ethical concerns or biases were found in the response ($D_{ethic}$), or when they had a lack of understanding of the internal logic of the generated response ($D_{trans}$), they often did not attempt to address the dissatisfaction or even terminated the conversation. For instance, one user explained their decision to end the conversation as follows: "I ended the conversation as I felt like there was no common understanding and was not looking forward to explaining myself any further than my original prompt." Through this, we can see that if the users experience such dissatisfaction, they not only have difficulty communicating with ChatGPT but also have no idea how to further improve their prompts, often terminating the conversation.

One notable point here is that $D_{ethic}$ and $D_{refuse}$ can be in a trade-off relationship. Including OpenAI [9], the company that developed ChatGPT, many companies have adopted a strategy where the LLM avoids answering when faced with potentially unethical or biased prompts, responding with statements like "As a language model, I am not capable of ... " [14, 98]. Although companies could avoid being embroiled in ethical issues, this approach might have introduced another dimension of dissatisfaction ($D_{refuse}$) for users. Self-reported scores on the level of dissatisfaction show the level of severity for both $D_{ethic}$ and $D_{refuse}$ are similar (Table 5). This suggests that the current approach of refusing to answer instead of giving responses with ethical concerns may not reduce users' overall dissatisfaction. Thus, it is necessary to find other measures that could also lower the users' dissatisfaction when faced with unethical or biased prompts.

*6.1.3 Differences in Dissatisfaction and Corresponding Tactics Across LLM Knowledge Levels.* Our result revealed that there exist significant differences in dissatisfaction and employed tactics between high- and low-knowledge user groups. We observed that the low-knowledge group reports higher occurrences of $D_{depth}$—dissatisfaction that ChatGPT's response is too general and lacks detail or originality—than the high-knowledge group (Table 5). One possible reason behind this is that the low-knowledge user group might have overestimated ChatGPT's creative capabilities. This could be because low-knowledge user groups may be more prone to unconditionally accepting media or news which states that Chat-GPT can perform creative tasks such as writing poetry and song lyrics [2, 23]. This may have led them to expect more creative responses, resulting in a higher possibility of feeling disappointment. In contrast, the high-knowledge group may have possessed a better understanding of ChatGPT's limitations. Knowing that ChatGPT's responses are based on trained patterns from existing datasets could have allowed them to be more generous towards the responses that lack originality. We speculate that the low-knowledge group might

---

[9]https://openai.com/

have a less accurate mental model of the capacity of LLM, misunderstand its capabilities, and experience more dissatisfaction in terms of $D_{depth}$.

Moreover, the tactics employed in response to these dissatisfactions differed between the two groups. Compared to the high-knowledge group, the low-knowledge group relied more on 'No tactic' and more frequently used $T_{repeat}$, which requires minimal effort for prompt writing (Table 6). This may be because the low-knowledge users may not know much about the various options of tactics they could take. Interestingly, however, although high-knowledge users used $T_{repeat}$ less, they found it more effective in solving their dissatisfaction. This may indicate that high-knowledge users tend to have a better sense of when is the right time to use $T_{repeat}$.

## 6.2 Design Implications for Building LLMs with Better Usability

Based on our study result, we suggest three design implications to enhance the usability of LLMs: (1) supporting users to represent their intent, (2) recommending effective multi-turn prompt tactics to users, and (3) providing personalized LLM experiences to users.

*6.2.1 Supporting users to represent their intent better.* We suggest a design that facilitates a better representation of the user's intent. In the current system interface, there is a lack of design support to help users' prompt writing process, and we found that users frequently face limitations in conveying their full intent in Sec 6.1.1. To address these challenges and facilitate a better representation of the user's intent, it is necessary to have a design that helps users refine their prompts to align them more precisely with their intent. This design could involve tokenizing user prompts and using this as a basis to offer keyword-specific suggestions. For example, if a user writes a prompt, "Explain recent issues related to Autonomous Vehicles (AVs) in simple terms.", keywords can tokenize the prompt, and the following keyword-specific suggestions can be provided: the types of AVs, the time frame for recent, the types of issues (e.g., ethical), and the appropriated level of simplicity for the terms used. Moreover, considering dissatisfaction arising from extensive and detailed responses ($D_{format}$), giving suggestions utilizing multi-modality, such as image and video, could enable a better user experience when they can succinctly represent the users' intent. This allows users to refine their prompts by selecting the suggestions, ensuring a more accurate alignment with their intent. Providing users with a range of suggestions and enabling them to select suggestions by reflecting their intent can empower users to express their intent effectively.

*6.2.2 Recommending effective multi-turn prompt tactics to users.* To enhance user satisfaction during multi-turn interactions with LLM, we suggest a design that recommends effective prompt tactics to users during the conversation. Our public dataset could be utilized for this process since it contains various prompt tactics (Table 3) and their effectiveness reported by users to address their dissatisfaction. For instance, an interaction can be envisioned where the LLM predicts the probability of user dissatisfaction with a generated response. If the probability is high, the system can proactively

guide users to employ some effective tactics in their subsequent prompt to address the anticipated dissatisfaction.

We also recommend evolving this design to incorporate effective prompt engineering techniques suitable for multi-turn interactions, such as Chain-of-Thought (CoT) [96]. While a thread of research has addressed effective prompt engineering techniques to get desired responses from LLMs, they usually focus on crafting one prompt. Moreover, there is a lack of research on prompt engineering techniques tailored to address or mitigate user dissatisfaction during conversations. By integrating our data-driven insights on users' effective prompt tactics with prompt engineering techniques, we propose that recommending tactics to users during multi-turn interactions will yield more favorable responses, enhancing their overall satisfaction.

*6.2.3 Providing personalized LLM experience.* We suggest the need for a design that provides personalized LLM experiences based on our finding that there exist differences in dissatisfaction and corresponding tactics depending on the user's level of knowledge about LLMs. One of the possible designs for personalized LLM experiences is to adjust the refusal policies or attitudes that LLM refuses to answer according to the user's knowledge levels about LLM. This is because our results show that the low-knowledge group experienced more dissatisfaction with ChatGPT's refusal to answer ($D_{refuse}$) than the high-knowledge group. This may be because the low-knowledge group tends to ask more questions that were limited for ChatGPT to answer without fully understanding ChatGPT's capabilities. Thus, rather than responding with a generic "As a language model, I am not capable of..." a more direct explanation addressing its limitations to better inform users of its capability may be required for low-knowledge users.

To facilitate personalized LLM experiences, we emphasize the need for user modeling based on prior sessions where LLM can gain information about the user's state before chatting. The user's state encompasses not only their knowledge level about LLM but also their usage purpose, specific task at hand, the language or proficiency level they used for chatting, and more. Such sessions serve to shape the user's mental model of LLM and vice versa, fostering a mutual understanding. Through this approach, users can benefit from customized interactions that consider their individual circumstances, ultimately improving their overall LLM experience.

## 6.3 Generalization of Results

While our study utilized ChatGPT as a case study, our research methodology and its implications can extend beyond ChatGPT. The seven categories of user-side dissatisfaction identified through our SLR (Table 1) encompassed references that span various LLMs. Hence, leveraging these categories and our analysis method, future research can apply similar investigations to different LLMs. In addition, the four categories of user tactics (Table 3) were derived from analyzing user behavior patterns in multi-turn conversations with LLM, based on ChatGPT user data. The consistent nature of user behavior across various LLMs with similar multi-turn chat-based interfaces suggests potential generalizability to other LLMs.

However, it is essential to consider potential variations that may arise due to specific features in LLMs, technical advancements, and changes in user perception towards LLMs. Even if the categories

remain constant, their distribution and severity may change. For instance, we can expect that while the proportion of dissatisfaction arising from accuracy ($D_{acc}$) might decrease as the performance of LLM improves, its perceived severity may intensify as user's expectations towards LLMs get higher. Furthermore, while many users may currently lack awareness of the ethical issues related to LLM, $D_{ethic}$ might increase as they become informed about the potential ethical threats posed by LLM. Therefore, extending our analysis to different LLMs or the same LLM over time allows for a comprehensive comparison of user dissatisfaction across various models and versions, providing insights into the direction of evolving LLMs. Our data analysis also showed that user dissatisfaction varied based on users' knowledge level regarding LLMs. From this, we may refer to the dissatisfaction distribution of the current high-knowledge group while inferring the dissatisfaction distribution of LLM users in the future. In terms of user tactics, the emergence of novel interaction components beyond chat-based interfaces may lead to different user behavior patterns, which would require further investigation.

## 6.4 Limitations and Future Work

We present the limitations of our work and possible future work.

First, our analysis is based on self-reported data from users. We tried to ensure the quality of the data by careful filtering and pre-processing of the data while checking on the actual conversation log. However, dissatisfaction levels and tactic effectiveness are based on participants' self-reported scores, which may suffer from subjectiveness and heavily rely on the participant's memory. We also tried to eliminate this problem by only collecting conversation logs within 1 month, but the problem may still linger.

Second, we investigated the difference in user dissatisfaction and tactics according to the difference in knowledge level of LLMs. Future work can expand on our work and further examine whether the differences in dissatisfaction and tactics exist according to other dimensions. For instance, since LLMs are chat-based, there may exist differences between those different English proficiency. Moreover, since users may have different expectations according to tasks, there may exist differences when given different tasks. For instance, fact-oriented tasks, such as finding information or explaining a real-world fact, will have more relevance with $D_{acc}$ since the user expects to get correct information. On the other hand, creative tasks, such as writing stories or scenarios, will have less relevance with $D_{acc}$ but more relevance with $D_{intent}$, since users will be interested in how well the LLM can understand their needed content or context of creating content to their situations.

Lastly, our analysis of user-side dissatisfaction and tactics was based on ChatGPT user data. Therefore, there may be some differences in how users undergo dissatisfaction depending on other LLMs. For instance, specific wordings used when LLM refuses to answer can affect how much users feel dissatisfaction regarding $D_{refuse}$. Moreover, since $D_{acc}$ is a category that is directly related to the performance of LLMs, users may face different levels of dissatisfaction for $D_{acc}$.

## 7 CONCLUSION

In this study, with ChatGPT as the case study, we explored user-side dissatisfaction and corresponding tactics during the conversation with chat-based LLM. Through a systematic literature review, we identified seven categories of user-side dissatisfaction from LLM-generated responses. Then, we collected data from 107 users conversing with ChatGPT, and uncovered prevalent, severe, and unaddressed dissatisfactions. We also analyzed four users' tactic categories to address their dissatisfaction and their prevalence and effectiveness. We also investigated how these vary depending on the users' knowledge level of LLMs. Our findings provide insights into how LLM and its interface can be further developed to aid people when they encounter dissatisfaction. One potential is user-side prompt engineering techniques that can be utilized in the middle of the conversation when dissatisfaction occurs. The pair of dissatisfactions and corresponding tactics can guide this prompt engineering. In addition to these contributions, we have made a publicly accessible dataset available, containing actual user conversation data related to dissatisfaction. This research deepens the understanding of user dissatisfaction in LLM interactions, providing a foundational knowledge base for future enhancements that can benefit users across knowledge levels.

## REFERENCES

[1] Accessed on 10/06/2023. LLM Jailbreak Study. https://sites.google.com/view/llm-jailbreak-study.
[2] Accessed on 10/08/2023. ChatGPT is a new AI chatbot that can answer questions and write essays. https://www.cnbc.com/2022/12/13/chatgpt-is-a-new-ai-chatbot-that-can-answer-questions-and-write-essays.html.
[3] Accessed on 10/08/2023. ChatGPT Masterclass: The Guide to AI & Prompt Engineering Udemy. https://www.udemy.com/course/chatgpt-ai-masterclass/.
[4] Accessed on 10/08/2023. gpt-4-system-card.pdf. https://cdn.openai.com/papers/gpt-4-system-card.pdf.
[5] Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 298–306.
[6] Badr AlKhamissi, Millicent Li, Asli Celikyilmaz, Mona Diab, and Marjan Ghazvininejad. 2022. A review on language models as knowledge bases. *arXiv preprint arXiv:2204.06031* (2022).
[7] Fares Antaki, Samir Touma, Daniel Milad, Jonathan El-Khoury, and Renaud Duval. 2023. Evaluating the performance of chatgpt in ophthalmology: An analysis of its successes and shortcomings. *Ophthalmology Science* (2023), 100324.
[8] Amos Azaria. 2022. ChatGPT Usage and Limitations. (Dec. 2022). https://hal.science/hal-03913837 working paper or preprint.
[9] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity. arXiv:2302.04023 [cs.CL]
[10] Ziv Bar-Yossef and Naama Kraus. 2011. Context-Sensitive Query Auto-Completion. In *Proceedings of the 20th International Conference on World Wide Web* (Hyderabad, India) *(WWW '11)*. Association for Computing Machinery, New York, NY, USA, 107–116. https://doi.org/10.1145/1963405.1963424

[11] Morteza Behrooz, William Ngan, Joshua Lane, Giuliano Morse, Benjamin Babcock, Kurt Shuster, Mojtaba Komeili, Moya Chen, Melanie Kambadur, Y-Lan Boureau, et al. 2023. The HCI Aspects of Public Deployment of Research Chatbots: A User Study, Design Recommendations, and Open Challenges. *arXiv preprint arXiv:2306.04765* (2023).

[12] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 610–623.

[13] Andrew Blair-Stanek, Nils Holzenberger, and Benjamin Van Durme. 2023. Can GPT-3 perform statutory reasoning? *arXiv preprint arXiv:2302.06100* (2023).

[14] Ali Borji. 2023. A categorical archive of chatgpt failures. *arXiv preprint arXiv:2302.03494* (2023).

[15] Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. 2022. What does it mean for a language model to preserve privacy?. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 2280–2292.

[16] Sarah Brown-Schmidt, Si On Yoon, and Rachel Anna Ryskin. 2015. People as contexts in conversation. In *Psychology of learning and motivation*. Vol. 62. Elsevier, 59–99.

[17] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712* (2023).

[18] Yihan Cao, Siyu Li, Yixin Liu, Zhiling Yan, Yutong Dai, Philip S Yu, and Lichao Sun. 2023. A comprehensive survey of ai-generated content (aigc): A history of generative ai from gan to chatgpt. *arXiv preprint arXiv:2303.04226* (2023).

[19] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. *Advances in neural information processing systems* 30 (2017).

[20] Hai Dang, Lukas Mecke, Florian Lehmann, Sven Goller, and Daniel Buschek. 2022. How to Prompt? Opportunities and Challenges of Zero-and Few-Shot Learning for Human-AI Interaction in Creative Applications of Generative Models. *arXiv preprint arXiv:2209.01390* (2022).

[21] Chenhe Dong, Yinghui Li, Haifan Gong, Miaoxin Chen, Junxin Li, Ying Shen, and Min Yang. 2022. A Survey of Natural Language Generation. *Comput. Surveys* 55, 8 (dec 2022), 1–38. https://doi.org/10.1145/3554727

[22] Dat Duong and Benjamin D Solomon. 2023. Analysis of large-language model versus human performance for genetics questions. *European Journal of Human Genetics* (2023), 1–3.

[23] Yogesh K Dwivedi, Nir Kshetri, Laurie Hughes, Emma Louise Slade, Anand Jeyaraj, Arpan Kumar Kar, Abdullah M Baabdullah, Alex Koohang, Vishnupriya Raghavan, Manju Ahuja, et al. 2023. "So what if ChatGPT wrote it?" Multidisciplinary perspectives on opportunities, challenges and implications of generative conversational AI for research, practice and policy. *International Journal of Information Management* 71 (2023), 102642.

[24] Yanai Elazar, Nora Kassner, Shauli Ravfogel, Abhilasha Ravichander, Eduard Hovy, Hinrich Sch"utze, and Yoav Goldberg. 2021. Measuring and improving consistency in pretrained language models. *Transactions of the Association for Computational Linguistics* 9 (2021), 1012–1031.

[25] Mohammadreza Farrokhnia, Seyyed Kazem Banihashem, Omid Noroozi, and Arjen Wals. 2023. A SWOT analysis of ChatGPT: Implications for educational practice and research. *Innovations in Education and Teaching International* (2023), 1–15.

[26] Luciano Floridi. 2023. AI as agency without intelligence: on ChatGPT, large language models, and other generative models. *Philosophy & Technology* 36, 1 (2023), 15.

[27] Vinitha Gadiraju, Shaun Kane, Sunipa Dev, Alex Taylor, Ding Wang, Emily Denton, and Robin Brewer. 2023. "" I wouldn't say offensive but...""": Disability-Centered Perspectives on Large Language Models. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 205–216.

[28] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. 2020. Realtoxicityprompts: Evaluating neural toxic degeneration in language models. *arXiv preprint arXiv:2009.11462* (2020).

[29] Sukhpal Singh Gill and Rupinder Kaur. 2023. ChatGPT: Vision and challenges. *Internet of Things and Cyber-Physical Systems* 3 (2023), 262–271.

[30] Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How Close is ChatGPT to Human Experts? Comparison Corpus, Evaluation, and Detection. arXiv:2301.07597 [cs.CL]

[31] Maanak Gupta, Charankumar Akiri, Kshitiz Aryal, Elisabeth Parker, and Lopamudra Praharaj. 2023. From ChatGPT to ThreatGPT: Impact of Generative AI in Cybersecurity and Privacy. *IEEE Access* 11 (2023), 80218–80245. https://api.semanticscholar.org/CorpusID:259316122

[32] Muhammad Usman Hadi, R Qureshi, A Shah, M Irfan, A Zafar, MB Shaikh, N Akhtar, J Wu, and S Mirjalili. 2023. A Survey on Large Language Models: Applications, Challenges, Limitations, and Practical Usage. *TechRxiv* (2023).

[33] Thomas M Holtgraves and Yoshihisa Kashima. 2008. Language, meaning, and social cognition. *Personality and Social Psychology Review* 12, 1 (2008), 73–94.

[34] Andreas Holzinger, Katharina Keiblinger, Petr Holub, Kurt Zatloukal, and Heimo M"uller. 2023. AI for life: Trends in artificial intelligence for biotechnology. *New Biotechnology* 74 (2023), 16–24.

[35] Jian Hu, Gang Wang, Fred Lochovsky, Jian-tao Sun, and Zheng Chen. 2009. Understanding User's Query Intent with Wikipedia. In *Proceedings of the 18th International Conference on World Wide Web* (Madrid, Spain) (*WWW '09*). Association for Computing Machinery, New York, NY, USA, 471–480. https://doi.org/10.1145/1526709.1526773

[36] Jie Huang and Kevin Chen-Chuan Chang. 2022. Towards reasoning in large language models: A survey. *arXiv preprint arXiv:2212.10403* (2022).

[37] Myeongjun Jang, Deuk Sin Kwon, and Thomas Lukasiewicz. 2022. BECEL: Benchmark for Consistency Evaluation of Language Models. In *International Conference on Computational Linguistics*. https://api.semanticscholar.org/CorpusID:252819451

[38] Myeongjun Jang and Thomas Lukasiewicz. 2023. Consistency analysis of chatgpt. *arXiv preprint arXiv:2303.06273* (2023).

[39] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *Comput. Surveys* 55, 12 (2023), 1–38.

[40] Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona Diab, and Bernhard Sch"olkopf. 2023. Can Large Language Models Infer Causation from Correlation? *arXiv preprint arXiv:2306.05836* (2023).

[41] Ishika Joshi, Ritvik Budhiraja, Harshal Dev, Jahnvi Kadia, M. Osama Ataullah, Sayan Mitra, Dhruv Kumar, and Harshal D. Akolekar. 2023. ChatGPT in the Classroom: An Analysis of Its Strengths and Weaknesses for Solving Undergraduate Computer Science Questions. https://api.semanticscholar.org/CorpusID:258417916

[42] Jean Kaddour, Joshua Harris, Maximilian Mozes, Herbie Bradley, Roberta Raileanu, and Robert McHardy. 2023. Challenges and Applications of Large Language Models. arXiv:2307.10169 [cs.CL]

[43] Enkelejda Kasneci, Kathrin Seßler, Stefan K"uchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan G"unnemann, Eyke H"ullermeier, et al. 2023. ChatGPT for good? On opportunities and challenges of large language models for education. *Learning and individual differences* 103 (2023), 102274.

[44] Rehan Ahmed Khan, Masood Jawaid, Aymen Rehan Khan, and Madiha Sajjad. 2023. ChatGPT-Reshaping medical education and clinical management. *Pakistan Journal of Medical Sciences* 39, 2 (2023), 605.

[45] Felipe C Kitamura. 2023. ChatGPT is shaping the future of medical writing but still requires human judgment. , e230171 pages.

[46] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems* 35 (2022), 22199–22213.

[47] Arun HS Kumar. 2023. Analysis of ChatGPT tool to assess the potential of its utility for academic writing in biomedical domain. *Biology, Engineering, Medicine and Science Reports* 9, 1 (2023), 24–30.

[48] Augustin Lecler, Loïc Duron, and Philippe Soyer. 2023. Revolutionizing radiology with GPT-based models: Current applications, future possibilities and limitations of ChatGPT. *Diagnostic and Interventional Imaging* 104, 6 (2023), 269–274.

[49] Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Evaluating the logical reasoning ability of chatgpt and gpt-4. *arXiv preprint arXiv:2304.03439* (2023).

[50] Yi Liu, Gelei Deng, Zhengzi Xu, Yuekang Li, Yaowen Zheng, Ying Zhang, Lida Zhao, Tianwei Zhang, and Yang Liu. 2023. Jailbreaking ChatGPT via Prompt Engineering: An Empirical Study. arXiv:2305.13860 [cs.SE]

[51] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. 2023. Trustworthy LLMs: a Survey and Guideline for Evaluating Large Language Models' Alignment. arXiv:2308.05374 [cs.AI]

[52] Li Lucy and David Bamman. 2021. Gender and representation bias in GPT-3 generated stories. In *Proceedings of the Third Workshop on Narrative Understanding*. 48–55.

[53] Ewa Luger and Abigail Sellen. 2016. Like Having a Really Bad PA" The Gulf between User Expectation and Experience of Conversational Agents. In *Proceedings of the 2016 CHI conference on human factors in computing systems*. 5286–5297.

[54] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2023. Self-refine: Iterative refinement with self-feedback. *arXiv preprint arXiv:2303.17651* (2023).

[55] Aman Madaan, Shuyan Zhou, Uri Alon, Yiming Yang, and Graham Neubig. 2022. Language models of code are few-shot commonsense learners. *arXiv preprint arXiv:2210.07128* (2022).

[56] Douglas L Mann. 2023. Artificial intelligence discusses the role of artificial intelligence in translational medicine: a JACC: basic to translational science interview with ChatGPT. *Basic to Translational Science* 8, 2 (2023), 221–223.

[57] Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2023. Recent advances in natural language processing via large pre-trained language models: A survey. *Comput. Surveys* 56, 2 (2023), 1–40.

[58] Chelsea Myers, Anushay Furqan, Jessica Nebolsky, Karina Caro, and Jichen Zhu. 2018. Patterns for how users overcome obstacles in voice user interfaces. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–7.

[59] Roberto Navigli, Simone Conia, and Bj̈orn Ross. [n. d.]. Biases in Large Language Models: Origins, Inventory and Discussion. *ACM Journal of Data and Information Quality* ([n. d.]).

[60] John J Nay. 2022. Law informs code: A legal informatics approach to aligning artificial intelligence with humans. *Nw. J. Tech. & Intell. Prop.* 20 (2022), 309.

[61] Saima Nisar and Muhammad Shahzad Aslam. 2023. Is ChatGPT a Good Tool for T&CM Students in Studying Pharmacology? *Available at SSRN 4324310* (2023).

[62] OpenAI. 2023. GPT-4 Technical Report. arXiv:2303.08774 [cs.CL]

[63] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems* 35 (2022), 27730–27744.

[64] Ethan Perez, Sam Ringer, Kamilė Lukošiūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Ben Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndousse, Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland, Nelson Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson, Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Jack Clark, Samuel R. Bowman, Amanda Askell, Roger Grosse, Danny Hernandez, Deep Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. 2022. Discovering Language Model Behaviors with Model-Written Evaluations. arXiv:2212.09251 [cs.CL]

[65] Martin Porcheron, Joel E Fischer, Stuart Reeves, and Sarah Sharples. 2018. Voice interfaces in everyday life. In *proceedings of the 2018 CHI conference on human factors in computing systems*. 1–12.

[66] Junaid Qadir. 2023. Engineering education in the era of ChatGPT: Promise and pitfalls of generative AI for education. In *2023 IEEE Global Engineering Education Conference (EDUCON)*. IEEE, 1–9.

[67] Jing Qian, Hong Wang, Zekun Li, Shiyang Li, and Xifeng Yan. 2022. Limitations of language models in arithmetic and symbolic induction. *arXiv preprint arXiv:2208.05051* (2022).

[68] Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a general-purpose natural language processing task solver? *arXiv preprint arXiv:2302.06476* (2023).

[69] Md. Mostafizer Rahman and Yutaka Watanobe. 2023. ChatGPT for Education and Research: Opportunities, Threats, and Strategies. *Applied Sciences* (2023). https://api.semanticscholar.org/CorpusID:258584102

[70] A Rao, J Kim, M Kamineni, M Pang, W Lie, and MD Succi. 2023. Evaluating ChatGPT as an adjunct for radiologic decision-making. medRxiv, 2023-02.

[71] Partha Pratim Ray. 2023. ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, limitations and future scope. *Internet of Things and Cyber-Physical Systems* (2023).

[72] Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–7.

[73] Soo Young Rieh and Hong (Iris) Xie. 2006. Analysis of multiple query reformulations on the web: The interactive information retrieval context. *Information Processing & Management* 42, 3 (2006), 751–768. https://doi.org/10.1016/j.ipm.2005.05.005

[74] Malik Sallam. 2023. ChatGPT utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In *Healthcare*, Vol. 11. MDPI, 887.

[75] Gaurav Sharma and Abhishek Thakur. 2023. ChatGPT in drug discovery. (2023).

[76] Chenglei Si, Dan Friedman, Nitish Joshi, Shi Feng, Danqi Chen, and He He. 2023. Measuring Inductive Biases of In-Context Learning with Underspecified Demonstrations. *arXiv preprint arXiv:2305.13299* (2023).

[77] Marita Skjuve, Ida Maria Haugstveit, Asbjørn Følstad, and Petter Brandtzaeg. 2019. Help! Is my chatbot falling into the uncanny valley? An empirical study of user experience in human–chatbot interaction. *Human Technology* 15, 1 (2019), 30–54.

[78] Anselm Strauss and Juliet Corbin. 1998. Basics of qualitative research techniques. (1998).

[79] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine* (2023), 1–11.

[80] H Holden Thorp. 2023. ChatGPT is fun, but not an author. , 313–313 pages.

[81] Shubo Tian, Qiao Jin, Lana Yeganova, Po-Ting Lai, Qingqing Zhu, Xiuying Chen, Yifan Yang, Qingyu Chen, Won Kim, Donald C. Comeau, Rezarta Islamaj, Aadit Kapoor, Xin Gao, and Zhiyong Lu. 2023. Opportunities and Challenges for ChatGPT and Large Language Models in Biomedicine and Health. arXiv:2306.10070 [cs.CY]

[82] Teun A Van Dijk. 2007. *Comments on context and conversation.* Citeseer.

[83] Shuai Wang, Harrisen Scells, Bevan Koopman, and Guido Zuccon. 2023. Can ChatGPT write a good boolean query for systematic review literature search? *arXiv preprint arXiv:2302.03495* (2023).

[84] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171* (2022).

[85] 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.

[86] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359* (2021).

[87] Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, et al. 2022. Taxonomy of risks posed by language models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 214–229.

[88] Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. *arXiv preprint arXiv:2302.11382* (2023).

[89] Robert Wolfe and Aylin Caliskan. 2022. American== white in multimodal language-and-image ai. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. 800–812.

[90] Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and Wei Cheng. 2023. Exploring the Limits of ChatGPT for Query or Aspect-based Text Summarization. arXiv:2302.08081 [cs.CL]

[91] Hongbin Ye, Tong Liu, Aijia Zhang, Wei Hua, and Weiqiang Jia. 2023. Cognitive Mirage: A Review of Hallucinations in Large Language Models. *ArXiv* abs/2309.06794 (2023). https://api.semanticscholar.org/CorpusID:261705916

[92] Yee Hui Yeo, Jamil S Samaan, Wee Han Ng, Peng-Sheng Ting, Hirsh Trivedi, Aarshi Vipani, Walid Ayoub, Ju Dong Yang, Omer Liran, Brennan Spiegel, et al. 2023. Assessing the performance of ChatGPT in answering questions regarding cirrhosis and hepatocellular carcinoma. *medRxiv* (2023), 2023–02.

[93] Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. 2023. How well do Large Language Models perform in Arithmetic tasks? *arXiv preprint arXiv:2304.02015* (2023).

[94] JD Zamfirescu-Pereira, Richmond Y Wong, Bjoern Hartmann, and Qian Yang. 2023. Why Johnny can't prompt: how non-AI experts try (and fail) to design LLM prompts. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–21.

[95] Chaoning Zhang, Chenshuang Zhang, Chenghao Li, Yu Qiao, Sheng Zheng, Sumit Kumar Dam, Mengchun Zhang, Jung Uk Kim, Seong Tae Kim, Jinwoo Choi, Gyeong-Moon Park, Sung-Ho Bae, Lik-Hang Lee, Pan Hui, In So Kweon, and Choong Seon Hong. 2023. One Small Step for Generative AI, One Giant Leap for AGI: A Complete Survey on ChatGPT in AIGC Era. arXiv:2304.06488 [cs.CY]

[96] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. 2022. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493* (2022).

[97] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A Survey of Large Language Models. arXiv:2303.18223 [cs.CL]

[98] Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang. 2023. Why Does Chat-GPT Fall Short in Providing Truthful Answers? https://api.semanticscholar.org/CorpusID:258865162

[99] Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto. 2023. Navigating the grey area: Expressions of overconfidence and uncertainty in language models. *arXiv preprint arXiv:2302.13439* (2023).

[100] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2022. Large language models are human-level prompt engineers. *arXiv preprint arXiv:2211.01910* (2022).

[101] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593* (2019).

| Category (7) | Code (19) | Literatures |
|---|---|---|
| Intent Understanding ($D_{intent}$) | C1. Response does not meet users' intent or instruction. | [14, 21, 42, 47, 70, 76] |
| | C2. Response is not aligned with the user's context. | [14, 15, 17, 21, 25, 30, 42, 44, 47, 70, 71, 74, 76, 86] |
| | C17. The tone or communication style is disappointing. | [8, 14, 30, 71, 74] |
| Content Depth and Originality ($D_{depth}$) | C3. Response is too general. | [14, 30, 47, 71, 75, 92, 97] |
| | C4. Response lacks originality. | [14, 17, 23, 30, 45, 48, 71, 74, 80] |
| | C5. Response lacks information. | [14, 17, 21, 25, 30, 43, 47, 48, 56, 61, 70, 71, 75, 92] |
| Information Accuracy ($D_{acc}$) | C6. The response contains incorrect information. | [6, 8, 14, 17, 18, 24, 26, 40, 43, 47, 51, 66, 74, 95, 99] [4, 22, 23, 25, 29, 32, 48, 57, 70, 71, 75, 80, 83, 86, 87] |
| | C7. Response is based on training data cut off at a certain date, and has limited access to newly created data. | [6, 21, 23, 25, 30, 42, 44, 74, 79, 92, 97] |
| | C8. Response is inconsistent. | [6, 17, 22, 24, 29, 34, 38, 49, 51, 83, 97, 99] |
| | C9. ChatGPT struggles with reasoning. | [6, 9, 14, 23, 25, 32, 40, 44, 51, 56, 68, 75, 95, 97] |
| | C10. (Hallucination) ChatGPT fabricates contents that conflict with the source content or cannot be verified from existing sources. | [4, 9, 14, 18, 29, 32, 39, 42, 51, 79, 81, 87, 97] |
| | C19. (Sycophancy) ChatGPT excessively conforms to the user. | [9, 30, 51, 64] |
| Transparency ($D_{trans}$) | C11. It's difficult to understand the reasons, criteria, logic, and evidence behind the responses. | [6, 9, 17, 23, 29, 32, 38, 51, 70, 71, 74, 75, 79, 83] |
| Refusing to Answer ($D_{refuse}$) | C12. ChatGPT avoids giving its own opinion by saying something similar to "As a language model, I am not capable …" | [14, 30] |
| | C13. ChatGPT avoids talking about difficult or controversial issues by saying something similar to "As a language model, I am not capable …" | [9, 30] |
| | C7. Response is based on training data cut off at a certain date, and has limited access to newly created data. | [6, 21, 23, 25, 30, 42, 44, 74, 79, 92, 97] |
| Content Ethics and Integrity ($D_{ethic}$) | C14. Response contains unlawful content | [34, 51] |
| | C15. Response contains unethical, harmful content. | [4, 12, 15, 28, 34, 51, 57, 74, 75, 79, 86, 86, 87] |
| | C16. Response contains biased content. | [5, 7, 14, 18, 25, 27, 29, 32, 43, 51, 52, 57, 59, 71, 74, 75, 79, 87, 89] |
| Response Format and Attitude ($D_{format}$) | C17. The tone or communication style is disappointing. | [8, 14, 30, 71, 74] |
| | C18. Response is overly detailed or too long | [17, 30, 70, 74, 90] |
| | C19. (Sycophancy) ChatGPT excessively conforms to the user. | [9, 30, 51, 64] |

Table 7: 7 category and corresponding 19 codes of user-side dissatisfaction from LLM Responses.

# A APPENDIX

## A.1 Systematic Literature Review Paper List

All paper lists corresponding user-side dissatisfaction codes are in Table 7.

## A.2 Data Filtering Criteria and Detailed Reason

### A.2.1 Conversation-Level Filtering. **Conversation older than 30 days.** We collected real-world experience data from individuals, which inherently consists of past data they have encountered. Therefore, in order to encourage respondents to recall these past experiences while responding to our data collection system, we restricted the chat dates to "previous 30 days" from the survey date. Although the survey included explicit instructions regarding this

matter, we identified four cases where participants reported chat dates older than 30 days, and excluded them.

***Conversation with a memory level of 3 or lower.*** Even if a conversation occurred within the previous 30 days, it was considered unreliable if the user had a low memory level regarding the conversation. Therefore, conversations where the user's memory level was rated 3 or lower on a 7-point scale were filtered out. This criterion led to the exclusion of five conversations.

***Conversation for fun or testing purposes.*** Our research focused on real-world experiences related to dissatisfaction encountered while using LLMs for practical purposes. Therefore, we do not delve into scenarios where users intentionally provoke dissatisfactory responses from LLMs, attempting to manipulate the model's behavior through techniques like jailbreaking [1, 50], using LLM solely for fun or testing. Despite the explicit instructions regarding this in the data collection system, seven conversations were identified as falling into this category and were filtered out.

***Conversation from versions other than GPT-3.5.*** Considering the significant differences in performance between GPT-3.5 and GPT-4 [4], we also considered the GPT version used in the conversation. Four conversations used GPT-4, while all others used GPT-3.5. To maintain data consistency, we filtered out the four conversations that used GPT-4.

*A.2.2 Response-Level Filtering.* **Dissatisfaction due to Chat-GPT's error messages** Dissatisfaction caused by ChatGPT responses being interrupted or encountering errors was not our research scope. Three responses fell under this category.

***Unconvincing dissatisfaction*** Seven cases were identified where it was challenging to understand why the user was dissatisfied when reviewing both the ChatGPT conversation and the user's dissatisfaction reasons.

***Mismatch between score and reason*** In one case, the effectiveness score for resolving dissatisfaction was 1 (indicating not effective), but the reason for that score was reported that the dissatisfaction was resolved by the prompt. This mismatch led to the exclusion of this case.

***No Correlation between selected dissatisfactions and subsequent prompts for resolving that dissatisfaction*** In five cases, we observed a lack of correlation between selected dissatisfaction categories and selected subsequent prompts to address such dissatisfaction. For example, it was the case a prompt that had nothing to do with the selected dissatisfaction was chosen to resolve the dissatisfaction.