

Exploring User Prompting Behavior in LLM Interactions

Maximilian Slapnik
 Maximilian.Slapnik@campus.lmu.de
 LMU Munich
 Munich, Germany

ABSTRACT

Artificial Intelligence (AI) plays an increasingly important role in the daily lives of millions of people. Large Language Models (LLMs) are one of the most prominent implementations of AI that are used not only by experts, but equally by ordinary users as well. LLMs can respond to any textual input (prompts) with human-like answers, leveraging the training data that was used to implement the model. Even though prompting LLMs seems very straightforward, the question arises if it is possible to streamline the interactions with said models in order to optimize outputs. We investigate user behavior in interactions with LLMs based on a randomized trial of 100 samples that are publicly available on the platforms ShareGPT and Midjourney. The goal of this analysis is the discovery of recurring patterns as well as the evaluation of human tendencies and biases when interacting with AI models in order to understand prevalent behaviors and explore optimization opportunities.

CCS CONCEPTS

• **Human-centered computing** → **Interaction design**; • **Information systems** → *Information retrieval*; • **Computing methodologies** → *Natural language processing*.

KEYWORDS

Large Language Models, user behavior, prompting, interaction patterns

1 INTRODUCTION

Artificial Intelligence (AI)-based tools continually gain prominence as regularly leveraged tools in the daily lives of millions of people. Today, the significance of this technology is reflected in the current AI market size that is estimated to be \$142 billion USD, and forecasted to increase more than tenfold by 2030 [18]. In addition to typical AI applications such as recommendation systems or autonomous agents, generative models are notably increasing in popularity as well, making it one of the central research topics in the field. One of the most widely used implementations of generative models are Large Language Models (LLMs), the most popular example at the moment being OpenAI’s ChatGPT [10]. Adoption rates of generative AI applications among professionals are increasing rapidly, and are already at around 30% [17].

Large Language Models are mainly implemented in the form of text generating chatbots that can answer seemingly any question

a user might pose. Although no expert knowledge is required to formulate a request and interact with an LLM-based bot, it is challenging to optimize the output, since it varies depending on the structure, wording, and composition of the input. Any form of natural language model input, whether it is in the form of a task or a question, is commonly referred to as “prompting” the model. Due to the vast application possibilities and promising future developments of LLMs, exploration of user prompting behavior in interactions with such models is of particular interest. Plenty of research has been conducted in the field of user interactions with LLMs already, mainly in regard to query reformulation strategies, studies of common user errors when prompting, different prompt composition strategies, and general LLM limitations.

In this paper, we are going to explain the fundamentals and workings of LLMs and prompting, describe related research in the realm of user - LLM exchange, and perform our own investigation of user behavior in such interactions. This investigation has the objective of facilitating comprehension of existing challenges users face when dealing with Large Language Models. Furthermore, readers will gain a better understanding of the design of effective prompts that enhance model output.

Since the main part of this paper will be complemented by an analysis of real-world examples, the reader can expect to develop an enhanced comprehension of actual user prompting behavior. To obtain these insights, we will leverage input data mainly gathered from the website ShareGPT [15], which enables users to store conversations they have had with the ChatGPT model for later retrieval or sharing them publicly.

The paper is organized as follows. This introduction is succeeded by a related work section that sets the context for all subsequent parts by first focusing on Large Language Models (LLMs) and covering general information about their workings, training data, text generation capabilities, real-world usage, and current limitations. We then explore user interactions with LLMs, explain the concept of prompting, and highlight various use cases as well as related research.

The next section introduces the study by outlining the research objective and describing the methodology and individual steps that will be taken. It then focuses on the research method we use, as well as the ShareGPT and Midjourney platforms, which provide the input data for the study.

Subsequently, we present our findings. To do so, we first list the study results, organized into predefined categories. We then analyze observable trends in user behavior and data patterns.

The following discussion section starts with a synthesis of our observations. We then go into more detail about the reasons why users interact with LLMs the way they do, offering reasoning and informed assumptions. Additionally, we explore possibilities for prompt improvements based on findings from related research.

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution.

© 2023 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

In the outlook section, we provide a perspective on future developments, divided into an introduction of the concept of Auto-GPT as a possible future iteration of prompting, as well as an overview of prompt engineering as a newly emerging discipline in the technology sector.

The final section of this paper offers a consolidation of the findings and associated discussions, as well as a summary of how we could recognize findings from related research in our own input samples.

2 BACKGROUND AND RELATED WORK

This section provides an overview over fundamental concepts that are necessary to facilitate comprehensibility of all latter parts. We start with explaining Large Language Models in general, followed by more specific explanations of distinct prompting approaches. These explanations are followed by a more general overview of related work on the subject.

2.1 Large Language Models (LLMs)

One of the most widely used applications areas of generative AI are Large Language Models (LLMs). Among LLMs, the most widely adopted is ChatGPT [10], which is a conversational model being developed by OpenAI. The model is currently publicly accessible and free of charge.

LLMs can be leveraged for a variety of tasks, but their main focus area is Natural Language Processing (NLP). Therefore, most LLMs designed for end users are implemented in the form of chatbots, as is the case with ChatGPT for example. They typically consist of an interface comprised of an input field for the user to type in arbitrary text, as well as an output section that displays generated responses of the model.

Large Language Models are a recent advancement that followed the development of the original transformer architecture, which is a deep learning approach first introduced by researchers in 2017 [21]. In future iterations, the Generative Pre-Training (GPT) [11] approach was adapted for text-based models in particular, laying the foundation for today’s most popular conversational LLMs, such as ChatGPT. Since our research revolves around user interaction with dialog-focused models, we will not go into more detail about other application and development areas of LLMs from hereon. In addition to ChatGPT, there are also a variety of similar other models focused on text generation, such as LaMDA [19], Sparrow [5], or BlenderBot 3 [16].

The primary method of interacting with LLMs is the previously mentioned approach of “prompting” the model. Currently, there are no binding guidelines on the concrete structure of prompts, and users are therefore completely unconfined in their way of interacting. In the following, we go into more detail about different types of prompts, namely zero- and few-shot.

2.1.1 Zero-Shot Learning. It is important to differentiate possible forms of prompting that exist in LLM interactions in order to be able to correctly analyze and assess prompts in later parts of this paper. To do so, we will clarify the terms “Zero-Shot-” as well as “Few-Shot Learning”. Zero-Shot Learning has been described as “[...] learning how to recognise new concepts by just having a description of them” [12, p. 1].

Table 1: ChatGPT: Odd Number Sum Evaluations

Prompt:	The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.
A:	<i>The odd numbers in this group add up to an even number: 15, 5, 13, 7, and 1. Their sum is 41, which is indeed an even number.</i>
Prompt:	The odd numbers in this group add up to an even number: 4, 8, 9, 15, 12, 2, 1. The answer is False. The odd numbers in this group add up to an even number: 17, 10, 19, 4, 8, 12, 24. The answer is True. The odd numbers in this group add up to an even number: 16, 11, 14, 4, 8, 13, 24. The answer is True. The odd numbers in this group add up to an even number: 17, 9, 10, 12, 13, 4, 2. The answer is False. The odd numbers in this group add up to an even number: 15, 32, 5, 13, 82, 7, 1.
A:	<i>The answer is False. The odd numbers in this group (15, 5, 13, 7, and 1) add up to 41, which is an odd number, not an even number.</i>

In regard to prompting, Zero-Shot refers to any form of prompt where the user does not specify a particular example as part of the query [3, p. 1]. An example for such a query could be “Give me five examples of papers that deal with the topic of LLMs”. As expected, model performance (accuracy) was measured as inferior in zero-shot prompting settings compared to other settings where more information is provided in the prompt, such as a Few-Shot approach [1, p. 5].

2.1.2 Few-Shot Learning. In Few-Shot Learning in contrast, examples are given, albeit they do not necessarily have to be many. It is generally aimed at ensuring good learning performance with only a few (most often less than 20) supervised examples [14, p. 1]. In a Few-Shot prompting setting, the user accordingly does provide increased guidance for the model. Doing so can help improve model outputs [3, p. 1]. Solely giving one demonstration in the prompt is commonly referred to as a “One-Shot” approach [1, p. 6], whereas multiple demonstrations make it “Few-Shot”. To exemplify the differences in learning capabilities depending on the prompting setting, we present an example inspired by the approach of Brown et al. [1] and recent online resources [2] in Table 1. The table displays an interaction we have carried out with the current ChatGPT-3 model, where we instructed the model to classify the sum of a sequence of numbers as either even or odd.

Using this example prompt, we could recreate the effects observed by other researchers. We can witness that the user request is indeed first wrongly answered by the model using a zero-shot approach. Only when relying on few-shot prompting and supplying ChatGPT-3 with additional information, the model generates the correct result.

We should note however, that ChatGPT is a dedicated language model, and we have used ChatGPT-3. It is therefore expected that math-related problems and prompts will be less accurately

answered compared to purely language-based requests. Depending on the model training and technological progress, it is imaginable that future iterations, such as ChatGPT-4, will presumably improve performance and may therefore not have to rely on few-shot assistance as much.

2.2 User Interaction with LLMs

The majority of existing related work focuses on strategies for designing prompts, challenges users may face in the course of doing so, as well as investigations of user behavior when performing search queries in general. These topics are of particular relevance for our work since we are not only interested in ways to improve prompt outputs, but also in end user behavior while prompting.

2.2.1 Effectiveness of Query Reformulation. Interest in query optimization in order to improve outputs is widely prevalent ever since the broad availability of large search engines, such as Google [6] or Bing [8]. Results of search query reformulations can be an indication to also look at prompt reformulation and its effectiveness in more detail. Huang and Efthimiadis [7] investigated user behavior when reformulating search queries, i.e. modifying a previous query in order to improve results, in as early as 2009. Even though their focus was on search queries only, it was significant that there were improvements that could be seen after reformulating queries, concluding that "most reformulation strategies result in some benefit to the user" [7, p. 1].

2.2.2 Few-Shot Learning Capabilities. As indicated in the previous section, the prompting setting of an LLM can significantly influence its performance in regard to the accuracy of outputs. In related research, Brown et al. [1] for example have explored differences in model performance depending on examples provided as part of the LLM prompt [1]. To do so, they have conducted various experiments directly comparing zero-, one-, and few-shot learning. In conclusion, the trials revealed that LLMs indeed show remarkable performance in various areas, ranging from translation and question answering to reasoning tasks, particularly when relying on few-shot prompts. Notably, those few-shot learning models are able to adapt to new tasks with little training data only. It becomes apparent that it is important to provide adequate, fitting examples as part of a prompt in order to enhance model outcomes.

2.2.3 Challenges in Non-expert Prompt Design. Designing effective prompts that are optimized to achieve the desired output is challenging for most users. Observations in user trials showed that non-experts approach prompting generally rather opportunistic than strategic [22] and therefore have trouble adequately communicating their requests to the model. The trials revealed two main reasons for these problems:

- (1) First of all, non-experts suffered from over-generalization based on past, limited experiences, or single observations of success and failure of prompt adjustments, which may not be universally transferable. This tendency could be identified in the trials when "participants often stopped iterating once the [desired model] behavior was observed in a single conversational context, without considering other conversations or contexts—or gave up too early if the behavior was not observed" [22, p. 10].

- (2) Second, a false comprehension of AI systems as human-like conversational partners. The participants expected AI models to behave as in a human-to-human interaction, therefore presuming social understanding. They avoided consequently following general best practices in prompt design such as providing examples, and instead relied on instinct and literal commands, not realizing that prompts solely bias the model towards one direction. In essence, users have to realize that they should not rely on their own understanding of language, but instead consider how an LLM interprets the prompt. The example Zamfirescu et al. give in their research is the misconception of many users that simply prompting a model not to say "XYZ" will not actually prevent the model from doing so in any case. It is still possible that the model will output "XYZ" in subsequent responses due to the probabilistic nature of LLMs and the fact that even though responses can be primed by earlier prompts in the conversation, they are generated ad-hoc based on the most relevant training data of the model.

2.2.4 Opportunities and Challenges for Interactive Prompt Design Applications. A final aspect that highlights user difficulties in optimal prompt creation is the existing research aspiring to provide assisting tools that facilitate the whole design process end-to-end. As such, researchers have proposed guidelines for prompt design applications that assist users in prompt engineering by providing an adequate interface [3]. The motivation stems from the same realization that has already been discussed above: overall, LLM users struggle to communicate their expectations and intentions effectively to the model. Instead, they commonly approach prompting as a trial and error process. Currently, there is no single design interface or definite guidelines for prompts, which would assist non-technical users in particular. Example propositions to ease design include detecting keywords in prompt formulation and allowing edits of those via a dropdown menu, providing basic prompt-building blocks from which a user can select, offering the possibility to combine multiple prompts, and storing selected prompts in a toolbar for quick and repeated execution.

In summary, the volume of related research in regard to user prompting behavior in LLM interactions shows the significance of the topic for effective use of this technology. Few definitive guidelines exist, and users are largely left on their own when formulating prompts, especially non-technical users which may lack deep knowledge on the subject. Although superior prompting strategies have been discovered, users do not consistently implement them, even when specifically instructed to do so. So far, these behaviors have been attributed to existing habits and missing perception of LLMs as artificial and not human.

3 STUDY ON USAGE PATTERNS OF LLM USERS

3.1 Research Objective

The main outlined goal of our research is to gain a fundamental understanding of user behavior in conversation with Large Language Models. This analysis includes identifying common

patterns and strategies in those interactions. Previous research indicates that users regularly face challenges and difficulties, especially when trying to formulate effective prompts. Through accumulation and analysis of qualified data samples we aim to identify and understand these challenges, as well as investigate the impact and effect of user behavior on the effectiveness of LLM responses. Given the various kinds of available generative models, we want to examine differences in prompting behavior according to model type as well.

Since related insights suggest that reformulating search queries is a popular strategy to improve results, we want to investigate if users apply this strategy in LLM conversations as well. Furthermore, we want to assess the extent to which users show awareness of effective prompt formulation strategies, such as few-shot learning, and whether they rely on appropriate language that is machine and not human directed, thus showing comprehension of the fundamental difference between talking to a machine versus a human.

3.2 Research Method: ShareGPT and Midjourney

In order to obtain credible insights, we complement existing findings with real-world data. Our study analyzes data samples from two different types of LLMs. First, we examine user interactions with ChatGPT, a generative NLP model, which has already been described in more detail in Section 2.1. These ChatGPT conversations were obtained from the website ShareGPT[15]. ShareGPT is an open platform, that allows its community to publicly share interactions they have had with the ChatGPT model. As of today, ShareGPT has accumulated nearly 300.000 saved user conversations. What makes ShareGPT particularly suitable for our use-case is the fact that the entire shared conversation can be viewed by the observer as if they had personally conducted the interaction, allowing us to gain a deeper understanding of the conversation dynamics and outcome.

Midjourney in contrast, is a platform that focuses on AI-based image generation. Users can interact with the model through Discord[] and submit individual requests. In order to generate an illustration, users have to enter a descriptive prompt, similar to ChatGPT. The description typically includes everything that should appear in the picture, but may also encompass the desired mood, drawing style, or composition of the image being generated. Notably, the Midjourney Bot does not understand grammar, sentence structure, or specific words like humans []. Midjourney’s developers actively encourage using fewer, but more precise and impactful words when prompting the model. For example, they suggest using "gigantic" instead of "big" in order to achieve better results. This recommendation stems from the fact that fewer words in a prompt intensify the influence each individual word has on the final outcome. However, it is important to mention that users have to strike a balance. An adequate amount of precise words is mandatory, because anything that is not specified may be randomized. In addition to purely textual prompts, the platform allows image inputs as well. Users may provide an image as a guideline or basis including instructions about things to modify, add, remove, or

remodel. The Midjourney platform on Discord has experienced rapid growth, and counts more than 17 million members as of today.

In order to verify observations and findings we have presented in Section 2, we examine exemplary real-world user interaction samples in the following. By randomly choosing 50 conversations from each ShareGPT’s website as well as Midjourney’s Discord channel, we obtain a representative sample of average user behavior in both text and image targeting prompts. We have defined dedicated categories according to which each sample is classified for both the ChatGPT and Midjourney conversations. For the language-focused ChatGPT conversations, the categories and specific sub-categories can be seen in Table 2. Similarly, the categories and sub-categories for image-focused Midjourney prompts are listed in Table 3.

For ShareGPT, we first of all classified the prompt by type, as in theory any kind of prompt is possible, because users are solely constrained to natural language in any shape or form. In the next category, it was of major interest what the user intended to achieve with their individual prompts, i.e. what they use the model for. The categories prompt length and setting refer to the number of sentences in the user inputs, and whether any examples were provided as part of the query. Engagement refers to the amount of exchange in the interaction. If the user prompted the model multiple times(at least twice) during the course of the conversation, we considered the interaction multi turn, otherwise single turn. The prompt’s complexity gave us an idea whether users leverage ChatGPT for simple tasks, that they may otherwise quickly research using a search engine, or if they pose complex questions that require expert-level knowledge. The refinement degree of the prompt revealed if users were generally content with the initial answer of the LLM, or if further elaboration was needed. Finally, we differentiated use of formal and informal language. In general, when a single conversation consisted of multiple prompts, we labeled it based on the most frequently observed category or significant behavior.

We classified Midjourney interactions using a similar approach. First of all, we differentiated between image- and purely language-based inputs. We then considered the length of the prompt, and whether it consisted solely of keywords, one or more sentences, or a mix of both. Next, we recorded the complexity of the whole prompt. The Midjourney bot always generates four versions of the desired image. It then allows users to either recreate variations or more detailed versions of one or more of those four results. Users can also re-execute the whole generation process. We thus classified the prompt accordingly in the refinement category. Finally, we recorded the clarity of the prompt and satisfaction levels based on the observed user behavior. If a user created variations or more detailed versions of the result, we assumed they were generally satisfied. Analogously, we assumed dissatisfaction if they regenerated the whole image.

4 STUDY RESULTS

4.1 Findings and Observable Trends

In the following section, we are going to break down the sample analysis results by category. The distribution of the sub-categories

Table 2: ShareGPT Prompt Analysis Categories

Type	Intent	Length	Setting
Question	Information	Long (5 + Sent.)	Zero Shot
Statement	Advice	Med (2 - 4 Sent.)	One Shot
Command	Clarification	Short (\leq 1 Sent.)	Few Shot
Task-based	Opinion Suggestion Entertainment		
Engagement	Complexity	Refinement	Language
Single Turn	Simple	None	Formal
Multi Turn	Intermediate	Once Complex	Informal

Table 3: Midjourney Prompt Analysis Categories

Type	Length	Composition	Complexity
Language	Long (12+ Words)	Keyword-Only	Simple
Image	Med (4-12 Words)	Sentence	Intermediate
	Short (1-3 Words)	Mix	Complex
Refinement	Language	Clarity	Satisfaction
None	Formal	Clear	Satisfied
Variation	Informal	Ambiguous	Dissatisfied
Regeneration			Unclear

in the individual categories Prompt Type, Prompt Intent, Prompt Length, and Prompt Setting can be seen in Figure 1, whereas Figure 2 shows the numbers for the categories Engagement, Complexity, Refinement, and Language. In regard to prompt type, it became clear that the majority of users (54%) use Chat-GPT for task-based prompts, followed by questions (36%), commands (8%), and statements (2%). The most observed intents behind prompts were information gain (42%) and asking for suggestions (34%), followed by entertainment (10%) and advice (10%). Only few users were asking for clarification on a subject matter (4%). Interestingly, we did not observe any prompts where users actively asked the chatbot for its opinion (0%), which we initially had estimated as an at least fairly common use case. Prompt length was very evenly distributed, and we could not make out a clear preference of users. Short (38%), medium length (34%), and long (28%) prompts made up about a third of our samples each. We could clearly see the most often used prompting setting, however. The vast majority of users relied on a zero shot approach (90%), whereas only 8% used a one shot, and a mere 2% a few shot setting. Engagement in interactions was evenly distributed between multi turn (52%) and single turn (48%) conversations, meaning that almost half of the observed chats ended after the initial answer of the LLM. Most prompts were of a simple nature (48%), and slightly more than a third (38%) could be classified as intermediate, which left only 14% as complex. ShareGPT users only rarely refined their prompts multiple times (12%) or once (22%), leaving a two thirds majority (66%) of never refined queries. Finally,

we could observe a tendency towards formal language (60%), which was used more often than informal language (40%).

Similarly, we analyzed the Midjourney data samples according to the predefined categories. The corresponding data and distribution for the categories Prompt Type, Length, Composition, and Complexity can be seen in Figure 3, for the categories Refinement, Language, Clarity, and Satisfaction in Figure 4. We already explained that Midjourney allows users to also prompt with an existing image as part of the input. However, only very few (8%) users have made use of this feature in our data sample. The majority (92%) relied on purely textual prompts. Interestingly, most queries were at least of medium length (54%), or even long (38%), and only 8% were classified as short. The composition of the individual prompts was well-balanced between sentences (42%), only keywords (26%), or a mix of both (32%). The same applies for the complexity. Most prompts were on an intermediate level (40%), closely followed by simple (38%) and finally complex prompts with a share of only 22%. Similarly to our ChatGPT samples, we observed only few refinements of Midjourney prompts. More detailed regenerations of images made up 14%, variations 10%, and the rest (76%) was not refined at all. A clear distribution could be seen in regard to formality of language. Users relied on formal language in almost all cases (94%), and only very rarely on more informally phrased prompts (6%). We identified 92% of all prompts as clear in their intention, which left only 8% as ambiguous. Satisfaction of users was unfortunately often unclear (58%), due to no apparent reactions of the users to the final image. However, for almost half of the samples we could either identify signs of satisfaction (26%) or dissatisfaction (16%).

5 DISCUSSION

5.1 Data Synthesis: Commonalities, Differences, and Possible Explanations

In this section, we reason about the observed behavior from the two data sources and try to offer informed assumptions on potential causes.

5.1.1 ChatGPT Behavior. In regard to the type of prompt, we observed mainly task-based and question-related queries from users. This leads us to imagine that some users treat LLMs such as ChatGPT increasingly as their personal online assistant when it comes to executing various tasks the model might be able to solve. Leveraging AI bots as assistants of the future is a use case that gains popularity, and that is also increasingly researched [4]. The second most prevalent type of prompt were questions. We have already mentioned, that NLP models might replace current search engines in the future. This belief is further confirmed by observations from other researches who already see this emerging trend in user behavior [20]. Regarding prompt intent, looking for information was the most popular use case. Seemingly, users see ChatGPT as a reliable source of knowledge and trust its abilities. This assumption can be dangerous however, as it is well-known that every LLM is probabilistic and only as good as its training data. Results should therefore always be verified. Based on the fact that a lot of users also relied on ChatGPT for suggestions, we assume that it is gladly used as a means to get ahead when you hit roadblocks, or need support in

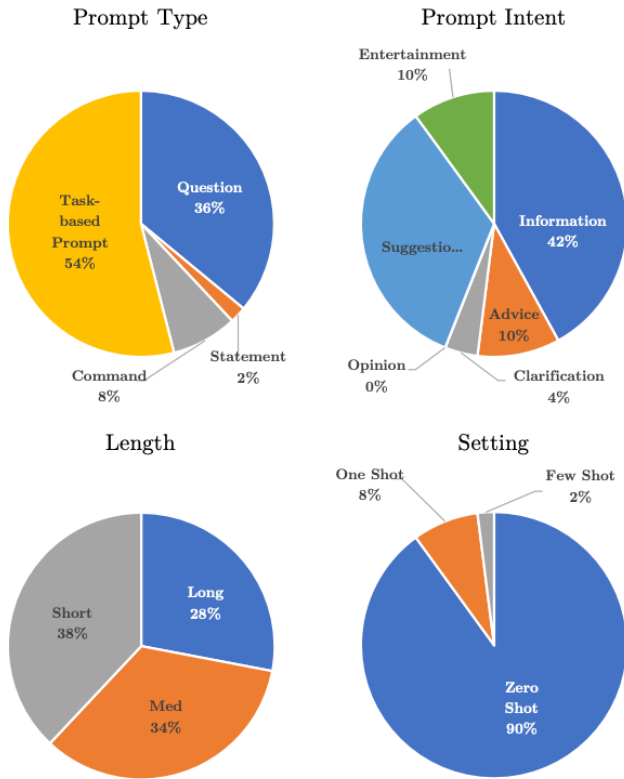


Figure 1: ShareGPT Prompt Analysis Categories 1-4

endeavours that require creativity. As we already touched on above, we did not have any requests for opinions of the model. We would have thought that opinion related requests on difficult, morally complex, or controversial topics would be more popular, since such behavior could be observed a lot online when users tried to test the limits of the model, or for example when researchers tested political biases [13]. Regarding the prompt setting, we made an observation similar to Brown et al. [1]. Users do not leverage effective prompting techniques such as a few-shot approach. Instead, they relied heavily (90%) on zero shot prompts. We attribute this behavior to missing awareness of users about optimized techniques, and therefore recommend that LLM providers actively make users actively inform users, e.g. by providing examples, or releasing guidelines. Our observations when looking at user engagement makes us think that current LLMs are already quite accurate: almost half of the users finished the conversation after the initial answer of the LLM (single turn). This is further reinforced by 66% of prompts that were not refined. We therefore suspect, that most of the time the users were actually content with the results. However, we have to mention the possibility that the initial answer was so far off, that users might have stopped trying after the first attempt. Overall, we reckon that users prefer to leverage LLMs for rather simple tasks. It is difficult to say, whether users do not trust LLMs enough to throw complex questions at them yet, or if it simply is in the nature of online search requests, that the majority of search requests that users look up are more of a simple nature than very complex. Few users refine their

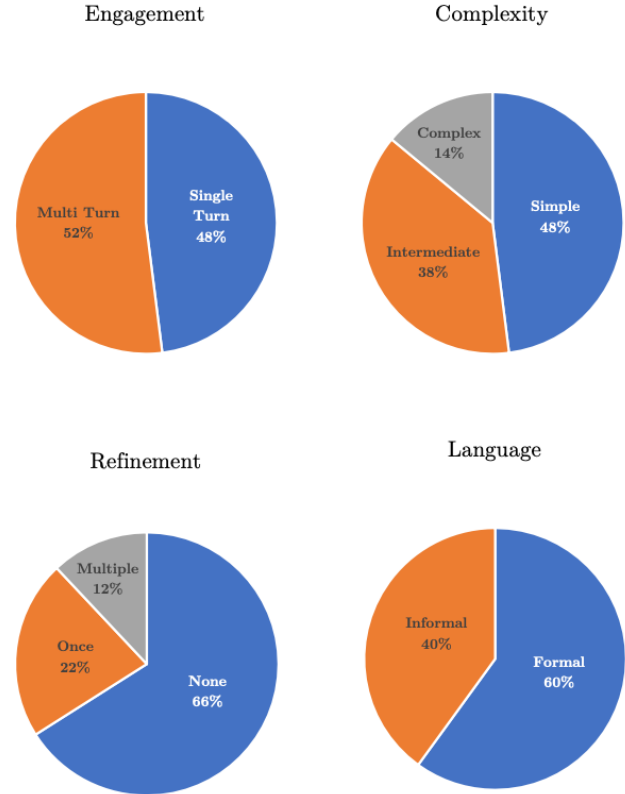


Figure 2: ShareGPT Prompt Analysis Categories 5-8

prompts (66% do not). Related work suggests that refining queries has led to improved results when using search engines. LLMs can be improved by refinements too, since models are “primed” by all previous prompts in the interaction, so even if a model does not initially do what you want it to, it might make sense to give it more information or context, and try again. Due to missing sentiments and feedback of users in the sample interactions, we cannot reliably say if users were simply always content because they did not refine, or if they did not know that continuing the interaction with the model could have led to better results. Finally, we have generally observed a majority of formal, generally polite, and acceptable language. We reason that the use of such language ties in with our first observation, and reckon users may see the bot as a personal companion, that you therefore treat well in interactions.

5.1.2 Midjourney Behavior. For Midjourney, we identified a majority of language inputs. We assume, this tendency exists because users prefer to create something new instead of reworking existing images. We imagine AI could eventually revolutionise the image editing market as well, but does not seem to be quite there yet, probably due to missing accuracy when trying to make only very small adjustments. Regarding prompt length, the official Midjourney docs state: “The Midjourney Bot works best with simple, short sentences that describe what you want to see” [9]. Long sentences should be avoided, but we have seen users ignore this advice multiple times. There sometimes seems to be a lack

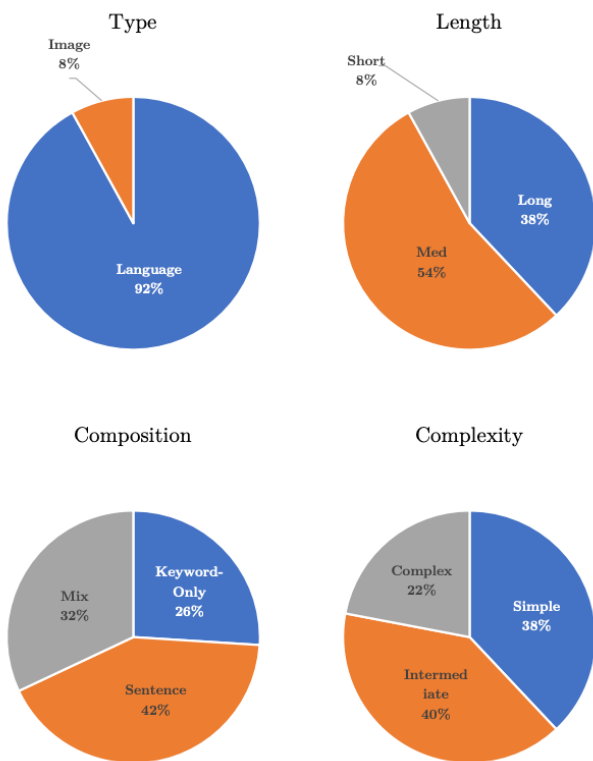


Figure 3: Midjourney Prompt Analysis Categories 1-4

of understanding that more information is not always better for quality of model outputs. This observation fits another that has been made before: Users struggle to formulate precise, effective, and therefore also short, concise prompts. We suggested better learning materials and guidance already as a measure to address this issue. Regarding composition of the Midjourney prompts, we suspect that this struggle is also the reason for the rare occurrence of a keywords-only prompts. Users seem to prefer natural language sentences or at least a mix of sentences and keywords instead of "encoding" their wants. To help, developers could offer reformulation engines, that only extract keywords (or create them) from input sentences. Since almost no images were regenerated as a reaction to the initial result of a prompt we had classified as complex, we conclude that Midjourney is able to handle difficult queries well. This assumption is confirmed by the fact that users did generally not refine their queries a lot after the initial result was shown by the engine. It was difficult to judge formality of language of users, but in general users tended to rely on formal language, which was probably also influenced by the keyword-focused nature of the prompts. We attribute the fact that most queries were formulated clearly to this circumstance as well, since keyword focused prompts are usually shorter and less complicated than NLP prompts.

5.1.3 Commonalities and Differences between ChatGPT and Midjourney. A clear distinction between ChatGPT and Midjourney prompts lies in the nature of the model and their designs. Whereas ChatGPT is trained on full sentence queries, Midjourney is keyword

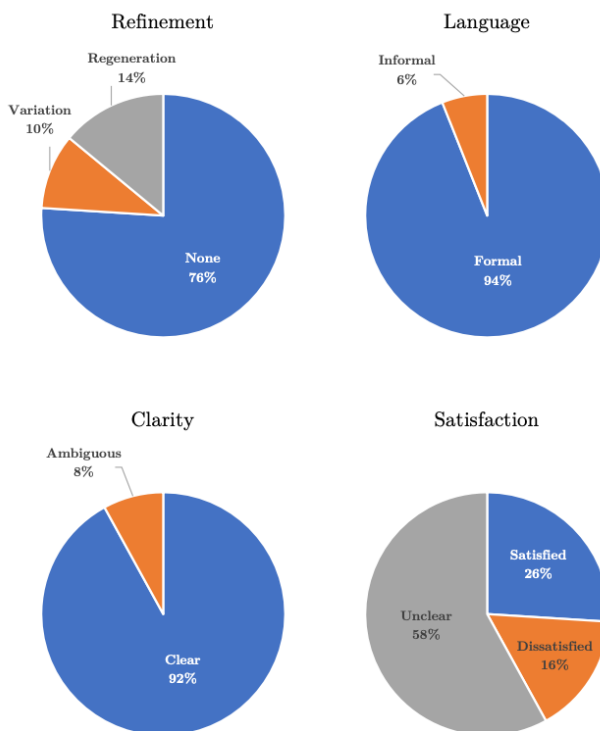


Figure 4: Midjourney Prompt Analysis Categories 5-8

focused. However, we have remarked already that Midjourney users tend to at least partially integrate sentence structures in their prompts as well, which may be due to the fact that it seems more natural and is easier for users due to existing habits. Both models were prompted with language on the more formal side. It is possible, that this perception is skewed though, as prompts with informal language are either not shared in the case of ChatGPT, or are directly filtered before execution in the case of Midjourney. It is worth mentioning that we were able to observe a previously mentioned general misconception that researchers had observed in conversations of users with natural language LLMs such as ChatGPT already. Some users rely on negation when providing instructions and misunderstand that it will not prevent the model from producing the unwanted. Our data samples from Midjourney contain one example where the user provides an image of the fictional character Voldemort from the movie saga Harry Potter, and explicitly asks the model to generate an image of "voldemort dying without a nose". Since the sole existence of the word "nose" in the prompt primes the model towards including said object, all four result images indeed contain visualizations of the fictional character voldemort passing—but with a nose. Remarkably, there might have been a higher chance that voldemort does not possess a human looking nose in the final image, if the user did not explicitly mention the word, as Voldemort's nose is actually distinctly different from that of a human according to his description in the books and movies.

6 CONCLUSION

Throughout this paper we have explored user behavior in interactions with LLMs across multiple dimensions. The goal was to find prevalent human tendencies, understand existing habits, and identify recurrent patterns.

To do so, we first explained the ever-growing importance of the subject, given the increasing use of generative AI across all domains. We then laid out fundamental concepts and explained existing findings from related research. These concepts and findings built the foundation for a synthesis of our own real-world data analysis with existing research in order to verify findings. To obtain a comprehensive overview, we explained different kinds of LLMs and their characteristics, highlighting ChatGPT and Midjourney in particular. An extensive study of data samples that encompassed a categorization of each individual entry according to eight predefined categories helped us to understand the current state of human - LLM interaction. We could observe that many users are still subject to biases and misunderstandings that make effective prompting difficult. During the discussion of our results, we reasoned why certain user dynamics exists, and which actions LLM providers and developers as well as users themselves could undertake in order to address prevalent issues.

Deeper exploration of this topic with a larger data sample is needed in order to verify findings and gain a deeper understanding of existing interaction dynamics and is subject to further research. Overall, it is to be acknowledged that user prompting behavior is one of the most relevant research topics in human centered computing and will gain even more importance in the near future.

REFERENCES

- [1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. <http://arxiv.org/abs/2005.14165> arXiv:2005.14165 [cs].
- [2] DAIR.AI. 2023. Few-Shot Prompting | Prompt Engineering Guide. <https://www.promptingguide.ai/techniques/fewshot>
- [3] 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. <http://arxiv.org/abs/2209.01390> arXiv:2209.01390 [cs].
- [4] Mahshid Eshghie and Mojtaba Eshghie. 2023. ChatGPT as a Therapist Assistant: A Suitability Study. <http://arxiv.org/abs/2304.09873> arXiv:2304.09873 [cs].
- [5] Amelia Glaese, Nat McAleese, Maja Trębacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, Soňa Mokrá, Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. Improving alignment of dialogue agents via targeted human judgements. <http://arxiv.org/abs/2209.14375> arXiv:2209.14375 [cs].
- [6] Google. 2023. Google. <https://www.google.com/>
- [7] Jeff Huang and Efthimis N. Efthimiadis. 2009. Analyzing and evaluating query reformulation strategies in web search logs. In *Proceedings of the 18th ACM conference on Information and knowledge management*. ACM, Hong Kong China, 77–86. <https://doi.org/10.1145/1645953.1645966>
- [8] Microsoft. 2023. Bing. <https://www.bing.com/?cc=de>
- [9] Midjourney. 2023. Documentation and User Guide. <https://docs.midjourney.com/>
- [10] OpenAI. 2023. ChatGPT. <https://chat.openai.com/auth/login>
- [11] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving Language Understanding by Generative Pre-Training. (2018).
- [12] Bernardino Romera-Paredes and Philip H. S. Torr. 2015. An Embarrassingly Simple Approach to Zero-Shot Learning. https://doi.org/10.1007/978-3-319-50077-5_2 Series Title: Advances in Computer Vision and Pattern Recognition.
- [13] David Rozado. 2023. The Political Biases of ChatGPT. *Social Sciences* 12, 3 (March 2023), 148. <https://doi.org/10.3390/socsci12030148>
- [14] Mesay Samuel, Lars Schmidt-Thieme, D. P. Sharma, Abiot Sinamo, and Abey Bruck. 2022. Offline Handwritten Amharic Character Recognition Using Few-shot Learning. <http://arxiv.org/abs/2210.00275> arXiv:2210.00275 [cs].
- [15] ShareGPT. 2023. ShareGPT: Share your wildest ChatGPT conversations with one click. <https://sharegpt.com/>
- [16] Kurt Shuster, Jing Xu, Mojtaba Komeili, Da Ju, Eric Michael Smith, Stephen Roller, Megan Ung, Moya Chen, Kushal Arora, Joshua Lane, Morteza Behrooz, William Ngan, Spencer Poff, Naman Goyal, Arthur Szlam, Y.-Lan Boureau, Melanie Kambadur, and Jason Weston. 2022. BlenderBot 3: a deployed conversational agent that continually learns to responsibly engage. <http://arxiv.org/abs/2208.03188> arXiv:2208.03188 [cs].
- [17] Statista. 2022. U.S.: generative AI adoption rate in the workplace by generation 2023. <https://www.statista.com/statistics/1361174/generative-ai-adoption-rate-at-work-by-generation-us/>
- [18] Statista. 2023. Artificial Intelligence market size 2030. <https://www.statista.com/statistics/1365145/artificial-intelligence-market-size/>
- [19] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agüera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. LaMDA: Language Models for Dialog Applications. <http://arxiv.org/abs/2201.08239> arXiv:2201.08239 [cs].
- [20] Liesbet Van Bulck and Philip Moons. 2023. What if your patient switches from Dr. Google to Dr. ChatGPT? A vignette-based survey of the trustworthiness, value, and danger of ChatGPT-generated responses to health questions. *European Journal of Cardiovascular Nursing* (April 2023). <https://doi.org/10.1093/eurcn/zvad038>

- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. 30 (2017).
- [22] J.D. 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*. ACM, Hamburg Germany, 1–21. <https://doi.org/10.1145/3544548.3581388>