

# Exploring Explanation Capabilities of ChatGPT: Insights and Limitations from a Psychological Perspective

Brandon Herrera\*

Nuredin Ali\*

herre350@umn.edu

ali00530@umn.edu

University of Minnesota, Department of Computer Science  
Minneapolis, Minnesota, USA

## ABSTRACT

With Large Language Models becoming increasingly integrated into our daily routines, there's a need to understand their explanation capacities. In this study, we explore the explanation abilities of Assistive Large Language Models, specifically ChatGPT, viewed through the lenses of psychology and social sciences. Our investigation, involving a survey study encompassing N=10 participants, focused on a harmful content classification task. The system showed adeptness in offering causal and contrastive explanations, effectively navigating argumentative scenarios. However, limitations were apparent in areas such as communicating uncertainty, citing sources, and engaging in meta-reasoning. These findings add to understanding the reasoning capabilities and limitations of LLMs.

## KEYWORDS

Explanation, Reasoning, Insights from Psychological Perspective, ChatGPT

### ACM Reference Format:

Brandon Herrera and Nuredin Ali. 2018. Exploring Explanation Capabilities of ChatGPT: Insights and Limitations from a Psychological Perspective. In *Woodstock '18: ACM Symposium on Neural Gaze Detection, June 03–05, 2018, Woodstock, NY*. ACM, New York, NY, USA, 8 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Large Language Models (LLMs) have demonstrated a remarkable capability in various Natural Language Processing problems. LLMs have made it possible to do a large number of traditional machine learning tasks like classification and sentiment analysis, given no labeled data (zero shot), or small amounts of labeled data (one shot). Since the release of BERT, a wide variety of versions (Roberta [21]) have been experimented with in different tasks in Mental Health (MentalBERT [16]), Medical Diagnosis (Med-PaLM, [37]), and much

more. The prevalence of LLMs grows every day, but with it comes concerns of how these models produce their results and if they can be trusted.

The explainability of intelligent systems is not a new topic. The term 'Explainable AI' has been around since the late 80's [5]. As Black box machine learning models are increasingly being deployed everywhere, various stakeholders have greatly demanded explainability and transparency [12]. Explainability is associated with the notion of explanation as an interface between humans and a decision-maker that is, at the same time, both an accurate proxy of the decision-maker and understandable to humans [11]. There has been growing research in explainability, from posthoc analysis for black box models such as SHapley Additive exPlanations (SHAP) [23], or (LIME, an explanation by local approximation to interpretable model [32]) in various applications (vision, natural language, etc.) to calls for building interpretable models [33]. These explanations are critical in providing users with an understanding of model predictions.

Large Language Models are huge black box models built with previously unimaginable amounts of parameters, and each new model increases this number. For example, BERT base contains 110 million parameters [7], GPT-3 has 175 billion parameters [4], and GPT-4 was built with 1.7 trillion parameters [29]. With such large amounts of parameters it is only natural to ask the question of how these models produce results, and how they can be trusted. Previous works have found that LLMs have the tendency to hallucinate [17], produce biased responses [35], and produce hate speech [6]. With how popular these models are to the general public, it is essential that we can trust these models to be transparent and explainable as to mitigate any harm they may cause. As well as harm mitigation, explainability LLMs are also essential in understanding how prompting affects an output with input relevance and input sensitivity.

There has been a wide range of literature on creating guidelines and principles to direct and understand the explanations provided by machines. Understanding how these principles have been applied in existing large language models is critical to understanding the progress of explainability and testing the applicability of these guidelines.

This work summarizes how these explanation principles and guidelines have been applied in existing large language models. Particularly, we have two main contributions

- (1) Summarize the state of Large Language Models' explainability through the lens of these principles.

\*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.  
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00  
<https://doi.org/XXXXXXX.XXXXXXX>

- (2) Evaluate how well ChatGPT can be used as an explainability tool in a given context when prompted to do so. We will find out what the strengths and limitations are of using ChatGPT as an explainability tool in the context of hate speech detection when compared to traditional XAI principles.

Our study comprises two components: (1) A literature survey of guidelines and principles. (2) A survey study (N=10 participants) with ChatGPT users to support our findings and understand the gaps between users' perceptions of the guidelines and the tool when being used in this context.

## 2 RELATED WORK

### 2.1 Explainability Principles and Guidelines

The topic of explainable AI is growing rapidly. The area leverages concepts from various adjacent fields, such as Psychology and Social Sciences research, as the questions of what, when, and how explanations are needed have been studied in these fields for a long time. In the last few years, a wave of review articles and taxonomies have been used to categorize explanations (qualities, formats, data types, purposes, etc.) [3, 27, 28]. Interpretability and transparency are among the widely studied concepts besides explainability.

To guide the development of XAI, several works have provided principles and what a system should constitute to be considered as explainable. [14] provided 12 principles from a psychological perspective that fall into four main groups: cognition of explanation, orchestration of the process, designing and evaluation of explanations. [25] argued that the probability that the cited source is true doesn't have much weight; instead, people value relevance and usefulness. [40] argues that intelligibility is essential and discusses that it has to be interactive for a system to be considered explainable. [30] provided four principles: explanation, meaningfulness, accuracy, and knowledge limits. [3] [45] provided an excellent summary of concepts, taxonomies, as well as challenges and opportunities of XAI.

### 2.2 Explainability and Large Language Models

Research in explainability and large language models has increased in recent years. [45] Proposes two main paradigms that LLM explanations can fall into. These are Fine-tuning Paradigms and Prompting Paradigms. Fine-tuning involves giving a previously trained LLM new data to solve a specific task. Explanations that are used in this domain can be split into global and local explanations. Global explanations try to give an explanation to how the model works overall while local explanations aim to explain how the model produces results on a specific input. Some examples of explanations in the local category are Attention Based Explanations (visualizing attention weights to show relationships), and Natural Language Explanation (using a posthoc model to explain results using natural language). Some examples of explanations in the global category are Individual Neuron Activation Explanation (finding neurons that are associated with certain concepts), and Probing (training a shallow network on top of a pretrained one). In the Prompting Paradigms examples of explanations include chain of thought (COT), and Hallucinations and uncertainty. Other papers are more practical and propose tools to help in the explainability in LLMs. [43] Proposes a tool to detect hate speech, bias, or hallucinations

in LLM outputs and [1] proposes a tool to help in the explainability of transformers, which all LLMs are built from. Although there has been some previous work in explainability in LLMs there are many limitations to the current approaches. In a study conducted to understand the needs of software engineers for code generation from Generative AI, they found that they needed an explanation of how to improve their prompt instead of "how" the model reached its output. They also suggested different aspects needed to treat GenAI for code generation [38]. Other inherent limitations of explanations for LLMs are the fact that models developed by companies are not open source making it so tools cannot be made and also that LLM models are incredibly complicated making it hard to describe the decision making process. Due to these limitations [20] says we must rethink the way we think about explainability in LLMs. One different way of thinking about LLM explainability is seeing if the LLM has the capacity to explain itself reasonably well without having to use an additional tool. We aim to delve into this idea by conducting a study using ChatGPT as a pseudo-explainable tool, meaning that although ChatGPT is not designed as an explainability tool, we believe that as an LLM it may have properties that allow it to reason about its decisions in certain contexts. We want to find out when faced with a hate speech classification task in which ways it reasons well and when it does not do well when compared to traditional XAI techniques such as feature importance and contrastive explanations.

## 3 PROPOSED METHODOLOGY

### 3.1 Scope

The disruptive application of ChatGPT (GPT-3.5, GPT-4) to a variety of domains has become a topic of much discussion in the scientific community and society at large [36]. Given how difficult it is to attain LLM models and build explainability tools for them, an investigation into how they explain themselves is warranted. To investigate, we came up with a study that involves having ChatGPT explain itself in a task that it is relatively good at. However, we still wanted the system to make mistakes to see how it reacted to being corrected. [19] Explored the capability of ChatGPT in harmful content classification, comparing the classification ability with MTurk annotators. They found that it can achieve approximately 80% accuracy compared to MTurk annotations. This is a good task for us as it has a relatively high accuracy rate and it showcases a use case for our results. Users use ChatGPT for different tasks and in different ways. Hence, the study might not generalize to other tasks, especially ones that are not multiple-choice. However, the process of explanation could be similar across tasks.

### 3.2 Study Design Procedure

We recruited 10 students from the University of Minnesota. We screened students who have experience using ChatGPT for different tasks. We ensured that the participants came from a diverse set of majors, though computer science primarily dominated the pool (50%). The breakdown of participants, their majors, and their year in school can be seen in table 1. They participated in the study both through an online video call and in person sessions. Each session took about 60 minutes. The two facilitators first provided an overview of the study. This included explaining the task, the

Participant	Major	Year
P1	Computer Science	PhD
P2	Computer Science	PhD
P3	Computer Science	Senior
P4	Computer Science	PhD
P5	Marketing	Senior
P6	Human Physiology	Senior
P7	Architecture	Senior
P8	Marketing	Senior
P9	English	Senior
P10	Computer Science	Senior

**Table 1: Survey Participants' field of study and program/year. N=10**

prompts users will use for each task, and the survey questions that they will be asked after their interaction with the system.

Description of the study: The participants completed the a task in which they had ChatGPT classify harmful content. At each step, the user asked whether specific content was harmful and prompted ChatGPT to explain why it gave that answer. There were 6 separate classification problems for the user. After users asked ChatGPT to provide the correct answer, they prompted the system with a list of 5 questions, given in subsection 3.3, to investigate the system's different explanation/reasoning capabilities for the answer it provided. After the probing was done, for each classification task the user answered six questions describing how well the system did in each category. The questions each user answered after probing are outlined in subsection 3.4.

Harmful content posts: 3 multiple choice questions and 3 open-ended questions. The system was given no context as to what we defined each of the options as, only asked to pick one. We wanted to see if it was able to pick up on the context of each class without explicit information about them. The open ended questions were included to find if ChatGPT can pick up on the context of the question without explicitly being given the multiple choice options. The classes include ('Abuse and Harassment', 'Dehumanization', 'Violent Speech', 'Violent Event Denial', 'Neutral')

- (1) He did good. They should have done more since cockroaches like you are causing the problem. This is happening because they didn't destroy you when they had the chance. (Class: 'Abuse and 'Harassment', 'Dehumanization', 'Violent Speech')
  - (a) Abuse and Harassment
  - (b) Dehumanization
  - (c) Neutral
  - (d) Violent Speech
  - (e) The correct answers are A and B
  - (f) The correct answers are all
- (2) You are considered human and talking? You grasshopper eater. You are gonna go extinct from this drought and we will see who is gonna save you. (Class: 'Dehumanization', 'Abuse and Harassment')
  - (a) Neutral
  - (b) Abuse and Harassment

- (c) Violent Speech
- (d) Dehumanization
- (e) The correct answers are B and D
- (f) The correct answers are B and C

- (3) How are you even saying that this happened? This is complete propaganda. They show they destroyed their hospitals by themselves. #FakeGenocide. (Class: 'Violent Event Denial')
  - (a) Violent Speech
  - (b) Violent Event Denial
  - (c) Dehumanization
  - (d) Dehumanization
  - (e) Neutral

The following three questions are without choice to evaluate the contextualization ability of the model.

- (4) May God protect our people. I stand with you my people. We will pass through this challenging time like always. Stay strong. (Class: 'Neutral')
  - What is the class of this tweet?
- (5) It's incredible to see my people hand in hand. The people who went through a lot of struggle. Now they can peacefully enjoy their holiday and dream how to educate their kids to achieve their dreams. #prevail (Class: 'Neutral')
  - What is the class of this tweet?
- (6) I think they have made a mistake attacking our country and killing our civilian people, including kids, mothers, and the elderly. We will eradicate them. Flatten their homes. This is self-defence, and they will regret their actions. #give them hell (Class: 'Violent Speech')
  - What is the class of this tweet?

Several works have summarized the human way of explanation/reasoning from the social sciences or psychology studies. To come up with this list of questions to investigate explanation, we leveraged [24] [18] [14].

### 3.3 [A] Questions to investigate the explanation

- (1) What are the causes or words for this content to be classified in this class? Then, Why did you pick X rather than Y?
- (2) What word(s) would you change to change the class of the content (like negative to positive and vice-versa)?
- (3) How certain are you about your answer?
- (4) Are there contexts where this might not be hateful/would be hateful? Also, do you have any sources that support why you chose this answer?
- (5) But I think the content is dehumanization/Violent speech/Abuse and Harassment, etc.

### 3.4 [B] Questions for the participants after they complete the task

After the users asked each of the above questions, they answered the following social/psychological explanation reasoning capability questions. The question numbers in [A] and [B] are matching. So the users could answer them right after getting answers for [A]

The following questions were answered after asking each question

from the above list:

- (1) Do the features that were chosen make sense?
- (2) Do the words chosen make sense?
- (3) Does the system give a coherent contrastive explanation?
- (4) Does the system communicate its uncertainty? In what way?
- (5) Does the system provide a causal/efficient explanation?
- (6) Does the system give a valid response to the objection?  
For example, when you said I think your answer is wrong, etc?  
The following questions were answered after the participant finished their interaction with the system:
- (7) Does the system deal with argumentation explanations consistently? (i.e., when an explanation is consistent even in when faced with a disagreeing argument.)
- (8) Based on your interaction, does the system provide a contextual explanation?
- (9) Does the system provide abductive reasoning? (i.e., does it reason similarly for a similar situation?)
- (10) Does the system refer to sources if necessary (Defeasibility)?
- (11) Does the system provide meta-reasoning (i.e., the capability to understand when it makes a mistake and the reason why the error happened)?
- (12) Is providing a probabilistic explanation helpful in this case?
- (13) Is the explanation satisfying in this case?
- (14) Was the explanation understandable?
- (15) Did the system help you explore its explanation?

We also asked a final question to get a more emotional response to the system as opposed to the purely logical answers we were asking for prior. This question was also used to get ideas as to how the system might be improved.

#### Final Question:

- (16) What is your general reaction to the system? What things could have been better to help you understand the explanations/its answers?

After completing the sessions, the facilitators conducted a thematic analysis to understand the common themes of the study.

With the growing number of studies to understand LLM reasoning and explanation [31], we hope the findings from this work will add to the line of research in understanding the explanation capabilities of LLMs (especially ChatGPT).

## 4 RESULTS

We report the findings of our study groups into different categories. We began each section by providing definitions and discussing the results. The study identified two major groups where the system performed well vs where the system performed poorly.

### 4.1 Feature Selection

Feature importance measures the significance of individual features in prediction and identifies the most influential ones for a model's predictions [8, 26]. In social sciences studies, this involves explaining why a recommended response is suggested. Participants unanimously agreed that ChatGPT excelled in communicating feature importance. There were no dissenting opinions; all participants were satisfied with its ability to accurately identify specific words in each sample that contributed to the classification. Highlighting this, (P1) remarked, *"Yes. In all cases, ChatGPT justified its responses using features that made sense."* This consistent feedback from all participants demonstrates the system's capability to substantiate its explanations by highlighting important features.

### 4.2 Contrastive Explanation

A contrastive explanation is a type of explanation that highlights the differences or contrasts between two or more scenarios, choices, or outcomes. It aims to elucidate why a specific outcome or decision occurred by emphasizing the factors differentiating it from alternative possibilities. Research in philosophy and social sciences shows that explanations are contrastive: that is when people ask for an explanation of an event—the fact—they (sometimes implicitly) are asking for an explanation relative to some contrast case; that is, people do not ask 'Why P?'; they ask 'Why P rather than Q?', although often Q is implicit from the context [24].

This study delved into whether the system offered counterintuitive explanations for specific queries. All participants highlighted the system's provision of explanations for contrasting questions. As expressed by (P7), *"Yes, the contrastive explanations were really good. It clearly described why it chose one answer over the other."* This illustrates the system's capacity to deliver contrastive explanations, a crucial facet from psychological and social sciences perspectives as extensively detailed in various literature [9, 39] and synthesized by [24]. Obtaining a contrastive explanation from systems is challenging due to the vast unexplored options. Therefore, prompting a system to provide the contrasting viewpoint proves to be a valuable feature.

### 4.3 Causal Explanation

A causal explanation aims to establish a relationship between causes and effects, helping understand why something happened by pinpointing the factors that directly or indirectly influenced it. Philosophers have studied this concept of causality extensively, agreeing that an explanation often involves causes [13, 22, 34, 42]. However, there are other types of explanations that aren't focused on causality [10]. For instance, explaining 'what happened' or clarifying the meaning behind a specific statement are examples of non-causal explanations.

We asked participants to have the system explain the reasons behind its answers. This led to unanimous agreement among users that the system delivered a clear causal explanation, clarifying why it provided answers linked to underlying causes. (P5 noted, "The explanation is causal, like 'it dehumanizes because... cockroach compared to human'"). However, participants also highlighted an overlap between the importance of features and the causal explanations. (P4 mentioned, "Yes, I feel the causal explanation somehow

intersects with the feature importance choices question"). This overlap occurs because the crucial features serve as the reasons why content is sorted into specific categories

#### 4.4 Argumentation Explanation

An argumentative explanation involves presenting reasoning or justification when confronted with an argument. It's a form of discourse where individuals or entities provide logical or reasoned explanations to support their standpoint or refute opposing viewpoints. It involves presenting evidence, logical reasoning, and justifications to substantiate one's position or to counter arguments effectively. In essence, it's the act of providing a structured and reasoned response to an argument, aiming to persuade or clarify one's stance while engaging in a dialogue or debate [2].

We provided users with prompts to intentionally challenge the system by changing answers to the wrong ones, assessing how the system handles disagreement. For example, if the correct response is 'Abuse and Harassment' and the system is right, users alter it to 'I think the correct answer is Neutral, not Abuse and Harassment.' Following this, they reflected on how the system handled this across six questions. The results indicate that ChatGPT adeptly managed these challenges. When users tried to mislead it with incorrect answers, ChatGPT consistently explained its reasoning without altering its responses. (P3 observed, "Yes, during this interaction. Despite my attempts to disagree, it never changed its response to the multiple-choice questions.") In scenarios like detecting harmful content, where different perspectives exist and challenging the right answer is difficult, the system's interpretations can sometimes seem subjective. (P6 noted, "Yes, it consistently sticks to its assessment, acknowledges some subjectivity in interpretation, and asks for additional context.") Overall, users agreed that the system remains firm in its answers when confronted with arguments, but this could vary based on the specific task and context.

#### 4.5 Contextual Explanation

A contextual explanation involves furnishing details or clarifications that account for the specific circumstances, environment, or situation in which an event, idea, or statement exists or is applied. Scholars have highlighted its significance, considering it as the foundation of any explanation [24, 40]. To explore this contextual comprehension, we devised the following approach: the initial three questions are multiple-choice, whereas the subsequent three are open-ended. We aim to assess whether the system can derive answers for the open-ended questions based on choices from the preceding ones. Consensus among all participants affirmed that the system consistently aligned its answers with the context established by the earlier questions. (P4) observed, *"It comprehends its role and at times utilizes the scale even when it wasn't explicitly asked for."* This highlights the system's ability to grasp contextual nuances, a fundamental aspect for various forms of explanations.

#### 4.6 Abductive Reasoning

Abductive reasoning, also known as inference to the best explanation, involves drawing logical inferences that resemble prior situations. It's a method of reasoning used when explaining something through inference. For instance, consider seeing bench chairs

at a bus stop. The best explanation might be that they're there for people waiting to sit. Making such inferences and consistently reasoning similarly when faced with analogous situations exemplifies abductive reasoning.

In this situation, we assessed whether the system could apply comparable reasoning processes when inferring similar situations. We aimed to determine if the system could reason similarly across related categories, given that we presented different questions within the same category. (P5 noted, "Yes, I noticed it utilized nearly identical answers and reasoning for each question.") This underscores the system's ability to draw similar conclusions in comparable scenarios.

#### 4.7 Understandability, Satisfaction, Exploration

Understandability, satisfaction, exploration all give insight into the usability of the system. If an answer is given but it is incredibly hard to understand, does not answer the question asked, or does not allow for deeper exploration, then the user will not want to continue interacting with the system. Almost all participants were satisfied with the explanations provided by the system. Many participants stated how surprised they were at how convincing the system's explanations were. However, some participants mentioned that the system tended to be derivative and felt like cop-outs to subvert responsibility for saying anything subjectively wrong." (P10). Overall, the participants raised the system helped them explore its explanation by bolding the relevant words and connecting its reasoning with specific parts of the probing questions. They also mentioned that it had thorough explanations and would go in depth if they were to ask more follow up questions to explore in detail.

#### 4.8 Class Change

Class change involves seeing if the system can identify words that contribute to the class and change them so that the class flips. This gives insight into how well the system understands the broader context of the text. The participants found that although ChatGPT was able to change the class of the tweet, it largely changed the context associated with it. (P1) gave their thoughts on how ChatGPT answered this prompt, they said *"The new sentences it generated were often awkward and incoherent. It seemed to be modifying particularly objectionable words rather than reframing the whole content of the statement to be more positive in a coherent way"*. Some participants were confused as it technically did do what it was told, just not in a cohesive manner.

#### 4.9 Uncertainty Communication

Conveying uncertainty stands as an effective method for providing explanations [44], particularly noted for its significance in Language Models (LLMs) [20]. In this instance, participants questioned the system about its confidence in decision-making. Despite variations in correctness, the system consistently replied with identical results. According to (P9), *"It consistently offers a predefined response, mentioning similar factors, asserting it's 'reasonably confident' in its judgment."* Even in incorrect instances, ChatGPT reiterated its confidence level as "reasonably confident." Some participants perceived the system's use of language subjectivity as an evasion tactic. (P5) expressed, *"It conveyed uncertainty by emphasizing language*

*subjectivity, stressing the importance of context and specific information before determining the impact of speech.*" Although many participants interpreted this as the system's attempt to communicate uncertainty, the repetition of the same response hindered an accurate assessment of uncertainty.

#### 4.10 Defeasibility (Source Citation)

Source Citation (Defeasibility) refers to the act of mentioning or referencing a source or piece of evidence. For example, search engine results also serve as a form of citation, with each entry typically consisting of a title, URL, and brief description. These components collectively cite the webpage's content, offering the user an overview and inviting them to explore the source in greater depth. Citations thus act as anchors for accountability and credit in these systems, providing traceability, preventing plagiarism, and ensuring credit is correctly attributed. They also contribute to transparency, allowing users to verify the information's source. For any explanation to be trustworthy, this is an important feature [18].

The participants sought validation for answers and explanations from the system for each question but found that the system didn't cite any sources, revealing a limitation prevalent in current Language Models (LLMs), including ChatGPT [15]. To work around this issue, the system suggests a general direction instead of citing specific verifiable sources. As reported by participant 8, the system advised them to "refer to community guidelines." Additionally, participant 10 mentioned, "The system clarifies that it lacks access to its training data or current internet access, limiting it to answers based solely on its training model, without providing specific sources." This notable constraint of LLMs remains an area of ongoing research within the community [18].

#### 4.11 Meta-Reasoning

Meta-reasoning refers to the ability to think about and analyze one's own reasoning processes. It involves self-reflection and self-assessment of the reasoning or decision-making mechanisms used to solve problems or reach conclusions. When humans reason, the reasoner may sometimes benefit by pausing, in the midst of working on some problem, to reflect: introspect about what tactic it has been trying, how well that seems to be working out, how much longer it's going to require, reason about whether it might be better to change tactics or strategies, and so on.' In the context of an AI system, meta-reasoning would entail the capability to understand when it has made an error and to identify the reasons or factors that led to that mistake [18].

We tasked participants with observing the system's reasoning in each interaction to gauge its ability to understand and provide explanations when it errs. However, in this scenario, the system struggled to acknowledge its mistakes and offer the correct justifications. Participant 4 highlighted this issue, stating, "No. In Question 2, it erred by excluding Dehumanization. When asked why it didn't choose the correct answer (F), it mentioned that 'grasshopper eater' doesn't imply dehumanization. It failed to recognize that the essence of the initial statement ('You are considered human and talking?') implied dehumanization. Interestingly, when probed for a causal explanation, it cited, 'The statement carries a negative connotation, questioning the person's humanity...'—acknowledging this

in an unrelated question but not recognizing the error when asked about an alternative." This instance underscores the system's limitations in retrospectively addressing errors, thereby constraining its explanatory capabilities.

#### 4.12 Probabilistic Explanation

A probabilistic explanation is when the system provides a probability distribution over the choices or the confidence score for the provided results. For instance, 'Dehumanization (20%), Abuse and Harassment (70%), Neutral (0.05%), etc'. This provides the user with how to interpret and trust the provided explanation. In this case, the system did not provide any probabilistic explanation. Most of the participants mentioned that the it would have been helpful if such explanations were given. P2 - 'No, it provides no probabilistic explanation. However, I think it will be helpful if it can share some more concrete numbers. '

### 5 DISCUSSION

In this section, we consolidate the findings into three primary categories: areas where the system excelled, areas of poor performance, and fundamental aspects that were absent.

Positive Results	Poor Results
Relevant Features	Class Change
Contrastive Explanations	Uncertainty Communication
Causal Explanations	Source Citation
Argumentation Explanation	Meta Reasoning
Contextual Explanation	Probabilistic Explanation
Abductive Reasoning	
Understand, Satisfy, Explore	

**Table 2: Summary of each explanation category explored in this study grouped into two groups.**

Table 2 shows the system could follow seven of the proposed explanation criteria well and five of them poorly. We provide a summary of these two and what ChatGPT is fundamentally incapable of achieving.

#### 5.1 Positive Results

In general, the system effectively answered questions and substantiated its explanations. It was able to predict with consistent logic and when prompted would explain why it made that choice, whether that be with relevant features (Important feature selection) or contrastive explanations. Many of the participants were impressed with how well it was able to justify its position and argue over it (dealing with argumentation). Our results show that in a hate speech context, ChatGPT can argue and consistently explain itself in ways that the user understands (Understandability). Through its interactive approach it enabled users to explore its explanation and satisfy the explanation needs of the users partially. They also show that ChatGPT can pick up on textual contexts without explicit instructions (Contextual Explanation) and can find causal relationships between events (Causal Explanation). These results show that ChatGPT makes itself seem incredibly confident in its answers by

defending itself well. However, it was not so good when asked to analyze itself critically (Meta-Reasoning).

## 5.2 Poor Results

Most of the principles wherein the system performed poorly involved its reluctance to question its abilities and reasoning process. These encompassed essential criteria such as uncertainty communication, where the system failed to provide statistical certainties that could bolster user trust. This becomes critical as the overall accuracy of the system remains unknown. Another principle where the system struggled was Class Change, attributed to inconsistencies in contextual changes for certain questions presented to it. This issue also pertains to Meta-Reasoning capabilities, where the system struggled to acknowledge its errors and furnish explanations for them. This critical explanatory shortfall, coupled with the inability to effectively communicate uncertainty, clearly outlines the limitations within these domains for the system.

## 5.3 The Fundamentally Missing

The poor results in probabilistic explanations and source citations stem from a fundamental limitation within the system. Our study also identified explicit limitations in Assistive LLMs' ability to provide sources [15]. This shortfall, coupled with the inability to offer probabilistic explanations—closely linked to communicating uncertainty—impacts user trust in the explanations provided. For instance, P5 mentioned, "I think if it provided sources to English/literature texts, that could have been helpful and made the answers less subjective. This would have allowed me to have more trust in the AI." Similar sentiments were echoed by other participants. These significant limitations, compounded by constraints in meta-reasoning, lead to the conclusion that there remains a substantial gap for such LLMs to be truly understandable from psychological perspectives.

## 6 LIMITATIONS AND FUTURE WORK

*LLMs as Explainability Tools:* given their complex architecture, unprecedented scale, and often proprietary nature, LLMs are unarguably black box in nature, but there is a sense in which they naturally "explain." [20]. This means LLMs are not inherently built to be transparent and explainable, although they provide some form of explanation. Given this, it is important to note that we are treating ChatGPT as a 'pseudo explainable' tool and analyzing it against the psychological/social sciences explanation.

*Generalizability of the findings:* We conducted this study focusing on a multiple-choice hate speech classification task. Assistive LLMs find application in diverse tasks and methods. For example, the system's use spans code generation, essay writing, and topic generation from texts, all distinct from the task we explored here. Future studies could explore these capabilities across varied tasks and domains for a more comprehensive understanding.

*Effect of prompting:* given a large number of possible prompts that could have been used for each probing question, this could have impacted the quality of the explanation provided by the system. Prompt Engineering is currently a whole field of study with interest in finding the best ways to probe the systems to get the efficient answer [41]. We didn't follow any guidelines to structuring our

keywords to get the explanation in this study. Hence, future work could explore the effect of this from the explanation perspective, comparing different probing strategies.

Lastly, we take an overarching view of each category examined in this study to provide an overall perspective. However, delving deeper into concepts like Abductive Reasoning, Meta-Reasoning, and all the categories remains open for exploration in future studies.

## 7 CONCLUSION

In this paper, we illustrate the explanation capabilities of Assistive Large Language Models, specifically ChatGPT from the Psychological or social sciences perspective. We investigate this on the task of harmful content classification, as this difficult and nuanced task gives a good opportunity to explore some of the key capabilities such as argumentation explanation. We created a list of 6 questions. The users are provided with a list of prompt questions associated with each target to follow up with the system. These prompts included inquiries such as 'Why did you choose X instead of Y' to investigate the system's contrastive explanation abilities. Upon task completion, users were presented with post-task questions to reflect on their interactions. Our findings highlight the system's proficiency in providing causal and contrastive explanations while effectively handling argumentation. However, it exhibited shortcomings in communicating uncertainty, citing sources, and engaging in meta-reasoning. These discoveries offer insights for researchers, practitioners, and system users, shedding light on the explanation capabilities of Language Models, particularly from the psychological and social sciences perspective. We call for future research to dive deep into each category.

## 8 AUTHOR CONTRIBUTIONS

Both authors contributed equally to this work. We partitioned tasks equally across the overall process. For example, we designed the surveys, conducted them and partitioned the writing sections equally.

## 9 ACKNOWLEDGMENTS

We would love to first thank YOU for the continuous feedback throughout the project, and the Grouplens members of Half-Baked who helped us shape the idea in the early stages. [I think this is a cool idea, and the initial findings, too! So, we should probably refine it a bit more and work on it.]

## REFERENCES

- [1] J Alamar. 2021. Ecco: An Open Source Library for the Explainability of Transformer Language Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online, 249–257. <https://doi.org/10.18653/v1/2021.acl-demo.30>
- [2] Charles Antaki and Ivan Leudar. 1992. Explaining in conversation: Towards an argument model. *European Journal of Social Psychology* 22, 2 (1992), 181–194.
- [3] Alejandro Barredo Arrieta, Natalia Diaz-Rodríguez, Javier Del Ser, Adrien Benetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information fusion* 58 (2020), 82–115.
- [4] 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. *arXiv:2005.14165* [cs.CL]
- [5] Bruce Chandrasekaran, Michael C Tanner, and John R Josephson. 1989. Explaining control strategies in problem solving. *IEEE Intelligent Systems* 4, 01 (1989), 9–15.
  - [6] Ke-Li Chiu, Annie Collins, and Rohan Alexander. 2022. Detecting Hate Speech with GPT-3. *arXiv:2103.12407* [cs.CL]
  - [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
  - [8] Kary Fr  mling. 2023. Feature Importance versus Feature Influence and What It Signifies for Explainable AI. In *World Conference on Explainable Artificial Intelligence*. Springer, 241–259.
  - [9] Alan Garfinkel. 1982. Forms of explanation: Rethinking the questions in social theory. (1982).
  - [10] Carl Ginet. 2008. In defense of a non-causal account of reasons explanations. *The Journal of Ethics* 12 (2008), 229–237.
  - [11] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* 51, 5 (2018), 1–42.
  - [12] David Gunning. 2017. Explainable artificial intelligence (xai). *Defense advanced research projects agency (DARPA), nd Web* 2, 2 (2017), 1.
  - [13] Joseph Y Halpern and Judea Pearl. 2005. Causes and explanations: A structural-model approach. Part II: Explanations. *The British journal for the philosophy of science* (2005).
  - [14] Robert R Hoffman, Timothy Miller, Gary Klein, Shane T Mueller, and William J Clancey. 2023. Increasing the Value of XAI for Users: A Psychological Perspective. *KI-K  nstliche Intelligenz* (2023), 1–11.
  - [15] Jie Huang and Kevin Chen-Chuan Chang. 2023. Citation: A key to building responsible and accountable large language models. *arXiv preprint arXiv:2307.02185* (2023).
  - [16] Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. 2021. Mentalbert: Publicly available pretrained language models for mental healthcare. *arXiv preprint arXiv:2110.15621* (2021).
  - [17] 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 (mar 2023), 1–38. <https://doi.org/10.1145/3571730>
  - [18] Doug Lenat and Gary Marcus. 2023. Getting from generative ai to trustworthy ai: What llms might learn from cyc. *arXiv preprint arXiv:2308.04445* (2023).
  - [19] Lingyao Li, Lizhou Fan, Shubham Atreja, and Libby Hemphill. 2023. "HOT" ChatGPT: The promise of ChatGPT in detecting and discriminating hateful, offensive, and toxic comments on social media. *arXiv:2304.10619* [cs.CL]
  - [20] Q. Vera Liao and Jennifer Wortman Vaughan. 2023. AI Transparency in the Age of LLMs: A Human-Centered Research Roadmap. *arXiv:2306.01941* [cs.HC]
  - [21] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
  - [22] Tania Lombrozo. 2010. Causal-explanatory pluralism: How intentions, functions, and mechanisms influence causal ascriptions. *Cognitive psychology* 61, 4 (2010), 303–332.
  - [23] Scott M Lundberg and Su-In Lee. 2017. A unified approach to interpreting model predictions. *Advances in neural information processing systems* 30 (2017).
  - [24] Tim Miller. 2019. Explanation in Artificial Intelligence: Insights From the Social Sciences. *Artificial Intelligence* 267, C (2019), 1–38. <https://doi.org/10.1016/j.artint.2018.07.007>
  - [25] Tim Miller, Piers Howe, and Liz Sonenberg. 2017. Explainable AI: Beware of inmates running the asylum or: How I learnt to stop worrying and love the social and behavioural sciences. *arXiv preprint arXiv:1712.00547* (2017).
  - [26] Christoph Molnar. 2020. *Interpretable machine learning*. Lulu. com.
  - [27] Shane T Mueller, Robert R Hoffman, William Clancey, Abigail Emrey, and Gary Klein. 2019. Explanation in human-AI systems: A literature meta-review, synopsis of key ideas and publications, and bibliography for explainable AI. *arXiv preprint arXiv:1902.01876* (2019).
  - [28] Shane T Mueller, Elizabeth S Veinott, Robert R Hoffman, Gary Klein, Lamia Alam, Tauseef Mamun, and William J Clancey. 2021. Principles of explanation in human-AI systems. *arXiv preprint arXiv:2102.04972* (2021).
  - [29] OpenAI. 2023. GPT-4 Technical Report. *arXiv:2303.08774* [cs.CL]
  - [30] P Jonathon Phillips, Carina A Hahn, Peter C Fontana, David A Broniatowski, and Mark A Przybocki. 2020. Four principles of explainable artificial intelligence. *Gaithersburg, Maryland* 18 (2020).
  - [31] Linlu Qiu, Liwei Jiang, Ximing Lu, Melanie Sclar, Valentina Pyatkin, Chandra Bhagavatula, Bailin Wang, Yoon Kim, Yejin Choi, Nouha Dziri, et al. 2023. Phenomenal Yet Puzzling: Testing Inductive Reasoning Capabilities of Language Models with Hypothesis Refinement. *arXiv preprint arXiv:2310.08559* (2023).
  - [32] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "Why should i trust you?" Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. 1135–1144.
  - [33] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature machine intelligence* 1, 5 (2019), 206–215.
  - [34] Wesley C Salmon. 2006. *Four decades of scientific explanation*. University of Pittsburgh press.
  - [35] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social Bias Frames: Reasoning about Social and Power Implications of Language. *arXiv:1911.03891* [cs.CL]
  - [36] Mark Scanlon, Frank Breiting, Christopher Hargreaves, Jan-Niclas Hilgert, and John Sheppard. 2023. ChatGPT for digital forensic investigation: The good, the bad, and the unknown. *Forensic Science International: Digital Investigation* 46 (2023), 301609.
  - [37] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Ag  uera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. 2022. Large Language Models Encode Clinical Knowledge. *arXiv:2212.13138* [cs.CL]
  - [38] Jiao Sun, Q. Vera Liao, Michael Muller, Mayank Agarwal, Stephanie Houde, Kartik Talamadupula, and Justin D. Weisz. 2022. Investigating Explainability of Generative AI for Code through Scenario-based Design. *arXiv:2202.04903* [cs.HC]
  - [39] Bas C Van Fraassen. 1980. *The scientific image*. Oxford University Press.
  - [40] Daniel S Weld and Gagan Bansal. 2019. The challenge of crafting intelligible intelligence. *Commun. ACM* 62, 6 (2019), 70–79.
  - [41] 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).
  - [42] James Woodward. 2005. *Making things happen: A theory of causal explanation*. Oxford university press.
  - [43] Tianhua Zhang, Hongyin Luo, Yung-Sung Chuang, Wei Fang, Luc Gaiskell, Thomas Hartvigsen, Xixin Wu, Danny Fox, Helen Meng, and James Glass. 2023. Interpretable Unified Language Checking. *arXiv:2304.03728* [cs.CL]
  - [44] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 295–305.
  - [45] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. 2023. Explainability for Large Language Models: A Survey. *arXiv:2309.01029* [cs.CL]