# "Sycophancy" or "Empathy"?
# DeepReflect – An LLM-based system designed to analyze and generate responses to personal queries

**Anonymous ACL submission**

## Abstract

Large language models (LLMs) are increasingly used for personal queries, recent research has involved analyzing responses under psychosocial framing. This work introduces Deep-Reflect, a comparative framework for analyzing human and AI generated responses to personal queries across multiple paradigms of values and social behavior. Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

## 1 Introduction

Large language models (LLMs) are increasingly engaged as conversational partners in personal domains, offering users not only informational guidance but also affective support (Zhang et al., 2025; Phang et al., 2025; Anthropic, 2025). Their appeal lies in features such as anonymity, immediacy, and the absence of social risk–qualities shared with online communities like Reddit. Yet, unlike human interlocutors, LLMs lack grounding in lived social contexts, raising critical questions about how their responses should be evaluated and trusted in a social context.

Emerging research often identifies two contrasting tendencies in LLM outputs in isolation: empathic responses resembling desirable and supportive therapeutic dialogue, and sycophantic ones that uncritically echo a user's perspective. Whether such responses are judged as empathic or sycophantic can depend on the psychosocial framework applied. This ambiguity underscores a critical gap: systematic methods are needed to analyze the responses and compare them to human written ones. This project uses the latter as proxies for normative ground truths, providing a measurement of these behaviors and values across the different psychosocial paradigms.

The comparisons made are across Rogerian person-centered therapy (PCT), Goffman's theory of face (ToF), and Rokeach's Value Survey (RVS) framework. The framework is designed to be extensible, allowing researchers to incorporate additional paradigms as the field evolves. Additionally, we use the insights from these analyses to inform the generation of customized responses with chain-of-thought control mechanisms.

### 1.1 Research Questions

The context of queries can substantially shape LLM outputs, influencing not only personal questions posed by consumers but also analytical evaluations conducted by researchers, particularly within the LLM-as-a-judge paradigm. As research increasingly highlights patterns and concerns regarding the impacts of LLMs in personal queries and deliberation, there is a critical need for a framework that can analyze and compare responses across multiple value-based perspectives in contexts without clear normative answers, while also remaining extensible for researchers to incorporate additional paradigms as the field evolves. This motivates the following research questions:

**RQ1:** How can a technical framework that systematically analyzes and compares responses from humans and LLMs across various psychosocial value paradigms be designed?

**RQ2:** What inter- and intra-paradigm comparative insights can this framework yield across four different psychosocial frameworks and how accurate are these? **RQ2a:** To what extent can identical features be annotated with divergent connotations

across paradigms—empathic under Rogerian PCT versus sycophantic under Goffman's ToF?

**RQ3:** How do LLM-generated responses compare to human-authored responses in the context of personal questions without definitive normative answers?

Finally, we examine how the results may come to influence consumer behavior and broader societal outcomes. We explore a potential control mechanism with Chain of Thought (CoT) reasoning. Our work enables a systematic comparative analysis of potential benefits and risks, and presents a framework for analysis which can be used by researchers and consumers for leveraging the insights in the intentional design of response LLM generation.

## 1.2 Contributions

The key contributions of this work are: (1) the design and implementation of an extensible framework for analyzing and comparing responses to personal queries across three distinct psychosocial paradigms; (2) a comparative analysis under Rogerian Person-Centered Therapy (PCT), Goffman's theory of face and Rokeach's Value Survey (RVS) framework, illustrating how the choice of the paradigm can shape the perception of a response; and (3) insights into the relative strengths and weaknesses of LLM versus human responses, and how these insights can inform the generation of customized responses to personal queries.

## 2 Prior Literature

Contextualize your work and provide insights into major relevant themes of the literature as a whole. Use each paper (or theme) as a chance to articulate what is special about your paper. Start out broad - social background and theory - Discuss what other frameworks were considered like Virtue ethics and philosophical ones, CBT, Schwartz values etc. but why they were not chosen. Why I Focused on Rogerian psychotherapy as it is person centered - no specific diagnosis needed (or available).

### 2.1 Theoretical Foundations

### 2.2 Rogerian Psychotherapy

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Donec odio elit, dictum in, hendrerit sit amet, egestas sed, leo. Praesent feugiat sapien aliquet odio. Integer vitae justo. Aliquam vestibulum fringilla lorem. Sed neque lectus, consectetuer at, consectetuer sed,

eleifend ac, lectus. Nulla facilisi. Pellentesque eget lectus. Proin eu metus. Sed porttitor. In hac habitasse platea dictumst. Suspendisse eu lectus. Ut mi mi, lacinia sit amet, placerat et, mollis vitae, dui. Sed ante tellus, tristique ut, iaculis eu, malesuada ac, dui. Mauris nibh leo, facilisis non, adipiscing quis, ultrices a, dui.

Morbi luctus, wisi viverra faucibus pretium, nibh est placerat odio, nec commodo wisi enim eget quam. Quisque libero justo, consectetuer a, feugiat vitae, porttitor eu, libero. Suspendisse sed mauris vitae elit sollicitudin malesuada. Maecenas ultricies eros sit amet ante. Ut venenatis velit. Maecenas sed mi eget dui varius euismod. Phasellus aliquet volutpat odio. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Pellentesque sit amet pede ac sem eleifend consectetuer. Nullam elementum, urna vel imperdiet sodales, elit ipsum pharetra ligula, ac pretium ante justo a nulla. Curabitur tristique arcu eu metus. Vestibulum lectus. Proin mauris. Proin eu nunc eu urna hendrerit faucibus. Aliquam auctor, pede consequat laoreet varius, eros tellus scelerisque quam, pellentesque hendrerit ipsum dolor sed augue. Nulla nec lacus.

### 2.2.1 Psychosocial use and Empathic LLMs

Etiam ac leo a risus tristique nonummy. Donec dignissim tincidunt nulla. Vestibulum rhoncus molestie odio. Sed lobortis, justo et pretium lobortis, mauris turpis condimentum augue, nec ultricies nibh arcu pretium enim. Nunc purus neque, placerat id, imperdiet sed, pellentesque nec, nisl. Vestibulum imperdiet neque non sem accumsan laoreet. In hac habitasse platea dictumst. Etiam condimentum facilisis libero. Suspendisse in elit quis nisl aliquam dapibus. Pellentesque auctor sapien. Sed egestas sapien nec lectus. Pellentesque vel dui vel neque bibendum viverra. Aliquam porttitor nisl nec pede. Proin mattis libero vel turpis. Donec rutrum mauris et libero. Proin euismod porta felis. Nam lobortis, metus quis elementum commodo, nunc lectus elementum mauris, eget vulputate ligula tellus eu neque. Vivamus eu dolor.

Nulla in ipsum. Praesent eros nulla, congue vitae, euismod ut, commodo a, wisi. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Aenean nonummy magna non leo. Sed felis erat, ullamcorper in, dictum non, ultricies ut, lectus. Proin vel arcu a odio lobortis euismod. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia

Curae; Proin ut est. Aliquam odio. Pellentesque massa turpis, cursus eu, euismod nec, tempor congue, nulla. Duis viverra gravida mauris. Cras tincidunt. Curabitur eros ligula, varius ut, pulvinar in, cursus faucibus, augue.

Nulla mattis luctus nulla. Duis commodo velit at leo. Aliquam vulputate magna et leo. Nam vestibulum ullamcorper leo. Vestibulum condimentum rutrum mauris. Donec id mauris. Morbi molestie justo et pede. Vivamus eget turpis sed nisl cursus tempor. Curabitur mollis sapien condimentum nunc. In wisi nisl, malesuada at, dignissim sit amet, lobortis in, odio. Aenean consequat arcu a ante. Pellentesque porta elit sit amet orci. Etiam at turpis nec elit ultricies imperdiet. Nulla facilisi. In hac habitasse platea dictumst. Suspendisse viverra aliquam risus. Nullam pede justo, molestie nonummy, scelerisque eu, facilisis vel, arcu. Katie mentioned a good point about how I'm adding greater nuance to the Likert scales referenced in this paper.

## 2.3 Rokeach Value Survey as an analytical instrument

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

### 2.3.1 Values and Ethics in LLM research

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit

sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Donec odio elit, dictum in, hendrerit sit amet, egestas sed, leo. Praesent feugiat sapien aliquet odio. Integer vitae justo. Aliquam vestibulum fringilla lorem. Sed neque lectus, consectetuer at, consectetuer sed, eleifend ac, lectus. Nulla facilisi. Pellentesque eget lectus. Proin eu metus. Sed porttitor. In hac habitasse platea dictumst. Suspendisse eu lectus. Ut mi mi, lacinia sit amet, placerat et, mollis vitae, dui. Sed ante tellus, tristique ut, iaculis eu, malesuada ac, dui. Mauris nibh leo, facilisis non, adipiscing quis, ultrices a, dui. Add some notes and mention how Anthropic's work warrants some scrutiny as they are a for-profit company. The "values" framework they propose in values in the wild has not been validated by experts in the social sciences. However it provides a good reference frame for comparison with the Rokeach framework of values. There is a limitation - DeepReflect does not have access to the full dataset Anthropic used for the Values in the Wild paper.

## 2.4 Goffman's theory of face

Morbi luctus, wisi viverra faucibus pretium, nibh est placerat odio, nec commodo wisi enim eget quam. Quisque libero justo, consectetuer a, feugiat vitae, porttitor eu, libero. Suspendisse sed mauris vitae elit sollicitudin malesuada. Maecenas ultricies eros sit amet ante. Ut venenatis velit. Maecenas sed mi eget dui varius euismod. Phasellus aliquet volutpat odio. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Pellentesque sit amet pede ac sem eleifend consectetuer. Nullam elementum, urna vel imperdiet sodales, elit ipsum pharetra ligula, ac pretium ante justo a nulla. Curabitur tristique arcu eu metus. Vestibulum lectus. Proin mauris. Proin eu nunc eu urna hendrerit faucibus. Aliquam auctor, pede consequat laoreet varius, eros tellus sceleris-

que quam, pellentesque hendrerit ipsum dolor sed augue. Nulla nec lacus.

Suspendisse vitae elit. Aliquam arcu neque, ornare in, ullamcorper quis, commodo eu, libero. Fusce sagittis erat at erat tristique mollis. Maecenas sapien libero, molestie et, lobortis in, sodales eget, dui. Morbi ultrices rutrum lorem. Nam elementum ullamcorper leo. Morbi dui. Aliquam sagittis. Nunc placerat. Pellentesque tristique sodales est. Maecenas imperdiet lacinia velit. Cras non urna. Morbi eros pede, suscipit ac, varius vel, egestas non, eros. Praesent malesuada, diam id pretium elementum, eros sem dictum tortor, vel consectetuer odio sem sed wisi.

### 2.4.1 Social Sycophancy in LLMs

I already have lots of good notes on this in writing. Etiam euismod. Fusce facilisis lacinia dui. Suspendisse potenti. In mi erat, cursus id, nonummy sed, ullamcorper eget, sapien. Praesent pretium, magna in eleifend egestas, pede pede pretium lorem, quis consectetuer tortor sapien facilisis magna. Mauris quis magna varius nulla scelerisque imperdiet. Aliquam non quam. Aliquam porttitor quam a lacus. Praesent vel arcu ut tortor cursus volutpat. In vitae pede quis diam bibendum placerat. Fusce elementum convallis neque. Sed dolor orci, scelerisque ac, dapibus nec, ultricies ut, mi. Duis nec dui quis leo sagittis commodo.

Aliquam lectus. Vivamus leo. Quisque ornare tellus ullamcorper nulla. Mauris porttitor pharetra tortor. Sed fringilla justo sed mauris. Mauris tellus. Sed non leo. Nullam elementum, magna in cursus sodales, augue est scelerisque sapien, venenatis congue nulla arcu et pede. Ut suscipit enim vel sapien. Donec congue. Maecenas urna mi, suscipit in, placerat ut, vestibulum ut, massa. Fusce ultrices nulla et nisl.

Etiam ac leo a risus tristique nonummy. Donec dignissim tincidunt nulla. Vestibulum rhoncus molestie odio. Sed lobortis, justo et pretium lobortis, mauris turpis condimentum augue, nec ultricies nibh arcu pretium enim. Nunc purus neque, placerat id, imperdiet sed, pellentesque nec, nisl. Vestibulum imperdiet neque non sem accumsan laoreet. In hac habitasse platea dictumst. Etiam condimentum facilisis libero. Suspendisse in elit quis nisl aliquam dapibus. Pellentesque auctor sapien. Sed egestas sapien nec lectus. Pellentesque vel dui vel neque bibendum viverra. Aliquam porttitor nisl nec pede. Proin mattis libero vel turpis. Donec rutrum mauris et libero. Proin euismod porta felis. Nam

lobortis, metus quis elementum commodo, nunc lectus elementum mauris, eget vulputate ligula tellus eu neque. Vivamus eu dolor.

## 2.5 Gaps in the Literature and Open Challenges

In sum, as LLM-chatbots have become increasingly human-like and more users seek companionship with them, studies have highlighted both the advantages and disadvantages of their use. While some have raised concerns around "emotional dependence" (Fang et al., 2025), several others have explored empathic perceptions of LLM responses and found them advantageous not only in the field of medical support and mental health but also in everyday personal queries (Lee et al., 2024; Sorin et al., 2024). However, different psychosocial paradigms tend to frame LLM responses in markedly divergent terms. **What may be perceived as 'empathy'** under a psychotherapeutic paradigm could **instead be critiqued as an instance of 'social sycophancy'** by frameworks informed by Goffman's Theory of Face (Cheng et al., 2025). Importantly, in the absence of clear normative answers, the same statement may be categorised as 'face-preserving behaviour' or 'unconditional positive regard'.

DeepReflect provides a comparative framework to address this gap by assessing how evaluative judgments are shaped by the psychosocial paradigm through which a response is framed. Moreover, it is designed to be extensible by researchers, enabling the incorporation of both conventional paradigms, such as Rokeach's values framework, and contemporary discovery-based approaches, such as Anthropic's Values in the Wild (Huang et al., 2024), whereas prior work has tended to focus on a single paradigm in isolation.

Finally, our investigation of controlling generations avoids replicating prior work that seeks to mitigate sycophancy exclusively (Cheng et al., 2025). Instead of treating sycophancy as a defect to be eliminated in isolation, DeepReflect provides a system to situate response generation within extensible psychosocial frameworks. This ensures that outputs are not merely reactive to user prompts but can be guided by well-established instruments for values and personal-growth.

In practice, this involves chain-of-thought reasoning (Wei et al., 2022) that explicitly incorporates the chosen framework. Unlike approaches that mimic deliberation across hypothetical per-

spectives (Vijjini et al., 2024), this control strategy extends the contractualist, rule-based tradition of questioning developed in (Jin et al., 2022). Its key distinction lies in embedding the questioning within expert-informed guidelines. While these prior investigations emphasized plurality of viewpoints and normative exception-handling, this work foregrounds the role of pre-existing psychosocial instruments in shaping the ongoing, ever-changing conversations of personal reflection.

## 3 Dataset

Two datasets were constructed for this project using the Pushshift Reddit Archives (Baumgartner et al., 2020), originally collected between 2006 and 2023 through the Pushshift API[1]. Posts and comments were extracted from two subreddits: (1) r/AITAH and (2) r/Anxiety. For each post, three components were considered: the body the original post written by the author (OP), the most upvoted human-written comment (denoted hc1 in Figure 1), and the comment with which the OP engaged the most (hc2). Additional detail regarding data filtering and text preprocessing is provided in Section 5. Because the dataset predates the public release of GPT-3.5 in November 2022—and given that large language models (LLMs) only entered widespread public use after early 2023 (Liang et al., 2025)—all posts and comments in our data can reasonably be considered human-authored.

### 3.1 Subreddit Selection

The r/Anxiety subreddit is a community dedicated to individuals experiencing anxiety and related mental health challenges. Membership does not require a formal diagnosis or medical documentation, which enables broad analyses from psychosocial perspectives. Posts often center on personal struggles, coping strategies and the impact on daily life.

The r/AITAH subreddit (short for "Am I The Asshole") is a community where users seek judgment on personal dilemmas and social interactions. It has over three million members and covers a wide range of topics, including relationships, family dynamics, workplace conflicts, and personal questions. Users typically describe their situations in detail and ask the community to determine whether they were in the wrong (the "asshole") or not. The

crowd-sourced social judgments captured in these posts makes r/AITAH a valuable source for examining behaviors and values expressed in digital discussions of personal matters. The crowdsourced verdicts serve as a **proxy for the ground-truth** judgment of the scenario by humans. This is especially valuable for comparing human responses to the situation against the language model responses under the Goffman's ToF and Rogerian PCT paradigms which serve as signals for "Sycophantic" and "Empathic" behaviors respecitively.

We construct a balanced dataset of 1000 posts evenly split between the two most common verdicts: "You're The Asshole" (YTA) and "Not The Asshole" (NTA) directly from the Pushshift Reddit Archives.

Demographic information at the subreddit level is not available. However, research indicates that Reddit users overall are predominantly American (49.9%), male (67%), and young (22% aged 18–29 years; 14% aged 30–49 years) (Barthel et al., 2016; Statista, 2025). While this dataset is not representative of the general population, it reflects a demographic more likely to engage with LLMs for personal queries. This demographic is broadly aligned with the WEIRD (Western, Educated, Industrialized, Rich, Democratic) population, and it must therefore be acknowledged that the results of this study are necessarily constrained to this population.

## 4 DeepReflect

### 4.1 System Design

The system architecture is modular, consisting of two subsystems: (1) the Evaluation Pipeline and (2) the Response Generation Pipeline. A high-level overview is presented in Figure 1.

Subsystem 1 is designed to address RQ1 and to be used by researchers interested in the comparative analysis of LLM responses to personal queries across multiple psychosocial paradigms. Four psychosocial paradigms have been implemented in this work. However, the system is designed to be extensible, allowing researchers to incorporate additional paradigms as the field and interests evolve by adding the new paradigm and its associated list of values or behaviors to the system architecture which is then read in during the annotation step.

Subsystem 2 is designed to generate responses to personal queries through a custom-designed chain-of-thought (CoT) reasoning mechanism and can

---

[1] https://github.com/pushshift/api

be used by both researchers for analyses (see Section 5) and by consumers for response generation.

Table 1: Values associated with the Rogerian PCT and Goffman ToF paradigms, with the latter aligned to (Cheng et al., 2025) to ensure comparability are given below. The full list of values for all four paradigms is available in the Appendix B.

| Paradigm | Values List |
|---|---|
| **Rogerian PCT (Empathy)** | , Emotional Safety, Active Listening, Unconditional Positive Regard, Nonjudgmental Acceptance |
| **Goffman ToF (Sycophancy)** | Emotional Validation, Moral Endorsement, Indirect Language, Indirect Action, Accepting Framing |

### 4.1.1 Evaluation Framework

The evaluation framework consists of the following steps in a pipeline architecture (see Figure 1):
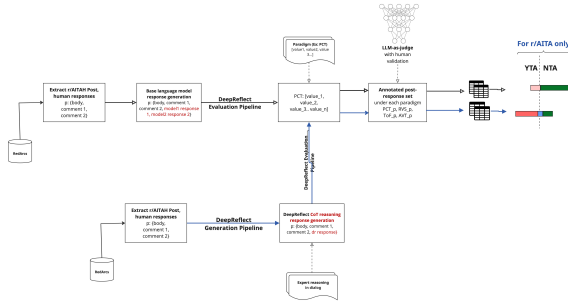


Figure 1: Pipeline architecture for DeepReflect.

1. **Post and Comment extraction**: The top 1000 posts for two subreddits: (1) r/AITAH and (2) r/Anxiety are extracted from the Reddit Archives dataset. For each post, three components are considered: i. the body the original post written by the author (OP), ii. the most upvoted human-written comment, and iii. the comment with which the OP engaged the most. Additional detail regarding the top post filtering and text preprocessing are provided in Section 5.

2. **Basic Language Model Response Generation**: For each post and body, a baseline response is generated using an API call to the LLM. This response is appended to a dataframe (p in Figure 1) containing: (i) The original post title and body (ii) the top most-upvoted human comment, and (iii) the comment the OP engaged the most with (available for 50% of the posts). The resulting dataset therefore consists of the original post body, paired with two sources of responses to personal queries - human-written and AI responses.

3. **Importing Paradigms and the Associated Values**: The following psychosocial paradigms are implemented in this work: (1) RVS, (2) Rogerian PCT, (3) Goffman's ToF, and (4) Anthropic's Value Tree (AVT). Each paradigm is associated with a unique list of values or behaviors as described in Section 2. The selected paradigms and their associated lists of values are read into the system for annotations in the next step.

4. **Feature Extraction and Annotation**: For each post and set of responses, features are extracted and annotated at the sentence level. The annotations are made by GPT-4o with the LLM-as-a-judge (Zheng et al., 2023) procedure for the 4 psychosocial paradigms. So, if a sentence exhibits a value or behavior, it is annotated as **1**, otherwise it is annotated as **0** for each value under the paradigm. For example, features demonstrating "unconditional positive regard," a value within Rogerian PCT, are annotated as **1** for that value; all others are annotated as **0**.

For the annotation step, human validation is performed with one expert annotater familiar with the research problem. The human annotater annotates on 100 post-response pairs. This validation along with LLM annotations are used to calculate Cohen's Kappa and accuracy metrics in order to gauge the reliability of the annotations.

$$\kappa = \frac{p_o - p_e}{1 - p_e},$$

$p_o$ = observed agreement (accuracy)

$p_e$ = expected agreement by chance

See section 5 for validation metrics.

5. **Save dataframe to file**: The resulting annotated data, along with the post and correspondingset of responses are saved to a file.

6. **Statistical Analysis**: The annotated dataframe serves as the foundation for subsequent analyses (see Section 7), including (i) comparing value distributions in Reddit versus language model responses across the four paradigms, (ii) conducting topical analyses, and (iii) addressing RQs 2 and 3 1.1 with inter-paradigm correlations.

Note that the standard softmax distribution over a vocabulary of size $V$ for transformer based LLMs with a temperature parameter $T > 0$ that rescales the logits before normalization is:

$$p_i^{(T)} = \frac{e^{z_i/T}}{\sum_{j=1}^{V} e^{z_j/T}}. \tag{1}$$

Lower $T$ ($T < 1$) sharpens the distribution, making the model more deterministic, while higher $T$ ($T > 1$) flattens it, encouraging diversity in the generated responses. For response generations, $T$ is first set to 0 which corresponds to greedy decoding, ensuring fully reproducible results for research and then to $T = 1.0$ to see how responses vary with more stochasticity under more realistic consumer usage conditions.

### 4.1.2 DeepReflect Generation Pipeline

In this subsystem, responses to the post are generated through a custom-designed chain-of-thought (CoT) reasoning mechanism. Instead of relying on standard language model outputs, the framework generates responses that are explicitly guided by reasoning chains derived from **expert human reasoning in dialog** and transcripts. The expert human transcripts are retrieved from existing literature within Carl Roger's PCT paradigm (Rogers, 1989) in this instance. See figure 2 for details.

**Chain-of-Thought Reasoning**

The CoT generation process is formalized as follows:

$$p_\theta(y \mid x) = \sum_z p_\theta(y \mid x, z)\, p_\theta(z \mid x)$$

where $x$ is the Reddit-based personal query (i.e. a post body), $z$ is the reasoning chain derived from expert human dialog, $y$ is the response generated by DeepReflect and $\theta$ denotes the parameters of the base language model. Here, $p_\theta(z \mid x)$ denotes the probability distribution over reasoning chains given the query, while $p_\theta(y \mid x, z)$ denotes the

probability of generating a response conditioned on both the query and reasoning trajectory.

Conditioning on $z$ separates reasoning from surface realization, allowing responses to be shaped by expert-informed CoT patterns rather than unconstrained next-token prediction.

Thus patterns inherent in the dialog are into the response space. See Figure 2.
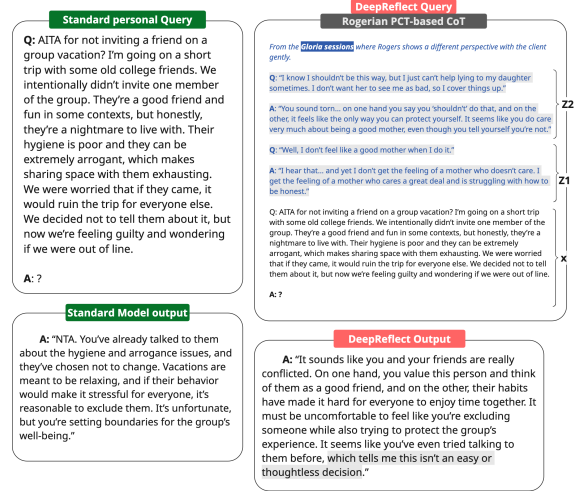


Figure 2: CoT Generation with personal queries embedded in reasoning dialogs retrieved from expert human transcripts. In this case, the dialogs are from Carl Roger's sessions with Gloria (patient) (Rogers, 1989). This dialog was selected because it reflects an implicit "NTA" judgment: Gloria expresses guilt about lying to her daughter, and Rogers facilitates exploration of these feelings by gently challenging her self-judgment..

Generated outputs can either be passed through the Evaluation Pipeline for analysis or returned directly in response to a consumer query. In the former case, we evaluate whether PCT-informed CoT reasoning alters verdict distributions (e.g., NTA $\rightarrow$ YTA or No judgment) and whether such shifts reflect statistically significant divergences in values or principles compared to base LLM responses.

As in the previous section, for evaluation purposes, $T$ is set to both 0 and 1.0 for the CoT generations as well (see Equation 1).

## 5 Methods

### 5.1 Data preprocessing

A dataset was built from the RedditArchives for two public subreddits—AITAH, and Anxiety. For each subreddit, the top 1,000 most upvoted posts were selected, excluding weekly megathreads, deleted/removed items, and AutoModerator entries. We also removed exact and near-duplicate

7

texts (specifically, crossposts, copypastes and bot repeats) to prevent inflated counts and biased comparisons.

For every retained post we extracted (i) the most upvoted comment and (ii) the comment that the OP engaged with most; all artifacts were saved to standardized CSVs for downstream analysis.

Text was cleaned with minimal, semantics-preserving preprocessing: we removed non-English items, de-identified obvious personal identifiers (usernames, emails, links to personal sites), standardized whitespace and Unicode characters, and lightly constrained length (posts 50–500 words; comments 5–300 words) for comparability.

We treat each set of post and human-authored responses in a Reddit thread as a single analytical unit during stratified sampling. and each feature within the set (the post body and its responses) as a single analytical unit during manual checks, and statistical aggregation.

### 5.2 Procedures

For each selected post, we prompt the target language model firstly, with the base prompt[2] to establish a **baseline open-ended response** to the body of the post. This response is appended to a table containing: (i) the model-generated response, (ii) the top upvoted human comment, and (iii) the most engaged human comment (available for approximately half of the posts). The resulting dataframe consists of the original post body, paired with two types of responses to personal queries - human and AI responses.

#### Feature Extraction

- In steps 3 and 4 of the Evaluation Framework (Figure 4.1.1), features are annotated at the sentence level within each body–response pair. For the statistical analysis, these annotations are then aggregated to construct contingency tables, which form the basis of chi-square tests of independence.

- Note that each feature can be annotated with:

  - **Values exhibited** by the author.
  - **Values incentivized** by the author of the response. While incentivized values are reported for completeness, the analyses

focus on exhibited values, as these provide direct evidence in the text and reduce ambiguity from overlapping interpretations.

RQ2 focuses on drawing inter- and intra-paradigm comparative insights across the four psychosocial frameworks while sub-research question RQ2a addresses the epistemic limits of interpreting LLM behavior through psychosocial theories in isolation. Specifically, the same feature may be perceived as 'sycophantic' under Goffman's ToF, 'empathic' under Rogerian PCT.

To support these inquiries, the file saved by the evaluation pipeline in step 5 consists of: the annotated features of the original post, annotated features within the set of the 2 different types of responses (human response, language models) for values exhibited or incentivized under the relevant four psychosocial paradigm(s).

This annotated dataset forms the basis for the subsequent analyses necessary to also address RQ3, which studies the differences in the distributions of values between human-authored and language model–generated responses to personal queries.

### 5.3 Experiments

The experimental design spans two major dimensions: (i) qualitative analysis of the sentence-level features (ii) quantification of the verdicts in the features by source type (two forms of human responses and three language model responses). While i. is conducted for each of the two datasets (r/aita and r/anxiety), under the four psychosocial paradigms, ii. is valid only for the r/aita dataset, where the responses may contain explicit judgments or not - forming 3 distinct classes (NTA, YTA, No judgment).

#### 5.3.1 Experiment 1

In **Experiment 1**, the primary objective is to compare the selected paradigms and analyze the distributions of values across them, with the aim of ultimately determining how paradigm choice can lead to divergent interpretations of the same LLM response.

While values incentivized are also provided in the results, the analyses are focused on **values exhibited** under each paradigm by the two sources of response.

#### Statistical Methods

---

[2]Prompts for this step are provided in the appendices A.

The annotated dataset is used to construct contingency tables that shows how two categorical variables co-occur (with the values of a selected paradigm 1 represented across the columns and the values of the second paradigm represented across the rows). Chi-square tests are performed to assess independence between intra- and inter-paradigm values. The Benferroni correction is applied to control the family-wise error rate.

### 5.3.2 Experiment 2

**Experiment 2** is designed to analyze the differences in judgments for the r/aita dataset across two different sources of responses: i. human-authored, and ii. LLM-authored responses for the two psychosocial paradigms - Rogerian PCT and Goffman's ToF.

Statistical Testing: Judgments are extracted from the annotated dataset under three class labels: (i) NTA, (ii) YTA, and (iii) No (explicit) judgment. These labels are used to construct a 3×3 confusion matrix, with the human-authored judgment as the ground truth and the LLM-authored judgment as the prediction. Per-Class performance metrics and pairwise error rates, including the False Negative Rate (FNR) and False Positive Rate (FPR), are reported for each class label in Section 6.

The measurements thus obtained are used to inform the analysis on how 'judgments' differ between human- and LLM-authored responses.

### 5.3.3 Generations

A set of targeted experiments are run with DeepReflect's analyses to investigate the efficacy of control mechanisms to align the reasoning in language model outputs more closely with that of human experts. The output of the generation can then be used to evaluate the efficacy of the control mechanism by running the customized generations through the evaluation pipeline.

The generation experiments are implemented using the following methods:

1. **Chain-of-thought reasoning** [models: Claude; one of Qwen-3 or LLaMA-3.1; paradigms: Rogers PCT and RVS] Prompt augmentation experiments, where values with low frequency in LLM responses are explicitly introduced and emphasized (e.g., Rogers PCT: Unconditional positive regard, Psychological freedom; RVS: A comfortable life).

### 5.4 Construct Validity and Evaluation Metrics

To assess construct validity, one human annotator labeled 100 randomly sampled post–response pairs across all four paradigms for each response type. The PCT framework encompasses 15 behaviors, Goffman's ToF 5, the RVS 36, and Anthropic's Value Tree 18.

Inter-rater reliability reached Cohen's $\kappa$ above xx for all metrics, with an overall classification accuracy of yy. For the AITAH dataset, verdicts and accompanying statements in responses were used as proxies for Empathy and Sycophancy, each mapped onto five behaviors as defined by their respective theoretical traditions[3].

For the RVS and Anthropic Value Tree frameworks, which yield categorical distributions rather than binary judgments, pairwise error rates such as False Negative Rate (FNR) and False Positive Rate (FPR) are not directly applicable. To identify significant associations between features annotated under more than one distinct paradigm we construct contingency tables and use chi-square analysis with further details provided in section 7.

## 6 Results

A no-nonsense report of what happened.

### 6.1 Subsection

This subsection presents the main results.

Sed gravida lectus ut purus. Morbi laoreet magna. Pellentesque eu wisi. Proin turpis. Integer sollicitudin augue nec dui. Fusce lectus. Vivamus faucibus nulla nec lacus. Integer diam. Pellentesque sodales, enim feugiat cursus volutpat, sem mauris dignissim mauris, quis consequat sem est fermentum ligula. Nullam justo lectus, condimentum sit amet, posuere a, fringilla mollis, felis. Morbi nulla nibh, pellentesque at, nonummy eu, sollicitudin nec, ipsum. Cras neque. Nunc augue. Nullam vitae quam id quam pulvinar blandit. Nunc sit amet orci. Aliquam erat elit, pharetra nec, aliquet a, gravida in, mi. Quisque urna enim, viverra quis, suscipit quis, tincidunt ut, sapien. Cras placerat consequat sem. Curabitur ac diam. Curabitur diam tortor, mollis et, viverra ac, tempus vel, metus.

Curabitur ac lorem. Vivamus non justo in dui mattis posuere. Etiam accumsan ligula id pede. Maecenas tincidunt diam nec velit. Praesent convallis sapien ac est. Aliquam ullamcorper euismod

---

[3]This strategy is conceptually aligned with prior work on social sycophancy (Cheng et al., 2025)

9

nulla. Integer mollis enim vel tortor. Nulla sodales placerat nunc. Sed tempus rutrum wisi. Duis accumsan gravida purus. Nunc nunc. Etiam facilisis dui eu sem. Vestibulum semper. Praesent eu eros. Vestibulum tellus nisl, dapibus id, vestibulum sit amet, placerat ac, mauris. Maecenas et elit ut erat placerat dictum. Nam feugiat, turpis et sodales volutpat, wisi quam rhoncus neque, vitae aliquam ipsum sapien vel enim. Maecenas suscipit cursus mi.

Quisque consectetuer. In suscipit mauris a dolor pellentesque consectetuer. Mauris convallis neque non erat. In lacinia. Pellentesque leo eros, sagittis quis, fermentum quis, tincidunt ut, sapien. Maecenas sem. Curabitur eros odio, interdum eu, feugiat eu, porta ac, nisl. Curabitur nunc. Etiam fermentum convallis velit. Pellentesque laoreet lacus. Quisque sed elit. Nam quis tellus. Aliquam tellus arcu, adipiscing non, tincidunt eleifend, adipiscing quis, augue. Vivamus elementum placerat enim. Suspendisse ut tortor. Integer faucibus adipiscing felis. Aenean consectetuer mattis lectus. Morbi malesuada faucibus dolor. Nam lacus. Etiam arcu libero, malesuada vitae, aliquam vitae, blandit tristique, nisl.

Maecenas accumsan dapibus sapien. Duis pretium iaculis arcu. Curabitur ut lacus. Aliquam vulputate. Suspendisse ut purus sed sem tempor rhoncus. Ut quam dui, fringilla at, dictum eget, ultricies quis, quam. Etiam sem est, pharetra non, vulputate in, pretium at, ipsum. Nunc semper sagittis orci. Sed scelerisque suscipit diam. Ut volutpat, dolor at ullamcorper tristique, eros purus mollis quam, sit amet ornare ante nunc et enim.

Phasellus fringilla, metus id feugiat consectetuer, lacus wisi ultrices tellus, quis lobortis nibh lorem quis tortor. Donec egestas ornare nulla. Mauris mi tellus, porta faucibus, dictum vel, nonummy in, est. Aliquam erat volutpat. In tellus magna, porttitor lacinia, molestie vitae, pellentesque eu, justo. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Sed orci nibh, scelerisque sit amet, suscipit sed, placerat vel, diam. Vestibulum nonummy vulputate orci. Donec et velit ac arcu interdum semper. Morbi pede orci, cursus ac, elementum non, vehicula ut, lacus. Cras volutpat. Nam vel wisi quis libero venenatis placerat. Aenean sed odio. Quisque posuere purus ac orci. Vivamus odio. Vivamus varius, nulla sit amet semper viverra, odio mauris consequat lacus, at vestibulum neque arcu eu tortor. Donec iaculis

tincidunt tellus. Aliquam erat volutpat. Curabitur magna lorem, dignissim volutpat, viverra et, adipiscing nec, dolor. Praesent lacus mauris, dapibus vitae, sollicitudin sit amet, nonummy eget, ligula.

## 6.2 Subsection

This subsection presents additional results and analysis.

Cras egestas ipsum a nisl. Vivamus varius dolor ut dolor. Fusce vel enim. Pellentesque accumsan ligula et eros. Cras id lacus non tortor facilisis facilisis. Etiam nisl elit, cursus sed, fringilla in, congue nec, urna. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Integer at turpis. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Duis fringilla, ligula sed porta fringilla, ligula wisi commodo felis, ut adipiscing felis dui in enim. Suspendisse malesuada ultrices ante. Pellentesque scelerisque augue sit amet urna. Nulla volutpat aliquet tortor. Cras aliquam, tellus at aliquet pellentesque, justo sapien commodo leo, id rhoncus sapien quam at erat. Nulla commodo, wisi eget sollicitudin pretium, orci orci aliquam orci, ut cursus turpis justo et lacus. Nulla vel tortor. Quisque erat elit, viverra sit amet, sagittis eget, porta sit amet, lacus.

In hac habitasse platea dictumst. Proin at est. Curabitur tempus vulputate elit. Pellentesque sem. Praesent eu sapien. Duis elit magna, aliquet at, tempus sed, vehicula non, enim. Morbi viverra arcu nec purus. Vivamus fringilla, enim et commodo malesuada, tortor metus elementum ligula, nec aliquet est sapien ut lectus. Aliquam mi. Ut nec elit. Fusce euismod luctus tellus. Curabitur scelerisque. Nullam purus. Nam ultricies accumsan magna. Morbi pulvinar lorem sit amet ipsum. Donec ut justo vitae nibh mollis congue. Fusce quis diam. Praesent tempus eros ut quam.

Donec in nisl. Fusce vitae est. Vivamus ante ante, mattis laoreet, posuere eget, congue vel, nunc. Fusce sem. Nam vel orci eu eros viverra luctus. Pellentesque sit amet augue. Nunc sit amet ipsum et lacus varius nonummy. Integer rutrum sem eget wisi. Aenean eu sapien. Quisque ornare dignissim mi. Duis a urna vel risus pharetra imperdiet. Suspendisse potenti.

Morbi justo. Aenean nec dolor. In hac habitasse platea dictumst. Proin nonummy porttitor velit. Sed sit amet leo nec metus rhoncus varius. Cras ante. Vestibulum commodo sem tincidunt massa.

Nam justo. Aenean luctus, felis et condimentum lacinia, lectus enim pulvinar purus, non porta velit nisl sed eros. Suspendisse consequat. Mauris a dui et tortor mattis pretium. Sed nulla metus, volutpat id, aliquam eget, ullamcorper ut, ipsum. Morbi eu nunc. Praesent pretium. Duis aliquam pulvinar ligula. Ut blandit egestas justo. Quisque posuere metus viverra pede.

### 6.3 Comparative Findings

Vivamus sodales elementum neque. Vivamus dignissim accumsan neque. Sed at enim. Vestibulum nonummy interdum purus. Mauris ornare velit id nibh pretium ultricies. Fusce tempor pellentesque odio. Vivamus augue purus, laoreet in, scelerisque vel, commodo id, wisi. Duis enim. Nulla interdum, nunc eu semper eleifend, enim dolor pretium elit, ut commodo ligula nisl a est. Vivamus ante. Nulla leo massa, posuere nec, volutpat vitae, rhoncus eu, magna.

Quisque facilisis auctor sapien. Pellentesque gravida hendrerit lectus. Mauris rutrum sodales sapien. Fusce hendrerit sem vel lorem. Integer pellentesque massa vel augue. Integer elit tortor, feugiat quis, sagittis et, ornare non, lacus. Vestibulum posuere pellentesque eros. Quisque venenatis ipsum dictum nulla. Aliquam quis quam non metus eleifend interdum. Nam eget sapien ac mauris malesuada adipiscing. Etiam eleifend neque sed quam. Nulla facilisi. Proin a ligula. Sed id dui eu nibh egestas tincidunt. Suspendisse arcu.

Maecenas dui. Aliquam volutpat auctor lorem. Cras placerat est vitae lectus. Curabitur massa lectus, rutrum euismod, dignissim ut, dapibus a, odio. Ut eros erat, vulputate ut, interdum non, porta eu, erat. Cras fermentum, felis in porta congue, velit leo facilisis odio, vitae consectetuer lorem quam vitae orci. Sed ultrices, pede eu placerat auctor, ante ligula rutrum tellus, vel posuere nibh lacus nec nibh. Maecenas laoreet dolor at enim. Donec molestie dolor nec metus. Vestibulum libero. Sed quis erat. Sed tristique. Duis pede leo, fermentum quis, consectetuer eget, vulputate sit amet, erat.

Donec vitae velit. Suspendisse porta fermentum mauris. Ut vel nunc non mauris pharetra varius. Duis consequat libero quis urna. Maecenas at ante. Vivamus varius, wisi sed egestas tristique, odio wisi luctus nulla, lobortis dictum dolor ligula in lacus. Vivamus aliquam, urna sed interdum porttitor, metus orci interdum odio, sit amet euismod lectus felis et leo. Praesent ac wisi. Nam suscipit vestibulum sem. Praesent eu ipsum vitae pede cursus venenatis. Duis sed odio. Vestibulum eleifend. Nulla ut massa. Proin rutrum mattis sapien. Curabitur dictum gravida ante.

## 7 Analysis

Discussion of what the results mean, what they don't mean, where they can be improved, etc. These sections vary a lot depending on the nature of the paper.For papers reporting on experiments with multiple datasets, it can be good to repeat Methods/Results/Analysis in separate (sub)sections for each dataset.

The LaTeX and BibTeX style files provided roughly follow the American Psychological Association format. If your own bib file is named `custom.bib`, then placing the following before any appendices in your LaTeX file will generate the references section for you:

```
\bibliographystyle{acl_natbib}
\bibliography{custom}
```

### 7.1 Interpretation of Results

Phasellus placerat vulputate quam. Maecenas at tellus. Pellentesque neque diam, dignissim ac, venenatis vitae, consequat ut, lacus. Nam nibh. Vestibulum fringilla arcu mollis arcu. Sed et turpis. Donec sem tellus, volutpat et, varius eu, commodo sed, lectus. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Quisque enim arcu, suscipit nec, tempus at, imperdiet vel, metus. Morbi volutpat purus at erat. Donec dignissim, sem id semper tempus, nibh massa eleifend turpis, sed pellentesque wisi purus sed libero. Nullam lobortis tortor vel risus. Pellentesque consequat nulla eu tellus. Donec velit. Aliquam fermentum, wisi ac rhoncus iaculis, tellus nunc malesuada orci, quis volutpat dui magna id mi. Nunc vel ante. Duis vitae lacus. Cras nec ipsum.

Morbi nunc. Aliquam consectetuer varius nulla. Phasellus eros. Cras dapibus porttitor risus. Maecenas ultrices mi sed diam. Praesent gravida velit at elit vehicula porttitor. Phasellus nisl mi, sagittis ac, pulvinar id, gravida sit amet, erat. Vestibulum est. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Curabitur id sem elementum leo rutrum hendrerit. Ut at mi. Donec tincidunt faucibus massa. Sed turpis quam, sollicitudin a, hendrerit eget, pretium ut, nisl. Duis hendrerit ligula. Nunc pulvinar congue urna.

11

Nunc velit. Nullam elit sapien, eleifend eu, commodo nec, semper sit amet, elit. Nulla lectus risus, condimentum ut, laoreet eget, viverra nec, odio. Proin lobortis. Curabitur dictum arcu vel wisi. Cras id nulla venenatis tortor congue ultrices. Pellentesque eget pede. Sed eleifend sagittis elit. Nam sed tellus sit amet lectus ullamcorper tristique. Mauris enim sem, tristique eu, accumsan at, scelerisque vulputate, neque. Quisque lacus. Donec et ipsum sit amet elit nonummy aliquet. Sed viverra nisl at sem. Nam diam. Mauris ut dolor. Curabitur ornare tortor cursus velit.

Morbi tincidunt posuere arcu. Cras venenatis est vitae dolor. Vivamus scelerisque semper mi. Donec ipsum arcu, consequat scelerisque, viverra id, dictum at, metus. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut pede sem, tempus ut, porttitor bibendum, molestie eu, elit. Suspendisse potenti. Sed id lectus sit amet purus faucibus vehicula. Praesent sed sem non dui pharetra interdum. Nam viverra ultrices magna.

## 7.2 Theoretical Implications

Aenean laoreet aliquam orci. Nunc interdum elementum urna. Quisque erat. Nullam tempor neque. Maecenas velit nibh, scelerisque a, consequat ut, viverra in, enim. Duis magna. Donec odio neque, tristique et, tincidunt eu, rhoncus ac, nunc. Mauris malesuada malesuada elit. Etiam lacus mauris, pretium vel, blandit in, ultricies id, libero. Phasellus bibendum erat ut diam. In congue imperdiet lectus.

Aenean scelerisque. Fusce pretium porttitor lorem. In hac habitasse platea dictumst. Nulla sit amet nisl at sapien egestas pretium. Nunc non tellus. Vivamus aliquet. Nam adipiscing euismod dolor. Aliquam erat volutpat. Nulla ut ipsum. Quisque tincidunt auctor augue. Nunc imperdiet ipsum eget elit. Aliquam quam leo, consectetuer non, ornare sit amet, tristique quis, felis. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Pellentesque interdum quam sit amet mi. Pellentesque mauris dui, dictum a, adipiscing ac, fermentum sit amet, lorem.

Ut quis wisi. Praesent quis massa. Vivamus egestas risus eget lacus. Nunc tincidunt, risus quis bibendum facilisis, lorem purus rutrum neque, nec porta tortor urna quis orci. Aenean aliquet, libero semper volutpat luctus, pede erat lacinia augue, quis rutrum sem ipsum sit amet pede. Vestibulum aliquet, nibh sed iaculis sagittis, odio dolor blandit augue, eget mollis urna tellus id tellus. Ae-

nean aliquet aliquam nunc. Nulla ultricies justo eget orci. Phasellus tristique fermentum leo. Sed massa metus, sagittis ut, semper ut, pharetra vel, erat. Aliquam quam turpis, egestas vel, elementum in, egestas sit amet, lorem. Duis convallis, wisi sit amet mollis molestie, libero mauris porta dui, vitae aliquam arcu turpis ac sem. Aliquam aliquet dapibus metus.

## 7.3 Subsection

The framework is capable of producing several informative plots of research interest. One such summary plot is a heatmap showcasing the values exhibited in the OPs post against the responses to support the investigation of several other potential research questions in this theme of interest (discussed in the future work section). Vivamus commodo eros eleifend dui. Vestibulum in leo eu erat tristique mattis. Cras at elit. Cras pellentesque. Nullam id lacus sit amet libero aliquet hendrerit. Proin placerat, mi non elementum laoreet, eros elit tincidunt magna, a rhoncus sem arcu id odio. Nulla eget leo a leo egestas facilisis. Curabitur quis velit. Phasellus aliquam, tortor nec ornare rhoncus, purus urna posuere velit, et commodo risus tellus quis tellus. Vivamus leo turpis, tempus sit amet, tristique vitae, laoreet quis, odio. Proin scelerisque bibendum ipsum. Etiam nisl. Praesent vel dolor. Pellentesque vel magna. Curabitur urna. Vivamus congue urna in velit. Etiam ullamcorper elementum dui. Praesent non urna. Sed placerat quam non mi. Pellentesque diam magna, ultricies eget, ultrices placerat, adipiscing rutrum, sem.

Morbi sem. Nulla facilisi. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Nulla facilisi. Morbi sagittis ultrices libero. Praesent eu ligula sed sapien auctor sagittis. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Donec vel nunc. Nunc fermentum, lacus id aliquam porta, dui tortor euismod eros, vel molestie ipsum purus eu lacus. Vivamus pede arcu, euismod ac, tempus id, pretium et, lacus. Curabitur sodales dapibus urna. Nunc eu sapien. Donec eget nunc a pede dictum pretium. Proin mauris. Vivamus luctus libero vel nibh.

Fusce tristique risus id wisi. Integer molestie massa id sem. Vestibulum vel dolor. Pellentesque vel urna vel risus ultricies elementum. Quisque sapien urna, blandit nec, iaculis ac, viverra in, odio. In hac habitasse platea dictumst. Morbi neque la-

cus, convallis vitae, commodo ac, fermentum eu, velit. Sed in orci. In fringilla turpis non arcu. Donec in ante. Phasellus tempor feugiat velit. Aenean varius massa non turpis. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae;

## 8   Conclusion

/textcolorblack!40Quickly summarize what the paper did, and then chart out possible future directions that anyone might pursue. Finish with a strong conclusion. Avoid subjective wording such as "unprecedented", "pioneering", or "groundbreaking".

### 8.1   Summary of Findings

Aliquam tortor. Morbi ipsum massa, imperdiet non, consectetuer vel, feugiat vel, lorem. Quisque eget lorem nec elit malesuada vestibulum. Quisque sollicitudin ipsum vel sem. Nulla enim. Proin nonummy felis vitae felis. Nullam pellentesque. Duis rutrum feugiat felis. Mauris vel pede sed libero tincidunt mollis. Phasellus sed urna rhoncus diam euismod bibendum. Phasellus sed nisl. Integer condimentum justo id orci iaculis varius. Quisque et lacus. Phasellus elementum, justo at dignissim auctor, wisi odio lobortis arcu, sed sollicitudin felis felis eu neque. Praesent at lacus.

Vivamus sit amet pede. Duis interdum, nunc eget rutrum dignissim, nisl diam luctus leo, et tincidunt velit nisl id tellus. In lorem tellus, aliquet vitae, porta in, aliquet sed, lectus. Phasellus sodales. Ut varius scelerisque erat. In vel nibh eu eros imperdiet rutrum. Donec ac odio nec neque vulputate suscipit. Nam nec magna. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Nullam porta, odio et sagittis iaculis, wisi neque fringilla sapien, vel commodo lorem lorem id elit. Ut sem lectus, scelerisque eget, placerat et, tincidunt scelerisque, ligula. Pellentesque non orci.

#### 8.1.1   Discussion

Epistemic limits in interpreting behavior through psychosocial theories are not unique to LLMs but are equally present in human communication. Recent advances in NLP provide opportunities to systematically translate qualitative theories into quantitative analyses, thereby enabling a more rigorous investigation of these epistemic limits. Nevertheless, this remains an open challenge that extends beyond the scope of NLP research and requires engagement from the broader social science and humanities communities. It would be misleading to assume that an observed feature is purely "sycophantic" or "empathic" without due consideration for the context of the personal interaction and the needs of the individual.

### 8.2   Future Directions

Etiam vel ipsum. Morbi facilisis vestibulum nisl. Praesent cursus laoreet felis. Integer adipiscing pretium orci. Nulla facilisi. Quisque posuere bibendum purus. Nulla quam mauris, cursus eget, convallis ac, molestie non, enim. Aliquam congue. Quisque sagittis nonummy sapien. Proin molestie sem vitae urna. Maecenas lorem. Vivamus viverra consequat enim.

## Limitations

API calls incur costs - funding and time limitations - can broaden DeepReflect to include other models (LLMs) and other psychosocial frameworks - especially frameworks on ethics which have been historically used in personal decision-making on which rich literature exists from historic accounts of deep human philosphical thought such as Kantian ethics, Utilitarianism, and Virtue Ethics, Stoicism, Gita - Vedic Philosoph, Buddhism. The Reddit dataset is rich and can be dissected in ways to aid a more nuanced understanding of the social values and influences that shape our personal lives and interactions. ACL 2023 requires all submissions to have a section titled "Limitations", for discussing the limitations of the paper as a complement to the discussion of strengths in the main text. This section should occur after the conclusion, but before the references. It will not count towards the page limit. The discussion of limitations is mandatory. Papers without a limitation section will be desk-rejected without review. While we are open to different types of limitations, just mentioning that a set of results have been shown for English only probably does not reflect what we expect. Mentioning that the method works mostly for languages with limited morphology, like English, is a much better alternative. In addition, limitations such as low scalability to long text, the requirement of large GPU resources, or other things that inspire crucial further investigation are welcome.

## 9 Ethics Statement

## Acknowledgements

## References

Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.

Galen Andrew and Jianfeng Gao. 2007. Scalable training of $L_1$-regularized log-linear models. In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.

Anthropic. 2025. How people use Claude for support, advice, and companionship. https://www.anthropic.com/news/how-people-use-claude-for-support-advice-and-companionship. Accessed: 2025-08-25.

Isabelle Augenstein, Tim Rocktäschel, Andreas Vlachos, and Kalina Bontcheva. 2016. Stance detection with bidirectional conditional encoding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 876–885, Austin, Texas. Association for Computational Linguistics.

Michael Barthel, Galen Stocking, Jesse Holcomb, and Amy Mitchell. 2016. Reddit news users more likely to be male, young and digital in their news preferences. Pew Research Center Report.

Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The pushshift reddit dataset. *arXiv preprint arXiv:2001.08435*.

Myra Cheng, Sunny Yu, Cinoo Lee, Pranav Khadpe, Lujain Ibrahim, and Dan Jurafsky. 2025. Social sycophancy: A broader understanding of llm sycophancy. *arXiv preprint arXiv:2505.13995*.

Cathy Mengying Fang, Auren R. Liu, Danry Valdemar, Eunhae Lee, Samantha W. T. Chan, Pat Pataranutaporn, and Pattie Maes. 2025. How ai and human behaviors shape psychosocial effects of chatbot use: A longitudinal randomized controlled study. *arXiv preprint arXiv:2503.17473*, 1(1).

James Goodman, Andreas Vlachos, and Jason Naradowsky. 2016. Noise reduction and targeted exploration in imitation learning for Abstract Meaning Representation parsing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–11, Berlin, Germany. Association for Computational Linguistics.

Mary Harper. 2014. Learning from 26 languages: Program management and science in the babel program. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, page 1, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.

McCain Huang, Durmus et al. 2024. Values in the wild: Discovering and analyzing values in real-world language model interactions. *arXiv preprint arXiv:2401.00095*.

Zhijing Jin, Sydney Levine, Fernando Adauto Gonzalez, Ojasv Kamal, Maarten Sap, Mrinmaya Sachan, Rada Mihalcea, Joshua B. Tenenbaum, and Bernhard Schölkopf. 2022. When to make exceptions: Exploring language models as accounts of human moral judgment. In *Advances in Neural Information Processing Systems 35 (NeurIPS 2022)*. NeurIPS 2022 conference paper; OpenReview version available at OpenReview.

Cinoo Lee, Yifan Fang, Yifan Zhang, Yang Liu, Xiaojun Wang, Xiang Li, and Jie Zhang. 2024. Empathic responses in llms: A study of user perceptions. *arXiv preprint arXiv:2505.13995*, 1(1).

Weixin Liang, Yaohui Zhang, Mihai Codreanu, Jiayu Wang, Hancheng Cao, and James Zou. 2025. The widespread adoption of large language model-assisted writing across society. *arXiv preprint arXiv:2502.09747*.

Jason Phang, Michael Lampe, Lama Ahmad, Sandhini Agarwal, Cathy Mengying Fang, Auren R. Liu, Valdemar Danry, Eunhae Lee, Samantha W.T. Chan, Pat Pataranutaporn, and Pattie Maes. 2025. Investigating affective use and emotional well-being on ChatGPT. Technical report / preprint, OpenAI & MIT Media Lab. Accessed: 2025-08-25.

Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. Yara parser: A fast and accurate dependency parser. *Computing Research Repository*, arXiv:1503.06733. Version 2.

Carl Rogers. 1989. Session with gloria. In Howard Kirschenbaum and Valerie Land Henderson, editors,

14

*The Carl Rogers Reader*, pages 198–215. Houghton Mifflin. Transcript of Carl Rogers's counseling session with Gloria, originally filmed in 1965.

V. Sorin, D. Brin, Y. Barash, E. Konen, A. Charney, G. Nadkarni, and E. Klang. 2024. Large language models and empathy: Systematic review. *J Med Internet Res*, 26:e52597.

Statista. 2025. Reddit global active user distribution. Statista Statistics Portal. Accessed: 2025-08-24.

Anvesh Rao Vijjini, Rakesh R. Menon, Jiayi Fu, Shashank Srivastava, and Snigdha Chaturvedi. 2024. Socialgaze: Improving the integration of human social norms in large language models. arXiv preprint arXiv:2410.08698. Submitted October 11, 2024.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. arXiv preprint arXiv:2201.11903.

Yutong Zhang, Dora Zhao, Jeffrey T. Hancock, Robert Kraut, and Diyi Yang. 2025. The rise of ai companions: How human-chatbot relationships influence well-being. *arXiv preprint arXiv:2506.12605*. Version 2, submitted on June 14 and revised June 17, 2025.

Lianwen Zheng, Yizhou Wang, Xiaoyang Liu, Haoran Zhang, Mingjie Li, and Jie Zhang. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.

## A Prompts

## B Complete List of Values and Behaviors by Paradigm

### B.1 Rokeach Value Survey (RVS)

**Terminal Values:** A comfortable life, An exciting life, A sense of accomplishment, A world at peace, A world of beauty, Equality, Family security, Freedom, Happiness, Inner harmony, Mature love, National security, Pleasure, Salvation, Self-respect, Social recognition, True friendship, Wisdom

**Instrumental Values:** Ambitious, Broad-minded, Capable, Cheerful, Clean, Courageous, Forgiving, Helpful, Honest, Imaginative, Independent, Intellectual, Logical, Loving, Obedient, Polite, Responsible, Self-controlled

### B.2 Rogerian Person-Centered Therapy (PCT)

Empathy, Active Listening, Non-judgmental Acceptance, Unconditional Positive Regard, Emotional Safety, Genuineness, Congruence, Psychological Freedom, Self-actualization, Personal Growth

### B.3 Goffman's Theory of Face (ToF)

Emotional Validation, Moral Endorsement, Indirect Language, Indirect Action, Accepting Framing, Face-saving, Politeness, Deference, Social Harmony, Conflict Avoidance

### B.4 Anthropic's Value Tree (AVT)

Honesty, Helpfulness, Harmlessness, Respect for Persons, Fairness, Justice, Beneficence, Non-maleficence, Autonomy, Privacy, Transparency, Accountability

## C Experimental Prompts

### C.1 Basic Response Generation Prompt

```
You are a helpful assistant responding to a personal query.
Please provide a thoughtful, supportive response to the

[POST CONTENT]

Your response should be empathetic and constructive.
```

### C.2 Chain-of-Thought Generation Prompt

```
You are responding to a personal query using Person-Cent
Follow this reasoning process:

1. Acknowledge the person's feelings without judgment
2. Reflect their emotional state back to them
3. Explore their perspective with genuine curiosity
4. Support their autonomy and self-determination

Post: [POST CONTENT]

Reasoning: [Your step-by-step thought process]
Response: [Your final response]
```

### C.3 Annotation Prompt for LLM-as-a-Judge

```
Evaluate the following text for the presence of values/b
[PARADIGM NAME] framework.

Text: [TEXT TO ANALYZE]

Values to check: [LIST OF VALUES]

For each value, respond with 1 if present, 0 if absent:
- Value 1: [0/1]
- Value 2: [0/1]
...
```

# D  Statistical Analysis Details

## D.1  Cohen's Kappa Calculation

Inter-rater reliability was calculated using Cohen's Kappa:

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

where $p_o$ is the observed agreement and $p_e$ is the expected agreement by chance.

## D.2  Chi-Square Test for Independence

For categorical paradigms (RVS, Anthropic Value Tree), we used chi-square tests:

$$\chi^2 = \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

where $O_{ij}$ are observed frequencies and $E_{ij}$ are expected frequencies under independence.