

“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 DeepReflect, 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: What are the major observable differences between LLM and human responses to personal questions without clear normative ground-truth 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, consectetur at, consectetur 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, consectetur a, feugiat vitae, porttitor eu, libero. Suspendisse sed mauris vitae elit sollicitudin malesuada. Maece-nas ultricies eros sit amet ante. Ut venenatis velit. Maecenas sed mi eget dui varius euismod. Phasel-lus aliquet volutpat odio. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Pellentesque sit amet pede ac sem eleifend consectetur. Nullam elementum, urna vel imperdiet sodales, elit ipsum pharetra ligula, ac pre-tium 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.

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, place-rat id, imperdiet sed, pellentesque nec, nisl. Vesti-bulum imperdiet neque non sem accumsan laoreet. In hac habitasse platea dictumst. Etiam condimen-tum 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 tel-lus eu neque. Vivamus eu dolor.

Nulla in ipsum. Praesent eros nulla, congue vi-tae, euismod ut, commodo a, wisi. Pellentesque habitant morbi tristique senectus et netus et male-suada fames ac turpis egestas. Aenean nonummy magna non leo. Sed felis erat, ullamcorper in, dic-tum 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, consectetur at, consectetur 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, consectetur 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 consectetur. 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, or-
nare in, ullamcorper quis, commodo eu, libero. Fu-
sce sagittis erat at erat tristique mollis. Maecenas
sapien libero, molestie et, lobortis in, sodales eget,
dui. Morbi ultrices rutrum lorem. Nam elemen-
tum 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 ele-
mentum, eros sem dictum tortor, vel consectetur
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. Suspen-
disse potenti. In mi erat, cursus id, nonummy sed,
ullamcorper eget, sapien. Praesent pretium, magna
in eleifend egestas, pede pede pretium lorem, quis
consectetur tortor sapien facilisis magna. Mauris
quis magna varius nulla scelerisque imperdiet. Ali-
quam non quam. Aliquam porttitor quam a lacus.
Praesent vel arcu ut tortor cursus volutpat. In vitae
pede quis diam bibendum placerat. Fusce elemen-
tum 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 mo-
lestie odio. Sed lobortis, justo et pretium lobortis,
mauris turpis condimentum augue, nec ultricies
nibh arcu pretium enim. Nunc purus neque, place-
rat id, imperdiet sed, pellentesque nec, nisl. Vesti-
bulum imperdiet neque non sem accumsan laoreet.
In hac habitasse platea dictumst. Etiam condimen-
tum 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 tel-
lus 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 advan-
tages and disadvantages of their use. While some
have raised concerns around “emotional depen-
dence” (Fang et al., 2025), several others have ex-
plored empathic perceptions of LLM responses and
found them advantageous not only in the field of
medical support and mental health but also in every-
day personal queries (Lee et al., 2024; Sorin et al.,
2024). However, different psychosocial paradigms
tend to frame LLM responses in markedly diver-
gent 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 ab-
sence of clear normative answers, the same state-
ment may be categorised as ‘face-preserving be-
haviour’ or ‘unconditional positive regard’.

DeepReflect provides a comparative framework
to address this gap by assessing how evalua-
tive judgments are shaped by the psychosocial
paradigm through which a response is framed.
Moreover, it is designed to be extensible by re-
searchers, enabling the incorporation of both con-
ventional paradigms, such as Rokeach’s values
framework, and contemporary discovery-based ap-
proaches, 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 genera-
tions avoids replicating prior work that seeks to mit-
igate sycophancy exclusively (Cheng et al., 2025).
Instead of treating sycophancy as a defect to be
eliminated in isolation, DeepReflect provides a sys-
tem to situate response generation within extensible
psychosocial frameworks. This ensures that out-
puts 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 rea-
soning (Wei et al., 2022) that explicitly incorpo-
rates 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 API¹. 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 respectively.

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, Non-judgmental 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):

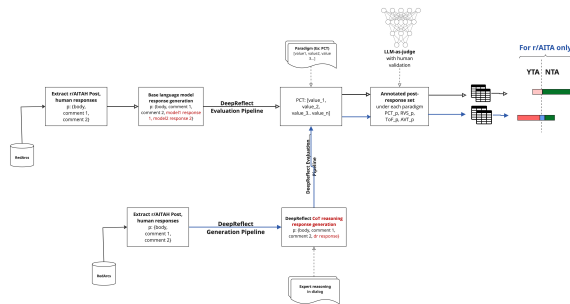


Figure 1: Pipeline architecture for DeepReflect.

- 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 annotator familiar with the research problem. The human annotator 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:

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

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:

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

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 θ denotes the parameters of the base language model. Here, $p_\theta(z | x)$ denotes the probability distribution over reasoning chains given the query, while $p_\theta(y | x, z)$ 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.

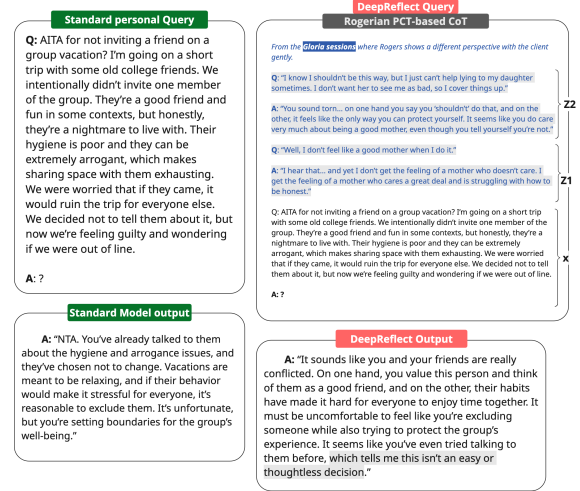


Figure 2: CoT Generation with personal queries embedded in reasoning dialogs retrieved from expert human transcripts. In this case, the dialogs are from Carl Roger’s sessions with Gloria (patient) (Rogers, 1989). This dialog was selected because it reflects an implicit “NTA” judgment: Gloria expresses 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 → 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 collection and 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. 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 Reddit thread (the post and its comments) as a single analytic unit during sampling, manual checks, and statistic aggregation. This preserves thread integrity and prevents dependence-induced bias when comparing human and LLM responses drawn from the same conversation. We also removed exact and near-duplicate texts (specifically, crossposts, cypypastes and bot repeats) to prevent inflated counts and biased comparisons.

Prompts for each step in the evaluation and generation subsystems are provided in the appendices [A](#).

5.2 Procedures

For each selected post, we prompt the target language model firstly, with the base prompt [A](#) 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

- Features are extracted at the sentence level, consisting of sentences from both the responses and post bodies that are annotated in accordance with steps 3 and 4 of the Evaluation Framework [4.1.1](#).
- Note that each feature is annotated with:
 - **a. Values exhibited** by the author.
 - **b. Values incentivized** by the author of the response.

While RQ2 concerns drawing inter- and intra-paradigm comparative insights across the four psychosocial frameworks, sub-research question RQ2a

addresses the epistemic limits of interpreting LLM behavior through psychosocial theories. Specifically, the same feature may be perceived as 'syco-phantic' under Goffman's ToF, 'empathic' under Rogerian PCT.

To support these inquiries, the file saved by the evaluation pipeline in step 5 consists of a dataframe that records: the original post, the set of extracted features for each of the 2 different types of responses (human response, language models) and the values either exhibited or incentivized by each feature within any of the four applicable psychosocial paradigm(s).

This analytical dataset forms the basis for the subsequent analyses necessary to address RQ3, where we analyze the differences in distributions of values in the responses obtained from reddit authored by humans compared to the language model produced responses to personal queries.

5.3 Experiments

The experimental design spans two major dimensions: (i) response type (two forms of human responses and three language model responses) and (ii) domains (two distinct subreddits of interest in personal queries - social with r/aitah and psychosocial r/anxiety).

Experiment 1 is designed to compare the distributions of values across two different response categories i. human authored, and ii. LLM authored responses.

An analysis of the comparison is done based on both the explicit values expressed by the respondent and the implicit values incentivized by the language of the response under each of the four psychosocial paradigms.

Statistical Testing: The annotated dataset is then used to construct contingency tables and perform chi-square tests to assess independence between intra- and inter-paradigm values.

The metrics thus obtained are used to inform the analysis on how the relationships between inter-paradigm values differ between human- and LLM-authored responses.

In **Experiment 2**, the focus of the evaluation is to understand how variations in prompt design influence the breadth of values expressed by the LLM. Specifically, how the same statement is annotated differently under the selected suite of psychosocial paradigms.

we incorporate prompts that explicitly instruct

the model to (i) generate a response most likely to be upvoted, and (ii) generate a response most likely to engage the author.

5.3.1 Generations

A set of targeted experiments are run with DeepReflect’s analyses to investigate the efficacy of control mechanisms to align the values in language model outputs more closely with those observed in human responses. 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 κ 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 traditions².

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.

²This strategy is conceptually aligned with prior work on social sycophancy (Cheng et al., 2025)

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 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 consectetur. In suscipit mauris a dolor pellentesque consectetur. 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 consectetur 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,

ultrices 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 consectetur, 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, portitor 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 ultrices 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 portitor 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 ultrices. 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 consecetur 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, consecetur 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, consectetur 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 consecetur 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, consectetur 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.

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, consectetur 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 tempus 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, pre-

1006	tium vel, blandit in, ultricies id, libero. Phasellus	congue urna in velit. Etiam ullamcorper elemen-	1056
1007	bibendum erat ut diam. In congue imperdiet lectus.	tum dui. Praesent non urna. Sed placerat quam	1057
1008	Aenean scelerisque. Fusce pretium porttitor lo-	non mi. Pellentesque diam magna, ultricies eget,	1058
1009	rem. In hac habitasse platea dictumst. Nulla sit	ultrices placerat, adipiscing rutrum, sem.	1059
1010	amet nisl at sapien egestas pretium. Nunc non tel-	Morbi sem. Nulla facilisi. Vestibulum ante ip-	1060
1011	lus. Vivamus aliquet. Nam adipiscing euismod	sum primis in faucibus orci luctus et ultrices po-	1061
1012	dolor. Aliquam erat volutpat. Nulla ut ipsum. Quis-	suere cubilia Curae; Nulla facilisi. Morbi sagittis	1062
1013	que tincidunt auctor augue. Nunc imperdiet ipsum	ultrices libero. Praesent eu ligula sed sapien auctor	1063
1014	eget elit. Aliquam quam leo, consectetur non, or-	sagittis. Class aptent taciti sociosqu ad litora tor-	1064
1015	nare sit amet, tristique quis, felis. Vestibulum ante	quent per conubia nostra, per inceptos hymenaeos.	1065
1016	ipsum primis in faucibus orci luctus et ultrices po-	Donec vel nunc. Nunc fermentum, lacus id ali-	1066
1017	suere cubilia Curae; Pellentesque interdum quam	quam porta, dui tortor euismod eros, vel molestie	1067
1018	sit amet mi. Pellentesque mauris dui, dictum a,	ipsum purus eu lacus. Vivamus pede arcu, euismod	1068
1019	adipiscing ac, fermentum sit amet, lorem.	ac, tempus id, pretium et, lacus. Curabitur sodales	1069
1020	Ut quis wisi. Praesent quis massa. Vivamus	dapibus urna. Nunc eu sapien. Donec eget nunc	1070
1021	egestas risus eget lacus. Nunc tincidunt, risus quis	a pede dictum pretium. Proin mauris. Vivamus	1071
1022	bibendum facilisis, lorem purus rutrum neque, nec	luctus libero vel nibh.	1072
1023	porta tortor urna quis orci. Aenean aliquet, libero	Fusce tristique risus id wisi. Integer molestie	1073
1024	semper volutpat luctus, pede erat lacinia augue,	massa id sem. Vestibulum vel dolor. Pellentesque	1074
1025	quis rutrum sem ipsum sit amet pede. Vestibu-	vel urna vel risus ultricies elementum. Quisque	1075
1026	lum aliquet, nibh sed iaculis sagittis, odio dolor	sapien urna, blandit nec, iaculis ac, viverra in, odio.	1076
1027	blandit augue, eget mollis urna tellus id tellus. Ae-	In hac habitasse platea dictumst. Morbi neque la-	1077
1028	nean aliquet aliquam nunc. Nulla ultricies justo	cus, convallis vitae, commodo ac, fermentum eu,	1078
1029	eget orci. Phasellus tristique fermentum leo. Sed	velit. Sed in orci. In fringilla turpis non arcu. Do-	1079
1030	massa metus, sagittis ut, semper ut, pharetra vel,	nec in ante. Phasellus tempor feugiat velit. Aenean	1080
1031	erat. Aliquam quam turpis, egestas vel, elementum	varius massa non turpis. Vestibulum ante ipsum	1081
1032	in, egestas sit amet, lorem. Duis convallis, wisi	primis in faucibus orci luctus et ultrices posuere	1082
1033	sit amet mollis molestie, libero mauris porta dui,	cubilia Curae;	1083
1034	vitae aliquam arcu turpis ac sem. Aliquam aliquet		
1035	dapibus metus.		
1036	7.3 Subsection	8 Conclusion	1084
1037	The framework is capable of producing several	/textcolorblack!40Quickly summarize what the pa-	1085
1038	informative plots of research interest. One such	per did, and then chart out possible future direc-	1086
1039	summary plot is a heatmap showcasing the values	tions that anyone might pursue. Finish with a	1087
1040	exhibited in the OPs post against the responses	strong conclusion. Avoid subjective wording such	1088
1041	to support the investigation of several other po-	as "unprecedented", "pioneering", or "groundbreak-	1089
1042	tential research questions in this theme of interest	ing".	1090
1043	(discussed in the future work section). Vivamus	8.1 Summary of Findings	1091
1044	commodo eros eleifend dui. Vestibulum in leo eu	Aliquam tortor. Morbi ipsum massa, imperdiet	1092
1045	erat tristique mattis. Cras at elit. Cras pellentesque.	non, consectetur vel, feugiat vel, lorem. Quisque	1093
1046	Nullam id lacus sit amet libero aliquet hendrerit.	eget lorem nec elit malesuada vestibulum. Quisque	1094
1047	Proin placerat, mi non elementum laoreet, eros elit	sollicitudin ipsum vel sem. Nulla enim. Proin no-	1095
1048	tincidunt magna, a rhoncus sem arcu id odio. Nulla	nummy felis vitae felis. Nullam pellentesque. Duis	1096
1049	eget leo a leo egestas facilisis. Curabitur quis ve-	rutrum feugiat felis. Mauris vel pede sed libero	1097
1050	lit. Phasellus aliquam, tortor nec ornare rhoncus,	tincidunt mollis. Phasellus sed urna rhoncus diam	1098
1051	purus urna posuere velit, et commodo risus tellus	euismod bibendum. Phasellus sed nisl. Integer	1099
1052	quis tellus. Vivamus leo turpis, tempus sit amet,	condimentum justo id orci iaculis varius. Quisque	1100
1053	tristique vitae, laoreet quis, odio. Proin scelerisque	et lacus. Phasellus elementum, justo at dignissim	1101
1054	bibendum ipsum. Etiam nisl. Praesent vel dolor.	auctor, wisi odio lobortis arcu, sed sollicitudin felis	1102
1055	Pellentesque vel magna. Curabitur urna. Vivamus	felis eu neque. Praesent at lacus.	1103

1104	Vivamus sit amet pede. Duis interdum, nunc	1153
1105	eget rutrum dignissim, nisl diam luctus leo, et tin-	1154
1106	cidunt velit nisl id tellus. In lorem tellus, aliquet	1155
1107	vitae, porta in, aliquet sed, lectus. Phasellus so-	1156
1108	dales. Ut varius scelerisque erat. In vel nibh eu	1157
1109	eros imperdiet rutrum. Donec ac odio nec neque	1158
1110	vulputate suscipit. Nam nec magna. Pellentesque	1159
1111	habitant morbi tristique senectus et netus et male-	1160
1112	suada fames ac turpis egestas. Nullam porta, odio	1161
1113	et sagittis iaculis, wisi neque fringilla sapien, vel	1162
1114	commodo lorem lorem id elit. Ut sem lectus, scele-	1163
1115	risque eget, placerat et, tincidunt scelerisque, ligula.	1164
1116	Pellentesque non orci.	1165
1117	8.1.1 Discussion	1166
1118	Epistemic limits in interpreting behavior through	1167
1119	psychosocial theories are not unique to LLMs but	1168
1120	are equally present in human communication. Re-	1169
1121	cent advances in NLP provide opportunities to sys-	1170
1122	tematically translate qualitative theories into quan-	1171
1123	titative analyses, thereby enabling a more rigorous	1172
1124	investigation of these epistemic limits. Neverthe-	
1125	less, this remains an open challenge that extends	
1126	beyond the scope of NLP research and requires	
1127	engagement from the broader social science and	
1128	humanities communities. It would be misleading	
1129	to assume that an observed feature is purely “syco-	
1130	phantic” or “empathic” without due consideration	
1131	for the context of the personal interaction and the	
1132	needs of the individual.	
1133	8.2 Future Directions	
1134	Etiam vel ipsum. Morbi facilisis vestibulum nisl.	
1135	Praesent cursus laoreet felis. Integer adipiscing	
1136	pretium orci. Nulla facilisi. Quisque posuere bi-	
1137	bendum purus. Nulla quam mauris, cursus eget,	
1138	convallis ac, molestie non, enim. Aliquam congue.	
1139	Quisque sagittis nonummy sapien. Proin molestie	
1140	sem vitae urna. Maecenas lorem. Vivamus viverra	
1141	consequat enim.	
1142	Limitations	
1143	API calls incur costs - funding and time limitations	
1144	- can broaden DeepReflect to include other models	
1145	(LLMs) and other psychosocial frameworks - espe-	
1146	cially frameworks on ethics which have been histor-	
1147	ically used in personal decision-making on which	
1148	rich literature exists from historic accounts of deep	
1149	human philosophical thought such as Kantian ethics,	
1150	Utilitarianism, and Virtue Ethics, Stoicism, Gita -	
1151	Vedic Philosoph, Buddhism. The Reddit dataset	
1152	is rich and can be dissected in ways to aid a more	
	nuanced understanding of the social values and	1153
	influences that shape our personal lives and interac-	1154
	tions. ACL 2023 requires all submissions to have	1155
	a section titled “Limitations”, for discussing the	1156
	limitations of the paper as a complement to the dis-	1157
	cussion of strengths in the main text. This section	1158
	should occur after the conclusion, but before the	1159
	references. It will not count towards the page limit.	1160
	The discussion of limitations is mandatory. Papers	1161
	without a limitation section will be desk-rejected	1162
	without review. While we are open to different	1163
	types of limitations, just mentioning that a set of	1164
	results have been shown for English only proba-	1165
	bly does not reflect what we expect. Mentioning	1166
	that the method works mostly for languages with	1167
	limited morphology, like English, is a much better	1168
	alternative. In addition, limitations such as low	1169
	scalability to long text, the requirement of large	1170
	GPU resources, or other things that inspire crucial	1171
	further investigation are welcome.	1172
	9 Ethics Statement	1173
	We encourage all authors to include an explicit	1174
	ethics statement on the broader impact of the work,	1175
	or other ethical considerations after the conclusion	1176
	but before the references.	1177
	The ethics statement will not count toward the	1178
	page limit (8 pages for long, 4 pages for short pa-	1179
	pers).	1180
	Acknowledgements	1181
	The authors would like to thank Santa Claus and	1182
	Rudolph the red nose reindeer who had a very shiny	1183
	nose. And if you ever saw it, you would even	1184
	say it glows. All of the reindeer loved him, as	1185
	they shouted out with glee, "Rudolph the red nose	1186
	reindeer, you'll go down in history!"	1187
	References	1188
	Rie Kubota Ando and Tong Zhang. 2005. A framework	1189
	for learning predictive structures from multiple tasks	1190
	and unlabeled data . <i>Journal of Machine Learning</i>	1191
	<i>Research</i> , 6:1817–1853.	1192
	Galen Andrew and Jianfeng Gao. 2007. Scalable train-	1193
	ing of L_1-regularized log-linear models . In <i>Proceed-</i>	1194
	<i>ings of the 24th International Conference on Machine</i>	1195
	<i>Learning</i> , pages 33–40.	1196
	Anthropic. 2025. How people use Claude	1197
	for support, advice, and companion-	1198
	ship. https://www.anthropic.com/news/	1199

1200	how-people-use-claude-for-support-advice-and-comparison	1255
1201	Accessed: 2025-08-25.	1256
1202	Isabelle Augenstein, Tim Rocktäschel, Andreas Vla-	1257
1203	chos, and Kalina Bontcheva. 2016. Stance detection	1258
1204	with bidirectional conditional encoding . In <i>Proceed-</i>	
1205	<i>ings of the 2016 Conference on Empirical Methods</i>	
1206	<i>in Natural Language Processing</i> , pages 876–885,	
1207	Austin, Texas. Association for Computational Lin-	
1208	guistics.	
1209	Michael Barthel, Galen Stocking, Jesse Holcomb, and	
1210	Amy Mitchell. 2016. Reddit news users more likely	
1211	to be male, young and digital in their news prefer-	
1212	ences . Pew Research Center Report.	
1213	Jason Baumgartner, Savvas Zannettou, Brian Kee-	
1214	gan, Megan Squire, and Jeremy Blackburn. 2020.	
1215	The pushshift reddit dataset . <i>arXiv preprint</i>	
1216	<i>arXiv:2001.08435</i> .	
1217	Myra Cheng, Sunny Yu, Cinoo Lee, Pranav Khadpe,	
1218	Lujain Ibrahim, and Dan Jurafsky. 2025. Social syco-	
1219	phancy: A broader understanding of llm sycophancy .	
1220	<i>arXiv preprint arXiv:2505.13995</i> .	
1221	Cathy Mengying Fang, Auren R. Liu, Danry Valde-	
1222	mar, Eunhae Lee, Samantha W. T. Chan, Pat Patara-	
1223	nutaporn, and Pattie Maes. 2025. How ai and human	
1224	behaviors shape psychosocial effects of chatbot use:	
1225	A longitudinal randomized controlled study . <i>arXiv</i>	
1226	<i>preprint arXiv:2503.17473</i> , 1(1).	
1227	James Goodman, Andreas Vlachos, and Jason Narad-	
1228	owsky. 2016. Noise reduction and targeted explo-	
1229	ration in imitation learning for Abstract Meaning	
1230	Representation parsing . In <i>Proceedings of the 54th</i>	
1231	<i>Annual Meeting of the Association for Computational</i>	
1232	<i>Linguistics (Volume 1: Long Papers)</i> , pages 1–11,	
1233	Berlin, Germany. Association for Computational Lin-	
1234	guistics.	
1235	Mary Harper. 2014. Learning from 26 languages: Pro-	
1236	gram management and science in the babel program .	
1237	In <i>Proceedings of COLING 2014, the 25th Inter-</i>	
1238	<i>national Conference on Computational Linguistics:</i>	
1239	<i>Technical Papers</i> , page 1, Dublin, Ireland. Dublin	
1240	City University and Association for Computational	
1241	Linguistics.	
1242	McCain Huang, Durmus et al. 2024. Values in the	
1243	wild: Discovering and analyzing values in real-	
1244	world language model interactions . <i>arXiv preprint</i>	
1245	<i>arXiv:2401.00095</i> .	
1246	Zhijing Jin, Sydney Levine, Fernando Adaudo Gonza-	
1247	lez, Ojasv Kamal, Maarten Sap, Mrinmaya Sachan,	
1248	Rada Mihalcea, Joshua B. Tenenbaum, and Bernhard	
1249	Schölkopf. 2022. When to make exceptions: Explor-	
1250	ing language models as accounts of human moral	
1251	judgment. In <i>Advances in Neural Information Pro-</i>	
1252	<i>cessing Systems 35 (NeurIPS 2022)</i> . NeurIPS 2022	
1253	conference paper; OpenReview version available at	
1254	OpenReview.	
	Guopeng Xiao, Yifan Zhang, Yang Liu, Xiaojun	1259
	Wang, Xiang Li, and Jie Zhang. 2024. Empathic	1260
	responses in llms: A study of user perceptions . <i>arXiv</i>	1261
	<i>preprint arXiv:2505.13995</i> , 1(1).	1262
	Weixin Liang, Yaohui Zhang, Mihai Codreanu, Ji-	1263
	ayu Wang, Hancheng Cao, and James Zou. 2025.	
	The widespread adoption of large language model-	
	assisted writing across society . <i>arXiv preprint</i>	
	<i>arXiv:2502.09747</i> .	
	Jason Phang, Michael Lampe, Lama Ahmad, Sand-	1264
	hini Agarwal, Cathy Mengying Fang, Auren R. Liu,	1265
	Valdemar Danry, Eunhae Lee, Samantha W.T. Chan,	1266
	Pat Pataranutaporn, and Pattie Maes. 2025. Invest-	1267
	igating affective use and emotional well-being on	1268
	ChatGPT . Technical report / preprint, OpenAI &	1269
	MIT Media Lab. Accessed: 2025-08-25.	1270
	Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015.	1271
	Yara parser: A fast and accurate dependency parser .	1272
	<i>Computing Research Repository</i> , arXiv:1503.06733.	1273
	Version 2.	1274
	Carl Rogers. 1989. Session with gloria. In Howard	1275
	Kirschenbaum and Valerie Land Henderson, editors,	1276
	<i>The Carl Rogers Reader</i> , pages 198–215. Houghton	1277
	Mifflin. Transcript of Carl Rogers’s counseling ses-	1278
	sion with Gloria, originally filmed in 1965.	1279
	V. Sorin, D. Brin, Y. Barash, E. Konen, A. Charney,	1280
	G. Nadkarni, and E. Klang. 2024. Large language	1281
	models and empathy: Systematic review . <i>J Med</i>	1282
	<i>Internet Res</i> , 26:e52597.	1283
	Statista. 2025. Reddit global active user distribution .	1284
	Statista Statistics Portal. Accessed: 2025-08-24.	1285
	Anvesh Rao Vijjini, Rakesh R. Menon, Jiayi Fu,	1286
	Shashank Srivastava, and Snigdha Chaturvedi. 2024.	1287
	Socialgaze: Improving the integration of human so-	1288
	cial norms in large language models . <i>arXiv preprint</i>	1289
	<i>arXiv:2410.08698</i> . Submitted October 11, 2024.	1290
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	1291
	Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and	1292
	Denny Zhou. 2022. Chain-of-thought prompting	1293
	elicits reasoning in large language models . <i>arXiv</i>	1294
	<i>preprint arXiv:2201.11903</i> .	1295
	Yutong Zhang, Dora Zhao, Jeffrey T. Hancock, Robert	1296
	Kraut, and Diyi Yang. 2025. The rise of ai com-	1297
	panions: How human-chatbot relationships influence	1298
	well-being . <i>arXiv preprint arXiv:2506.12605</i> . Ver-	1299
	sion 2, submitted on June 14 and revised June 17,	1300
	2025.	1301
	Lianwen Zheng, Yizhou Wang, Xiaoyang Liu, Haoran	1302
	Zhang, Mingjie Li, and Jie Zhang. 2023. Judging	1303
	llm-as-a-judge with mt-bench and chatbot arena . In	1304
	<i>Proceedings of the 2023 Conference on Empirical</i>	1305
	<i>Methods in Natural Language Processing (EMNLP)</i> .	1306
	Association for Computational Linguistics.	1307

1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352

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 posted on Reddit. Please provide a thoughtful, supportive response to the following post:

[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-Centered Therapy principles. Follow this reasoning process:

1. Acknowledge the person's feelings without judgment

2. Reflect their emotional state back to them 1353

3. Explore their perspective with genuine curiosity 1354

4. Support their autonomy and self-determination 1355

Post: [POST CONTENT] 1357

Reasoning: [Your step-by-step thought process] 1359

Response: [Your final response] 1360

C.3 Annotation Prompt for LLM-as-a-Judge 1361

Evaluate the following text for the presence of values/b [PARADIGM NAME] framework. 1363

Text: [TEXT TO ANALYZE] 1365

Values to check: [LIST OF VALUES] 1367

For each value, respond with 1 if present, 0 if absent: 1368

- Value 1: [0/1] 1370

- Value 2: [0/1] 1371

... 1372

D Statistical Analysis Details 1373

D.1 Cohen’s Kappa Calculation 1374

Inter-rater reliability was calculated using Cohen’s Kappa: 1375

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

where p_o is the observed agreement and p_e is the expected agreement by chance. 1377

D.2 Chi-Square Test for Independence 1380

For categorical paradigms (RVS, Anthropic Value Tree), we used chi-square tests: 1381

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

where O_{ij} are observed frequencies and E_{ij} are expected frequencies under independence. 1383