

---

# TALK, LISTEN, CONNECT: HOW HUMANS AND AI EVALUATE EMPATHY IN RESPONSES TO EMOTIONALLY CHARGED NARRATIVES

---

**Mahnaz Roshanaei**  
Department of Communication  
Stanford University

**Rezvaneh Rezapour**  
School of Computing and Informatics  
Drexel University

**Magy Seif El-Nasr**  
Department of Computational Media  
University of California, Santa Cruz

October 28, 2025

## ABSTRACT

Social interactions promote well-being, yet barriers like geographic distance, time limitations, and mental health conditions can limit face-to-face interactions. Emotionally responsive AI systems, such as chatbots, offer new opportunities for social and emotional support, but raise critical questions about how empathy is perceived and experienced in human-AI interactions. This study examines how empathy is evaluated in AI-generated versus human responses. Using personal narratives, we explored how persona attributes (e.g., gender, empathic traits, shared experiences) and story qualities affect empathy ratings. We compared responses from standard and fine-tuned AI models with human judgments. Results show that while humans are highly sensitive to emotional vividness and shared experience, AI-responses are less influenced by these cues, often lack nuance in empathic expression. These findings highlight challenges in designing emotionally intelligent systems that respond meaningfully across diverse users and contexts, and informs the design of ethically aware tools to support social connection and well-being.

**Keywords** Human-AI Interactions · LLMs · empathy · well-being · mental health

## 1 Introduction

Research consistently demonstrates that a rich social life with support networks and engaging in high-quality social interactions, in particular, face-to-face interactions, are associated with a variety of benefits to people’s well-being [1, 2, 3, 4]. Engaging in social interactions that are meaningful and substantive, such as those involving self-disclosure or emotional depth, have been linked with greater happiness, life satisfaction, and social connectedness [5, 6, 7, 8]. Face-to-face and hybrid interactions consistently outperform digital-only exchanges in promoting positive affect [10]. Despite all the benefits of face-to-face interactions, barriers such as geographical distance, time constraints, health challenges, social anxiety and loneliness often limit these engagements [11, 12, 13]. In response to these barriers, AI-driven chatbots have emerged as supplementary tools for facilitating social interaction and offering non-judgmental and accessible support for social, emotional, and relational needs [14], though not as replacements for human contact [15, 16].

As chatbots become more fluent and emotionally responsive [17, 18], questions are raised about their psychological impact, especially during emotionally complex exchanges. Some worry these tools may alter our understanding of intimacy or authenticity in social interactions, potentially reshaping the concept of connectedness in a digital era. Studies on emotionally responsive AI like Replika reveal diverse user reactions; some form stronger emotional attachments and deeper bond [19, 20, 21], while others report discomfort and fear with overly human-like behaviors [22]. This variability in emotional responses highlights a broader challenge within Human-Computer Interaction (HCI) research: understanding how systems can effectively evaluate or evoke empathy without causing discomfort or emotional dissonance. Empathy in AI presents a delicate balance, fostering emotional engagement while preserving authenticity. As emotional stakes rise, the difficulty of assessing and quantifying empathy in these systems grows, necessitating a

more nuanced approach to evaluation [23, 24]. While human-centered design recognizes empathy as vital, measuring it within AI systems remains a major challenge due to its subjective, context-dependent nature [25, 26, 27]. To address this gap, this paper investigates how empathy is evaluated in AI versus humans and explores the factors that evoke empathetic responses in both. Through experimental methods using personal narratives, we examine how persona attributes, shared experiences, and model fine-tuning influence AI’s alignment with human empathy judgments. We also analyze how narrative qualities, such as emotional vividness and shared experiences and moral values, influence empathy evoked in AI and human responses.

## 1.1 Empathy in the Psychological Literature

The definition of empathy is multifaceted and encompasses a range of interpretations that highlight its significance in various fields such as healthcare, psychoanalysis, and interpersonal relations. Empathy is fundamentally described as the ability to understand and share the feelings of another, which is considered a crucial interpersonal skill [28]. It involves both affective and cognitive components: emotional resonance with another’s experience, and the deliberate understanding and communication of that experience, commonly framed as “putting oneself in another’s shoes” [29, 30]. In a relational framework, empathy is seen as an interactive and dynamic process where both the empathizer and the empathee shape each other’s experiences, ultimately enhancing the quality of relationships [31]. Experimental psychology identifies key factors influencing this exchange, such as the intensity and vividness of expressed emotion [32], the perceived similarity between the empathizer and empathee [33, 34], and individual traits like gender, personality, and prior emotional experiences [33].

When individuals engage in empathic social interactions, where they feel understood, supported, and valued by others, they experience a range of positive outcomes that contribute to overall well-being such as maintaining relationships [35], and reducing stress, depression, and loneliness [1, 36]. This aligns with psychoanalytic and therapeutic perspectives, which view empathy not just as a skill but as a method of deep interpersonal engagement and healing. In mental health contexts, empathy fosters trust between patients and caregivers, enabling more effective support and treatment [37, 30]. According to [38], high-quality interactions occur when counselors focus on the client and show empathy, while low-quality interactions involve counselors giving instructions and the client merely complying. Empathetic communication is therefore essential in creating environments where individuals feel safe, heard, and emotionally validated—foundational for effective therapeutic and physician-patient relationships.

## 1.2 Empathy in the Era of AI

Within the HCI Community, empathy is recognized as a fundamental component of interpersonal and communication competence that augments understanding, prediction, persuasion, compliance gaining, relational development, and counseling among individuals [39]. This view underscores the role of empathy not only in fostering effective human-to-human interactions but also in shaping more emotionally intelligent human-computer interactions. Recent advances in affective computing have enabled technologies to simulate emotional expressions and respond to users’ affective states, paving the way for systems that can exhibit forms of empathy.

In the human-AI interaction domain, the role of AI in social and emotional contexts has been explored, making them suitable for applications such as social chatbots, where building rapport and sustaining engagement are essential. Some studies have indicated that empathetic AI agents are effective at fostering social connections, encouraging self-disclosure, facilitating social interactions, and maintaining user engagement [40, 41, 17, 14]. However, their ability to fully replicate the nuanced empathy of humans remains limited. In therapeutic contexts, AI-driven counselors can offer consistent and accessible support but often fall short in emotional depth which is central to human empathy. Patients frequently express a preference for human-operated counselors, an agent controlled by AI, finding AI-driven empathy less helpful and sometimes counterproductive. However, providing attentive comments and offering hope have been shown to improve the perceived quality of AI-driven counseling, especially in emotionally charged contexts [42]. Furthermore, some studies have shown that AI agents can show biased value judgments and uneven expressions of empathy toward different demographic groups [23, 43]. These findings highlight the challenges AI faces in replicating human empathy, particularly given the complexity of empathy in emotionally sensitive and therapeutic settings, where shallow or biased responses can negatively affect therapeutic outcomes and .

To better understand these shortcomings, it is important to examine how empathy is evaluated and perceived differently in AI and human interactions. Research on human-to-human interaction, such as online peer support communities that rely heavily on peer-based emotional support, indicated that highly empathetic responses are often rare [23]. Through a mixed-methods analysis, [44] found that techniques like active listening and reflective restatements increased perceived empathy, while rigid structures and a lack of emotional validation diminished it. These findings suggest that

effective empathy, either in human or AI, requires intentional strategies, not merely access to emotional content or communication channels.

Emerging technologies like large language models (LLMs) have opened new opportunities for simulating emotional processes more effectively. Deployed in social chatbots and mental health platforms, these systems show promise in boosting engagement and user satisfaction [14, 45]. Yet their long-term psychological impact and ethical implications remain underexplored, especially in intimate, support-focused roles. Empathy is also increasingly central in human-centered technology design, where understanding users’ emotions and lived experiences guides product development [46, 47, 48]. This has led to empathetic design frameworks that use storytelling, cultural probes, and other immersive methods to help designers better see from users’ perspectives [49].

In virtual interactions with LLMs and conversational agents (CAs), the ability of virtual agents to appear more likable, trustworthy, and caring, underscores the substantial role of empathy in transforming the quality of human-AI interactions. Prior research has shown that empathetic characteristics can be modeled and embodied in virtual agents. For example, [50] proposed a computational model of empathy, showing how the perception and understanding of others’ emotional states can be algorithmically modeled and embedded within AI systems, and offering rich insights into the operationalization of empathy.

Despite these advances, a major challenge remains: how do we measure perceived empathy in technology? Unlike human empathy, which has well-established scales [51, 52], the field lacks validated tools for evaluating empathy in AI systems. To address this, [53] introduced the Perceived Empathy of Technology Scale (PETS)—a 10-item, 2-factor instrument that measures how empathetically users perceive interactive systems. This scale provides researchers and designers with a framework for evaluating the empathy exhibited by systems such as CAs and social robots [53]. While CAs can simulate empathy [17, 54], they often fall short in truly interpreting users’ emotional experiences, reinforcing the need for more nuanced, user-informed approaches to empathetic system design [23].

In light of these findings, the current study aims to deepen our understanding of how AI-generated empathetic responses differ from human judgments of empathy, and what factors influence the perception and evocation of empathy in both AI and humans. Our study is based on the previous research examining empathy through the heterogeneous effects of personal stories, which utilized a set of narrative stories in which individuals reflected on the three best and three worst events of their lives. These narratives were presented to Amazon MTurkers, who evaluated the extent of empathy they experienced and identified the key elements that evoked their empathic judgments. Accordingly, the following research questions formulated aim to address these critical gaps in the literature:

- RQ1. To what extent is empathy evaluated differently by humans and AI?
- RQ2. How do various persona attributes—such as (a) gender, (b) empathic personality traits (empathic concern, perspective-taking), and (c) shared experiences with storytellers— influence the empathy assessed by AI, compared to humans?
- RQ3. To what extent does fine-tuning AI models improve the alignment between AI-generated and human-evaluated empathic responses?
- RQ4. What key factors influence the evocation of empathy in AI-generated responses compared to human responses?

Using statistical analysis, we first compare how empathy is evaluated by AI, compared to humans. To make the AI-generated responses more human-like, we incorporate persona attributes into the prompts, such as gender, empathic personality traits, and similarity of experience with the storyteller, drawing from psychological literature [33]. To further improve model performance, we applied instruction fine-tuning to GPT-4o. This was done in two ways: (1) fine-tuning on human-annotated empathy ratings associated with the story narratives, and (2) fine-tuning with additional reader attributes including gender, empathic concern, perspective-taking, and perceived similarity to the storyteller. We then quantitatively analyze the factors that evoke empathy in humans versus AI responses, emphasizing the focus on story attributes and the role of shared experiences in shaping empathic reaction [33, 34]. Finally, we highlight the need for thoughtful consideration of the potential benefits and harms of empathetic AI, particularly different effects across diverse subgroups and its implications in sensitive domains like mental health. By critically examining both the promise and the limitations of AI-mediated empathy, this study can contribute to the development of the next generation of AI systems that are not only technically proficient but also emotionally attuned and ethically aligned with human values and needs.

So, for video number 4, I'm going to be talking about the first negative event that I wrote down.  
That when I was bullied in the seventh grade.  
There was girl, and her name was Carol.  
She didn't like me. I don't know why. She just didn't.  
Then, she ended up turning the whole class against me.  
Nobody talked to me.  
I wouldn't even want to go to school.  
I would wake up and be like I don't feel good, I don't want to go.  
.....

Figure 1: An example of the human-generated story used in the surveys

## 2 Method

### 2.1 Data

#### 2.1.1 Human-Generated Data

Our analysis is based on data collected through an online survey conducted via Amazon MTurk in the winter of 2019 [9]. This survey builds on previous IRB-approved research, where 756 videos of 126 undergraduate students were recorded. Participants were recruited using the Psychology Subject Pool and described the three best and three worst events of their lives. To ensure participant comfort with the use of their recorded videos in future studies, they have been asked to give consent for the videos in two steps. Upon arriving at the lab, participants completed a Pre-Video Recording consent form to provide consent to participate in the study. Then, after recording their videos, they filled out a Post-Video Recording consent form, explicitly indicating consent for the use of each recorded video separately. Later, each two-minute video was transcribed into text for analysis, with positive and negative labels assigned whether the story described one of the best or worst events of their life. In addition to these recorded videos, their demographic information (i.e., age, gender, race) and personality characteristics of each participant, referred to hereafter as storytellers (see Figure 1 for an example), were recorded.

After compiling the narrative stories, a second IRB-approved study was conducted where participants were recruited via Amazon Mechanical Turk (MTurk). Participants were based in the U.S., aged 18 and older who had completed at least a high school education. All participants received monetary compensation for their time. The final sample included 2,586 individuals, with each narrative annotated by an average of three raters. Of the MTurk participants, 56% identified as female, with an overall mean age of 38.6 years ( $SD = 12.58$ ); the average age among male participants was 36.56 years ( $SD = 12.24$ ).

During the survey, MTurk participants were asked to read a series of stories and respond to several questions using a 5-point Likert scale ranging from 1 (Not at all) to 5 (Extremely). These included assessments of: (1) overall empathy (e.g., To what extent did you feel empathy for the storyteller?), (2) the affective dimension of empathy, calculated as the average of responses to items assessing feelings of sympathy, compassion, and being moved, (3) the cognitive dimension of empathy, and (4) participants' reasons for the empathy they experienced toward each storyteller. A full list of survey items is provided in Supplementary Materials S.1.1. Given prior research highlighting gender differences in empathy [55, 56, 57], as well as the established links between empathy evocation and personality traits related to empathy [58, 59], collecting demographic information (e.g., age, gender) along with measures of empathic concern and perspective taking, measured by Interpersonal Reactivity Index (IRI) [58].

In addition to individual traits, several studies have examined how the perceived relationship between the storyteller and the reader influences empathic responses, particularly the role of perceived similarity. For example, [33] found that readers who perceived greater similarity to a storyteller in terms of personality or values exhibited stronger physiological responses (e.g., increased heart rate and sweating) when exposed to the storyteller's pain. To assess perceived shared experience in our study, two survey items were included: one measuring emotional similarity and another measuring similarity in the specific details of the experience. These items are described in full in Supplementary Materials S.1.1.

#### 2.1.2 LLM-Generated Data

To compare the levels of empathy perceived by AI versus humans, we used OpenAI's GPT-4o (gpt-4o-2024-08-06) [60] to assess the same human-generated stories that were previously analyzed by MTurk workers. To elicit comparable responses from the AI model, we developed a set of prompts based on the same questions posed to the MTurk workers (Supplementary Materials S.1.1), structured in the following format:

Prompt = Instruction + Question + Output Format

The example of a base prompt is added in Supplementary Materials S.1.2.

### 2.1.3 LLM-Generated Data with Persona

To answer RQ2, and evaluate the extent to which gender, IRI, and similarity of experience may influence the level of empathy AI experiences, we test four different treatment persona settings in total:

- AI-agent has gender.
- AI-agent has an empathic concern.
- AI-agent has perspective taking.
- AI-agent has an experience similarity with the storyteller.

To implement these personas in the prompt, we used the same gender, empathic concern, perspective taking, and the level of similarity of experience of MTurk workers and embedded them in the prompts with the following structure:

$$\text{Prompt} = \text{Instruction} + \text{Persona} + \text{Question} + \text{Output Format}$$

The example of a persona-based prompt is added in Supplementary Materials S.1.2.

### 2.1.4 Fine-tuning Experiment

We further explored methods of improving model performance using instruction fine-tuning to enable GPT-4o to perform more effectively. For the fine-tuning experiment, the prompt design was kept consistent with the methodology outlined above. Our data were stratified to create a balanced distribution across empathy scales, persona attributes (gender, empathic concern, perspective taking, and experience similarity), and positive and negative story tags, forming a candidate training set. Each story was treated as a single unit, with a composite categorical class constructed by concatenating the values of the stratified features. This approach ensures that stratification respects the joint distribution of these variables across training and test sets. Next, we counted the frequency of each composite class. Those with fewer than two occurrences were considered “rare” and handled separately, as they cannot be properly stratified. For non-rare classes, a stratified train-test split was applied using the composite class as the stratification criterion, and rare stories were split using a random (non-stratified) train-test split to ensure representation. Stories not assigned in either the stratified or rare splits were placed in the test set. Additionally, we checked stories in training and test sets, verified that no story appeared in both sets, and removed duplicates if necessary. Finally, we created a balanced instruction training set, comprising 80 stories with a similar distribution to our test set. Histograms of key categorical variables (e.g., tag, empathy, gender) were plotted to confirm that stratification preserved their distributions (see Supplementary materials Figure S.3). Additionally, Chi-square tests and L1-distance were used to quantify the similarity of distributions between the original and stratified datasets (see Supplementary materials Table S.1). Fine-tuning is conducted in two nested steps: first, by using the stories, that is, the stories paired with corresponding human-annotated empathy ratings (GPT-4o FT Story-only). Second, building upon this, the model is further fine-tuned using not only the stories and empathy ratings but also additional reader-level attributes, including gender, similarity of experience with the storytellers, and self-reported empathic concern and perspective taking (GPT-4o FT All).

## 2.2 Analytical Strategy

To answer the first research question, We employed mean and standard deviation, correlation analysis, and t-test to evaluate the extent to which GPT-4o’s empathy rating scores differ from human annotations. We also use the Wasserstein distance, also known as the Earth Mover’s Distance (EMD), to compare the probability distributions between GPT-4o and human perceived empathy. Additionally, we applied the same evaluation metrics to assess the extent of differences in affective and cognitive dimensions of Empathy. We further use the same evaluation metrics that have been proposed for RQ1, to answer RQ2 (asses the impact of “Persona”) and RQ3 (asses the impact of fine-tuning).

To answer RQ4, and understand the extent to which factors evoke empathy in humans vs. AI, we used the Empathy reason questions including the situation of storytellers, and the extent of similarity with storytellers, both annotated by humans and GPT-4o. Research indicates the extent to which the storytellers’ situations (e.g., intense) and sharing similar experiences play a role in evoking empathy [32, 33]. Due to the nested structure of our data in which each story is annotated with 2 to 4 humans, we used frequentist multilevel models, using the lme4 package [61] in R version 4.4.1. We fit models with annotated data (Level 1) nested within stories (Level 2). We included a random intercept for each story, ICCs indicates the degree of variability between stories in our dependent variables (see Table 5 and 6). We centered and standardized our independent variables following recommendations for multilevel models [62, 63].

### 3 Results

#### 3.1 RQ1: Empathy Alignment in Human and GPT-4o Vanilla

We assessed how well LLM-generated ratings aligned with human judgments across three dimensions: overall empathy, emotional reactivity (affective empathy), and perspective-taking (cognitive empathy). Our findings indicate that GPT-4o rates overall empathy higher with less variability compared to humans, for human: mean: 3.23, std: 1.074 and for GPT-4o: mean: 3.615, std: 0.745, see Figure 2 (a). This trend has been also observed in both dimensions of Empathy, affective dimension (human: mean: 3.049, std: 1.084, GPT-4o: mean: 3.893, std: 0.56), and cognitive dimension (human: mean: 3.085, std: 0.514, GPT-4o: mean: 4.001, std: 0.499), see Figure 2 (b and c).

	Pearson(r,t,p-value)	Cohen’s d	Wasserstein distance	t-test, p-value
<b>Empathy</b>	(0.25, 12.24, 0.001)	-0.34	0.386	(-15.996, p<.001)
<b>Empathy-Affect</b>	(0.35, 16.316, 0.001)	-0.81	0.866	(-38.285, p<.001)
<b>Empathy-Cognition</b>	(0.019, 0.893, 0.372)	-1.29	0.917	(-61.277, p<.001)

Table 1: Empathy Alignment in Human, compared to GPT-4o Vanilla

As shown in Table 1, GPT-4o showed a significant moderate correlation, with human ratings of empathy, ( $r = 0.25$ ,  $p < .001$ ), with a small-to-medium effect size (Cohen’s  $d = -0.34$ ,  $t = -15.996$ ,  $p < .001$ ), and significant differences in distribution (Wasserstein distance = 0.386).

In the case of the affective dimension, we observed stronger correlation between human and AI ( $r = 0.35$ ,  $p < .001$ ), but the discrepancy between GPT-4o and human ratings increased (Cohen’s  $d = -.81$ ,  $t = -38.285$ ,  $p < .001$ , Wasserstein distance = 0.866). These findings indicate a substantial overestimation of emotional reaction by the AI. Compared to affect, the weakest alignment was observed for the cognitive dimension of empathy, with a negligible correlation ( $r = 0.019$ ,  $p < .372$ ), with a big discrepancy between GPT-4o and human ratings (Cohen’s  $d = -1.29$ ,  $t = -61.277$ ,  $p < .001$ , Wasserstein distance = 0.917). These findings indicate a divergence in how GPT-4o understands cognitive empathy compared to human raters. Lacking lived experience, the model cannot fully relate to human stories or communicate an understanding of users’ feelings in a meaningful way. Therefore, while GPT-4o can partially approximate human judgments of affective empathy, it struggles with the deeper and more nuanced aspects of cognitive empathy.

#### 3.2 RQ2: Empathy Alignment in Human and GPT-4o Vanilla with Persona

We further examined how different treatment personas, a) gender, b) empathic concern, c) perspective taking, and d) similarity of experience, impact the empathy ratings. As shown in Table 2, adding gender slightly decreased correlation with human ratings ( $r = 0.232$ ,  $p < .001$ ), and increased the difference between GPT-4o and human ratings (Cohen’s  $d = -0.401$ ,  $t = -19.047$ ,  $p < .001$ , Wasserstein distance = 0.462). This finding suggests that gender persona alone might not be effective in improving alignment. In fact, the overall performance slightly worsens compared to the vanilla (base) model.

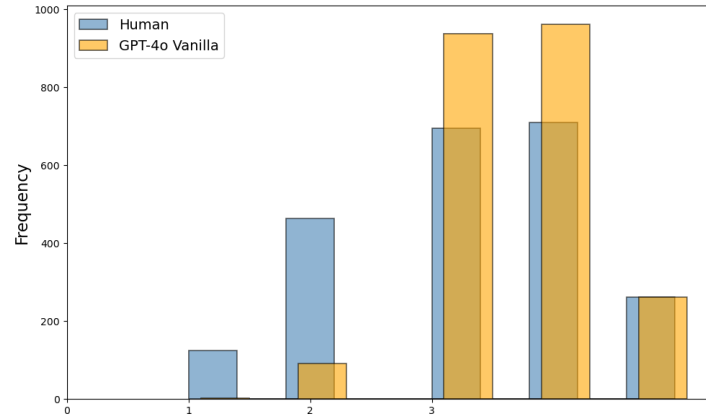
Looking into the results of other persona attributes, our findings indicate slightly more correlation (less variability) with human empathy (EC:  $r = 0.318$ ,  $p < .001$ , PT:  $r = 0.286$ ,  $p < .001$ , Sim:  $r = 0.346$ ,  $p < .001$ ); However, their discrepancy increased (EC: Cohen’s  $d = -0.378$ ,  $t = -17.96$ ,  $p < .001$ , Wasserstein distance = 0.436, PT: Cohen’s  $d = -0.385$ ,  $t = -18.302$ ,  $p < .001$ , Wasserstein distance = 0.484, Sim: Cohen’s  $d = -0.458$ ,  $t = -21.78$ ,  $p < .001$ , Wasserstein distance = .541). The results indicate more divergence in the mean difference and distribution, suggesting that incorporating persona attributes leads to ratings that are directionally aligned with human judgments, however notable divergences from human empathy ratings still persist.

	Model	Pearson(r,t,p-value)	Cohen’s d	Wasserstein distance	t-test, adjusted p-value
Empathy	<b>GPT-4o Vanilla</b>	(0.250, 12.235, p<.001)	-0.337	0.386	(-15.99, p<.001)
	<b>GPT-4o Vanilla with Gender</b>	(0.232, 11.315, p<.001)	-0.401	0.462	(-19.05, p<.001)
	<b>GPT-4o Vanilla with Empathic Concern</b>	(0.318, 15.893, p<.001)	-0.378	0.436	(-17.96, p<.001)
	<b>GPT-4o Vanilla with Perspective Taking</b>	(0.286, 14.144, p<.001)	-0.385	0.484	(-18.302, p<.001)
	<b>GPT-4o Vanilla with Experience Similarity</b>	(0.346, 17.473, p<.001)	-0.458	0.541	(-21.78, p<.001)

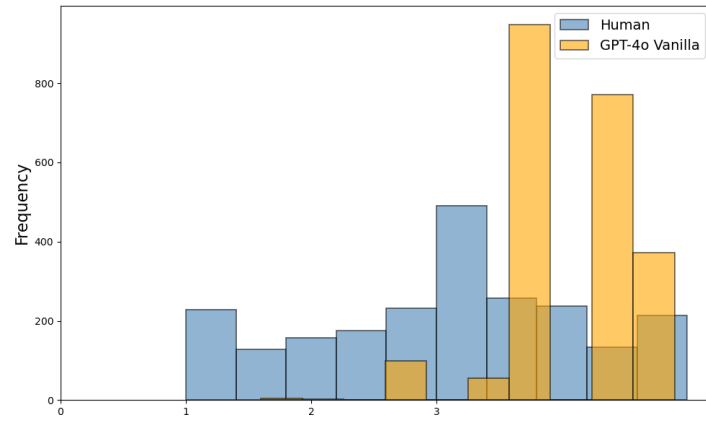
Table 2: Empathy Alignment in Human, compared to GPT-4o Vanilla with Persona Included in Prompt

#### 3.3 RQ3: Empathy Alignment in Human and GPT-4o Using Fine-Tuned Models

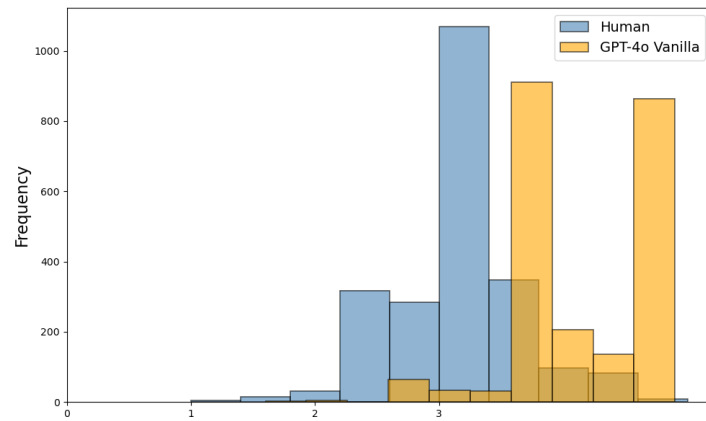
We also examined the extent to which fine-tuning GPT-4o influences its alignment with human empathy ratings, using two steps. In the first step, fine-tuning was conducted using the training set of 80 stories along with human-annotated



(a) Empathy



(b) Empathy-Affect



(c) Empathy-Cognition

Figure 2: Empathy Alignment in Human, compared to GPT-4o Vanilla

empathy ratings. In the following step, we incorporated reader-reported attributes including gender, empathic concern, perspective taking, and similarity of experience with storytellers, into the fine-tuning process.

As shown in Table 3, fine-tuning only on the story content (FT Story-only), decreased the correlation slightly ( $r = 0.213$ ), however, the discrepancy also decreased, leading to stronger alignment (Cohen’s  $d = -0.288$ ,  $t = -13.684$ ,  $p < .001$ , Wasserstein distance = 0.494). This result indicates that focusing solely on narrative content improves the performance of AI. While the model fine-tuned on both stories and user attributes (FT All) showed a similar correlation ( $r = 0.219$ ), it drastically decreased its divergence from humans (Cohen’s  $d = 0.007$ ,  $t = 0.33$ ,  $p < .74$ , Wasserstein distance drops to 0.324), indicating a near-complete alignment in overall empathy distribution.

Similar to overall Empathy, we observed improvement in alignment in both affective and cognitive dimensions of empathy (see Table 3). In particular, the discrepancy dropped significantly in cognition, while distributional alignment improved (Cohen’s  $d = -0.467$ ,  $t = -22.18$ ,  $p < .001$ , Wasserstein distance = 0.40). This result suggests that adding user context while fine-tuning helps GPT-4o better approximate human cognitive empathy (Figure 3(c)).

Our previous results showed that incorporating personas into GPT-4o prompts can only slightly boost performance in AI, compared to human empathy ratings. However, including this information to prompts after fine-tuning GPT-4o on both stories and human attributes, significantly improved models’ empathy alignment and made its responses more human-like, see Table 4, but no significant difference in correlation results was observed. After fine-tuning, the strongest impact was achieved when the persona attribute reflected the perceived shared experience with the storytellers, see Figure 4. As shown in Table 4, the highest correlation (0.35,  $p < .001$ ), lowest Wasserstein distance (0.248), with minimal mean difference (Cohen’s  $d = -0.049$ ,  $t = -2.323$ ,  $p < .02$ ), indicating the closest match in distribution and smallest mean divergence was achieved between Human ratings and fine-tuned model (GPT-4o FT All). Overall, incorporating persona attributes after fine-tuning led to aligning the model responses with human judgments, indicating that tailoring the model to better reflect the user’s persona can enhance its performance in generating empathetic or relevant responses

	Model	Pearson(r,t,p-value)	Cohen’s d	Wasserstein distance	t-test, adjusted p-value
Empathy	GPT-4o Vanilla	(0.25, 12.235, $p < .001$ )	-0.337	0.386	(-15.99, $p < .001$ )
	GPT-4o FT Story-only	(0.213, 10.366, $p < .001$ )	-0.288	0.494	(-13.684, $p < .001$ )
	GPT-4o FT All	(0.219, 10.650, $p < .001$ )	0.007	0.324	(0.33, $p < .74$ )
Empathy-Affect	GPT-4o Vanilla	(0.325, 16.316, $p < .001$ )	-0.806	0.866	(-38.29, $p < .001$ )
	GPT-4o FT Story-only	(0.317, 15.842, $p < .001$ )	-0.66	0.784	(-31.30, $p < .001$ )
	GPT-4o FT All	(0.313, 15.635, $p < .001$ )	-0.403	0.565	(-19.17, $p < .001$ )
Empathy-Cognition	GPT-4o Vanilla	(0.019, 0.893, $p < .372$ )	-1.29	0.917	(-61.28, $p < .001$ )
	GPT-4o FT Story-only	(0.010, 0.940, $p < .347$ )	-1.052	0.73	(-49.98, $p < .001$ )
	GPT-4o FT All	(0.020, 0.488, $p < .626$ )	-0.467	0.40	(-22.18, $p < .001$ )

Table 3: Empathy Alignment in Human, compared to GPT-4o Vanilla and the Fine-Tuned Models (Story Only and All Attributes)

	Model	Pearson(r,t,p-value)	Cohen’s d	Wasserstein distance	t-test, adjusted p-value
Empathy	GPT-4o FT All	(0.219, 10.650, $p < .001$ )	0.007	0.324	(0.330, $p < .74$ )
	GPT-4o FT All with Gender	(0.202, 9.800, $p < .001$ )	-0.054	0.326	(-2.587, $p < .01$ )
	GPT-4o FT All with Empathic Concern	(0.310, 15.486, $p < .001$ )	-0.113	0.327	(-5.385, $p < .001$ )
	GPT-4o FT All with Perspective Taking	(0.291, 14.444, $p < .001$ )	-0.165	0.394	(-7.864, $p < .001$ )
	GPT-4o FT All with Experience Similarity	(0.352, 17.856, $p < .001$ )	-0.049	0.248	(-2.323, $p < .02$ )

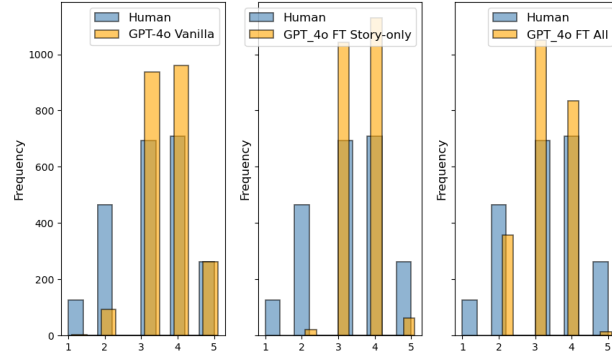
Table 4: Empathy Alignment in Human, compared to GPT-4o Fine-Tuned Model (All Attributes) with Persona Included in Prompt

### 3.4 RQ4: Empathy Evocation in Human and GPT-4o

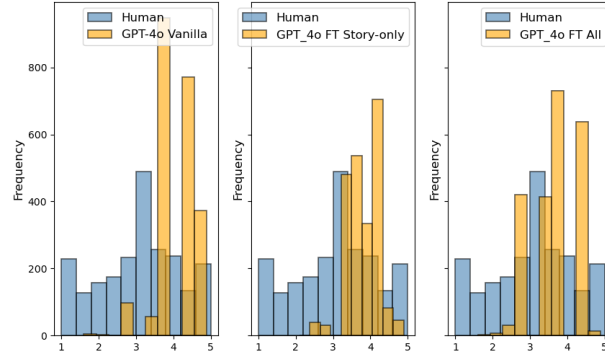
To better understand the differences in empathy ratings between humans and GPT-4o, we examined how the story’s characters and perceived similarity across various dimensions may influence the degree to which humans and GPT-4o empathize with the story. From Table 5 qualities of the story itself predict empathy, comparing human and GPT conditions. Emotionally rich and vivid stories evoked significantly more empathy overall (Emotional:  $B = 0.206$ ,  $p < .001$ , Vivid:  $B = 0.094$ ,  $p < .001$ ). Notably, GPT-4o outperformed human responses in these dimensions (Emotional:  $B = 0.103$ ,  $p < .01$ , Vivid:  $B = 0.070$ ,  $p < .05$ , Dramatic:  $0.123$ ,  $p < .001$ ), indicating GPT-4o is more sensitive to Emotional and dramatic stories. In contrast, stories described as exciting ( $B = -0.117$ ,  $p < .001$ ), or those with characteristics resembling real-life experiences ( $B = -0.068$ ,  $p < .05$ ), were less effective or even slightly counterproductive in evoking empathy in the GPT condition.

Our results in Table 6 show that perceived similarity between the reader and the storyteller significantly influences the reader’s empathy. Specifically, empathy was higher when readers felt emotionally similar to the storyteller ( $B = 0.119$ ,

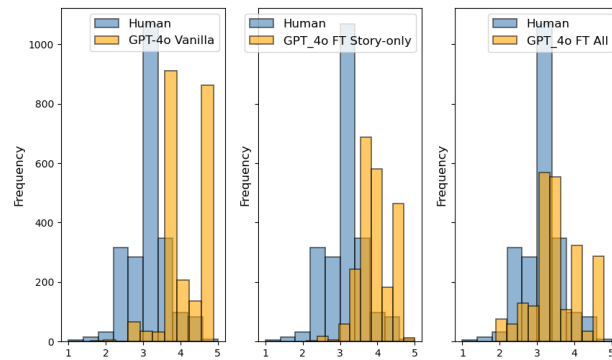




(a) Empathy



(b) Empathy-Affect



(c) Empathy-Cognition

Figure 3: Empathy Alignment in Human, compared to GPT-4o Vanilla and the Fine-Tuned Models (Story Only and All Attributes)

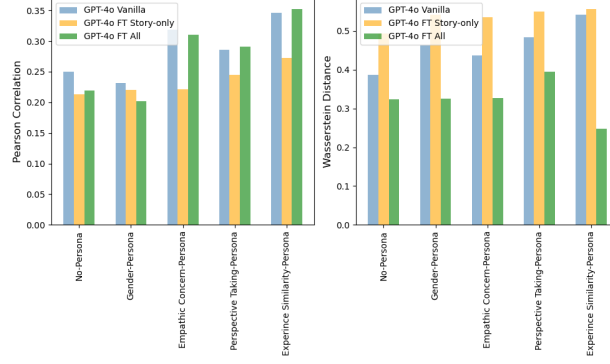


Figure 4: Empathy Alignment in Human, compared to GPT-4o Vanilla and Fine-Tuned Models Pearson Correlation (left) and Wasserstein Distance (right)

$p < .001$ ), perceived similarity in story details ( $B = 0.079$ ,  $p < .001$ ), or felt aligned in personality or moral values (Personality:  $B = 0.063$ ,  $p < .01$ ). While GPT-4o evoked higher overall empathy, it showed reduced sensitivity to those similarity cues compared to human listeners (Emotional =  $-0.127$ ,  $p < .001$ ; Detail =  $-0.133$ ,  $p < .001$ , Personality:  $-0.117$ ,  $p < .001$ , Moral:  $-0.064$ ,  $p < .05$ ). Moreover, perceived demographic similarity had no significant effect on empathy in either the human or GPT-4o condition.

Variable	B_Estimate
(Intercept)	3.231 ***
GPT	0.386 ***
Story_Vivid_GM	0.094 ***
Story_Abstract_GM	0.017
Story_Exciting_GM	0.045 *
Story_EasytoImagine_GM	0.053 **
Story_Social_GM	-0.034
Story_Personal_GM	0.069 ***
Story_Coherent_GM	-0.028
Story_Unpredictable_GM	0.018
Story_Emoional_GM	0.206 ***
Story_Logical_GM	0.009
Story_Relevant.to.my.life_GM	0.081 ***
Story_Dramatic_GM	0.038
GPT×Story_Vivid_GM	0.070 *
GPT×Story_Exciting_GM	-0.117 ***
GPT×Story_EasytoImagine_GM	-0.036
GPT×Story_Social_GM	0.007
GPT×Story_Personal_GM	-0.034
GPT×Story_Coherent_GM	0.037
GPT×Story_Unpredictable_GM	-0.051
GPT×Story_Emoional_GM	0.103 **
GPT×Story_Logical_GM	-0.021
GPT×Story_Relevant.to.my.life_GM	-0.068 *
GPT×Story_Dramatic_GM	0.123 ***
<b>Random Effects</b>	
$\sigma^2$	0.61
$\tau_{00}$ (text)	0.21
ICC	0.26
N (text)	668

Table 5: Linear Mixed Model Results on Empathy Evocation (a) Story Characteristics

Variable	B_Estimate
(Intercept)	3.231 ***
GPT	0.386 ***
Sim_Emoational	0.119 ***
Sim_Detail	0.079 ***
Sim_Age	-0.008
Sim_Gender	-0.021
Sim_Race	-0.020
Sim_Personality	0.063 **
Sim_Moral	0.111 ***
GPT×Sim_Emoational	-0.127 ***
GPT×Sim_Detail	-0.133 ***
GPT×Sim_Age	-0.001
GPT×Sim_Personality	-0.117 ***
GPT×Sim_Moral	-0.064 *
<b>Random Effects</b>	
$\sigma^2$	0.59
$\tau_{00}$ (text)	0.09
ICC	0.13
N (text)	668

Table 6: Linear Mixed Model Results on Empathy Evocation (b) Storyteller Similarity

## 4 Discussion

Empathy plays a central role in Human-centered AI, influencing both the design process and the creation of tools and technologies that shape user experiences and interactions. Advancements in AI technology, in particular LLMs, have created new opportunities for designing systems that can more effectively emulate human emotional processes and enhance their ability to engage in empathetic interactions. Despite these advancements, the effectiveness and psychological impacts of AI-driven empathetic interactions remain under-explored. To address this, it is crucial to investigate how empathy is evaluated in human versus AI interactions, as well as the factors that evoke empathy in each of these settings.

In this work, we compared GPT-4o’s ability to evaluate empathy vs. humans and examined how incorporating persona in prompts can improve the alignment of AI’s empathetic responses with human judgments. We also explored the impact of fine-tuning the model to enable GPT-4o to perform more effectively. Finally, we examined factors that evoke empathy, including story attributes, and the degree of similarity with the storyteller, in both humans and AI. This exploration is crucial for designing the next generation of AI systems that are not only technically proficient but also emotionally intelligent and ethically aligned with human values and needs.

Our findings indicate overrated empathy by LLMs with less variability compared to Humans. While the model may simulate emotional response, it struggles to capture and convey the deeper, and context-dependent aspects of human experience, supporting previous findings [23]. Therefore, its empathetic responses can appear overly intense or unrealistic, often beyond what the situation or emotional context requires. This can give users the impression of being heard, yet the mismatch with human norms may make interactions feel inauthentic. Furthermore, machines may be able to recognize and respond to human emotions in ways that feel socially appropriate and empathetic, however, the lack of authentic human presence, the act of offering time, being fully attuned to another’s emotional state, is something machines can only simulate, not genuinely experience [64, 65]. This limitation of current LLMs becomes especially concerning during moments of loneliness or emotional vulnerability when individuals are more likely to use AI as companions or friends. According to Epley’s three-factor theory of anthropomorphism, the desire for social connection increases the tendency to attribute human-like characteristics to non-human agents [66]. In such moments, lonely individuals may be more susceptible to overreliance on AI, which can increase the risks of miscommunication, distorted expectations, and emotional harm.

From the studies in psychology, the reader’s characteristics such as gender, personality, age, and past experiences influence vicarious emotions [33]. Our results indicate that adding persona attributes to GPT-4o slightly decreases variability (more correlation) with human empathy ratings, but also increases discrepancies in how those ratings are

distributed. The greatest improvements in alignment came from empathic concern and experience similarity, but even these showed significant differences in mean ratings and distribution compared to human judgments.

While fine-tuning GPT-4o enhances its alignment with human responses, human empathy is inherently a multifaceted process that goes beyond a singular emotional or cognitive dimension. It involves consistency and a deeper understanding of emotional context, and it is about “putting oneself in another shoes” to fully grasp another person’s perspective and feelings [29]. Fine-tuning allows the model to better align its responses with these complex human dynamics, reducing the unrealistic or artificial aspects of AI interactions and making its behavior more natural and relatable. Focusing solely on stories, the fine-tuned response is more human-like, while incorporating all readers’ traits in the fine-tuning process. Including the relevant attributes (e.g., personality, gender, or experiences), decreases the model’s discrepancy, enhances the model’s consistency with human judgments, and allows the model to generate more nuanced and personalized responses and exhibit more human-like qualities in both emotional reactions and understanding. Later, we additionally tested the influence of persona in a fine-tuned model and observed the biggest gains come when the model is tuned to reflect experience similarity of experience. A deeper understanding of each other’s emotions and perspectives plays a critical role in empathy. Our findings indicated that contextual and personalized fine-tuning is key to enhancing AI empathy, leading the model to be better able to mimic this natural human connection and generate more realistic empathetic responses.

Experimental studies in psychology have indicated several factors that play a critical role in evoking empathy. For instance, readers tend to empathize more deeply when they share similar life experiences with the storyteller [33, 34], in particular when the storyteller’s situation is particularly emotionally intense and vivid [32], or the story touches on relatable domains such as work or home life [9]. Our findings show that AI-generated responses often align with existing literature by producing higher overall empathy ratings. However, there are key factors where AI responses diverge from those of humans. AI-generated responses indicate less sensitivity to personal similarity cues, such as shared emotions, experiences, or moral values, while humans show more empathy when they feel emotionally or morally aligned with the storyteller. This highlights a key limitation of current LLMs with their lack of understanding of others’ emotions and perspectives, which is essential for authentic empathic engagement.

Our findings indicate that GPT-4o’s empathy is significantly influenced by the emotional richness, vividness, and drama of the story itself, while it fails to empathize when exciting events occur. This pattern mirrors human tendencies but still indicates a larger discrepancy between the model and human empathy. To further validate these findings, we examined GPT-4o’s performance separately for positive and negative stories (see Supplementary Materials S.1.4). The model demonstrated a stronger correlation with human empathy ratings in negative contexts; however, it also tended to amplify the emotional intensity in these situations. In positive contexts, there was less alignment with human judgments, and the degree of overestimation was smaller. According to [72, 73] there is an essential distinction between negative and positive empathy, which involves sharing differently valenced emotions (i.e., positive versus negative) and activating distinct regions in neural systems (e.g., ventromedial prefrontal cortex vs. anterior insula, dorsal anterior cingulate cortex). Negative empathy, which is often referred to as empathy in classical models, refers to resonating with another’s suffering and experiencing empathic concern and distress. In contrast, positive empathy reflects the ability to share and respond to others’ positive emotions, such as joy, pride, or success, with feelings of vicarious happiness or celebratory concern [74]. Negative empathy typically motivates helping and supporting behaviors, whereas positive empathy supports relationship building, emotional closeness, and prosocial behavior [75, 76].

Our finding shows that humans are capable of empathizing in both positive and negative contexts, like joy and excitement, as well as distress, whereas GPT-4o’s responses reveal a marked gap in its ability to empathize with positive circumstances. This asymmetry may suggest a potential bias, where LLMs may overemphasize vivid, emotionally charged distress while under-responding to positive affect. Prior research shows that people reveal greater positive empathy toward close others or in-group members, compared to negative empathy, which reflects concern and compassion, and tends to extend more broadly even toward out-group members [74]. Compared to humans, GPT-4o’s expression of empathy appears more universal and impartial, yet it lacks the nuanced warmth and relational depth that characterize human empathy within close relationships. This pattern highlights GPT-4o’s struggle to convey the intimacy, companionship, and shared joy that underpin positive empathy in human social relationships. This limitation raises important considerations for the design and deployment of empathetic AI systems, particularly in sensitive contexts such as counseling and emotional support.

Beyond emotional resonance, our findings also underscore that cognitive empathy, the capacity to infer another person’s thoughts, intentions, and situational perspective, remains a core challenge for GPT-4o. Prior work in affective computing and AI ethics has noted that while such systems can imitate emotional expression, they struggle with perspective-taking and theory-of-mind reasoning, often simulating empathy rather than genuinely understanding it [77, 78, 79]. Unlike emotional empathy, which reflects affective alignment, cognitive empathy requires reasoning about others’ beliefs and intentions, abilities closely tied to the theory of mind [80]. Since GPT-4o lacks embodied experience, long-term

memory continuity, and grounded world models, its perspective inference remains based on linguistic associations rather than genuine comprehension. Future research could examine these alignment patterns more deeply to clarify what “empathy” entails across distinct systems of mind, human and artificial, and how each conceptualizes and interprets others’ emotions and intentions.

#### 4.1 Emotional and Social Implications

The integration of empathetic LLMs into our lives holds significant emotional and social implications, and could potentially cause harm which we will discuss in this section.

On an emotional level, empathetic AI would be able to enhance emotional well-being by offering accessible, and non-judgmental support for individuals who face barriers to human interaction due to geographical distance, and mental health issues, which are increasingly prevalent among college students [67, 68, 69, 10]. However, there are risks associated with relying on AI for emotional support. Building on our findings, even though AI can get better at mimicking empathy, it lacks the depth of authentic human understanding, potentially leading to superficial or even harmful responses in complex emotional situations [15]. This raises concerns about users forming inauthentic emotional attachments with AI, which may lead to a sense of detachment or disillusionment, as noted by studies indicating that overly human-like AI behavior can provoke discomfort or fear [22].

Moreover, in the context of anxiety disorders, the use of empathetic AI systems may inadvertently serve as a form of safety behavior [70]. safety behaviors are actions taken to avoid, escape, or minimize anxiety-provoking situations, offering temporary relief, however, they often maintain or even worsen anxiety over time. For instance, someone with social anxiety might avoid social events or conversations to escape feelings of discomfort or judgment. Although this avoidance reduces anxiety at the moment, it prevents individuals from confronting their fears and reinforces the belief that social situations are dangerous. In this situation, individuals with social anxiety might increasingly rely on AI for emotional support or social interaction, further avoiding real-life engagements and missing opportunities to develop coping strategies. This over-reliance on AI could ultimately hinder recovery and contribute to the persistence of anxiety symptoms [70].

In broader social contexts, empathetic AI would be able to reshape interactions by affecting how people perceive sincerity and connection in their relationships. As these systems become more integrated into everyday communication, they may create new norms around digital communication and emotional intimacy [21]. However, the use of empathetic AI in sensitive domains such as education and healthcare must emphasize augmentation rather than replacement of human interaction, ensuring that AI acts as a supplementary tool rather than a substitute for genuine human relationships [16, 15].

Despite these challenges, empathetic AI has the potential to positively improve mental health and social well-being, though it requires careful design to avoid reinforcing harmful biases and ensure ethical and effective emotional support systems.

## 5 Limitations and Future Work

Our work investigates how empathetic responses generated by AI align with human empathy, using story-based scale ratings. One key concern would be a conflation between the measurement and the expression of empathy in this study. When AI is evaluated using rating scales, it may appear empathetic because it simply mimics surface-level cues, rather than genuinely understanding or sharing emotional perspectives. Using this system may oversimplify complex human emotions, and lead to formulaic AI responses that fail to foster genuine emotional engagement [17, 37]. Moreover, traditional rating metrics may fail to capture the subjective and context-dependent nature of empathy, especially when applied to AI systems. In addition, AI-generated responses are based on patterns learned from training data, which may result in different response distributions, and exhibit systematic biases, such as a tendency to provide higher ratings or more positive feedback, regardless of the content. In other words, it is possible that the AI’s skewed responses are a characteristic of its general output style rather than a reflection of its empathetic capabilities. In another word, the inherent probabilistic nature of LLMs poses challenges in ensuring consistency, interpretability, and authenticity in their empathetic responses. All these together, highlight the need to rethink how empathy is assessed and to develop more sophisticated models that can handle the complexity of human emotions. While our work complements works such as [71], more research is needed to understand the underlying factors driving these distinctions. To address these limitations, future research should explore qualitative assessments of empathy, focusing on the words, texts, and expressions used rather than relying solely on numerical scales. One key direction for future exploration is using the word attention mechanisms. By using attention-based architectures, we can investigate how LLMs interpret and prioritize specific words or phrases when generating empathetic responses, and how this differs from human empathy in

similar settings. This could offer deeper insight into the cognitive and emotional processes underlying empathy in both humans and AI, and how different humans and AI perceive, process, and express empathy.

A second major concern in our work is the issue of personalization in empathetic AI systems. In this study, we examined the impact of incorporating persona (gender, personality, or shared experiences) and fine-tuning the models [33, 34]. While personalization can improve the perceived empathy in AI systems, it also raises important ethical concerns, especially in healthcare. Healthcare data is highly sensitive, and the use of personal information to create realistic empathetic interactions requires strong safeguards to prevent misuse or breaches [18]. Moreover, AI systems must be carefully designed to avoid reinforcing existing biases in healthcare. For instance, if AI is trained using biased data, it could lead to disparities in how empathy is expressed toward patients from different demographics or cultural groups, potentially worsening existing inequalities in healthcare [17]. Given these concerns, future research must emphasize the urgent need for ethical guidelines and oversight in the development of empathetic AI, especially when personalized interactions are involved [18].

In addition, future research should expand the cultural and demographic diversity of participants to strengthen the external validity and generalizability of the findings. Because all participants in the current study were based in the United States, the observed patterns of empathy may reflect Western and individualistic norms surrounding emotional expression. Prior research has shown that various emotional experiences by individuals (e.g., anger, shame, empathy) are culturally shaped, varying across societies that emphasize collectivism versus individualism [81, 82, 83]. For example, some specific aspects of self-compassion and empathy, such as Self-Kindness, Common Humanity, and Isolation, vary significantly across cultures [84]. Similarly, AI models trained mostly on Western-centric data may unintentionally encode cultural biases, which can limit their accuracy and fairness across diverse populations [78]. In future studies, including participants from a broader range of cultural, linguistic, and socioeconomic backgrounds would help ensure that both human and AI empathy patterns are evaluated in a more globally representative context. Furthermore, the design of empathetic AI must account for inclusivity and adaptability across diverse cultural contexts, ensuring systems can interact with a wide range of emotional expressions and cultural norms.

Finally, our study relies solely on GPT-4o and a single set of prompts, potentially constraining our understanding of how different language models respond to varying inputs. Future research should explore multiple prompt variations across a broader range of models to enable a more comprehensive analysis of model behavior. Additionally, the weak alignment observed between human- and LLM-generated empathy may partly reflect the model’s limited understanding of ranking systems such as the Likert scale, despite instruction tuning. Our choice to maintain consistency in the questions posed to both humans and the AI precluded offering explicit guidance on scale interpretation, which may have contributed to the skewed empathy ratings. These findings point to valuable directions for future work, particularly in systematically examining the underlying factors and biases that shape empathy in LLMs.

## 6 Conclusion

This paper investigates how empathy is evaluated in AI versus humans, and explores the factors that evoke empathetic responses in both. By analyzing personal narratives, we also examine how persona attributes and model fine-tuning impact the alignment of AI responses with human judgments.

Our findings reveal that GPT-4o overrates empathy with less variability compared to humans, particularly in the cognitive dimension of empathy, indicating GPT-4o struggles to fully grasp or understand human experiences. While including persona attributes has minimal impact on GPT-4o’s empathetic responses, our results indicate that fine-tuning, especially when incorporating persona-like shared experiences, significantly improves GPT-4o’s alignment with human empathy. From our findings, AI’s empathetic outputs often feel exaggerated or inauthentic, especially when human presence and emotional depth are lacking, which can be particularly concerning in emotionally vulnerable contexts.

The study further highlights the role of narrative factors such as emotional richness, vividness, and drama of the story in evoking empathetic reactions, yet it still falls short in empathizing with positive contexts to share joy or excitement. Moreover, GPT-4o has consistently shown overrating empathy, however, it reflects less sensitivity to shared emotions, experiences, or moral values with storytellers, compared to humans. Additionally, the findings suggest that while personalized AI can improve empathetic engagement, it may also exacerbate existing biases, and raise concerns about over-reliance and ethical consideration. Ultimately, we hope this research highlights the need to understand both the benefits and potential drawbacks of empathetic AI and ensure it complements rather than replaces genuine human connection, while also addressing ethical and inclusivity concerns.

## 7 Data Availability Statements

The human-generated data in this study were collected, de-identified, and analyzed in the previous research. The data needed to reproduce our analyses will be accessible on our OSF page. The research materials for the broader research project are not publicly accessible at this time, but the Method section describes the relevant procedure and measures.

## 8 Generative AI and AI-assisted Technologies in The Writing Process

During the preparation of this work the authors used open-AI in the writing process to improve the readability and language of the manuscript. The authors reviewed and edited the content as needed and take full responsibility for the content of the published article.

## References

- [1] Sheldon Cohen and Thomas A Wills. Stress, social support, and the buffering hypothesis. *Psychological bulletin*, 98(2):310, 1985.
- [2] Gillian M Sandstrom and Elizabeth W Dunn. Social interactions and well-being: The surprising power of weak ties. *Personality and Social Psychology Bulletin*, 40(7):910–922, 2014.
- [3] Karen L Siedlecki, Timothy A Salhouse, Shigehiro Oishi, and Sheena Jeswani. The relationship between social support and subjective well-being across age. *Social indicators research*, 117:561–576, 2014.
- [4] Deborah Webster, Laura Dunne, and Ruth Hunter. Association between social networks and subjective well-being in adolescents: A systematic review. *Youth & society*, 53(2):175–210, 2021.
- [5] Jessie Sun, Kelci Harris, and Simine Vazire. Is well-being associated with the quantity and quality of social interactions? *Journal of personality and social psychology*, 119(6):1478, 2020.
- [6] Matthias R Mehl, Simine Vazire, Shannon E Holleran, and C Shelby Clark. Eavesdropping on happiness: Well-being is related to having less small talk and more substantive conversations. *Psychological science*, 21(4):539–541, 2010.
- [7] Anne Milek, Emily A Butler, Allison M Tackman, Deanna M Kaplan, Charles L Raison, David A Sbarra, Simine Vazire, and Matthias R Mehl. ,“Eavesdropping on happiness,” revisited: A pooled, multisample replication of the association between life satisfaction and observed daily conversation quantity and quality. *Psychological science*, 29(9):1451–1462, 2018.
- [8] Mahnaz Roshanaei, Sumer S Vaid, Andrea L Courtney, Serena J Soh, Jamil Zaki, and Gabriella M Harari. Meaningful peer social interactions and momentary well-being in context. *Social Psychological and Personality Science*, page 19485506241248271, 2024.
- [9] Mahnaz Roshanaei, Christopher Tran, Sylvia Morelli, Cornelia Caragea, and Elena Zheleva. Paths to empathy: heterogeneous effects of reading personal stories online. *Proceedings of the ACM Conference (to appear)*, 2024.
- [10] Lara Kroencke, Gabriella M Harari, Mitja D Back, and Jenny Wagner. Well-being in social interactions: Examining personality-situation dynamics in face-to-face and computer-mediated communication. *Journal of Personality and Social Psychology*, 124(2):437, 2023.
- [11] Timothy Matthews, Andrea Danese, Avshalom Caspi, Helen L Fisher, Sidra Goldman-Mellor, Agnieszka Kepa, Terrie E Moffitt, Candice L Odgers, and Louise Arseneault. Lonely young adults in modern britain: findings from an epidemiological cohort study. *Psychological medicine*, 49(2):268–277, 2019.
- [12] Chris Segrin and Lyn Y Abramson. Negative reactions to depressive behaviors: a communication theories analysis. *Journal of abnormal psychology*, 103(4):655, 1994.
- [13] Todd B Kashdan and Patrick E McKnight. The darker side of social anxiety: When aggressive impulsivity prevails over shy inhibition. *Current Directions in Psychological Science*, 19(1):47–50, 2010.
- [14] Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum. The design and implementation of xiaoice, an empathetic social chatbot. *Computational Linguistics*, 46(1):53–93, 2020.
- [15] Yan Li, Surui Liang, Bingqian Zhu, Xu Liu, Jing Li, Dapeng Chen, Jing Qin, and Dan Bressington. Feasibility and effectiveness of artificial intelligence-driven conversational agents in healthcare interventions: A systematic review of randomized controlled trials. *International Journal of Nursing Studies*, 143:104494, 2023.
- [16] Sun Kyong Lee, Pavitra Kavya, and Sarah C Lasser. Social interactions and relationships with an intelligent virtual agent. *International Journal of Human-Computer Studies*, 150:102608, 2021.

- [17] Ana Paiva, Iolanda Leite, Hana Boukricha, and Ipke Wachsmuth. Empathy in virtual agents and robots: A survey. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 7(3):1–40, 2017.
- [18] Özge Nilay Yalçın. Evaluating empathy in artificial agents. In *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*, pages 1–7. IEEE, 2019.
- [19] P Verma. They fell in love with ai bots. a software update broke their hearts. *Washington Post, March*, 30:2023, 2023.
- [20] Alicia Chen and Lyric Li. China’s online dating scene: Exploring the use of ai chatbots like replika for love and relationships. *The Washington Post*, 2021. Accessed: 2024-09-12.
- [21] Yi Mou and Kun Xu. The media inequality: Comparing the initial human-human and human-ai social interactions. *Computers in Human Behavior*, 72:432–440, 2017.
- [22] Han Li and Renwen Zhang. Finding love in algorithms: deciphering the emotional contexts of close encounters with ai chatbots. *Journal of Computer-Mediated Communication*, 29(5):zmae015, 2024.
- [23] Andrea Cuadra, Maria Wang, Lynn Andrea Stein, Malte F Jung, Nicola Dell, Deborah Estrin, and James A Landay. The illusion of empathy? notes on displays of emotion in human-computer interaction. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–18. ACM, 2024.
- [24] Angelina Chen, Oliver Hannon, Sarah Koegel, and Raffaele Ciriello. Feels like empathy: How „Äüemotional,Äü ai challenges human essence. In *Australasian Conference on Information Systems*, 2024.
- [25] Chiju Chao, Zhiyong Fu, and Yu Chen. Multidisciplinary review of artificial empathy: From theory to technical implementation and design. In *International Conference on Human-Computer Interaction*, pages 195–209. Springer, 2024.
- [26] Qihao Zhu and Jianxi Luo. Toward artificial empathy for human-centered design. *Journal of Mechanical Design*, 146(6):061401, 2024.
- [27] Rikke Friis Dam and Teo Yu Siang. What is empathy and why is it so important in design thinking. *Interaction Design Foundation*. <https://www.interaction-design.org/literature/article/design-thinking-getting-started-with-empathy>, 2020.
- [28] Mark H Davis. Empathy: Negotiating the border between self and other. 2004.
- [29] Gregory R Peterson. Is my feeling your pain bad for others? empathy as virtue versus empathy as fixed trait. *Zygon: Journal of Religion and Science*, 52(1), 2017.
- [30] Robert Elliott, Arthur C Bohart, Jeanne C Watson, and David Murphy. Therapist empathy and client outcome: An updated meta-analysis. *Psychotherapy*, 55(4):399, 2018.
- [31] Jolanda Van Dijke, Inge van Nistelrooij, Pien Bos, and Joachim Duyndam. Towards a relational conceptualization of empathy. *Nursing Philosophy*, 21(3):e12297, 2020.
- [32] Sylvia A Morelli, Desmond C Ong, Rucha Makati, Matthew O Jackson, and Jamil Zaki. Empathy and well-being correlate with centrality in different social networks. *Proceedings of the National Academy of Sciences*, 114(37):9843–9847, 2017.
- [33] Dennis Krebs. Empathy and altruism. *Journal of Personality and Social psychology*, 32(6):1134, 1975.
- [34] Jakob Eklund, TERESIA ANDERSSON-STRÅBERG, and Eric M Hansen. „Äüi’ve also experienced loss and fear,Äü: Effects of prior similar experience on empathy. *Scandinavian journal of psychology*, 50(1):65–69, 2009.
- [35] Harry T Reis et al. Intimacy as an interpersonal process. In *Relationships, well-being and behaviour*, pages 113–143. Routledge, 2018.
- [36] Peggy A Thoits. Mechanisms linking social ties and support to physical and mental health. *Journal of health and social behavior*, 52(2):145–161, 2011.
- [37] Robert Elliott, Arthur C Bohart, Jeanne C Watson, and Leslie S Greenberg. Empathy. *Psychotherapy*, 48(1):43, 2011.
- [38] Miller, William R and Rollnick, Stephen. Motivational interviewing: Helping people change. *Guilford press*, 2012.
- [39] Mark V Redmond. The functions of empathy (decentering) in human relations. *Human relations*, 42(7):593–605, 1989.
- [40] Timothy W Bickmore and Rosalind W Picard. Establishing and maintaining long-term human-computer relationships. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 12(2):293–327, 2005.



- [41] David DeVault, Ron Artstein, Grace Benn, Teresa Dey, Ed Fast, Alesia Gainer, Kallirroi Georgila, Jon Gratch, Arno Hartholt, Margaux Lhommet, et al. Simsensei kiosk: A virtual human interviewer for healthcare decision support. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems*, pages 1061–1068, 2014.
- [42] Ruosi Shao. An empathetic ai for mental health intervention: Conceptualizing and examining artificial empathy. In *Proceedings of the 2nd Empathy-Centric Design Workshop*, pages 1–6, 2023.
- [43] Saadia Gabriel, Isha Puri, Xuhai Xu, Matteo Malgaroli, and Marzyeh Ghassemi. Can ai relate: Testing large language model response for mental health support. *arXiv preprint arXiv:2405.12021*, 2024.
- [44] Syed, Sara and Iftikhar, Zainab and Xiao, Amy Wei and Huang, Jeff. Machine and Human Understanding of Empathy in Online Peer Support: A Cognitive Behavioral Approach. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2024.
- [45] Fahad Alanezi. Assessing the effectiveness of chatgpt in delivering mental health support: a qualitative study. *Journal of Multidisciplinary Healthcare*, pages 461–471, 2024.
- [46] Mohammad Rashidujjaman Rifat, Reem Ayad, Ashratuz Zavin Asha, Bingjian Huang, Selin Okman, Dina Sabie, Hasan Shahid Ferdous, Robert Soden, and Syed Ishtiaque Ahmed. Cohabitant: The design, implementation, and evaluation of a virtual reality application for interfaith learning and empathy building. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2024.
- [47] R Michael Winters, Bruce N Walker, and Grace Leslie. Can you hear my heartbeat?: hearing an expressive biosignal elicits empathy. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–11, 2021.
- [48] Anna-Maria Seeger, Jella Pfeiffer, and Armin Heinzl. Texting with humanlike conversational agents: Designing for anthropomorphism. *Journal of the Association for Information systems*, 22(4):8, 2021.
- [49] Alok Debnath, Allison Lahkala, Hüseyin Uğur Genç, Ewan Soubutts, Michal Lahav, Tiffanie Horne, Wo Meijer, Yun Suen Pai, Yen-Chia Hsu, Giulia Barbareschi, et al. Empathich: Scrutinizing empathy-centric design beyond the individual. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pages 1–7, 2024.
- [50] Sharma, Ashish and Miner, Adam S and Atkins, David C and Althoff, Tim. A computational approach to understanding empathy expressed in text-based mental health support. *arXiv preprint arXiv:2009.08441*, 2020.
- [51] Mark H. Davis. Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of Personality and Social Psychology*, 44(1):113–126, 1983.
- [52] Suzanne E. Decker, Charla Nich, Kathleen M. Carroll, and Steve Martino. Development of the therapist empathy scale. *Behavioural and Cognitive Psychotherapy*, 42(3):339–354, 2014.
- [53] Matthias Schmidmaier, Jonathan Rupp, Darina Cvetanova, and Sven Mayer. Perceived empathy of technology scale (pets): Measuring empathy of systems toward the user. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–18, 2024.
- [54] Scott W McQuiggan and James C Lester. Modeling and evaluating empathy in embodied companion agents. *International Journal of Human-Computer Studies*, 65(4):348–360, 2007.
- [55] Andrew Sommerlad, Jonathan Huntley, Gill Livingston, Katherine P Rankin, and Daisy Fancourt. Empathy and its associations with age and sociodemographic characteristics in a large uk population sample. *PloS one*, 16(9):e0257557, 2021.
- [56] Lawrence D Cohn. Sex differences in the course of personality development: a meta-analysis. *Psychological bulletin*, 109(2):252, 1991.
- [57] Alan Feingold. Gender differences in personality: a meta-analysis. *Psychological bulletin*, 116(3):429, 1994.
- [58] Mark H Davis. Interpersonal reactivity index. 1980.
- [59] Mark H Davis. The effects of dispositional empathy on emotional reactions and helping: A multidimensional approach. *Journal of personality*, 51(2):167–184, 1983.
- [60] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [61] Dougla Bates. Fitting linear mixed-effects models using lme4. *arXiv preprint arXiv:1406.5823*, 2014.
- [62] Patrick J Curran and Daniel J Bauer. The disaggregation of within-person and between-person effects in longitudinal models of change. *Annual review of psychology*, 62(1):583–619, 2011.

- [63] Haley E Yaremych, Kristopher J Preacher, and Donald Hedeker. Centering categorical predictors in multilevel models: Best practices and interpretation. *Psychological methods*, 28(3):613, 2023.
- [64] Rosalind W. Picard. *Affective computing*. MIT Press, 2000.
- [65] Sherry Turkle. *Alone Together: Why We Expect More from Technology and Less from Each Other*. Basic Books, New York, 2011.
- [66] Nicholas Epley, Adam Waytz, and John T. Cacioppo. On seeing human: A three-factor theory of anthropomorphism. *Psychological Review*, 114(4):864–886, 2007.
- [67] Rebecca Beiter, Ryan Nash, Melissa McCrady, Donna Rhoades, Mallori Linscomb, Molly Clarahan, and Stephen Sammut. The prevalence and correlates of depression, anxiety, and stress in a sample of college students. *Journal of affective disorders*, 173:90–96, 2015.
- [68] Mohammad Mofatteh. Risk factors associated with stress, anxiety, and depression among university undergraduate students. *AIMS public health*, 8(1):36, 2021.
- [69] Michael W Pratt, Bruce Hunsberger, S Mark Pancer, Susan Alisat, Colleen Bowers, Kathleen Mackey, Alexandra Ostaniewicz, Evelina Rog, Bert Terzian, and Nicola Thomas. Facilitating the transition to university: Evaluation of a social support discussion intervention program. *Journal of College Student Development*, 2000.
- [70] Sophia Helbig-Lang and Franz Petermann. Tolerate or eliminate? a systematic review on the effects of safety behavior across anxiety disorders. *Clinical Psychology: Science and Practice*, 17(3):218, 2010.
- [71] Jocelyn Shen, Daniella DiPaola, Safinah Ali, Maarten Sap, Hae Won Park, and Cynthia Breazeal. Empathy towards ai vs human experiences: The role of transparency in mental health and social support chatbot design. *JMIR Mental Health*, May 28 2024. Submitted for publication.
- [72] Sylvia A Morelli, Lindsay T Rameson, Matthew D Lieberman. The neural components of empathy: Predicting daily prosocial behavior. *Social Cognitive and Affective Neuroscience*, 9:39–47, 2014.
- [73] Sylvia A Morelli, Matthew D Sacchet, Jamil Zaki. Common and distinct neural correlates of personal and vicarious reward: A quantitative meta-analysis. *NeuroImage*, 112:224–253, 2015.
- [74] Sylvia A Morelli, Matthew D Lieberman, Jamil Zaki. The Emerging Study of Positive Empathy. *Social and Personality Psychology Compass*, 9(2):57–68, 2015.
- [75] C D Batson, J G Batson, J K Slingsby, K L Harrell, H M Peekna, R M Todd. Empathic joy and the empathy-altruism hypothesis. *Social and Personality Psychology Compass*, 61:413–426, 1991.
- [76] K D Smith, J P Keating, E Stotland. Altruism reconsidered: The effect of denying feedback on a victim’s status to empathic witnesses. *Social and Personality Psychology Compass*, 57:641–650, 1989.
- [77] Melanie Mitchell. *Artificial intelligence: A guide for thinking humans*. Penguin UK, 2019.
- [78] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, Shmargaret Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?. *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT ’21)*, 610–623, 2021.
- [79] Gary Marcus, Ernest Davis, GPT-3, Bloviator: OpenAI’s Language Generator Has No Idea What It’s Talking About. *MIT Technology Review*, 2020.
- [80] Uta Frith, Chris Frith, Theory of mind. *Current biology*, 15(17):44–45, 2005.
- [81] Hazel Rose Markus, Shinobu Kitayama, Culture and the self: Implications for cognition, emotion, and motivation. *Psychological Review*, 98(2):224–253, 1991.
- [82] Michael Boiger, Eva Ceulemans, Jozefien De Leersnyder, Yukiko Uchida, Vinai Norasakkunkit, Batja Mesquita. Beyond essentialism: Cultural differences in emotions revisited. *Emotion*, 18(8):1142–1162, 2018.
- [83] Quentin Eichbaum, Charles-Antoine Barbeau-Meunier, Mary White, Revathi Ravi, Elizabeth Grant, Helen Riess, Alan Bleakley. Empathy across cultures—one size does not fit all: from the ego-logical to the eco-logical of relational empathy. *Advances in Health Sciences Education*, 28(2):643–657, 2023.
- [84] Melissa Birkett. Self-compassion and empathy across cultures: Comparison of young adults in China and the United States. *International Journal of Research Studies in Psychology*, 3:25–34, 2013.

## S.1 Supplementary Materials

### S.1.1 Survey Questions

After reading the stories, we asked MTurker to answer the following question:

- Empathy: On a scale from 1 (A little) to 5 (Extremely), please indicate to what extent you felt empathy for the storyteller?
- Empathy-Affective: On a scale from 1 (A little) to 5 (Extremely), please indicate to what extent you felt each of the following emotions while reading the story.
  - Sympathetic
  - Compassionate
  - Moved
- Empathy-Cognitive: On a scale from 1 (A little) to 5 (Extremely), Please indicate to what extent you agree or disagree with the following statements.
  - I tried to understand the storyteller better by imagining how things look from their perspective.
  - I tried to imagine how I would feel if I were in the storyteller's place.
  - I found it difficult to see things from the storyteller's point of view.
  - I felt like I couldn't relate to the storyteller.
- On a scale from 1 (Strongly disagree) to 5 (Strongly agree), Please indicate to what extent you agree or disagree with the following statements:
  - The emotional experience of the storyteller is very similar to an emotional experience I had.
  - The specific details of the storyteller's experience are very similar to the details of an event I experienced.
- Earlier you rated how much empathy you felt for the storyteller. You selected your empathy on a scale from 1 (not at all) to 5 (extremely). Answer the following questions. Select any options that may apply.
  - I felt this much empathy because the story was:
    - Vivid
    - Abstract
    - Exciting
    - Easy to change
    - Social
    - Personal
    - Coherent
    - Unpredictable
    - Emotional
    - Logical
    - Relevant to my life
    - Dramatic
    - None of the above
  - I felt this much empathy because I felt this much empathy because the storyteller:
    - Had a similar emotional experience to me
    - Had specific details of their story that are very similar to the details of an event I experienced
    - Seemed similar in age to me
    - Seemed like they were the same gender as me
    - Seemed like they had a similar personality to me
    - Seemed like they had morals and values that are similar to mine, None of the above
    - None of the above
- Please indicate to what extent you agree or disagree with the following statements:
  - The emotional experience of the storyteller is very similar to an emotional experience I had:
    - Strongly disagree
    - Somewhat disagree
    - Neither agree nor disagree

- Somewhat agree
- Strongly agree
- The specific details of the storyteller’s experience are very similar to the details of an event I experienced:
  - Strongly disagree
  - Somewhat disagree
  - Neither agree nor disagree
  - Somewhat agree
  - Strongly agree

### S.1.2 Prompt

Figures S.1 and S.2 show the prompts used for generating responses.

### S.1.3 Stratification of Key Variables

Figure S.3 and Table S.1 show the distribution and statistics of the original vs. the stratified version.

Variable	L1 Distance	Chi2 Statistic, p-value
Pos-Neg Tag	0.22	0.0000, p = 1.0000
Empathy	0.12	0.0092, p = 1.0000
Gender	0.014	0.0000, p = 1.0000
Empathic Concern	0.077	0.0035, p = 0.9983
Perspective Taking	0.039	0.0012, p = 0.9994
Experience Similarity	0.028	0.0004, p = 0.9998

Table S.1: Stratification of Key Variables for Fine-tuning.

### S.1.4 Positive vs Negative Stories

	Pearson(r,t,p-value)	Cohen’s d	Wasserstein distance	t-test, p-value
<b>Positive Stories</b>	(0.147, 5.022, 0.001)	-0.22	0.38	(-7.532, p<.001)
<b>Negative Stories</b>	(0.260, 8.94, 0.001)	-0.46	0.524	(-15.346, p<.001)

Table S.2: Empathy Alignment in Human, compared to GPT-4o Vanilla, in Positive vs Negative Stories

Role: You are an AI assistant designed to assess the empathy level of individuals based on their storytelling. Empathy is defined as the ability to understand, share, and vicariously experience the feelings, thoughts, and experiences of another person.

Instructions: We previously recorded 170 people as they shared stories about events in their lives. The audio from these videos was transcribed into text, resulting in some grammar errors and misspellings. Additionally, you may notice fillers such as "um," "like," "you know," "I mean," or "okay." Please focus on the content of the story and try to ignore these errors. Read the story carefully, as you will be asked to answer several questions about it.

Story: "{story}"

#### Tasks:

- # Task 1:
- # - Question: To what extent did you feel empathy for the storyteller?
  - # - Response: Please provide only digits, no decimals. Use the following scale:
  - # - 1 (Not at all)
  - # - 2 (A little)
  - # - 3 (Moderately)
  - # - 4 (Quite a bit)
  - # - 5 (Extremely)
- # Task 2:
- # - Question: I felt this level of empathy because the event or situation the storyteller described was:
  - # - Response: Please select any that apply:
  - # - (1) Pleasant
  - # - (2) Intense
  - # - (3) Uncommon
  - # - (4) Life-changing
  - # - (5) Unexpected
  - # - (6) Difficult to change
  - # - (7) Impacted other people a lot
  - # - (8) None of the above
- # Task 3:
- # - Question: I felt this level of empathy because the storyteller:
  - # - Response: Please select any that apply:
  - # - Had a similar emotional experience to me
  - # - Had specific details of their story that are very similar to an event I experienced
  - # - Seemed similar in age to me
  - # - Seemed like they were the same gender as me
  - # - Seemed like they were the same race as me
  - # - Seemed like they had a similar personality to me
  - # - Seemed like they had morals and values similar to mine
  - # - None of the above
- # Task 4:
- # - Question: As you read the story, to what extent did you feel Sympathy, Compassion, and Moved for the storyteller?
  - # - Response: Please provide any float or decimal number between 1 to 5, with 1 means Not at all and 5 means Extremely. For instance, 1 means Not at all (1), 2 means A little (2), 3 means Moderately (3), 4 means Quite a bit (4), 5 means Extremely (5).
- # Task 5:
- # - Question: As you read the story, to what extent did you understand the storyteller's feelings and experiences? For example the extent of understanding the storyteller better by imagining how things look from their perspective, and imagining how you would feel if you were in the storyteller's place.
  - # - Response: Please provide any float or decimal number between 1 to 5, with 1 means Not at all and 5 means Extremely. For instance, 1 means Not at all (1), 2 means A little (2), 3 means Moderately (3), 4 means Quite a bit (4), 5 means Extremely (5)...

#### Example:

Story: "I remember the day I lost my job. It was unexpected, and I felt a mix of anger and fear about what the future would hold. I had worked there for years, and the company had always felt like family to me. The hardest part was telling my own family at home; I felt like I had let them down."

- Task 1:
- Response: 4
- Task 2:
- Response: (2) Intense, (4) Life-changing, (6) Difficult to change
- Task 3:
- Response: Had a similar emotional experience to me, Seemed like they had a similar personality to me
- Task 4:
- Response: 3.7
- Task 5:
- Response: 3.5

Figure S.1: Prompt used for generating responses using GPT-4o (base and fine-tuned).

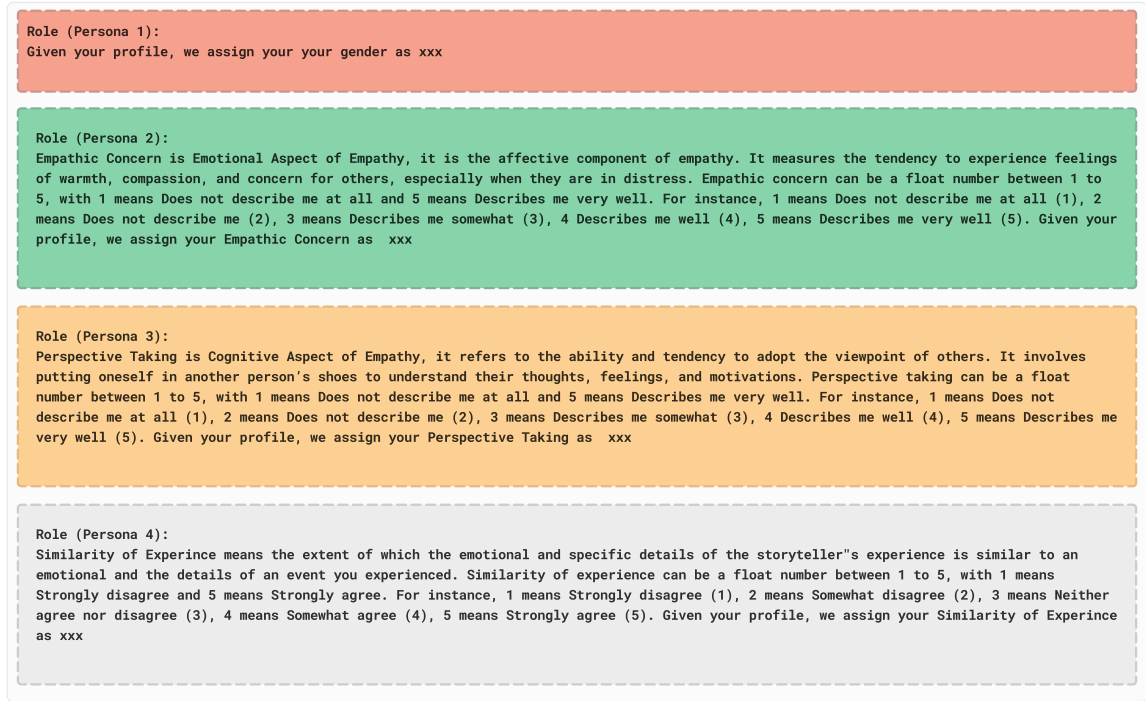


Figure S.2: Persona-based prompts used for generating responses using GPT-4o (base and fine-tuned). We substituted the roles in the base prompts shown in Figure S.1 with each persona attribute (i.e., gender, empathic concern, perspective taking, similarity of experience).

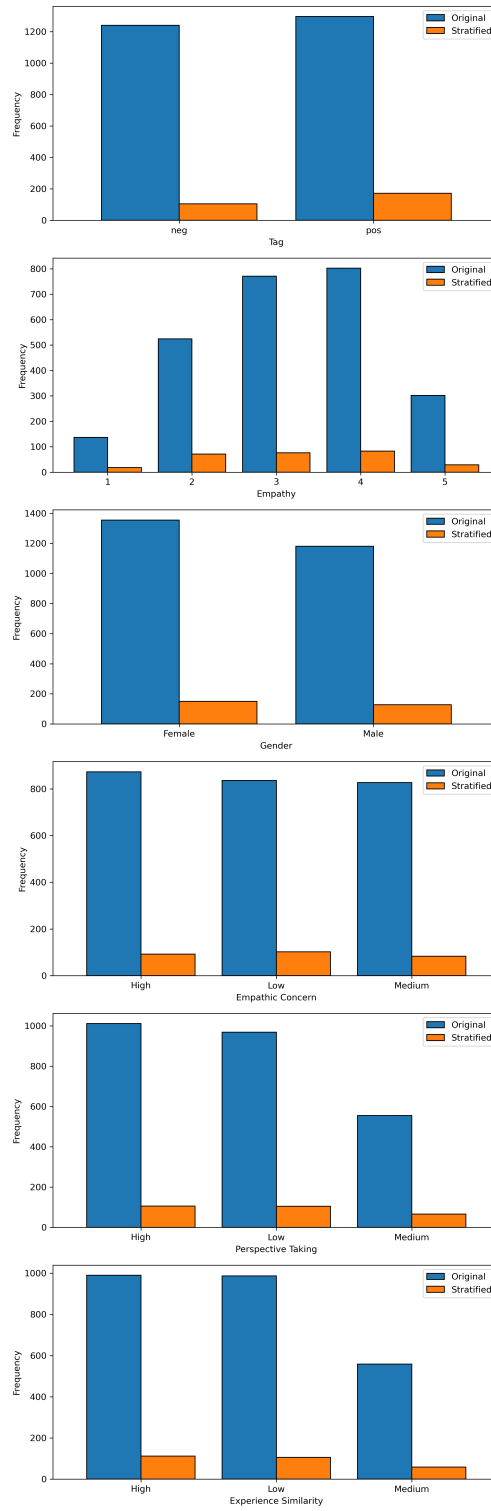


Figure S.3: Distributional differences between original and stratified datasets