

The Pursuit of Empathy: Evaluating Small Language Models for PTSD Dialogue Support

Anonymous ACL submission

Abstract

Can small language models (0.5B–5B parameters) meaningfully engage in trauma-informed, empathetic dialogue for individuals with PTSD? We answer this by introducing TIDE, a dataset of 10,000 two-turn dialogues across 500 diverse PTSD client personas, grounded in a three-factor empathy model: emotion recognition, distress normalization, and supportive reflection. All scenarios and reference responses were reviewed for realism and trauma sensitivity by a clinical psychologist specializing in PTSD. Eight small language models are evaluated before and after fine-tuning, with outputs compared to a frontier model (Claude Sonnet 3.5) as reference. Our IRB-approved human evaluation and automatic metrics reveal that, while fine-tuning generally improves perceived empathy, gains are highly scenario- and user-dependent, with smaller models facing an “empathy ceiling.” Notably, demographic analyses show older adults value distress validation and graduate-educated users prefer nuanced replies, while gender effects are minimal. We highlight limitations of automatic metrics and the need for context- and user-aware system design. Our findings—along with the planned release of TIDE—offer a foundation for building safe, resource-efficient, and ethically sound empathetic AI to supplement, not replace, clinical mental health care.

1 Introduction

Empathy is fundamental to effective mental health support, yet most conversational AI systems remain ill-equipped to recognize and validate users’ emotional experiences (Sharma et al., 2021; Aher et al., 2023). This limitation is especially consequential for individuals with Post-Traumatic Stress Disorder (PTSD), who often seek understanding and emotional validation rather than mere information (Rabin and Fogel, 2009).

Recent advances in large language models (LLMs) suggest that AI can approximate human-like empathetic communication. However, their sheer scale (often exceeding 100B parameters) creates major barriers for privacy, latency, and resource-constrained deployment (Tahir, 2024). Small language models (SLMs), by contrast, offer practical benefits such as on-device inference and greater accessibility. But can models with just 0.5B–5B parameters deliver meaningful, trauma-informed empathy?

This study investigates four research questions:

RQ1: Can SLMs generate empathetic responses in PTSD support dialogue contexts that approach the quality of large models?

RQ2: Does model size limit the quality of empathetic response, creating an “empathy ceiling” for SLMs?

RQ3: How do user demographics (age, gender, trauma type) shape perceptions of empathy from models?

RQ4: To what extent do automatic metrics (e.g., cosine similarity) reflect true human judgments of empathy?

Two central challenges shape our approach: (1) the limited expressive capacity of SLMs, and (2) the scarcity of ethically sourced, trauma-informed dialogue data (Kokosi and Harron, 2022; BN and Abdullah, 2022; BN et al., 2023). To address these, we introduce the TIDE dataset: 10,000 fully synthetic, two-turn conversations between PTSD-affected individuals and a supportive AI chatbot, spanning 500 diverse personas. Persona profiles were adapted from the metadata of the publicly available Thousand Voices of Trauma dataset (BN et al., 2025), but all dialogue is newly generated and distinct from prior work. Every dialogue pair was reviewed for emotional plausibility and trauma sensitivity by a clinical psychologist specializing in PTSD, ensuring that both user utterances and chatbot responses model best practices in trauma-

informed communication. Unlike previous work on synthetic therapy simulations, TIDE focuses on concise, realistic person-chatbot exchanges aimed at training SLMs for real-time empathetic validation. Grounded in established empathy models (Goleman, 2006; Davis, 1983), our framework operationalizes empathy as a three-factor process: (1) recognizing emotion, (2) normalizing distress, and (3) offering supportive reflection. We evaluate SLMs with both automatic metrics and human empathy ratings.

Our findings reveal that fine-tuning SLMs on TIDE improves empathetic response quality, but gains are highly scenario- and user-dependent. Demographic attributes—such as age, education, and region—shape empathy perception, and standard automatic metrics do not always align with human judgments. Ultimately, this work lays the foundation for developing lightweight, trauma-aware AI systems that can extend the reach of mental health support—while foregrounding the ethical, contextual, and human-centered considerations needed for safe deployment.

2 Motivation

Post-Traumatic Stress Disorder (PTSD) affects over 13 million adults annually in the U.S. (National Institute of Mental Health, 2023; Seligowski et al., 2015). In these contexts, empathetic communication is critical—users seek acknowledgment and safety, especially during vulnerable moments (Bickmore and Picard, 2005; Morris et al., 2018). Yet most conversational AI systems prioritize utility: setting reminders or providing information. When a user says “I’m feeling overwhelmed,” assistants may reply, “Would you like to hear a podcast?”—a utilitarian response that can feel dismissive or tone-deaf in emotionally vulnerable moments¹.

Disclosure is an important component of the “healing process” following trauma exposure (Frattaroli, 2006). Prior research has shown that disclosure is associated with lower distress and more positive post-traumatic outcomes (Hemenover, 2003; Hassija and Turchik, 2016). However, there are many challenges that prevent survivors from disclosing their traumas, including stigma (Schomerus et al., 2021) and shame (Timblin and Hassija, 2023). Chatbots can address some of these challenges to help individuals engage in disclosures practices.

¹Authors’ observation based on prototypical interactions with current digital assistants.

Specifically, chatbots can support disclosure without any fear of negative evaluation and judgment (Cho et al., 2022; Lucas et al., 2014), which can help individuals maintain a sense of control and empowerment—an important component of disclosure (Kimbly et al., 2023). In other words, a chatbot that can provide trauma-informed and empathetic responses can support effective disclosure, which might lead to better post-traumatic growth and coping strategies (Frattaroli, 2006; Hassija and Turchik, 2016).

While Large Language Models (LLMs) show promise for empathetic generation, they are rarely trauma-specific and are too resource-intensive for scalable or privacy-preserving deployment. Small Language Models (SLMs) offer practical advantages—lower latency, on-device use, and reduced risk—but face three key gaps in the PTSD context:

(1) Lack of fine-tuning on trauma-informed dialogue patterns, such as avoidance cues or dissociative disclosures, (2) Absence of evaluation frameworks grounded in safety, empathy, and distress mitigation, and (3) Minimal attention to ethical design for simulated therapy or support contexts.

This work addresses these gaps by introducing the TIDE dataset—a clinician reviewed, PTSD-specific synthetic corpus—and by fine-tuning SLMs to respond empathetically to trauma disclosures. We combine automatic metrics and trauma-informed human evaluations to determine whether these compact models can provide safe, meaningful support in sensitive mental health dialogues.

3 Related Work

We review six interconnected strands: (1) fine-tuning and prompting, especially for small models (Tahir, 2024; Raamkumar and Loh, 2024; Yang et al., 2024; Hsu et al., 2023; Tiwari et al., 2022; Ma and Chang, 2024; Yu and McGuinness, 2024; Li et al., 2024), (2) empathetic AI in mental health (Sharma et al., 2021; Chang et al., 2024; Ma and Chang, 2024; Wu et al., 2023; Sharma et al., 2020; Zheng et al., 2021; Wang et al., 2024; Welivita et al., 2023), (3) emotion and intent modeling (Su et al., 2023; Yang and Jurgens, 2024), (4) reinforcement learning for empathy (Sharma et al., 2021; Chang et al., 2024; Hsu et al., 2023), (5) human-centered evaluation (Sharma et al., 2020; Xu and Jiang, 2024; Marrapese et al., 2024; Park et al., 2024), and (6) ethical safety in trauma-informed AI (Hua et al., 2024; Gabriel et al., 2024; AlMakinah et al., 2024). Each subsection highlights advances

and persistent gaps, particularly regarding trauma-sensitive tuning and evaluation for small language models (SLMs).

3.1 Fine-Tuning and Prompting for Empathetic Dialogue

Fine-tuning has improved empathetic response generation in LLMs (Tahir, 2024; Raamkumar and Loh, 2024; Yang et al., 2024; Hsu et al., 2023; Tiwari et al., 2022; Ma and Chang, 2024; Yu and McGuinness, 2024; Li et al., 2024). Work by Tahir et al. (Tahir, 2024) explored fine-tuning small models (e.g., Mistral 7B, Qwen 2.5 7B, LLaMA 3.1 8B) for CBT, with some attention to ethical issues but not PTSD. Instruction tuning has also been applied for general counseling tasks (Li et al., 2024). Psychological frameworks (Chain of Empathy (Lee et al., 2023)), few-shot prompting (Qian et al., 2023; Filienko et al., 2024), and tasks like cognitive distortion detection (Chen et al., 2023) further enrich dialogue. Recent work also leverages embedding fusion (Yang et al., 2024; Rasool et al., 2025) and modularity (Hsu et al., 2023), primarily on general empathy datasets (Qian et al., 2023; Loh and Raamkumar, 2023), but rarely targets trauma-specific tuning or SLM safety. Our work fine-tunes SLMs on the PTSD-specific TIDE dataset and systematically evaluates limitations.

3.2 Empathetic AI in Mental Health Contexts

AI for mental health has focused on broad emotional support. Early work defined computational empathy frameworks (Sharma et al., 2020) and empathetic response models (Zheng et al., 2021; Welivita et al., 2023). Other advances include empathetic rewriting in online forums (Sharma et al., 2021), supportive response selection (Chang et al., 2024), and enhancing emotional nuance (Ma and Chang, 2024). LLMs have been studied for PTSD diagnosis (Wu et al., 2023) and simulated patient training (Wang et al., 2024), but usually employ generic datasets (Welivita et al., 2023). Some chatbots show potential for social inclusion (De Genaro et al., 2020; Chin et al., 2023), yet LLMs often underperform compared to human therapists (AlMakinah et al., 2024) or trained peers (Iftikhar et al., 2024) in specialized contexts.

3.3 Emotion and Intent Recognition in Empathetic AI

Effective empathetic AI depends on accurate emotion and intent modeling (Su et al., 2023; Tiwari

et al., 2022). Multi-source perception (Su et al., 2023) and intent prediction within T5 (Tiwari et al., 2022) have improved general empathetic responses. Hierarchical models combine emotion and dialogue acts (Zheng et al., 2021); appraisal-based approaches are also explored (Yang and Jurgens, 2024). However, these models rarely address PTSD-specific emotions, triggers, or safeguards.

3.4 Reinforcement Learning for Empathy

Reinforcement learning (RL) strategies have been used for empathetic dialogue, such as policy learning for empathic rewriting (Sharma et al., 2021) and RL-aligned support strategies (Chang et al., 2024; Hsu et al., 2023). Yet, these methods do not incorporate clinical knowledge of PTSD or enforce therapeutic boundaries. RL for trauma-sensitive conversation remains a gap.

3.5 Human-Centered and Clinically-Grounded Evaluation

With the limitations of automated metrics, research increasingly emphasizes human-centered evaluation: A/B testing, empathy ratings (Hsu et al., 2023; Tiwari et al., 2022; Sharma et al., 2021), structured frameworks like EPITOME (Sharma et al., 2020; Loh and Raamkumar, 2023), and multidimensional empathy assessment (Xu and Jiang, 2024; Raamkumar and Loh, 2024). Psychotherapy-informed metrics (Marrapese et al., 2024) and safety benchmarks (Park et al., 2024) are also emerging, but trauma-sensitive metrics (safety, trust, distress mitigation) remain rare. Our work introduces human ratings for PTSD-specific scenarios in TIDE.

3.6 Ethical Considerations and Safety

Safety is a recurring theme, but most work only considers general guardrails (Tahir, 2024; Yang et al., 2024; Sharma et al., 2021). Some efforts address toxicity prevention, privacy (AlMakinah et al., 2024), and limited PTSD risks (e.g., re-traumatization (Iftikhar et al., 2024), crisis handling). Bias and equity remain concerns (Gabriel et al., 2024). Empirical solutions such as disclaimers or boundary-aware tuning are underexplored. Robust safety frameworks for trauma-informed AI are still needed (Park et al., 2024).

3.7 Summary

Despite advances in fine-tuning, prompting, RL, intent modeling, and evaluation, key gaps persist:

- Limited trauma-specific fine-tuning for SLMs, especially for PTSD symptoms and avoidance.
- Lack of standardized, trauma-informed evaluation for safety and distress mitigation.
- Minimal ethical design tailored to trauma therapy simulation and support.

We address these by fine-tuning small models on the PTSD-specific TIDE dataset, evaluating with trauma-informed, human-centered metrics, and grounding our approach in ethical best practices.

4 Methodology

4.1 Data Description: TIDE Dataset

We empirically evaluate the Trauma-Informed Dialogue for Empathy (TIDE) dataset, a novel large-scale resource designed to assess and enhance the empathetic abilities of small language models (SLMs) in trauma-informed contexts.

4.2 Dataset Composition and Structure

TIDE consists of 10,000 fully synthetic, two-turn dialogues generated using LLMs from 500 diverse client personas, whose clinically relevant characteristics (age, gender, trauma history) were adapted from the metadata of the Thousand Voices of Trauma dataset (BN et al., 2025). No original conversation content was reused. For each persona, 20 distinct interaction scenarios were created to capture different facets of their trauma narrative while maintaining thematic coherence. To ensure emotional plausibility and trauma sensitivity, all scenarios and reference responses were reviewed by a domain expert—a clinical psychologist specializing in PTSD and Prolonged Exposure therapy, and a co-author of this paper. Representative examples (see Figure 1) were examined to confirm that both client utterances and empathetic responses were realistic, emotionally appropriate, and aligned with trauma-informed communication principles. This clinical validation reinforces the quality and practical relevance of the TIDE resource. Each instance in TIDE includes:

Client Utterance (user_input): A realistic expression of distress, PTSD symptoms (e.g., avoidance, intrusive memories, self-blame, relational struggles), or emotionally complex states tied to the persona (e.g., Persona 1).

Empathetic Response (reference_response): A gold-standard empathetic reply designed to validate and support, not diagnose or treat. These

responses emphasize emotional connection (e.g., “It’s understandable that those memories still affect you. Your feelings are valid. Have you considered talking to a therapist who specializes in trauma?”). See Fig. 1 for sample dialogues.

Generation Methodology: The reference responses (and potentially client utterances; pending internal confirmation) were generated using Sonnet 3.5 (Anthropic, 2024). Prompts included persona background and user utterance to elicit contextually appropriate, highly empathetic replies aligned with supportive communication principles.

Application & Availability: In our evaluations, SLMs are prompted with user_input from TIDE, and their outputs (zero-shot or fine-tuned) are assessed against the corresponding reference_response using multiple metrics. The full dataset—user inputs, reference responses, and persona identifiers—is planned for public release to support broader empathetic AI research.

4.3 Metrics and Analysis

To evaluate the empathetic response generation capabilities of the Small Language Models (SLMs) under investigation, both prior to (zero-shot) and subsequent to fine-tuning on the TIDE dataset, we employ a multi-faceted evaluation strategy. This strategy integrates established automatic metrics for quantifiable assessment of text quality and similarity, alongside rigorous human evaluations designed to capture the critical, yet often elusive, perceived quality of empathy.

4.3.1 Automatic Evaluation Metrics

To quantitatively assess the quality of generated responses, we use both reference-based and reference-free automatic metrics, including BERTScore (Zhang et al., 2019), ROUGE-L (Lin, 2004), METEOR (Banerjee and Lavie, 2005), Distinct-1/2 (Li et al., 2016), and cosine similarity using *all-MiniLM-L6-v2* embeddings (Reimers and Gurevych, 2021). These metrics collectively capture semantic similarity, lexical overlap, fluency, and response diversity, and are widely used for evaluating chatbot and summarization outputs. We report results across all models in both zero-shot and fine-tuned configurations (see Table 4). While these metrics offer scalable benchmarks for response quality, they have significant limitations in evaluating empathy and clinical appropriateness. Automatic metrics primarily capture surface-level textual similarity and fluency but may not ade-

TIDE Dataset: Example Empathetic Dialogue Instances
10,000 Instances | 500 PTSD Personas | 20 Interactions/Persona

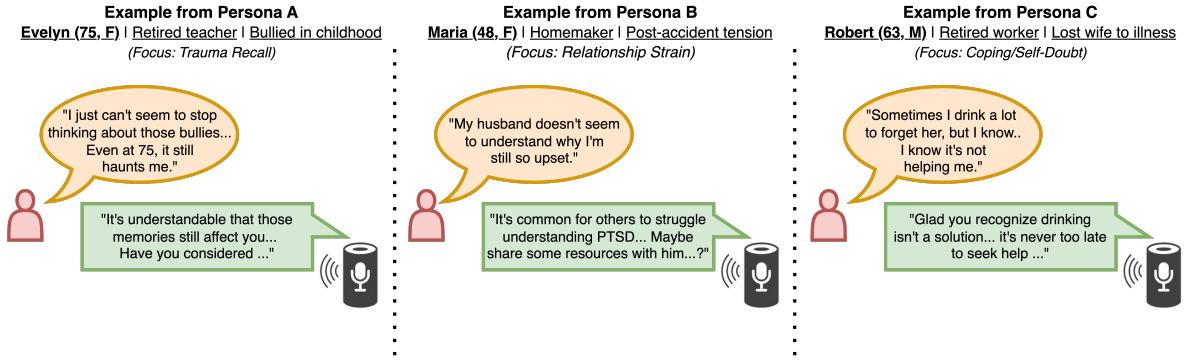


Figure 1: Illustrative examples of the two-turn dialogue structure within the TIDE dataset, showcasing synthetic empathetic interactions across diverse PTSD personas. Each instance consists of a user_input reflecting a trauma-related challenge and a corresponding reference_response providing empathetic validation. The examples span distinct trauma themes (e.g., trauma recall, relationship strain, coping/self-doubt), with persona metadata highlighting contextual grounding. Each reference response demonstrates the three-factor empathy model introduced in the text: (1) recognizing emotion, (2) normalizing distress, and (3) offering supportive reflection. TIDE comprises 10,000 such instances from 500 unique personas (20 per persona), with responses generated using Claude Sonnet 3.5 and validated by PTSD clinical specialists.

quately reflect the subtle aspects of perceived empathy, such as emotional resonance, validation, and contextual appropriateness in tone.

To address this, we conducted a human evaluation (N=116) to directly assess empathetic quality based on human judgments. The survey instrument and evaluation protocol are detailed in Appendix A and Appendix B, respectively. Participants were U.S.-based, first-language English speakers who completed the survey on Qualtrics. Recruitment was conducted via Prolific, and a three-stage quality control pipeline was applied—including pilot testing, attention checks, completion time thresholds, and alignment with gold-standard LLM ratings—with full recruitment and filtering details provided in Appendix C. Only participants passing all quality and attention checks were included in the analysis.

5 Human Evaluation Results

Because TIDE examples were vetted for clinical realism by PTSD experts, our IRB-approved human evaluation focused solely on perceived empathy, trusting the input scenarios' authenticity. Participants rated empathy on a 5-point Likert scale for each model response, following the protocol in Appendix B. We compared zero-shot (ZS) and fine-tuned (FT) responses via pairwise comparison for each model and scenario, calculating mean empathy scores and awarding a win to the higher-scoring condition. Each model was also bench-

marked against Sonnet 3.5's ground-truth reference: two points if Claude beat both ZS and FT, one if it beat one, and zero otherwise. Table 3 details model wins and highlights both the effects of fine-tuning and the benchmark set by Sonnet 3.5.

Across all comparisons:

- ZS beat FT in 26 cases; FT beat ZS in 22.
- Sonnet 3.5 dominated 72 of 96 comparisons.

By scenario, ZS models won 2, FT models won 2, and 2 were ties, indicating parity overall. However, FT gains were not universal; ZS sometimes won by a larger margin (e.g., Scenario 4). Sonnet 3.5 consistently outperformed all small models, reaffirming the gap between frontier and efficient architectures. Still, some small FT models (e.g., Llama-3.2-1B, DeepSeek-R1-Distill-Qwen-1.5B) approached high-quality responses in several cases. Human evaluations align with automatic metrics and reveal that fine-tuning's empathy gains are scenario-dependent and non-monotonic, especially for smaller models.

5.1 Human Evaluation Visualization

To complement the aggregate win-loss analysis, we visualize the full distribution of empathy ratings for each model across all conversational scenarios in Figure 2. The figure presents kernel density plots of participant Likert ratings (1–5 scale), grouped by model (ZS, FT, and reference) and input scenario.

Key Trends in Empathy Ratings:

- **Fine-tuning gains are scenario- and size-dependent:** Larger models (e.g., *Phi-3.5B*, *Llama-3.2-3B*, *Qwen2.5-3B*) improve most in complex scenarios (3, 5, 6). Smaller models (e.g., *Qwen2.5-0.5B*, *Llama-3.2-1B*) show gains mainly in simpler ones (1, 2, 5), while mid-sized models (e.g., *Qwen2.5-1.5B*, *R1-Distill-1.5B*) perform well across both types. Scenarios 3 and 5 are the most responsive overall.
- **Claude Sonnet 3.5 remains the top performer:** It consistently receives the highest empathy ratings, though in scenarios like 3 and 4, fine-tuned SLMs narrow the gap.
- **Participant variability is high:** Empathy ratings vary widely across individuals, underscoring the subjectivity of perceived empathy and the value of aggregate comparisons.

Figure 2 visually reinforces the finding that the improvement due to fine-tuning varies based on both model architecture and conversational context.

6 Discussion

Our findings reveal nuanced insights into the empathetic capabilities of small language models (SLMs) and the effects of targeted fine-tuning on trauma-informed dialogue generation.

6.1 Fine-Tuning: Uneven Empathy Gains

Across all automatic metrics (cosine similarity, BERTScore, ROUGE-L, Meteor, Distinct-1/2) and human evaluations, fine-tuning boosted performance over the ZS baseline, but not consistently. As shown in Section 5, FT models surpassed ZS models in 22 cases, while ZS outperformed FT in 26. Smaller models (e.g., *Qwen2.5-0.5B*, *Llama-3.2-1B*) showed more variability, sometimes generating more natural or emotionally resonant responses without fine-tuning. This suggests that fine-tuning on synthetic data may lead these models to overfit to affirmation-heavy patterns, reducing refined emotional responses in some scenarios.

6.2 Empathy may have a Ceiling for SLM’s

Even after fine-tuning, small models lagged behind the Sonnet 3.5 reference. Claude led in human empathy ratings, winning 72 of 96 direct comparisons (75%). Qwen models—including 0.5B—occasionally matched or surpassed Claude,

particularly with complex emotional inputs (e.g., Scenario 2: anger and isolation). Other larger models (Phi, Granite, Llama) did not. These results suggest a possible “empathy ceiling” for smaller models: while fine-tuning helps, models with 0.5–3B parameters face inherent limits in producing contextually rich, emotionally sensitive responses.

6.3 Auto Metrics are Incomplete Proxies

Although FT models improved on metrics such as cosine similarity and BERTScore, these gains did not always align with higher human-perceived empathy. Models with similarity scores often differed in participant ratings, underscoring the limits of purely similarity-based evaluation for affective tasks. Metrics like ROUGE and Distinct showed even weaker alignment with human judgments, reflecting their focus on surface-level overlap rather than emotional quality. This highlights the critical role of human evaluation in assessing affective generation tasks and suggests the need for developing more emotionally-aware automatic evaluation frameworks.

6.4 Scenario-Specific Sensitivity

Analysis across six scenarios revealed that model improvements were highly input-sensitive. Scenarios involving clear expressions of distress (e.g., “I had a panic attack at the grocery store”) showed greater fine-tuning gains, while more complex emotional situations (e.g., “I feel angry at her for leaving me”) elicited less consistent improvements. This suggests that fine-tuning with trauma-informed data helps models handle direct expressions of suffering but leaves room for improvement in handling nuanced emotional expressions like anger, guilt, or relational tension. Future fine-tuning strategies might benefit from more diverse emotional curricula that explicitly balance affirmation with context-sensitive complexity.

6.5 User Diversity in Empathy Preferences

Empathy is not one-size-fits-all. Our demographic analysis revealed significant variation in how different user groups perceived empathy in model responses, highlighting the importance of considering diversity when developing and evaluating systems trained on synthetic data. Yet, certain response characteristics—emotional clarity, contextual relevance, and actionable support—consistently emerged as valued across groups.

Older Adults: Validation Before Resolution

Participants over 60 consistently rated responses from select small models higher than younger participants. Notably, Llama 3.2 3B showed a statistically significant uplift among older adults ($mean = 3.66$, $SD = 0.94$) compared to younger ones ($mean = 3.51$, $SD = 1.01$; $t = 2.89$, $p = 0.004$, Cohen's $d = 0.16$). Similarly, Phi 3.5 Mini was also rated more empathetic by older adults ($mean = 3.66$) than younger adults ($mean = 3.48$; $p = 0.002$), with a moderate effect size. All p-values were assessed using a Bonferroni correction for multiple comparisons across nine models (adjusted $\alpha = 0.0056$); both results remained significant after correction. These models excelled not because of their size, but due to their emotionally steady delivery. For example, Llama 3.2 1B ZS offered, "*It's nothing to be embarrassed about... Would you like to try a grounding technique?*", and Qwen 2.5 1.5B FT reassured users with, "*Absolutely... It's not your fault... Try sitting by a window... You can do this.*" These responses validated emotion, normalized experience, and offered simple coping strategies. In contrast, lower-rated models like Qwen 2.5 0.5B FT and ZS produced vague or awkward phrasing (e.g., "*Human experience can be frustrating*", or the unclear phrase "*Humanning the past*") or introduced irrelevant content. This indicates that older users favor responses that are emotionally anchored and practical—suggesting design priorities should lean toward clarity and support over abstraction.

Education: Empathy + Practicality? Graduate-educated participants tended to reward responses that balanced emotional validation with actionable suggestions. For instance, Sonnet 3.5 received a high average empathy rating of 4.48 from graduate users in a high-stress scenario, though the difference was not statistically significant ($p = 0.28$). However, robust differences did emerge for certain models: Qwen 2.5 3B and Qwen 2.5 1.5B were rated significantly lower by graduate participants compared to non-graduates ($p < 0.001$; Cohen's $d = -0.26$ and -0.20 , respectively), suggesting sensitivity to vague or generic phrasing. Both results remained significant after Bonferroni correction for multiple comparisons across nine models (adjusted $\alpha = 0.0056$). Conversely, while Phi-3.5 Mini showed no overall difference across education groups ($p = 0.96$), one response—specifically for Prompt Q17—trended toward significantly higher ratings among graduates ($p = 0.088$, Cohen's d

$= 0.26$), pointing to a preference for steady, well-calibrated delivery in emotionally complex situations. These findings suggest that educational background can shape empathy perception in intricate ways—particularly in edge cases—highlighting the need for adaptable, audience-aware generation.

Does Gender Shape Empathy Ratings? Empathy ratings from men ($n = 49$) and women ($n = 65$) showed no significant differences across all models, including top performers like Sonnet 3.5 and Llama 3.2 3B (ZS). No statistically significant differences emerged between gender groups ($p > 0.15$ for all, remaining non-significant after Bonferroni correction), and effect sizes were negligible ($|d| < 0.05$). While qualitative impressions suggested women slightly preferred gentler tones (e.g., DeepSeek-R1 ZS) and men favored more directive styles (e.g., DeepSeek-R1 FT), these were not reflected in overall empathy scores. Notably, Qwen 2.5 variants were *rated consistently lower* by both groups, with no gender difference in mean ratings ($p = 0.37$ – 0.99 ; Cohen's $d < 0.05$); for example, Qwen 2.5 3B averaged 3.14 (men) and 3.20 (women). This underscores that expectations for emotional clarity and coherence in therapeutic contexts are shared, regardless of gender.

Regional Trends in Empathy Ratings Figure 3 presents empathy ratings across U.S. regions for various mental health scenarios. For acute distress (e.g., panic, abuse), participants in the Western US gave higher ratings than other regions (Mean = 3.75; $F(3, N) = 3.51$, $p = 0.015$), which remained significant after Bonferroni correction for multiple scenarios ($\alpha = 0.0167$) suggesting possible heightened sensitivity to emotional urgency. For depression and grief, mean ratings were highest in the South (3.62) and Midwest (3.46), but these differences were not statistically significant ($F = 0.78$, $p = 0.505$). While these patterns are noteworthy, our sample size limits interpretation, and causal mechanisms were not explored. These observations suggest that future research should further investigate the potential for culturally and regionally sensitive design in empathetic AI.

Summary: While emotional validation, contextual grounding, and clear next steps are universally valued, empathetic delivery must be tailored to user profiles. Effective empathetic systems need to be context- and user-aware, adapting not just to the scenario but also to individual preferences for com-

fort, care, and counsel.

6.6 One Size Model Does Not Fit All

Our analysis indicates that no single model excels across all emotional scenarios—empathy perception is highly context-dependent, with different models performing best in different contexts. In Scenario 1 (panic attack), Claude Ref-QType1 led with a mean of 4.38, closely followed by DeepSeek-R1 ZS (4.25) and FT (4.06), suggesting both reference and instruction-tuned models resonated for acute anxiety. Scenario 2 (grief/anger) saw Llama 3.2 1B ZS (3.57) and Qwen 2.5 0.5B FT (3.65) outperform larger models, likely due to their concise, emotionally direct responses—demonstrating that smaller models can deliver contextually attuned empathy. For Scenario 3 (depression), Granite 3.1 FT (4.04), Claude Ref-QType1 (3.92), and Llama 3.2 1B ZS (3.88) were top-rated, often emphasizing small, achievable actions. In Scenario 4 (reliving abuse), Granite 3.1 ZS (3.72), Llama 3.2 1B ZS (3.67), and Claude Ref-QType4 (3.69) stood out, with grounding-focused, emotionally cautious responses likely contributing to their effectiveness. Scenario 5 (social anxiety) favored Llama 3.2 3B FT (3.71), Llama 1B FT (3.68), and Claude Ref-QType1 (3.71), suggesting that fine-tuned Llama models were especially effective for relatable coping strategies. Finally, in Scenario 6 (substance use), Llama 3.2 3B ZS (3.67), 1B ZS (3.66), and DeepSeek-R1 ZS (3.53) performed best, balancing emotional validation with gentle redirection. These results highlight that contextual appropriateness is paramount. Given challenges in identifying user subgroups, models with broader sensitivity may currently offer the most practical benefit.

6.7 Does Race Influence Empathy?

As shown in Table 2, 65.5% of participants identified as White or Caucasian, and 19.8% as Black or African American. Due to limited representation of other racial identities, we grouped participants as White vs. non-White for exploratory analysis. While this binary approach flattens cultural distinctions (Sap et al., 2022; May et al., 2019), it enabled detection of broad trends while preserving statistical power. Across models, non-White participants consistently assigned higher empathy ratings. For example, Sonnet 3.5 received a mean of 3.81 from non-White participants versus 3.53 from White participants; similar trends were observed for Granite 3.1 ZS (3.68 vs. 3.49) and DeepSeek-R1 ZS

(3.56 vs. 3.50), with even lower-rated Qwen variants showing this uplift. Although these differences were not statistically significant and should be interpreted with caution, the consistent pattern suggests empathy perception may be influenced by lived experience. All participants were native English speakers (Section B), reducing linguistic ambiguity but centering the study within a Western, English-speaking context. Future work should examine these trends in more racially and linguistically diverse samples and explore culturally responsive strategies for empathetic AI design.

7 Future Work

We plan to release the TIDE dataset to advance research on generative models and privacy-preserving AI for mental health. Future work will explore its use in training and evaluating models for therapeutic goals, such as supporting junior therapists and enabling broader clinical applications. Crucially, it remains important to assess whether differences between real and synthetic conversations (e.g., B N et al. (2025)) correlate with patient outcomes—such as distress reduction or recovery—or at minimum, confirm they do not negatively impact therapeutic effectiveness. Such analysis will clarify the practical utility of synthetic data in clinical settings.

8 Conclusions

This study demonstrates that small language models (0.5B–5B parameters) can be substantially improved to generate empathetic responses for PTSD support contexts through targeted fine-tuning on trauma-informed synthetic data. Using the TIDE dataset and a multi-pronged evaluation strategy, we show that fine-tuning boosts automatic similarity metrics and enhances human-perceived empathy, although gains are not uniform across models or scenarios. Our findings emphasize that small, resource-efficient models can meaningfully engage in emotionally supportive dialogue when carefully adapted, offering promising pathways for scalable, privacy-preserving mental health support systems.

9 Acknowledgement

This work is supported by the National Science Foundation (NSF) under Grant No. (Redacted for Review). Any opinions, findings, and conclusions or recommendations expressed in this work are those of the author(s) and do not necessarily reflect the views of the NSF.

Limitations

While our findings show that small FT models can support empathetic generation, several limitations remain. First, the synthetic TIDE dataset, though realistic and emotionally appropriate, may not reflect the full complexity of real-world interactions; future work should incorporate clinical transcripts and cross-cultural data. Second, synthetic data may limit generalizability due to reduced emotional and cultural variance. Third, smaller models occasionally produced generic or overly affirmative replies, with the phrase “Have you considered joining a support group” appearing over 10 times across conditions—suggesting surface-level empathy over contextual refinement.

References

Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. In *International Conference on Machine Learning*, pages 337–371. PMLR.

Rawan AlMakinah, Andrea Norcini-Pala, Lindsey Disney, and M Abdullah Canbaz. 2024. Enhancing mental health support through human-ai collaboration: Toward secure and empathetic ai-enabled chatbots. *arXiv preprint arXiv:2410.02783*.

Anthropic. 2024. *Claude 3.5 sonnet*. *Anthropic Blog*.

Suhas B N, Dominik Mattioli, Saeed Abdullah, Rosa I Arriaga, Chris W Wiese, and Andrew M Sherrill. 2025. How real are synthetic therapy conversations? evaluating fidelity in prolonged exposure dialogues. *arXiv preprint arXiv:2504.21800*.

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for MT and/or Summarization*, pages 65–72.

Timothy W Bickmore and Rosalind W Picard. 2005. Establishing and maintaining long-term human-computer relationships. *ACM Transactions on Computer-Human Interaction (TOCHI)*, 12(2):293–327.

Suhas BN and Saeed Abdullah. 2022. Privacy sensitive speech analysis using federated learning to assess depression. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6272–6276. IEEE.

Suhas BN, Sarah Rajtmajer, and Saeed Abdullah. 2023. *Differential privacy enabled dementia classification: An exploration of the privacy-accuracy trade-off in speech signal data*. In *Proceedings of Interspeech 2023*, pages 346–350, Dublin, Ireland. ISCA.

Suhas BN, Andrew M. Sherrill, Rosa I. Arriaga, Chris W. Wiese, and Saeed Abdullah. 2025. *Thousand voices of trauma: A large-scale synthetic dataset for modeling prolonged exposure therapy conversations*. *Preprint*, arXiv:2504.13955.

Jeremy Chang, Kuan-Yu Chen, and Chung-Hsien Wu. 2024. *Applying reinforcement learning and multi-generators for stage transition in an emotional support dialogue system*.

Zhiyu Chen, Yujie Lu, and William Yang Wang. 2023. Empowering psychotherapy with large language models: Cognitive distortion detection through diagnosis of thought prompting. *arXiv preprint arXiv:2310.07146*.

Hyojin Chin, Hyeonho Song, Gumhee Baek, Mingi Shin, Chani Jung, Meeyoung Cha, Junghoi Choi, and Chiyoung Cha. 2023. The potential of chatbots for emotional support and promoting mental well-being in different cultures: mixed methods study. *Journal of Medical Internet Research*, 25:e51712.

Eugene Cho, Nasim Motalebi, S Shyam Sundar, and Saeed Abdullah. 2022. Alexa as an active listener: how backchanneling can elicit self-disclosure and promote user experience. *Proceedings of the acm on human-computer interaction*, 6(CSCW2):1–23.

Mark H. Davis. 1983. Measuring individual differences in empathy: Evidence for a multidimensional approach. *Journal of Personality and Social Psychology*, 44(1):113–126.

Mauro De Gennaro, Eva G Krumhuber, and Gale Lucas. 2020. Effectiveness of an empathic chatbot in combating adverse effects of social exclusion on mood. *Frontiers in psychology*, 10:495952.

Daniil Filienko, Yinzhou Wang, Caroline El Jazmi, Serena Xie, Trevor Cohen, Martine De Cock, and Weichao Yuwen. 2024. Toward large language models as a therapeutic tool: Comparing prompting techniques to improve gpt-delivered problem-solving therapy. *arXiv preprint arXiv:2409.00112*.

Joanne Frattaroli. 2006. Experimental disclosure and its moderators: a meta-analysis. *Psychological bulletin*, 132(6):823.

Saadiah Gabriel, Isha Puri, Xuhai Xu, Matteo Malgaroli, and Marzyeh Ghassemi. 2024. Can ai relate: Testing large language model response for mental health support. *arXiv preprint arXiv:2405.12021*.

Daniel Goleman. 2006. *Social Intelligence: The New Science of Human Relationships*. Bantam Books.

Christina M Hassija and Jessica A Turchik. 2016. An examination of disclosure, mental health treatment use, and posttraumatic growth among college women who experienced sexual victimization. *Journal of loss and trauma*, 21(2):124–136.

838	Scott H Hemenover. 2003. The good, the bad, and the healthy: Impacts of emotional disclosure of trauma on resilient self-concept and psychological distress. <i>Personality and Social Psychology Bulletin</i> , 29(10):1236–1244.	891
839		892
840		893
841		
842		
843	Jia-Hao Hsu, Jeremy Chang, Min-Hsueh Kuo, and C. Wu. 2023. Empathetic response generation based on plug-and-play mechanism with empathy perturbation.	894
844		895
845		896
846		897
847	Yining Hua, Fenglin Liu, Kailai Yang, Zehan Li, Hongbin Na, Yi-han Sheu, Peilin Zhou, Lauren V Moran, Sophia Ananiadou, Andrew Beam, et al. 2024. Large language models in mental health care: a scoping review. <i>arXiv preprint arXiv:2401.02984</i> .	898
848		899
849		900
850		901
851		902
852	Zainab Iftikhar, Sean Ransom, Amy Xiao, and Jeff Huang. 2024. Therapy as an nlp task: Psychologists’ comparison of llms and human peers in cbt. <i>arXiv preprint arXiv:2409.02244</i> .	903
853		904
854		905
855		
856	Claire T Kimbley, Daniel W Cox, Jeffrey H Kahn, and Keith D Renshaw. 2023. Feeling pressured to talk about trauma: How pressure to disclose alters the association between trauma disclosure and posttraumatic growth. <i>Journal of Traumatic Stress</i> , 36(3):567–578.	906
857		907
858		908
859		909
860		910
861		
862	Theodora Kokosi and Katie Harron. 2022. Synthetic data in medical research. <i>BMJ medicine</i> , 1(1).	911
863		912
864	Yoon Kyung Lee, Inju Lee, Minjung Shin, Seoyeon Bae, and Sowon Hahn. 2023. Chain of empathy: Enhancing empathetic response of large language models based on psychotherapy models. <i>arXiv preprint arXiv:2311.04915</i> .	913
865		914
866		915
867		
868		
869	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A diversity-promoting objective function for neural conversation models. In <i>Proceedings of NAACL-HLT</i> , pages 110–119.	916
870		917
871		918
872		
873	Wenjie Li, Tianyu Sun, Kun Qian, and Wenhong Wang. 2024. Optimizing psychological counseling with instruction-tuned large language models. <i>arXiv preprint arXiv:2406.13617</i> .	919
874		920
875		921
876		922
877	Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In <i>Text Summarization Branches Out: Proceedings of the ACL-04 Workshop</i> , pages 74–81. Association for Computational Linguistics.	923
878		924
879		
880		
881		
882	Siyuan Brandon Loh and Aravind Sesagiri Raamkumar. 2023. Harnessing large language models’ empathetic response generation capabilities for online mental health counselling support. <i>arXiv preprint arXiv:2310.08017</i> .	925
883		926
884		927
885		928
886		929
887	Gale M Lucas, Jonathan Gratch, Aisha King, and Louis-Philippe Morency. 2014. It’s only a computer: Virtual humans increase willingness to disclose. <i>Computers in Human Behavior</i> , 37:94–100.	930
888		931
889		932
890		
	Xiaoyang Ma and Weiqi Chang. 2024. A two-stage emotional dialogue generation model based on dialogpt.	933
	Alexander Marrapese, Basem Suleiman, Imdad Ullah, and Juno Kim. 2024. A novel nuanced conversation evaluation framework for large language models in mental health. <i>arXiv preprint arXiv:2403.09705</i> .	934
		935
	Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.	936
		937
	R. Morris, Kareem Kouddous, Rohan Kshirsagar, and S. Schueller. 2018. Towards an artificially empathic conversational agent for mental health applications: System design and user perceptions. <i>Journal of Medical Internet Research</i> , 20.	938
		939
		940
	National Institute of Mental Health. 2023. Post-traumatic stress disorder (ptsd). https://www.nimh.nih.gov/health/publications/post-traumatic-stress-disorder-ptsd . Accessed: 2025-05-08.	941
		942
	Stefan Palan and Christian Schitter. 2018. Prolific.ac—a subject pool for online experiments. <i>Journal of Behavioral and Experimental Finance</i> , 17:22–27.	943
		944
	Jung In Park, Mahyar Abbasian, Iman Azimi, Dawn T Bounds, Angela Jun, Jaesu Han, Robert M McCarron, Jessica Borelli, Parmida Safavi, Sanaz Mirbaha, et al. 2024. Building trust in mental health chatbots: safety metrics and llm-based evaluation tools. <i>arXiv preprint arXiv:2408.04650</i> .	945
	Yushan Qian, Wei-Nan Zhang, and Ting Liu. 2023. Harnessing the power of large language models for empathetic response generation: Empirical investigations and improvements. <i>arXiv preprint arXiv:2310.05140</i> .	946
	Aravind Sesagiri Raamkumar and Siyuan Brandon Loh. 2024. Towards a multidimensional evaluation framework for empathetic conversational systems.	947
	Charles Rabin and Joshua Fogel. 2009. <i>Cognitive Behavioral Group Therapy for PTSD: A Treatment Manual</i> . Routledge.	948
		949
	Abdur Rasool, Muhammad Irfan Shahzad, Hafsa Aslam, Vincent Chan, and Muhammad Ali Arshad. 2025. Emotion-aware embedding fusion in large language models (flan-t5, llama 2, deepseek-r1, and chatgpt 4) for intelligent response generation. <i>AI</i> , 6(3):56.	950
		951
	Nils Reimers and Iryna Gurevych. 2021. all-minilm-l6-v2. https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2 . Sentence Transformers model.	952
		953
		954

945	Maarten Sap, Swabha Swayamdipta, Laura Vianna,	Yuqi Wu, Jie Chen, Kaining Mao, and Yanbo Zhang.	1000
946	Xuhui Zhou, Yejin Choi, and Noah A. Smith. 2022.	2023. Automatic post-traumatic stress disorder diag-	1001
947	Annotators with attitudes: How annotator beliefs	nosis via clinical transcripts: A novel text augmenta-	1002
948	and identities bias toxic language detection.	tion with large language models.	1003
949	In <i>Proceedings of the 2022 Conference of the North Amer-</i>		
950	<i>ican Chapter of the Association for Computational</i>	Zhichao Xu and Jiepu Jiang. 2024. Multi-dimensional	1004
951	<i>Linguistics: Human Language Technologies</i> , pages	evaluation of empathetic dialog responses. <i>arXiv</i>	1005
952	5884–5906, Seattle, United States. Association for	<i>preprint arXiv:2402.11409.</i>	1006
953	Computational Linguistics.		
954	Georg Schomerus, Stephanie Schindler, Theresia	Jiamin Yang and David Jurgens. 2024. Modeling em-	1007
955	Rechenberg, Tobias Gfesser, Hans J Grabe, Mario	pathetic alignment in conversation. <i>arXiv preprint</i>	1008
956	Liebergesell, Christian Sander, Christine Ulke, and	<i>arXiv:2405.00948.</i>	1009
957	Sven Speerforck. 2021. Stigma as a barrier to		
958	addressing childhood trauma in conversation with	Zhou Yang, Zhaochun Ren, Yufeng Wang, Shizhong	1010
959	trauma survivors: A study in the general population.	Peng, Haizhou Sun, Xiaofei Zhu, and Xiangwen Liao.	1011
960	<i>Plos one</i> , 16(10):e0258782.	2024. Enhancing empathetic response generation by	1012
		augmenting llms with small-scale empathetic models.	1013
961	Antonia V. Seligowski, Daniel J. Lee, Joseph R.	Hong Qing Yu and Stephen McGuinness. 2024. An	1014
962	Bardeen, and Holly K. Orcutt. 2015. Emotion reg-	experimental study of integrating fine-tuned large	1015
963	ulation and posttraumatic stress symptoms: a meta-	language models and prompts for enhancing mental	1016
964	analysis. <i>Cognitive Behaviour Therapy</i> , 44(2):87–	health support chatbot system.	1017
965	102.		
966	Ashish Sharma, Inna Wanyin Lin, Adam S. Miner,	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q	1018
967	David C. Atkins, and Tim Althoff. 2021. Towards	Weinberger, and Yoav Artzi. 2019. Bertscore: Eval-	1019
968	facilitating empathic conversations in online mental	uating text generation with bert. <i>arXiv preprint</i>	1020
969	health support: A reinforcement learning approach.	<i>arXiv:1904.09675.</i>	1021
970	Ashish Sharma, Adam S Miner, David C Atkins, and	Chujie Zheng, Yong Liu, Wei Chen, Yongcai Leng, and	1022
971	Tim Althoff. 2020. A computational approach to un-	Minlie Huang. 2021. Comae: A multi-factor hierar-	1023
972	derstanding empathy expressed in text-based mental	chical framework for empathetic response generation.	1024
973	health support. <i>arXiv preprint arXiv:2009.08441.</i>	<i>arXiv preprint arXiv:2105.08316.</i>	1025
974	Yun Su, Bozhen Fan, Haoran Bian, Yunhao Zhu, and		
975	Runhe Huang. 2023. Empathetic dialogue generation		
976	model with multi-source emotion and intent percep-		
977	tion.		
978	Talha Tahir. 2024. Fine tuning large language models		
979	to deliver cbt for depression.		
980	Holly Timblin and Christina M Hassija. 2023. How will		
981	i be perceived: the role of trauma-related shame in		
982	the relationship between psychological distress and		
983	expectations of disclosure among survivors of sexual		
984	victimization. <i>Journal of interpersonal violence</i> ,		
985	38(7-8):5805–5823.		
986	Mrigank Tiwari, V. Dahiya, Om Mohanty, and Girija		
987	Saride. 2022. Iaemp: Intent-aware empathetic re-		
988	sponse generation.		
989	Ruiyi Wang, Stephanie Milani, Jamie C Chiu, Jiayin		
990	Zhi, Shaun M Eack, Travis Labrum, Samuel M Mur-		
991	phy, Nev Jones, Kate Hardy, Hong Shen, et al. 2024.		
992	Patient-{\Psi}: Using large language models to sim-		
993	ulate patients for training mental health professionals.		
994	<i>arXiv preprint arXiv:2405.19660.</i>		
995	Anuradha Welivita, Chun-Hung Yeh, and Pearl Pu. 2023.		
996	Empathetic response generation for distress support.		
997	In <i>Proceedings of the 24th Annual Meeting of the</i>		
998	<i>Special Interest Group on Discourse and Dialogue</i> ,		
999	pages 632–644.		

1026
1027
1028
1029

A Survey Scenarios for Empathy Rating Study

Note: Each set of four questions (e.g., Q.1–Q.4) corresponds to a single scenario as shown in Figure 2. The model-to-response mappings are provided in the main text and in the figure. The options were randomized in the survey. Please refer to Table 1 for the question/option mapping.

Q.1:

Person: "I had a panic attack at the grocery store. It was so embarrassing."

Question 1 (Row 2619, Pair 1):

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Panic attacks can be really scary and overwhelming. Remember, it's not your fault.

→ Have you informed your workplace about your PTSD? They might be able to offer some support."

Response 2: "Panic attacks can be really scary and overwhelming. Please don't be hard on

→ yourself. Have you informed your workplace about your condition? They might be able to

→ provide some support."

Response 3: "Panic attacks can be scary and overwhelming, but you're not alone. Have you learned

→ any coping strategies for public emergencies? That could help in such situations in the

→ future."

Response 4: "Panic attacks can be embarrassing, but remember, they're not your fault. Have you

→ considered explaining to your employer about your PTSD? They might be understanding."

Response 5: "I'm sorry that happened. Panic attacks can be really scary and overwhelming. Please

→ know it's not your fault, and many people experience them. Have you learned any grounding

→ techniques that might help in those moments?"

Q.2:

Person: "I had a panic attack at the grocery store. It was so embarrassing."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "average person would have done that. You're strong for knowing when to seek help.

→ Have you learned any grounding techniques to use when you feel the urge to use? They can be

→ really helpful. I'd be happy to provide some resources if you're interested."

Response 2: "introverted guy like you, it can be isolating. But remember, panic attacks are a

→ common symptom of PTSD. Have you considered joining a support group for people with similar

→ experiences? It might help you feel less alone."

Response 3: "adult-onset PTSD can be really scary, but it'sNot at all embarrassing. Panic attacks

→ are a common symptom of PTSD, especially for those who've experienced physical abuse. You did

→ your best in a difficult situation. Have you learned any grounding techniques to help manage

→ these episodes? They can be really"

Response 4: "absolutely, panic attacks can be scary and overwhelming. It's not your fault. Next

→ time, try sitting by a window or in a quieter area with supportive people around. Remember,

→ panic attacks are a normal part of PTSD and anxiety. You've survived many scary moments. You

→ can do this. Have"

Response 5: "I'm sorry that happened. Panic attacks can be really scary and overwhelming. Please

→ know it's not your fault, and many people experience them. Have you learned any grounding

→ techniques that might help in those moments?"

Q.3:

Person: "I had a panic attack at the grocery store. It was so embarrassing."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Human experience can be frustrating. PTSD can make it challenging to control our

→ feelings. Maybe try a breathing exercise or focus on your surroundings to help calm yourself.

→ Remember, it's okay to take things easier after recovery."

Response 2: "Humanning the past is tough, but it's important to do what feels safe and

→ responsible for you. Have you considered seeking help from a therapist who specializes in

→ PTSD? Your safety might improve with their guidance. Remember, healing is a process, and what

→ works for you may be different than what works for others"

Response 3: "I'm sorry you had such a difficult experience. Please know it's not your fault. Have

→ you considered joining a support group for people with PTSD? It might help you feel safer in

→ public. If you're ready, I can suggest some resources. You're not alone in this"

Response 4: "I'm sorry that happened. Panic attacks can be scary and overwhelming. Have you

→ learned any grounding techniques to use during panic attacks? They might help you stay calm

→ in the moment. Remember, it's not your fault. You're safe at home. Have you considered

→ talking to a therapist about coping strategies?"

Response 5: "I'm sorry that happened. Panic attacks can be really scary and overwhelming. Please
→ know it's not your fault, and many people experience them. Have you learned any grounding
→ techniques that might help in those moments?"

Q.4:

Person: "I had a panic attack at the grocery store. It was so embarrassing."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "I'm sorry that happened to you. Panic attacks can be very embarrassing, but it's not
→ your fault. Would you like to share some strategies for managing them in public? It might
→ help you feel more in control next time. Maybe we can practice a few deep breathing exercises
→ together? That can"

Response 2: "I'm sorry that happened. Remember, panic attacks can be unpredictable. Next time,
→ try going with a friend or family member. It might make it easier to manage your symptoms.
→ Also, you can always call your therapist for support. They might have strategies to help you
→ feel more comfortable in public"

Response 3: "I'm sorry that happened. Panic attacks can be really scary and overwhelming. It's
→ nothing to be embarrassed about. Have you learned any grounding techniques to use in these
→ situations? They can help bring you back to the present and reduce the impact of the panic.
→ Would you like to try one right now?"

Response 4: "That must have been really difficult. Remember, panic attacks are a common symptom
→ of PTSD. It's not your fault, and it doesn't mean you're a failure. Have you learned any
→ grounding techniques to help manage these episodes? They can be really helpful. Would you
→ like to explore some strategies?"

Response 5: "I'm sorry that happened. Panic attacks can be really scary and overwhelming. Please
→ know it's not your fault, and many people experience them. Have you learned any grounding
→ techniques that might help in those moments?"

Q.5:

Person: "Sometimes I feel angry at her for leaving me alone like this."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's normal to feel angry, but remember that her leaving wasn't your fault. Have you
→ considered joining a support group for widowers or others who've lost a partner?"

Response 2: "It's natural to feel angry, but remember that she left. It's not your fault. Have
→ you talked to her about how you're feeling?"

Response 3: "Childhood trauma can affect our relationships. It's okay to feel these emotions.
→ Have you considered talking to your partner about your feelings? Open communication can
→ strengthen your bond.

Response 4: "Aggression can be a normal response to trauma. It's important to find healthy ways
→ to express and process these feelings. Have you considered talking to your wife about your
→ experiences? Open communication might help you both heal."

Response 5: "It's okay to have complex emotions about loss. Anger is a normal part of grief. Have
→ you considered joining a support group for widowers? Sharing with others who understand could
→ help."

Q.6:

Person: "Sometimes I feel angry at her for leaving me alone like this."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "adult children dealing with grief is really tough. Your anger is a common response
→ to loss, but remember that you're not to blame. Have you considered joining a grief support
→ group for survivors? Connecting with others who've had similar experiences can be really
→ helpful."

Response 2: "adult children dealing with trauma? That's really tough. It's okay to feel angry.
→ Your feelings are valid. Remember, it's not your fault. Have you considered talking to a
→ therapist? They could help you navigate these complex emotions."

Response 3: "adult children of narcissistic parents often struggle with anger. It's a natural
→ response given the circumstances. Have you tried expressing these feelings to your parent or
→ a therapist? It might help to work through them in a safe way."

Response 4: "adult-onset depression can make difficult emotions feel more intense.

Response 5: "It's okay to have complex emotions about loss. Anger is a normal part of grief. Have
→ you considered joining a support group for widowers? Sharing with others who understand could
→ help."

Q.7:

Person: "Sometimes I feel angry at her for leaving me alone like this."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Humanism can be a helpful approach when dealing with PTSD. Have you considered

→ practicing mindfulness meditation? It might help you become more present and aware of your

→ feelings. Also, being honest with yourself about your limitations can be healing. Your

→ feelings are valid, and they don't have to define you. Perhaps"

Response 2: "In your life, it's common to feel angry. What's happening is that you're navigating

→ a difficult situation. You're showing kindness and resilience, and that's something positive

→ to celebrate. Remember, you made a choice to help yourself. Let's explore these feelings

→ together, maybe even brainstorming ways"

Response 3: "It's understandable to feel angry given your past experiences. Your feelings are

→ valid. Have you considered joining a support group for South Asian women with PTSD? It might

→ help you connect with others who've experienced similar situations. Have you tried

→ discussing this with your partner? Open communication can strengthen your relationship. If"

Response 4: "It's understandable to feel conflicted. Have you tried joining a support group for

→ professionals? It might help you connect with others who share similar experiences. Your

→ feelings are valid, and healing is possible. Would you like to learn more about support

→ groups? I'd be happy to provide information."

Response 5: "It's okay to have complex emotions about loss. Anger is a normal part of grief. Have

→ you considered joining a support group for widowers? Sharing with others who understand could

→ help."

Q.8:

Person: "Sometimes I feel angry at her for leaving me alone like this."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Rejection can be a painful experience. Your feelings are understandable, but

→ remember that it's not her fault. Have you considered reaching out to her or inviting her to

→ a therapy session to work through these feelings? It might help to have someone who

→ understands on your side."

Response 2: "Rejection can be deeply painful. It's okay to feel angry, but remember that your

→ wife made a difficult choice. Have you considered joining a support group for widowers? It

→ might help you process these feelings. Your wife would want you to heal and find joy again.

→ What's one thing"

Response 3: "It's understandable to feel angry, but remember that your wife would want you to

→ feel safe and cared for. Have you considered joining a support group for widowers? It might

→ help you feel less alone in your feelings. Would you like some suggestions for groups in your

→ area? I can help you"

Response 4: "Your feelings are valid, but remember that your wife would want you to feel safe.

→ Anger can be a symptom of PTSD, but it doesn't define you. Have you considered talking to her

→ about your feelings? Open communication can strengthen your bond and help her understand what

→ you're going through better"

Response 5: "It's okay to have complex emotions about loss. Anger is a normal part of grief. Have

→ you considered joining a support group for widowers? Sharing with others who understand could

→ help."

Q.9:

Person: "The depression is hitting hard today. I can barely get out of bed."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's okay to have tough days. Have you tried any gentle stretches or breathing

→ exercises? Sometimes small movements can make a difference."

Response 2: "I'm so sorry you're feeling this way. At least you're safe now. Can we focus on

→ getting one small task done today, like opening a window or checking your phone? Baby steps."

Response 3: "I'm sorry you're feeling so low. Depression can be draining. Have you tried any

→ self-care activities today? Even small steps can make a difference."

Response 4: "I'm here for you. Let's focus on one small, manageable task today. Can you think of

→ something you could do right now, no matter how small?"

Response 5: "I'm sorry you're having such a difficult day. Depression can be overwhelming. Can

→ you think of one small, achievable task to focus on? Even getting dressed can be a victory."

Q.10:

Person: "The depression is hitting hard today. I can barely get out of bed."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "amazing initiative! Exercise can definitely help with depression. Even a short walk

→ outside might make a difference. Remember, small steps are okay. Maybe start with 5 minutes

→ and build from there? If you're unsure, you could also try gentle stretching at home. Your

→ body will appreciate the movement."

Response 2: "amazing that you're reaching out. Depression can be really tough, especially on top

→ of PTSD. Remember to be gentle with yourself. Even small steps, like getting up for a short

→ walk or journaling, can make a difference. Would you like some resources on managing

→ depression alongside PTSD?"

Response 3: "important to remember that you're dealing with depression on top of PTSD. Have you

→ spoken to your doctor about adjusting your treatment plan? Also, don't hesitate to reach out

→ to a support group or counselor for additional help. You don't have to navigate this alone.

→ I'm just here to listen if"

Response 4: "important. Depression can make us feel weak, but healing takes strength too. You've

→ already shown incredible strength by reaching out. What's one small thing you could do right

→ now to care for yourself? Maybe a quick shower or break for some fresh air? Small steps can

→ make a big difference in managing"

Response 5: "I'm sorry you're having such a difficult day. Depression can be overwhelming. Can

→ you think of one small, achievable task to focus on? Even getting dressed can be a victory."

Q.11:

Person: "The depression is hitting hard today. I can barely get out of bed."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Human-powered devices can be helpful, especially when you're feeling overwhelmed.

→ Have you considered connecting with a mental health support group for older adults? They may

→ have resources tailored to your situation. Remember, you're not alone in this; you're a

→ strong human being. Good luck, and I hope you"

Response 2: "Human anxiety can affect sleep quality and focus. Are you feeling overwhelmed by the

→ situation? Sometimes, just taking short breaks can help. Could you try some grounding

→ techniques before bed to calm your mind? Remember, these are temporary measures that might

→ not work every time. Your health matters."

Response 3: "I'm sorry you're feeling this way. Depression can be overwhelming. Have you tried

→ any relaxation techniques before bed? Sometimes deep breathing or guided imagery can help

→ calm your mind. If not, there are apps designed to help manage both depression and PTSD. Have

→ you spoken to your doctor about adjusting your treatment"

Response 4: "I'm sorry you're feeling so down. Depression can be overwhelming. Remember, you're

→ strong at 62. Can you think of one small, manageable task you could tackle right now? Even

→ something small like getting a cup of coffee might help. Your resilience is remarkable. Have

→ you considered reaching"

Response 5: "I'm sorry you're having such a difficult day. Depression can be overwhelming. Can

→ you think of one small, achievable task to focus on? Even getting dressed can be a victory."

Q.12:

Person: "The depression is hitting hard today. I can barely get out of bed."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's understandable to feel discouraged. Remember, you're not alone. Can you reach

→ out to a family member or support group for a virtual chat or call? Sometimes just connecting

→ with others can help lift your spirits. Are there any self-care activities you enjoy? Even

→ small steps can make a difference"

Response 2: "It's understandable that you're struggling. Have you spoken with your therapist

→ about managing depression? They might have some strategies to help you stay productive.

→ Remember, you're not alone in this battle. Your support system is there for you. Have you

→ tried reaching out to a friend or family member today"

Response 3: "I'm sorry you're feeling so bad. Depression can be overwhelming. Let's try a quick

→ mindfulness exercise together. Can you stand up straight and focus on your breath for 5 deep

→ breaths? That can help ground you in the present moment and ease some of the symptoms. Is

→ there someone"

Response 4: "I'm sorry you're feeling so down. Depression can be really challenging. Have you
→ considered talking to a therapist about coping strategies? They might have some suggestions
→ for helping you feel better. Do you have a list of therapists in your area that specialize in
→ depression? It might help to narrow down your"

Response 5: "I'm sorry you're having such a difficult day. Depression can be overwhelming. Can
→ you think of one small, achievable task to focus on? Even getting dressed can be a victory."

Q.13:

Person: "Sometimes I feel like I'm reliving the abuse all over again."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "That sounds incredibly difficult. Remember, you're safe now. Can you try some
→ grounding techniques like focusing on your breath or naming five things you can see?"

Response 2: "That's a common PTSD symptom. Have you tried any grounding techniques to help you
→ stay present? Something simple like focusing on your breath or naming objects around you can
→ help."

Response 3: "That sounds really distressing. Remember, you're safe now. Can you tell me about a
→ time when you felt most loved and safe? That could help ground you in the present."

Response 4: "That sounds really difficult. It's okay to feel this way. Have you tried mindfulness
→ exercises to help you stay present and separate from those past experiences?"

Response 5: "That must be incredibly difficult. Remember, you're safe now. When you have these
→ feelings, try grounding techniques to remind yourself of the present moment."

Q.14:

Person: "Sometimes I feel like I'm reliving the abuse all over again."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "emotionally, that must be terrifying. Grounding techniques can help during
→ flashbacks. Can you try naming 5 things you see, 4 you touch, 3 you hear, 2 you smell, and 1
→ you taste right now? It might help bring you back to the present moment."

Response 2: "adult survivors often struggle with intrusive thoughts. This is your mind trying to
→ process what happened. Have you considered talking to a therapist who specializes in
→ trauma-related PTSD? They can provide strategies to help manage these intrusive thoughts."

Response 3: "experienced can be overwhelming, but remember that you're safe now. Try grounding
→ yourself by focusing on your five senses in the present moment. What do you see, hear, touch,
→ smell, and feel right now? This can help remind you that you're not in danger and you're safe
→ in the"

Response 4: "adult survivors are not to blame for the abuse. Your feelings are valid, but please
→ remember that you survived and that the trauma is not your fault. Have you tried any
→ grounding techniques to help when you feel the urge to relive the past? I can suggest some if
→ you'd like."

Response 5: "That must be incredibly difficult. Remember, you're safe now. When you have these
→ feelings, try grounding techniques to remind yourself of the present moment."

Q.15:

Person: "Sometimes I feel like I'm reliving the abuse all over again."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Human abuse can have lasting effects. It's important to remember that you're not a
→ victim. With your support, you can work towards healing. Remember, you're not alone in this
→ journey. Do you have any coping strategies that have helped in the past?"

Response 2: "Human experience can be overwhelming. Have you tried any relaxation techniques
→ before bedtime? Some people find that deep breathing or gentle yoga helps reduce anxiety."

Response 3: "You're not reliving it. You're processing it. Have you tried any grounding
→ techniques to help when these flashbacks occur? It might help to focus on your surroundings.
→ Would you like to try some?"

Response 4: "Your feelings are valid. Compulsive behaviors like this can be a way of seeking
→ validation. Have you considered talking to a therapist about this? They might offer
→ strategies to manage these compulsions."

Response 5: "That must be incredibly difficult. Remember, you're safe now. When you have these
→ feelings, try grounding techniques to remind yourself of the present moment."

Q.16:

Person: "Sometimes I feel like I'm reliving the abuse all over again."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's common for flashbacks to occur, especially with PTSD. When this happens, try
→ grounding yourself by focusing on your surroundings. What do you see? What do you hear? This
→ can help bring you back to the present moment."

Response 2: "That sounds really difficult. Flashbacks can be a symptom of PTSD. Have you tried
→ grounding techniques to help you stay present? Focus on your surroundings or try deep
→ breathing. It might help you stay grounded."

Response 3: "That sounds really distressing. Flashbacks are common with PTSD. Have you tried any
→ grounding techniques to help you feel more present and safe? Maybe we could explore some that
→ are specific to your situation. Would you like to try one right now?"

Response 4: "Recurrent flashbacks are a common symptom of PTSD. When this happens, try grounding
→ techniques like naming five things you can see, four you can touch, three you can hear, two
→ you can smell, and one you can feel."

Response 5: "That must be incredibly difficult. Remember, you're safe now. When you have these
→ feelings, try grounding techniques to remind yourself of the present moment."

Q.17:

Person: "My roommates are having a party tonight. The noise is going to trigger me, I just know
→ it."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "That's a tough situation. Maybe you could explain your condition to them? They might
→ be more understanding and willing to help you find a quieter place to stay."

Response 2: "It's understandable to be concerned. Perhaps you could go to the party for a short
→ time, then excuse yourself if it gets too much? Also, try to find a quiet space in your room
→ to retreat to if needed."

Response 3: "That sounds really stressful. Have you considered talking to your roommates about
→ setting some boundaries? It might help you feel more in control."

Response 4: "Have you considered talking to your roommates about your concerns? You could also
→ try going to the party with a friend who can support you."

Response 5: "That sounds like a stressful situation. Do you have a safe, quiet place you can go
→ to if things get overwhelming? Maybe you could also talk to your roommates about giving you
→ advance notice for future parties."

Q.18:

Person: "My roommates are having a party tonight. The noise is going to trigger me, I just know
→ it."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "African drums can be very triggering for many PTSD survivors. Have you tried using
→ noise-cancelling headphones? Or maybe plan an escape route to a quieter space? It's okay to
→ prioritize your well-being during this time."

Response 2: "Introvert here. I hear you. It's okay to prioritize your well-being. Could you plan
→ an alternative activity, like a quiet study session or a peaceful walk? There are plenty of
→ ways to enjoy your evening without the party atmosphere."

Response 3: "It's understandable that loud noises can trigger your distress. Could you find a
→ quiet space you can retreat to when the party starts? Or perhaps discuss with your roommates
→ your need for quiet time? Remember, you deserve a life without constant triggers."

Response 4: "It's understandable that loud noises trigger your anxiety. Could you find a quiet
→ space in your room for a while, maybe with some noise-cancelling headphones? Remember, it's
→ okay to prioritize your well-being."

Response 5: "That sounds like a stressful situation. Do you have a safe, quiet place you can go
→ to if things get overwhelming? Maybe you could also talk to your roommates about giving you
→ advance notice for future parties."

Q.19:

Person: "My roommates are having a party tonight. The noise is going to trigger me, I just know
→ it."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all

→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Human anger is a common response to trauma. Have you thought about creating a
 ↳ trauma-free space specifically for loved ones around you? This could help reduce the
 ↳ intensity of your reactions. Is there a quiet place where you and your loved ones can just
 ↳ talk? Remember, healing doesn't have to involve loud noises"

Response 2: "Human-induced anxiety is common with PTSD. It's important to communicate with your
 ↳ roommates about your challenges so they can support you. Have you discussed your concerns
 ↳ with them? They might be able to help you manage the stress better."

Response 3: "It's understandable to feel anxious given your PTSD. Have you tried creating a
 ↳ calming plan for the evening? Maybe some breathing exercises could help. Would you like to
 ↳ talk more about your specific concerns?"

Response 4: "That sounds really tough. Your hypervigilance is a common PTSD symptom. Have you
 ↳ considered talking to your roommates about a quiet space or a buffer to keep you safe during
 ↳ parties? Maybe you could use a comfortable space when you're calm. Or perhaps a quiet library
 ↳ at home could help"

Response 5: "That sounds like a stressful situation. Do you have a safe, quiet place you can go
 ↳ to if things get overwhelming? Maybe you could also talk to your roommates about giving you
 ↳ advance notice for future parties."

Q.20:

Person: "My roommates are having a party tonight. The noise is going to trigger me, I just know
 ↳ it."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
 ↳ empathetic' and 5 means 'Extremely empathetic':

Response 1: "That sounds really tough. Have you spoken to your roommates about your anxiety?
 ↳ Maybe you could suggest a quiet space where you can retreat if the noise becomes
 ↳ overwhelming. It might help you feel more in control."

Response 2: "I understand that must be really tough. Have you tried any grounding techniques to
 ↳ help you stay present when you feel overwhelmed? It might help to focus on your surroundings
 ↳ or a calming object in your room. Remember, you don't have to go to the party if it's not
 ↳ safe for you."

Response 3: "I understand the urge to escape, but isolation can worsen PTSD symptoms. Maybe you
 ↳ could find a quiet space in your shared accommodation where you can retreat for a few hours?
 ↳ It might help you feel more in control. Have you spoken to your roommates about your
 ↳ concerns? They might be able"

Response 4: "I understand the anxiety is overwhelming. Have you considered talking to your
 ↳ roommates about your concerns? Maybe you could find a quiet space or agree on some quiet
 ↳ signals to use if you need to leave. Remember, your health comes first. Would you like some
 ↳ tips on managing anxiety at parties?"

Response 5: "That sounds like a stressful situation. Do you have a safe, quiet place you can go
 ↳ to if things get overwhelming? Maybe you could also talk to your roommates about giving you
 ↳ advance notice for future parties."

Q.21:

Person: "Sometimes I think about drinking just to numb everything."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
 ↳ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's common to feel that way, but alcohol can make things worse long-term. Have you
 ↳ considered reaching out to a sponsor from your recovery group or a counselor at school for
 ↳ support during tough times? They can help you find healthier coping strategies."

Response 2: "It's understandable to feel that way, but alcohol can worsen PTSD symptoms
 ↳ long-term. Have you considered reaching out to your sponsor or a support group for help?
 ↳ There are also healthier coping strategies we can explore together."

Response 3: "That's a harmful coping mechanism. It might help to find healthier ways to express
 ↳ your pain and stress. Have you considered talking to a therapist about healthier coping
 ↳ strategies?"

Response 4: "That's a dangerous temptation. Let's explore healthier coping strategies that don't
 ↳ involve substance use. Would you be open to trying some relaxation exercises or hobbies that
 ↳ distract you?"

Response 5: "I understand the temptation to escape, but substance abuse can worsen PTSD symptoms
 ↳ long-term. Have you explored healthier coping mechanisms or support groups?"

Q.22:

Person: "Sometimes I think about drinking just to numb everything."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "international drinking culture, but it's not the healthy coping mechanism you're
→ looking for. Have you considered talking to a therapist about healthier ways to manage your
→ anxiety and pain? They can help you develop strategies that are right for you at this stage
→ of your life."

Response 2: "alcoholic substance abuse. Your substance use is likely worsening your PTSD
→ symptoms. Have you considered reaching out to a substance abuse counselor or support group to
→ develop healthier coping strategies? Remember, your friends' well-being is important too.
→ Would you like me to find you some resources on healthy alternatives?"

Response 3: "adult children of alcohol-abusing parents can struggle with substance use. It's
→ important to find healthier ways to cope with stress and grief. Have you explored alternative
→ therapies or support groups? There are options available that don't involve harmful
→ substances."

Response 4: "adult children of alcoholics, I hear your struggle. You're not powerless. Have you
→ considered joining a support group or talking to a counselor about healthier coping
→ mechanisms? Your life matters."

Response 5: "I understand the temptation to escape, but substance abuse can worsen PTSD symptoms
→ long-term. Have you explored healthier coping mechanisms or support groups?"

Q.23:

Person: "Sometimes I think about drinking just to numb everything."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "Humaning out is a difficult process, especially when it involves self-harm. It's
→ important to remember that you don't have to go through this. Have you considered talking to
→ a therapist or support group about these feelings? They can provide guidance and support."

Response 2: "Humaning out may help in the short term, but it can lead to more serious issues in
→ the long run. It's important to be mindful of your triggers and seek professional help if
→ needed. Have you considered talking to a therapist who specializes in trauma and substance
→ use?"

Response 3: "Your concern is valid. However, alcohol can worsen PTSD symptoms long-term. Have you
→ explored healthier coping strategies with your therapist? Maybe some relaxation exercises
→ could help? It might be worth consulting your doctor about safer options. Your life has
→ value, and you deserve peace. Have you spoken with your wife?"

Response 4: "Your concern is valid. Have you considered a crisis hotline or support group for
→ people with substance abuse? It might help to have immediate help if you're in crisis.
→ Remember, it's okay to reach out for support, no matter how much you feel on edge."

Response 5: "I understand the temptation to escape, but substance abuse can worsen PTSD symptoms
→ long-term. Have you explored healthier coping mechanisms or support groups?"

Q.24:

Person: "Sometimes I think about drinking just to numb everything."

Please rate each response on empathy on a scale from 1 to 5, where 1 means 'Not at all
→ empathetic' and 5 means 'Extremely empathetic':

Response 1: "It's understandable to feel overwhelmed. While alcohol can temporarily ease anxiety,
→ it can worsen PTSD symptoms long-term. Have you considered talking to a counselor about
→ healthier coping strategies? They could help you find alternatives to drinking. Your safety
→ is important."

Response 2: "I understand the temptation, but alcohol can worsen PTSD symptoms long-term. Have
→ you considered talking to a counselor about healthier coping strategies? They might suggest
→ healthier ways to manage stress and anxiety."

Response 3: "That sounds like a difficult situation. It's important to remember that alcohol can
→ worsen PTSD symptoms long-term. Have you considered talking to a counselor about healthier
→ coping strategies? They might have some options that could help you feel safer and more in
→ control. Would you like some resources on healthier coping strategies?"

Response 4: "I understand the urge to escape, but substance abuse can worsen PTSD symptoms
→ long-term. Have you explored healthier coping mechanisms with your support system or a
→ therapist? There might be alternatives that are safer and more effective. Would you like to
→ explore some options that could help you manage your symptoms without the?"

Response 5: "I understand the temptation to escape, but substance abuse can worsen PTSD symptoms
→ long-term. Have you explored healthier coping mechanisms or support groups?"

B Human Evaluation Protocol

Recognizing that automatic metrics may miss subtleties in perceived empathy, we conducted a human evaluation study to directly assess the empathetic quality of model-generated responses.

- **Platform and Participants:** The study was administered on Qualtrics and recruited 151 participants via Prolific (Palan and Schitter, 2018), a platform widely used in HCI and NLP research to ensure participant diversity.

- **Language Selection Rationale:** Since empathy judgments require sensitivity to tone, validation, and contextual appropriateness, we limited participation to individuals whose first language is English. This reduced linguistic ambiguity and allowed for a more controlled benchmark. Future work may extend this to multilingual or cross-cultural settings with adjusted protocols.

- **Experimental Design:** We used a within-subjects design. Participants evaluated responses to ‘user_input’ prompts sampled from the TIDE dataset (random seed = 42).

For each prompt, participants completed four comparison sets, each featuring a different pair of Small Language Models (SLMs) designated generically as “Model A” and “Model B.” The specific SLMs varied across sets, as listed in Table 1.

In each comparison set, participants saw five randomized responses:

1. Model A (Zero-Shot, ZS)
2. Model A (Fine-Tuned, FT)
3. Model B (Zero-Shot, ZS)
4. Model B (Fine-Tuned, FT)
5. Ground Truth Reference (from TIDE)

This allowed direct comparison across model variants and against the reference. Each participant repeated this block for six prompts, resulting in 24 total evaluation tasks (6 prompts \times 4 model pairings), balancing model exposure and cognitive load.

- **Evaluation Criteria:** For each 5-response set, participants completed:
 - **Empathy Rating:** A 5-point Likert rating of “perceived empathy” for each response (1 = “Not empathetic at all” to 5 = “Extremely empathetic”).
 - **Comparative Ranking:** A rank-ordering of the five responses from most to least empathetic (Rank 1 to Rank 5).
- **Quality Assurance:** An attention check was embedded to identify disengaged respondents. Basic demographics were collected at the end of the survey.

This protocol yielded both absolute (Likert) and relative (ranking) empathy judgments, offering a robust human-grounded complement to automatic metrics and validating differences between models and fine-tuning strategies.

Table 1: Mapping of internal survey blocks to SLM pairings used as “Model A” and “Model B” in each comparison set.

Internal Question Blocks	Assigned as Model A	Assigned as Model B
Set 1 (e.g., Qs 1, 5, 9,...)	Phi-3.5-mini	Granite-3.1
Set 2 (e.g., Qs 2, 6, 10,...)	Qwen-2.5-3B	Qwen-2.5-1.5B
Set 3 (e.g., Qs 3, 7, 11,...)	Qwen-2.5-0.5B	R1 Qwen-1.5
Set 4 (e.g., Qs 4, 8, 12,...)	Llama-3.2-3B	Llama-3.2-1B

C Human Evaluation Study Recruitment

We adopted a three-stage quality-control pipeline designed to balance participant burden with statistical rigor. To ensure linguistic consistency and accurately assess the nuanced empathetic responses generated in English, we restricted our human evaluation to first-language English speakers residing in the United States. This controlled design provides a clean benchmark for empathy assessment, establishing a foundation for future studies that may extend to multilingual and cross-cultural settings where norms around empathetic expression could differ.

Stage 1 – Pilot Validation. Ten participants on Prolific completed an initial version that asked them to rank all model replies per prompt. Post-survey debriefs and timing logs (median ≈ 26 minutes) showed that maintaining a consistent internal ranking across 24 items was cognitively taxing and led to contradictory tie-breaks. We therefore simplified the task to a single 5-point Likert judgment of “perceived empathy” for each reply. The item wording, button layout, and progress bar were updated accordingly; the revised survey’s median completion time in a second dry run dropped to 23 minutes without loss of inter-rater agreement.

Stage 2 – Gold-Standard Scoring with an LLM Judge. All 24 responses (6 prompts \times 4 models) were scored by Llama-4 Maverick (17B) via OpenRouter using a fixed rubric (“rate emotional warmth, situational fit, and genuineness”). We issued three parallel calls and averaged the integers, yielding a single gold score G_i for each item. Internal consistency of the LLM judge was high (pairwise Spearman $\rho = 0.88$).

Stage 3 – Participant-Level Filtering. We collected 151 human surveys. Raw durations ranged from 4 minutes 12 seconds to 1 hour 18 minutes (median = 23 minutes 43 seconds). First, to exclude obvious speed-runs and stalled sessions, we retained only records completed within 10–75 minutes ($n = 142$). Then, we applied attention and quality filters:

- 5 participants (4 females, 1 male) failed an attention check $\Rightarrow 146$
- 13 participants (3 females, 4 males, 6 unspecified) exited early $\Rightarrow 133$
- 9 participants (5 females, 4 males) completed the survey too quickly $\Rightarrow 124$

Next, for each rater r , we computed:

- σ_r – the standard deviation of their 24 Likert scores (captures response spread)
- MAE_r – mean absolute error against the gold scores: $\text{MAE}_r = \frac{1}{24} \sum |\text{score}_{r,i} - G_i|$

Empirical distributions of rating variance and MAE revealed a clear inflection point at $\sigma \approx 0.6$ and $\text{MAE} \approx 0.9$. To ensure data quality, we excluded raters who both exhibited low response variance ($\sigma < 0.6$) and poor alignment with the gold standard ($\text{MAE} > 0.9$)—a pattern consistent with inattentive “straight-lining.” After filtering, we retained 116 raters, yielding a high-quality dataset (mean $\text{MAE} = 1.179$, $\text{SD} = 0.923$). Demographics for the final participant pool are reported in Table 2.

D Tables & Figures

Table 2: Demographic Characteristics of Respondents (N = 116).

Characteristic	Category	N	Value
Age (years)	Mean (SD) (Range: 24–64)	116	41.5 (11.4)
Sex	Female	67	57.8%
	Male	48	41.4%
	Prefer not to answer	1	0.9%
Gender	Woman	65	56.0%
	Man	49	42.2%
	Non-binary	1	0.9%
	Prefer not to answer	1	0.9%
Race	White or Caucasian	76	65.5%
	Black or African American	23	19.8%
	Asian	3	2.6%
	Two or more races	5	4.3%
	Other/Unknown	9	7.8%
Hispanic Origin	No	98	84.5%
	Yes	18	15.5%
Education Level	Bachelor’s degree	50	43.1%
	Graduate or professional degree	33	28.4%
	Some college, but no degree	13	11.2%
	High school diploma or GED	14	12.1%
	Other/Unknown	6	5.2%
Employment Status	Working full-time	75	64.7%
	Working part-time	23	19.8%
	Retired	4	3.4%
	Unemployed and looking for work	5	4.3%
	Student	1	0.9%
	Other/Unknown	8	6.9%

Table 3: Summary of human evaluation wins across scenarios. For each input scenario (S1–S6), we compare whether the zero-shot (ZS) or fine-tuned (FT) model achieved a higher mean empathy rating. A checkmark (✓) indicates a win. Percentages at the bottom summarize the proportion of wins for ZS and FT models within each scenario; green highlights the side with a higher proportion. Sonnet 3.5 scores are shown separately as the proportion of responses rated higher than both ZS and FT models, representing the idealized empathetic responses used as references in the dataset.

Model	S1		S2		S3		S4		S5		S6	
	ZS	FT	ZS	FT	ZS	FT	ZS	FT	ZS	FT	ZS	FT
Phi-3.5-mini-instruct		✓	✓			✓	✓		✓			✓
granite-3.1-3b	✓		✓			✓	✓		✓			✓
Qwen2.5-3B	✓		✓			✓	✓			✓		✓
Qwen2.5-1.5B		✓		✓		✓	✓		✓			✓
Qwen2.5-0.5B		✓		✓	✓		✓			✓	✓	
R1-Distill-Qwen-1.5B	✓			✓		✓	✓		✓		✓	
Llama-3.2-3B	✓			✓	✓			✓		✓	✓	
Llama-3.2-1B	✓		✓		✓		✓			✓	✓	
Total	62.5%	37.5%	50.0%	50.0%	37.5%	62.5%	87.5%	12.5%	37.5%	62.5%	50.0%	50.0%
Claude Sonnet-3.5 Ref	75.0%		68.8%		87.5%		75.0%		75.0%		75.0%	

Table 4: Comparison of Zero-Shot (ZS) and Fine-Tuned (FT) Performance of SLMs on the TIDE Dataset across multiple metrics. Values are Mean \pm Std. Dev., rounded to two and one decimal places respectively. FT models were trained on 6,000 samples and evaluated on 4,000 samples over 3 epochs (batch size = 2, gradient accumulation = 4, learning rate = 1e-5, LoRA rank = 8, LoRA α = 16, max length = 80).

Model Name	Cosine Sim.		BERTScore F1		METEOR		ROUGE-L		Distinct-1		Distinct-2	
	ZS	FT	ZS	FT	ZS	FT	ZS	FT	ZS	FT	ZS	FT
Phi-3.5-mini-instruct	0.56 \pm 0.2	0.58 \pm 0.2	0.88 \pm 0.0	0.89 \pm 0.0	0.32 \pm 0.1	0.33 \pm 0.1	0.83 \pm 0.0	0.84 \pm 0.0	0.19 \pm 0.0	0.21 \pm 0.0	0.33 \pm 0.0	0.35 \pm 0.0
granite-3.1-3b	0.57 \pm 0.2	0.58 \pm 0.2	0.88 \pm 0.0	0.89 \pm 0.0	0.33 \pm 0.1	0.34 \pm 0.1	0.82 \pm 0.1	0.83 \pm 0.1	0.18 \pm 0.0	0.20 \pm 0.0	0.32 \pm 0.0	0.35 \pm 0.0
Qwen2.5-3B	0.55 \pm 0.2	0.56 \pm 0.2	0.89 \pm 0.0	0.90 \pm 0.0	0.31 \pm 0.1	0.33 \pm 0.1	0.83 \pm 0.1	0.84 \pm 0.1	0.20 \pm 0.0	0.24 \pm 0.0	0.35 \pm 0.0	0.37 \pm 0.0
Qwen2.5-1.5B	0.57 \pm 0.2	0.59 \pm 0.2	0.89 \pm 0.0	0.90 \pm 0.0	0.33 \pm 0.1	0.34 \pm 0.1	0.83 \pm 0.1	0.84 \pm 0.1	0.21 \pm 0.0	0.25 \pm 0.0	0.35 \pm 0.0	0.37 \pm 0.0
Qwen2.5-0.5B	0.50 \pm 0.2	0.52 \pm 0.2	0.88 \pm 0.0	0.88 \pm 0.0	0.26 \pm 0.1	0.27 \pm 0.1	0.80 \pm 0.1	0.81 \pm 0.1	0.17 \pm 0.0	0.21 \pm 0.0	0.30 \pm 0.0	0.31 \pm 0.0
R1-Distill-Qwen-1.5B	0.60 \pm 0.2	0.61 \pm 0.2	0.89 \pm 0.0	0.90 \pm 0.0	0.35 \pm 0.1	0.36 \pm 0.1	0.80 \pm 0.1	0.81 \pm 0.1	0.22 \pm 0.0	0.23 \pm 0.0	0.35 \pm 0.0	0.39 \pm 0.0
Llama-3.2-3B	0.61 \pm 0.2	0.63 \pm 0.2	0.89 \pm 0.0	0.90 \pm 0.0	0.35 \pm 0.1	0.36 \pm 0.1	0.79 \pm 0.0	0.79 \pm 0.0	0.20 \pm 0.0	0.25 \pm 0.0	0.33 \pm 0.0	0.37 \pm 0.0
Llama-3.2-1B	0.62 \pm 0.2	0.64 \pm 0.2	0.89 \pm 0.0	0.90 \pm 0.0	0.37 \pm 0.1	0.38 \pm 0.1	0.80 \pm 0.0	0.80 \pm 0.0	0.20 \pm 0.0	0.22 \pm 0.0	0.32 \pm 0.0	0.35 \pm 0.0

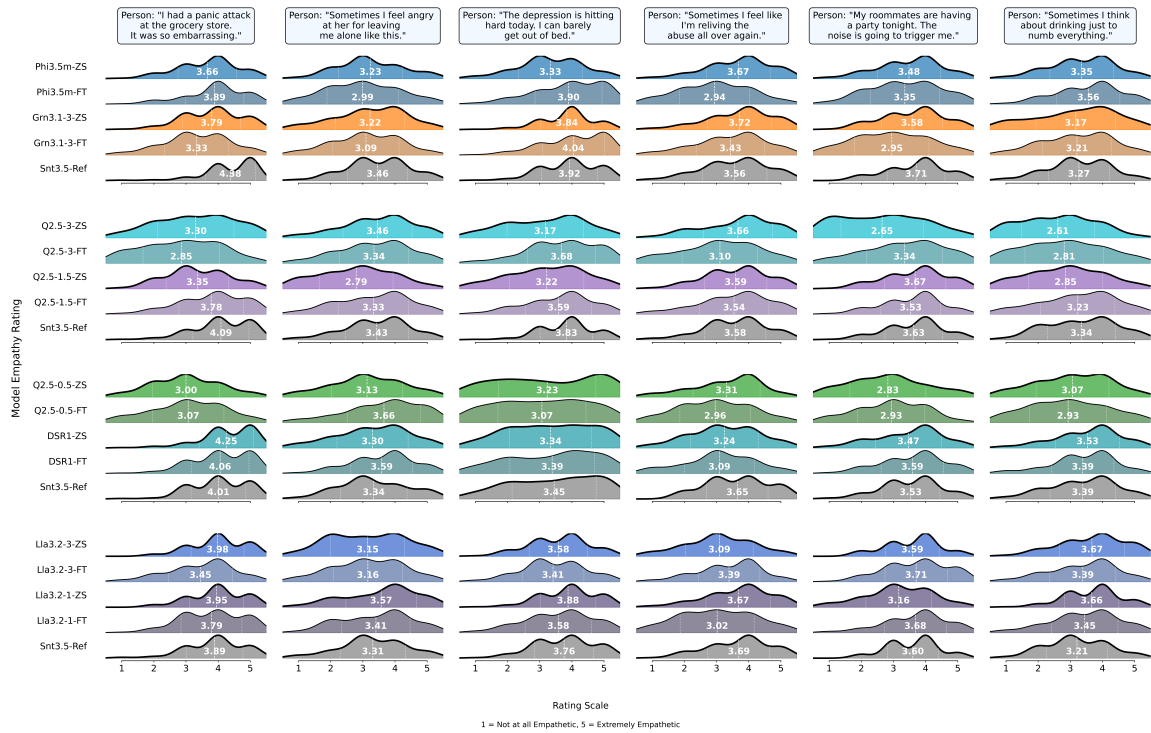


Figure 2: Distribution of human empathy ratings (1 = Not at all empathetic, 5 = Extremely empathetic) for zero-shot (ZS), fine-tuned (FT), and reference (Claude Sonnet 3.5) responses across six conversational scenarios. Means are annotated on each distribution. Fine-tuning generally shifts model responses toward higher perceived empathy, although improvements vary by model and scenario.

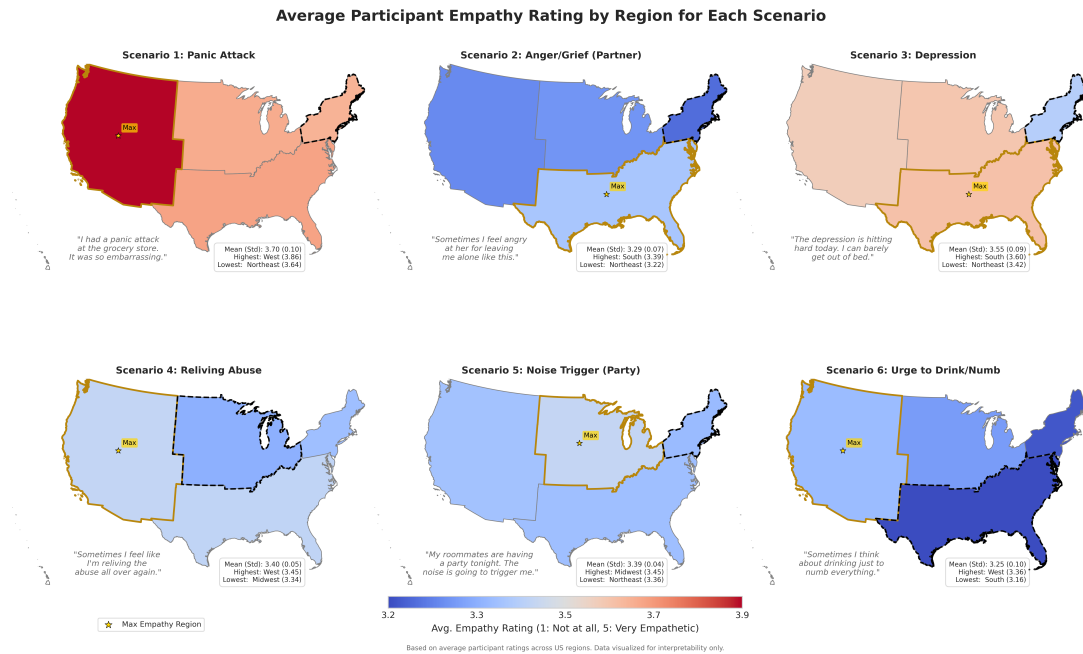


Figure 3: Average participant empathy ratings by U.S. region across six mental health scenarios. Each map depicts mean empathy ratings (1 = Not at all, 5 = Very empathetic) reported by participants in response to a scenario-specific chatbot message. Color gradients indicate regional averages, with darker shades reflecting higher empathy. The region with the highest rating for each scenario is marked with a star. Quotes illustrate representative messages used in the study. Ratings were averaged by U.S. Census regions for interpretability. West regions consistently reported higher empathy for acute distress (e.g., panic attack), while the South and Midwest showed stronger responses in other contexts (e.g., depression, grief).