

From Heart to Words: Generating Empathetic Responses via Integrated Figurative Language and Semantic Context Signals

Anonymous ACL submission

Abstract

Although generically expressing empathy is straightforward, effectively conveying empathy in specialized settings presents nuanced challenges. We present a conceptually motivated investigation into the use of figurative language and causal semantic context to facilitate targeted empathetic response generation within a specific mental health support domain, studying how these factors may be leveraged to promote improved response quality. Our approach achieves a 7.6% improvement in BLEU, a 36.7% reduction in Perplexity, and a 7.6% increase in lexical diversity (D-1 and D-2) compared to models without these signals, and human assessments show a 24.2% increase in empathy ratings. These findings provide deeper insights into grounded empathy understanding and response generation, offering a foundation for future research in this area.

1 Introduction

Whether through a gentle idiom like “I’ve got your back” or a careful choice of phrasing, people often reach beyond straightforward language to convey empathetic support (Barak et al., 2008; Naslund et al., 2016; Sharma et al., 2020a). In particular, employing figurative language (e.g., metaphors or idioms) is a prominent tool used to this effect (Lee et al., 2024b), and is commonly understood to enhance emotional expression by making the abstract more vivid and relatable (Fussell and Moss, 2014). Consider the two statements: “*I understand it’s tough*” (literal) and “*I understand it feels like fighting an endless battle*” (figurative). The latter conveys stronger emotional resonance, fostering deeper connection with those seeking support.

Despite the intuitive advantages of generating more rhetorically and contextually targeted empathy, existing empathy generation research predominantly focuses on understanding the emotions of the speaker (the individual sharing their struggles)

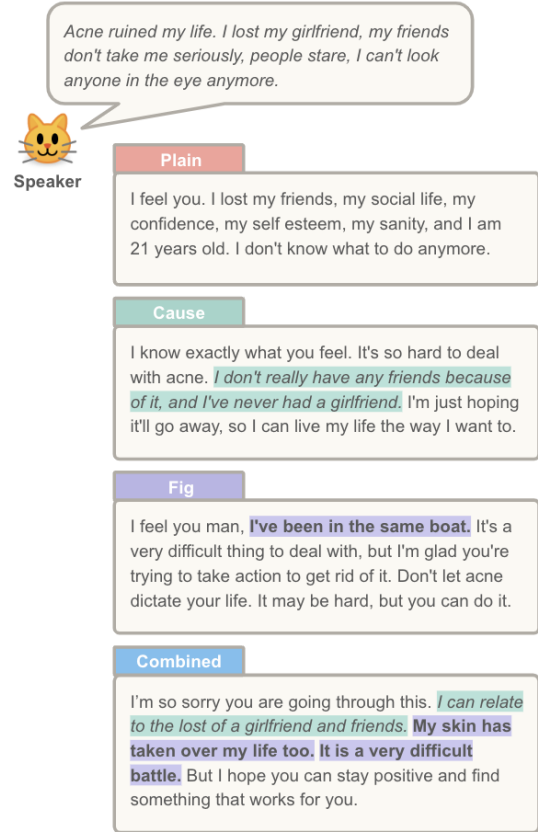


Figure 1: Illustration of empathetic responses generated using different LLM fine-tuning approaches: plain, figurative language, cause annotations, and combined.

(Rashkin et al., 2019; Welivita and Pu, 2023), rather than on *how* the responder (the individual providing support) conveys empathy. While some studies, such as work by Welivita et al. (2023), explore structured response intents (e.g., agreeing or suggesting), these approaches align responses with emotions rather than leveraging deeper rhetorical tools. This highlights a gap in exploring how these more complex language phenomena can be used to effectively convey empathy, particularly when informed by the speaker’s context.

We address this gap by investigating how both rhetorical craft (focusing on figurative language) and semantic context (focusing on empathy cause) contribute to empathetic response generation in a specialized mental health setting. We summarize our contributions as:

- **We enrich Lee and Parde (2024)’s AcnEmpathize corpus** with figurative language metadata and manual empathy cause annotations, supplying a useful additional layer of data.
- **We propose prompting methods to integrate figurative language metadata and empathy cause data** into the empathetic response generation pipeline.
- **We show that the joint integration of these components significantly improves empathy response generation performance** across a broad range of automated and human evaluation metrics in a specialized setting.

Through this work, we aim to inspire further exploration of targeted, nuanced signals in empathetic communication. While our findings are drawn from a specialized domain, we hope that they pave the way for broader study in diverse support settings.

2 Related Work

The importance of effectively communicating empathy has been demonstrated across multiple domains (Decety and Jackson, 2004; Green et al., 2005; Riess and Kraft-Todd, 2014). Computational approaches have specifically sought to automate empathetic response generation by modeling speaker emotions and generating appropriate replies, often using the popular EmpatheticDialogues dataset (Rashkin et al., 2019) which contains emotion-labeled conversations (Lin et al., 2020; Majumder et al., 2020).

Beyond emotion recognition, researchers have studied how the causes behind a speaker’s emotions can better contextualize responses. For example, Gao et al. (2021) and Li et al. (2021) analyzed emotion causes in the EmpatheticDialogues dataset to identify triggers behind the speaker’s emotions. Similarly, Qian et al. (2023) integrated emotion cause recognition into large language models (LLMs) for empathetic response generation. While these works focus on understanding the triggers (e.g., “I failed an exam”) behind the speaker’s own

emotions (e.g., sadness), our work introduces *empathy causes*, which identify specific parts of the speaker’s text that evoke empathy *from a responder* in interpersonal communication (see Figure 2 for an example and Section 3.2 for further details).

Recent works have also started considering communication strategies to guide empathetic response generation. For instance, Welivita et al. (2021) introduced a dataset combining emotion labels and response intents (such as agreeing or suggesting), and in follow-up work they demonstrated how these intents could guide the generation of emotionally supportive and empathetic responses (Welivita et al., 2023). Similarly, Saha et al. (2022) incorporated reinforcement learning for empathetic rewriting.

While these communication strategies can help shape empathetic responses, they do not explore more complex rhetoric. Recent work on empathetic storytelling (Shen et al., 2024) shows how narrative style elements—such as tone and phrasing in the speaker’s text—can influence perceived empathy, demonstrating the potential for leveraging rhetorical device for empathetic response generation. Figurative language is a powerful such tool that enriches emotional expression (Fussell and Moss, 2014; Citron and Goldberg, 2014). Computational studies have shown that incorporating figurative language—specifically metaphors, idioms, and hyperbole—can improve predictions of both emotion (Lee et al., 2024a) and empathy (Lee et al., 2024b). Despite the clear value of figurative language to empathetic expression, as highlighted by these works, figurative language remains unexplored in empathetic response generation.

Building on the success of emotion cause annotations in improving response generation (Gao et al., 2021), we integrate *figurative language* and *empathy cause annotations* to address both the rhetorical and contextual aspects of empathetic response generation. We aim to create responses that are not only emotionally engaging but also well-aligned with the speaker’s concerns.

3 Data

3.1 Source Dataset

We use the publicly available AcnEmpathize (Lee and Parde, 2024), which captures authentic emotional exchanges from an online acne support community. Unlike general-domain empathy datasets, such as *EmpatheticDialogues* (Rashkin et al., 2019) (multi-topic emotional dialogues) and *EPITOME*

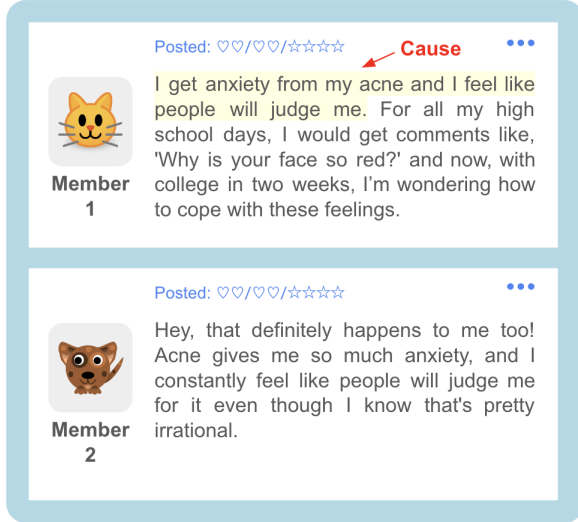


Figure 2: Example of an initial (speaker) post and an empathetic reply in the AcnEmpathize dataset. The highlighted portion in the speaker’s post indicates the annotated cause that evokes empathy in the given reply.

(Sharma et al., 2020b) (conversations spanning multiple mental health topics), AcnEmpathize represents a focused domain-specific peer-support community. It features over 12K posts categorized into *initial posts*, written by individuals seeking support, and *responses* from others responding to emotional challenges. We sampled a subset of 2,492 posts in speaker-response pair format, where all responses contain empathy. These pairs correspond to 1,110 unique speaker posts, as many posts received multiple empathetic replies.

3.2 Cause Annotations

We annotated empathy causes in our sampled speaker posts (see Figure 2 for an example of an annotated cause in a speaker-response pair), and we release these annotations publicly. By creating an explicit signal for which parts of a speaker’s text require empathetic acknowledgment, we anticipate that these annotations can help models directly address relevant points of distress, rather than engaging with the post more generically. We annotated cause sentences, using the collaborative tool INCEpTION (Klie et al., 2018). We recruited three graduate student volunteers with formal training in natural language processing at a U.S.-based institution. Annotators were instructed to highlight cause sentence(s) in each speaker post that were most likely to prompt the corresponding empathetic reply across three rounds (see Appendix A.1 for annotation details).

| Language Type | # Posts (%) |
|------------------|----------------|
| Idiom | 1,225 (49.16%) |
| Metaphor | 887 (35.59%) |
| Hyperbole | 559 (22.43%) |
| Total Figurative | 1,723 (69.14%) |

Table 1: Distribution of figurative language type (idioms, metaphors, and hyperbole) in responses within the cause-annotated AcnEmpathize dataset. The counts represent the number of response posts that contain each type of figurative language. Each response may contain more than one type.

In **Round 1**, annotators independently labeled 10 identical conversations and participated in a discussion afterward to resolve disagreements, resulting in eventual perfect inter-annotator agreement (IAA) using Krippendorff’s alpha (Krippendorff, 1970). In **Round 2**, annotators labeled 90 additional identical conversations, resulting in an initial IAA of 0.70. Disagreements in this round were also resolved through discussion to reach full consensus. Pre-consensus pairwise IAA scores for the 100 triple-annotated samples ranged from 0.67 to 0.73, consistent with the IAA score of 0.68 reported in an external empathy study (Sharma et al., 2020b). For **Round 3**, the remaining conversations were divided among the annotators in a ratio of 476:476:150.¹

3.3 Figurative Language Metadata

We automatically identified the presence of metaphors, idioms, and hyperbole in AcnEmpathize using an externally validated multitask framework proposed by Lai et al. (2023) and built on top of mT5 (Xue, 2020). This method identifies figurative language through template-based prompt learning. We used the detection prompt:

Which figure of speech does this text contain? (A) Literal (B) [Task] | Text: [Text]

[Task] corresponds to one of the figurative language types: idiom, metaphor, or hyperbole. Each sentence in the response text was iteratively assessed for each figurative language type, and the results were recorded as binary indicators.

As shown in Table 1, approximately 69% of empathetic responses in our data contain figurative language, with idioms being the most common

¹One annotator had an unavoidable and unexpected schedule constraint that prevented equal distribution.

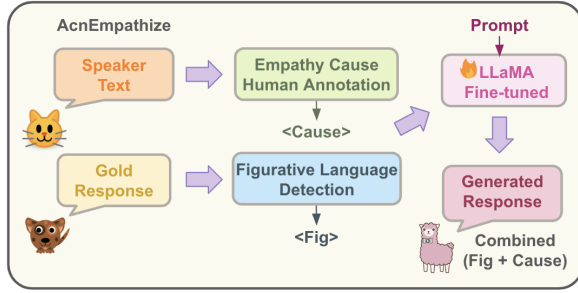


Figure 3: Overview of our pipeline for empathetic response generation. Speaker texts are manually annotated with empathy causes, while figurative language in responses is identified using a detection method. These elements are integrated during the fine-tuning of LLaMA, guided by prompts, to generate empathetic responses that are both contextually relevant and linguistically expressive.

(49.16%), followed by metaphors (35.59%) and hyperbole (22.43%). This highlights the frequent use of figurative language in empathetic responses, further justifying its incorporation alongside cause annotations to enrich our generation process.

4 Empathetic Response Generation

Using the new cause annotations and figurative language metadata, we sought to generate empathetic responses that meaningfully address the struggles expressed in speaker posts through both emotional engagement and contextual alignment. We study these factors both independently and in concert with one another. An overview of our full pipeline is presented in Figure 3. Although gold empathy cause labels are used in this study to demonstrate proof of concept, a promising future direction (and one facilitated by the new cause labels) involves the automated detection of empathy cause.

4.1 Modeling Framework

We systematically finetune the LLaMA-3-8B model (Touvron et al., 2023) using our gold cause annotations and figurative language metadata, and also compare to a zero-shot baseline. The fine-tuned models (*plain*, *cause*, *figurative*, and *combined*) are trained to learn response patterns and expressions present in the dataset. Each approach employs different prompts tailored to the specific objectives of the fine-tuning strategy.

4.2 Plain

In plain fine-tuning, the model is trained without cause or figurative language annotations. Thus,

the model learns to generate empathetic responses based on the natural patterns and style of replies present in the dataset. During generation, the following prompt is used:

Given the input text, generate an empathetic response.

We also used this prompt for the zero-shot baseline, which by definition did not involve any fine-tuning and instead used the base LLaMA-3-8B model for inference directly.

4.3 Cause

In this approach, we fine-tuned the model using cause annotations in the speaker posts, as described in Section 3.2. We specify gold-standard causes by wrapping them in `<cause>` tags. The prompt correspondingly acknowledges these tags:

Given the input text with `<cause>` tags, generate a targeted empathetic response that acknowledges the specific concerns expressed.

See Appendix A.1 for a detailed example of how the training data was formatted for fine-tuning.

4.4 Figurative Language (Fig)

In this approach, we fine-tuned the model using figurative language identified in the response texts, based on the detection method outlined in Section 3.3. We classify sentences containing idioms, metaphors, or hyperbole, and mark the classified sentences with `<idiom>`, `<metaphor>`, and `<hyperbole>` tags. These tags expose the model to examples of how figurative expressions are employed in empathetic replies. For generation, we design the prompt to flexibly include figurative language where appropriate, without explicitly specifying particular expressions:

Given the input text, generate an empathetic response that uses figurative language where appropriate, specifically `<idiom>`, `<metaphor>`, or `<hyperbole>`.

4.5 Combined (Fig + Cause)

Our final approach integrates both figurative language and cause annotations to enrich responses both rhetorically and contextually. Speaker texts

are tagged with <cause> tags around labeled empathy causes, while response texts are tagged with <idiom>, <metaphor>, and <hyperbole> during training. During generation, we use the prompt:

Given the input text, generate an empathetic response that very strongly emphasizes the use of figurative language, specifically <idiom>, <metaphor>, and <hyperbole>, optimizing <idiom> and <metaphor> to maximize emotional support, while addressing the concerns indicated by the <cause> tags.

The emphasis on idioms and metaphors in the prompt is motivated by a prior study on empathy detection using AcnEmpathize (Lee et al., 2024b), which shows their statistically significant association with empathy labels. While hyperbole remains a useful linguistic device, idioms and metaphors are prioritized to maximize the emotional supportiveness of the responses, with cause annotations incorporated to maintain contextual relevance.

5 Evaluation

5.1 Experimental Setup

All experiments utilized a 4-bit quantized version of the LLaMA-3-8B model, implemented with the FastLanguageModel framework² to optimize memory usage and computational efficiency. For the fine-tuning approaches, the model was trained on both speaker posts and responses to generate empathetic outputs. The model was fine-tuned using a batch size of 1, a learning rate of 5e-5, and the AdamW optimizer in 8-bit mode. Training was conducted for three epochs using three NVIDIA 2080 TI GPUs with FP16 or BF16 support, utilizing the PEFT framework with a LoRA (Hu et al., 2021) configuration (rank = 16, alpha = 16, dropout = 0). All models and implementation details will be made available via GitHub after publication.

5.2 Evaluation Frameworks

We compare conditions using both automated metrics and human evaluation to provide a well-rounded assessment of the generated responses, encompassing linguistic quality, lexical alignment, and empathetic support. We describe our evaluation frameworks below.

²We implemented it using the unsloth GitHub repository <https://github.com/unslothai/unsloth>.

5.2.1 Automatic Evaluation

We assess different aspects of response quality using perplexity (PPL), BLEU, and Distinct-1 (D-1) and Distinct-2 (D-2) scores. PPL is measured using the Hugging Face transformers library (Wolf et al., 2020) to evaluate the likelihood of gold responses under the model’s probability distribution, indicating fluency and coherence. BLEU evaluates lexical overlap between generated and gold responses (Papineni et al., 2002). We compute it using sentence_bleu from the NLTK library (Bird et al., 2009), averaging across multiple gold responses, with a smoothing function applied to handle sparsity (Chen and Cherry, 2014). Finally, we calculate the ratio of unique unigrams (D-1) and bigrams (D-2) to the total number of tokens (Li et al., 2015) (also using the NLTK library) to measure the lexical diversity of generated responses.

5.2.2 Human Evaluation

Additionally, we conducted a human evaluation to provide a more holistic assessment of the generated responses. The evaluation was performed by three graduate student volunteers with formal training in NLP at a U.S.-based institution.³ The annotators rated the responses based on the following criteria, which are recognized as the most commonly used in empathetic conversational systems (Raamkumar and Yang, 2022), following Rashkin et al. (2019):⁴

- **Empathy:** The response’s ability to demonstrate understanding of the speaker’s feelings.
- **Relevance:** The extent to which the response is appropriate and on-topic.
- **Fluency:** The ease of understanding and linguistic clarity.

These aspects collectively capture the effectiveness of an empathetic response, aligning with prior work in evaluating empathetic dialogue systems. We randomly sampled 111 sets of four generated responses, each corresponding to a single speaker post.⁵ To minimize potential bias, the responses

³Two annotators were not involved in the cause annotation process, while one annotator participated in both tasks.

⁴Agreement scores are typically not reported for subjective dimensions in empathetic response generation, as seen in prior work (Rashkin et al., 2019; Majumder et al., 2020). Rather than enforcing uniform agreement, our evaluations were designed to capture diverse perspectives, aligning with real-world variability in how individuals perceive empathy, relevance, and fluency.

⁵These 111 sets represent 10% of the 1,110 unique speaker posts used for generation.

| Model | PPL (↓) | BLEU (↑) | D-1 (↑) | D-2 (↑) |
|-----------|--------------|--------------|--------------|--------------|
| Zero-shot | 14.944 | 0.058 | 0.587 | 0.847 |
| Plain | 14.182 | 0.764 | 0.515 | 0.751 |
| Cause | 13.990 | 0.772 | 0.523 | 0.755 |
| Fig | 9.100 | 0.775 | 0.561 | 0.814 |
| Combined | 8.980 | 0.822 | 0.569 | 0.814 |

Table 2: Performance of the zero-shot baseline and fine-tuned approaches on automated metrics for empathetic response generation. *Combined* shows the best overall performance. While *Zero-shot* achieves the highest D-1 and D-2 scores, it is excluded from further evaluation due to its extremely low BLEU score (0.058). (↑) means higher is better, (↓) means lower is better.

| Model | E | R | F | Most Supportive (%) |
|----------|--------------|--------------|--------------|---------------------|
| Plain | 3.565 | 3.631 | 3.207 | 3.19% |
| Cause | 3.889 | 4.024 | 3.799 | 11.70% |
| Fig | 4.195 | 4.021 | 4.132 | 32.98% |
| Combined | 4.426 | 4.135 | 4.189 | 52.13% |

Table 3: Performance on the fine-tuned approaches in human evaluation for empathetic response generation. The values represent the average scores for Empathy (E), Relevance (R), and Fluency (F) across all evaluated samples, along with the Most Supportive (%) column reflecting the percentage of responses selected as Most Supportive by the majority of annotators. *Combined* achieves the best overall performance across all metrics.

generated from each of the four approaches were shuffled before being presented to annotators. Annotators were also asked to select the response they consider **Most Supportive**, beyond the individual scores for Empathy, Relevance, and Fluency. These counts were weighted based on majority agreement (i.e., responses selected by at least two annotators) for evaluation. Complete evaluation guidelines, including definitions and examples for each criterion, are provided in Appendix A.2.

6 Results

6.1 Automatic Evaluation

Table 2 summarizes the results of evaluating the Zero-shot baseline and fine-tuned approaches us-

ing the automated metrics. *Combined* (*Fig + Cause*) achieves the best overall performance, balancing highest BLEU (0.822) and lowest Perplexity (8.980) with competitive D-1 (0.569) and D-2 (0.814) scores. This suggests that balancing figurative language and cause annotations yields responses that are not only linguistically diverse but also coherent and contextually aligned.

After *Combined*, *Fig* achieves the best overall performance, reducing PPL by 35.8% (14.182→9.100), increasing D-1 by 8.9% (0.515→0.561), D-2 by 8.4% (0.751→0.814), and BLEU by 1.4% (0.764→0.775) compared to *Plain*. While both *Fig* and *Cause* enhance response quality, *Fig* has a more pronounced impact on lexical diversity (D-1, D-2) and overall coherence (PPL), whereas *Cause* demonstrates modest improvements in semantic alignment, with a 1.1% increase in BLEU (0.764→0.772).

In contrast, the *Zero-shot* baseline, despite achieving high diversity scores (D-1: 0.587, D-2: 0.847), shows extremely poor performance in other key metrics. Its BLEU score is especially low at 0.058 and it has the highest PPL (14.944), reflecting poor fluency and alignment with gold responses. Due to these limitations, we excluded *Zero-shot* from further human evaluation and analysis, as its low-quality outputs would not provide meaningful insights into the effectiveness of different fine-tuned approaches. Manual inspection further confirmed that many responses were frequently incoherent or contextually misaligned.

6.2 Human Evaluation

The results of human evaluation on the fine-tuned approaches are summarized in Table 3. Similar to the performance on automated metrics, *Combined* (*Fig + Cause*) demonstrates the best overall performance across all criteria, achieving the highest average scores for Empathy (4.426), Relevance (4.135), and Fluency (4.189). *Fig* follows closely, with strong scores for Empathy (4.195) and Fluency (4.132). It also performs comparably to *Cause* in Relevance (4.021 vs. 4.024), indicating that figurative language alone can align responses well with the speaker’s context. *Cause* excels in Relevance (4.024), confirming that it effectively addresses the content of the speaker’s text, with less pronounced scores for Empathy (3.889) and Fluency (3.799) compared to *Fig* and *Combined*. In contrast, *Plain* lags behind across all three metrics, with the lowest scores for Empathy (3.565), Relevance (3.631),

| Model | tone_pos | pro-social | cog-proc | adj |
|----------|--------------|--------------|---------------|--------------|
| Gold | 3.495 | 0.966 | 15.940 | 6.910 |
| Plain | 2.686 | 0.892 | 15.147 | 6.249 |
| Cause | 2.478 | 0.814 | 15.347 | 6.181 |
| Fig | 3.560 | 1.414 | 16.304 | 6.718 |
| Combined | 3.603 | 1.425 | 16.136 | 6.927 |

Table 4: LIWC analysis of *Gold* & generated responses for different approaches. Features include tone_pos, reflecting positive tone; prosocial, capturing supportive language; cogproc, representing cognitive engagement and contextual reasoning; and adj, measuring linguistic richness through descriptive adjectives. For consistency, multiple gold responses corresponding to a single speaker post were aggregated to a single representation.

and Fluency (3.207).

For the Most Supportive metric, which reflects perceived supportiveness (as described in Section 5), 84.7% (94 out of 111) of evaluated samples reach majority agreement, with at least two annotators selecting the same response. Among these, the responses generated by *Combined (Fig + Cause)* are selected the most frequently as being the Most Supportive, taking up 52.13% (49 responses out of 94) of evaluated samples. *Fig* follows, with 32.98% (31 responses out of 94), while *Cause* and *Plain* are selected less frequently, with 11.7% (11 responses) and 3.19% (3 responses) of evaluated samples, respectively.

7 Analysis of Generated Responses

In this section, we analyze the generated responses to gain deeper insights into various dimensions of empathetic expression.

Psycholinguistic Insights

Using LIWC (Tausczik and Pennebaker, 2010) psycholinguistic features, we examine emotional, social, cognitive, and linguistic aspects related to empathy in generated and gold responses (*Gold*). We use the LIWC 2022 edition⁶ to extract and select four psycholinguistic features from each response:

- **tone_pos**: Encompasses words related to positive emotions (Tausczik and Pennebaker, 2010). Their presence can contribute to creating uplifting and supportive responses.

- **prosocial**: Captures social supportiveness, reflecting language that signals a willingness to help or show care (Pennebaker et al., 2015).
- **cogproc**: Indicates cognitive engagement, such as reasoning and understanding. Ensures that the response is relevant and thoughtful.
- **adj**: Measures the use of descriptive adjectives, capturing the vividness and expressiveness of the responses.

Table 4 provides results for these selected features. *Combined (Fig + Cause)* demonstrates the most well-rounded performance, surpassing both *Gold* and other generated methods in most metrics (tone_pos: 3.603, prosocial: 1.425, adj: 6.927), except in cogproc, where it ranks second. Overall, it effectively balances positive tone, social supportiveness, cognitive engagement, and linguistic richness in the generated responses.

Fig also excels, achieving the highest cognitive engagement score in cogproc (16.304 vs. *Gold*: 15.940). This demonstrates how figurative language can enhance reasoning and thoughtful engagement beyond emotion expression. It also significantly boosts the score for prosocial (1.414 vs. *Gold*: 0.966) which makes it particularly effective in shaping socially supportive responses. While *Fig* does not surpass *Gold* in adj (6.718 vs. *Gold*: 6.910), it remains competitive in its role to leverage descriptive adjectives in empathetic text.

Cause shows nuanced results, with a slight improvement in cogproc (15.347, an increase from *Plain*: 15.147) but a lower tone_pos score (2.478 vs. *Plain*: 2.686). This suggests that while *Cause* enhances reasoning, it may benefit from complementary strategies to elevate positive tone and social supportiveness. These findings highlight the synergy between *Cause* and *Fig*, as evidenced by *Combined*, which effectively balances their strengths to enhance empathetic responses.

What Makes a Response Supportive?

We extend our analysis beyond all generated responses to focus on the Most Supportive responses, identified by majority agreement during human evaluation (see Section 5). These responses were compared against others to explore the role of Empathy, Relevance, and Fluency in determining perceived supportiveness.

Our analysis reveals that a balance among Empathy, Relevance, and Fluency is critical for per-

⁶<https://www.liwc.app/>

| Response Type | E | R | F |
|-----------------|-------------|-------------|-------------|
| Most Supportive | 4.67 | 4.37 | 4.40 |
| Other Responses | 3.78 | 3.80 | 3.64 |

Table 5: Average scores for Empathy (E), Relevance (R), and Fluency (F) in Most Supportive and Other Responses. All differences were tested using a paired t -test and found statistically significant ($p < 0.001$).

ceived supportiveness (See Table 5). While Empathy scores were consistently high for Most Supportive responses (average: 4.67), high empathy alone was insufficient. When we observed responses that received a perfect empathy score (5) but weren’t selected as being the Most Supportive, 64.71% (22 out of 34) of such responses had lower fluency scores (average: 3.18 vs. Most Supportive: 4.40). Similarly, 35.29% (12 out of 34) of responses with perfect empathy scores that were not selected as Most Supportive had lower relevance (average: 3.78 vs. Most Supportive: 4.37). While low relevance may have some influence, it appears to be a less critical breaking factor than fluency. This is supported by our effect size analysis using Cohen’s d (fluency: 1.24, relevance: 0.85), aligning with research that frames supportiveness as a multidimensional construct requiring high empathy, contextual alignment, and linguistic clarity (Cutrona, 1990; Halpern, 2001; Burleson, 2003).

8 Conclusion

In this study, we explored a novel approach to empathetic response generation by integrating figurative language and semantic context signals via manually annotated empathy causes. This integration significantly improves the response quality across emotional, contextual, and linguistic dimensions, as demonstrated by automated metrics: BLEU improves by 7.6% (0.764 \rightarrow 0.822), PPL decreases by 36.7% (14.182 \rightarrow 8.980), and lexical diversity increases by 7.6% (D-1: 0.515 \rightarrow 0.569, D-2: 0.751 \rightarrow 0.814) from *Plain* fine-tuning. Human evaluations affirm that *Combined* (Fig + Cause) achieves the highest ratings for Empathy (4.426, +24.2%), Relevance (4.135, +13.9%), and Fluency (4.189, +30.6%) out of 5. These findings, supported by our psycholinguistic analysis and exploration of the Most Supportive responses, underscore the synergy between rhetoric and contextual alignment when generating empathetic responses. Overall, this study advances empathetic response genera-

tion by investigating this balance, moving beyond conventional approaches focused solely on understanding speaker’s emotions. It also offers valuable insights into what makes a response truly supportive and engaging. Our work provides a more holistic approach to empathetic communication, addressing underlying factors that drive emotionally and socially supportive interactions.

9 Limitations

Our study is limited in several aspects. Human evaluation inherently involves subjectivity, which could introduce variability in assessing empathy, relevance, and fluency. Efforts were made to mitigate this by carefully crafting an evaluation guideline and shuffling responses; however, subjectivity remains a potential limitation. Additionally, the findings are based on the AcnEmpathize dataset, which focuses on an acne support community. The results may not necessarily generalize to other contexts. In future work, researchers are encouraged to test and adapt these strategies to diverse domains that require support, while also comparing with state-of-the-art (SOTA) methods and more advanced models to fully explore their potential.

10 Ethical Considerations

This study utilizes the AcnEmpathize dataset, which is based on publicly available and anonymized data, ensuring compliance with ethical standards for research involving online communities. The dataset does not include any personal identifying information, and all annotation tasks were conducted by volunteers who were informed about the research goals and methods. The dataset and annotations are intended solely for research purposes, with the aim of advancing empathetic communication through computational methods.

Acknowledgments

Writing quality in early drafts of some portions of this manuscript was checked and occasionally revised using ChatGPT; in later drafts, writing quality was manually reviewed and edited.

References

- Azy Barak, Meyran Boniel-Nissim, and John Suler. 2008. Fostering empowerment in online support groups. *Computers in human behavior*, 24(5):1867–1883.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. "O'Reilly Media, Inc."
- Brant R Burleson. 2003. The experience and effects of emotional support: What the study of cultural and gender differences can tell us about close relationships, emotion, and interpersonal communication. *Personal relationships*, 10(1):1–23.
- Boxing Chen and Colin Cherry. 2014. A systematic comparison of smoothing techniques for sentence-level bleu. In *Proceedings of the ninth workshop on statistical machine translation*, pages 362–367.
- Francesca MM Citron and Adele E Goldberg. 2014. Metaphorical sentences are more emotionally engaging than their literal counterparts. *Journal of cognitive neuroscience*, 26(11):2585–2595.
- CE Cutrona. 1990. Type of social support and specific stress: Toward a theory of optimal matching. *Social support: An interactional view/Wiley*.
- Jean Decety and Philip L Jackson. 2004. The functional architecture of human empathy. *Behavioral and cognitive neuroscience reviews*, 3(2):71–100.
- Susan R Fussell and Mallie M Moss. 2014. Figurative language in emotional communication. In *Social and cognitive approaches to interpersonal communication*, pages 113–141. Psychology Press.
- Jun Gao, Yuhan Liu, Haolin Deng, Wei Wang, Yu Cao, Jiachen Du, and Ruifeng Xu. 2021. Improving empathetic response generation by recognizing emotion cause in conversations. In *Findings of the association for computational linguistics: EMNLP 2021*, pages 807–819.
- David Green et al. 2005. *Troubled talk: Metaphorical negotiation in problem discourse*, volume 15. Walter de Gruyter.
- Jordi Halpern. 2001. From detached concern to empathy: Humanizing medical practice.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Jan-Christoph Klie, Michael Bugert, Beto Boullosa, Richard Eckart De Castilho, and Iryna Gurevych. 2018. The inception platform: Machine-assisted and knowledge-oriented interactive annotation. In *Proceedings of the 27th international conference on computational linguistics: System demonstrations*, pages 5–9.
- Klaus Krippendorff. 1970. Estimating the reliability, systematic error and random error of interval data. *Educational and psychological measurement*, 30(1):61–70.
- Huiyuan Lai, Antonio Toral, and Malvina Nissim. 2023. Multilingual multi-figurative language detection. *arXiv preprint arXiv:2306.00121*.
- Gyeongeeun Lee and Natalie Parde. 2024. Acnempathize: A dataset for understanding empathy in dermatology conversations. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 143–153.
- Gyeongeeun Lee, Zhu Wang, Sathya N Ravi, and Natalie Parde. 2024a. Empatheticfig at wassa 2024 empathy and personality shared task: Predicting empathy and emotion in conversations with figurative language. In *Proceedings of the 14th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis*, pages 441–447.
- Gyeongeeun Lee, Christina Wong, Meghan Guo, and Natalie Parde. 2024b. Pouring your heart out: Investigating the role of figurative language in online expressions of empathy. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 519–529.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Yanran Li, Ke Li, Hongke Ning, Xiaoqiang Xia, Yalong Guo, Chen Wei, Jianwei Cui, and Bin Wang. 2021. Towards an online empathetic chatbot with emotion causes. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pages 2041–2045.
- Zhaojiang Lin, Peng Xu, Genta Indra Winata, Farhad Bin Siddique, Zihan Liu, Jamin Shin, and Pascale Fung. 2020. Caire: An end-to-end empathetic chatbot. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 13622–13623.
- Navonil Majumder, Pengfei Hong, Shanshan Peng, Jiansun Lu, Deepanway Ghosal, Alexander Gelbukh, Rada Mihalcea, and Soujanya Poria. 2020. Mime: Mimicking emotions for empathetic response generation. *arXiv preprint arXiv:2010.01454*.
- John A Naslund, Kelly A Aschbrenner, Lisa A Marsch, and Stephen J Bartels. 2016. The future of mental health care: peer-to-peer support and social media. *Epidemiology and psychiatric sciences*, 25(2):113–122.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318.

| | | | |
|-----|--|---|-----|
| 714 | James W Pennebaker, Ryan L Boyd, Kayla Jordan, and | Anuradha Welivita and Pearl Pu. 2023. Use of a tax- | 767 |
| 715 | Kate Blackburn. 2015. The development and psycho- | onomy of empathetic response intents to control and | 768 |
| 716 | metric properties of liwc2015. | interpret empathy in neural chatbots. <i>arXiv preprint</i> | 769 |
| | | <i>arXiv:2305.10096</i> . | 770 |
| 717 | Yushan Qian, Wei-Nan Zhang, and Ting Liu. 2023. | Anuradha Welivita, Yubo Xie, and Pearl Pu. 2021. A | 771 |
| 718 | Harnessing the power of large language models | large-scale dataset for empathetic response gener- | 772 |
| 719 | for empathetic response generation: Empirical in- | ation. In <i>Proceedings of the 2021 Conference on</i> | 773 |
| 720 | vestigations and improvements. <i>arXiv preprint</i> | <i>Empirical Methods in Natural Language Processing</i> , | 774 |
| 721 | <i>arXiv:2310.05140</i> . | pages 1251–1264. | 775 |
| 722 | Aravind Sesagiri Raamkumar and Yinping Yang. 2022. | Anuradha Welivita, Chun-Hung Yeh, and Pearl Pu. 2023. | 776 |
| 723 | Empathetic conversational systems: A review of cur- | Empathetic response generation for distress support. | 777 |
| 724 | rent advances, gaps, and opportunities. <i>IEEE Trans-</i> | In <i>Proceedings of the 24th Annual Meeting of the</i> | 778 |
| 725 | <i>actions on Affective Computing</i> , 14(4):2722–2739. | <i>Special Interest Group on Discourse and Dialogue</i> , | 779 |
| | | pages 632–644. | 780 |
| 726 | Hannah Rashkin, Eric Michael Smith, Margaret Li, and | Thomas Wolf, Lysandre Debut, Victor Sanh, Julien | 781 |
| 727 | Y-Lan Boureau. 2019. Towards empathetic open- | Chaumond, Clement Delangue, Anthony Moi, Pier- | 782 |
| 728 | domain conversation models: A new benchmark and | ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, | 783 |
| 729 | dataset . In <i>Proceedings of the 57th Annual Meet-</i> | et al. 2020. Transformers: State-of-the-art natural | 784 |
| 730 | <i>ing of the Association for Computational Linguistics</i> , | language processing. In <i>Proceedings of the 2020 con-</i> | 785 |
| 731 | pages 5370–5381, Florence, Italy. Association for | <i>ference on empirical methods in natural language</i> | 786 |
| 732 | Computational Linguistics. | <i>processing: system demonstrations</i> , pages 38–45. | 787 |
| 733 | Helen Riess and Gordon Kraft-Todd. 2014. Empathy: a | L. Xue. 2020. mt5: A massively multilingual pre- | 788 |
| 734 | tool to enhance nonverbal communication between | trained text-to-text transformer. <i>arXiv preprint</i> | 789 |
| 735 | clinicians and their patients. <i>Academic Medicine</i> , | <i>arXiv:2010.11934</i> . | 790 |
| 736 | 89(8):1108–1112. | | |
| 737 | Tulika Saha, Vaibhav Gakhreja, Anindya Sundar Das, | | |
| 738 | Souhitya Chakraborty, and Sriparna Saha. 2022. To- | | |
| 739 | wards motivational and empathetic response gener- | | |
| 740 | ation in online mental health support. In <i>Proceedings</i> | | |
| 741 | <i>of the 45th international ACM SIGIR conference on</i> | | |
| 742 | <i>research and development in information retrieval</i> , | | |
| 743 | pages 2650–2656. | | |
| 744 | Ashish Sharma, Monojit Choudhury, Tim Althoff, and | | |
| 745 | Amit Sharma. 2020a. Engagement patterns of peer- | | |
| 746 | to-peer interactions on mental health platforms. In | | |
| 747 | <i>Proceedings of the International AAAI Conference on</i> | | |
| 748 | <i>Web and Social Media</i> , volume 14, pages 614–625. | | |
| 749 | Ashish Sharma, Adam S Miner, David C Atkins, and | | |
| 750 | Tim Althoff. 2020b. A computational approach to un- | | |
| 751 | derstanding empathy expressed in text-based mental | | |
| 752 | health support. <i>arXiv preprint arXiv:2009.08441</i> . | | |
| 753 | Jocelyn Shen, Joel Mire, Hae Won Park, Cynthia | | |
| 754 | Breazeal, and Maarten Sap. 2024. Heart-felt narra- | | |
| 755 | tives: Tracing empathy and narrative style in personal | | |
| 756 | stories with llms. <i>arXiv preprint arXiv:2405.17633</i> . | | |
| 757 | Yla R Tausczik and James W Pennebaker. 2010. The | | |
| 758 | psychological meaning of words: Liwc and comput- | | |
| 759 | erized text analysis methods. <i>Journal of language</i> | | |
| 760 | <i>and social psychology</i> , 29(1):24–54. | | |
| 761 | Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier | | |
| 762 | Martinet, Marie-Anne Lachaux, Timothée Lacroix, | | |
| 763 | Baptiste Rozière, Naman Goyal, Eric Hambro, | | |
| 764 | Faisal Azhar, et al. 2023. Llama: Open and effi- | | |
| 765 | cient foundation language models. <i>arXiv preprint</i> | | |
| 766 | <i>arXiv:2302.13971</i> . | | |

A Appendix

A.1 Empathy Cause Annotation

The annotation process for empathy cause identification was designed to be straightforward and intuitive. While no formal instruction document was provided, the following general guidelines were followed during annotation:

- **Task Overview:**

- For each empathetic reply, annotators were instructed to highlight at least one sentence in the corresponding initial post that was perceived as the cause of empathy in that reply.

- **Annotation Criteria:**

- Cause sentences typically contained expressions of frustration, sadness, or distress (e.g., "I feel so sad and hopeless").
- Statements describing personal experiences that others may relate to were also considered (e.g., "I have cystic acne on my body...it's genetic, so it's extra difficult.").

- **Annotation Process and Adjustments:**

- Annotators initially discussed definitions of empathy and its potential triggers to ensure consistency, following [Lee and Parde \(2024\)](#) to maintain uniformity within the dataset.
- Due to frequent issues with automatic sentence detection in INCEpTION ([Klie et al., 2018](#)), annotators treated spans ending in punctuation as full sentences.
- To simplify the process, annotators excluded threads with no replies, posts lacking textual content (e.g., only emojis), replies quoting other posts without independent content, and extremely long posts over 1K sentences to prevent annotation errors.

A.2 Example of Data Formatting for Fine-Tuning

All training data was formatted by marking specific linguistic features in speaker posts using XML-style tags before fine-tuning. These tags were applied in the same way across different experimental settings:

- **Cause Annotations:** `<cause>...</cause>` (to mark text likely to trigger empathy)

- **Figurative Language:**

- `<idiom>...</idiom>` (to mark idiomatic expressions)
- `<metaphor>...</metaphor>` (to mark metaphorical comparisons)
- `<hyperbole>...</hyperbole>` (to mark exaggeration)

- **Combined (Figurative + Cause):** Both cause and figurative language tags were used together.

This tagged version of the data was fed into the model for fine-tuning, following the same format as the original dataset. The only modification was the addition of cause and figurative language tags, allowing the model to learn associations between these markers and appropriate empathetic responses.

A.3 Human Evaluation Guideline

Human annotators were provided with a common file containing 111 samples, each consisting of a speaker text and four generated responses (see Table 6). Each annotator was given an individual evaluation file (see Table 7) to record scores for the generated responses with the following instruction:

- **Instruction:** For each entry, read the "Speaker text" and four responses. Rate empathy (E), relevance (R), and fluency (F) for each response on a scale of 1-5 (1: not at all, 3: somewhat, 5: very likely), using the format ERF. Finally, choose the response that feels most supportive.

They were also provided with the definitions and examples of each measure, as detailed below.

A.3.1 Empathy (E)

- **Definition:** Does the response show understanding of the speaker's feelings?

*Note: Empathy doesn't necessarily involve having the exact same experience or simply agreeing to the speaker. If a response includes any of the following empathy communication mechanisms (adapted from [Sharma et al. \(2020b\)](#)), you can assign at least a 3.

| Index # | Speaker text | Response 1 | Response 2 | Response 3 | Response 4 |
|---------|---|--|---|---|---|
| 1 | At the end of June I'm moving in with my current boyfriend. It makes me sad, however, that the ONLY thing that is bothering me about moving in with him is my acne. | I totally understand how you feel. My boyfriend has been so supportive, but I think it's just natural for me to feel embarrassed about my skin around him. | I feel the same way. My boyfriend is always telling me how beautiful I am and that he doesn't care about my acne, but I still feel like he deserves better. | I know exactly how you feel. I've been living with my boyfriend for 3 years now and was really self-conscious about my skin. It's tough, but I've learned that you can't let it get to you. | I can relate to this. I've been seeing my boyfriend for 4 months now and he has seen me at my worst. It's a tough situation. I'm sure you'll figure it out. |
| : | : | : | : | : | : |

Table 6: Example common file entries provided to annotators, showing the speaker text and four generated responses used for evaluation.

| Index # | Speaker text | Response 1 | Response 2 | Response 3 | Response 4 | Most Supportive |
|---------|---|------------|------------|------------|------------|-----------------|
| 1 | At the end of June I'm moving in with my current boyfriend. It makes me sad, however, that the ONLY thing that is bothering me about moving in with him is my acne. | ERF | ERF | ERF | ERF | Choose from 1-4 |
| : | : | : | : | : | : | : |

Table 7: Example evaluation entries provided to annotators for scoring Empathy (E), Relevance (R), Fluency (F), and selecting the Most Supportive response from the four responses.

- **Emotional Reactions:** Does the response express or allude to warmth, compassion, concern, or similar feelings of the responder towards the seeker? (e.g., *Everything will be fine; I feel really sad for you.*)
 - **Interpretations:** Does the response communicate an understanding of the seeker's experiences and feelings? In what manner? (e.g., *I understand how you feel; This must be terrifying; I also have anxiety attacks at times which makes me really terrified.*)
 - **Explorations:** Does the response make an attempt to explore the seeker's experiences and feelings? (e.g., *What happened?; Are you feeling alone right now?*)
 - **Example:**

Speaker text: I have acne and worry that my boyfriend will think it's gross.

Responses:

 - *Just get over it.* (1)
 - *A lot of people worry about their acne around others.* (3)
 - *I completely understand feeling self-conscious about acne, especially around people who matter to you. I've felt that way too.* (5)
- A.3.2 Relevance (R)**
- **Definition:** Is the response appropriate to the conversation? Is it on-topic?
 - **Example:**

Speaker text: I have acne and worry that my boyfriend will think it's gross.

Responses:

 - *I hope to get hired soon.* (1)
 - *A lot of people feel self-conscious about their skin.* (3)
 - *It's understandable to feel self-conscious about acne around someone you care about, like your boyfriend.* (5)
- A.3.3 Fluency (F)**
- **Definition:** Is the response easy to understand? Does it flow smoothly?

927 • **Example:**

928 **Speaker text:** I hate acne.

929 **Responses:**

- 930 – *I acne understand your concerns about.*
931 (1)
- 932 – *Acne is annoying. It is tiring. It is bad.*
933 (3)
- 934 – *I understand your frustration with acne.*
935 *It's tough to deal with every day, and it*
936 *can be tiring. (5)*

937 **A.3.4 Most Supportive**

- 938 • **Definition:** Imagine you are the person who
939 shared the concerns in the “Speaker text” col-
940 umn. Which of the four responses (“Response
941 1”, “Response 2”, “Response 3”, “Response
942 4”) would make you feel the most supported?