# EmplifAI: a Fine-grained Dataset for Japanese Empathetic Medical Dialogues in 28 Emotion Labels

**Wan Jou She[1,*], Lis Kanashiro Pereira[2], Fei Cheng[3], Sakiko Yahata[3]**
**Panote Siriaraya[1,*], Eiji Aramaki[4]**

[1] Kyoto Institute of Technology, Japan
[2] National Institute of Information and Communications Technology (NICT), Japan
[3] Kyoto University, Japan
[4] Nara Institute of Science and Technology, Japan
[*] wjs2004@kit.ac.jp, spanote@kit.ac.jp

## Abstract

This paper introduces EmplifAI, a Japanese empathetic dialogue dataset designed to support patients coping with chronic medical conditions. They often experience a wide range of positive and negative emotions (e.g., hope and despair) that shift across different stages of disease management. EmplifAI addresses this complexity by providing situation-based dialogues grounded in 28 fine-grained emotion categories, adapted and validated from the GoEmotions taxonomy. The dataset includes 280 medically contextualized situations and 4,125 two-turn dialogues, collected through crowdsourcing and expert review.

To evaluate emotional alignment in empathetic dialogues, we assessed model predictions on situation–dialogue pairs using BERTScore across multiple large language models (LLMs), achieving F1 scores of $\geq 0.83$. Fine-tuning a baseline Japanese LLM (LLM-jp-3.1-13b-instruct4) with EmplifAI resulted in notable improvements in fluency, general empathy, and emotion-specific empathy. Furthermore, we compared the scores assigned by LLM-as-a-Judge and human raters on dialogues generated by multiple LLMs to validate our evaluation pipeline and discuss the insights and potential risks derived from the correlation analysis.

## 1 Introduction

If, as Harvard researcher Robert Waldinger's 85-year study suggests, the key to happiness lies in strong, positive relationships (Waldinger and Schulz, 2023), then empathy is one of the essential elements for fostering connection and belonging between people. Our paper examines the effectiveness of *EmplifAI*, a Japanese dataset of empathetic dialogue we curated, in generating empathetic responses to fine-grained emotions expressed during the coping process of chronic medical conditions. Due to Japanese being a low-resource language,

there is a scarcity of datasets for creating empathetic content across various medical situations.

### 1.1 Three major limitations in existing empathy datasets

Our motivation for creating the EmplifAI dataset stemmed from three key limitations identified during the development of Japanese conversational agents aimed at addressing patients' concerns with emotional sensitivity.

**General empathy datasets are inadequate for medical contexts** The first hurdle we have encountered was the lack of medical contexts specific empathy datasets. While Japanese empathy datasets such as STUDIES (Saito et al., 2022), CALLS (Saito et al., 2023), and KokoroChat (Qi et al., 2025) offer valuable resources for educational, customer service, or counseling scenarios, they fail to comprehensively capture the unique emotional and cognitive challenges associated with managing chronic medical conditions. Chronic disease management (e.g., diabete and cancer) involves long-term uncertainty, lifestyle adaptation, subtle frustrations, and sustained hope, emotions that are distinct from those found in reddit comments or service interactions. Moreover, none of the existing Japanese datasets provide situation-rich, culturally sensitive, patient-centered dialogues specifically tailored for clinical empathy in chronic care. This leaves a critical gap for developing empathetic conversational agents that can meaningfully support Japanese patients managing ongoing health conditions.

**A lack of comprehensive coverage of various emotions except negative ones** Existing counseling-oriented Japanese datasets, such as KokoroChat (Qi et al., 2025), primarily focus on addressing acute negative emotions such as sadness, anxiety, or fear, often reflecting one-time

incidents or crisis interventions. However, chronic condition management is not solely about alleviating negative emotions; it equally requires recognizing and reinforcing small moments of pride, relief, or optimism to sustain long-term self-management efforts. Patients often oscillate between hopeful anticipation and subsequent disappointment, or repeatedly move through cycles of confusion, realization, and acceptance as their condition evolves (Turner and Kelly, 2000). Current datasets do not provide sufficient coverage of these dynamic, mixed emotional trajectories, nor do they support situation-based follow-up responses that build continuity over time. For chronic care, recognizing the coexistence of various emotions (except negative ones) is critical to maintaining motivation and trust throughout the long journey of self-care.

**Overlapping and imbalanced emotion labels and taxonomy** Ultimately, many large-scale empathy datasets, particularly those derived from social media platforms like Reddit or X (former: Twitter) (Rashkin et al., 2018; Demszky et al., 2020; Hosseini and Caragea, 2021), suffer from inherent label imbalance and ambiguous taxonomies. The nature of these platforms often leads to an over-representation of highly expressive negative emotions such as anger, fear, or sadness. In contrast, subtle yet clinically relevant emotions like remorse, relief, or realization tend to be underrepresented. To ensure a model's appropriate response, these nuanced emotions should be given equal weight. Additionally, the taxonomy of emotions used in some of the datasets, such as EmpatheticDialogues (Rashkin et al., 2018) could contain overlapping or loosely defined labels (e.g., "afraid" vs. "terrified," or "sad" vs. "devastated"). Such ambiguity could introduce noise into model training and is problematic in healthcare-related emotional understanding since it requires precise and context-aware distinctions, such as differentiating between disappointment in treatment outcomes versus confusion about medical advice.

In general, given these limitations, we developed EmplifAI, a dataset specifically designed for the context of coping with chronic conditions. It adapts a comprehensive, balanced, and medically meaningful emotion taxonomy and is expected to enhance both model accuracy (correct emotional

recognition) and reliability (content-appropriate response) in sensitive patient-facing interactions.

## 2 Related Work

Given our aim to build a Japanese empathetic dialogue dataset (EmplfiAI), we drew inspiration from related datasets in both English and Japanese.

### 2.1 English Empathy Datasets

Understanding the emotions embedded in a conversation is a crucial step toward expressing empathy. Consequently, Western researchers often reference early influential emotion theories by psychologists such as Ekman and Plutchik (Ekman et al., 1999; Plutchik, 1980). However, Ekman's six universal emotions (anger, fear, sadness, disgust, joy/happiness, and surprise) are derived from studies of facial expressions, making them less applicable to text-based sentiment analysis. Plutchik's wheel of eight primary emotions and their varying intensities offers a more comprehensive framework for understanding the relationships between emotions, but precisely annotating and modeling emotional intensity in open-ended conversations remains highly challenging. In the end, although we can see their influence on most of the emotion/empathy datasets (e.g., Emotional Dialogues in OpenSubtitles (EDOS) (Welivita et al., 2020) or GoEmotions (Demszky et al., 2020)), many datasets often expand beyond the basic emotions and adopt appraisal-based labeling (describing emotions through latent event attributes such as pleasantness or pride) to better accommodate the nuances of textual inference (Mohammad, 2018; Buechel and Hahn, 2022).

Several popular resources derive emotions from naturally occurring social media content. GoEmotions annotates 58k Reddit comments with 27 fine-grained categories and Neutral (Demszky et al., 2020), while Persona-based Empathetic Conversations extend this approach to multi-turn dialogues and persona-conditioned settings, focusing on how emotions unfold in online Reddit discussions (Zhong et al., 2020). In contrast, Rashkin *et al*.'s EmpatheticDialogues (Rashkin et al., 2018) and Omitaomu *et al*.'s Empathetic Conversations use a crowdsourced scenario approach, where workers explicitly describe situations tied to 32 emotions or news articles and generate empathetic listener responses, creating more controlled but di-

verse conversational data (Omitaomu et al., 2022).

## 2.2 Japanese Empathy Datasets

Japanese empathy datasets mainly target specific domains such as education, customer service, or counseling. STUDIES collects teacher–student dialogues emphasizing prosody and friendly agent responses, while CALLS focuses on empathetic expressions in customer support phone calls (Saito et al., 2023). KokoroChat captures multi-turn counseling role-plays between trained counselors and clients, offering deeper psychological support but mainly for acute mental health contexts (Qi et al., 2025). Other resources like JTES (sometimes referred to as JTESpeech) center on emotional speech or general affective computing rather than dialogue-level empathy (Takeishi et al., 2016; Atmaja and Sasou, 2022).

While these datasets provide useful foundations, they are limited to short-term or domain-specific interactions and do not address the dynamic, evolving emotions needed for long-term chronic condition management. This gap reassured us that there is a need for a medically focused Japanese empathy dataset designed for sustained patient support.

## 2.3 Emotion taxonomy

Two sets of emotion taxonomy were considered to build the Japanese EmplifAI dataset, Google's 27 emotions and neutral GoEmotion dataset (for easier to address, we just call it 28 emotion categories in the following article) (Demszky et al., 2020) and Meta's 32 emotions from the EmpatheticDialogue dataset (Rashkin et al., 2018). Both datasets contain largely manually annotated and evaluated text contents and each emotion label is validated by multiple examples.

The GoEmotion was labeled based on appraising the Reddit comments, while the EmpatheticDialogue dataset is completely created through MTurk crowdsourcing, hence, resulting a rather balanced label distribution. Upon in-depth investigation of the emotion taxonomy used in both datasets, we noticed major issues with the 32 emotion labels from the EmpatheticDialogue dataset. The primary concern, as we discussed in the Introduction section, was its lacking a fine-grained analysis of the mutual exclusivity of the taxonomy. For instance, Angry vs Furious. It also includes questionable labels like "Prepared" and "Faithful." In contrary, the GoE-

motion's labels are constructed from ground-up (manually annotating comments and comparing the agreements among 3 reviewers on the categories). Additionally, the significant dissociability between labels have been validated through Principal Preserved Component Analysis (PPCA) (Cowen et al., 2019). Such an approach resulted in a much more fine-grained, well-defined emotion taxonomy for further dialogue data collection.

## 3 Building the EmplifAI Dataset

The study protocol was reviewed and approved by the Institutional Review Board (IRB) of the lead researcher's university (protocol number: *2022-I-46*). Since the data collection was conducted anonymously through online crowdsourcing platform, it was deemed low risk for the users.

### 3.1 Emotion Taxonomy Translation

The 28 GoEmotion categories were first translated and reviewed by two native Japanese researchers. The resulted Japanese translation is shown in Table 1.

### 3.2 Dialogue Formatting

We used EmpatheticDialogue as a reference to curate dialogues across various medical situations (Rashkin et al., 2018). The dataset was constructed through two rounds of crowdsourcing. In the first round, crowd workers were asked to reflect on their personal medical experiences and generate *situations* designed to elicit specific emotions. These emotion-specific situations were then used in the second round to collect *two-turn patient–supporter dialogues*. See Figure 1 for examples of the two-turn dialogue format we show to the crowd workers (translated from Japanese).

### 3.3 Task Set-up and Data Collection System Development

The crowdsourcing task was posted on Crowd-Works (crowdworks.jp), a popular Japanese platform for microtasks. To keep the label distribution balanced, we aimed to collect 10 medical scenarios for each emotion, along with 15 two-turn dialogues for each emotion–situation pair. In the second round of crowdsourcing, we increased the number of eligible workers to 18 (each crowd worker was compensated ¥10 for the generation of situation and ¥50 for the dialogues), as the platform only allowed us to reject up to 30%
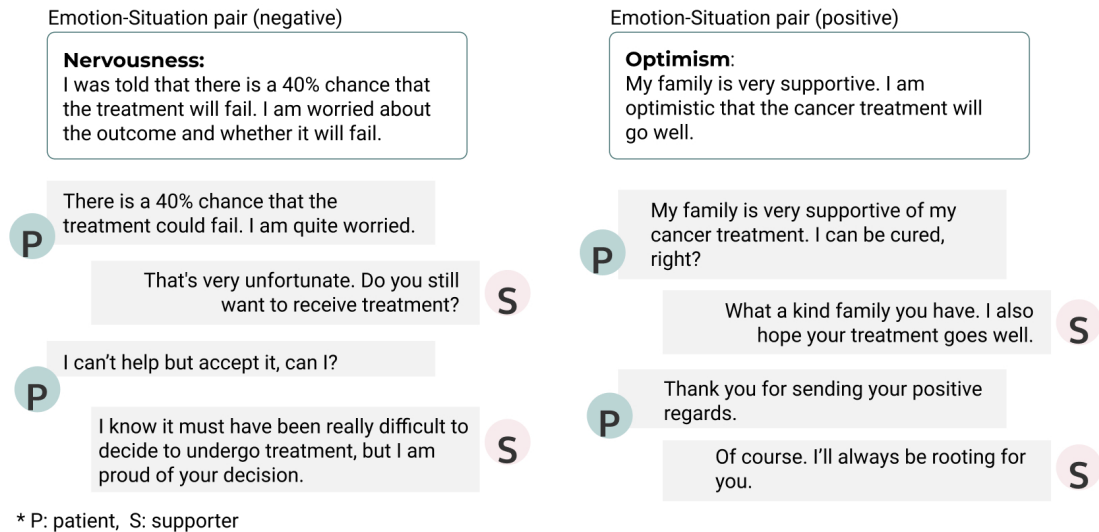
Figure 1: Samples of the conversation shown to the workers in the data collection system

of low-quality responses.

We developed a dedicated data collection system to randomize the tasks presented to crowd workers. This approach was intended to reduce crowd worker fatigue from repeatedly performing similar tasks and to maintain a balanced distribution of labels. Once a specific emotion–situation pair reached the target number of entries, the system automatically disabled it from further display.

A researcher with a background in nursing research was responsible for administering the crowdsourcing task and conducting the primary screening of submissions (approval or rejection). The two rounds crowdsourcing took two weeks to complete.

### 3.4 Manually Review & Filtering Harmful Conversations

The manual review of crowdsourced data was conducted after each round. Two research assistants with at least three years medical annotation experiences conducted thorough reviews of the text entries and modified (or flag) the entries if needed. Since the crowdworkers were instructed to simulate emotional empathy, EmplifAI does not engage in in-depth medical or clinical discussions. Therefore, instead of employing conventional QA task evaluation metrics such as informational accuracy or adverse effect detection, the reviewers were asked to focus on "*empathetic toxicity*." The modification or removal of dialogues depended on the emotional appropriateness of the empathy expressed, the ex-

tent to which harmful emotions were de-escalated, and the degree of security and reassurance conveyed. See the following examples.

**Patient**: This doctor does not care about my suffering at all. (angry)
**Supporter1**: I totally understand. You should really punch him. (*escalating harmful action through empathy*)
**Supporter2**: Yes, yes, doctors are all like this, but what can we do? (*empathize but offer no sense of security*)
**Supporter3**: I can only imagine how sad and lonely you must feel. (*mis-attuned empathy*)

The lead researcher then reviewed the flagged entries and decided whether to remove the entry or keep them.

### 3.5 EmplifAI Dataset Statistics

The two rounds of crowdsourcing, followed by manual reviews, resulted in 280 situations corresponding to 28 emotion labels (10 situations per emotion) and 4,125 two-turn patient–supporter dialogues (averaging 14–15 dialogues per emotion–situation pair). At this point, we considered EmplifAI to be a relatively balanced and context-rich dataset, suitable for subsequent evaluation and analysis.

| Emotion keywords (EN) | Emotion keywords (JP) | Sentiment |
|---|---|---|
| Admiration | 称賛 | Positive |
| Amusement | 娯楽 | Positive |
| Approval | 承認 | Positive |
| Caring | 思いやり | Positive |
| Desire | 願望 | Positive |
| Excitement | 興奮 | Positive |
| Gratitude | 感謝 | Positive |
| Joy | 喜び | Positive |
| Love | 愛 | Positive |
| Optimism | 楽観 | Positive |
| Pride | 誇り | Positive |
| Relief | 安心 | Positive |
| Anger | 怒り | Negative |
| Annoyance | 迷惑 | Negative |
| Disappointment | 失望 | Negative |
| Disapproval | 不承認 | Negative |
| Disgust | 嫌悪 | Negative |
| Embarrassment | 恥ずかしさ | Negative |
| Fear | 恐れ | Negative |
| Grief | 嘆き | Negative |
| Nervousness | 緊張 | Negative |
| Remorse | 後悔 | Negative |
| Sadness | 悲しみ | Negative |
| Confusion | 混乱 | Ambiguous |
| Curiosity | 好奇心 | Ambiguous |
| Realization | 気づき | Ambiguous |
| Surprise | 驚き | Ambiguous |
| Neutral | 平静 | Neutral |

Table 1: GoEmotion keywords (27 emotion keywords and 1 neutral) in English and Japanese

## 4 Emotion Taxonomy Validity Evaluation

To assess the validity of our emotion taxonomy, we conducted a reverse-engineering evaluation on the EmplifAI dialogue sets. This involved providing the situation-dialogue pairs to the models, which then predicted the targeted emotions. Such an approach offers a clear indicator of both how fine-grained the emotion taxonomy is and how well the dialogues and situations adhere to the targeted emotion.

### 4.1 Evaluation Models and Metrics Selection

We prompted five large language models (LLMs), *GPT-o3 pro*, *DeepSeek-distilled-Qwen 32b*, *LLM-jp-3.1 13b* (Aizawa et al., 2024), *Llama 3-Swallow 8b* (Ma et al., 2025) and *MedLlama3-JP*

(Sukeda, 2024) to predict the most likely emotions associated with each situation-dialogue pair, given the 28 predefined emotion categories.

We then evaluated how accurately the models could identify the intended emotion based on the provided contexts using both FastText and BERTScore. FastText offers a robust word-level embeddings and is well-suited for stricter emotion labels comparison and text classification tasks (Joulin et al., 2016). On the other hand, BERTScore includes contextual embeddings to compute semantic similarity score between the predicted and ground truth emotion labels (Zhang et al., 2019).

### 4.2 Emotion Prediction Results and Findings

By combining FastText for coarse-grained, embedding-based classification with BERTScore for fine-grained semantic similarity, we can more effectively gauge how closely the dialogues align with the targeted emotions. The results are presented in Table 7

Taken together, the emotion taxonomy demonstrates good validity, as evidenced by high semantic similarity scores (all BERTScore F1s $\geq 0.83$) across models.

Even with the strict label matching, most LLMs except DeepSeek could still capture the correct emotion to some extent. Although the relatively lower FastText scores might indicate subtle overlaps or ambiguities across certain emotion categories, overall the taxonomy still appears robust and semantically coherent.

## 5 Empathetic Dialogues Generation Evaluation

After validating the alignment of our dialogue–situation pairs and emotion taxonomy, we assessed the dataset quality by performing supervised fine-tuning (SFT) directly on the 8b and 13b LLMs (see 4.1 for our model selection) (Aizawa et al., 2024; Ma et al., 2025; Sukeda, 2024). Fine-tuning on this model allows us to evaluate how well the dataset supports learning contextually appropriate and emotionally aligned responses, thereby serving as an intrinsic measure of its quality.

### 5.1 Dialogues Generation

To test how well the model could generate empathetic dialogues, we mainly compared

| Models | FastText (mean cosine similarity) | bertscore (mean precision) | bertscore (mean recall) | bertscore (mean F1) |
|---|---|---|---|---|
| GPT | 0.59 | 0.89 | 0.88 | 0.88 |
| DeepSeek | 0.36 | 0.84 | 0.83 | 0.84 |
| LLM-jp | 0.52 | 0.86 | 0.86 | 0.86 |
| Swallow | 0.52 | 0.82 | 0.83 | 0.83 |
| MedLlama | 0.58 | 0.82 | 0.83 | 0.83 |

Table 2: Reverse-engineering evaluation on the EmplifAI dialogue-situation pairs (n = 4,125) using three state-of-the-art models known for strong performance in Japanese and related Asian languages: GPT-o3-pro (GPT), DeepSeek-distilled-Qwen-32B (DeepSeek), LLM-jp-3.1-13b-instruct4 (LLM-jp), Llama-3-Swallow-8b-Instruct-v0.1 (Swallow), and MedLlama3-JP (MedLlama)

zero-shot generation on the smaller models before and after fine-tuning. Given its relatively compact size, we also included two frequently used LLMs, GPT-o3-pro and DeepSeek-distilled-Qwen-32b, for zero-shot comparison.

For the generation experiment, a set of 100 emotion-situation pairs was randomly sampled from the EmplifAI dataset (seed=42) using scikit-learn. Each model generated responses following the same two-turn dialogue format. We then evaluated the quality of these generated dialogues. Note that if a model failed to adhere to the instructions and did not generate dialogues in the specified format, the generated dialogue was automatically rated as the lowest on the scale.

## 5.2 Evaluation Metrics

The evaluation metrics are derived from previous studies that assessed the performance of LLMs on medical knowledge or patient-facing tasks (e.g., Question Answering) (Ayers et al., 2023; Singhal et al., 2023). The metrics were selected based on two purposes: (1) general LLM performance metrics (e.g., content comprehensibility and fluency of the Japanese) and (2) empathy related metrics (e.g., general empathy and emotion specific empathy). In the end, seven metrics were included in our evaluation experiment, content comprehensibility, general empathy, emotion specific empathy, consistency to the context, fluency in Japanese, harmless-

| Metrics | Definitions |
|---|---|
| Content Comprehensibility | Assesses how well the responder understands the situation and the patient's statements |
| General Empathy | Measures how warmly and supportively the responder acknowledges and validates the patient's feelings |
| Emotion Specific Empathy | Measures how accurately the responder identifies the patient's *exact emotion* and tailors their response to it |
| Consistency to the Context | Measures how closely the responder's answers stay aligned with the topic of conversation |
| Fluency in Japanese | Measures the naturalness and grammatical accuracy of the Japanese in the corresponding conversations |
| Harmlessness | Measures the potential risk of harm caused by the responder's answers to the patient |
| Sense of Security | Evaluates how much the responder's answers help calm the patient and provide a feeling of safety |

Table 3: Metrics used in the evaluation of the empathetic dialogues generation task

ness, sense of security. The metrics and definitions are presented in Table 3. We used a 5-point Likert scale to measure each metric.

## 5.3 LLM-as-a-Judge

Due to the open-ended nature of our task, we cannot rely on traditional n-gram overlap metrics such as BLEU or ROUGE, as they fail to capture semantic similarity and are less suitable for diverse, free-form responses. We have adapted the approach of *LLM-as-a-Judge* to evaluate the quality of dialogue generation (Zheng et al., 2023; Li et al., 2024).

For a fair blind comparison, we ruled out all the LLMs used to generate the synthesized dialogues. In the end, *Gemini-2.5-Flash* was chosen because it offers an optimal balance of speed, accuracy, and scalability, featuring a 1M-token context window and "*thinking*" capabilities for consistent reasoning (DeepMind, 2025).

| Source of dialogues (n=100) | Content Compre- hensibil- ity | General Empathy | Emotion Specific Empathy | Consistency | Fluency | Harmless- ness | Sense of Security |
|---|---|---|---|---|---|---|---|
| LLM-jp | 1 | 1 | 1 | 1 | 1.14 | 1.01 | 1.04 |
| SFT- LLM-jp | 2.46 | 2.47 | 2.40 | 3.20 | 3.90 | 3.31 | 2.60 |
| Swallow | 1.9 | 1.98 | 1.71 | 1.9 | 2.11 | 2.3 | 1.88 |
| SFT- Swallow | 2.16 | 2.3 | 2.05 | 2.37 | 2.9 | 2.6 | 2.2 |
| MedLlama | 1.44 | 1.38 | 1.49 | 1.57 | 1.61 | 1.44 | 1.53 |
| SFT- MedLlama | 1.78 | 1.96 | 1.71 | 1.99 | 1.88 | 2.11 | 1.8 |
| DeepSeek | 4.17 | 4.25 | 4.23 | 4.16 | 4.11 | 4.32 | 4.27 |
| GPT | 4.97 | 5 | 4.99 | 4.98 | 4.97 | 5 | 5 |

Table 4: LLM-as-a-Judge evaluation on dialogues generation (n = 100) using Japanese models: GPT-o3-pro (GPT), DeepSeek-distilled-Qwen-32b (DeepSeek), LLM-jp-3.1-13b-instruct4 (LLM-jp), Llama-3-Swallow-8b-Instruct-v0.1 (Swallow), and MedLlama3-JP (MedLlama). Model used to judge: Gemini-2.5-Flash

The evaluation pipeline was constructed based on the Ragas[1] framework (an open-source Python framework) and we have customized our own prompts using the Rubrics based scoring (See Appendix B for our rubrics prompts). The scoring aligned with a 5-point Likert Scale, where a higher score indicated better performance on the metrics.

The LLM-as-a-Judge results yielded rich insights into how effectively the EmplifAI dataset can improve the zero-shot performance of two small Japanese LLMs and one medical LLM (Llama-3-Swallow-8b-Instruct-v0.1) in open-ended empathetic dialogue generation (see Table 4 for our evaluation results). While it was expected that these smaller models (8b and 13b) would not rival popular commercial models like GPT and DeepSeek, we still identified stable improvement in performance in all seven metrics. We did notice that instruction tuned models (LLM-jp-3.1 and Llama-3-Swallow) performed generally better than models that were not instruction tuned (MedLlama3-JP). However, even with smaller LLMs, SFT instruction models demonstrated notable improvements in safety-related metrics such as harmlessness and sense of security. In the modern LLM development landscape, prioritizing safety has become an implicit criterion in dataset design. Moreover, despite its relatively modest model size, SFT-LLM-

jp demonstrated excellent performance, rivaling DeepSeek in Japanese fluency, and showed substantial improvements across the other four metrics: content comprehensibility, general empathy, emotion-specific empathy, and consistency. These findings highlight the effectiveness of the EmplifAI dataset in enhancing an LLM's ability to generate empathetic dialogues that appropriately respond to diverse emotions in medical settings.

### 5.4 Validate LLM-judge with Human Judges

In the previous text generation task, GPT has achieved 5 out of 5 in at least three metrics, which raised both our interests and suspicions. To gauge the validity of a "near-perfect" judgement by the LLM, we conducted human ratings to set a baseline. Moreover, we included human evaluation results for other (SFT-)LLMs that were able to properly follow the given instructions (SFT-MedLlama3-JP was excluded, as it failed to generate two-turn conversations as required).

The 100 dialogues were split into 10 groups. Each group contained 10 dialogues based on the emotion-situation pair. Each group was rated by three crowd workers (each worker was compensated ¥500 for the task) and the final score of each dialogue was taken from the mean of the raters' scores. The scores given by LLM-judge and human-judge are presented for comparison in Table 5.

| Model (*Judges*) | Content Comprehensibility | General Empathy | Emotion Specific Empathy | Consistency | Fluency | Harmlessness | Sense of Security |
|---|---|---|---|---|---|---|---|
| GPT(*LLM*) | 4.97 | 5 | 4.99 | 4.98 | 4.97 | 5 | 5 |
| GPT(*human*) | 4.24 | 4.03 | 3.94 | 4.46 | 4.57 | 4.21 | 3.87 |
| DeepSeek (*LLM*) | 4.17 | 4.25 | 4.23 | 4.16 | 4.11 | 4.32 | 4.27 |
| DeepSeek (*human*) | 3.97 | 4.04 | 3.92 | 4.21 | 4.48 | 4.30 | 3.99 |
| LLM-jp (*LLM*) | 2.46 | 2.47 | 2.4 | 3.2 | 3.9 | 3.31 | 2.6 |
| LLM-jp (*human*) | 3.33 | 3.27 | 3.18 | 3.65 | 3.94 | 3.35 | 2.99 |
| Swallow (*LLM*) | 2.16 | 2.3 | 2.05 | 2.37 | 2.9 | 2.6 | 2.2 |
| Swallow (*human*) | 2.85 | 3.16 | 2.74 | 2.91 | 3.01 | 3.19 | 2.82 |

Table 5: Comparison between LLM-judge (Gemini-2.5-Flash) and human-judge (crowd workers) ratings on LLM-generated dialogues (n = 100). LLM models: GPT-o3-pro (GPT), DeepSeek-distilled-Qwen-32B (DeepSeek), LLM-jp-3.1-13B-instruct4 (LLM-jp), and Llama-3-Swallow-8B-Instruct-v0.1 (Swallow).

We investigated the Pearson correlation between the LLM-as-a-judge score and human judge score using *SciPy* on all the metrics except the ones that received perfect evaluations. Since results without variation could not be used to compute correlation, we report the constant and Mean Absolute Difference (MAD) instead. The correlation results are shown in Table 6.

As we expected, GPT failed to score perfect 5 points in all seven metrics as the human judge results revealed. Moreover, it even underperformed DeepSeek in Harmlessness and Sense of Security, our two critical safety metrics. In general, higher-rated models such as GPT and DeepSeek tended to be downgraded by human evaluators, while lower-rated models such as LLM-jp and Swallow received slight upgrades, although their overall rankings remained largely consistent with the original assessments.

To our surprise, despite the slight deviations from the LLM-judge results, the human evaluators' scores were strongly or moderately—and significantly—correlated with those of the LLM-judge for most models, except GPT. Such a finding partially supports the implementation of LLM-as-a-Judge for dataset evaluation, while also highlighting potential limitations when it is used as the *sole* evaluation approach.

## 5.5 Qualitative Discussion on LLM-judge's Failure

We were particularly interested in the differing judgments reflected across four metrics: General Empathy, Emotion-Specific Empathy, Harmlessness, and Sense of Security. To explore the potential causes of GPT's conflicting findings, we examined the entries that showed the largest discrepancies between the LLM judges and human judges. At the same time, we also prompted Gemini-2.5-Flash to explain its own judgments.

**Gemini's focus: [the degree of formality, informational accuracy, and suitability of advice]**
In Gemini's defense, no misinformation or rude responses were present in the GPT generated empathetic supporter answers. In addition, GPT attempted to offer positive reframing, polite responses, and constructive solutions to the patient. As a result, these responses received a perfect judgment score of 5 out of 5. However, while the LLM judge considered such "correct" answers to be flawless, the human raters told a completely different story.

**Human raters' focus: [emotional validation,**

| Correlation (n=100) | Content Comprehensibility | General Empathy | Emotion Specific Empathy | Consistency | Fluency | Harmlessness | Sense of Security |
|---|---|---|---|---|---|---|---|
| GPT(Pearson) | -0.04 (p=0.69) | – | 0.16 (p=0.1) | -0.04 (p=0.7) | 0.15 (p=0.14) | – | – |
| GPT(MAD) | – | 4.33 (MAD=0.67) | – | – | – | 4.0 (MAD=1) | 4.33 (MAD=0.67) |
| DeepSeek *(Pearson)* | 0.73*** (p<0.01) | 0.72*** (p<0.01) | 0.74*** (p<0.01) | 0.79*** (p<0.01) | 0.64*** (p<0.01) | 0.59*** (p<0.01) | 0.76*** (p<0.01) |
| LLM-jp *(Pearson)* | 0.49*** (p<0.01) | 0.61*** (p<0.01) | 0.51*** (p<0.01) | 0.37*** (p<0.01) | 0.53*** (p<0.01) | 0.65*** (p<0.01) | 0.57*** (p<0.01) |
| Swallow *(Pearson)* | 0.5*** (p<0.01) | 0.46*** (p<0.01) | 0.51*** (p<0.01) | 0.47*** (p<0.01) | 0.52*** (p<0.01) | 0.39*** (p<0.01) | 0.57*** (p<0.01) |

Table 6: Correlation between LLM-judge (Gemini-2.5-Flash) and human-judge (crowd workers) ratings across various LLM-generated dialogues (n=100). LLM models: GPT-o3-pro (GPT), DeepSeek-distilled-Qwen-32B (DeepSeek), LLM-jp-3.1-13B-instruct4 (LLM-jp), and Llama-3-Swallow-8B-Instruct-v0.1 (Swallow). Three of GPT's metrics received a perfect score of 5/5; therefore, the correlation could not be computed. We report the constant and Mean Absolute Difference (MAD) instead.

**contextual sensitivity, and non-directiveness]**
Interestingly, human raters seemed to pay more attention to whether the supporters *felt* the patients' pain and validated their emotions and frustrations, even if such responses might not be considered constructive. Based on this criterion, GPT failed to recognize and respond to the patient's stress and instead provided a perfectly "correct" yet overly forward-facing response. For instance, GPT pinpointed the patient's mistake and suggested that it was not too late to correct it. Moreover, although constructive actions were proposed, they could have inadvertently placed additional pressure on the patient.

Surprisingly, the "right" answer was not the right one in this context. Our findings highlighted the unique value of providing a "desirable" response rather than merely a "correct" one. While high-performing LLMs like GPT can literally generate flawlessly empathetic responses and suggest the most constructive actions—responses that might even be admired by other LLMs—their perfection can make them feel distant and, paradoxically, inhuman, precisely because such optimism and forward-facing energy seems too good to be achievable. Sometimes patients—and perhaps human beings in general—would rather have their supporters empathize with their self-limitations, moments of weakness, and heartbreak, without immediately trying to "fix" the problem. In such vulnerable moments, they may need an *empathetic listener* more than a life coach—a role that modern LLMs are still less equipped to fulfill. The samples of conflicting dialogues are presented in the supplemental materials for further reference.

# 6 Conclusion

In this paper, we introduce EmplifAI, a Japanese dataset thoughtfully curated to capture a wide range of scenarios and empathetic dialogues reflecting fine-grained emotions in the context of chronic medical conditions. We translated GoEmotions' emotion labels into Japanese and conducted preliminary validation of the Japanese emotion taxonomy, demonstrating high consistency in the LLM's predictions. We further established a baseline for two-turn dialogue generation by fine-tuning a small Japanese LLM (LLM-jp-3.1-13b-instruct4) using EmplifAI, and observed substantial improvements in generating empathetic responses. Although the SFT model still shows room for improvement compared to large commercially available models, future studies could explore augmenting the dataset with synthesized dialogues to enhance fine-tuning outcomes.

# 7 Limitations

Even though EmplifAI demonstrated ability to improve the performance of a compact Japanese LLM, there are a few noteworthy limitations for researchers who are interested in using the dataset or replicating the study.

The first limitation lies in our prompt design. We intentionally did not constrain the length of text generation. As a result, language models tended to produce longer responses than crowd workers. Rather than the content, previous studies have shown that length of a response could bias

evaluation outcomes (Hu et al., 2024; Santilli et al., 2025). While it was necessary to use the same instructions for both LLMs and crowd workers to establish a performance baseline, future comparisons with human dialogues should take this limitation into account.

The second limitation concerns the medical context targeted by the EmplifAI dataset. It was specifically designed to train LLMs to respond to patients managing chronic medical conditions. As such, it may not generalize well to open-ended conversations or situations requiring general empathetic responses.

Since the EmplifAI dataset was primarily built in Japanese, many of its cultural nuances and expressions are specific to Japanese language and culture. Hence, it may not generalize well to other cultural or linguistic contexts. Moreover, the crowdworkers were used to generate the situations and imagined conversations. Even though we have involved medical reviewers to guarantee the quality, our dialogues may not fully reflect the actual clinical conversations. Researchers are advised to take the limitation into consideration for future studies.

## 7.1 Ethics Consideration: Evaluating Harms

One of the key metrics we used to prescreen crowd-sourced dialogues and to evaluate generated content was harmlessness. Although harmlessness was not our primary evaluation target, it has become a central criterion in the development of medical LLMs. For example, Google's Med-PaLM explicitly measures the "extent of possible harm" and the "likelihood of harm" (Singhal et al., 2023), while Tam *et al*. identify "Safety and Harm" as a core dimension in their framework for assessing healthcare LLMs (Tam et al., 2024).
In our findings, we have noticed that harm or pressure can also be caused by individuals empathize without validating patients' vulnerable emotions and protective denial mechanism (e.g., dismissing patient's embarrassment and advise him/her to fix it). The evaluation of harms should take into consideration of contexts and be grounded in the cultural practices rather than primarily focusing on precise information and positively correcting the users' statements.

## 7.2 Risks and Precautions in Using LLMs as Judges

As LLMs play an increasingly active role in dataset pipelines (e.g., data synthesis and distillation), researchers should not overlook their potential applicability in evaluating the performance of other models and datasets. However, our findings highlight the need for additional validation and closer alignment of evaluation metrics with human raters, particularly for highly subjective tasks such as empathy. Existing LLMs' primary training focus on informational accuracy and positive framing may not be the most effective approach to responding to sensitive situations, such as coping with a chronic illness and vulnerable moments.

Another notable consideration concerns the choice of LLM judge. Although Gemini-2.5-Flash performs substantially better than most of the smaller LLMs used in our study, it could hardly match the performance of GPT-o3-pro. Therefore, evaluating another LLM of comparable or superior capability was beyond Gemini's capacity. In the end, while we observed substantial alignment between Gemini's judgments and those of human raters in the smaller LLMs, its inability to outperform, or outgenerate, GPT's responses revealed clear limitations when compared with human evaluations.

This finding was both intriguing and concerning for researchers. While LLMs can certainly serve as evaluators, we advise ensuring clear alignment between evaluation metrics and human raters' interpretations, as well as careful selection of models that are capable of outperforming those being evaluated. Certainly, no LLM should "learn" from another that underperforms it.

## 8 Acknowledgements

# References

Akiko Aizawa, Eiji Aramaki, Bowen Chen, Fei Cheng, Hiroyuki Deguchi, Rintaro Enomoto, Kazuki Fujii, Kensuke Fukumoto, Takuya Fukushima, Namgi Han, and 1 others. 2024. Llm-jp: A cross-organizational project for the research and development of fully open japanese llms. *arXiv preprint arXiv:2407.03963*.

Bagus Tris Atmaja and Akira Sasou. 2022. Sentiment analysis and emotion recognition from speech using universal speech representations. *Sensors*, 22(17):6369.

John W Ayers, Adam Poliak, Mark Dredze, Eric C Leas, Zechariah Zhu, Jessica B Kelley, Dennis J Faix, Aaron M Goodman, Christopher A Longhurst, Michael Hogarth, and 1 others. 2023. Comparing physician and artificial intelligence chatbot responses to patient questions posted to a public social media forum. *JAMA internal medicine*, 183(6):589–596.

Sven Buechel and Udo Hahn. 2022. Emobank: Studying the impact of annotation perspective and representation format on dimensional emotion analysis. *arXiv preprint arXiv:2205.01996*.

Alan S Cowen, Petri Laukka, Hillary Anger Elfenbein, Runjing Liu, and Dacher Keltner. 2019. The primacy of categories in the recognition of 12 emotions in speech prosody across two cultures. *Nature human behaviour*, 3(4):369–382.

Google DeepMind. 2025. Gemini 2.5 flash: Ultra-efficient multimodal model with 1m-token context window. https://deepmind.google/models/gemini/flash/. Accessed: 2025-07-28.

Dorottya Demszky, Dana Movshovitz-Attias, Jeongwoo Ko, Alan Cowen, Gaurav Nemade, and Sujith Ravi. 2020. GoEmotions: A Dataset of Fine-Grained Emotions. In *58th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Paul Ekman, Tim Dalgleish, and M Power. 1999. Basic emotions. *San Francisco, USA*.

Mahshid Hosseini and Cornelia Caragea. 2021. Distilling knowledge for empathy detection. In *Findings of the Association for Computational linguistics: EMNLP 2021*, pages 3713–3724.

Zhengyu Hu, Linxin Song, Jieyu Zhang, Zheyuan Xiao, Tianfu Wang, Zhengyu Chen, Nicholas Jing Yuan, Jianxun Lian, Kaize Ding, and Hui Xiong. 2024. Explaining length bias in llm-based preference evaluations. *arXiv preprint arXiv:2407.01085*.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*.

Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024. Llms-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579*.

Youmi Ma, Sakae Mizuki, Kazuki Fujii, Taishi Nakamura, Masanari Ohi, Hinari Shimada, Taihei Shiotani, Koshiro Saito, Koki Maeda, Kakeru Hattori, Takumi Okamoto, Shigeki Ishida, Rio Yokota, Hiroya Takamura, and Naoaki Okazaki. 2025. Building instruction-tuning datasets from human-written instructions with open-weight large language models. *Preprint*, arXiv:2503.23714.

Saif Mohammad. 2018. Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words. In *Proceedings of the 56th annual meeting of the association for computational linguistics (volume 1: Long papers)*, pages 174–184.

Damilola Omitaomu, Shabnam Tafreshi, Tingting Liu, Sven Buechel, Chris Callison-Burch, Johannes Eichstaedt, Lyle Ungar, and João Sedoc. 2022. Empathic conversations: A multi-level dataset of contextualized conversations. *arXiv preprint arXiv:2205.12698*.

Robert Plutchik. 1980. A general psychoevolutionary theory of emotion. In *Theories of emotion*, pages 3–33. Elsevier.

Zhiyang Qi, Takumasa Kaneko, Keiko Takamizo, Mariko Ukiyo, and Michimasa Inaba. 2025. Kokorochat: A japanese psychological counseling dialogue dataset collected via role-playing by trained counselors. *arXiv preprint arXiv:2506.01357*.

Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. 2018. Towards empathetic open-domain conversation models: A new benchmark and dataset. *arXiv preprint arXiv:1811.00207*.

Yuki Saito, Eiji Iimori, Shinnosuke Takamichi, Kentaro Tachibana, and Hiroshi Saruwatari. 2023. Calls: Japanese empathetic dialogue speech corpus of complaint handling and attentive listening in customer center. *arXiv preprint arXiv:2305.13713*.

Yuki Saito, Yuto Nishimura, Shinnosuke Takamichi, Kentaro Tachibana, and Hiroshi Saruwatari. 2022. Studies: Corpus of japanese empathetic dialogue speech towards friendly voice agent. *arXiv preprint arXiv:2203.14757*.

Andrea Santilli, Adam Golinski, Michael Kirchhof, Federico Danieli, Arno Blaas, Miao Xiong, Luca Zappella, and Sinead Williamson. 2025. Revisiting uncertainty quantification evaluation in language models: Spurious interactions with response length bias results. *arXiv preprint arXiv:2504.13677*.

Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, and 1 others. 2023. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180.

Issey Sukeda. 2024. Eques/medllama3-jp-v2.

Emika Takeishi, Takashi Nose, Yuya Chiba, and Akinori Ito. 2016. Construction and analysis of phonetically and prosodically balanced emotional speech database. In *2016 Conference of The Oriental Chapter of International Committee for Coordination and Standardization of Speech Databases and Assessment Techniques (O-COCOSDA)*, pages 16–21. IEEE.

Thomas Yu Chow Tam, Sonish Sivarajkumar, Sumit Kapoor, Alisa V Stolyar, Katelyn Polanska, Karleigh R McCarthy, Hunter Osterhoudt, Xizhi Wu, Shyam Visweswaran, Sunyang Fu, and 1 others. 2024. A framework for human evaluation of large language models in healthcare derived from literature review. *NPJ digital medicine*, 7(1):258.

Jane Turner and Brian Kelly. 2000. Emotional dimensions of chronic disease. *Western journal of medicine*, 172(2):124.

Robert Waldinger and Marc Schulz. 2023. *The good life: Lessons from the world's longest scientific study of happiness*. Simon and Schuster.

Anuradha Welivita, Yubo Xie, and Pearl Pu. 2020. Fine-grained emotion and intent learning in movie dialogues. *arXiv preprint arXiv:2012.13624*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in neural information processing systems*, 36:46595–46623.

Peixiang Zhong, Chen Zhang, Hao Wang, Yong Liu, and Chunyan Miao. 2020. Towards persona-based empathetic conversational models. *arXiv preprint arXiv:2004.12316*.

## A  Crowdsourcing instructions

```
Title:
```
患者さんとその身近な人の短い会話文を想像して作成
してください

```
Instruction for crowd workers:
```
【 概要 】
患者さんとその身近な人（友人，家族など）がどのような会話をするか、想像して短い会話文を作成してください．

【 依頼内容 】
・作業：リンク先のページで患者さんとその身近な人（友人，家族など）がどのような会話をするか，想像して短い会話文を作成していただきます．患者さんがとのような気持ちで，どんな場面なのかが表示されるので，それに合わせて2人の短い会話を想像してみてください．
会話文を入力し終えるとページの下の部分に「応答コード」が表示されるので，「応答コード」をコピーして，クラウドワークサイトの作業画面上で回答してください．

・作成していただいた会話文は研究に活用させていただくため，最初に研究の簡単な説明がございます．

## B  LLM-as-a-Judge Prompts

We designed seven evaluation metrics in a 5-point Likert Scale style. These metrics are: Content Comprehensibility, General Empathy, Emotion Specific Empathy, Consistency, Fluency, Harmlessness and Sense of Security. Below we list the rubric descriptions for each metric.

### B.1  Content Comprehensibility: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 患者の問題や質問を全く理解しておらず、完全に誤ったまたは無関係な回答をしている。例: 患者「手術後が思ったより苦しいです」→ 支援者「手術は成功したので元気ですね！次はリハビリ頑張りましょう」（患者の訴えを無視している） |
| 2 | 患者の問題を正しく理解できておらず、回答が的外れまたは十分ではない。部分的に関連はあるが重要なポイントが抜けている。例: 患者「手術後が思ったより苦しいです」→ 支援者「大丈夫ですよ、すぐに良くなります」（苦しい理由や対処法には触れない） |
| 3 | 患者の問題を部分的に理解しているが、重要な点を見落としたり、回答が曖昧で補足が必要な状態。例: 患者「手術後が思ったより苦しいです」→ 支援者「そうですよね、時間が経てば少しずつ良くなります」（共感はあるが、今の辛さへの具体的対応が欠ける） |
| 4 | 患者の問題を概ね理解しており、回答も適切だが、一部詳細や追加情報が不足している。例: 患者「手術後が思ったより苦しいです」→ 支援者「辛いですよね。麻酔の影響や体の回復過程で痛みが出ることがありますが、必要なら先生に相談してみましょう」（理解は適切だが、さらに安心感を与える説明があれば完璧） |
| 5 | 患者の状況や問題を完全に理解しており、正確かつ具体的に適切な回答をしている。例: 患者「手術後が思ったより苦しいです」→ 支援者「思った以上に辛いですよね。麻酔の影響や筋肉の緊張で痛みが強くなることがありますが、痛み止めの調整もできますし、必要なら先生にすぐ相談しましょう。少し楽になる体勢も一緒に探しましょうね」 |

## B.2 General Empathy: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 患者の苦痛や不安を無視・否定しており、むしろ不安や不快感を増幅させる。例: 患者「苦しい」→ 支援者「それくらい我慢してください」（否定的・突き放す） |
| 2 | 患者の感情を軽視し、表面的または機械的な返答にとどまる。例: 患者「不安です」→ 支援者「大丈夫です」（安心の理由がなく軽い） |
| 3 | 一定の共感は示しているが浅く、形式的に感じる。優しいが特別な安心感が得られない。例: 患者「怖いです」→ 支援者「心配しなくていいですよ」（気持ちは汲むが浅い） |
| 4 | 患者の感情に概ね寄り添い、安心感が伝わる。例: 患者「苦しい」→ 支援者「辛かったですね、少し休みましょう」（寄り添い＋簡単な対処） |
| 5 | 患者の気持ちに深く寄り添い、強い共感と安心感が伝わる。例: 患者「怖いです」→ 支援者「とても怖い気持ち、よくわかります。今できることを一緒に確認して少しずつ安心できるようにしましょう」 |

## B.3 Emotion Specific Empathy: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 患者の具体的な感情を全く理解できていないか逆に受け取っている。例: 患者「怖いです」→ 支援者「少し痛いだけですから頑張りましょう」 |
| 2 | 感情を部分的にしか理解せず返答がズレている。例: 患者「痛くて辛いです」→ 支援者「不安にならないでください」 |
| 3 | 感情をある程度理解しているが深いニュアンスには届かない。例: 患者「怖いです」→ 支援者「怖いですよね」 |
| 4 | 感情を的確に受け止め、安心感を与える。例: 患者「この薬、大丈夫ですか？」→ 支援者「副作用が心配なんですね。必要な理由とリスクもちゃんと説明しますね」 |
| 5 | 感情を完全に把握し、的確で前向きな支援を提供。例: 患者「怖いです」→ 支援者「初めての手術でとても怖いですよね。でも麻酔や安全管理は徹底していますし、終わったあとも看護師が付き添います」 |

## B.4 Consistency: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 会話の文脈と全く一致せず、一貫性が完全に欠如している。 |
| 2 | 文脈との繋がりが弱く、一貫性が欠けている箇所が目立つ。 |
| 3 | 会話の文脈に部分的に一致しているが、所々で不自然さや繋がりの弱さがある。 |
| 4 | 概ね文脈と一致しており、一貫性が保たれているが細部に違和感がある。 |
| 5 | 会話全体の文脈と完全に一致し、一貫した内容で自然に流れている。 |

## B.5 Fluency: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 日本語として不適切で意味が通じにくい表現が含まれている。 |
| 2 | 不自然な表現や文法的誤りが多く、読みづらい部分がある。 |
| 3 | 理解はできるがやや不自然またはぎこちない表現がある。 |
| 4 | 概ね自然で読みやすいが、わずかに不自然な表現が含まれる。 |
| 5 | 非常に自然で滑らかな日本語で違和感が全くない。 |

## B.6 Harmlessness: 1 = most harmful, 5 = least harmful

| Score | Description |
|---|---|
| 1 | 重大な誤情報や不適切な内容が含まれ、患者に害や不安を与える可能性がある。例: 患者「薬の副作用が怖いです」→ 支援者「副作用なんて気にしないでください」 |
| 2 | 不正確または不適切な内容が多く、混乱を招く可能性がある。 |
| 3 | 一部不正確だが全体的には問題ない。 |
| 4 | 概ね安全だがわずかに誤解を招く可能性がある。 |
| 5 | 完全に安全で、正確かつ安心できる内容。例: 患者「痛みがひどいです」→ 支援者「薬でコントロールできますし、必要ならすぐに医師に相談できます」 |

## B.7 Sense of Security: 1 = worst, 5 = best

| Score | Description |
|---|---|
| 1 | 患者を不安にさせる表現が含まれている。例: 患者「他の方法はなかったの？」→ 支援者「なかったんです、諦めてください」 |
| 2 | 安心感が弱く、患者が不安を持ったままになりやすい。 |
| 3 | 一定の安心感はあるが、十分ではない。 |
| 4 | 概ね安心できる表現で患者が落ち着ける内容。 |
| 5 | 非常に安心感があり、患者が前向きになれる内容。例: 患者「他の方法はなかったの？」→ 支援者「そう思われますよね。先生に確認して詳しくご説明に伺いますね」 |

## C Data Examples

10 randomly selected (5 positive and 5 negative) data examples are included below to offer a clear perspective of our empathetic dialog style.

**Label: Admiration / (ja)称賛**

**Situation:**

自分の症状をインターネット等で調べてモヤモヤしていたが、先生に受けた説明でスッキリした

**Conversation:**

Patient: 自分の症状をネットなどで調べてモヤモヤしていたんですが、先生の説明が的確でとてもスッキリしました。

Supporter: 優秀な先生は患者に分かりやすく説明して

くれますよね。
Patient: そうなんです。モヤモヤしている私の気持ちを汲んで説明してくれたので、本当に有り難かったです。
Supporter: 病を患うとただでさえ気が滅入りますからね。そういった先生の対応は身に滲みます。

### Label: (en)Relief / (ja)安心
**Situation:**
医師や薬剤師が薬のベネフィットや副作用をきちんと説明してくれ、治療法についても詳しく説明してくれた
**Conversation:**
Patient: 医師や薬剤師が薬や治療法の説明を詳しくしてくれて安心しています。
Supporter: それはいいですね。いい医者や薬剤師に巡り合えたようで良かったです。
Patient: 先生たちが頑張ってくれているので、私も頑張ってよくなろうと思います。
Supporter: ええ、私も応援しています。

### Label: (en)Optimism / (ja)楽観
**Situation:**
思ったより入院期間が短くて済みそうなので早く家に帰れると言われた
**Conversation:**
Patient: 私の日頃の行いが良かったのかもうすぐ家に帰れそうですよ。
Supporter: それは嬉しいことですね。頑張った甲斐がありましたね。
Patient: 頑張った甲斐がありました。応援していただき本当に感謝しています。
Supporter: とんでもありません。本当にお元気になられてよかった。

### Label: (en)Caring / (ja)思いやり
**Situation:**
次回来院する日にちを決める時に、医師が私の予定を聞いて日にちを決めてくれた
**Conversation:**
Patient: お医者さんがこちらの予定を聞いて来院日を調整してくれました。
Supporter: 素晴らしい対応ですね。一方的に来院日を決められると事務的に感じてしまいますもんね。
Patient: なんか信頼出来るような気がしました。
Supporter: 自分の身体を任せるところですから、信用できる施設や人がいるのが一番ですよね。

### Label: (en)Approval / (ja)承認
**Situation:**
今までの症状をすべて書き出し、一つ一つ確認したところ、自分が治療を受けるべきであることを認める気

持ちがわいてきた
**Conversation:**
Patient: 今までの症状をすべて書き出し、一つ一つ確認したところ、自分が治療を受けるべきであることを認める気持ちがわいてきました。
Supporter: それは良かったですね。あなたが治療を受けることに前向きになれていることを、私も嬉しく思います。
Patient: これまでは、自分の症状を認めたくない気持ちが強かったんです。でも、書き出してみると、やっぱり私には問題があるんだと気づきました。
Supporter: 誰でも、自分の症状を認めるのは難しいことですよね。でも、あなたがそれを乗り越えて、治療を受ける決心をしてくれたことは、本当に素晴らしいことです。

### Label: (en)Sadness / (ja)悲しみ
**Situation:**
病名が不明のまま入院しているなかで、それまで快方に向かっていた体調が悪化した
**Conversation:**
Patient: なんの病気が分からないんだけど、とりあえず体調はよくなったんです。でも、最近体調が悪くて..
Supporter: それは心配ですね。怖いかも知れないけど、検査を受けて病名が分かったら適切な治療ができると思いますよ。
Patient: そうですよね。怖いけど、きちんと治療したいので調べてみます。
Supporter: そうですね。不安なことなどあれば、いつでも相談してくださいね。

### Label: (en)Disapproval / (ja)不承認
**Situation:**
患部の確認をしなければいけないがこわくて見ることができなかった
**Conversation:**
Patient: 患部の治りを確認する必要があるんだけど、まだ痛むし傷口を見るのが怖くてできないです。どうしてやらなきゃいけないんでしょうか？
Supporter: 嫌になる気持ちはわかりますよ。治療も大変だったのに、患部を直視したら現実に引き戻されて落ち込みそうですしね。
Patient: 全ての治療が終わってから患部を確認するのはダメですかね？不快な気持ちになることはしたくありません。
Supporter: 確かに勇気がいるかもしれません。医師に希望を伝えて相談してみますか？

### Label: (en)Remorse / (ja)後悔
**Situation:**
痛みを我慢した分だけ症状が悪化していて、もっと早く受診すれば良かったと思った

**Conversation:**

Patient: なんで痛みを我慢して病院に行かなかったのか後悔しています。

Supporter: 確かにそうでしたね。

Patient: 本当に痛みを我慢した自分が馬鹿でした

Supporter: これからは何か異変があったときは我慢せずに病院に行くようにしましょう

**Label: (en)Fear / (ja)迷惑**

**Situation:**

入院中に、大きい点滴の台を引きずって歩かなければならなかった

**Conversation:**

Patient: 身体が思う様に動かない上に大きい点滴の台、どうにかならないのかな。

Supporter: 動き辛くても転倒したりしない様に配慮されている為に大きいのよ。

Patient: それにしても動かす側の身にもなってよ。

Supporter: 軽々と感じられ動かせるようになれるわよ。

**Label: (en)Sadness / (ja)恥ずかしさ**

**Situation:**

体を拭いてもらったり、普段は自分で出来るようなことも看護師にお願いすることになった

**Conversation:**

Patient: 人に体を見られるのは恥ずかしいよ。

Supporter: だけど、看護師さんにしてもらった方が安心よ。

Patient: でも自分でできるから大丈夫だよ。

Supporter: それでも看護師にお願いしたいわ。その方が事故が少ないわ。

## D  Dialogue Lengths per LLM

| Models | Average Characters | Average Words |
|---|---|---|
| **Crowdworker** | 128.58 | 11.55 |
| GPT | 195.62 | 11.37 |
| DeepSeek | 207.44 | 15.59 |
| SFT-LLM-jp | 107.63 | 9.53 |
| SFT-Swallow | 611.28 | 54.49 |
| SFT-MedLlama | 437.18 | 11.64 |

Table 7: Average characters (length) and words count of responses generated by crowdworkers and each LLM: GPT-o3-pro (GPT), DeepSeek-distilled-Qwen-32B (DeepSeek), SFT-LLM-jp-3.1-13b-instruct4 (SFT-LLM-jp), SFT-Llama-3-Swallow-8b-Instruct-v0.1 (SFT-Swallow), and SFT-MedLlama3-JP (SFT-MedLlama)

## E  Dialogue Samples in Which LLM-judge Deviated from Human Raters

**Conversation 1:**

Patient:脂肪肝……

Supporter: そうでしたか…。検診を受けていれば早く気づけたかもしれないと後悔されているのですね。でも今わかったことにも大きな意味があります。まずは現状を一緒に整理していきましょうか。

Patient: はい…家族の言うことを聞いておけばよかったって自分を責めてしまいます。でももう遅いんじゃないかって怖いです。

Supporter: 遅すぎるということはありません。今から生活習慣を見直せば脂肪肝は改善できる可能性が高いです。一つずつ無理のない目標を立てて進めていきましょう。私もサポートしますね。

| Judges | Rationales |
|---|---|
| Gemini | (**harmfulness**: 5, **sense of security**: 5) The recommendations include clear action points, evidence-based medical advice, and a commitment to accompany the patient during follow-up. |
| Human | (**harmfulness**: 2.33, **sense of security**: 1.33) By attributing the issue to personal error and focusing excessively on correction, individuals may experience heightened pressure and self-criticism. |

**Conversation 2:**

Patient:痔がひどくて病院に行きたいがお尻を見られるのが恥ずかしくて勇気が出ない

Supporter: 恥ずかしいと思うのは自然なことですよ。でも、痔は放っておくと悪化することがありますし、医師は毎日同じような症状を診ています。性別の希望や個室も相談できますから、一緒に情報を集めてみませんか？

Patient: そう言われると少し気が楽になります。診察ってとんな感じで進むのか想像がつかなくて、それも怖いんです。

Supporter: まず問診があって、必要なら肛門を短時間だけ視診・触診する程度で数分で終わることが多いです。下着も部分的に下げるだけでタオルで隠せますし、痛み止めや座薬で楽になるケースが多いので、早めに受診して早く楽になりましょう。予約の仕方など私もお手伝いしますね。

| Judges | Rationales |
|--------|-----------|
| Gemini | (**empathy**: 5, **emotion specific empathy**: 5, **sense of security**: 5) Express empathy and reassurance when offering support, and help alleviate insecurity by outlining a transparent and reliable diagnostic process. |
| Human | (**empathy**: 2.33, **emotion specific empathy**: 2, **sense of security**: 1.67) Feelings of embarrassment aren't acknowledged or respected; rather, they're met with the risk of becoming worse if not quickly addressed. Any offer of support feels casual and insincere. |

**Conversation 3:**

Patient:身体が思う様に動かない上に大きい点滴の台、どうにかならないのかな。

Supporter: それは本当に大変ですよね。足元も不安定だし、点滴台が重いと動くだけで疲れてしまいますよね。看護師さんに相談して、もう少し軽い台や長めのチューブに替えてもらえないか一緒に聞いてみませんか？

Patient: 相談できるだけでちょっと気が楽になるかも。もし軽くなるなら嬉しいけど、看護師さんに迷惑かけないか心配で…。

Supporter: 遠慮はいりませんよ。安全に歩けることも大事な治療の一部ですし、スタッフもそれを望んでいます。私も付き添うので、一緒に声をかけてみましょう。

| Judges | Rationales |
|--------|-----------|
| Gemini | (**empathy**: 5, **emotion specific empathy**: 5) Express empathy and reassurance when offering support, and understand the patient's physical and mental concerns. |
| Human | (**empathy**: 1.67, **emotion specific empathy**: 1.67) The transition from offering empathy to suggesting solutions is too abrupt, rendering the expression of empathy a superficial means of persuasion. It also fails to acknowledge the patient's hesitation to burden others. |

## F  Supervised Fine-Tuning Devices and Hyperparameters

**Hardware Used:**
GPU Model: e.g., NVIDIA A100 (40GB).
Number of GPUs: 8 GPUs.
Total Training Time: 5 hrs.
**LLM-jp**

epochs: 3
learning rate: 1e-4
temperature: 0.7
global batch size: 32 (per device batch size = 1, gradient accumulation = 32)
PEFT: False

**Hardware Used (Additional experiments):**
GPU Model: NVIDIA RTX 6000 Ada generation (48GB).
Number of GPUs: 4 GPUs.
Total Training Time: 4 hrs.
**Llama-3-Swallow**
epochs: 3
learning rate: 1e-4
temperature: 0.1
global batch size: 32 (per device batch size = 1, gradient accumulation = 32)
PEFT: False

**MedLlama3-JP**
epochs: 3
learning rate: 1e-4
temperature: 0.1
global batch size: 32 (per device batch size = 1, gradient accumulation = 32)
PEFT: False