

Empathy Applicability Modeling for General Health Queries

Anonymous ACL submission

Abstract

LLMs are increasingly being integrated into clinical workflows, yet they often lack clinical empathy, an essential aspect of effective doctor–patient communication. Existing NLP frameworks focus on reactively labeling empathy in doctors’ responses but offer limited support for anticipatory modeling of empathy needs, especially in general health queries. We introduce the Empathy Applicability Framework (EAF), a theory-driven approach that classifies patient queries in terms of the applicability of emotional reactions and interpretations, based on clinical, contextual, and linguistic cues. We release a benchmark of real patient queries, dual-annotated by Humans and GPT-4o. In the subset with human consensus, we also observe substantial human–GPT alignment. To validate EAF, we train classifiers on human-labeled and GPT-only annotations to predict empathy applicability, achieving strong performance and outperforming the heuristic and zero-shot LLM baselines. Error analysis highlights persistent challenges: implicit distress, clinical-severity ambiguity, and contextual hardship, underscoring the need for multi-annotator modeling, clinician-in-the-loop calibration, and culturally diverse annotation. EAF provides a framework for identifying empathy needs *before* response generation, establishes a benchmark for anticipatory empathy modeling, and enables supporting empathetic communication in asynchronous healthcare.

1 Introduction

Clinical empathy integrates cognitive (understanding), emotional (resonating), and action-oriented (expressing) components (Guidi and Traversa, 2021). It is indispensable for clinical care, deepening therapeutic relationships and improving outcomes such as patient satisfaction, care effectiveness, reduced distress, and hospital length of stay (Guidi and Traversa, 2021; Olson, 1995; Hoffstädt et al., 2020); yet clinicians miss 70-90% of

empathic opportunities during patient interactions (Morse et al., 2008; Hsu et al., 2012).

Large Language Models (LLMs) are increasingly integrated into healthcare workflows and patient interactions, with major EHR vendors such as EPIC adopting them for clinical messaging and nearly half of physicians reporting patients consult ChatGPT before visits (Antoniak et al., 2024; Sermo Team, 2025). While these trends highlight rapid adoption of LLMs in healthcare, they also raise concerns of lacking empathy crucial for asynchronous physician-patient interactions (Koranteng et al., 2023). However, effective empathy requires discernment, not just fluency. This highlights a critical, antecedent challenge: How can we systematically model the applicability of empathy, allowing systems to recognize the specific clinical and linguistic cues that necessitate an emotional response?

Modeling empathy in text is inherently difficult without non-verbal cues, and NLP research has historically over-weighted emotional aspects while overlooking cognitive empathy (Lahnala et al., 2022), even though both are vital in clinical care. To redress this imbalance, EPITOME (Sharma et al., 2020) captures the multidimensionality of empathy through emotional reactions, interpretations, and explorations, offering an empathetic lens on warmth, understanding, and curiosity in mental health support. Online Empathy (Chai et al., 2019) also addresses multidimensionality, classifying responses as Informational and Emotional. However, both EPITOME and Online Empathy assess empathy post hoc, labeling support-giver responses after they appear and thus offering no guidance while a clinician is composing a response to the patient query.

Lahnala et al. attempt to solve this particular problem with with an Appraisal Framework that annotates empathic opportunities and clinician elicitation and response as functions of (affect | judgment

l appreciation) in breaking-bad-news oncology dialogues (Lahnala et al., 2024). This discourse analysis lens excels at teaching stance shifts over multi-turn synchronous conversations, yet is not suited to single-turn, asynchronous general health queries: it classifies stance, not what the patient needs (cognitive clarification vs emotional warmth). Thus, it remains need-blind, providing little actionable help for one-off replies.

To address this gap, we propose the Empathy Applicability Framework (EAF), a theoretically grounded method to proactively identify when and what type of clinical empathy should be expressed in response to patient queries. EAF operationalizes empathy along two key dimensions: affective (emotional reactions) and cognitive (interpretations) – labeling each as *Applicable* or *Not Applicable* based on clinical, contextual, and linguistic cues within patient queries. Unlike prior work that evaluates empathy *in the response itself* (Chai et al., 2019; Sharma et al., 2020), EAF analyzes the patient’s query before any response, enabling anticipatory reasoning. Evidence for the value of anticipation comes from Sibyl (Wang et al., 2025), which shows that anticipating user’s emotional and contextual trajectory improves empathetic response generation. EAF can also mitigate limitations of reactive frameworks that overrely on surface lexical cues. For instance, EPITOME misclassified generic phrases as empathetic in nonsensical contexts, with false-positive rates exceeding 95% (Lee et al., 2023). By shifting empathy assessment from *post-hoc* response scoring to *pre-response* query analysis, EAF is methodologically significant: it enables anticipatory inference of emotional and interpretive support needs, helping providers and LLMs detect and act on empathic opportunities in general health queries.

We make three primary contributions: (i) **Framework Design**: we introduce and theoretically ground the EAF in clinical empathy literature, clearly differentiating our anticipatory model from prior post-hoc approaches; (ii) **Annotated and analyzed Benchmark**: a novel dataset of 1,300 patient queries annotated by humans and GPT-4o (included in the supplementary materials), demonstrating EAF’s reliability and interpretability; and (iii) **Operationalization Challenges**: we identify and systematically analyze specific contexts where anticipatory empathy annotations diverge, highlighting opportunities for future research in multi-annotator modeling, clinician-in-the-loop systems,

and culturally sensitive annotation strategies.

2 Empathy Applicability Framework and Theoretical Grounding

The EAF identifies empathetic needs proactively by assessing patient queries along two dimensions adapted from EPITOME (Sharma et al., 2020) and informed by Chai et al.’s distinction between emotional and informational support (Chai et al., 2019): *Emotional Reactions* and *Interpretations*. Table 1 summarizes the EAF, detailing applicable and non-applicable cues for each dimension.

To develop EAF, we performed inductive thematic coding on 300 randomly selected patient queries from the HealthcareMagic and iCliniq datasets (Li et al., 2023), identified themes, formed subcategories (cues), and finally, we iteratively refined to comprehensively and distinctly capture empathy applicability.

Additionally, we ground EAF cues in Patient-Centred Care (PCC), a widely used framework for clinician–patient communication, focusing on the PCC functions of responding to emotions and managing uncertainty to ensure clinically grounded expressions of empathy (Epstein and Street Jr, 2007; McCormack et al., 2011).

3 Methods

To determine whether EAF is reliably interpretable across a range of clinical queries and to identify any systematic challenges, we curated a diverse dataset of health-related queries and annotated them using the EAF, employing both human annotators and an LLM. To assess whether these annotations exhibit learnable patterns, indicating the internal consistency of EAF, we trained classifiers on the EAF-labeled data. The following subsections detail the annotation and modeling procedures.

3.1 Data Source

We sampled 9,500 patient queries from two publicly available datasets (HealthCareMagic and iCliniq) released by Li et al. (Li et al., 2023). We sampled 4,750 queries each from HealthCareMagic ($\approx 100k$ dialogues) and iCliniq ($\approx 10k$), to maximize linguistic and contextual diversity and avoid overfitting to a single source. As these datasets are publicly available and anonymized, our IRB determined that this study was exempt from human subjects review. The datasets do not carry an explicit license; therefore, we use them exclusively

Dimensions	Applicable cues	Not Applicable cues
Emotional Reactions Expressions of warmth, compassion, concern, or similar feelings conveyed by a doctor in response to a patient’s query.	<ul style="list-style-type: none"> • Severe Negative Emotion • Inferred Negative State • Seriousness of Symptoms • Concern for Relations <i>Rationale:</i> Signals reflect distinct pathways of emotional distress, guiding when emotional reactions are warranted.	<ul style="list-style-type: none"> • Routine Health Management • Purely Factual Medical Queries • Neutral Symptom Descriptions • Hypothetical Queries <i>Rationale:</i> Signals no emotional content; omit reactions to maintain factual medical focus.
Interpretations Communication of an understanding of the patient’s feelings (expressed or implied) and/or experiences (including contextual factors) inferred from the patient’s query.	<ul style="list-style-type: none"> • Expression of Feeling • Experiences or Context Affecting Emotional State • Symptoms with an Emotional Impact • Distressing Uncertainty About Health <i>Rationale:</i> Signals lived burden, context, or uncertainty requiring interpretive acknowledgment.	<ul style="list-style-type: none"> • Emotional-Reactions N/A cues +: with absence of distressing contextual or experiential details. <i>Rationale:</i> Signals absence of both emotional and contextual cues, preventing over-empathizing and maintaining focus on informational needs.

Table 1: Empathy Applicability Framework (EAF). Each dimension lists cues for when an empathic dimension is *Applicable* or *Not Applicable*; Brief rationales explaining what each cue set captures follow the cues. Detailed description of the EAF and its cues with examples is provided in the Appendix A. Also, see Appendix Table 4 for concrete query scenarios illustrating cues usage and EAF operationalization.

for non-commercial research, in line with the authors’ stated intention and public availability, and will release our de-identified EAF benchmark under the same non-commercial terms. To balance rigor and cost, 1,500 of the queries were earmarked for dual annotation by humans and GPT-4o to support reliability and error analyses, while the remaining 8,000 were annotated only by GPT-4o for predictive validity testing.

3.2 Annotation Task

The annotation task required using EAF to label patient queries as applicable or not applicable (see Table 1) on two dimensions of empathy: Emotional Reactions (EA) and Interpretations (IA). Human annotators were instructed to identify at least one best-fitting subcategory per dimension to justify their labels (they mostly selected a single best-fitting subcategory). The GPT annotations listed all relevant subcategories supporting labeling decisions.

3.2.1 Annotator Recruitment, Training and Calibration

Due to empathy annotation subjectivity, we prioritized consistency by avoiding crowdsourcing and instead recruited and trained two annotators from Pakistan with high English proficiency: HA1, a female with an MS in Linguistics, and HA2, a male with a BS in Computer Science. We recruited two annotators via departmental channels for about a

one-month engagement. Informed consent to use the annotated dataset to train large language models was collected from the annotators prior to the start of the annotation process. Annotators were compensated US\$360 (equivalent to a local monthly research salary). Annotators underwent three-stage training on 200 queries (50 + 50 + 100) from a subset of 1,500, with convergence meetings after each stage to clarify misunderstandings and align labeling. Training queries were excluded from later experiments. Annotators then independently labeled the remaining 1,300 queries following procedures in Section 3.2. Annotation instructions are detailed in Appendix B.

We *intentionally* employed lay annotators to capture the patient’s perspective. Prior research shows that empathy is ‘in the eye of the beholder’ (Bernardo et al., 2018), and given that the empathy levels in the clinician’s response will be perceived by the patient, lay annotators whose judgments reflect the patient/recipient experience are *better-suited* for this task. Additionally, prior studies show that clinicians often overlook empathic opportunities in favor of diagnostic focus (Hsu et al., 2012) and that patients’ ratings of clinicians’ empathy often diverge from clinicians’ assessments (Bernardo et al., 2018; Hermans et al., 2018).

3.2.2 GPT Annotations

To scale the data set and enable comparison with human annotations, we used GPT-4o via the OpenAI API, prompted to act as an expert annotator using contrastive prompting (Gao and Das, 2024). The model was given definitions of EA and IA, subcategory descriptions with examples, and labels indicating whether each subcategory was Applicable or Not Applicable. Then it returned the matching subcategories, with the format inherently indicating the applicability class (annotation scripts included in the supplementary software). Complete prompt specifications, including with and without our framework are included in Appendix F.

For the 1,300 human-annotated queries, GPT-4o generated five annotation passes per query, with final labels determined by majority vote¹. For the remaining 8,000 queries, a single-pass annotation was used due to cost constraints. This yielded two subsets: 1,300 queries labeled by both humans and GPT (with majority-voted GPT labels) and 8,000 labeled solely by GPT (single-pass annotation). **Note:** Throughout the remainder of this text, all references to GPT refer specifically to GPT-4o.

3.3 Modeling Task and Approach

We frame empathy applicability prediction as two independent binary classification tasks. Given a patient query P_i , the objective is to predict, for each empathy dimension $d \in \{EA, IA\}$, whether that dimension is *Applicable* (1) or *Not Applicable* (0), denoted A_{id} . For each dimension, we fine-tune a distinct RoBERTa-based classifier (Liu et al., 2019). Full architectural details, including the attention mechanism, the pooling operation, and the model diagram, are provided in the appendix E.

4 Evaluation Setup and Experiments

This section details the evaluation setup and model training configurations used in our experiments.

4.1 Annotator Agreement

We assessed human annotation reliability using raw agreement and Cohen’s Kappa across the 1,300 independently labeled queries. For GPT-generated annotations, we compared majority-voted GPT labels with a subset of human-annotated queries:

¹Majority voting ensured consistency across passes. More than 94% of queries received the same label on the first pass and as the majority vote for both empathy dimensions, indicating minimal divergence. Hence, we report evaluation metrics only with the majority-voted labels.

queries where both human annotators reached an agreement. This allows us to evaluate GPT performance without confounding disagreement over error or subjectivity.

4.2 Conceptual Alignment

To examine whether humans and GPT rely on similar rationales, we performed an UpSet plot analysis (Figure 1). This analysis was limited to queries where humans and GPT agreed on the overall applicability label, allowing us to assess alignment in subcategory reasoning rather than outcome. A match is coded as *Full* if GPT includes both subcategories selected by the two human annotators, *Partial* if GPT’s subcategories overlap with only one human’s subcategory label and *No match* if GPT matches neither human subcategory.

4.3 Divergence Bar and Qualitative Analysis

Given the subjective nature of empathy, we analyze mismatches as directional divergences rather than strict errors. To characterize disagreement, we use three-way divergence bars (Figure 2) that decompose label mismatches within each subcategory into *Annotator Spread* (one human labeled Applicable, the other Not), *LLM-Adds* (GPT labeled Applicable, humans Not), and *LLM-Omits* (GPT labeled Not, humans Applicable). Furthermore, we performed qualitative analysis on a subset of queries where GPT labeled differently, and identified thematic patterns that highlight the different labeling.

4.4 Model Evaluation

We evaluated the performance of the classifiers trained to predict empathy applicability (Applicable vs. Not Applicable). Reported metrics include accuracy, weighted F1 score, and macro-averaged F1 score across both dimensions (EA and IA). To contextualize classifier performance, we compared results against four baselines: Random Guessing (assigns labels at random), Always Applicable, Always Not Applicable, and o1-Zero-Shot (based on OpenAI’s reasoning model, without invoking empathy applicability framework). For the o1 baseline, we provide only the definition of the target dimension (EA or IA) and prompt it to classify each patient query as ‘Applicable’ or ‘Not Applicable’, preserving the zero-shot setting without framework cues. These baselines help determine whether our trained models learn meaningful patterns beyond simple heuristics or zero-shot LLM reasoning. In

Dimension	Human–Human	Human–GPT
	κ (agree / disagree)	κ (agree/ disagree)
EA	0.521 (981 / 315)	0.617 (668 / 152)
IA	0.404 (898 / 398)	0.652 (678 / 142)

Table 2: Cohen’s κ with agreement counts: human–human agreement on the full set and human–GPT alignment on the human-consensus subset.

addition, following best practices for sanity checking and overfitting control, we include *classical* text classifiers trained and evaluated on the Human Set (Section 4.5): Logistic Regression (LR) and Linear SVM with TF–IDF features(1–3 grams). These linear models provide a transparent reference for what is learnable from local lexical features.

4.5 Model Training and Training Sets

Each classifier for the EA and IA tasks is based on RoBERTa-base (≈ 125 M parameters) and was trained on two distinct datasets (data and scripts included in the supplementary material): **Human Set**: Contains only queries where both human annotators reach consensus on a label for a given dimension, serving as a high-fidelity benchmark aligned with human judgment. **Autonomous Set**: Consists of GPT-labeled data from the 8,000-query pool, with no human supervision. This tests whether models trained solely on GPT output can approximate human consensus.

For the Human Set, we split the data into subsets of training (75%), validation (5%), and test (20%). For the Autonomous Set, training was done entirely on GPT-labeled data, but testing used the same human-consensus test set as the Human Set to enable consistent evaluation relative to human agreement. Training used a single NVIDIA A40 GPU per run. A Human-Set run finished in ≈ 15 min GPU time, while an Autonomous-Set run took ≈ 40 min; thus the total compute budget per dimension is < 1 GPU-hour. All models were trained for 10 epochs using a learning rate of 2×10^{-5} and a batch size of 8. To ensure comparability, all models shared the same architecture and hyperparameters.

5 Results

In this section, we present our findings related to the reliability of the EAF and the challenges in operationalizing it. Additional dataset characterization beyond agreement and modeling results is provided in Appendix G.

5.1 Reliability of the EAF

We evaluated reliability along three axes of Consistency, Predictive validity and Conceptual alignment.

Consistency. We first assess agreement on Applicable/Not Applicable labeling between human annotators across 1,300 queries, and between GPT-4o and the *human consensus* on a subset of 820² queries. As shown in Table 2, human annotators achieved moderate agreement on both empathy dimensions, with an overall Cohen’s κ of 0.46. This falls within the typical range for empathy annotation tasks; for example, Sibyl (Wang et al., 2025) reported scores between 0.4 and 0.6. Notably, agreements outnumbered disagreements *by a factor of two to three*, suggesting that the EAF supports relatively consistent human labeling despite the inherent subjectivity of empathy.

GPT aligned well with the *human consensus dataset*, queries where both humans agreed, achieving three-way agreement. For both EA and IA, Cohen’s κ exceeded 0.6 and raw agreement was about 80% (Table 2). These results reflect agreement on human-aligned cases, demonstrating EAF’s effectiveness in guiding GPT to anticipate empathy applicability in clearer contexts, excluding more ambiguous or complex queries (see section 5.2).

Predictive Validity. We next evaluated whether EAF annotations are machine-learnable. As shown in Table 3, classifiers trained on human consensus data achieved high performance: LR attains **0.84** Macro-F1 for EA and **0.80** for IA (SVM: 0.83/0.77), establishing a strong classical reference. Our transformer (RoBERTa-base) exceeds these values significantly (**for EA: (0.92 vs. 0.84); and for IA (0.87 vs. 0.80)**). Models trained on GPT-only annotations (the Autonomous set) also performed well, achieving around **0.85** for EA and **0.77** for IA on the same held-out human-consensus test set, reflecting expected loss from noisier labels or differences from human labeling. Our models also significantly outperformed the trivial baselines (random guessing, always applicable, always not applicable, and o1 Zero-Shot), which yielded substantially lower scores. McNemar’s test (McNemar, 1947) confirmed statistical significance of the transformers over the trivial baselines (max $p < 10^{-4}$)

²For the full set, Appendix Table 5 shows comparable GPT agreement with each human annotator on affective EA, but substantially more variable agreement on cognitive IA.

Table 3: Classification results across training sets and baselines (single run on the human-consensus test set). Bold indicates best performance. Classical baselines (TF-IDF+LR/SVM) are trained on the Human Set only. Our transformer significantly outperforms all the baselines.

Training Set / Model	EA			IA		
	Acc	Macro-F1	Wtd-F1	Acc	Macro-F1	Wtd-F1
Random	0.47	0.47	0.47	0.44	0.43	0.44
Always Applicable	0.52	0.34	0.36	0.53	0.35	0.37
Always Not Applicable	0.48	0.32	0.31	0.47	0.32	0.30
o1 Zero-Shot	0.55	0.40	0.41	0.62	0.53	0.54
<i>Human-supervised models (train and tested on human-consensus set)</i>						
Logistic Regression	0.84	0.84	0.84	0.80	0.80	0.80
Linear SVM	0.83	0.83	0.83	0.77	0.77	0.77
Transformer (RoBERTa-base)	0.92	0.92	0.92	0.87	0.87	0.87
<i>Autonomous-supervised model (train on GPT labels, test on human-consensus test set)</i>						
Transformer (RoBERTa-base)	0.85	0.85	0.85	0.78	0.77	0.77

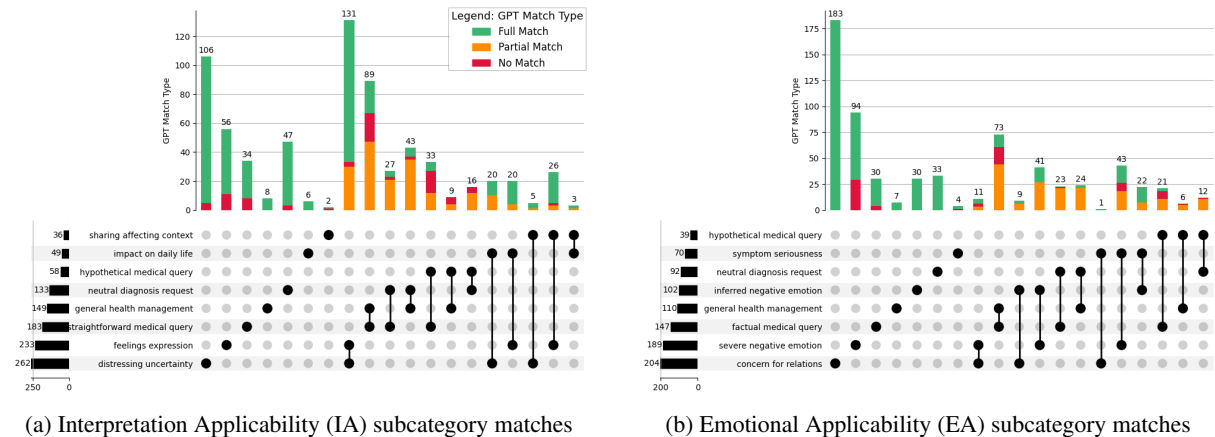


Figure 1: UpSet plots comparing GPT and human rationales for (a) IA and (b) EA subcategories. For each query, each human annotator selects one best-fit subcategory for their rationale; thus the human set is either a *single-dot* combination (both humans chose the same subcategory) or a *two-dot* combination (humans chose different subcategories). Horizontal bars show how often each subcategory appears in the human annotation set across all queries. Each vertical bar shows the frequency of a unique human-combination and is split by GPT agreement: **Full** (GPT’s subcategory set covers the entire human set), **Partial** (GPT matches only one of the two human subcategories), and **No match** (GPT matches neither human subcategory).

and over the classical baselines (max $p \leq 0.02$). Taken together, strong linear performance indicates consistent linguistic realizations of the constructs, while the transformer’s margin suggests benefits from broader context rather than overfitting. Overall, these results show that EAF-labeled data encode structured and learnable patterns.

Conceptual Alignment. We further examined whether humans and GPT rely on similar reasoning when assigning EAF labels. UpSet plot analysis (Figure 1) shows strong conceptual alignment. In many cases, both human annotators independently selected the same subcategory and GPT matched it, especially for both applicability and non-applicability cues such as *Severe Emotion* or *Factual Queries*. These matches indicate that the EAF defines meaningful categories that are consistently identifiable by both humans and LLMs.

tently identifiable by both humans and LLMs.

When annotators selected different subcategories for the same label, GPT often matched both. For example, in queries involving both *Expression of Feeling* and *Distressing Uncertainty*, GPT cited both reasons, suggesting that GPT can reconcile diverse human rationales and underscores the framework’s breadth in conceptualizing clinical empathy. No-match cases are rare, and GPT typically overlaps with at least one human subcategory. Appendix G.4 quantifies and corroborates these trends, showing match rates (match vs. miss, conditioned on humans using that subcategory) above 80% for most subcategories.

Collectively, these results establish that EAF supports consistent human judgments, yields learnable patterns, and promotes interpretable reasoning

453 across both humans and LLMs, making it well
454 suited for anticipatory empathy modeling in clinical
455 settings.

456 5.2 Systematic Challenges in Operationalizing 457 Anticipatory Empathy

458 Divergence bar analysis (Section 4) revealed that
459 inter-human agreement is significantly lower for
460 interpretations (IA) than for Emotional Reactions
461 (EA) (Table 2), and that despite moderate overall
462 human-GPT agreement (Table 2), there is diver-
463 gence at the subcategory level. Subsequent qual-
464 itative analysis revealed three key challenges in
465 applying the EAF, with implications for any clinical
466 empathy framework in NLP.

467 5.2.1 Challenge 1: Subjectivity in Identifying 468 Implied Distress

469 The categories *Inferred Negative State* (EA) and
470 *Distressing Uncertainty* (IA) show substantial di-
471 vergence in inter-human and human-GPT annota-
472 tions (Figure 2).

473 A qualitative review of 50 randomly selected
474 cases³ (25 each for *Distressing Uncertainty* and *In-*
475 *ferred Negative State*)⁴ by the first author acting as
476 adjudicator revealed that in more than 50% of the
477 queries, one could reasonably infer implied emo-
478 tional distress *or* determine that the query is driven
479 by factual intent. For instance, the female anno-
480 tator labeled a pain-and-menstrual-cycle query as
481 *Distressing Uncertainty*, while the male annotator
482 treated it as a factual diagnostic request, illustrating
483 the subjectivity of distress inference.

484 5.2.2 Challenge 2: Clinical-Severity 485 Ambiguity

486 In the category *Serious Symptoms* (EA), GPT la-
487 beled 100 queries as requiring emotional reactions
488 when humans did not (Figure 2). Qualitative anal-
489 ysis of 25 randomly selected cases³ where only
490 GPT had labeled empathy applicability revealed
491 three patterns: (1) In 40% of the cases, GPT appro-
492 priately flagged empathy needed for patients with
493 chronic or life-threatening conditions (e.g., post-
494 liver transplant complications) that human annota-
495 tors with no medical background had overlooked
496 (2) borderline cases with reasonable disagreement

³Detailed patient queries, mis-aligned labels, and qualita-
tive interpretations are included in the supplementary material
as the misalignment_analysis.csv file

⁴a sample size consistent with prior clinical-NLP error
analyses; (Hu et al., 2024)

(16%), such as prolonged low-grade fever after kid-
ney stones, and (3) GPT overgeneralization of vivid
but non-serious pain symptoms (44%) that did not
meet the EAF criteria of chronic or life-threatening
severity (for example, lip numbness after dental
problems).

503 5.2.3 Challenge 3: Contextual Hardship

504 GPT frequently over-applied *Symptoms Emotional*
505 *Impact* (SEI) and *Context Sharing* (CS) tags com-
506 pared to humans (Figure 2). An analysis of 25 ran-
507 domly selected³ mismatched labels in SEI category,
508 and all 20 mismatches in CS revealed that while
509 GPT sometimes correctly identified complex dis-
510 tress signals humans missed (20-25% of the cases),
511 it more often (75-80% of the cases) equated physi-
512 cal discomfort with emotional distress – potentially
513 reflecting Western-centric training biases (Johnson
514 et al., 2022; Cao et al., 2023).

515 These challenges, rooted in subjective inference,
516 clinical ambiguity, and cultural variation, highlight
517 the complexity of implementing clinical empathy.
518 Addressing them requires moving beyond single-
519 annotator consensus toward frameworks that em-
520 brace interpretive pluralism, clinical expertise, and
521 cultural sensitivity.

522 6 Discussion and Conclusion

523 Asynchronous patient communication requires *an-*
524 *ticipatory* mechanisms that can signal empathic
525 needs *before* a response is written. EAF addresses
526 this gap by assigning applicability labels to patient
527 queries, indicating whether empathy is warranted
528 and which dimension (emotional versus interpre-
529 tive) should be expressed. This framing comple-
530 ments recent advances in empathetic response gen-
531 eration, which enrich responses by reasoning about
532 emotional causes, rhetorical framing, and likely fu-
533 ture trajectories. For example, Lee et al. integrate
534 figurative language and empathy-cause signals to
535 generate responses that are both emotionally en-
536 gaging and contextually grounded. Complementar-
537 ily, Chen et al. use cause-aware chain-of-thought
538 prompting: they extract the emotional cause from
539 the dialogue and insert generated commonsense
540 relations (xEffect, xReact, xIntent, xNeed, xWant)
541 conditioned on the extracted cause into the prompt,
542 which significantly improves empathy in human
543 evaluations (Chen et al., 2024). Likewise, Sibyl
544 (Wang et al., 2025) helps models anticipate the
545 seeker’s likely next move and improves empa-
546 thetic generation. However, these systems are de-

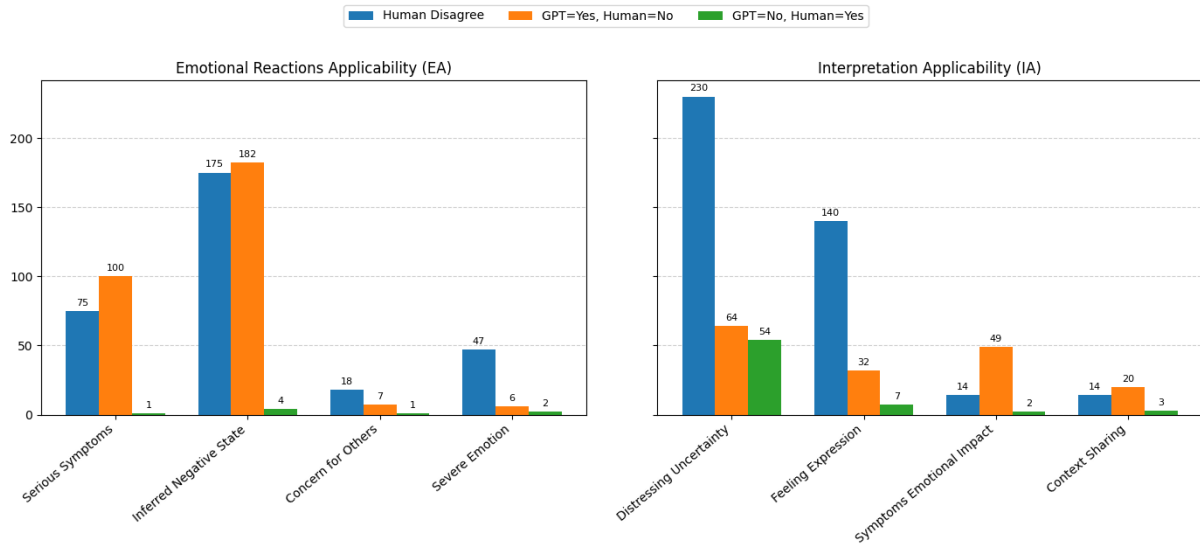


Figure 2: Three-way divergence for every subcategory. Orange = *Annotator Spread in Humans* (One Applicable, other not); Blue = *LLM-Adds Empathy Dimension* (GPT Applicable, Humans Not); Green = *LLM-Omits Empathy Dimension* (GPT Not, Humans Applicable).

signed primarily for open-domain emotional support and often assume the user’s emotional state is explicit or readily inferable. In contrast, general health queries frequently contain ambiguous, subtle, or medically grounded cues about whether emotional reactions or interpretive support are appropriate. Our anticipatory assessment framework identifies the *empathetic needs embedded in the patient’s query* prior to generation, and can be integrated with cause-aware, figurative-language, and commonsense-augmented methods to guide empathetic response generation in clinical and general-health settings.

However, EAF faces challenges from **subjective inference**, particularly when cues are implicit (e.g., *Inferred Negative Emotional State*). As appraisal theory suggests, divergent interpretations of distress often reflect genuine ambiguity rather than noise (Wondra and Ellsworth, 2015). Furthermore, these appraisals are shaped by **cultural priors of emotion**. Eichbaum et al. (2023) warn that Western-centric empathy models can misfire cross-culturally. Indeed, GPT-4o—trained on predominantly Western data (Johnson et al., 2022)—often labeled minor inconveniences as empathy-worthy where our South Asian annotators did not. This highlights a **cultural bias** inherent to LLMs that encode American norms (Cao et al., 2023).

The NLP community increasingly embraces this variability through multi-annotator models and annotator-aware representations that yield cali-

brated uncertainty estimates and capture interpretive styles (Davani et al., 2022; Mokhberian et al., 2023). Gordon et al.’s (Gordon et al., 2021) *jury learning* further shows how selecting annotator subsets aligned with demographic perspectives can preserve pluralism (Gordon et al., 2021). In clinical empathy contexts, retaining subjective variability can help anticipate diverse patient needs, including domain-specific perspectives (e.g., oncologists who prioritize emotional support as central to care) (Dekker et al., 2020). Building on this, we advocate extending disagreement-aware and annotator-aware frameworks toward **diversity-aware modeling** that explicitly accounts for culturally patterned differences in what counts as an empathy need, rather than collapsing them into a single consensus label.

This work makes three contributions to clinical empathy in NLP. First, we introduce the Empathy Applicability Framework (EAF), shifting from reactive to anticipatory modeling. Second, we establish a benchmark of 1,300 patient queries demonstrating reliable EAF labels. Third, our analysis identifies challenges, namely subjective inference, clinical-severity ambiguity, and contextual hardship, as opportunities to embrace interpretive pluralism via multi-annotator frameworks. By combining a practical framework with empirical operationalization, this work advances empathy modeling that respects interpretive complexity while remaining computationally tractable.

7 Limitations

Our study faces three key constraints, the first two mirroring limitations reported by Ali et al. (2025). First, we relied on only two human annotators, neither of whom had clinical training, which limited the range of perspectives represented; expanding the size, clinical expertise, and cultural diversity of the annotator pool would better capture the variability of empathy judgments. Second, all automatic annotations were produced with GPT-4o—selected for its widespread availability through ChatGPT—but this exclusive focus on the GPT series limits the generalization of our findings to other model architectures (e.g., Gemini, Claude, GPT reasoning models, or open-source alternatives). Third, human annotators selected a single most-salient subcategory per dimension, while GPT-4o returned multiple subcategories; this procedural mismatch hinders direct comparison of disagreement patterns, and aligning the guidelines would allow for more rigorous evaluation. Future work should therefore involve a more diverse set of human annotators, evaluate multiple LLM families trained under different specifications, and standardize annotation procedures between humans and models to obtain broader insights for improving empathy modeling in NLP for clinical contexts.

8 Ethical considerations

We developed the EAF to augment not replace clinician empathy judgments. Deploying EAF therefore requires close attention to several intertwined ethical risks that must be mitigated through thoughtful design and implementation.

A primary concern is the moral and social impact of artificial empathy. Because LLMs lack authentic emotional experience, we must ask whether the ‘applicable emotional reactions’ they generate can truly convey warmth or connection. If users perceive these reactions as hollow or manipulative, an *uncanny valley* effect could ensue, in which attempted comfort backfires by appearing inauthentic. Determining *whether, when, and how* automated empathy should be implemented, and addressing potential deception or user discomfort, requires a systematic study of user perceptions of authenticity versus artificiality.

A second mirror image danger arises from the same gap between simulated language and genuine feeling. As *Empathic AI Can’t Get Under the Skin* discussed, LLMs lack the biological and

psychological underpinnings that ground human empathy, yet their empathic language can evoke real emotional responses (Nature Machine Intelligence, 2024). Kirk et al. warn that users may form perceived emotional bonds with such systems, risking unhealthy attachment or disclosure of sensitive information (Nature Machine Intelligence, 2024). Thus, rejection born of perceived inauthenticity and devotion born of mistaken authenticity are twin failure modes rooted in the same ontological limitation.

For these reasons, we insist that the EAF be used strictly within a *human-in-the-loop* pipeline. Clinicians must retain final authority over how and when empathy is expressed, supported by transparent rationales and safeguards that guard against both deceptive alienation and false intimacy, thus protecting patients from the dual harms of artificial empathy.

References

- Iqra Ali, Jesse Atuhurra, Hidetaka Kamigaito, and Taro Watanabe. 2025. Hlu: Human vs llm generated text detection dataset for urdu at multiple granularities. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3495–3510.
- Maria Antoniak, Aakanksha Naik, Carla S Alvarado, Lucy Lu Wang, and Irene Y Chen. 2024. Nlp for maternal healthcare: Perspectives and guiding principles in the age of llms. In *Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency*, pages 1446–1463.
- Monica Oliveira Bernardo, Dario Cecílio-Fernandes, Patrício Costa, Thelma A Quince, Manuel João Costa, and Marco Antonio Carvalho-Filho. 2018. Physicians’ self-assessed empathy levels do not correlate with patients’ assessments. *PloS one*, 13(5):e0198488.
- Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. Assessing cross-cultural alignment between chatgpt and human societies: An empirical study. *arXiv preprint arXiv:2303.17466*.
- Yibo Chai, Fengyang Wu, Rui Sun, Zhongliang Zhang, Jie Bao, Runxin Ma, Qizhou Peng, Danqin Wu, Yexing Wan, and Keyu Li. 2019. Predicting future alleviation of mental illness in social media: an empathy-based social network perspective. In *2019 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom)*, pages 1564–1571. IEEE.

711	Xinhao Chen, Chong Yang, Man Lan, Li Cai, Yang	Ian Hsu, Somnath Saha, Phillip Todd Korthuis, Victo-	766
712	Chen, Tu Hu, Xinlin Zhuang, and Aimin Zhou.	ria Sharp, Jonathon Cohn, Richard D Moore, and	767
713	2024. Cause-aware empathetic response genera-	Mary Catherine Beach. 2012. Providing support to	768
714	tion via chain-of-thought fine-tuning. <i>arXiv preprint</i>	patients in emotional encounters: a new perspective	769
715	<i>arXiv:2408.11599</i> .	on missed empathic opportunities. <i>Patient education</i>	770
716	Aida Mostafazadeh Davani, Mark Díaz, and Vinodku-	<i>and counseling</i> , 88(3):436–442.	771
717	mar Prabhakaran. 2022. Dealing with disagreements:	Yan Hu, Qingyu Chen, Jingcheng Du, Xueqing Peng,	772
718	Looking beyond the majority vote in subjective an-	Vipina Kuttichi Keloth, Xu Zuo, Yujia Zhou, Zehan	773
719	notations. <i>Transactions of the Association for Com-</i>	Li, Xiaoqian Jiang, Zhiyong Lu, et al. 2024. Im-	774
720	<i>putational Linguistics</i> , 10:92–110.	proving large language models for clinical named	775
721	Joost Dekker, Jeanet Karchoud, Annemarie MJ	entity recognition via prompt engineering. <i>Journal</i>	776
722	Braamse, Hilde Buiting, Inge RHM Konings, Myra E	<i>of the American Medical Informatics Association</i> ,	777
723	van Linde, Claudia SEW Schuurhuizen, Mirjam AG	31(9):1812–1820.	778
724	Sprangers, Aartjan TF Beekman, and Henk MW Ver-	Rebecca L Johnson, Giada Pistilli, Natalia Menéndez-	779
725	heul. 2020. Clinical management of emotions in	González, Leslye Denisse Dias Duran, Enrico Panai,	780
726	patients with cancer: introducing the approach “emo-	Julija Kalpokiene, and Donald Jay Bertulfo. 2022.	781
727	tional support and case finding”. <i>Translational be-</i>	The ghost in the machine has an american ac-	782
728	<i>havioral medicine</i> , 10(6):1399–1405.	cent: value conflict in gpt-3. <i>arXiv preprint</i>	783
729	Quentin Eichbaum, Charles-Antoine Barbeau-Meunier,	<i>arXiv:2203.07785</i> .	784
730	Mary White, Revathi Ravi, Elizabeth Grant, Helen	Erica Koranteng, Arya Rao, Efren Flores, Michael Lev,	785
731	Riess, and Alan Bleakley. 2023. Empathy across	Adam Landman, Keith Dreyer, and Marc Succi. 2023.	786
732	cultures—one size does not fit all: from the ego-logi-	Empathy and equity: Key considerations for large	787
733	cal to the eco-logical of relational empathy. <i>Advances in</i>	language model adoption in health care. <i>JMIR Medi-</i>	788
734	<i>Health Sciences Education</i> , 28(2):643–657.	<i>cal Education</i> , 9:e51199.	789
735	Ronald M Epstein and Richard L Street Jr. 2007. Patient-	Allison Lahnala, Charles Welch, David Jurgens, and	790
736	centered communication in cancer care: promoting	Lucie Flek. 2022. A critical reflection and forward	791
737	healing and reducing suffering.	perspective on empathy and natural language process-	792
738	Xiang Gao and Kamalika Das. 2024. Customizing lan-	ing. <i>arXiv preprint arXiv:2210.16604</i> .	793
739	guage model responses with contrastive in-context	Allison Claire Lahnala, Béla Neuendorf, Alexander	794
740	learning. In <i>Proceedings of the AAAI Conference</i>	Thomin, Charles Welch, Tina Stibane, and Lucie Flek.	795
741	<i>on Artificial Intelligence</i> , volume 38, pages 18039–	2024. Appraisal framework for clinical empathy: A	796
742	18046.	novel application to breaking bad news conversa-	797
743	Mitchell L Gordon, Kaitlyn Zhou, Kayur Patel, Tat-	sations. In <i>Proceedings of the 2024 Joint International</i>	798
744	sunori Hashimoto, and Michael S Bernstein. 2021.	<i>Conference on Computational Linguistics, Language</i>	799
745	The disagreement deconvolution: Bringing machine	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,	800
746	learning performance metrics in line with reality. In	pages 1393–1407.	801
747	<i>Proceedings of the 2021 CHI Conference on Human</i>	Andrew Lee, Jonathan K Kummerfeld, Larry An, and	802
748	<i>Factors in Computing Systems</i> , pages 1–14.	Rada Mihalcea. 2023. Empathy identification sys-	803
749	Clarissa Guidi and Chiara Traversa. 2021. Empathy	tems are not accurately accounting for context. In	804
750	in patient care: from ‘clinical empathy’ to ‘empathic	<i>Proceedings of the 17th Conference of the European</i>	805
751	concern’. <i>Medicine, Health Care and Philosophy</i> ,	<i>Chapter of the Association for Computational Lin-</i>	806
752	24:573–585.	<i>guistics</i> , pages 1686–1695.	807
753	Lianne Hermans, Tim Olde Hartman, and Patrick W	Gyeongun Lee, Zhu Wang, Sathya N Ravi, and Natalie	808
754	Dielissen. 2018. Differences between gp perception	Parde. 2025. From heart to words: Generating em-	809
755	of delivered empathy and patient-perceived empa-	pathetic responses via integrated figurative language	810
756	thy: a cross-sectional study in primary care. <i>British</i>	and semantic context signals. In <i>Findings of the As-</i>	811
757	<i>Journal of General Practice</i> .	<i>sociation for Computational Linguistics: ACL 2025</i> ,	812
758	Hinke Hoffstädt, Jacqueline Stouthard, Maartje C Mei-	pages 4490–4502.	813
759	jers, Janine Westendorp, Inge Henselmans, Peter	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, Steve	814
760	Spreeuwenberg, Paul de Jong, Sandra van Dulmen,	Jiang, and You Zhang. 2023. Chatdoctor: A medical	815
761	and Liesbeth M van Vliet. 2020. Patients’ and clin-	chat model fine-tuned on a large language model	816
762	icians’ perceptions of clinician-expressed empathy	meta-ai (llama) using medical domain knowledge.	817
763	in advanced cancer consultations and associations	<i>Cureus</i> , 15(6).	818
764	with patient outcomes. <i>Palliative Medicine Reports</i> ,	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	819
765	1(1):76–83.	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	820
		Luke Zettlemoyer, and Veselin Stoyanov. 2019.	821
		Roberta: A robustly optimized bert pretraining ap-	822
		proach. <i>arXiv preprint arXiv:1907.11692</i> .	823

824	Lauren A McCormack, Katherine Treiman, Douglas Rupert, Pamela Williams-Piehot, Eric Nadler, Neeraj K Arora, William Lawrence, and Richard L Street Jr. 2011. Measuring patient-centered communication in cancer care: a literature review and the development of a systematic approach. <i>Social science & medicine</i> , 72(7):1085–1095.	A Empathy Applicability Framework Detail	874
825			875
826		A.1 Emotional Reactions in General Health Queries	876
827			877
828		A.1.1 Definition	878
829		Emotional Reactions refer to expressions of warmth, compassion, concern, or similar feelings conveyed by a doctor in response to a patient’s query. These reactions aim to provide emotional support and reassurance to the patient.	879
830			880
831	Quinn McNemar. 1947. Note on the sampling error of the difference between correlated proportions or percentages. <i>Psychometrika</i> , 12(2):153–157.		881
832			882
833		A.1.2 Emotional Reactions Not Applicable (N/A)	884
834	N Mokhberian, MG Marmarelis, FR Hopp, V Basile, F Morstatter, and K Lerman. 2023. Capturing perspectives of crowdsourced annotators in subjective learning tasks. arXiv preprint. <i>arXiv preprint arXiv:2311.09743</i> .	Emotional reactions are not necessary or expected in the doctor’s response when the patient’s query is factual, neutral, or a simple advice request, without expressing emotional distress. Below are detailed categories reflecting when emotional reactions are not applicable:	885
835			886
836		1. Purely Factual Medical Queries Description: The patient requests specific medical information, including explanations of medical concepts, without emotional distress or underlying distressing uncertainty.	887
837		Examples:	888
838		• "What is the use of Tylenol?"	889
839	Diane S Morse, Elizabeth A Edwardsen, and Howard S Gordon. 2008. Missed opportunities for interval empathy in lung cancer communication. <i>Archives of internal medicine</i> , 168(17):1853–1858.	• "Is it possible to outgrow a seafood allergy?"	890
840		2. General Health Management Without Emotional Involvement Description: The patient seeks guidance on health management, follows up on prior advice, or requests basic guidance on minor health issues, without expressing emotional distress or underlying distressing uncertainty. Here the guidance is on what the patient should do.	891
841		Examples:	892
842		• "I’m managing diabetes with insulin. How often should I check my blood sugar levels?"	893
843	Nature Machine Intelligence. 2024. Empathic ai can’t get under the skin . <i>Nature Machine Intelligence</i> , 6:495.	• "I have swelling in my ankle after a long walk. Should I be concerned?"	894
844		• "I had an X-ray for a fracture; should it be strapped or cast right away?"	895
845		3. Diagnosis Requests with Neutral Symptom Descriptions Description: The patient describes symptoms neutrally without expressing emotional distress or underlying distressing uncertainty. Here the request is about asking what the doctor thinks the issue is.	896
846	Joanne K Olson. 1995. Relationships between nurse-expressed empathy, patient-perceived empathy and patient distress. <i>Image: The Journal of Nursing Scholarship</i> , 27(4):317–322.	Examples:	897
847		• "I have intermittent knee pain from working out. How would I know if I tore cartilage?"	898
848			899
849			900
850	Torkel Richert, Björn Johnson, and Bengt Svensson. 2018. Being a parent to an adult child with drug problems: Negative impacts on life situation, health, and emotions. <i>Journal of Family Issues</i> , 39(8):2311–2335.		901
851			902
852			903
853			904
854			905
855	Sermo Team. 2025. Can physicians and patients trust AI doctor apps like ChatGPT? https://www.sermo.com/resources/ai-doctor-app/ . Blog post; accessed 22 July 2025.		906
856			907
857			908
858			909
859	Ashish Sharma, Adam S Miner, David C Atkins, and Tim Althoff. 2020. A computational approach to understanding empathy expressed in text-based mental health support. <i>arXiv preprint arXiv:2009.08441</i> .		910
860			911
861			912
862			913
863	Lanrui Wang, Jiangnan Li, Chenxu Yang, Zheng Lin, Hongyin Tang, Huan Liu, Yanan Cao, Jingang Wang, and Weiping Wang. 2025. Sibyl: Empowering empathetic dialogue generation in large language models via sensible and visionary commonsense inference. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 123–140.		914
864			915
865			916
866			917
867			918
868			919
869			920
870	Joshua D Wondra and Phoebe C Ellsworth. 2015. An appraisal theory of empathy and other vicarious emotional experiences. <i>Psychological review</i> , 122(3):411.		921
871			922
872			
873			

923	• "Hello. I am having pain in my jaw area, immediately in front of my left ear. The pain is random. My feeling is it is somehow related to sinus but that's just a gut feeling."		
924			
925			
926			
927	4. Hypothetical Medical Queries Without Emotional Concern		
928	Description: The patient inquires about hypothetical situations without emotional involvement.		
929			
930			
931	Examples:		
932	• "If someone has XYZ symptoms, what might be the cause?"		
933			
934	• "What would happen if a person skipped their medication?"		
935			
936	A.1.3 Emotional Reactions Applicable		
937	Definition: Emotional reactions are necessary or expected in the doctor's response when:		
938			
939	• The patient expresses emotions like fear, worry, frustration, or distress.		
940			
941	• The patient implies emotional distress over symptoms affecting their well-being.		
942			
943	• The patient's tone suggests a need for reassurance or emotional support.		
944			
945	• The patient is expressing concern for a close relation (e.g., a child, spouse).		
946			
947	Below are detailed categories reflecting when emotional reactions are applicable:		
948			
949	1. Seriousness of Symptoms Definition: The patient describes symptoms that suggest a life-threatening or chronic health condition significantly impacting long-term health or quality of life. This includes diseases like cancer, heart disease, mental health issues, or chronic conditions leading to disability. The symptoms suggest a life-threatening or serious health condition that could significantly impact long-term health or quality of life.		
950			
951			
952			
953			
954			
955			
956			
957			
958			
959	Examples:		
960	• "My father has been having severe chest pains and shortness of breath. Could it be a heart attack?"		
961			
962			
963	• "I've been experiencing numbness and weakness in my limbs for months. Could this be multiple sclerosis?"		
964			
965			
966	• "I'm 78 and have been told I have a floating hernia after bowel cancer surgery. Can it be cured?"		
967			
968	2. Severe Negative Emotion Expressed Definition: The patient explicitly states intense emotions such as fear, frustration, or anger regarding their health.		
969			
970			
971			
	Examples:		972
	• "I feel depressed and anxious like never before. I cannot sleep at night."		973
			974
	• "I am scared and plan on taking my son to the doctor. Should I be overly worried?"		975
			976
	• "I'm terrified about my recent diagnosis of cancer."		977
			978
	3. Underlying Negative Emotional State Inferred Definition: The patient implies emotional distress that isn't explicitly stated but can be inferred from their tone or descriptions, such as subtle signs of emotional worry, frustration, or distress about delays or uncertainties. Focus on emotional worry, not the medical concern.		979
			980
			981
			982
			983
			984
			985
	Examples:		986
	• "I am starting to get a little alarmed by this spotting after ovulation. Is this cause for concern?" (Worry inferred)		987
			988
			989
	• "I have been trying to conceive, and the report does not look right to me. I just want to take a second opinion." (Anxiety inferred)		990
			991
			992
	• "I need to be a bit more at ease after what I read about diabetic enteropathy. I was a bit scared if it might be fatal." (Fear inferred)		993
			994
			995
	4. Concern Severity for Close Relations Definition: The patient is asking on behalf of someone with whom they share a close, protective relationship, implying heightened emotional concern.		996
			997
			998
			999
	Examples:		1000
	• "Hello, I am the mother of a five-year-old. He has a small lump that hasn't gone away. Should I take him to a dermatologist?"		1001
			1002
			1003
	• "My son recently started daycare and has gotten sick. His fever was 102.9. Should I take him to the hospital?"		1004
			1005
			1006
	A.2 Interpretations in General Health Queries		1007
			1008
	A.2.1 Definition		1009
	Interpretations refer to the communication of an understanding of the patient's feelings (expressed or implied) and/or experiences (contextual factors) inferred from the patient's query. It's about recognizing and articulating what the patient is feeling and why, based on their situation, concerns, and history.		1010
			1011
			1012
			1013
			1014
			1015
			1016
	A.2.2 Interpretations Applicable		1017
	Interpretations are necessary when the patient's query requires the doctor to communicate an un-		1018
			1019

1020	derstanding of the patient’s feelings (expressed or implied) and/or experiences (contextual factors).	1069
1021		1070
1022	This involves acknowledging emotions, underlying concerns, or contextual elements that influence the patient’s emotional state. Below are detailed categories reflecting when interpretations are applicable:	1071
1023		1072
1024		1073
1025		1074
1026		1075
1027	1. Expression of Feelings (Explicit or Implicit)	1076
1028	Description:	1077
1029	The patient expresses emotions directly or implies them through language or tone. This includes feelings such as fear, anxiety, frustration, sadness, or hopelessness.	1078
1030		1079
1031		1080
1032	Examples:	1081
1033	• Explicit Expression:	1082
1034	– "I’m really scared about these chest pains."	1083
1035	– "I’m frustrated because my symptoms aren’t improving."	1084
1036	– "I have been in severe pain. It hurts so bad getting out of bed."	1085
1037		1086
1038		1087
1039		1088
1040	• Implicit Expression:	1089
1041	– "I guess I have to accept this is how things will be now."	1090
1042	– "Nothing seems to be helping."	1091
1043	– "I don’t know what to do anymore."	1092
1044		1093
1045	2. Sharing Experiences or Contextual Factors Affecting Emotional State and Well-being	1094
1046	Description:	1095
1047	The patient shares personal experiences, contextual factors, or circumstances that influence their health and emotional state. These include social, environmental, or personal situations beyond medical concerns that affect their emotional state.	1096
1048		1097
1049		1098
1050	Examples:	1099
1051	• "With my father’s illness and financial stress, I’m feeling overwhelmed."	1100
1052	• "I’ve been under a lot of pressure at work, and now I’m having trouble sleeping."	1101
1053	• "Ever since the accident, I can’t stop thinking about what happened."	1102
1054	• "I recently moved to a different state, haven’t found a general practitioner, and haven’t paid my high deductible for the year."	1103
1055		1104
1056		1105
1057		1106
1058		1107
1059		1108
1060		1109
1061		1110
1062		1111
1063	3. Expressions of Distressing Uncertainty About Health or Treatment	1112
1064	Description:	1113
1065	Uncertainties, confusion, or mistrust about their health status, treatment, or future are leading to emotional distress. This includes questions about prognosis, treatment effectiveness, or doubt about potential outcomes that indicate or imply underlying emotional distress. The focus should not be on uncertainty alone but specifically on uncertainty that reflects or suggests emotional distress in the patient.	1114
1066		
1067	Examples:	
1068	• "I’m not sure if this treatment is really working for me."	
	• "Do you think I should get a second opinion?"	
	• "Will chemo be fatal?"	
	• "Should my wife also get examined?"	
	• "Is this something that sounds like I should consider doing?"	
	• "I am wondering if I should see a doctor."	
	4. Symptoms Significantly Affecting Emotional Well-being or Daily Life	
	Description:	
	The patient describes symptoms that significantly impact their emotional well-being or daily functioning, and they express or imply emotional distress because of these symptoms. The key is the emotional impact of the symptoms, not just the symptoms themselves.	
	Examples:	
	• "My symptoms have been affecting my job for months."	
	• "I’m so tired all the time that I can’t take care of my kids properly."	
	• "These migraines are making it impossible to enjoy my hobbies."	
	• "The pain is getting worse every day, and it’s really wearing me down."	
	A.2.3 Interpretations Not Applicable	
	Interpretations are not necessary when the patient’s query does not require the doctor to communicate an understanding of the patient’s feelings or experiences. This occurs when:	
	• The query is straightforward, factual, or routine.	
	• There are no expressed or implied feelings needing acknowledgment.	
	• There are no contextual factors (experiences) or underlying uncertainty concerns leading to emotional distress that require understanding.	
	Below are detailed categories reflecting when interpretations are not applicable:	

1115	1. Straightforward Medical Queries Lacking Emotion, Distressing Uncertainty, and Context			
1116				
1117	Description:	The patient requests specific medical information or explanations of medical concepts without expressing emotional distress, underlying distressful uncertainty, or providing context (social, environmental, or personal situations) implying an emotional state. These queries are strictly informational and lack emotional or experiential elements requiring interpretation.		
1118				
1119				
1120				
1121				
1122				
1123				
1124				
1125	Examples:			
1126	• "What is the use of Tylenol?"			
1127	• "Hello doctor, I would like to get an opinion regarding the attached chest radiograph. I wish to know if there are any abnormalities like scarring."			
1128				
1129				
1130				
1131	2. General Health Management Requests Without Emotion, Context, and Distressing Uncertainty			
1132				
1133				
1134	Description:	The patient seeks guidance on health management, follows up on prior advice, or requests basic guidance on minor health issues without expressing emotional distress, underlying distressful uncertainty, or providing contextual factors (social, environmental, or personal situations) that imply an emotional state. Here the guidance is on what the patient should do.		
1135				
1136				
1137				
1138				
1139				
1140				
1141				
1142	Examples:			
1143	• "I'm managing diabetes with insulin. How often should I check my blood sugar levels?"			
1144	• "I have intermittent knee pain from working out. How would I know if I tore cartilage?"			
1145	• "I had an X-ray for a fracture; should it be strapped or cast right away?"			
1146				
1147				
1148				
1149	3. Diagnosis Requests with Neutral Symptom Descriptions Lacking Distressing Uncertainty and Context			
1150				
1151				
1152	Description:	The patient describes symptoms neutrally without expressing emotional distress or underlying distressful uncertainty. They provide necessary details without implying feelings or contextual factors (social, environmental, or personal situations) that need acknowledgment. These descriptions are straightforward and lack emotional or experiential content requiring interpretation. Here the request is about asking what the doctor thinks the issue is.		
1153				
1154				
1155				
1156				
1157				
1158				
1159				
1160				
1161				
1162	Examples:			
1163	• "I have swelling in my ankle after a long walk. Should I be concerned?"			
1164				
	• "Hello doctor, I am suffering from pain in my mouth. It feels like sensitivity pain. I cannot say it is pain exactly; it is irritating a lot. No pain in teeth. It feels like itching in my gums (middle of the teeth). Please tell me what I can do."			
	4. Hypothetical Medical Queries With No Emotions, Context, and Distressing Uncertainty			
	Description:	The patient inquires about hypothetical situations or general medical information without expressing or implying personal feelings or contextual factors (social, environmental, or personal situations) that need acknowledgment.		
		These queries are theoretical and lack emotional or experiential aspects requiring interpretation.		
	Examples:			
	• "If someone has XYZ symptoms, what might be the cause?"			
	• "What would happen if a person skipped their medication?"			
	B Annotation Instructions for Human Annotators			
		Annotators received an Excel workbook containing the patient queries and a fixed header with the instructions shown in Figure 3. For each pat_query, they assigned <i>Emotional Reactions</i> and <i>Interpretations</i> labels (Applicable / Not Applicable) and selected the justifying sub-category, as defined in Appendix A. The header also links to a Google Doc, reproduced verbatim in Appendix A, that provides the full framework details for reference during annotation.		
		For accessibility and clarity, we restate the instructions here in text form.		
	B.1 Instructions Given to Annotators			
	Instructions:			
	1. Read the Document:	Access and thoroughly review the following document containing the Framework Details: defined in Appendix A. Focus on understanding the details outlined below.		
	2. Understand Emotional Reactions:			
	• Emotional Reactions Definition:	Learn what emotional reactions are and their role in doctor-patient communication.		
	• Understand when emotional reactions are applicable or not applicable by reviewing:			

1212	Sub-definitions, Subcategories and Examples that illustrate their use in the relevant scenarios.		
1213			
1214			
1215	3. Classify Emotional Reactions: For each patient query, follow these steps:		
1216			
1217	• Determine Emotional Reactions Applicability or Not Applicability: Decide whether emotional reactions are applicable or not applicable in response to the patient query.		
1218			
1219			
1220			
1221			
1222	• Select a Subcategory:		
1223	– If applicable, choose the subcategory that best explains why emotional reactions are needed in response to the patient query.		
1224			
1225			
1226	– If not applicable, select the subcategory that justifies why emotional reactions are not necessary in response to the patient query.		
1227			
1228			
1229			
1230	4. Understand Interpretations:		
1231	• Interpretations Definition: Learn what interpretations are and their role in doctor-patient communication.		
1232			
1233			
1234	• Understand when interpretations are applicable or not applicable by reviewing: Sub-definitions, Subcategories and Examples that illustrate their use in the relevant scenarios.		
1235			
1236			
1237			
1238			
1239	5. Classify Interpretations: For each patient query, follow these steps:		
1240			
1241	• Determine Interpretations Applicability or Not Applicability: Decide whether interpretations are applicable or not applicable in response to the patient query.		
1242			
1243			
1244			
1245	• Select a Subcategory:		
1246	– If applicable, choose the subcategory that best explains why interpretations are needed in response to the patient query.		
1247			
1248			
1249	– If not applicable, select the subcategory that justifies why interpretations are not necessary in response to the patient query.		
1250			
1251			
1252	B.2 Additional Verbal Clarifications		
1253	During training sessions, annotators received the following clarifications:		
1254			
1255	• If they could not understand whether the symptoms or medical issue is severe , they were allowed to briefly search online (e.g., Google) to check whether the condition is typically serious.		
1256			
1257			
1258			
		• If they were unsure whether the query was emotionally significant for the patient, they were encouraged to go with what they felt and believe their own judgment.	1259 1260 1261 1262
		• For each dimension (EA and IA), annotators were instructed to:	1263 1264
		1. Read the patient query.	1265
		2. While annotating a dimension, first read the <i>Applicable</i> definition. If they believe it fits, go through the <i>Applicable</i> subcategories one by one and tag at least the one they think fits best.	1266 1267 1268 1269 1270
		3. If the <i>Applicable</i> definition does not feel like it fits, they should still briefly review the <i>Applicable</i> subcategories to verify this.	1271 1272 1273
		4. Then move to the <i>Not Applicable</i> definition and repeat the same process with the <i>Not Applicable</i> subcategories.	1274 1275 1276
	B.3 Boundary Cases: Subjectivity and Lack of Medical Expertise		1277 1278
	Empathy applicability judgments are inherently subjective, and some patient queries lie at the boundary between emotional and informational intent. Such disagreements often reflect legitimate interpretive variability rather than simple annotation error. These challenges are explored in detail in Section 5.2 (Systematic Challenges in Operationalizing Anticipatory Empathy); here, we briefly revisit representative cases to make these boundary conditions more transparent.		1279 1280 1281 1282 1283 1284 1285 1286 1287 1288
	To make these cases transparent, we provide a supplementary file (<i>misalignment_analysis.csv</i>) listing detailed patient queries, mis-aligned labels, and qualitative interpretations. The queries highlighted as exhibiting <i>Divergent Interpretation</i> correspond to <i>Reasonable Disagreement</i> around both symptom severity and whether emotional content is present.		1289 1290 1291 1292 1293 1294 1295 1296
	C Illustrative Scenarios for EAF Operationalization		1297 1298
	See Table 4 for illustrative scenarios demonstrating the operationalization of the EAF.		1299 1300
	D Appendix: Human-GPT Agreement Analysis		1301 1302
	Table 5 presents pairwise agreement between GPT and each human annotator. “Agreed” and “Disagreed” columns denote the number of queries		1303 1304 1305

Instructions:		
<p>1. Read the Document: Access and thoroughly review the following document containing the Framework Details: https://docs.google.com/document/d/1XQZL4i2lcsQZqVNFdHXQ_XMRwanqjLUOVIZD_WiNE9I/edit?tab=t.0 Focus on understanding the details outlined below.</p> <p>2. Understand Emotional Reactions: Emotional Reactions Definition: Learn what emotional reactions are and their role in doctor-patient communication. Applicability and Not applicability of Emotional Reactions: Understand when emotional reactions are applicable or not applicable by reviewing: Sub-definitions, Subcategories and Examples that illustrate their use in the relevant scenarios.</p> <p>3. Classify Emotional Reactions: For each patient query, follow these steps: Determine Emotional Reactions Applicability or Not Applicability: Decide whether emotional reactions are applicable or not applicable in response to the patient query. Select a Subcategory: If applicable, choose the subcategory that best explains why emotional reactions are needed in response to the patient query. If not applicable, select the subcategory that justifies why emotional reactions are not necessary in response to the patient query.</p> <p>4. Understand Interpretations: Interpretations Definition: Learn what interpretations are and their role in doctor-patient communication. Applicability and Not applicability of Interpretations: Understand when interpretations are applicable or not applicable by reviewing: Sub-definitions, Subcategories and Examples that illustrate their use in the relevant scenarios.</p> <p>5. Classify Interpretations: For each patient query, follow these steps: Determine Interpretations Applicability or Not Applicability: Decide whether interpretations are applicable or not applicable in response to the patient query. Select a Subcategory:</p>		
Patient Query	Emotional	Emotional Reactions
My blood pressure has been running 91/66 to 93/62 is that low, i am 32 years old, my weight is 180. I am tired all the time. I feel weak and I never have any energy. I was also diagnosed with Situs Inversus. Should I see a doctor for my blood pressure and should i worry about it?	Not Applicable	Purely Factual Medical Queries

Figure 3: Screenshot of the annotation spreadsheet provided to annotators. The header shows the instructions and links to the framework document.

1306 where both annotators assigned the same or dif-
1307 ferent labels of Applicable or Not Applicable, re-
1308 spectively.

1309 E Model Architecture Details

1310 Each empathy dimension—Emotional Reactions
1311 (EA) and Interpretations (IA)—is modeled inde-
1312 pendently. We fine-tune a pretrained RoBERTa-
1313 based model (Liu et al., 2019) separately for each
1314 dimension, while maintaining the same overall ar-
1315 chitecture. “Independently” means each classifier
1316 learns to predict the applicability of one dimension
1317 without sharing parameters or optimization across
1318 tasks. For fine-tuning, we incorporate an attention
1319 mechanism based on a feed-forward network. The
1320 model architecture is illustrated in Figure 4.

1321 The model follows an attention-based pooling
1322 approach built on top of a pretrained RoBERTa en-
1323 coder. The encoder converts patient queries into
1324 contextualized token embeddings, capturing the
1325 meaning of each word based on its surrounding con-
1326 text. When a sentence is processed by RoBERTa, it
1327 generates a hidden representation for each token, re-
1328 flecting its contextual meaning. Unlike traditional
1329 methods that rely solely on the [CLS] token or an
1330 average of all embeddings, this model applies a
1331 learned attention mechanism to identify the most

relevant tokens for classification.

1332 Specifically, the model uses a feed-forward neu-
1333 ral network to compute attention scores for each
1334 token. A linear transformation first maps each to-
1335 ken embedding to a scalar score, which then passes
1336 through a Tanh activation to constrain values be-
1337 tween [1,1] and avoid extremes. Since not all to-
1338 kens contribute equally to classification, the model
1339 converts these raw scores into attention weights
1340 using a softmax function across the sequence. This
1341 normalization ensures that important words receive
1342 higher weights, while less relevant words are as-
1343 signed lower importance.

1345 After computing attention weights, the model
1346 performs a weighted sum of token embeddings.
1347 Tokens with higher attention scores contribute
1348 more significantly to the final pooled representa-
1349 tion, highlighting the most relevant parts of the
1350 query. This pooled vector is then passed through
1351 a classification-linear layer, which outputs logits
1352 representing the likelihood of belonging to either
1353 the "Not Applicable" or "Applicable" class. Dur-
1354 ing training, the model optimizes both the attention
1355 mechanism and the classification layer via cross-
1356 entropy loss, thereby improving accuracy in empa-
1357 thy classification.

1358 Training separate models for EA and IA avoids
1359 crosstalk between tasks. Each classifier learns

Empathy Dimension	Scenario Type	Scenario	Applicability	Explanation
Emotional Reaction	Explicit Need	<i>"Hello doctor, I am having constant eye floaters, low back and hip pain, and also my rib cage hurts. I feel depressed and anxious like never before. I cannot sleep at night. An MRI of my brain shows a tiny flare, but radiologists say it's nothing to worry about. What should I do?"</i>	Applicable	The patient explicitly expresses intense negative emotions—feeling depressed and anxious—and states an inability to sleep. An emotional reaction from the doctor is necessary to provide support and reassurance.
Emotional Reaction	Implicit Need	<i>"Hello doctor, my son has been experiencing frequent headaches over the past week. We've tried over-the-counter medications, but there's no improvement. What should we do?"</i>	Applicable	Emotional reactions are applicable here because, as Richert et al. (Richert et al., 2018) find, parents of children with health (drug) issues often experience significant distress and negative mental health effects. The mother may be experiencing worry and anxiety about her child's well-being, even if she doesn't explicitly express it.
Emotional Reaction	Not Needed	<i>"Hello doctor, I was suffering from an infection in my tonsil for the past four days. I went to an ENT specialist who prescribed antibiotics. Now my tonsil pain has subsided, but I still feel something stuck on the left side of my throat where the pain was. I have no problem swallowing. Kindly advise me on what to do next."</i>	Not Applicable	The patient provides a neutral description of symptoms without expressing emotional concern or distress. The primary need is factual medical advice. An emotional reaction from the doctor is not necessary in this case.
Interpretation	Explicit Need	<i>"Hello doctor, I am feeling extremely anxious about my upcoming surgery. I can't stop worrying about the possible complications."</i>	Applicable	The patient explicitly expresses feelings of anxiety and worry. The doctor should communicate an understanding of these feelings, acknowledging the patient's emotional state and providing appropriate support.
Interpretation	Implicit Need	<i>"Hello doctor, I've been taking the medication as prescribed, but I'm not seeing any improvement. Is there something I'm doing wrong?"</i>	Applicable	The patient implies feelings of frustration and possibly self-blame. The doctor should interpret and acknowledge these underlying feelings, demonstrating understanding and support.
Interpretation	Not Needed	<i>"I was playing with my sister's boyfriend's brother and I swung to hit him like I said we were playing around and I my wrist hit his elbow really hard when it happened my hand got really numb and my vein was hurting really bad and it's 6 hours later and my vein still hurts what should I do"</i>	Not Applicable	The query is a straightforward request for diagnosis with neutral symptom descriptions. It does not express emotions or distressing contextual factors that require acknowledgment. The doctor's response should focus solely on providing a factual diagnosis.

Table 4: Empathy Dimensions, Scenarios, Applicability, and Explanations

dimension-specific patterns from the data, resulting in a simple and modular approach that enables focused analysis of empathy applicability in patient queries.

F Prompt Design for LLM Annotations

For more detail on the prompt design used for LLM (GPT-4o, o1) based annotations, we provide here the exact prompts used in our experiments.

We used two styles of prompts:

- With-framework (contrastive) prompts:** the LLM received the full Empathy Applicability Framework for a given dimension (Emotional Reactions or Interpretations), including both Applicable and Not Applicable subcategories with examples for each. This creates a contrastive in-context signal: the model must decide between multiple subcategories across both classes.
- Without-framework prompts:** the LLM only received a short task definition (dimension defini-

Table 5: Cohen’s κ agreement scores and confusion matrix counts between GPT-4o and each human annotator for Emotional Reactions (EA) and Interpretations (IA)

Annotator 1	Annotator 2	Kappa EA	Kappa IA	Agreed EA	Disagreed EA	Agreed IA	Disagreed IA
HA2	GPT	0.4402	0.5306	917	379	988	308
HA1	GPT	0.4096	0.3612	940	356	890	406

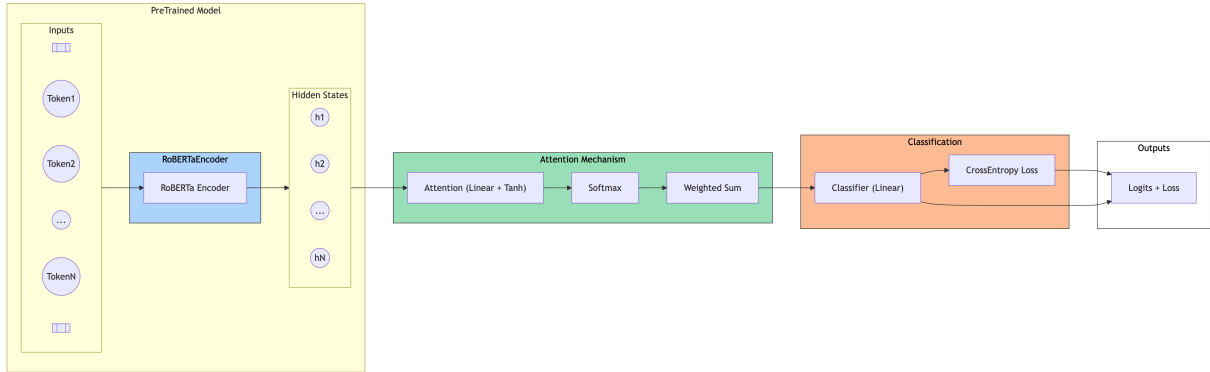


Figure 4: Empathy Dimension Applicability Model Architecture

tion + binary Applicability decision), without any subcategories or examples. This approximates a generic zero-shot setup without our framework.

F.1 Emotional Reactions: With-Framework Contrastive Prompt

The core schema passed to the LLM for Emotional Reactions with the full framework was:

```

ANNOTATION_SCHEMA = {
  "instruction": "Annotate emotional reactions in general health queries based on the following schema. For each query, return the matching subcategories. Think logically and ensure to revisit your annotation for each query",
  "definitions": {
    "Emotional Reactions": {
      "description": "Expressions of warmth, compassion, concern, or similar feelings conveyed by a doctor in response to a patient's query.",
      "categories": [
        {
          "name": "Purely Factual Medical Queries",
          "description": "The patient requests specific medical information, including explanations of medical concepts, without emotional distress or underlying distressing uncertainty.",
          "examples": ["What is the use of Tylenol?", "Is it possible to outgrow a seafood allergy?"],
          "class": "Emotional Reactions Not Applicable"
        },
        {
          "name": "General Health Management Without Emotional Involvement",
          "description": "The patient seeks guidance on health management, follows up on prior advice, or requests basic guidance on minor health issues, without expressing emotional distress or underlying distressing uncertainty. Here the guidance is on what the patient should do.",
          "examples": ["I'm managing diabetes with insulin. How often should I check my blood sugar levels?", "I have swelling

```

```

in my ankle after a long walk. Should I be concerned?"],
"class": "Emotional Reactions Not Applicable"
},
{
  "name": "Diagnosis Requests with Neutral Symptom Descriptions",
  "description": "The patient describes symptoms neutrally without expressing emotional distress or underlying distressing uncertainty. Here the request is about asking what the doctor thinks the issue is.",
  "examples": ["I have intermittent knee pain from working out. How would I know if I tore cartilage?", "Hello. I am having pain in my jaw area, immediately in front of my left ear."],
  "class": "Emotional Reactions Not Applicable"
},
{
  "name": "Hypothetical Medical Queries Without Emotional Concern",
  "description": "The patient inquires about hypothetical situations without emotional involvement.",
  "examples": ["If someone has XYZ symptoms, what might be the cause?", "What would happen if a person skipped their medication?"],
  "class": "Emotional Reactions Not Applicable"
},
{
  "name": "Seriousness of Symptoms",
  "description": "The patient describes symptoms that suggest a life-threatening or chronic health condition significantly impacting long-term health or quality of life. This includes diseases like cancer, heart disease, mental health issues, or chronic conditions leading to disability.",
  "examples": ["My father has been having severe chest pains and shortness of breath. Could it be a heart attack?", "I've been experiencing numbness and weakness in my limbs for months."],
  "class": "Emotional Reactions Applicable"
}
},

```

1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477

```

1478 {
1479   "name": "Severe Negative Emotion Expressed",
1480   "description": "The patient explicitly states
1481     intense emotions such as fear,
1482     frustration, or anger regarding their
1483     health.",
1484   "examples": ["I feel depressed and anxious
1485     like never before. I cannot sleep at
1486     night.", "I'm terrified about my recent
1487     diagnosis of cancer."],
1488   "class": "Emotional Reactions Applicable"
1489 },
1490 {
1491   "name": "Underlying Negative Emotional State
1492     Inferred",
1493   "description": "The patient implies emotional
1494     distress that isn't explicitly stated
1495     but can be inferred from their tone or
1496     descriptions, such as subtle signs of
1497     emotional worry, frustration, or
1498     distress about delays or underlying
1499     distressing uncertainties. Focus on
1500     emotional worry, not the medical
1501     concern.",
1502   "examples": ["I am starting to get a little
1503     alarmed by this spotting after
1504     ovulation. Is this cause for concern?",
1505     "I need to be a bit more at ease after
1506     what I read about diabetic enteropathy
1507     ."],
1508   "class": "Emotional Reactions Applicable"
1509 },
1510 {
1511   "name": "Concern Severity for Close Relations
1512     ",
1513   "description": "The patient is asking on
1514     behalf of someone with whom they share
1515     a close, protective relationship,
1516     implying heightened emotional concern.",
1517
1518   "examples": ["Hello, I am the mother of a
1519     five-year-old. He has a small lump that
1520     hasn't gone away.", "My son recently
1521     started daycare and has gotten sick.
1522     His fever was 102.9. Should I take him
1523     to the hospital?"],
1524   "class": "Emotional Reactions Applicable"
1525 }
1526 ]
1527 }
1528 },
1529 "output_format": "json",
1530 "example_query": {
1531   "query": "I'm scared and plan on taking my son to the
1532     doctor. Should I be overly worried?",
1533   "annotations": [
1534     {
1535       "subcategories": [
1536         {"name": "Severe Negative Emotion Expressed"},
1537         {"name": "Concern Severity for Close Relations"}
1538       ],
1539       "class": "Emotional Reactions Applicable",
1540       "reason": "The patient explicitly expresses fear
1541         regarding their son's health and shows
1542         heightened emotional concern for a close
1543         relation."
1544     }
1545 ]
1546 }
1547 }

```

```

"description": "Interpretations refer to the
1559 communication of an understanding of the
1560 patients feelings (explicit or implied) and/or
1561 experiences (contextual factors) inferred
1562 from their query. It's about recognizing and
1563 articulating what the patient is feeling and
1564 why, based on their situation, concerns, and
1565 history.",
1566 "categories": [
1567   {
1568     "name": "Expression of Feelings (Explicit or
1569     Implicit)",
1570     "description": "The patient expresses
1571     emotions directly or implies them
1572     through language or tone. This includes
1573     feelings such as fear, anxiety,
1574     frustration, sadness, or hopelessness.",
1575
1576     "examples": [
1577       "I'm really scared about these chest
1578       pains.",
1579       "I'm frustrated because my symptoms aren't
1580       improving.",
1581       "I guess I have to accept this is how
1582       things will be now.",
1583       "Nothing seems to be helping.",
1584       "I don't know what to do anymore."
1585     ],
1586     "class": "Interpretations Applicable"
1587   },
1588   {
1589     "name": "Sharing of Experiences or Contextual
1590     Factors Affecting Emotional State and
1591     Well being",
1592     "description": "The patient shares personal
1593     experiences, contextual factors, or
1594     circumstances that influence their
1595     health and emotional state. These
1596     include social, environmental, or
1597     personal situations, beyond medical
1598     concerns, that affect their emotional
1599     state.",
1600     "examples": [
1601       "With my fathers illness and financial
1602       stress, Im feeling overwhelmed.",
1603       "I've been under a lot of pressure at
1604       work, and now I'm having trouble
1605       sleeping.",
1606       "Ever since the accident, I can't stop
1607       thinking about what happened.",
1608       "I recently moved to a different state,
1609       haven't found a general practitioner
1610       , and haven't paid my high
1611       deductible for the year."
1612     ],
1613     "class": "Interpretations Applicable"
1614   },
1615   {
1616     "name": "Expressions of Distressing
1617     Uncertainty About Health or Treatment",
1618     "description": "Uncertainties, confusion, or
1619     mistrust expressed by a patient about
1620     their health status, treatment, or
1621     future that lead to significant
1622     emotional distress. This includes
1623     statements involving questions about
1624     prognosis, treatment effectiveness, or
1625     doubt about potential outcomes,
1626     specifically when accompanied by
1627     explicit or implicit signs of emotional
1628     distress.",
1629     "examples": [
1630       "Im not sure if this treatment is really
1631       working for me, and its making me
1632       anxious.",
1633       "Is this something that sounds like I
1634       should consider doing? Im so
1635       confused about whats right.",
1636       "I feel lost. Should my wife also get
1637       examined?",
1638       "Do you think theres any hope for me
1639       after trying this?"
1640     ],
1641     "class": "Interpretations Applicable"
1642   }
1643 ],
1644 "class": "Interpretations Applicable"

```

1548 F.2 Interpretations: With-Framework

1549 Contrastive Prompt

1550 The corresponding schema for Interpretations (IA)

1551 with the full framework was:

```

1552 ANNOTATION_SCHEMA = {
1553   "instruction": "Annotate interpretations in general health
1554     queries based on the following schema. For each query,
1555     return the matching subcategories. Think logically
1556     and ensure to revisit your annotation for each query",
1557   "definitions": {
1558     "Interpretations": {

```

```

1645 "name": "Symptoms Significantly Affecting          emotional distress or underlying          1732
1646 Emotional Well-being or Daily Life",           distressful uncertainty. They provide          1733
1647 "description": "The patient describes           necessary details without implying          1734
1648 symptoms that significantly impact             feelings or contextual factors. These          1735
1649 their emotional well-being or daily           descriptions are straightforward and          1736
1650 functioning, and they express or imply        lack emotional or experiential content          1737
1651 emotional distress because of these           requiring interpretation. Here the          1738
1652 symptoms. The key is the emotional           request is about asking what the doctor          1739
1653 impact of the symptoms, not just the         thinks the issue is.",                      1740
1654 symptoms themselves.",                       "examples": [                               1741
1655 "examples": [                                 "I have swelling in my ankle after a long          1742
1656 "These migraines are making it impossible    walk. Should I be concerned?",             1743
1657 to enjoy my hobbies.",                      "Hello doctor, I am suffering from pain          1744
1658 "I'm so tired all the time that I can't      in my mouth. It feels like                  1745
1659 take care of my kids properly.",           sensitivity pain. I cannot say it is          1746
1660 "My symptoms have been affecting my job      pain exactly; it is irritating a            1747
1661 for months.",                               lot. No pain in teeth. It feels like          1748
1662 "The pain is getting worse every day, and    itching in my gums (middle of the            1749
1663 it's really wearing me down."               teeth). Please tell me what I can do          1750
1664 ],                                           .",                                           1751
1665 ],                                           ],                                           1752
1666 "class": "Interpretations Applicable"        ],                                           1753
1667 },                                           "class": "Interpretations Not Applicable"    1754
1668 {                                           },                                           1755
1669 "name": "Straightforward Medical Queries      {                                           1756
1670 Lacking Emotion, Distressing                 "name": "Hypothetical Medical Queries with no          1757
1671 Uncertainty, and Context",                   Emotions, Context, and Distressing          1758
1672 "description": "The patient requests specific  Uncertainty",                               1759
1673 medical information or explanations of        "description": "The patient inquires about          1760
1674 medical concepts without expressing          hypothetical situations or general          1761
1675 emotional distress, underlying               medical information without expressing          1762
1676 distressful uncertainty or providing         or implying personal feelings or            1763
1677 context (social, environmental, or           contextual factors that need                  1764
1678 personal situations) implying an             acknowledgment. These queries are            1765
1679 emotional state. These queries are           theoretical and lack emotional or            1766
1680 strictly informational and lack             experiential aspects requiring              1767
1681 emotional or experiential elements           interpretation.",                             1768
1682 "examples": [                                 "examples": [                               1769
1683 "What is the use of Tylenol?",              "If someone has XYZ symptoms, what might          1770
1684 "Hello doctor, I would like to get an        be the cause?",                             1771
1685 opinion regarding the attached chest         "What would happen if a person skipped          1772
1686 radiograph. I wish to know if there        their medication?"                          1773
1687 are any abnormalities like scarring        ],                                           1774
1688 .",                                           "class": "Interpretations Not Applicable"    1775
1689 ],                                           ]                                           1776
1690 "class": "Interpretations Not Applicable"    },                                           1777
1691 },                                           }                                           1778
1692 {                                           },                                           1779
1693 "name": "General health management requests  "output_format": "json",                    1780
1694 Without Emotion, Context, and               "example_query": {                            1781
1695 Distressing Uncertainty",                   "query": "I'm not sure if this treatment is really          1782
1696 "description": "The patient seeks guidance on working for me.",                            1783
1697 health management, follows up on prior      "annotations": [                             1784
1698 advice, or requests basic guidance on      {                                           1785
1699 minor health issues without expressing      "subcategories": [                            1786
1700 emotional distress, underlying              {                                           1787
1701 distressful uncertainty, or providing        "name": "Expressions of Distressing          1788
1702 contextual factors that imply an            Uncertainty About Health or                  1789
1703 emotional state. Additionally, they may     Treatment"                                   1790
1704 include personal medical context, such     },                                           1791
1705 as test results, medications taken,        ],                                           1792
1706 and previous medical consultations.         "class": "Interpretations Applicable",        1793
1707 Here the guidance is on what the           "reason": "The patient explicitly expresses doubt          1794
1708 patient should do.",                       about the effectiveness of the treatment,          1795
1709 "examples": [                               which requires interpretation."              1796
1710 "I have intermittent knee pain from         ],                                           1797
1711 working out. How would I know if I         ],                                           1798
1712 tore cartilage?",                           }                                           1799
1713 "I had an X-ray for a fracture; should it   }                                           1800
1714 be strapped or cast right away?",          }                                           1801
1715 "Hi, my husband is 39, and his SGPT and    }                                           1802
1716 SGOT levels in a recent test were          }                                           1803
1717 101 and 98 respectively. His               }                                           1804
1718 triglycerides are 280, which is high        }                                           1805
1719 . His height is 168 cm and weight is       }                                           1806
1720 79 kg. What does a rise in these          }                                           1807
1721 values indicate? What precautions          }                                           1808
1722 should he take?"                            }                                           1809
1723 ],                                           }                                           1810
1724 "class": "Interpretations Not Applicable"    }                                           1811
1725 },                                           }                                           1812
1726 {                                           }                                           1813
1727 "name": "Diagnosis Requests with Neutral     }                                           1814
1728 Symptom Descriptions Lacking               }                                           1815
1729 Distressing Uncertainty and Context",       }                                           1816
1730 "description": "The patient describes        }                                           1817
1731 symptoms neutrally without expressing        }                                           1818

```

F.3 Prompts Without the Framework (Definition-Only) 1800

For the *without-framework* condition, the LLM received only a short task description and the names of the Applicability labels. No subcategories or examples were provided. 1802-1805

F.3.1 Emotional Reactions (without framework). 1806

ANNOTATION_SCHEMA = { 1808

```

1809     "instruction": "Read the patient query and decide whether
1810         emotional reactions are necessary in the response. "
1811         "Emotional reactions refer to the expressions
1812         of warmth, compassion, concern, or
1813         similar feelings "
1814         "conveyed by a doctor in response to a patient
1815         's query. "
1816         "If emotional reactions are necessary, mark it
1817         as 'Emotional Reactions Applicable'. "
1818         "If not, mark it as 'Emotional Reactions Not
1819         Applicable'. Think carefully and be
1820         consistent.",
1821     "output_format": "json"
1822 }

```

1823 F.3.2 Interpretations (without framework).

```

1824 ANNOTATION_SCHEMA = {
1825     "instruction": "Read the patient query and decide whether
1826         interpretations are necessary in the response. "
1827         "Interpretations refer to the communication of
1828         an understanding of the patients
1829         feelings "
1830         "(explicit or implied) and/or experiences (
1831         contextual factors) inferred from their
1832         query. "
1833         "If interpretations are necessary, mark it as
1834         'Interpretations Applicable'. "
1835         "If not, mark it as 'Interpretations Not
1836         Applicable'. Think carefully and be
1837         consistent.",
1838     "output_format": "json"
1839 }

```

1840 Together, these listings document the exact
1841 prompts used in both the with-framework (con-
1842 trastive) and without-framework settings for LLM
1843 annotations.

1844 G Dataset Analyses

1845 To characterize the released benchmark beyond
1846 agreement and modeling results, Figure 5 summa-
1847 rizes label base rates (H1/H2/GPT), EA–IA co-
1848 applicability patterns, and the heavy-tailed distri-
1849 bution of query lengths, providing a high-level
1850 view of dataset variability. In the following sec-
1851 tions, we present additional analyses of annotation
1852 consistency, framework coherence, and stability
1853 across the labeling process. Specifically, we ana-
1854 lyze: (i) subcategory usage and co-occurrence pat-
1855 terns across human annotators and GPT rationales;
1856 (ii) length effects on applicability with confidence
1857 intervals; (iii) run-order drift checks to assess sta-
1858 bility over the labeling sequence; and (iv) match
1859 vs. miss portions of subcategory rationales when
1860 humans and GPT agree on the overall applicability
1861 label, to assess alignment in rationales.

1862 G.1 Subcategory Prevalence and EA×IA 1863 Co-occurrence (Humans vs. GPT)

1864 Table 6 and Table 7 report the prevalence of sub-
1865 categories along with the **Applicable** or **Not Ap-**
1866 **licable** classification of each subcategory under
1867 the EAF. This makes it possible to assess the *in-*
1868 *tuitive coherence* of the framework: subcategories

1869 that encode similar affective cues or uncertainty
1870 (Applicable) should align with each other in all di-
1871 mensions, while informational or routine requests
1872 (Not Applicable) should be clustered separately.

1873 Because each query receives exactly one subcat-
1874 egorY from each human annotator, human counts
1875 sum to $N = 1296$ per dimension. In contrast, GPT
1876 may assign multiple subcategories per query; thus
1877 GPT totals exceed N (EA: 2100; IA: 2572), corre-
1878 sponding to an average of 1.62 EA subcategories
1879 and 1.98 IA subcategories per query.

1880 The prevalence distributions show both stable
1881 structure and meaningful variability in how cues
1882 are operationalized. For EA, H1 frequently uses
1883 the Not Applicable subcategory *Factual Medical*
1884 *Query* (31.8%), while H2 rarely uses it (2.9%) and
1885 instead assigns both Applicable and Not Applicable
1886 categories such as *Concern for Relations* (20.6%)
1887 and *General Health Management* (19.4%). A simi-
1888 lar pattern appears in IA: H1 frequently assigns
1889 the Not Applicable subcategory *Straightforward*
1890 *Medical Query* (33.3%), whereas H2 more often
1891 assigns the Applicable subcategory *Distressing Un-*
1892 *certainty* (40.8%). GPT assigns Applicable EA
1893 categories at higher rates, especially *Inferred Nega-*
1894 *tive Emotion* (59.4%) and *Symptom Seriousness*
1895 (27.5%), reflecting broader recall for empathic-
1896 need cues. Across dimensions, GPT again as-
1897 signs multiple IA cues per query and frequently
1898 marks Applicable interpretive cues (*Distressing*
1899 *Uncertainty*, *Feelings Expression*, *Impact on Daily*
1900 *Life*, *Sharing Affecting Context*), consistent with its
1901 broader recall of applicability signals.

1902 Figure 6 further supports the framework’s *intu-*
1903 *itive coherence* when interpreted through the Appli-
1904 cable vs. Not Applicable split of subcategories (Ta-
1905 bles 6–7). Not Applicable request-types tend to pair
1906 with Not Applicable interpretations (Not×Not),
1907 while affective/uncertainty cues (Applicable) more
1908 often co-occur with Applicable interpretive signals
1909 (App×App). Concretely, in the human heatmaps,
1910 EA *Factual Medical Query* aligns strongly with IA
1911 *Straightforward Medical Query* (H1: 30%), and
1912 *Neutral Diagnosis Request* co-occurs with its IA
1913 counterpart (H1: 19%; H2: 9%), both canonical
1914 Not×Not pairings. In contrast, Applicable EA cues
1915 align with Applicable IA cues: *Severe Negative*
1916 *Emotion* co-occurs with *Feelings Expression* (H1:
1917 15%), and *Concern for Relations* frequently pairs
1918 with *Distressing Uncertainty* (H2: 16%), reflecting
1919 clinically intuitive links between expressed (or in-
1920 ferred) affect and corresponding interpretive needs.

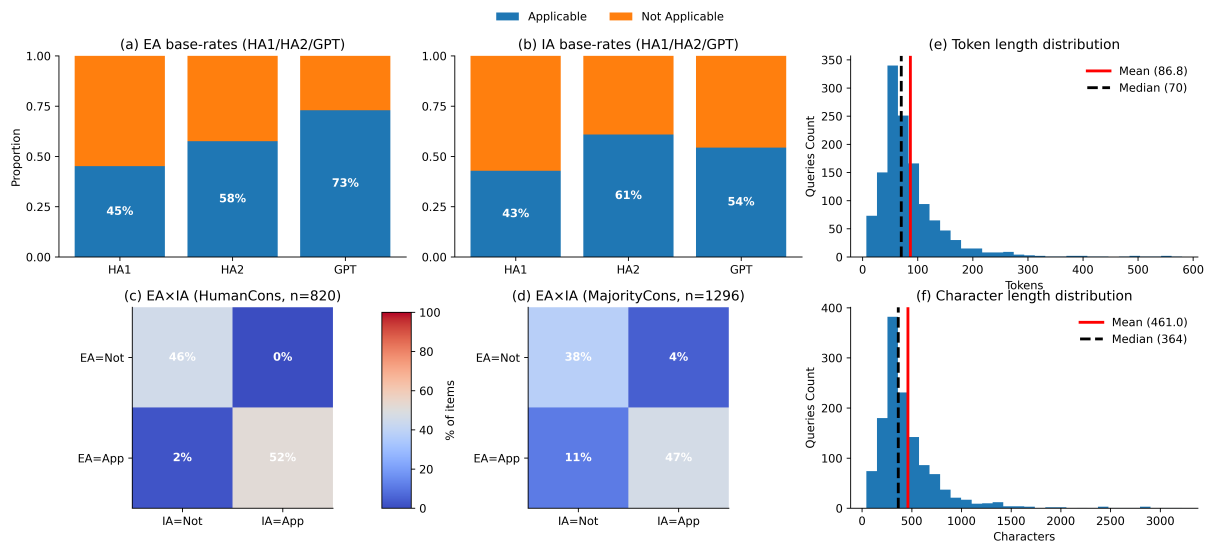


Figure 5: **Dataset overview panel: base rates, EA–IA coupling, and query length distributions.** Panels (a–b) report the binary label base rates for *Emotional Reactions* (EA) and *Interpretations* (IA) from Human Annotator 1 (HA1), Human Annotator 2 (HA2), and GPT. Bars are shown as stacked proportions of *Applicable* vs. *Not Applicable*. Panels (c–d) show EA×IA co-occurrence as 2×2 heatmaps for (c) *Human consensus* (only items where HA1 and HA2 agree on both EA and IA) and (d) *Majority consensus* (majority vote over HA1, HA2, and GPT), with each cell annotated by the percentage of items in that consensus subset; Panels (e–f) summarize query length: (e) token-count histogram (simple word tokenization) and (f) character-count histogram; vertical reference lines mark the mean (solid red) and median (dashed black). Together, the figure summarizes label prevalence, the empirical coupling between EA and IA decisions, and the distribution of textual input lengths in patient queries.

GPT exhibits broader App×App cross-pairings, most prominently EA *Inferred Negative Emotion* co-occurring with multiple Applicable IA cues including *Distressing Uncertainty* (38%), *Feelings Expression* (27%), *Impact on Daily Life* (24%), and *Sharing Affecting Context* (20%), which is expected given GPT’s multi-label rationale annotations. Overall, these structured pairings indicate that cross-dimensional co-occurrence is not arbitrary: it aligns with the EAF’s applicability semantics while also highlighting how multi-cue rationales (GPT) differ from single-label human assignments.

Additionally, these prevalence and co-occurrence results provide actionable signals for refining the EAF and improving annotation practice. First, the strong Not×Not and App×App structure suggests the framework’s applicability split is broadly coherent, but the sharp annotator skews in how “routine/informational” vs. “affective/uncertainty” cues are operationalized (e.g., heavier use of *Factual Medical Query* and *Straightforward Medical Query* versus greater use of *Concern for Relations* and *Distressing Uncertainty*) highlight subcategories that may be very broad or boundary-sensitive. These

are prime candidates for guideline refinement (clearer decision rules, additional contrastive examples, or merging/splitting categories), while consistently rare categories can be reconsidered for consolidation if they contribute limited discriminative value. Second, the co-occurrence matrices can be leveraged to diagnose systematic annotation patterns and potential drift: stable, clinically intuitive pairings (e.g., *Severe Negative Emotion* with *Feelings Expression*) indicate consistent interpretation, whereas unexpected or diffuse pairings can reveal where annotators may be diverging due to inconsistent application rather than subjectivity, and targeted review could mitigate those differences. In this sense, co-occurrence structure is not only a validation of the EAF semantics, but also a practical tool for targeted adjudication, annotator training, and future modeling choices.

G.2 Length Effects on Applicability (with Confidence Intervals)

To assess whether query length is associated with applicability judgments, we stratify items into token-length deciles and plot EA/IA applicability rates for each label source (HA1, HA2, GPT, Hu-

EA subcategory	Shortname	Class	HA1	HA2	GPT
Underlying Negative Emotional State Inferred	inferred negative emotion	App	121 (9.3)	197 (15.2)	770 (59.4)
Concern Severity For Close Relations	concern for relations	App	246 (19.0)	267 (20.6)	277 (21.4)
Severe Negative Emotion Expressed	severe negative emotion	App	209 (16.1)	132 (10.2)	177 (13.7)
Seriousness Of Symptoms	symptom seriousness	App	9 (0.7)	150 (11.6)	357 (27.5)
Diagnosis Requests With Neutral Symptom Descriptions	neutral diagnosis request	Not	262 (20.2)	165 (12.7)	226 (17.4)
Purely Factual Medical Queries	factual medical query	Not	412 (31.8)	38 (2.9)	108 (8.3)
General Health Management Without Emotional Involvement	general health management	Not	37 (2.9)	252 (19.4)	156 (12.0)
Hypothetical Medical Queries Without Emotional Concern	hypothetical medical query	Not	0 (0.0)	95 (7.3)	29 (2.2)

Table 6: EA subcategory prevalence with applicability class (App vs. Not). Values are count (percent of $N = 1296$). Humans assign one subcategory per query (counts sum to N). GPT may assign multiple subcategories per query; thus GPT counts can exceed N .

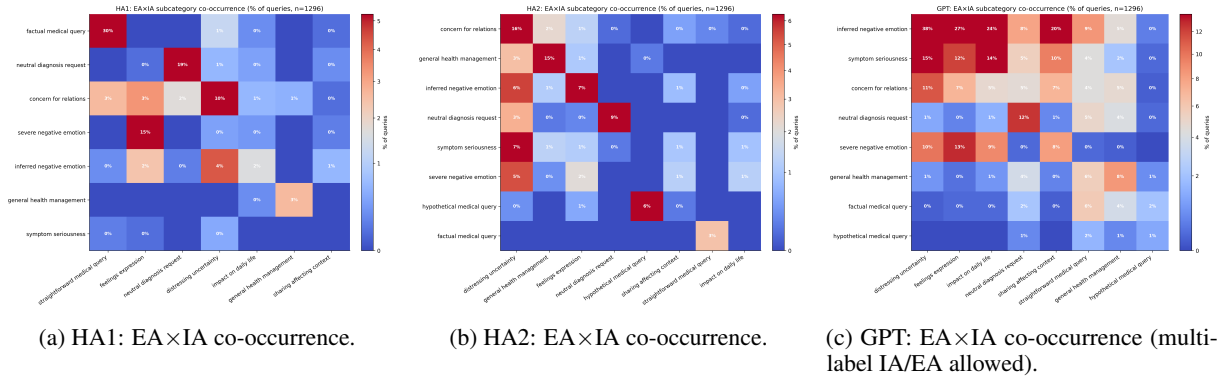


Figure 6: EA x IA subcategory co-occurrence (percent of queries, $N = 1296$). Humans provide a single EA and IA subcategory per query, yielding sharper pairings; GPT may assign multiple subcategories per query, producing broader co-occurrence patterns.

manCons, MajorityCons). The x-axis reports the mean token length per decile (μ), and we annotate the number of items per decile (n). Error bars denote confidence intervals for each applicability estimate.

Across both dimensions, applicability increases with length: shorter queries (D1–D2) receive substantially lower applicability rates, while longer queries (D9–D10) show consistently higher applicability across all sources. Although absolute rates differ by source (e.g., GPT tends to assign applicability more frequently than humans), the upward trend is shared, suggesting that length (and the additional context typically present in longer queries) is systematically associated with applicability rather than reflecting annotator idiosyncrasies alone.

G.3 Run-order Drift Checks

Finally, we evaluate whether labeling behavior drifts over the annotation sequence (e.g., fatigue or order effects). We plot the rolling mean of pairwise differences in applicability rates (window size = 32) across run index for each pair (H1–H2, H1–GPT, H2–GPT), and apply ADWIN change detection ($\delta = 0.002$) to flag potential change points.

Across $N = 1296$ items, ADWIN identifies only a small number of localized change points in each rolling-difference stream. For EA, ADWIN flags 3 points for H1–H2 (at indices 767, 1023, 1279), 2 for H1–GPT (223, 1151), and 3 for H2–GPT (191, 479, 959). For IA, ADWIN flags 3 points for H1–H2 (415, 639, 1279), 3 for H1–GPT (191, 863, 1151), and 4 for H2–GPT (223, 415, 607, 1087). Because these indices correspond to positions in

IA subcategory	Shortname	Class	HA1	HA2	GPT
Expressions Of Distressing Uncertainty About Health Or Treatment	distressing uncertainty	App	222 (17.1)	529 (40.8)	556 (42.9)
Expression Of Feelings (Explicit Or Implicit)	feelings expression	App	274 (21.1)	176 (13.6)	423 (32.6)
Symptoms Significantly Affecting Emotional Well-Being Or Daily Life	impact on daily life	App	44 (3.4)	36 (2.8)	382 (29.5)
Sharing Experiences Or Contextual Factors Affecting Emotional State And Well Being	sharing affecting context	App	15 (1.2)	48 (3.7)	308 (23.8)
Straightforward Medical Queries Lacking Emotion, Distressing Uncertainty, And Context	straightforward medical query	Not	431 (33.3)	41 (3.2)	285 (22.0)
Diagnosis Requests With Neutral Symptom Descriptions Lacking Distressing Uncertainty And Context	neutral diagnosis request	Not	269 (20.8)	118 (9.1)	339 (26.2)
General Health Management Requests Without Emotion, Context, And Distressing Uncertainty	general health management	Not	41 (3.2)	264 (20.4)	242 (18.7)
Hypothetical Medical Queries With No Emotions, Context, And Distressing Uncertainty	hypothetical medical query	Not	0 (0.0)	84 (6.5)	37 (2.9)

Table 7: IA subcategory prevalence with applicability class (App vs. Not). Values are count (percent of $N = 1296$). Humans assign one subcategory per query (counts sum to N). GPT may assign multiple subcategories per query; thus GPT counts can exceed N .

the rolling stream, each detection reflects a local shift over the surrounding ~ 32 -item neighborhood rather than a single item. Overall, the sparse detections and the absence of sustained shifts in the rolling trajectories suggest broadly stable annotation behavior over the run, with only occasional local fluctuations in pairwise disagreement.

G.4 Match vs. Miss by subcategory

In this section, we quantify subcategory-level rationale overlap in cases where GPT agrees with the human consensus label (Applicable / Not Applicable) for the given dimension. For each query, each human provides one best-fit subcategory; the human rationale set is therefore either a singleton (both humans chose the same subcategory) or a size-two set (humans chose different subcategories). We then compare this human set to GPT’s provided subcategories as rationales for labeling the query as Applicable or not. For each human-selected subcategory occurrence, we count a **Match** if GPT includes that same subcategory, and a **Miss** otherwise. Figure 9 reports, for each subcategory, the % Match vs. % Miss, along with the total number of such occurrences (N).

Across both EA and IA, match rates are generally high, indicating that when GPT agrees with humans on the binary applicability label, it often

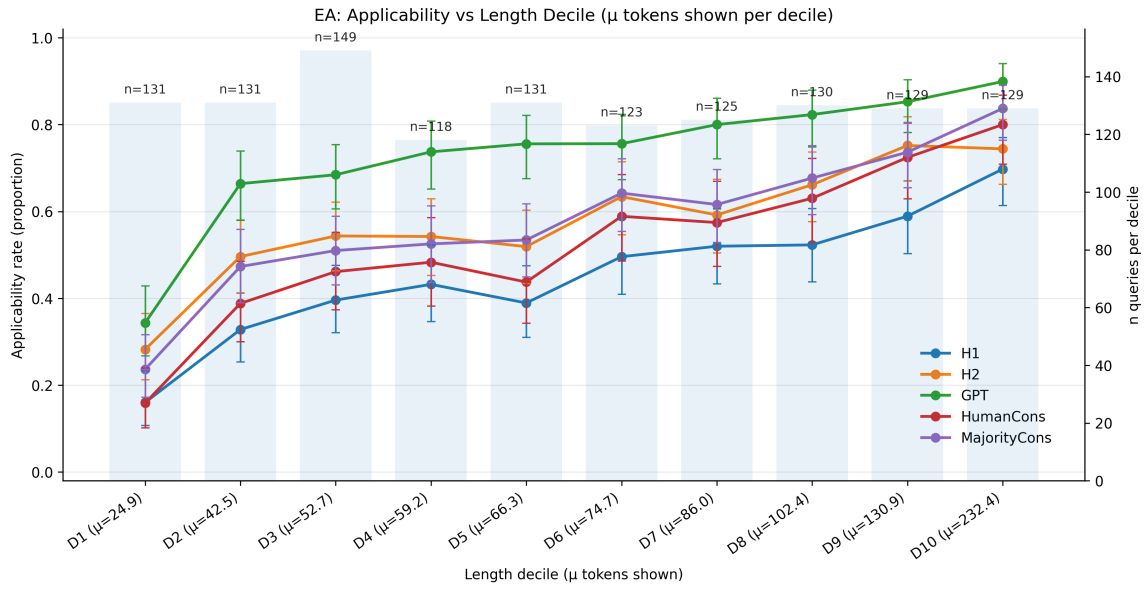
also identifies the same underlying rationale cues. However, several subcategories exhibit higher miss rates, suggesting systematic differences in how GPT justifies an agreed-upon label.

For **EA** (Figure 9a), GPT shows near-complete alignment for *Concern for Relations* ($N = 204$) and *Inferred Negative Emotion* ($N = 102$), and very strong alignment for *Neutral Diagnosis Request* ($N = 92$) and *Symptom Seriousness* ($N = 70$). In contrast, *Hypothetical Medical Query* exhibits the largest miss portion ($N = 39$), and both *Severe Negative Emotion* ($N = 189$) and *Factual Medical Query* ($N = 147$) show notable misses, indicating that GPT sometimes agrees on EA applicability while grounding its justification in different cues than the humans.

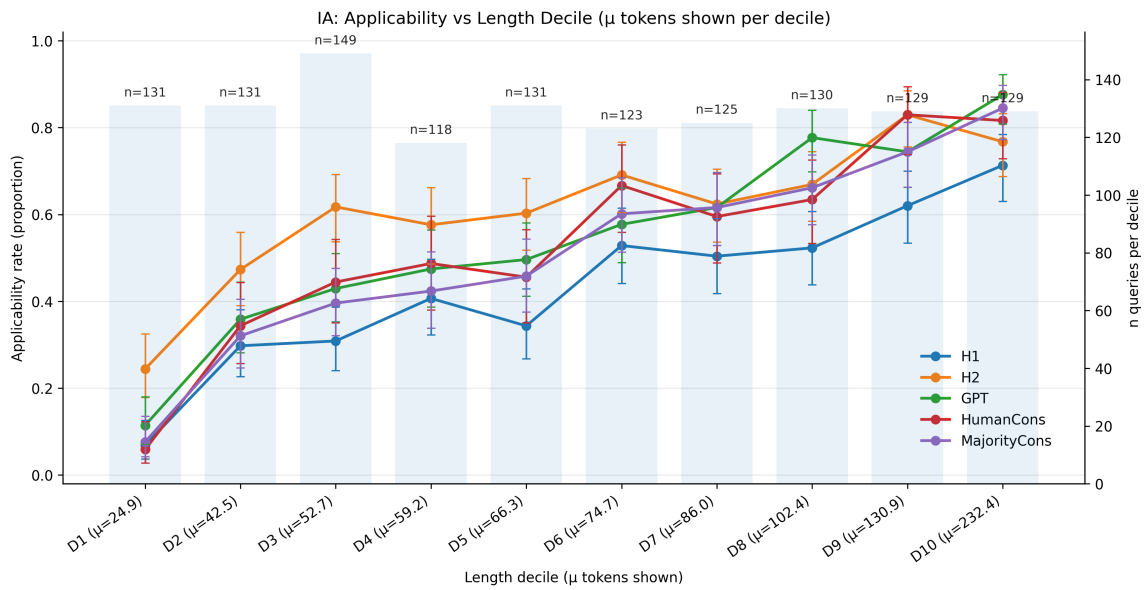
For **IA** (Figure 9b), alignment is strongest for *Distressing Uncertainty* ($N = 262$) and *Feelings Expression* ($N = 233$), and is also high for *Neutral Diagnosis Request* ($N = 133$). *Impact on Daily Life* shows a perfect match in this subset ($N = 49$), though this category is comparatively smaller. The most challenging IA category is again *Hypothetical Medical Query* ($N = 58$), which has the largest miss fraction; *Straightforward Medical Query* ($N = 183$) also shows a higher miss rate than the more affective/uncertainty categories.

Together, these patterns indicate substantial ra-

2060 tionale overlap when GPT agrees with the human
2061 consensus applicability label, but alignment varies
2062 by subcategory. Misses occur not only for hypothet-
2063 ical and some informational categories, but also for
2064 certain affective categories (e.g., *Severe Negative*
2065 *Emotion*). This suggests that mismatches reflect
2066 less a single cue-type split and more the prevalence
2067 of multi-cue queries, where multiple plausible ratio-
2068 nales can support the same applicability judgment.

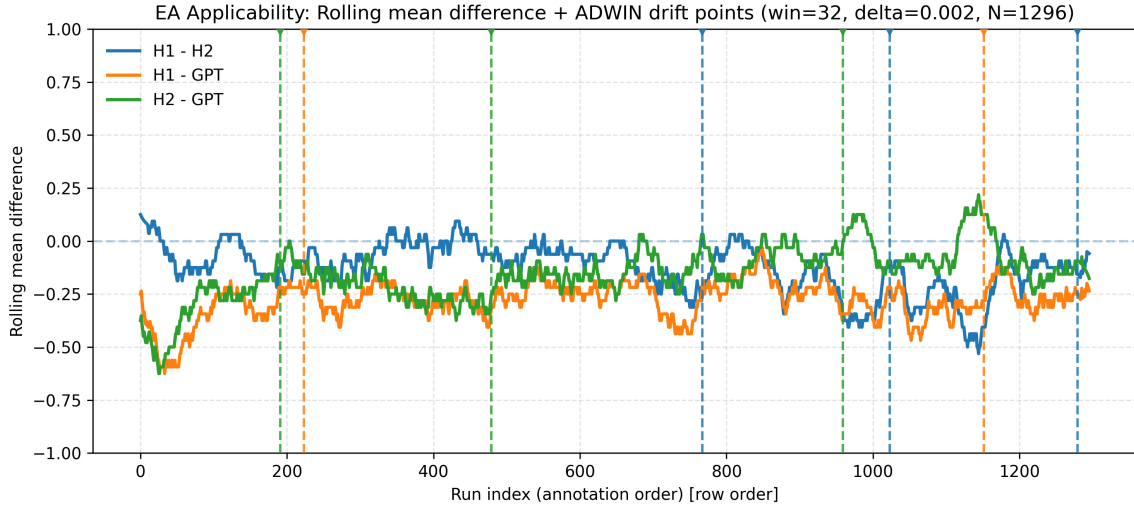


(a) EA: Applicability vs. token-length decile (μ tokens shown per decile; error bars are confidence intervals).

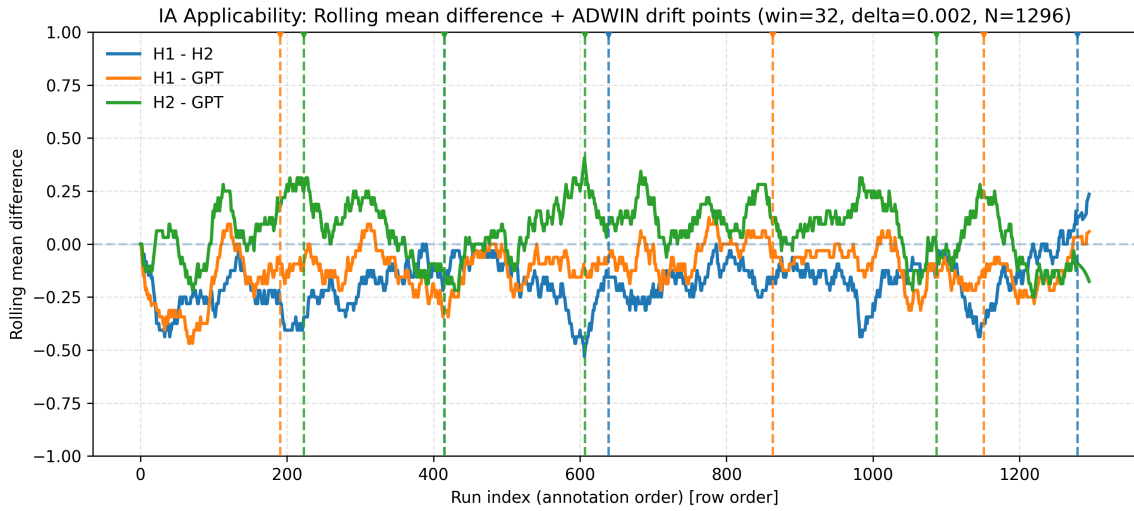


(b) IA: Applicability vs. token-length decile (μ tokens shown per decile; error bars are confidence intervals).

Figure 7: Applicability rates by query-length decile. Bars indicate the number of queries per decile (n). Error bars denote confidence intervals.



(a) EA: rolling mean pairwise difference + ADWIN drift points (win=32, $\delta = 0.002$).



(b) IA: rolling mean pairwise difference + ADWIN drift points (win=32, $\delta = 0.002$).

Figure 8: Run-order drift diagnostics using rolling mean differences in applicability rates and ADWIN change detection. Vertical dashed lines indicate detected change points.

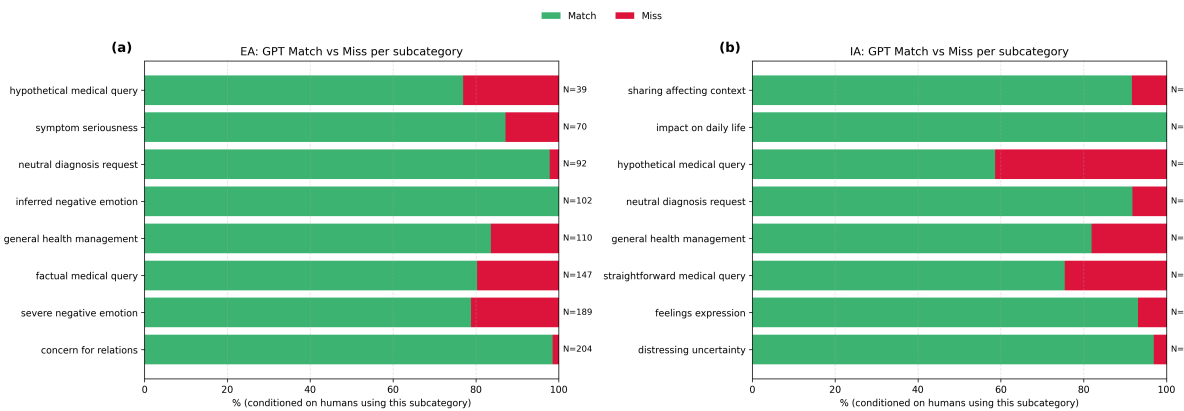


Figure 9: **Rationale Match vs. Miss by subcategory** Stacked bars show, for each subcategory, the percentage of **Match** (green) vs. **Miss** (red) between GPT and the human-selected subcategory rationale, computed only on queries where GPT agrees with the human consensus applicability label. The N label denotes the number of times the subcategory appears in the human rationale set within this agreement subset. Panels show (a) EA and (b) IA.