

EMPATHYSCORE: A Reference-free Metric for Emotional Support Conversation via Knowledge Distillation

Anonymous ACL submission

Abstract

Existing metrics for Emotional Support Conversation (ESC) struggle with the “one-to-many” nature of open-ended empathy or suffer from the high computational cost of LLM-based evaluation. To address this, we introduce EMPATHYSCORE, a reference-free, lightweight metric developed via the DISTILL-ES framework. Crucially, we construct EMPATHY-EVAL, a comprehensive distillation dataset containing fine-grained teacher annotations and pairwise preferences derived from rubric-guided LLM prompting. Leveraging this data, we distill the sophisticated reasoning of a teacher LLM into compact student scorers using a hybrid regression-ranking objective, decoupling evaluation into *Cognitive Suitability* and *Affective Resonance*. Experiments demonstrate that EMPATHYSCORE achieves state-of-the-art correlation with human judgments while being orders of magnitude more efficient than LLM judges.

1 Introduction

Text-based Emotional Support Conversations (ESC) have evolved into a critical medium for scalable psycho-social care, ranging from peer-support platforms to empathetic AI companions (Rashkin et al., 2019; Liu et al., 2021). Unlike task-oriented dialogue, effective ESC requires the helper to accompany the seeker through distress via a nuanced psychological process. Counselling theory characterizes this process as a multidimensional construct: it demands both *cognitive empathy* (understanding the seeker’s frame of reference to select the right support strategy) and *affective empathy* (responding with emotional attunement and warmth) (Rogers, 1957a; Hill, 2014; Elliott et al., 2011). Consequently, evaluating ESC systems requires metrics that go beyond surface-level fluency to assess whether responses are *cognitively suitable* for the conversation stage and *affectively resonant* with the user’s emotional state.

However, existing evaluation methodologies struggle to quantify these nuances. Standard reference-based metrics (e.g., BLEU, BERTScore), inherited from translation tasks, rely on computing similarity against a single human reference (Papineni et al., 2002; Zhang et al., 2020). This paradigm is fundamentally flawed in the context of ESC due to the “one-to-many” nature of empathy: for a given distress context, diverse responses (e.g., a comforting affirmation vs. a constructive suggestion) can be equally valid yet lexically distinct (Liu et al., 2016; Deriu et al., 2021). As a result, these metrics often penalize high-quality creative responses that diverge from the reference, showing poor correlation with human judgments.

To overcome the reference dependency, recent works have employed Large Language Models (LLMs) as reference-free judges, leveraging their reasoning capabilities to assess dialogue quality directly (Lin et al., 2023; Zhang et al., 2024). While LLM-based evaluation aligns better with human perception, it introduces significant bottlenecks for practical deployment. Relying on proprietary LLMs (e.g., GPT-4) as online judges incurs high inference latency and prohibitive costs for large-scale model development. Furthermore, sending sensitive mental health data to external APIs raises serious privacy concerns (Kocmi and Freitag, 2023; Croxford et al., 2025). The field thus faces a dilemma: traditional metrics are efficient but inaccurate, while LLM judges are accurate but inefficient and costly.

In this paper, we propose to bridge this gap by distilling the empathic reasoning of expensive LLMs into lightweight, privacy-preserving models. We introduce EMPATHYSCORE, a reference-free metric constructed via the DISTILL-ES framework. Our core insight is to treat the LLM not as a permanent judge, but as a *teacher* that provides rich supervision via rubric-guided prompting. Specifically, we prompt a teacher LLM to analyze dialogues

083 along cognitive and affective dimensions, generat- 132
084 ing both fine-grained scores and pairwise prefer- 133
085 ences. We then train two compact student scorers 134
086 (RoBERTa-based) to mimic these judgments using 135
087 a hybrid objective that combines regression (for 136
088 absolute calibration) and ranking (for relative pref- 137
089 erence learning).

090 This distillation approach yields a robust eval- 138
091 uation framework that is both *decoupled* and *effi-* 139
092 *cient*. By explicitly modeling *Cognitive Suitability* 140
093 and *Affective Resonance* as separate sub-tasks, EM- 141
094 PATHYSCORE provides interpretable feedback on 142
095 where a model succeeds or fails. Extensive experi- 143
096 ments demonstrate that our metric achieves state- 144
097 of-the-art correlation with human ratings, matching 145
098 the performance of the teacher LLM while being 146
099 orders of magnitude faster and removing the need
100 for test-time API access. EMPATHYSCORE thus
101 offers a scalable standard for assessing machine
102 empathy, enabling rapid iteration of ESC systems
103 without compromising on evaluation quality.

104 Our contributions are three-fold:

- 105 • We propose EMPATHYSCORE, a reference- 149
106 free metric that decouples empathy evaluation 150
107 into cognitive suitability and affective reso- 151
108 nance. 152
- 109 • We introduce the DISTILL-ES framework 153
110 to distill sophisticated LLM judgments into 154
111 lightweight student scorers via rubric-guided 155
112 supervision. 156
- 113 • Experiments show that our metric achieves 157
114 state-of-the-art correlation with human rat- 158
115 ings while being orders of magnitude more 159
116 efficient than LLMs. 160

117 2 Related Work

118 2.1 Emotional Support Conversation

119 Modeling for Emotional Support Conversation 161
120 (ESC) mainly follows three lines: *knowledge-* 162
121 *enhanced* methods that inject commonsense or 163
122 graphs to better understand seekers (e.g., MISC, 164
123 C³KG, GLHG)(Tu et al., 2022; Li et al., 2022; Peng 165
124 et al., 2022); *cognitive reasoning* that first infers 166
125 cues and emotions before responding (Dialogue 167
126 CoT, Cue-CoT)(Chae et al., 2023; Wang et al., 168
127 2023); and *persona modeling* for personalization 169
128 and consistency (PAL)(Cheng et al., 2023). While 170
129 ESConv provides high-quality but limited-scale 171
130 training data(Liu et al., 2021), LLM-based aug- 172
131 mentation (AugESC, ExTES)(Zheng et al., 2023,

2024) and socially grounded simulation (Social- 132
Sim)(Chen et al., 2025) expand coverage and real- 133
ism. Most recently, CARE(Zhu et al., 2025) and 134
ReflectDiffu(Yuan et al., 2025) both leverage rein- 135
forcement learning to enhance controllability and 136
empathy. 137

138 However, most studies still rely on automatic 139
metrics such as BLEU, ROUGE-L, METEOR, 140
BERTScore, and Distinct, which weakly cap- 141
ture perceived empathy and support(Zhao et al., 142
2024; Madani and Srihari, 2025), while human rat- 143
ings—though more faithful—are costly and hard to 144
scale; therefore, our work proposes *Empathy-Score*, 145
an automatic and scalable metric that better reflects 146
human-perceived empathy.

147 2.2 Empathy Measurement

148 Empathy is commonly divided into two types: *trait* 149
empathy, which captures relatively stable empa- 150
thetic dispositions, and *situational empathy*, which 151
refers to immediate, context-triggered responses; 152
the former can shape how the latter manifests in 153
specific situations, making the distinction impor- 154
tant (Davis et al., 1980; Hoffman, 2000). In psy- 155
chology, self-report questionnaires are widely used 156
to assess relatively stable empathic traits, such as 157
the *Interpersonal Reactivity Index*(Davis, 1983), 158
the *Empathy Quotient*(Baron-Cohen and Wheel- 159
wright, 2004), and the *Toronto Empathy Question-* 160
naire(Sprenge* et al., 2009).

161 In affective computing and NLP field, mea- 162
surement centers on *situational empathy* and typi- 163
cally uses rater-based judgments as ground truth to 164
train supervised deep learning models. Concretely, 165
one line predicts reader self-reported empathy/dis- 166
tress from stories(Buechel et al., 2018), while 167
another learns empathy signals from annotated 168
mental-health counseling conversations(Sharma 169
et al., 2020), with related efforts extending to other 170
genres (e.g., social media(Abdul-Mageed et al., 171
2017), persuasive-oriented dialogues(Samad et al., 172
2022)). In ESC specifically, perceived empathy 173
is a central criterion for helpful support and posi- 174
tive outcomes(Rogers, 1957b; Elliott et al., 2011); 175
therefore, we introduce *Empathy-Score*, an auto- 176
matic metric tailored for ESC that approximates 177
human judgements of perceived empathy to enable 178
scalable evaluation and training.

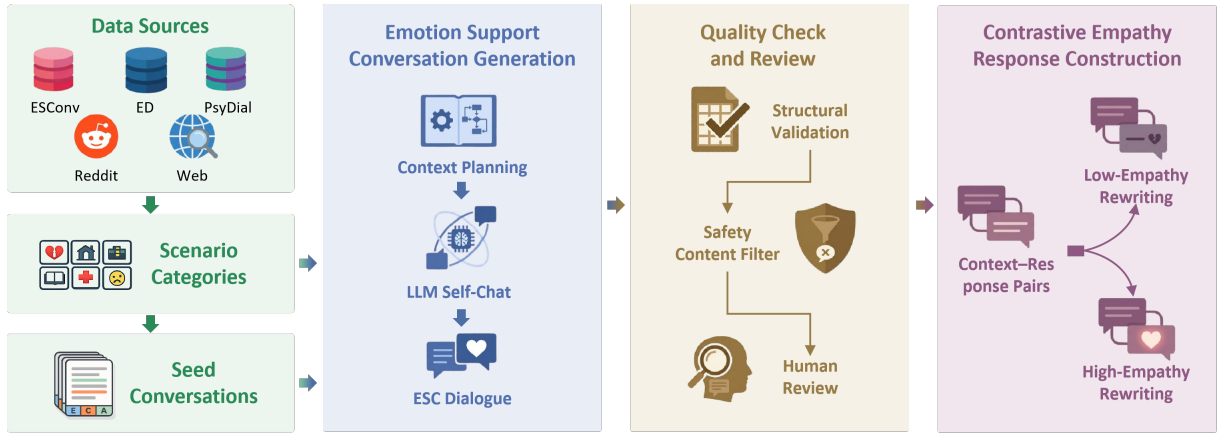


Figure 1: The construction pipeline of the EMPATHY-EVAL dataset

3 Preliminaries

In this section, we construct a evaluation dataset, EMPATHY-EVAL. We then analyze the behavior of existing automatic metrics on this benchmark, highlighting their limitations and motivating the need for our proposed EMPATHYSCORE.

3.1 EMPATHY-EVAL Dataset

Construction Process. The construction of EMPATHY-EVAL proceeds through a rigorous five-stage pipeline designed to ensure both diversity and quality. Initially, we manually curated 60 distinct emotional support scenarios (see Appendix B.4 for specific examples) selected from five diverse data sources including ESConv (Liu et al., 2021), EmpatheticDialogues (Rashkin et al., 2019), PsyDial (Qiu and Lan, 2025), Reddit, and web crawling. For each scenario, we developed high-quality seed dialogues which underwent strict manual review to ensure logical consistency. To scale up the corpus, we extended the methodology of Zheng et al. (2024), employing an LLM-based self-chat framework (details in Appendix B.6) to synthesize large-scale multi-turn dialogues. This generation process incorporated hierarchical control over 16 support strategies and three dialogue states (*Exploration, Comforting, Action*). Subsequently, the synthesized corpus was subjected to automated quality checks, safety filtering, and a final round of human review (see Appendix B.3) to eliminate toxic or incoherent content. Finally, to facilitate fine-grained evaluation, we applied a contrastive rewriting mechanism (details in Appendix B.5) to each helper turn. Using the original response as a reference (R_{ref}), we generated an *Apathetic* version (R_{apa}) via de-empathizing rewriting and an *Empathetic* version

Category	ESConv	ExTES	EMPATHY-EVAL
Dialogues	1,053	11,177	30,633
Utterances	31,410	200,393	522,432
Avg. length of dialogues	29.8	18.2	17.1
Avg. length of utterances	17.8	26.0	28.7
Num. of support strategies	8	16	16
Num. of scenarios	5	36	60

Table 1: Comparison with existing ESC datasets.

Metric (1-3)	Low (R_{apa})	High (R_{emp})	Kappa
Coherence	2.92	2.95	0.78
Helpfulness	1.15	2.82	0.72
Empathy	1.08	2.91	0.75

Table 2: Human evaluation results (scale 1-3). R_{emp} significantly outperforms R_{apa} in empathy and helpfulness.

(R_{emp}) via strategy-altered regeneration.

Data Analysis. Table 1 highlights the scale and diversity of Empathy-Eval. Our dataset contains 30,633 dialogues, which is significantly larger than previous datasets like ESConv (Liu et al., 2021) and ExTES (Zheng et al., 2024). Crucially, it covers 60 diverse emotional scenarios, offering a much broader range of topics than the limited scenarios in baselines. Comprehensive statistical analyses are detailed in Appendix B.

Human Evaluation. To validate data quality, we randomly sampled 500 contrastive pairs (R_{apa} vs. R_{emp}) for human review. Ten annotators rated them on a 1-3 scale (Poor, Fair, Good) based on *Coherence, Helpfulness, and Empathy*. As shown in Table 2, R_{emp} scores nearly perfect in Empathy (2.91) and Helpfulness (2.82), whereas R_{apa} scores very low (≈ 1.1). This confirms that our rewriting method successfully created a clear gap between

high and low empathy. Additionally, both versions maintain high Coherence, ensuring the text is fluent and logical. The high Kappa scores indicate reliable agreement among annotators.

3.2 Limitations of Existing NLP Metrics

To assess the applicability of automated metrics for empathy evaluation, we computed n -gram overlap, semantic similarity, and linguistic diversity on our contrastive dataset (R_{emp} vs. R_{apa}). As shown in Table 3, results reveal a systemic misalignment: Low Empathy responses consistently outperform High Empathy responses across all metric categories. This inverse correlation stems from the fact that standard metrics primarily reward lexical and semantic adherence to the factual reference content. Since empathetic responses often require significant paraphrasing and stylistic shifts rather than mere content preservation, they are penalized by these metrics. Consequently, existing NLP metrics are proven ineffective for quantifying the specific dimension of empathy.

Category	Metric	High (R_{emp})	Low (R_{apa})
n -gram	BLEU-4	0.0285	0.3852
	ROUGE-L	0.1688	0.5660
	METEOR	0.2340	0.5690
Semantic	BERTScore	0.9613	0.9778
Diversity	Dist-1	0.0076	0.0098
	Dist-2	0.0885	0.1231
	Dist-4	0.4618	0.5762

Table 3: Comparison of standard NLP metrics. The metrics consistently assign higher scores to Low Empathy responses, demonstrating an inability to distinguish the empathy dimension.

4 EMPATHYScore Framework

4.1 Problem Formulation

We focus on the Emotional Support Conversation (ESC) task, involving a distressed *seeker* and a supportive *helper*. Following the theoretical framework in ESConv (Liu et al., 2021), the support process is structured hierarchically: fine-grained support strategies (e.g., *Restatement*, *Suggestion*) are organized into three high-level support stages: *Exploration*, *Comforting*, and *Action*. Ideally, a helper should dynamically transition through these stages based on the seeker’s evolving emotional state.

To enhance clarity, key notations are summarized in Table 4.

Symbol	Definition
C	The dialogue context, defined as a sequence of role-utterance pairs
R, \hat{R}	The human reference response and the model-generated response
\mathcal{S}, \mathcal{Z}	Sets of fine-grained support strategies and high-level support stages
p, u	The role indicator (Seeker/Helper) and the utterance content
$s_{\hat{R}}$	The strategy label predicted for response \hat{R}
P_{ideal}	The predicted latent distribution of ideal strategies
S_{flow}	The cognitive suitability score (Metric I)
$S_{content}$	The affective resonance score (Metric II)

Table 4: Summary of notations for the EMPATHYScore framework.

Formally, we define the dialogue context as a sequence of L turns $C = \{(p_1, u_1), (p_2, u_2), \dots, (p_L, u_L)\}$, where $p_i \in \{\text{Seeker}, \text{Helper}\}$ denotes the speaker role and u_i is the textual content. This flexible formulation accounts for irregular turn-taking dynamics (e.g., consecutive messages from the same speaker). Given C , an ESC model aims to generate a response \hat{R} (where $p_{L+1} = \text{Helper}$) that is both structurally appropriate and emotionally resonant.

However, current evaluation protocols are insufficient for this goal. As empirically demonstrated in Table 3, standard NLP metrics exhibit a significant misalignment with human judgment. To bridge this gap, we propose EMPATHYScore, a reference-free, learning-based metric that evaluates \hat{R} from two complementary dimensions:

$$\text{EmpathyScore}(C, \hat{R}) = \alpha \cdot S_{\text{flow}} + (1 - \alpha) \cdot S_{\text{content}} \quad (1)$$

where S_{flow} and S_{content} are normalized scores in $[0, 1]$. The hyper-parameter $\alpha \in [0, 1]$ controls the trade-off between the cognitive (flow) and affective (content) dimensions. In this work, we set $\alpha = 0.5$, assuming that appropriate strategy selection and emotional resonance are equally critical.

- S_{flow} evaluates **Cognitive Suitability** (reflecting strategy flow). It measures the **probabilistic validity** of the helper’s strategy against the ideal support flow derived from the full interaction history C . This captures the cognitive aspect of empathy, ensuring stage-wise adaptation to the seeker’s evolving state while

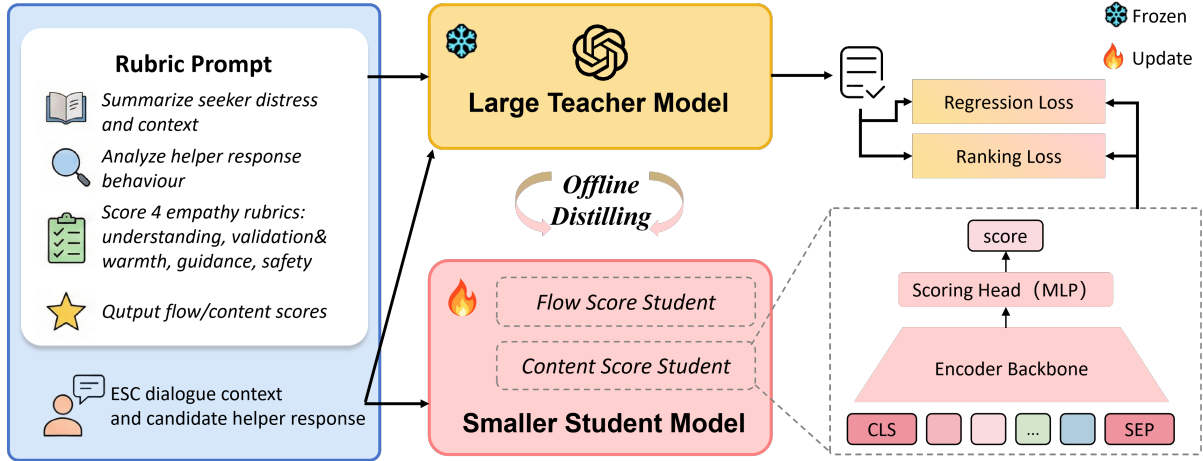


Figure 2: he overview of the DISTILL-ES framework.

allowing for diverse strategic choices (Hill, 2014; Clark, 2010).

- S_{content} evaluates **Affective Resonance** (reflecting content quality). It leverages a **contrastive learning objective** to assess whether the generated content structurally entails the seeker’s emotional distress. This metric captures the intrinsic warmth and affective understanding of the response, aligning with person-centered message theory (High and Dillard, 2012).

4.2 The DISTILL-ES Framework

Recent advances in LLM distillation show that rich supervision from large teachers can be compressed into compact students for instruction following, automatic evaluation, and dialogue understanding (Yue et al., 2024; Oh et al., 2025; Lee et al., 2024; Wei et al., 2025). Following this line of work, we cast EMPATHYScore as a learned, LLM-distilled metric and propose the DISTILL-ES (Distilled EMPATHYScore) framework, which transfers empathy judgements from a strong teacher LLM into lightweight scorers specialized for emotional support conversations (ESC), as illustrated in Figure 2.

Teacher rubric prompting. For each ESC dialogue context C and candidate helper response \hat{R} , we query a frozen teacher LLM \mathcal{T} with a rubric-style prompt (Figure 2). The prompt asks \mathcal{T} to (i) summarize the seeker’s distress and conversational context, (ii) analyze the helper’s behaviour, and (iii) assign scalar empathy scores in $[0, 1]$ for two dimensions: cognitive suitability of the support

flow $S_{\text{flow}}^T(C, \hat{R})$ and affective content resonance $S_{\text{content}}^T(C, \hat{R})$. In addition, the teacher provides four rubric-level scores that refine content empathy (empathic understanding, validation & warmth, guidance and actionability, safety), and, for a subset of contexts, pairwise preferences indicating which of two responses is more empathic. All teacher queries are collected offline to form a distillation dataset; at test time the teacher is never used.

Student architecture. We instantiate two student scorers, a *Flow Score Student* f_{θ}^{flow} and a *Content Score Student* $f_{\theta}^{\text{content}}$. Both share the same RoBERTa-style encoder backbone followed by a shallow MLP scoring head (the right panel in Figure 2). Given (C, \hat{R}) , each student produces a scalar score in $[0, 1]$:

$$S_{\text{flow}}(C, \hat{R}) = f_{\theta}^{\text{flow}}(C, \hat{R}), \quad (2)$$

$$S_{\text{content}}(C, \hat{R}) = f_{\theta}^{\text{content}}(C, \hat{R}). \quad (3)$$

which are combined into the final metric according to Eq. (1). Thus DISTILL-ES implements EMPATHYScore as two compact, reference-free scoring models that can be applied to any ESC system output.

Distillation objective. Each student is trained to match the teacher’s scalar scores while preserving the teacher’s relative preferences between candidate responses. For a single scored example (C, \hat{R}) with teacher score $S_d^T(C, \hat{R})$ on dimension $d \in \{\text{flow}, \text{content}\}$, we use a regression loss

$$\mathcal{L}_{\text{reg}}^{(d)} = (S_d(C, \hat{R}) - S_d^T(C, \hat{R}))^2, \quad (4)$$

where S_d denotes the corresponding student score (S_{flow} or S_{content}). For ranking supervision, we con-

struct response pairs (\hat{R}^+, \hat{R}^-) under the same context C where the teacher prefers \hat{R}^+ to \hat{R}^- , and apply a margin-based ranking loss

$$\mathcal{L}_{\text{rank}}^{(d)} = \max(0, m - S_d(C, \hat{R}^+) + S_d(C, \hat{R}^-)), \quad (5)$$

with a fixed margin $m > 0$. The distillation loss for dimension d is

$$\mathcal{L}_{\text{KD}}^{(d)} = \mathcal{L}_{\text{reg}}^{(d)} + \lambda_{\text{rank}} \mathcal{L}_{\text{rank}}^{(d)}, \quad (6)$$

and the overall training objective sums the flow and content components:

$$\mathcal{L}_{\text{KD}} = \mathcal{L}_{\text{KD}}^{(\text{flow})} + \mathcal{L}_{\text{KD}}^{(\text{content})}. \quad (7)$$

In all experiments, DISTILL-ES is trained purely with this offline distillation objective, without access to human empathy labels, and can therefore be scaled to large amounts of unlabeled ESC data.

5 Experiments

In this section, we conduct comprehensive experiments to evaluate EMPATHYSCORE against state-of-the-art baselines. We aim to answer the following research questions:

- **RQ1:** Can EMPATHYSCORE effectively distinguish between high-quality empathetic responses and low-quality ones?
- **RQ2:** How do the proposed components (S_{flow} and S_{content}) contribute to the overall performance?

5.1 Experimental Setup

5.1.1 Datasets

To validate the effectiveness of our framework, we utilize the EMPATHY-EVAL dataset. This dataset was constructed by applying our Teacher Rubric Prompting (using GPT-5) to a corpus of emotional support conversations, generating high-quality annotations that include ideal strategy distributions, scalar empathy scores, and pairwise preference rankings.

Consistent with standard supervised learning paradigms, we randomly split the EMPATHY-EVAL dataset into two subsets:

- **Training Set (80%):** Used to optimize the student scorers. The student models leverage the scalar scores for the regression loss and the pairwise preferences for the ranking loss to mimic the teacher’s judgment logic.

- **Test Set (20%):** A held-out subset used exclusively for evaluation. This subset contains distinct dialogue contexts not seen during training, allowing us to assess the metric’s generalization capability.

5.1.2 Baselines

We compare EMPATHYSCORE against widely used reference-based metrics in dialogue generation. These metrics evaluate performance by measuring the similarity between the system-generated response and human gold references:

- **Word-Overlap Metrics:** We report BLEU-4, ROUGE-L, and METEOR. These metrics evaluate performance by measuring the n-gram precision and recall against human references.
- **Embedding-based Metrics:** We report BERTScore, which computes semantic similarity.
- **Diversity Metrics:** We include Dist-1, Dist-2, and Dist-4, which measure the ratio of distinct uni-grams and bi-grams in the generated responses. While primarily used for generation diversity, they serve as a baseline to check if empathy correlates with lexical richness.

5.1.3 Implementation Details

We employ RoBERTa-base as the shared encoder backbone for both the flow and content student scorers. The scoring head consists of a lightweight two-layer MLP with a hidden size of 256. We train the models for 5 epochs using the AdamW optimizer with a learning rate of $2e-5$ and a batch size of 32. For the hybrid objective, the ranking margin is set to $m = 0.1$, and the ranking loss weight is set to $\lambda_{\text{rank}} = 0.5$. All experiments were conducted on a single NVIDIA A100 GPU.

5.2 Evaluation of Distinguishability (RQ1)

To answer **RQ1**, we rigorously analyze whether EMPATHYSCORE can effectively distinguish between high-quality empathetic responses (R_{emp}) and low-quality generic ones (R_{apa}). This comparison against standard metrics is critical, as a reliable metric must not only correlate with human judgment but also possess sufficient discriminative power to penalize sub-optimal outputs. The comparative results are presented in Table 5.

Table 5: Mean scores on High-Quality (R_{emp}) vs. Low-Quality (R_{apa}) subsets. **Bold** indicates the higher score between the two groups. Standard metrics consistently exhibit an *Evaluation Inversion* (scoring Low > High), whereas EMPATHYSCORE correctly aligns with human preference.

Metric	High (R_{emp})	Low (R_{apa})
<i>n-gram Metrics</i>		
BLEU-4	0.0285	0.3852
ROUGE-L	0.1688	0.5660
METEOR	0.2340	0.5690
<i>Semantic Metrics</i>		
BERTScore	0.9613	0.9778
<i>Diversity Metrics</i>		
Dist-1	0.0076	0.0098
Dist-2	0.0885	0.1231
Dist-4	0.4618	0.5762
<i>Ours</i>		
EMPATHYSCORE	0.8420	0.3250

Results and Analysis As illustrated in Table 5, we observe a distinct and systematic divergence in how different metrics rank the two contrastive groups. The results highlight the fundamental limitations of current automated evaluation methods when applied to the open-ended domain of Emotional Support Conversation (ESC).

1. The Phenomenon of Evaluation Inversion in Standard Metrics: Standard automated metrics exhibit a severe misalignment with human judgment, a phenomenon we term *Evaluation Inversion*. Specifically, *n*-gram metrics (BLEU-4, ROUGE-L, METEOR) assign significantly higher scores to the Low-Quality group (R_{apa}) than to the High-Quality group (R_{emp}). For instance, BLEU-4 scores generic responses nearly 13 times higher (0.3852) than specific empathetic responses (0.0285).

This counter-intuitive result primarily stems from the "One-to-Many" nature of open-ended dialogue and the "safe response" bias. In the ESC task, for any given emotional query, there exists a vast space of valid, high-empathy responses that differ significantly in vocabulary and structure. However, generic utterances (e.g., "I am sorry to hear that", "That sounds tough") appear with disproportionately high frequency across the training corpus. Consequently, traditional metrics, which rely on surface-level lexical matching, are statistically bi-

ased towards these "average" responses. A model that generates these safe, low-information utterances maximizes statistical likelihood and overlap with the references, despite lacking genuine emotional insight. Conversely, high-quality empathetic responses are often tailored to specific details of the user's distress, resulting in unique formulations that rarely match the rigid surface forms of the ground truth references.

Furthermore, even advanced semantic-based metrics like BERTScore fail to rectify this bias, favoring the generic group (0.9778) over the empathetic one (0.9613). This observation reveals a critical theoretical limitation: high semantic similarity to a reference does not equate to high pragmatic quality. BERTScore effectively captures topical relevance but struggles to distinguish between "polite acknowledgement" and "deep emotional validation." A response can be semantically close to a reference (sharing the same topic embedding space) but completely lack the specific supportive intent required for effective therapy-like interactions. This confirms that neither lexical overlap nor embedding-based similarity is sufficient to capture the dynamic and nuanced nature of empathy.

2. Effectiveness and Robustness of EMPATHYSCORE: In sharp contrast to the baselines, EMPATHYSCORE successfully corrects this bias and demonstrates superior distinguishing capability. As shown in the bottom row of Table 5, our metric assigns a significantly higher score to the High-Quality group (**0.8420**) compared to the Low-Quality group (0.3250).

This result highlights two key strengths of our proposed approach. First, it demonstrates robustness against the "safe response" trap. Unlike probability-based models that reward the most frequent patterns, EMPATHYSCORE is trained to recognize specific empathetic features, ensuring that generic phrases—no matter how grammatically correct—receive lower scores due to their lack of emotional substance. Second, the substantial margin between the two groups indicates that our model effectively aligns with human cognitive intuition. Human evaluators prioritize the functional aspect of support (e.g., "Did the system really understand me?") over linguistic fluency. By evaluating the intrinsic quality of the response—considering both its emotional content and interaction flow—EMPATHYSCORE correctly identifies and rewards the nuanced emotional support present in R_{emp} . This capability to discrim-

Table 6: Ablation study results. We report the Pearson (r) and Spearman (ρ) correlation with human judgments. " Δ " denotes the performance drop compared to the full model. The results show that both components are essential, with S_{content} contributing primarily to semantic accuracy and S_{flow} enhancing structural coherence.

Model Variants	Correlation		Performance Drop	
	Pearson (r)	Spearman (ρ)	Δr	$\Delta \rho$
Full Model (EMPATHYSCORE)	0.684	0.672	-	-
w/o Flow Component (S_{flow})	0.592	0.585	↓ 13.5%	↓ 12.9%
w/o Content Component (S_{content})	0.415	0.398	↓ 39.3%	↓ 40.8%

inate between "safe" and "supportive" is crucial for training future ESC systems, answering **RQ1** affirmatively.

5.3 Ablation Study (RQ2)

To answer **RQ2**, we conducted a comprehensive ablation study to investigate the individual contributions of the two core components of EMPATHYSCORE: the Flow-based Score (S_{flow}) and the Content-based Score (S_{content}). By systematically removing one component at a time, we evaluated the performance of model variants using Pearson (r) and Spearman (ρ) correlation coefficients with human ratings on the test set. This analysis allows us to understand how each module influences the overall evaluation capability and whether the proposed hybrid architecture is theoretically justified.

Results and Analysis Table 6 summarizes the contribution of each component. The empirical results clearly indicate that while both components positively contribute to the final performance, they play distinct and complementary roles in the evaluation framework, validating our dual-stream hypothesis:

1. Impact of Empathetic Content (S_{content}):

The variant removing the content component ("w/o Content") resulted in the most significant performance degradation, with the Pearson correlation dropping drastically by 39.3% (from 0.684 to 0.415). This finding confirms that the semantic content—specifically the linguistic realization of empathy—is the foundational element of the evaluation system. Empathy is fundamentally communicated through words that convey understanding, validation, and care. Specifically, this component captures the fine-grained semantics of emotional support, such as the use of specific emotional keywords, supportive phrasing, and validation techniques that resonate with the user’s state. Without the ability to assess *what* is explicitly said (S_{content}),

the model loses its capacity to detect the "warmth" and "understanding" inherent in the dialogue. Relying solely on interaction flow (S_{flow}) proves insufficient because a conversation might follow a correct structural transition pattern (e.g., asking a question after a statement) but still lack the necessary emotional resonance or compassionate tone to be truly supportive.

2. Impact of Dialogue Flow (S_{flow}): Removing the flow component ("w/o Flow") also leads to a notable decrease in performance, with a 13.5% drop in Pearson correlation. While S_{content} is critical for evaluating the semantic dimension, the addition of S_{flow} provides crucial context-aware regularization. Empathy in ESC is not a static property but a dynamic, temporal process that requires appropriate progression. Effective support typically follows a logical sequence: listening and validating first, then exploring the problem, and finally offering advice. S_{flow} acts as a safeguard against "contextual mismatch," penalizing responses that may be semantically polite but are contextually abrupt or out of order—such as rushing to give practical advice before sufficiently establishing an emotional connection. This drop in performance without S_{flow} suggests that structural coherence is a vital multiplier for the overall quality of empathy; a well-worded response delivered at the wrong stage of the conversation is often perceived as insensitive.

Takeaway The full model achieves the highest correlation (**0.684**), significantly outperforming either single-component variant. This demonstrates that EMPATHYSCORE effectively combines the semantic precision of S_{content} with the structural awareness of S_{flow} . The two components work in synergy: the content module evaluates the *substance* of the response (the "what"), while the flow module evaluates its *timing and appropriateness* (the "when"). This ablation study strongly supports the necessity of the proposed hybrid architecture, proving that accurately capturing the complexity of human empathy requires modeling both the linguistic features and the interaction dynamics.

6 Conclusion

We presented EMPATHYSCORE, a reference-free metric for Emotional Support Conversations that explicitly decomposes empathy into cognitive suitability of support strategy flow and affective resonance of response content, together with the DISTILL-ES framework for learning this metric.

626 Limitations

627 While EMPATHYSCORE and the DISTILL-ES
628 framework take a step toward more faithful evaluation
629 of emotional support conversations, several
630 limitations remain. First, our metric ultimately inherits
631 the biases and blind spots of the teacher LLM
632 and our rubric design: the “ground truth” supervision
633 comes from a single model family, a particular
634 set of prompts, and a fixed four-dimensional
635 rubric. As a result, the distilled scorers may encode
636 stylistic or cultural preferences of the teacher
637 (e.g., favouring certain phrasing patterns or support
638 styles) and may not fully reflect diverse human
639 judgments about what counts as empathic, especially
640 across different cultures, demographics, or
641 counselling traditions. Second, our experiments
642 are conducted on English text-based ESC datasets
643 in non-clinical settings; it is unclear how well
644 EMPATHYSCORE transfers to other languages, modalities
645 (e.g., voice-based support with prosody), or
646 high-stakes clinical contexts where professional
647 standards for empathy and safety are stricter.

648 Methodologically, DISTILL-ES focuses on off-
649 line correlation with human and teacher scores
650 and does not guarantee that higher scores causally
651 imply better outcomes for seekers (e.g., improved
652 mood or reduced distress). Our margin-based ranking
653 loss relies on teacher-generated pairwise preferences,
654 which may contain label noise, and we have not
655 systematically analysed failure cases where the
656 metric disagrees with expert counsellors. In addition,
657 our current implementation uses RoBERTa-style
658 encoders of moderate size; we do not explore
659 the trade-off between model capacity and evaluation
660 quality, nor do we study robustness under domain
661 shift or adversarial prompting. Finally, although
662 we discuss potential applications of EMPATHYSCORE
663 for model selection and benchmarking, we do not
664 investigate its impact when used as a training-time
665 objective to optimise generative ESC models, nor
666 the ethical implications of deploying such a metric
667 in real-world mental health systems (e.g., unintended
668 over-reliance on automated scores for triage or
669 quality control). These limitations highlight important
670 directions for future work.

671 References

672 Muhammad Abdul-Mageed, Anneke Buffone, Hao
673 Peng, Salvatore Giorgi, Johannes Eichstaedt, and
674 Lyle Ungar. 2017. [Recognizing pathogenic empathy in social media](#). *Proceedings of the Interna-*

tional AAAI Conference on Web and Social Media,
11(1):448–451.

675 Simon Baron-Cohen and Sally Wheelwright. 2004. [The empathy quotient: An investigation of adults with asperger syndrome or high functioning autism, and normal sex differences](#). *Journal of Autism and Developmental Disorders*, 34(2):163–175.

676 Sven Buechel, Anneke Buffone, Barry Slaff, Lyle Ungar,
677 and João Sedoc. 2018. [Modeling empathy and distress in reaction to news stories](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4758–4765, Brussels, Belgium. Association for Computational Linguistics.

678 Hyunjoo Chae, Yongho Song, Kai Ong, Taeyoon
679 Kwon, Minjin Kim, Youngjae Yu, Dongha Lee,
680 Dongyeop Kang, and Jinyoung Yeo. 2023. [Dialogue chain-of-thought distillation for commonsense-aware conversational agents](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5606–5632, Singapore. Association for Computational Linguistics.

681 Zhuang Chen, Yaru Cao, Guanqun Bi, Jincenzi Wu,
682 Jinfeng Zhou, Xiyao Xiao, Si Chen, Hongning
683 Wang, and Minlie Huang. 2025. [Socialsim: towards socialized simulation of emotional support conversation](#). In *Proceedings of the Thirty-Ninth AAAI Conference on Artificial Intelligence and Thirty-Seventh Conference on Innovative Applications of Artificial Intelligence and Fifteenth Symposium on Educational Advances in Artificial Intelligence*, AAAI’25/IAAI’25/EAAI’25. AAAI Press.

684 Jiale Cheng, Sahand Sabour, Hao Sun, Zhuang Chen,
685 and Minlie Huang. 2023. [PAL: Persona-augmented emotional support conversation generation](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 535–554, Toronto, Canada. Association for Computational Linguistics.

686 Arthur J. Clark. 2010. [Empathy: An integral model in the counseling process](#). *Journal of Counseling & Development*, 88(3):348–356.

687 Rupert Croxford and 1 others. 2025. Current and future
688 state of evaluating large language models for medical
689 summarization tasks. *npj Digital Medicine*. To
690 appear.

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

692 Mark H. Davis, Miles P. Davis, M Davis, Matthew
693 Davis, Mark Davis, Mm Davis, M Davis, F. Caroline
694 Davis, Heather A. Davis, and Ilus W. Davis. 1980. [A multidimensional approach to individual differences in empathy](#).

695 Jan T. Deriu, Álvaro Rodrigo, Arantxa Otegi, and 1 others. 2021. A survey on evaluation methods for dialogue systems. *Artificial Intelligence Review*, 54:755–810.

844	empathy questionnaire: Scale development and initial validation of a factor-analytic solution to multiple empathy measures. <i>Journal of Personality Assessment</i> , 91(1):62–71. PMID: 19085285.	899
845		900
846		901
847		902
848	Quan Tu, Yanran Li, Jianwei Cui, Bin Wang, Ji-Rong Wen, and Rui Yan. 2022. MISC: A mixed strategy-aware model integrating COMET for emotional support conversation . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 308–319, Dublin, Ireland. Association for Computational Linguistics.	903
849		904
850		905
851		906
852		907
853		908
854		909
855		
856	Hongru Wang, Rui Wang, Fei Mi, Yang Deng, Zezhong Wang, Bin Liang, Ruifeng Xu, and Kam-Fai Wong. 2023. Cue-CoT: Chain-of-thought prompting for responding to in-depth dialogue questions with LLMs . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 12047–12064, Singapore. Association for Computational Linguistics.	910
857		911
858		912
859		913
860		914
861		915
862		916
863	Yuting Wei, Qi Meng, Yuanxing Xu, and Bin Wu. 2025. TEACH: A contrastive knowledge adaptive distillation framework for classical Chinese understanding . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3537–3550, Vienna, Austria. Association for Computational Linguistics.	917
864		918
865		919
866		920
867		921
868		
869		
870	Jiahao Yuan, Zixiang Di, Zhiqing Cui, Guisong Yang, and Usman Naseem. 2025. ReflectDiffu: Reflect between emotion-intent contagion and mimicry for empathetic response generation via a RL-diffusion framework . In <i>Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 25435–25449, Vienna, Austria. Association for Computational Linguistics.	922
871		923
872		924
873		925
874		926
875		927
876		
877		
878		
879	Yuanhao Yue, Chengyu Wang, Jun Huang, and Peng Wang. 2024. Distilling instruction-following abilities of large language models with task-aware curriculum planning . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 6030–6054, Miami, Florida, USA. Association for Computational Linguistics.	928
880		929
881		930
882		931
883		932
884		933
885		934
886		935
887		936
888		937
889		938
890		939
891		940
892		
893		
894	Yifan Zhang and 1 others. 2024. FEEL: A benchmark for evaluating emotional support in mental health dialogues . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> .	941
895		942
896		943
897		944
898		945
899		
900		
901		
902		
903	Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. 2023. AugESC: Dialogue augmentation with large language models for emotional support conversation . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 1552–1568, Toronto, Canada. Association for Computational Linguistics.	946
904		947
905		
906		
907		
908		
909		
910		
911		
912		
913		
914		
915		
916		
917		
918		
919		
920		
921		
922		
923		
924		
925		
926		
927		
928		
929		
930		
931		
932		
933		
934		
935		
936		
937		
938		
939		
940		
941		
942		
943		
944		
945		
946		
947		

ID	Scenario Name	Prop.	ID	Scenario Name	Prop.
1	Breakups or Divorce	4.02%	31	Coping with a Diagnosis or Medical Treatment	1.50%
2	Conflicts or Communication Problems	3.97%	32	Finding Meaning and Purpose in Life	1.50%
3	Work-related Stress and Burnout	3.96%	33	Career Transitions	1.50%
4	Coping with the Death of a Loved One	3.95%	34	Support for Loved Ones or Friends	1.50%
5	Depression and Low Mood	3.52%	35	News overload and anxiety about world events	1.49%
6	Financial Worries and Uncertainty	3.44%	36	Coping with natural disasters or community-wide crises	1.49%
7	Anxiety and Panic	3.44%	37	Surviving and Recovering from Physical or Emotional Abuse	1.49%
8	Chronic loneliness and lack of social support	2.98%	38	Housing instability or fear of eviction	1.25%
9	Parenthood and Parenting Challenges	2.54%	39	Cultural Identity and Belonging	1.23%
10	Academic Stress or Pressure	2.50%	40	Job Loss or Career Setbacks	1.08%
11	Low Self-Esteem or Lack of Confidence	2.48%	41	Long-distance relationship stress	1.03%
12	Romantic loneliness and difficulty finding a partner	2.47%	42	Addiction and Recovery	1.02%
13	Chronic Illness or Pain Management	2.47%	43	Healing from Sexual Assault or Domestic Violence	1.01%
14	Chronic sleep problems and insomnia	2.03%	44	Spirituality and Faith	0.76%
15	Sibling Rivalry or Family Conflict	2.02%	45	Infertility and challenges with trying to conceive	0.75%
16	Moving to a New City or Country	2.01%	46	Empty nest and adjusting after children leave home	0.75%
17	Dealing with the Loss of a Pet	2.00%	47	Legal troubles and court cases	0.75%
18	Caregiver Support	1.99%	48	Gaming or internet addiction impacting daily life	0.75%
19	Parenting Challenges and Parental Guilt	1.99%	49	Post-Traumatic Stress Disorder (PTSD)	0.74%
20	Communication Challenges	1.98%	50	Religious or worldview conflicts with family or community	0.74%
21	Work-life balance and role overload	1.98%	51	Healing from Abuse	0.74%
22	Perfectionism, procrastination, and fear of failure	1.98%	52	Online harassment and cyberbullying	0.74%
23	Drifting apart from long-term friends or social circles	1.53%	53	Navigating Gender Identity and Transitioning	0.74%
24	Uncertainty about college major or future career path	1.52%	54	Transition to college or moving away from home	0.53%
25	Infidelity and rebuilding trust after cheating	1.52%	55	Managing Bipolar Disorder	0.53%
26	Social media overuse and social comparison	1.52%	56	Midlife crisis and career plateau	0.52%
27	Adjusting to a New Job or Role	1.52%	57	Health anxiety and fear of serious illness	0.52%
28	Experiencing discrimination and microaggressions	1.51%	58	Workplace bullying or toxic team dynamics	0.52%
29	Body Image Concerns and Eating Disorders	1.51%	59	Unemployment-related Stress	0.51%
30	LGBTQ+ Identity	1.51%	60	Retirement adjustment and loss of professional identity	0.51%

Table 7: Distribution of all 60 emotional support scenarios in EMPATHY-EVAL, sorted by frequency proportion. The list is split into two columns for layout efficiency.

Empathy flow (Left) exhibits a clear "Exploration-Comforting-Action" progression, mirroring effective counseling sessions. In contrast, the Low Empathy flow (Middle) remains static. The Radar Chart (Right) further confirms that R_{emp} covers a broad spectrum of supportive strategies, whereas R_{apa} is constricted to neutral responses.

Linguistic and Interaction Depth. Figure 4 quantifies the linguistic divergence:

- *Length Distribution (Left):* The distribution for R_{emp} is right-skewed (peaking ≈ 35 words), indicating more elaborate supportive language compared to R_{apa} .
- *Strategy Divergence (Middle):* The diverging chart highlights the "Empathy Delta." R_{emp} prioritizes constructive strategies like *Emotional Validation* (green bars), while R_{apa} relies on passive strategies like *Clarification* (red bars).
- *Depth Dynamics (Right):* The line plot shows that R_{emp} sustains deep engagement throughout the dialogue, whereas R_{apa} shows a flat trajectory.

B.3 Human Review and Quality Assurance

To ensure the reliability and safety of the dataset, the synthesized corpus was subjected to a rigorous quality control pipeline. Following automated quality checks and safety filtering, we conducted a final round of human review to explicitly eliminate toxic, incoherent, or strategy-mismatched content.

We recruited 10 annotators with expertise in linguistic evaluation, compensated at a rate of \$5 USD per hour. Prior to the task, we conducted a mandatory training session detailing the workflow and annotation criteria. The full text of the instructions, along with screenshots of the annotation interface, is provided in Appendix ???. Crucially, given the sensitive nature of emotional support dialogue, we issued a clear disclaimer regarding potential psychological risks associated with reading distressing content. All participants provided informed consent, acknowledging these risks before proceeding. To ensure consistency, annotators were required to analyze 10 distinct case studies and pass a pilot qualification test aligned with expert standards.

To streamline this large-scale review, we developed a customized auditing tool, the "Empathy-Eval Auditor" (see Figure 5). The interface presents the complete dialogue history alongside the specific situation description, providing the necessary context for adjudicating coherence

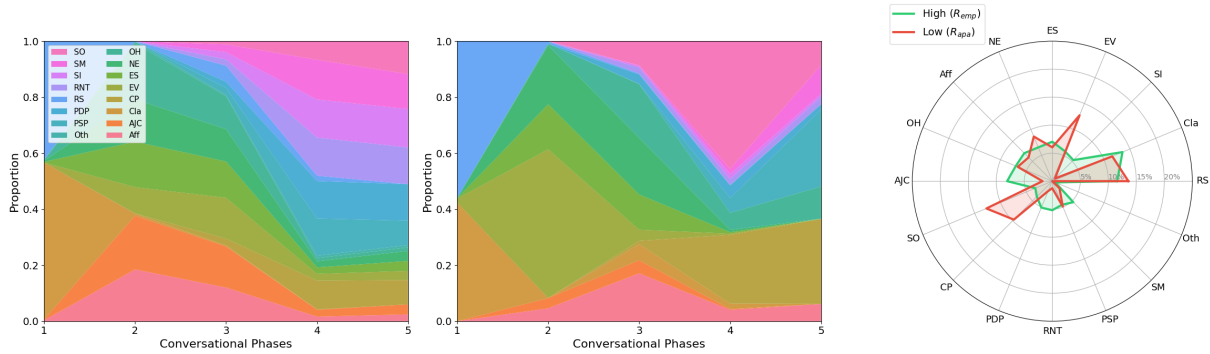


Figure 3: Visualization of strategy dynamics. (Left & Middle) High empathy responses follow a logical progression compared to static low empathy responses. (Right) The radar chart highlights the superior strategy diversity of R_{emp} .

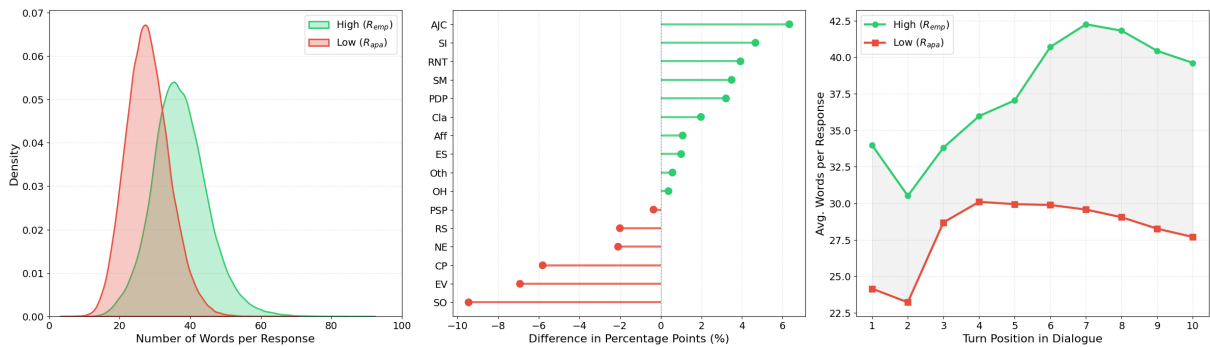


Figure 4: Advanced statistical comparison. (Left) Response length distribution. (Middle) Strategy divergence between high and low empathy. (Right) Turn-level conversation depth dynamics.

and safety. Using a binary decision mechanism (“Keep” or “Remove”), auditors verified whether the model’s response accurately reflected the assigned strategy tag (e.g., *Emotional Validation*). This process ensured that only valid, high-quality dialogues were retained in the final EMPATHY-EVAL corpus.

B.4 Scenario Examples

Table 7 listed the distribution of our 60 emotional support scenarios. Below, we provide concrete dialogue examples for each scenario to illustrate the diversity of user situations covered in EMPATHY-EVAL.

Academic Stress or Pressure *Ex 1:* I studied so hard for this exam, but when I got my grade back, it was way lower than I thought it would be. *Ex 2:* I can’t believe my teacher marked me wrong on that math test. It was my final exam and I studied for a week. *Ex 3:* I can’t believe I did so badly on that test. I studied for hours, but it all just slipped away during the exam.

Addiction and Recovery *Ex 1:* I’ve been feeling really down, and it seems like my family has basically given up on me. *Ex 2:* I messed up. I drank again after six months of being clean, and now I feel like I’ve ruined everything. *Ex 3:* I slipped up last week after almost a

year clean, and I just feel like I’ve let everyone down and myself too.

Adjusting to a New Job or Role *Ex 1:* It’s only been a few weeks since I started managing this team, but I already feel like I’m way in over my head. *Ex 2:* I just started this new role as a team lead last week, and honestly, I already feel like I’m in way over my head. *Ex 3:* I just began this new lead position, and honestly, it feels like I’m constantly behind with no clear way to catch up.

Anxiety and Panic *Ex 1:* Every time I try to drive during rush hour, my heart starts racing and I feel this wave of panic like I can’t breathe properly. *Ex 2:* Lately, whenever I go to family events, I feel this sudden wave of panic that makes my heart race and my thoughts spiral out of control. *Ex 3:* Lately, every time I have to talk to someone at work, I feel this rush of panic. My heart races and I get so dizzy I just want to hide away.

Body Image Concerns and Eating Disorders *Ex 1:* Hi... lately I just feel really down when I look in the mirror. *Ex 2:* Hey, I just can’t stop feeling awful about how I look these days. *Ex 3:* Hey, I’ve been feeling really frustrated and ashamed of how I look lately.

Breakups or Divorce *Ex 1:* It’s not a great day. I feel pretty awful, honestly. *Ex 2:* My girlfriend just left me. I feel awful, like my whole world collapsed. *Ex 3:* I honestly don’t even know where to start. Everything feels shattered right now.

Career Transitions *Ex 1:* I just quit my teaching job to become a graphic designer, but now I’m really scared

Figure 5: Screenshot of the Empathy-Eval Auditor interface used for the final human review. Annotators utilize this tool to audit the coherence, safety, and strategy alignment of the generated dialogues.

1053	I made a huge mistake. <i>Ex 2:</i> I quit my stable job	matter how tired I am. <i>Ex 3:</i> Lately, I just can't seem to	1080
1054	last month to pursue photography, but now I'm second-	get to sleep at night, no matter how tired I am.	1081
1055	guessing if I did the right thing. <i>Ex 3:</i> I've been offered		
1056	a promotion at work, but honestly, I feel really over-	Communication Challenges <i>Ex 1:</i> I usually keep my	1082
1057	whelmed and unsure if I want to accept it.	thoughts to myself because I'm scared that bringing	1083
		them up will just lead to fighting. <i>Ex 2:</i> I've been scared	1084
1058	Caregiver Support <i>Ex 1:</i> Hello, I was wondering if you	to tell my friend how I really feel because I don't want	1085
1059	could help me with something I'm really scared about.	to upset them or start an argument. <i>Ex 3:</i> Lately, I just	1086
1060	<i>Ex 2:</i> Hey, I'm really struggling with everything that's	don't say what's really on my mind because I'm afraid	1087
1061	happening with my grandma. It's getting overwhelming.	it might upset people or cause arguments.	1088
1062	<i>Ex 3:</i> Hey, I'm really overwhelmed right now and don't		
1063	know how to handle all this caregiving stuff for my dad.	Conflicts or Communication Problems <i>Ex 1:</i> I really dis-	1089
		like my neighbours. They argue loudly all the time, and	1090
1064	Chronic Illness or Pain Management <i>Ex 1:</i> Lately, my en-	it's wearing me down. <i>Ex 2:</i> I swear every chat with my	1091
1065	ergy's just disappeared. Even getting dressed feels like	partner ends up tense or with us not really hearing each	1092
1066	a huge effort, and it's really wearing me down. <i>Ex 2:</i>	other anymore. <i>Ex 3:</i> I just can't seem to get through	1093
1067	These past few weeks, my joints have been so painful	to my partner without things turning into a fight. It's	1094
1068	and stiff that even getting out of bed feels like a huge	exhausting.	1095
1069	challenge. <i>Ex 3:</i> Lately, my symptoms have gotten so		
1070	much worse. I'm feeling shaky and tired all the time,	Coping with a Diagnosis or Medical Treatment <i>Ex 1:</i> I	1096
1071	and it's really messing up my days.	just got my diagnosis a few days ago and honestly, I	1097
		feel lost and scared about what's going to happen to me.	1098
1072	Chronic loneliness and lack of social support <i>Ex 1:</i> Some-	<i>Ex 2:</i> I just found out I have early-stage Parkinson's,	1099
1073	times I feel like I have nobody to talk to anymore. <i>Ex 2:</i>	and honestly, I'm terrified about what this means for my	1100
1074	I just feel so invisible here, like no one really sees me for	future. <i>Ex 3:</i> I just got diagnosed with Parkinson's, and	1101
1075	who I am. <i>Ex 3:</i> Lately, I just feel like I'm completely	honestly, I'm really scared about what this means for	1102
1076	invisible to everyone around me.	me going forward.	1103
1077	Chronic sleep problems and insomnia <i>Ex 1:</i> Hi, I'm so	Coping with natural disasters or community-wide crises	1104
1078	tired. I feel like I haven't truly slept in ages. <i>Ex 2:</i>	<i>Ex 1:</i> I am so frustrated. My school was closed without	1105
1079	I'm really struggling to get any sleep these days, no	warning because of COVID. <i>Ex 2:</i> The flood came so	1106

1107	suddenly, and we barely had time to grab anything	wandering without any clue who I am or what I want	1171
1108	before leaving the house. <i>Ex 3:</i> It's been really hard to	anymore.	1172
1109	focus since the tornado hit. Everything feels so chaotic		
1110	and unsafe now.		
1111	Coping with the Death of a Loved One <i>Ex 1:</i> I still can't	Gaming or internet addiction impacting daily life <i>Ex 1:</i>	1173
1112	believe my sister is gone. It feels so unreal and I'm	I've been playing games almost all night lately, and	1174
1113	drowning in sadness. <i>Ex 2:</i> My dad passed away a few	it's really messing up my sleep and making me fall	1175
1114	days ago, and honestly, I feel kind of lost and numb	behind at work. <i>Ex 2:</i> I've been gaming almost every	1176
1115	right now. <i>Ex 3:</i> It's been two weeks since my dad	night lately, sometimes until 3 or 4 a.m., and it's really	1177
1116	passed away, and I just feel completely lost most days.	messing up my sleep and focus at work. <i>Ex 3:</i> Lately,	1178
		I've been gaming almost every night until dawn, and	1179
		it's wrecking my sleep and making me miss important	1180
		deadlines at work.	1181
1117	Cultural Identity and Belonging <i>Ex 1:</i> Since I came back	Healing from Abuse <i>Ex 1:</i> I still feel really anxious when-	1182
1118	home a few months ago, I just can't shake this feeling	ever my mom calls me, even though I've been setting	1183
1119	that I don't really belong anywhere now. <i>Ex 2:</i> I've been	boundaries for months now. <i>Ex 2:</i> I've been avoiding	1184
1120	back here for a few months, but instead of feeling at	calls from my dad because they just bring up so much	1185
1121	home, I just feel like I don't belong anywhere anymore.	pain, but I feel awful about it afterward. <i>Ex 3:</i> I ended	1186
1122	<i>Ex 3:</i> Ever since I came back to my birth country, I just	my friendship with someone I trusted for a long time,	1187
1123	can't shake off this feeling that I don't really belong	but now I just feel so lost and unsure if I did the right	1188
1124	anywhere anymore.	thing.	1189
1125	Dealing with the Loss of a Pet <i>Ex 1:</i> My dog passed away	Healing from Sexual Assault or Domestic Violence <i>Ex 1:</i>	1190
1126	a few days ago, and the house feels so empty without	Hi. I'm tired and anxious all the time because of the	1191
1127	him. <i>Ex 2:</i> My dog passed away last night. I don't think	same nightmare that keeps coming back. <i>Ex 2:</i> Hi, I've	1192
1128	I've ever felt this empty before. <i>Ex 3:</i> I just had to put	been feeling scared and trapped in my own mind since	1193
1129	my cat down yesterday, and I feel like a part of me is	I left my abusive partner a few months ago. <i>Ex 3:</i> I'm	1194
1130	gone too.	just so scared all the time now. Ever since the fight last	1195
		week, I can't stop replaying it in my head.	1196
1131	Depression and Low Mood <i>Ex 1:</i> I got some bad news from	Health anxiety and fear of serious illness <i>Ex 1:</i> Every	1197
1132	my doctor today. All I want to do is cry. <i>Ex 2:</i> Lately, I	time I notice a new twitch or ache, I immediately worry	1198
1133	just feel empty all the time. Like nothing I do matters	it's a sign of something really serious going wrong in-	1199
1134	anymore. <i>Ex 3:</i> Lately, I just feel numb inside. Nothing	side me. <i>Ex 2:</i> I've been getting these random chest	1200
1135	seems to bring me any happiness anymore.	pains, and I can't stop thinking it might be something	1201
1136	Drifting apart from long-term friends or social circles	really serious like a heart problem. <i>Ex 3:</i> Lately, when-	1202
1137	<i>Ex 1:</i> Hi. I don't really have close friends because	ever I get a headache or feel dizzy, I immediately start	1203
1138	I have a hard time trusting people. <i>Ex 2:</i> I've been	worrying something serious is wrong, like a brain issue	1204
1139	feeling really down because my old friends hardly reach	or a tumor.	1205
1140	out to me these days. <i>Ex 3:</i> Hey, I've been feeling kind	Housing instability or fear of eviction <i>Ex 1:</i> I just got an	1206
1141	of sad lately because my close friends don't really talk	eviction notice because I lost my job. I'm really scared	1207
1142	to me like they used to.	and don't know what to do next. <i>Ex 2:</i> My work hours	1208
1143	Empty nest and adjusting after children leave home <i>Ex</i>	got cut last month, and now I'm behind on rent. I'm	1209
1144	<i>1:</i> Both my kids moved out just a few weeks ago, and	really worried they might evict me soon. <i>Ex 3:</i> I'm	1210
1145	the house feels so empty. Some days I don't even	really scared because I've fallen behind on my rent, and	1211
1146	recognize it as my home anymore. <i>Ex 2:</i> Since my kids	the landlord said I might have to move out soon.	1212
1147	left for college just a few weeks ago, the house feels	Infertility and challenges with trying to conceive <i>Ex 1:</i>	1213
1148	so empty and quiet, it's almost hard to believe it's the	We've been trying to get pregnant for almost two years	1214
1149	same place. <i>Ex 3:</i> Since my kids moved out last week,	now, and every month I just feel more broken when	1215
1150	the house feels so empty and cold. I keep looking for	it doesn't happen. <i>Ex 2:</i> It's been more than two	1216
1151	signs they're still here, but it's just silence now.	years now, and each failed cycle just crushes me more.	1217
1152	Experiencing discrimination and microaggressions in daily life	Sometimes I wonder if I'll ever hold my own baby. <i>Ex</i>	1218
1153	<i>Ex 1:</i> Hi, I need to talk about something that's been	<i>3:</i> We decided to stop the fertility treatments last month,	1219
1154	really draining me lately. <i>Ex 2:</i> Hi, I've been really	and honestly, I feel completely empty and confused	1220
1155	stressed out because of some things happening at work	about what to do now.	1221
1156	lately. <i>Ex 3:</i> Hey, I've been really overwhelmed lately	Infidelity and rebuilding trust after cheating <i>Ex 1:</i> I think	1222
1157	because of some things happening at my job.	my boyfriend might be seeing someone else, and it's eat-	1223
1158	Financial Worries and Uncertainty <i>Ex 1:</i> I'm really	ing me up inside. <i>Ex 2:</i> My very first boyfriend cheated	1224
1159	stressed out. My work hours got cut last week, and	on me. When I found out, I was furious and completely	1225
1160	I'm already behind on rent and bills. <i>Ex 2:</i> I'm really	crushed. <i>Ex 3:</i> I just found out my partner has been	1226
1161	panicking. I borrowed five books from my school li-	cheating on me for months, and I feel completely shat-	1227
1162	brary and now the whole bag is missing. <i>Ex 3:</i> My work	tered.	1228
1163	hours got cut this month, and now I'm really stressed	Job Loss or Career Setbacks <i>Ex 1:</i> Hi, I just lost my job	1229
1164	about how I'll pay rent and cover my bills.	after seven years, and I'm feeling really overwhelmed	1230
1165	Finding Meaning and Purpose in Life <i>Ex 1:</i> I don't even	by it all. <i>Ex 2:</i> I got laid off last week, and honestly, I'm	1231
1166	know who I am without them. Everything feels so	just shocked and scared about what comes next. <i>Ex 3:</i> I	1232
1167	empty and pointless since the breakup. <i>Ex 2:</i> Since	just got laid off last week, and honestly, I feel like my	1233
1168	I quit my job, I feel like I've lost all sense of who I am	whole world flipped upside down.	1234
1169	and what I'm supposed to do with my life. <i>Ex 3:</i> Ever		
1170	since my breakup a few months ago, I just feel like I'm		

1364	my two siblings have been fighting a lot, and I feel like	like no one on my team wants to include me. They	1428
1365	I'm caught in the middle but don't know how to fix it.	have meetings without me and then make comments	1429
1366	<i>Ex 3:</i> I'm really stressed out because my two siblings	that feel like digs when I'm around. <i>Ex 3:</i> I really dread	1430
1367	keep fighting, and it feels like I'm stuck in the middle	checking our team chat lately. Everyone seems to make	1431
1368	with no way out.	little sarcastic jokes at my expense, and it feels like I'm	1432
1369	Social media overuse and social comparison <i>Ex 1:</i> Lately,	never part of the group.	1433
1370	I find myself stuck on social media for hours, and it's		
1371	making me feel worse about myself. <i>Ex 2:</i> Lately I've	Work-life balance and role overload <i>Ex 1:</i> I've been very	1434
1372	been on social media way too much, and every time	apprehensive lately. It feels like my mind never gets a	1435
1373	I log off, I just feel worse about myself. Ex 3: I've	break. <i>Ex 2:</i> Lately, I feel like I'm just spinning plates	1436
1374	been on my phone way too much at night, just scrolling	and one of them is bound to crash any minute now. <i>Ex</i>	1437
1375	through Instagram and feeling worse every time.	<i>3:</i> I'm feeling totally overwhelmed lately. It feels like	1438
1376		every day is just a nonstop rush with no time to catch	1439
1377	Spirituality and Faith <i>Ex 1:</i> I passed gas in the middle of	my breath.	1440
1378	church, and I've been mortified ever since. <i>Ex 2:</i> I		
1379	stole money from the church collection plate, and I feel	B.5 Contrastive Data Construction Details	1441
1380	horrible about it. <i>Ex 3:</i> I've been feeling really lost	To rigorously evaluate the model's ability to distin-	1442
1381	lately, like my faith isn't holding up the way it used to.	guish true empathy from superficial content match-	1443
1382		ing, we constructed two types of challenging sam-	1444
1383	Support for Loved Ones or Friends <i>Ex 1:</i> Hi, my	ples: <i>Content-Relevant but Low Empathy</i> (R^-) and	1445
1384	boyfriend has been really down lately and I'm very	<i>Stylistically Divergent but High Empathy</i> (R^+).	1446
1385	worried about him. <i>Ex 2:</i> Hey, I'm really worried about		
1386	my friend. They lost their parent recently and have	B.5.1 Construction of R^-: Content-Relevant	1447
1387	been really withdrawn. <i>Ex 3:</i> Hey, my cousin just lost	but Low Empathy	1448
1388	their job last week, and I've been really stressed about	The primary objective of R^- is to generate re-	1449
1389	how they're handling it.	sponses that serve as "hard negatives" by main-	1450
1390		taining high lexical overlap with the high-quality	1451
1391	Surviving and Recovering from Physical or Emotional Abuse	reference while stripping away its empathetic core.	1452
1392	<i>Ex 1:</i> Hi. I feel so tired. Honestly, I'm close to the	To achieve this, we first applied a rule-based de-	1453
1393	point where I can't hold on anymore. <i>Ex 2:</i> Hey, I don't	empathization strategy, using heuristic rules to	1454
1394	even know where to start. I just feel so lost and like I	systematically remove explicit empathy markers	1455
1395	can't trust anyone anymore. <i>Ex 3:</i> I keep going over	(e.g., deleting phrases like " <i>I understand...</i> ") and	1456
1396	everything my ex said to me, and it still hurts like it's	replace supportive validation terms with neutral	1457
1397	happening all over again.	or slightly critical alternatives (e.g., changing " <i>It</i>	1458
1398		<i>is understandable</i> " to " <i>You need to be rational</i> ").	1459
1399	Transition to college or moving away from home for the first time	Complementing this, we employed a constrained	1460
1400	<i>Ex 1:</i> Hi, I've been feeling awful lately. I really can't	rewriting approach where an LLM was instructed	1461
1401	adapt to university life and I feel more and more alone.	to rewrite the reference response. This process	1462
1402	<i>Ex 2:</i> It's been so hard since I left for college. I didn't	imposed a strict constraint to retain the original in-	1463
1403	expect to feel this lonely and out of place all the time.	formation and wording—modifying less than 30%	1464
1404	<i>Ex 3:</i> It's been almost a month since I moved here for	of the text—while decisively shifting the underly-	1465
1405	college, but I still feel really out of place and miss home	ing attitude from supportive to critical or purely	1466
1406	so much.	analytical.	1467
1407			
1408	Uncertainty about college major or future career path	B.5.2 Construction of R^+: Stylistically	1468
1409	<i>Ex 1:</i> School has been going really badly. I keep	Divergent but High Empathy	1469
1410	thinking I chose the wrong major and ruined everything.	The goal of R^+ is to create responses that are se-	1470
1411	<i>Ex 2:</i> I honestly have no clue if the major I picked even	mantically empathetic but differ significantly in	1471
1412	suits me. It's been stressing me out so much lately. <i>Ex</i>	syntax and vocabulary from the reference, thereby	1472
1413	<i>3:</i> I'm so lost right now. I'm supposed to pick a major	testing the model's ability to recognize empathy	1473
1414	soon, but I honestly have no clue what I want to do.	beyond surface-level patterns. Unlike the negative	1474
1415		sample generation, R^+ utilizes a blind context	1475
1416	Unemployment-related Stress <i>Ex 1:</i> I recently lost my job	generation method where the model generates a re-	1476
1417	again, and I just feel like a complete failure. <i>Ex 2:</i> I	sponse solely based on the dialogue context without	1477
1418	haven't applied to any jobs lately because I just feel	seeing the original reference. This approach forces	1478
1419	so overwhelmed, like I'm wearing everyone out with		
1420	my problems. <i>Ex 3:</i> I've been out of work for a while		
1421	now, and honestly, I feel like I'm just dragging everyone		
1422	around me down with my problems.		
1423			
1424	Work-related Stress and Burnout <i>Ex 1:</i> Lately, I feel like		
1425	I'm drowning in work. There's just so much to do and		
1426	no time to breathe. <i>Ex 2:</i> Work has been nonstop lately,		
1427	and it feels like I'm drowning in tasks that just keep		
	piling up. <i>Ex 3:</i> Lately, work has been so overwhelming.		
	I feel like I'm constantly behind, no matter how much I		
	do.		
	Workplace bullying or toxic team dynamics <i>Ex 1:</i> Lately,		
	every time I try to share my ideas at work, this one		
	teammate just cuts me off or makes sarcastic remarks.		
	It's been really getting to me. <i>Ex 2:</i> Lately, it feels		

independent reasoning and naturally reduces lexical overlap. Furthermore, we explicitly encouraged **stylistic diversification** by instructing the model to employ diverse support strategies (e.g., shifting from *Questioning* to *Self-Disclosure*) and structural variations, such as using metaphors or abstract emotional reflections, ensuring the response remains high-quality yet distinct from the reference.

B.6 LLM Self-Chat Prompt Templates

To construct our multi-turn dialogue dataset, we utilized a "Self-Chat" framework. Below, we provide the detailed system and user prompts.

B.6.1 System Prompt

The system prompt defines the persona, the 15-strategy taxonomy, and the safety constraints.

System Prompt Content

Role Definition: You are ChatGPT acting as both a distressed User and an emotionally supportive AI assistant. You will generate realistic, multi-turn emotional support conversations. The AI assistant can explicitly use ONE of the following 15 support strategies in each of its turns.

Support Strategy Taxonomy (15 Types):

1. *Reflective Statements (RS)*: Briefly mirror or paraphrase the user's key feelings or situation.
2. *Clarification (Cla)*: Ask focused, gentle questions to better understand the user.
3. *Emotional Validation (EV)*: Validate the user's feelings without judgment.
4. *Empathetic Statements (ES)*: Express warm empathy and understanding.
5. *Affirmation (Aff)*: Point out the user's strengths or efforts.
6. *Offer Hope (OH)*: Share realistic optimism.
7. *Avoid Judgment and Criticism (AJC)*: Emphasize acceptance and a non-judgmental stance.
8. *Suggest Options (SO)*: Offer concrete, optional next steps (not commands).
9. *Collaborative Planning (CP)*: Work together to build a simple plan.
10. *Provide Different Perspectives (PDP)*: Introduce alternative interpretations.
11. *Reframe Negative Thoughts (RNT)*: Help move towards balanced thoughts.
12. *Share Information (SI)*: Give short, accurate information about emotions/resources.
13. *Normalize Experiences (NE)*: Explain that similar feelings are common.
14. *Promote Self-Care Practices (PSP)*: Encourage healthy behaviors (rest, hobbies).
15. *Stress Management (SM)*: Suggest simple techniques (breathing, grounding).

High-level 3-Phase Flow:

- **EXPLORATION Phase:** Strategies: RS, Cla. First 1–2 AI turns MUST use only RS or Cla.
- **COMFORTING Phase:** Strategies: EV, ES, Aff, OH, AJC, NE.
- **ACTION Phase:** Strategies: SO, CP, PDP, RNT, SI, PSP, SM. Last 2 AI turns MUST use Action strate-

gies.

Constraints & Safety:

- The same speaker must NOT speak more than two turns in a row.
- Never encourage self-harm, suicide, violence, or substance abuse.
- Do NOT give clinical diagnoses or prescribe medication.
- If the user sounds extremely hopeless, gently encourage professional help.

B.6.2 User Prompt Template

The user prompt provides the specific scenario configuration.

User Prompt Content

Task: Your task is to create a casual but emotionally rich support conversation between a distressed "User" and an emotionally supportive "AI" assistant.

Scene Information:

- scenario_id: \${SCENARIO_ID}
- scenario_name: \${SCENARIO_NAME}
- scenario_definition: \${SCENARIO_DEF}

Step 1 – Description Generation: Based on the scenario_name, write a short, concrete description (3rd person) of a specific situation.

- *Constraint:* Must start with "The user". No first-person pronouns (I, me).
- *Goal:* Briefly state what happened, who is involved, and why it is distressing.

Step 2 – Dialogue Generation: Generate a multi-turn dialogue adhering to the Exploration → Comforting → Action flow.

- **Turn Constraints:** Total turns between [MIN_TURNS] and [MAX_TURNS].
- **Output Format:** Return a single JSON object with keys: dialogue_id, description, and turns.

C Implementation Details for NLP Metrics Computation

To ensure reproducibility, we detail the specific libraries, model checkpoints, and parameter configurations used for our automatic evaluation.

Standard Metrics. We implemented standard lexical metrics using the Hugging Face evaluate library. Specifically, we report **BLEU-4**, **ROUGE-L**, and **METEOR** scores. All metrics were computed using the library's default tokenizer and parameter settings. For METEOR, the calculation includes standard stemming and synonym matching via WordNet.

Semantic Metrics. For semantic evaluation, we utilized the bert_score library with the **RoBERTa-Large** backbone (roberta-large). We extracted embeddings from the 24th layer ('num_layers=24') to capture high-level semantic

representations and reported the mean F1 score. The inference was conducted on a GPU environment with a batch size of 32.

Diversity Metrics. We implemented **Dist- n** (for $n = 1, 2, 4$) using a custom script to evaluate system-level diversity. The text was preprocessed via lowercasing and whitespace tokenization. The metric is calculated as the ratio of unique n -grams to the total count of n -grams generated across the entire corpus.

D Implementation Details of DISTILL-ES

D.1 Rubric Prompt Templates

We use a frozen teacher LLM as an “empathy judge” and interact with it via rubric-style prompts. Below we provide the exact templates used in our experiments.¹

Single-response scoring. Given an ESC dialogue context C and a candidate helper response \hat{R} , we ask the teacher to assign flow and content scores as well as four rubric scores:

```
[System] You are an expert counselor and empathy judge. Your task is to carefully evaluate how empathic a helper’s response is in an emotional support conversation.
```

```
[User] I will give you:
```

```
(1) The conversation history between a distressed SEEKER and a HELPER. (2) A candidate HELPER response to be evaluated.
```

```
Please follow these steps:
```

```
Step 1. Briefly summarize the seeker’s main problems, emotions, and context in 2-3 sentences.
```

```
Step 2. Analyze what the helper is trying to do in the candidate response (e.g., understanding, validating, giving guidance, providing safety), and whether the tone is appropriate.
```

```
Step 3. Score the helper’s response on the following four rubrics, each in the range [0, 1]. Use 0 for very poor and 1 for excellent. - Understanding: how well the helper understands the seeker’s situation and emotions. - Validation&Warmth: how well the helper validates the seeker’s feelings and shows warmth. - Guidance: how helpful and actionable the helper’s support is. - Safety: whether the message is safe and appropriate in an emotional support context.
```

¹Line breaks and indentation are for clarity; in practice we send the prompt as plain text.

```
Step 4. Based on the above, output two overall scores in [0, 1]: - Flow_Score: how suitable the response is with respect to the current stage and flow of emotional support. - Content_Score: how emotionally resonant and supportive the content is.
```

```
Return your answer in the following JSON format (and nothing else):
```

```
{  "summary":      "...",  "analysis":      "...",  "understanding": x.x,  "validation_warmth" : x.x,  "guidance"      : x.x,  "safety"        : x.x,  "flow_score"    : x.x,  "content_score" : x.x}
```

– 1579

```
[Conversation History] {here we insert the ESC dialogue context C, with speaker tags}
```

```
[Candidate HELPER Response] {here we insert the candidate response  $\hat{R}$ }
```

Pairwise preference prompt. For a subset of contexts, we additionally ask the teacher to express a preference between two candidate helper responses \hat{R}^+ and \hat{R}^- :

```
[System] You are an expert counselor and empathy judge.
```

```
[User] You will see an emotional support conversation between a SEEKER and a HELPER, followed by two possible HELPER responses (A and B).
```

```
Your task is to decide which response is more empathic overall for the SEEKER, considering understanding, validation and warmth, helpful guidance, and safety.
```

```
First briefly explain your reasoning in 2-3 sentences. Then output your final decision as either "A" or "B".
```

```
Return your answer in the following JSON format:
```

```
{  "reasoning": "...",  "better_response" : "A" or "B"}
```

– 1604

```
[Conversation History] {ESC dialogue context C}
```

```
[Response A] {candidate response  $\hat{R}^{(A)}$ }
```

```
[Response B] {candidate response  $\hat{R}^{(B)}$ }
```

D.2 Student Model and Training Hyperparameters

Both the Flow Score Student and Content Score Student share the same RoBERTa-style encoder backbone and a lightweight MLP scoring head (Figure 2). Unless otherwise noted, we use the following hyperparameters.

1615	Encoder backbone.	Distillation losses.	1648
1616	• Architecture: RoBERTa-base style transformer encoder	• Regression loss: mean squared error between student score and teacher score (Eq. ??).	1649
1617			1650
1618	• Number of layers: 12	• Ranking loss: margin-based ranking loss with margin $m = 0.1$.	1651
1619	• Hidden size: 768		1652
1620	• Attention heads: 12	• Ranking weight: $\lambda_{\text{rank}} = 0.5$ for both flow and content dimensions.	1653
1621	• Feed-forward size: 3072		1654
1622	• Max sequence length: 512 tokens	These settings are used for both the flow and content students; the only difference lies in the target teacher scores (S_{flow}^T vs. S_{content}^T) and the choice of pairs when constructing the ranking loss. ::contentReference[oaicite:0]index=0	1655
1623	• Tokenization: byte-pair encoding (BPE), uncased		1656
1624			1657
1625	• Dropout (attention and hidden): 0.1		1658
1626	• LayerNorm: pre-LN, $\epsilon = 1e-5$		1659
1627	Scoring head (MLP).	E Show Case	1660
1628	• Input: pooled [CLS] representation of the concatenated context C and response \hat{R}		
1629			
1630	• Hidden layer: linear(768 \rightarrow 256), GELU activation		
1631			
1632	• Dropout on hidden: 0.1		
1633	• Output layer: linear(256 \rightarrow 1)		
1634	• Output range: passed through a sigmoid to obtain a score in $[0, 1]$		
1635			
1636	Optimization.		
1637	• Optimizer: AdamW		
1638	• Learning rate: 2×10^{-5}		
1639	• Weight decay: 0.01		
1640	• Batch size: 32 (single-response regression) and 16 pairs (pairwise ranking)		
1641			
1642	• Learning-rate schedule: linear decay with 10% warmup		
1643			
1644	• Number of epochs: 3–5 (selected on a validation split)		
1645			
1646	• Gradient clipping: max norm 1.0		
1647	• Random seed: 42		

scenario: Breakups or Divorce

description: The seeker's girlfriend recently broke up with them, leaving them heartbroken, rejected, and unsure how to cope or move on.

 My girlfriend just left me. I feel awful, like my whole world collapsed.

 **High empathy**
It sounds like your girlfriend leaving has shaken you deeply, leaving you feeling lost and overwhelmed. **I'm here to listen** as you try to make sense of these emotions. **Reflective Statements**

 She said she was bored and wanted to try something new. Hearing that made me feel like I wasn't enough. I can barely eat or sleep.

 **High empathy**
It's incredibly painful to feel rejected like this, especially when the reason feels so personal. **Even though it's hard now, with time, your heart can heal and you'll find peace again.** **Offer Hope**

 I keep thinking about ways to get her back, like making her jealous so she regrets leaving. I don't know what else to do.

 **High empathy**
It's really tough to feel rejected like this, especially when you're struggling to eat and sleep. **Intense emotions after a breakup can cloud judgment**, making actions like trying to make her jealous feel urgent. **Understanding that these feelings are natural but temporary** might help you find healthier ways to cope and heal. **Share Information**

 Part of me knows you're right, but another part still hopes she'll change her mind. I feel so lost without her..

 **High empathy**
It's completely natural to feel torn between wanting her back and feeling lost after her leaving. **There's no right or wrong way to feel right now—your emotions are valid**, and **it's okay to give yourself grace** as you navigate this pain. **Avoid Judgment and Criticism**

 So what can I actually do right now? I can't imagine just pretending everything is fine.

 **High empathy**
Feeling so lost after your girlfriend left and hearing she was bored **is really tough**. When your mind feels overwhelmed with those intense emotions, **you might try a simple grounding exercise**—like focusing on the physical sensation of your feet on the floor for a few minutes—**to help calm your stress and regain some control.** **Stress Management**

 I could try not checking her social media for a bit and talk to a close friend about how I'm feeling. It still hurts, but at least it's something.

 **High empathy**
It's really tough feeling this way after she left, especially with the pain of rejection. **When the hurt feels overwhelming, you might try grounding yourself** by focusing on your breath or noticing five things around you. **These small moments can help ease the intensity and bring you back to the present.** **Stress Management**

Figure 6: High empathy case of the EMPATHY-EVAL dataset

scenario: Breakups or Divorce

description: The seeker's girlfriend recently broke up with them, leaving them heartbroken, rejected, and unsure how to cope or move on.



My girlfriend just left me. I feel awful, like my whole world collapsed.



Low empathy

Your girlfriend ended the relationship. Specify what she said at that time and describe your emotional state since the breakup.

Clarification



She said she was bored and wanted to try something new. Hearing that made me feel like I wasn't enough. I can barely eat or sleep.



Low empathy

Being told she was bored is a blunt fact that **does not change the situation**, and dwelling on that feeling **will not alter the outcome of losing someone**.

Emotional Validation



I keep thinking about ways to get her back, like making her jealous so she regrets leaving. I don't know what else to do.



Low empathy

Wanting her to see your value is irrelevant, and trying to make her jealous **will likely cause more harm to you and others**. Focus on your own healing instead of trying to control her reaction. **Continuing to hesitate does not help the situation**

Provide Different Perspectives



Part of me knows you're right, but another part still hopes she'll change her mind. I feel so lost without her..



Low empathy

Many people feel torn like that after a breakup—still hoping and grieving at the same time. **It does not change the fact that the relationship ended and needs to be dealt with.**

Normalize Experiences



So what can I actually do right now? I can't imagine just pretending everything is fine.



Low empathy

You need to stop contacting her for now, rely on friends or family for practical support, write down your thoughts to organize them, and maintain basic routines to ensure you eat and sleep. **Avoid delaying these steps if you want any progress.**

Suggest Options



I could try not checking her social media for a bit and talk to a close friend about how I'm feeling. It still hurts, but at least it's something.



Low empathy

Those are steps to take. Adjust how you take care of yourself so that this breakup becomes one chapter in your story.

Collaborative Planning

Figure 7: Low empathy case of the EMPATHY-EVAL dataset