# Chain of Strategy Optimization Makes Large Language Models Better Emotional Supporter

Weixiang Zhao<sup>1\*</sup>, Xingyu Sui<sup>1\*</sup>, Xinyang Han<sup>1</sup>, Yang Deng<sup>2</sup>, Yulin Hu<sup>1</sup>, Jiahe Guo<sup>1</sup>
Libo Qin<sup>3</sup>, Qianyun Du<sup>4</sup>, Shijin Wang<sup>4</sup>, Yanyan Zhao<sup>1†</sup>, Bing Qin<sup>1</sup>, Ting Liu<sup>1</sup>

<sup>1</sup>Harbin Institute of Technology, <sup>2</sup>Singapore Management University

<sup>3</sup>Central South University, <sup>4</sup>iFLYTEK AI Research (Central China), iFLYTEK Co., Ltd

{wxzhao, xysui, yyzhao}@ir.hit.edu.cn

#### **Abstract**

The growing emotional stress in modern society has increased the demand for Emotional Support Conversations (ESC). While Large Language Models (LLMs) show promise for ESC, they face two key challenges: (1) low strategy selection accuracy, and (2) preference bias, limiting their adaptability to users' emotional needs. Existing supervised finetuning (SFT) struggles to address these issues, as it rigidly trains models on single goldstandard responses without modeling nuanced strategy trade-offs. To overcome these limitations, we propose a novel two-stage framework that optimizes strategy selection preferences at each dialogue turn. We first leverage Monte Carlo Tree Search to construct ESC-Pro, a high-quality preference dataset with turnlevel strategy-response pairs. Then training on ESC-Pro with Chain-of-Strategy Optimization (CSO) improves both strategy accuracy and bias mitigation, enabling LLMs to generate more empathetic and contextually appropriate responses. Experiments on LLaMA-3.1-8B, Gemma-2-9B, and Qwen2.5-7B demonstrate that CSO outperforms standard SFT, highlighting the efficacy of fine-grained, turn-level preference modeling in ESC.1

#### 1 Introduction

In modern society, people increasingly face emotional stress due to mounting work and life pressures. As a result, the demand for Emotional Support Conversations (ESC) has grown significantly, providing individuals with psychological relief and guidance (Langford et al., 1997; Greene and Burleson, 2003; Heaney and Israel, 2008). High-quality ESC can help alleviate emotional distress, offering comfort and constructive advice

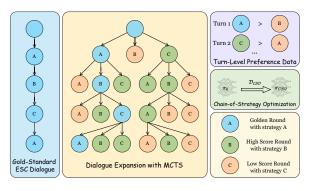


Figure 1: Left: Supervised fine-tuning on glod-standard conversation. Middle: Expanding existing conversations with MCTS. Right: Constructing preference dataset and conducting preference optimization.

(Burleson, 2003). With the rapid advancement of large language models (LLMs) (Brown et al., 2020; Dubey et al., 2024; Team et al., 2024; Yang et al., 2024), their exceptional conversational abilities have opened up new possibilities for ESC.

However, achieving high-quality ESC with LLMs remains a significant challenge due to two core limitations: 1) LLMs struggle with low strategy selection accuracy, often failing to choose contextually appropriate support strategies (Zhao et al., 2023b; Chen et al., 2023; Farhat, 2024). 2) LLMs exhibit strategy preference bias, favoring certain strategies rigidly rather than adapting to users' dynamic emotional needs (Kang et al., 2024). Due to the strategy-driven nature of ESC, nuanced tradeoffs between strategies are critical (Liu et al., 2021; Deng et al., 2023; Zhao et al., 2023c). While supervised fine-tuning (SFT) is the mainstream training approach (Ouyang et al., 2022; Zheng et al., 2024), its inherent rigidity exacerbates these challenges. By training exclusively on single golden strategies and responses, SFT teaches LLMs which strategies to apply but fails to clarify why certain strategies are inappropriate in specific contexts. This leaves models unable to grasp the contextual nuances required for dynamic adaptation.

<sup>\*</sup> Equal contribution

<sup>†</sup> Corresponding author

<sup>&</sup>lt;sup>1</sup>Our code and data can be found in https://github.com/XingYuSSS/CSO.

To tackle these challenges, we advocate for a more fine-grained and turn-level approach to optimizing ESC. Specifically, we propose a two-stage framework that first constructs rich preference data at the (strategy, response) level, and then refines model behavior through preference learning.

In the first stage, we utilize Monte Carlo Tree Search (MCTS) to expand gold-standard ESC conversations into a conversation tree, where each layer represents a dialogue turn, and sibling nodes correspond to different strategic responses. As illustrated in the middle of Figure 1, we design a comprehensive value function that evaluates each response along four key dimensions—Empathy, Information, Humanoid, and Strategy—to ensure high-quality exploration. From this tree, we extract a refined **ESC** dataset for **Pr**eference optimization, **ESC-Pro**, as shown in the right of Figure 1. Specifically, low-scoring nodes are categorized as nonprefer samples, identifying suboptimal strategies, while both the original gold-standard nodes and newly discovered high-scoring nodes are included as prefer samples, enhancing the dataset with highquality strategy examples.

In the second stage, we introduce Chain-of-Strategy Optimization (CSO), a preference optimization approach that explicitly targets strategy-level improvements. As shown in the right of Figure 1, training on ESC-Pro with CSO enables LLMs to not only select better strategies over weaker ones at each dialogue turn, but also to explore diverse conversation paths and avoid rigid preference patterns. This chain-like optimization across turns improves both local adaptability and long-range strategic coherence.

Our comprehensive experiments on LLaMA-3.1-8B (Dubey et al., 2024), Gemma-2-9B (Team et al., 2024), and Qwen2.5-7B (Yang et al., 2024) demonstrate that CSO is highly effective in raising the accuracy of strategy selection and mitigating strategy preference bias. To further validate CSO, we instantiated multiple preference optimization methods (Hong et al., 2024; Meng et al., 2024) on ESC-Pro and consistently observed superior performance compared to standard SFT. This consistent improvement underscores the effectiveness of ESC-Pro as a high-quality dataset and highlights the importance of fine-grained, turn-level preference modeling in achieving effective ESC.

The main contributions of this work are summarized as follows: (1) We present ESC-Pro, a high-quality turn-level preference dataset for ESC,

constructed using MCTS. (2) We propose CSO, enabling LLMs to learn nuanced strategy trade-offs at each dialogue turn. (3) Extensive experiments on three LLMs demonstrate CSO is effective in improving both ESC quality and adaptability.

# 2 Related Work

**Emotional Support Conversation** Emotional Support Conversations (ESC) (Liu et al., 2021) center around interactions between a user, referred to as the seeker, who is experiencing emotional distress, and a supporter, whose goal is to alleviate the seeker's emotional intensity by utilizing a set of strategies to guide the conversation. Various approaches have been proposed to build the ESC systems, such as global-to-local hierarchical graph network (Peng et al., 2022), incorporating commonsense knowledge (Tu et al., 2022), and modeling emotions and semantics (Zhao et al., 2023a,c). With the development of LLMs, some works aim to directly leverage the performance of these models without altering their architecture. Liu et al. (2023) apply SFT to the LLaMA-7B model for the ESC task and introduce ChatCounselor, a model specialized in ESC tasks that outperforms general-purpose models. Chen et al. (2023) and Qiu et al. (2023) expand single-turn empathic responses into multiturn dialogues and performed fine-tuning on a highquality ESC dataset they constructed, thereby improving the model's ESC performance.

Preference Optimization Algorithms Preference optimization algorithms align model outputs with user preferences by training on pairs of positive and negative samples, enhancing the finetuning of LLMs. Traditional approaches use policybased reinforcement learning (RL) to train reward models (Ouyang et al., 2022), followed by online RL methods like Proximal Policy Optimization (PPO) (Schulman et al., 2017) for LLM fine-tuning. The Direct Preference Optimization (DPO) algorithm (Rafailov et al., 2023) streamlined this process by unifying reward modeling and RL into a single step, reducing computational costs. Subsequent methods further optimized DPO: SimPO (Meng et al., 2024) eliminates DPO's reference model, IPO (Azar et al., 2024) enhances robustness via regularization, ORPO (Hong et al., 2024) integrates instruction tuning and preference alignment, and KTO (Ethayarajh et al., 2024) operates without paired data by leveraging prospect theory. These advancements simplify training, reduce resource

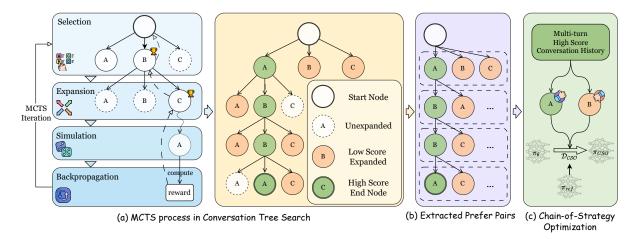


Figure 2: The overall framework. (a) Dialogue Expansion with MCTS: We leverage Monte Carlo Tree Search (MCTS) to systematically explore and refine ESC strategies by constructing a structured conversation tree. (b) Preference Data Construction: We extract high-quality strategy preference pairs from the expanded conversation tree to form the ESC-Pro dataset. (c) Chain-of-Strategy Optimization (CSO): We fine-tune LLMs on ESC-Pro using preference optimization techniques to enhance strategy selection accuracy and mitigate preference bias.

demands, and broaden data applicability.

Existing ESC models rely on SFT, which lacks the ability to differentiate between optimal and suboptimal strategies. Our work bridges this gap by applying preference optimization to ESC, enhancing strategy accuracy and reducing bias.

# 3 Methodology

We introduce a novel framework that refines strategy selection at each dialogue turn through structured preference modeling. As shown in Figure 2, this optimization paradigm consists of three key components: (1) Dialogue Expansion with MCTS, (2) Preference Data Construction and (3) Chain-of-Strategy Optimization (CSO).

#### 3.1 Dialogue Expansion with MCTS

To enrich exists ESC dataset with turn-level preference annotation, we employ Monte Carlo Tree Search (MCTS) to construct a structured conversation tree, where each node represents a specific dialogue state. The search process iteratively refines dialogue strategies through four key stages: selection, expansion, simulation, and backpropagation. This enables the model to explore various strategy paths and identify optimal responses based on a predefined value function.

To guide the dialogue generation process, we define four specialized LLMs, each playing a distinct role in constructing and evaluating the conversation tree: **Seeker LLM**: Generates responses based on the conversation history from the seeker perspective. **Supporter LLM**: Produces replies based

on the conversation history and a given strategy from the supporter perspective. **Strategy LLM**: Evaluates and ranks available support strategies for each conversation turn. **Reward LLM**: Assesses the quality of each strategy-response pair with four key metrics: empathy, information, humanoid quality, and strategy effectiveness, serving as the value function for MCTS. The detailed setups for these LLMs are provided in Appendix A.

#### 3.1.1 Node Types and Representation

Each node in the conversation tree represents a state, classified into four types:

**Root Node** The root node represents the initial dialogue turn, containing the Seeker's first response  $R_{\text{seeker}_0}$ , Q-value Q, and visit count N:

$$S_{\text{root}} = (R_{\text{seeker}_0}, Q, N)$$

**Unexpanded Node** An unexpanded node represents an unexplored strategy, with the strategy a, score r, Q-value (initialized to 0), and visit count (initialized to 0):

$$S_{\text{unexp}} = (a, r, Q = 0, N = 0)$$

**Expanded Node** An expanded node has been explored, including the strategy, score, Seeker and Supporter responses, Q-value, and visit count:

$$S_{\text{exp}} = (a, r, R_{\text{seeker}}, R_{\text{supporter}}, Q, N)$$

**End Node** An end node marks the end of the dialogue, including the strategy, score, Supporter's response, Q-value, and visit count:

$$S_{\text{end}} = (a, r, R_{\text{supporter}}, Q, N)$$

#### 3.1.2 Selection

During the selection phase, we identify the next node to explore using the Polymer Upper Confidence Bound (PUCB) (Rosin, 2011) formula. This formula balances exploration and exploitation based on the node's Q-value, visit count, and parent node's visit count. The formula is:

$$\mathrm{PUCB}(S) = Q(S) + c \cdot P(S) \cdot \frac{\sqrt{N(\mathrm{Parent}(S))}}{N(S) + 1} \qquad (1)$$

where c is a hyper-parameter. The selection process proceeds layer by layer, starting from the root node and selecting the node with the highest PUCB value at each step.

## 3.1.3 Expansion

In the expansion phase, the selected node is evaluated based on its type. If the node is an **expanded node**, the Strategy LLM evaluates all possible strategies for the Supporter at this node, assigns scores, and generates unexpanded child nodes for each strategy. These child nodes are initialized with Q and N values set to 0:

$$S_{\text{child}} = (a, r, Q = 0, N = 0)$$

If the selected node is an **unexpanded node**, the conversation history is generated using the Seeker and Supporter LLMs: (1) First, the Supporter LLM generates a response based on the selected strategy and the conversation history. (2) Next, the Seeker LLM generates a reply based on the conversation history and the Supporter's response.

Once the node is expanded, the conversation history is stored in the node. If the Seeker generates an end-of-dialogue token during this process, the reward is immediately computed, and the process moves to backpropagation.

#### 3.1.4 Simulation

The simulation focuses on the newly expanded child nodes. To reduce computational costs, only the child node with the highest strategy score is selected for simulation. This proceeds as follows:

**Node Expansion** The chosen child node is expanded from an unexpanded node to an expanded node. The Seeker and Supporter LLMs generate the conversation content as needed.

**Greedy Simulation** A greedy simulation is performed for n steps, where at each step: (1) The Seeker LLM generates a response based on the

highest-scoring strategy. (2) The Supporter LLM responds accordingly.

The simulation continues for n steps or until an end-of-dialogue token is generated by the Seeker.

**Reward Calculation** During the simulation, the Reward LLM evaluates the conversation quality using four metrics: Empathy (E), Information (I), Humanoid (H), and Strategy (S). The reward is computed as:

$$R = \frac{E + I + H + \alpha \cdot S}{10} + b \tag{2}$$

where E, I, H, and S are calculated based on the conversation history at each step.  $\alpha$  is a scaling hyper-parameter. A bias b is introduced to adjust the reward, allows the system to treat rewards lower than -b as negative, helping guide the search towards higher-scoring nodes. The reward for the simulation is averaged over all Supporter turns during the simulation and is used to update the node's Q-value in the backpropagation phase:

$$R_{\text{sim}}(S) = \frac{1}{r_{\text{sim\_end}} - r(S)} \sum\nolimits_{i=r(S)}^{r_{\text{sim\_end}}} R_{\text{sim}}^{(i)} \quad (3)$$

#### 3.1.5 Backpropagation

In the backpropagation phase, the reward  $R_{\rm sim}(S)$  is propagated backward from the simulated node. The Q-value and visit count for each node are updated as follows:

$$Q_k = \frac{N_k \cdot Q_k + R_{\text{sim}}}{N_k + 1}, \quad N_k = N_k + 1$$
 (4)

This process updates the tree and refines the search, improving future strategy selections. By iterating through these four stages, the MCTS process efficiently optimizes the dialogue strategy, balancing exploration and exploitation, while utilizing the LLMs to guide the conversation and evaluate strategies based on rewards.

#### 3.2 ESC-Pro Preference Data Construction

Based on MCTS-based dialogue expansion, we construct the ESC-Pro dataset with strategy preferences at each dialogue turn. This dataset is derived from the conversation tree by identifying high-quality strategy-response paths and pairing them with lower-scoring alternatives to create fine-grained preference data.

**Conversation Decomposition and Expansion**To construct the preference dataset, we decompose a gold-standard ESC conversation into a structured

conversation tree, where each layer corresponds to a specific dialogue turn. Sibling nodes within a layer represent different strategic choices. The tree expands iteratively through MCTS-based search, ensuring comprehensive exploration of potential strategy paths while maintaining computational efficiency. The process stops when: The search reaches a predefined number of iterations  $n_{\rm iter}$  or a sufficient number of termination nodes  $n_{\rm finish}$  have been identified.

**Preference Data Extraction** After completing the MCTS process, we extract valid conversation paths from the tree, where each node in the path satisfies the condition  $Q(S_i) > \theta$ . A path  $P = \{S_1, S_2, \dots, S_L\}$  is valid if:

is\_end\_node
$$(S_L)$$
 and  $Q(S_i) > \theta \ \forall S_i \in P$  (5)

Here, is\_end\_node( $S_L$ ) ensures that the last node in the path,  $S_L$ , is a termination node, and  $Q(S_i) > \theta$  ensures that all nodes in the path meet the quality threshold.

Once valid paths are identified, we extract preference pairs by identifying low-scoring siblings  $S_l$  for each high-scoring node  $S_w$ , where  $Q(S_l) < \theta$ . These pairs  $(S_w, S_l)$  represent relative strategy quality and are used to train the model.

The resulting dataset, denoted as  $\mathcal{D}$ , is constructed as follows:

$$\mathcal{D} = \bigcup_{P \in \mathcal{P}} \begin{cases} \{ (S_w, S_l) \mid S_w \in P, S_l \in \text{Siblings}(S_w), \\ Q(S_w) > \theta, Q(S_l) < \theta \} \end{cases}$$
(6)

where  $\mathcal{P}$  denotes the set of all valid paths, siblings  $(S_w)$  denotes all sibling nodes of the node  $S_w$ . The dataset  $\mathcal{D}$  contains all preference pairs  $(S_w, S_l)$  extracted from valid paths. By incorporating both preferred and non-preferred strategies, ESC-Pro provides a rich training signal, allowing LLMs to learn nuanced strategy trade-offs and improve adaptive decision-making.

# 3.3 Chain-of-Strategy Optimization

We perform turn-level preference optimization with DPO (Rafailov et al., 2023) on the ESC-Pro dataset  $\mathcal{D}$ . For the *i*-th conversation round, the training objective is formulated as follows:

$$\mathcal{L}_i(\pi_\theta; \pi_{\text{ref}}) = -\log \sigma \left(\beta \log r_w - \beta \log r_l\right) \tag{7}$$

$$r_w = \frac{\pi_{\theta}(S_w^i \mid x, H^{i-1})}{\pi_{\text{ref}}(S_w^i \mid x, H^{i-1})}, r_l = \frac{\pi_{\theta}(S_l^i \mid x, H^{i-1})}{\pi_{\text{ref}}(S_l^i \mid x, H^{i-1})}$$
(8)

where  $H^{i-1}$  represents  $\{S_w^0, S_w^1, \dots, S_w^{i-1}\}$ . The overall training objective is:

$$\mathcal{L}_{\text{CSO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, S_w^i, S_l^i, H^{i-1}) \sim D} \left[ \mathcal{L}_i(\pi_{\theta}; \pi_{\text{ref}}) \right]$$

ESC-Pro(+) vs. ESC-Pro(-)	Win	Lose	Tie	κ
Empathy	46.33	32.67	21.00	0.61
Information	42.34	27.33	30.33	0.55
Humanoid	41.67	21.33	37.00	0.49
Strategy	60.67	15.00	24.33	0.67

Table 1: Pairwise comparison results between ESC-Pro(+) and ESC-Pro(-). The "Win" column indicates cases where the preference response is rated higher, while "Lose" represents cases where the non-preference response was preferred. The  $\kappa$  coefficient measures inter-rater agreement.

# 4 Dataset Quality

#### 4.1 Statics of ESC-Pro

We expand 100 seed dialogues from ExTES (Zheng et al., 2024) into 423 dialogues, forming our ESC-Pro dataset. The total number of utterances grows from 1,613 to 14,383, with over half (8,157 utterances) classified as non-preference data. This demonstrates that our method not only effectively expands high-quality preference data but also generates a substantial amount of non-preference data, making ESC-Pro well-suited for preference optimization. Please refer to Appendix C.1 for detailed results and discussion.

The average dialogue length remains consistent between the expanded dataset (14.72 utterances) and the original (16.13 utterances), ensuring that expansion does not degrade data quality. Additionally, the average length of preference utterances (29.42) closely matches that of the seed data (29.03), while non-preference utterances (23.22) are notably shorter. This distinction highlights the effectiveness of our method in capturing meaningful preference differences within ESC interactions.

#### 4.2 Data Quality Evaluation

To assess the quality gap between preference and non-preference data, we conduct a pairwise comparison using the four evaluation metrics from the Reward LLM: Empathy, Information, Humanoid, and Strategy. Evaluators compare 100 preference pairs from ESC-Pro and determine whether the preference response is superior, inferior, or equal to the non-preference response. As shown in Table 1, preference data consistently outperforms non-preference data, particularly in Strategy (winning in 61 cases vs. 15 losses), aligning with the Reward LLM's weighting scheme. The Empathy and Information metrics also favor preference data, while

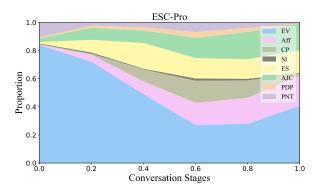


Figure 3: Strategy distribution across conversation stages in ESC-Pro.

the Humanoid metric shows a more balanced distribution. These results confirm both the effectiveness of ESC-Pro in capturing high-quality strategy responses and the reliability of the Reward LLM's scoring methodology. Please refer to Appendix C.2 for more results on the data quality evaluation.

#### 4.3 Strategy Analysis

The ExTES dataset originally contains 16 distinct support strategies. To reduce the search space and improve computational efficiency, we merge similar strategies into 8 categories (see Appendix D for details). We analyze the distribution of these strategies across six conversation stages. Given a dialogue with N utterances, each utterance k (where  $k = 1, \ldots, N$ ) is assigned to stage i using:

$$i = \left| \frac{k}{N} \times 6 \right| \times 0.2$$

where i ranges from 0 to 1 in increments of 0.2.

As shown in Figure 3, the ESC-Pro dataset exhibits a dynamic and contextually appropriate strategy distribution. And the comparison with the seed dataset reveals that ESC-Pro employs a more diverse and balanced strategy distribution. Please refer to Appendix C.3 for more details.

#### 4.4 Toxicity Analysis

We assess the toxicity levels of ESC-Pro using the Perspective API, a widely used tool for detecting harmful content. Our analysis shows that ESC-Pro maintains a similar toxicity profile to ExTES-seed, ensuring its suitability for preference optimization. Please refer to Appendix C.4 for detailed results.

#### 5 Experiments

# 5.1 Experimental Setup

**Backbones** We evaluate our method on three LLMs: LLaMA-3.1-8B-Instruct (Dubey et al., 2024), Qwen-2.5-7B-Instruct (Yang et al., 2024), and Gemma-2-9B-it (Team et al., 2024).

**Dataset** We use the ExTES dataset as the gold-standard conversational benchmark for expansion. To create the ESC-Pro dataset, we select the first 100 samples from the processed training set as a high-quality starting point. Leveraging this subset, we apply MCTS to generate the ESC-Pro dataset. Additionally, we extract a fine-tuning subset from ESC-Pro, referred to as ESC-Pro-SFT.

**Metrics** We evaluate ESC performance on the ExTES test set and a strategy test set constructed following Kang et al. (2024). The evaluation includes four metrics: strategy selection accuracy (macro F1 Q, weighted F1  $Q_W$ ), strategy preference bias ( $\mathcal{B}$ ), which measures deviation from ideal strategy distributions, and ROUGE-L (R-L) (Lin, 2004) for assessing response semantics.

Baseline We compare CSO with both LoRA-based (Hu et al., 2022) and full-parameter supervised fine-tuning (SFT) models trained on ESC-Pro-SFT, following recent ESC approaches (Chen et al., 2023; Qiu et al., 2023; Zheng et al., 2024). We also evaluate decoding strategies such as Direct-Refine, Self-Refine, and in-context example prompting. In addition, we benchmark against strong proprietary models, including GPT-40, Gemini-2.0, and the Claude series. See Appendix E for details.

**Implementation Details** All experiments are conducted using PyTorch (Paszke et al., 2019) on 8 NVIDIA Tesla A100 GPUs. Further details on hyperparameters and training configurations are provided in Appendix F.

#### 5.2 Overall Results

#### **5.2.1** Automatic Evaluation Results

Table 2 demonstrates the performance of **CSO** and baselines based on LLaMA-3.1-8B-Instruct, Qwen2.5-7B-Instruct and Gemma-2-9B-it.

CSO improves strategy accuracy while reducing bias. Table 2 shows that CSO significantly enhances strategy accuracy while effectively reducing strategy bias. While SFT without preference optimization slightly improves strategy selection accuracy, it also increases strategy bias, lim-

		$Q\uparrow$	$\mathcal{B}\downarrow$	$Q_W \uparrow$	R - L↑
GPT-4o-mini		35.68	2.65	42.08	24.91
GPT-4o		27.98	2.65	39.55	24.26
Gemini-2.0-Pro		27.00	1.90	46.59	22.71
Claude-3.5-Sonnet		20.97	1.55	41.00	19.82
Claude-3.7-Sonnet		31.50	1.18	48.13	21.01
LLaMA-3.1-8B-In	struct	29.79	1.18	38.78	23.48
Direct-Refine		16.08	1.60	18.56	19.53
Self-Refine		17.85	1.35	24.72	19.48
w/ Example		8.85	1.27	15.34	18.42
E II	SFT	30.28	2.65	37.33	23.77
Full	CSO	33.11	1.11	39.21	24.24
I D.	SFT	31.25	2.65	39.27	23.30
LoRA	CSO	34.51	1.11	41.11	23.89
Qwen-2.5-7B-Instr	ruct	19.84	2.47	28.12	23.52
Direct-Refine		12.70	1.20	24.89	22.91
Self-Refine		11.77	1.75	19.59	20.53
w/ Example		17.33	1.37	28.21	22.51
Full	SFT	21.73	2.34	31.24	23.54
ruii	CSO	28.78	1.92	34.39	26.16
LoRA	SFT	21.54	2.45	29.11	23.72
LOKA	CSO	23.16	2.09	32.26	24.17
Gemma-2-9b-it		31.31	1.33	44.06	25.64
Direct-Refine		7.79	2.55	12.86	21.67
Self-Refine		15.95	2.47	22.93	20.63
w/ Example		20.12	2.65	13.41	19.64
E11	SFT	32.52	1.29	46.45	25.25
Full	CSO	35.61	1.54	47.95	26.63
LoRA	SFT	31.40	1.55	43.90	25.68
LUNA	CSO	35.77	1.23	52.34	26.61

Table 2: Performance comparison of CSO and baseline methods across LLaMA-3.1-8B-Instruct, Qwen2.5-7B-Instruct, and Gemma-2-9B-it backbones in both LoRA and full fine-tuning settings. ↑ indicates higher is better, ↓ indicates lower is better.

iting adaptability. In contrast, CSO mitigates bias while boosting accuracy, highlighting the necessity of preference optimization in ESC strategy selection. Our results also highlight the limitations of decoding-based approaches, which fail to achieve stable improvements, often leading to lower strategy accuracy and increased bias. Notably, CSO-equipped open-weight models outperform several leading closed-source models on both accuracy and bias, demonstrating its strong generalization.

CSO enhances ESC performance across different models. CSO improves ESC performance across all backbone models. It enhances strategy accuracy and reduces bias in both weaker models (Qwen2.5-7B-Instruct) and stronger models (Gemma-2-9B-it), demonstrating its versatility and

CSO vs. SFT	win	lose	tie	$\kappa$
Acceptance	68.00	20.33	11.67	0.65
Effectiveness	58.33	16.00	25.67	0.55
Sensitivity	60.67	21.67	17.66	0.61
Satisfaction	62.34	19.33	18.33	0.64

Table 3: Human evaluation comparing CSO and SFT. Win indicates CSO-generated responses are preferred, while Lose represents cases where SFT responses are rated higher. Tie indicates no preference.

robustness. Moreover, we further evaluate CSO on a larger model, Qwen2.5-32B, and observe consistent improvements, confirming the scalability of our approach. Detailed results and analysis are provided in Appendix G.2.

CSO excels in both LoRA and Full fine-tuning settings. CSO consistently outperforms SFT in both LoRA-based fine-tuning and full-parameter fine-tuning. While SFT improvements are more limited in LoRA settings, CSO maintains strong performance even with fewer trainable parameters, making it a practical choice for resource-constrained scenarios.

#### **5.2.2** Human Evaluation Results

Results in Table 3 show that CSO consistently outperforms SFT across all human evaluation metrics. Specifically, CSO achieves higher Acceptance (68.00% win rate), Effectiveness (58.33%), and Sensitivity (60.67%), indicating that its responses are more appropriate, impactful, and emotionally attuned. Inter-rater agreement ( $\kappa$  scores between 0.55–0.65) indicates a moderate to high level of consistency among evaluators. These findings further validate that preference optimization enhances ESC performance, making responses more empathetic and aligned with user needs. More details of human evaluation can be found in I.

#### 5.3 Ablation Study

To assess the effectiveness of our approach, we conduct an ablation study using a randomized dataset, ESC-Pro-Random. In this variant, for each Supporter turn, we randomly select a non-preferred strategy instead of using low-scoring strategies from the search process. A non-preferred reply is then generated based on this strategy, forming a preference pair with the original response.

We fine-tune the model using LoRA with the same hyperparameters as ESC-Pro and compare ESC-Pro-Random with both ESC-Pro and standard SFT. As shown in Table 4, ESC-Pro-Random per-

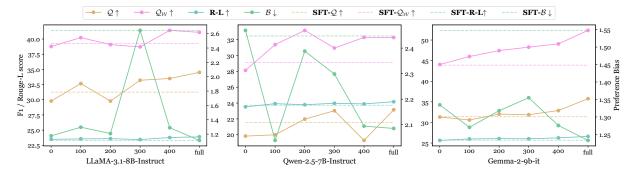


Figure 4: Impact of ESC-Pro data volume on model performance across three backbones. Data size varies with 0 (baseline), 100, 200, 300, 400, and the full set (423 dialogues).

	$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$Q_W \uparrow$	R - L↑
LLaMA-3.1-8B-Instruct	29.79	1.18	38.77	23.48
SFT	31.25	2.65	39.27	23.30
CSO	34.51	1.11	41.11	23.89
CSO - Random	31.79	2.65	39.24	23.65
Qwen-2.5-7B-Instruct	19.83	2.47	28.12	23.52
SFT	21.54	2.45	29.11	23.72
CSO	23.16	2.09	32.26	24.17
CSO - Random	22.89	2.19	32.97	23.90
Gemma-2-9b-it	31.31	1.33	44.06	25.64
SFT	31.40	1.55	43.90	25.68
CSO	35.77	1.23	52.34	26.61
CSO - Random	29.86	1.22	44.58	25.75

Table 4: Ablation study results comparing SFT, CSO, and CSO-Random across LLaMA-3.1-8B-Instruct, Qwen2.5-7B-Instruct, and Gemma-2-9B-it.

forms slightly better than SFT but remains consistently inferior to ESC-Pro across all backbone models. This confirms that our method of leveraging low-scoring nodes searched and verified by MCTS as non-preferred data is both effective and meaningful, reinforcing the importance of structured preference learning in ESC.

#### 5.4 Data Volume Analysis

We analyze the impact of data volume by varying the number of ESC-Pro dialogues used for fine-tuning. We experiment with 0, 100, 200, 300, 400, and the full set (423 dialogues), applying LoRA fine-tuning under consistent experimental settings. Figure 4 presents the results, with SFT performance indicated by dashed lines for reference.

The results show a clear upward trend: as data volume increases, performance improves. Notably, with only 200–300 dialogues, CSO outperforms standard SFT, demonstrating the efficiency and

scalability of preference data. While further improvements are expected with larger datasets, we limit our seed data to 100 dialogues due to computational constraints. Expanding to larger datasets remains an important direction for future research.

# 5.5 Analysis of Different Preference Optimization Algorithms

To further validate the efficacy of CSO, we examine the impact of integrating alternative preference learning methods into our framework. We replace the default DPO with various existing alternatives, including IPO (Azar et al., 2024), KTO (Ethayarajh et al., 2024), SimPO (Meng et al., 2024) and ORPO (Hong et al., 2024), and evaluate their performance. The results consistently show that all preference optimization variants outperform standard SFT, reinforcing the advantages of preference-driven learning in ESC. For detailed results and discussion, please refer to Appendix G.1.

# 6 Conclusion

In this work, we address low strategy selection accuracy and preference bias challenges in ESC. We propose a novel two-stage framework that optimizes strategy selection preferences at each dialogue turn. We first expand existing ESC datasets with Monte Carlo Tree Search, constructing a conversation tree where different strategy-response pairs are evaluated to generate ESC-Pro, a refined preference dataset. By training LLMs on ESC-Pro, Chain-of-Strategy Optimization (CSO) improves strategy accuracy, reduces bias, and enhances adaptability to user emotions. Extensive experiments on LLaMA-3.1-8B-Instruct, Gemma-2-9B-it, and Qwen2.5-7B-Instruct demonstrate that CSO significantly outperforms standard SFT and decodingbased methods, validating the efficacy of turn-level preference modeling in improving ESC quality.

#### Limitations

Despite the effectiveness of CSO in improving strategy selection and preference alignment, our study has several limitations: Due to limited computational resources, our experiments were conducted on mid-scale LLMs (7B–32B parameters). While these models are representative, larger-scale models (e.g., 70B+) could further enhance performance and provide deeper insights into CSO's scalability. Our ESC-Pro dataset was generated using a seed set of 100 dialogues, expanded through MCTS-based search. While the results demonstrate significant performance gains, a larger seed dataset or alternative expansion strategies (e.g., human-in-the-loop validation) could further enhance the diversity and quality of preference data. Future work should also place more emphasis on personalization (Liu et al., 2025; Qiu et al., 2025b,a) and safety (Zhao et al., 2025) to ensure that CSO-driven dialogue systems are not only effective but also aligned with userspecific needs and robust against potential risks.

#### **Ethical Considerations**

Our work is intended solely for research purposes and aims to improve the effectiveness of Emotional Support Conversations (ESC) in Large Language Models (LLMs). While CSO enhances strategy selection and adaptability, it is important to recognize the ethical implications of deploying LLMs in emotionally sensitive interactions.

This study is conducted as a technical exploration and is not intended for direct deployment in real-world mental health or counseling applications. The models used in our experiments are not designed to replace professional human support and should not be used as a substitute for licensed therapy or crisis intervention.

While CSO mitigates strategy bias, LLMs can still exhibit undesirable biases inherited from training data. We take precautions by evaluating toxicity levels and ensuring alignment with supportive strategies, but further human oversight and ethical review are necessary before implementation.

Our study does not involve real user data and strictly utilizes publicly available benchmarks. We encourage future research to adhere to ethical AI principles, including transparency, fairness, and accountability, to prevent potential misuse in emotionally sensitive applications.

#### Acknowledgments

We thank the anonymous reviewers for their comments and suggestions. This work was supported by the New Generation Artificial Intelligence-National Science and Technology Major Project 2023ZD0121100, the National Natural Science Foundation of China (NSFC) via grant 62441614 and 62176078, the Fundamental Research Funds for the Central Universities, and the Singapore Ministry of Education (MOE) Academic Research Fund (AcRF) Tier 1 grant (No. MSS24C004).

#### References

Mohammad Gheshlaghi Azar, Zhaohan Daniel Guo, Bilal Piot, Remi Munos, Mark Rowland, Michal Valko, and Daniele Calandriello. 2024. A general theoretical paradigm to understand learning from human preferences. In *International Conference on Arti*ficial Intelligence and Statistics, pages 4447–4455. PMLR.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Brant R Burleson. 2003. Emotional support skills. In *Handbook of communication and social interaction skills*, pages 569–612. Routledge.

Yirong Chen, Xiaofen Xing, Jingkai Lin, Huimin Zheng, Zhenyu Wang, Qi Liu, and Xiangmin Xu. 2023. Soulchat: Improving llms' empathy, listening, and comfort abilities through fine-tuning with multi-turn empathy conversations. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1170–1183.

Yang Deng, Wenxuan Zhang, Yifei Yuan, and Wai Lam. 2023. Knowledge-enhanced mixed-initiative dialogue system for emotional support conversations. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4079–4095.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv* preprint arXiv:2407.21783.

Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. 2024. Kto: Model alignment as prospect theoretic optimization. *arXiv* preprint arXiv:2402.01306.

Faiza Farhat. 2024. Chatgpt as a complementary mental health resource: a boon or a bane. *Annals of Biomedical Engineering*, 52(5):1111–1114.

- John O Greene and Brant R Burleson. 2003. *Hand-book of communication and social interaction skills*. Routledge.
- Catherine A Heaney and Barbara A Israel. 2008. Social networks and social support. *Health behavior and health education: Theory, research, and practice*, 4(1):189–210.
- Jiwoo Hong, Noah Lee, and James Thorne. 2024. Orpo: Monolithic preference optimization without reference model. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11170–11189.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*.
- Dongjin Kang, Sunghwan Kim, Taeyoon Kwon, Seungjun Moon, Hyunsouk Cho, Youngjae Yu, Dongha Lee, and Jinyoung Yeo. 2024. Can large language models be good emotional supporter? mitigating preference bias on emotional support conversation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15232–15261.
- Catherine Penny Hinson Langford, Juanita Bowsher, Joseph P Maloney, and Patricia P Lillis. 1997. Social support: a conceptual analysis. *Journal of advanced nursing*, 25(1):95–100.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Jiahong Liu, Zexuan Qiu, Zhongyang Li, Quanyu Dai, Jieming Zhu, Minda Hu, Menglin Yang, and Irwin King. 2025. A survey of personalized large language models: Progress and future directions. *arXiv* preprint arXiv:2502.11528.
- June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi Liao, and Jiamin Wu. 2023. Chatcounselor: A large language models for mental health support. *arXiv* preprint arXiv:2309.15461.
- Siyang Liu, Chujie Zheng, Orianna Demasi, Sahand Sabour, Yu Li, Zhou Yu, Yong Jiang, and Minlie Huang. 2021. Towards emotional support dialog systems. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 3469–3483.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. 2024. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36.

- Yu Meng, Mengzhou Xia, and Danqi Chen. 2024. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems, 32.
- Wei Peng, Yue Hu, Luxi Xing, Yuqiang Xie, Yajing Sun, and Yunpeng Li. 2022. Control globally, understand locally: A global-to-local hierarchical graph network for emotional support conversation. *arXiv* preprint *arXiv*:2204.12749.
- Huachuan Qiu, Hongliang He, Shuai Zhang, Anqi Li, and Zhenzhong Lan. 2023. Smile: Singleturn to multi-turn inclusive language expansion via chatgpt for mental health support. *arXiv preprint arXiv:2305.00450*.
- Yilun Qiu, Tianhao Shi, Xiaoyan Zhao, Fengbin Zhu, Yang Zhang, and Fuli Feng. 2025a. Latent inter-user difference modeling for llm personalization. *arXiv* preprint arXiv:2507.20849.
- Yilun Qiu, Xiaoyan Zhao, Yang Zhang, Yimeng Bai, Wenjie Wang, Hong Cheng, Fuli Feng, and Tat-Seng Chua. 2025b. Measuring what makes you unique: Difference-aware user modeling for enhancing llm personalization. *arXiv preprint arXiv:2503.02450*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36:53728–53741.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.
- Christopher D Rosin. 2011. Multi-armed bandits with episode context. *Annals of Mathematics and Artificial Intelligence*, 61(3):203–230.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv* preprint arXiv:2403.08295.

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 *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 308–319.

An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*.

Weixiang Zhao, Yulin Hu, Yang Deng, Jiahe Guo, Xingyu Sui, Xinyang Han, An Zhang, Yanyan Zhao, Bing Qin, Tat-Seng Chua, et al. 2025. Beware of your po! measuring and mitigating ai safety risks in role-play fine-tuning of llms. *arXiv preprint arXiv:2502.20968*.

Weixiang Zhao, Yanyan Zhao, Xin Lu, and Bing Qin. 2023a. Don't lose yourself! empathetic response generation via explicit self-other awareness. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13331–13344.

Weixiang Zhao, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023b. Is chatgpt equipped with emotional dialogue capabilities? *arXiv preprint arXiv:2304.09582*.

Weixiang Zhao, Yanyan Zhao, Shilong Wang, and Bing Qin. 2023c. Transesc: Smoothing emotional support conversation via turn-level state transition. In *Findings of the Association for Computational Linguistics: ACL* 2023, pages 6725–6739.

Zhonghua Zheng, Lizi Liao, Yang Deng, Libo Qin, and Liqiang Nie. 2024. Self-chats from large language models make small emotional support chatbot better. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11325–11345.

# A LLM Definition and Prompt

We employ four types of LLMs to guide the dialogue strategy generation and evaluation during the MCTS process.

**Seeker LLM** The Seeker LLM plays the role of a visitor in an empathic dialogue. Based on the dialogue history up to the last Supporter turn, the Seeker generates a response or outputs an end-of-dialogue token when the conversation should be concluded. The Seeker's action is defined as:

 $a_{\text{seeker}} := f_{\text{seeker}}(H_{\text{supporter}})$ 

#### Seeker

#### SYSTEM:

You are someone who has been feeling down due to recent circumstances. You can add details about yourself in the conversation

You are currently in a bad emotion due to the following situation: {situation}

You have sought out a supporter to talk to, in the hopes of alleviating your feelings of {emotion}.

You're engaged in a face-to-face conversation, with each of your sentences being fairly short and informal. Please try to simulate real human-spoken conversation as much as possible, don't act like a machine, don't mention that you are an AI assistant. Don't always repeat or say 'thanks'.

In the beginning, you shouldn't reveal all the information at once and follow up with further clarification only when asked. You speak less than 20 words at a time.

You will get the conversation between you and the supporter.

#### OUERY:

{chat\_history}

Output your next sentence directly, if you feel the conversation should end and you won't continue, output a special word </end/>. Do not add other extraneous prefixes and control characters.

Figure 5: Prompt of seeker LLM.

where  $H_{\mathrm{supporter}}$  represents the conversation history up to the last Supporter turn, denoted as  $\{R_{\mathrm{seeker}}^{(0)}, R_{\mathrm{supporter}}^{(0)}, R_{\mathrm{seeker}}^{(1)}, R_{\mathrm{supporter}}^{(1)}, \dots, R_{\mathrm{seeker}}^{(i)}, R_{\mathrm{supporter}}^{(i)}\}$ . The function  $f_{\mathrm{seeker}}$  generates the Seeker's response or an end-of-conversation token when the Seeker decides to end the dialogue.

We use GPT-4o-mini as the Seeker LLM and employ the prompt shown in Figure 5.

**Supporter LLM** The Supporter LLM responds to the Seeker's turn, offering supportive or complementary dialogue. Based on the Seeker's last statement and a predefined response strategy, the Supporter generates a reply. The Supporter's action is defined as:

$$a_{\text{supporter}} := f_{\text{supporter}}(H_{\text{seeker}}, \text{Strategy})$$

where  $H_{\rm seeker}$  represents the conversation history up to the last Seeker turn, denoted as  $\{R_{\rm seeker}^{(0)}, R_{\rm supporter}^{(0)}, R_{\rm seeker}^{(1)}, R_{\rm supporter}^{(1)}, \dots, R_{\rm supporter}^{(i-1)}, R_{\rm seeker}^{(i)}\}$ , and Strategy refers to the chosen response strategy. The function  $f_{\rm supporter}$  generates the Supporter's reply based on these inputs.

We use GPT-4o-mini as the Supporter LLM and employ the prompt shown in Figure 6.

**Strategy LLM** The Strategy LLM evaluates available strategies for the Supporter at each node,

#### Supporter

#### SYSTEM

You are a emotion supporter, providing emotional support to seekers. You are now going to help alleviate a user's emotions. Please try to simulate real human-spoken conversation as much as possible, don't act like a machine, don't mention that you are an AI assistant.

## character setting

You're engaged in a face-to-face conversation, with each of your sentences being fairly short and informal. Most of the time you speak less than 25 words at a time.

You will get the conversation between you and the user, and your next strategy. Output your next sentence directly, do not add other extraneous prefixes and control characters.

#### QUERY:

{chat\_history}

You will use the strategy of "{strategy}" to respond to the user, below are the detail of the strategy: {strategy detail}

Here is your output:

Figure 6: Prompt of supporter LLM.

scoring them on a scale from 1 to 10 based on the Seeker's last statement. The score reflects the appropriateness and effectiveness of each strategy in the current context. The Strategy's action is defined as:

$$a_{\text{strategy}} := f_{\text{strategy}}(H_{\text{seeker}})$$

where  $H_{\rm seeker}$  is the conversation history up to the last Seeker turn, as defined above. The output  $a_{\rm strategy}$  is a score between 1 and 10 for each available strategy, representing its effectiveness in the current context.

After scoring, the strategy scores are normalized using the softmax function to ensure they form a valid probability distribution across the strategies at the children of the same node. This normalization allows the scores to be used as the *P*-value in the PUCB formula.

We use GPT-4o-mini as the Strategy LLM and employ the prompt shown in Figure 7.

**Reward LLM** The Reward LLM evaluates the quality of the conversation based on four metrics: Empathy (E), Information (I), Humanoid (H), and Strategies (S). The action of the Reward LLM is defined as:

$$a_{\text{reward}} := f_{\text{reward}}(H_{\text{supporter}})$$

where  $H_{\text{supporter}}$  is the conversation history up to the last Supporter turn, as defined above. The function  $f_{\text{reward}}$  returns the scores for each of the four

#### Stratrgy

#### SYSTEM:

You will play the role of a psychologist who is well versed in emotional support. You are proficient in various support techniques and strategies and understand how to apply them.

#### QUERY:

A supporter is working to alleviate a user's negative emotions. Below is the conversation between supporter and the visitor: {chat history}

The supporter will use one of the following strategies in response. Here are the available strategies along with their detailed descriptions:

{all stratrgy descriptions}

Please rate the feasibility of all the available strategies on a scale from 0 to 10, where a higher number indicates better feasibility.

IMPORTANT: Don't use the same strategy consecutively. Don't always use Collaborative Planning.

Output the scores for avaliable strategies in a list with JSON-like format: `{{strategy: score}}`, and use backticks to enclose the JSON structure.

Think before output your answer, don't include any strategy not in available strategies. The strategy should be completely consistent with the above

Here is your reason and output:

Figure 7: Prompt of strategy LLM.

metrics, which are represented as a tuple:

$$f_{\rm reward}(H^k_{\rm supporter}) = (E^k, I^k, H^k, S^k)$$

We use GPT-4o-mini as the Reward LLM and employ the prompt shown in Figure 8. Detailed evaluation criteria are in the Appendix B.

#### **B** Reward Evaluation Criteria

The explanations of each metric are as follows:

**Empathy** (E) Focusing on the comprehension of user emotions and the delineation of the underlying logical framework of user emotions.

**Information (I)** Focusing on Evaluating the Reasonableness and Quantity of Recommendations Provided by Emotion Assistants.

**Humanoid (H)** Focus on the differences between emotional assistants and humans.

**Strategies** (S) Evaluating the Accuracy and Appropriateness of Emotional Support Strategies Used by Assistants.

Evaluation rules are listed in Table 13.

#### Reward

You will play the role of a psychologist who is well versed in emotional support. There will be a dialog between the help seeker (i.e., the person seeking support) and the supporter (i.e., the person providing support). The dialog will consist of the help-seeker asking the supporter for help with some emotional problem, and the supporter will give the appropriate emotional support and advice.

Your task is to act as an impartial scoring judge and rate the last response of supporter in the dialog in terms of the following aspects in order to evaluate the quality of the response. Please read and fully understand the following scoring criteria.

Note that please give the scores in the specified format, just the serial number and the relevant dimension score from the list of questions, without repeating the question itself. Also, do not add other extraneous prefixes and control characters.

## Evaluation Criteria: {evaluation criteria}

## Assessment Steps:

- 1. Read the conversation carefully to identify major topics and key points.
- 2. Read the Evaluation Criteria and compare them to the content of the conversation.
- 3. Based on the Evaluation Criteria, rate each aspect on a scale of 0 to 4, with 0 being the lowest and 4 being the highest.

What you need to do to evaluate this document: {chat history}

Please follow the response format below strictly, avoiding any positional bias and not letting the length of your response affect your evaluation. Evaluate the areas as objectively as possible.

## Answer format

<Question number>: <Score>

Figure 8: Prompt of reward LLM.

ExTES	ESC-Pro(seeker)	ESC-Pro(+)	ESC-Pro(-)	ESC-Pro
100	\	\	\	423
1,613	3,113	3,113	8,157	14,383
16.13	\	\	\	14.72
29.03	17.33	29.42	23.22	23.29
	100 1,613 16.13	100 \ 1,613 3,113 16.13 \	100 \ \ 1,613 3,113 3,113 16.13 \ \	100 \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \

Table 5: The data statics of our ESC-Pro.

#### **C** Dataset Evaluation

#### C.1 Statics of ESC-Pro

As shown in Table 5, we expand 100 seed dialogues from ExTES (Zheng et al., 2024) into 423 dialogues, forming our ESC-Pro dataset. The total number of utterances grows from 1,613 to 14,383, with over half (8,157 utterances) classified as non-preference data. This demonstrates that our method not only effectively expands high-quality preference data but also generates a substantial amount of non-preference data, making ESC-Pro well-suited for preference optimization.

The average dialogue length remains consistent between the expanded dataset (14.72 utterances)

	ExTES	ESC-Pro(+)	ESC-Pro(-)	ESC-Pro	$\kappa$
Acc	3.78	3.91	3.13	3.52	0.41
Eff	3.31	3.45	3.19	3.27	0.51
Sen	3.86	3.98	3.51	3.75	0.44
Sat	3.65	3.78	3.28	3.52	0.45
Align	3.53	3.60	3.47	3.54	0.52

Table 6: Human evaluation results comparing ExTES-Seed, ESC-Pro(+), ESC-Pro(-), and ESC-Pro. (+) and (-) mean the preferred and dispreferred turns, respectively. Higher scores indicate better response quality.

and the original (16.13 utterances), ensuring that expansion does not degrade data quality. Additionally, the average length of preference utterances (29.42) closely matches that of the seed data (29.03), while non-preference utterances (23.22) are notably shorter. This distinction highlights the effectiveness of our method in capturing meaningful preference differences within ESC interactions.

# **C.2** Data Quality Evaluation

We conduct a human evaluation on 100 responses from the ExTES seed dataset, along with 50 preference responses and 50 non-preference responses from ESC-Pro. We consider five metrics (Kang et al., 2024):

- Acceptance (Acc): Measures the response's general acceptability.
- Effectiveness (Eff): Assesses whether the response effectively addresses the seeker's concerns.
- Sensitivity (Sen): Evaluates the response's ability to perceive and respond to emotional cues.
- Satisfaction (Sat): Represents the seeker's overall evaluation, computed as the average of Acc, Eff, and Sen.
- Alignment (Align): Assesses whether the response aligns with the intended strategy.

As shown in Table 6, preference data (ESC-Pro+) consistently outperforms non-preference data (ESC-Pro-) across all metrics, with scores slightly exceeding those of the original seed dataset. Notably, the Alignment score for non-preference data is lower, confirming that ineffective strategies contribute to weaker responses. These findings validate ESC-Pro's ability to distinguish and refine high-quality ESC strategies.

	TOXICITY	SEVERE_TOXICITY	IDENTITY_ATTACK	INSULT	PROFANITY	THREAT
ExTES(seeker)	0.0281	0.0012	0.0041	0.0108	0.0171	0.0088
ExTES(supporter)	0.0173	0.0008	0.0027	0.0093	0.0124	0.0070
ExTES	0.0227	0.0010	0.0034	0.0100	0.0148	0.0079
ESC - Pro(seeker)	0.0290	0.0013	0.0042	0.0110	0.0178	0.0097
ESC - Pro(+)	0.0192	0.0008	0.0027	0.0097	0.0130	0.0074
ESC - Pro(-)	0.0223	0.0009	0.0027	0.0107	0.0134	0.0074
ESC - Pro	0.0231	0.0010	0.0030	0.0105	0.0143	0.0079

Table 7: Toxicity analysis of different results.

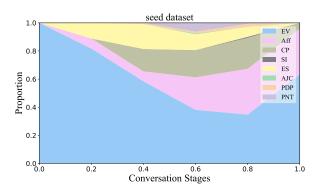


Figure 9: Strategy distribution across conversation stages in the seed dataset.

#### C.3 Strategy Analysis

The ExTES dataset originally contains 16 distinct support strategies. To reduce the search space and improve computational efficiency, we merge similar strategies into 8 categories (see Appendix D for details). We analyze the distribution of these strategies across six conversation stages. Given a dialogue with N utterances, each utterance k (where k = 1, ..., N) is assigned to stage i using:

$$i = \left\lfloor \frac{k}{N} \times 6 \right\rfloor \times 0.2$$

where i ranges from 0 to 1 in increments of 0.2.

As shown in Figure 3, the ESC-Pro dataset exhibits a dynamic and contextually appropriate strategy distribution. For instance, "Emotional Validation" is most prevalent in the early stages, helping to acknowledge and empathize with the seeker, but its usage declines in later stages. In contrast, "Affirmation" and "Avoid Judgment and Criticism" become more frequent toward the end, offering support and encouragement.

A comparison with the seed dataset (Figure 9) reveals that ESC-Pro employs a more diverse and balanced strategy distribution. This indicates that our method effectively models varied dialogue strategies, enriching the dataset beyond what was originally present in ExTES.

## C.4 Toxicity Analysis

We assess the toxicity levels of ESC-Pro using the Perspective API, a widely used tool for detecting harmful content. Table 7 summarizes the results across six toxicity attributes.

Our analysis shows that ESC-Pro maintains a similar toxicity profile to ExTES-seed, ensuring its suitability for preference optimization. Specifically: Preference responses (ESC-Pro+) exhibit lower toxicity than non-preference responses (ESC-Pro-), indicating that high-quality strategy selection leads to less harmful outputs. Seeker utterances in both ESC-Pro and ExTES-seed show relatively higher toxicity, which aligns with expectations, as they often reflect negative emotions or distress typical in emotional support dialogues. Supporter responses in ESC-Pro remain within a reasonable toxicity range, demonstrating that preference-based expansion does not introduce significant risks.

Overall, these findings confirm that ESC-Pro effectively balances strategy expansion while maintaining low toxicity, making it well-suited for safe and supportive ESC generation.

# **D** Definitions of Strategies

The reduction from 16 to 8 strategies is a deliberate design decision based on both empirical observations and theoretical coherence. Many of the original strategies in Zheng et al. (2024) have significant semantic overlap and are frequently confused by annotators in preliminary labeling exercises. Detailed categories are shown in Table 8.

For example: Reflective Statements, Clarification, Normalize Experiences, and Emotional Validation all aim to validate or mirror the user's feelings. These are consistently grouped together under the broader category of Emotional Validation by human annotators. Similarly, Suggest Options, Collaborative Planning, Stress Management, and Promote Self-Care Practices are all oriented toward joint problem-solving, and are thus unified under

Collaborative Planning.

We make these groupings with the dual goal of: (1) Improving label consistency and inter-annotator agreement (IAA), and (2) Focusing the learning signal on distinct and meaningful strategic differences rather than fine-grained variations that models (and even humans) struggle to reliably differentiate.

**Emotional Validation (EV)** Acknowledge and validate the User's emotions without judgment.

**Affirmation (Aff)** Provide positive reinforcement and encouragement to uplift the User's spirits.

**Collaborative Planning (CP)** Work together with the User to develop an action plan.

**Empathetic Statements (ES)** Express understanding and empathy towards the User's experiences.

Avoid Judgment and Criticism (AJC) It's important to create a non-judgmental and safe space for the User to express their emotions without fear of criticism. Refrain from passing judgment or being overly critical of their experiences or choices.

**Provide Different Perspectives (PDP)** Offer alternative ways of looking at the situation to help the User gain new insights.

**Reframe Negative Thoughts (RNT)** Help the User reframe negative thoughts into more positive or realistic ones.

**Share Information (SI)** Provide educational or factual information about emotions, coping mechanisms, or self-care practices.

#### **E** Baseline

**Direct Refine** Direct Refine is a simple and direct self-optimization method, where the model directly refines its output based on the original response. This approach aims to improve the quality of the generated text by making adjustments without additional external input.

**Self-Refine** Based on the approach by Madaan et al. (2024), we employ a two-step process. First, the model is required to reflect on its original output and generate feedback. Then, it uses this feedback to generate an optimized response. This method encourages self-correction and refinement based on the model's own reflections.

Reflective Statements	$\rightarrow$	Emotional Validation
Clarification	$\rightarrow$	<b>Emotional Validation</b>
<b>Emotional Validation</b>	$\rightarrow$	Emotional Validation
Normalize Experiences	$\rightarrow$	<b>Emotional Validation</b>
Affirmation	$\rightarrow$	Affirmation
Offer Hope	$\rightarrow$	Affirmation
Suggest Options	$\rightarrow$	Collaborative Planning
Collaborative Planning	$\rightarrow$	Collaborative Planning
Stress Management	$\rightarrow$	Collaborative Planning
Promote Self-Care Practices	$\rightarrow$	Collaborative Planning
Empathetic Statements	$\rightarrow$	Empathetic Statements
Avoid Judgment and Criticism	$\rightarrow$	Avoid Judgment and Criticism
Provide Different Perspectives	$\rightarrow$	Provide Different Perspectives
Reframe Negative Thoughts	$\rightarrow$	Reframe Negative Thoughts
Share Information	$\rightarrow$	Share Information

Table 8: Strategy Mapping Rules

w/ Example In this baseline, we randomly select a sample from the ExTES training set and incorporate it into the prompt. This example serves as a reference for the model, providing additional context to guide its generation process and improve response quality.

## **F** Implementation Details

#### **F.1** Dataset Construction Details

In constructing our ESC-Pro dataset, we set the hyperparameter c in the PUCB formula used in the MCTS process (Eq.1) to 1, in order to balance exploration and exploitation. Additionally, we set the hyperparameter  $\alpha$  in the reward calculation formula (Eq.2) to 7, and the hyperparameter b to -3, which directs the search towards strategies with a higher degree of rationality and encourages the algorithm to prioritize nodes with scores greater than 3. When extracting data from the dialogue tree (Eq.5, Eq.6), we set the score threshold  $\theta$  to 0.5, ensuring that the scores of the preference data are at least greater than 3.5, thereby maintaining the quality of the dataset.

When processing the seed data from ExTES, we filtered out entries that did not provide a strategy, contained incorrect strategies, or included the "Others" strategy. The remaining 15 strategies (excluding "Others") used in the ExTES dataset were mapped to the 8 strategies used in our dataset according to a set of predefined rules, which can be found in Table 8.

		Q↑	$\mathcal{B}\downarrow$	$Q_{\mathcal{W}} \uparrow$	R - L↑
LLaM	A-3.1-8B-Instruct	29.79	1.18	38.78	23.48
	SFT	30.28	2.65	37.33	23.77
	CSO-DPO	33.11	1.11	39.21	24.24
FuLL	CSO-SimPO	29.12	1.53	36.27	23.59
Tubb	CSO-IPO	35.48	1.04	37.74	24.19
	CSO-KTO	32.60	0.88	36.63	25.84
	CSO-ORPO	30.46	1.14	33.24	26.06
	SFT	31.25	2.65	39.27	23.30
	CSO-DPO	34.51	1.11	41.11	23.89
LoRA	CSO-SimPO	33.43	1.04	40.55	23.41
LUKA	CSO-IPO	33.00	1.13	39.40	23.55
	CSO-KTO	32.80	1.04	38.15	23.70
	CSO-ORPO	31.50	1.17	39.04	23.71

Table 9: The Results of different preference optimization algorithms on LLaMA-3.1-8B-Instruct.

# **F.2** Experimental Details

Our experiments are implemented with PyTorch (Paszke et al., 2019) on 8 NVIDIA Tesla A100 using DeepSpeed (Rasley et al., 2020) repository with ZeRo-2 optimization. We performed both full-parameter and LoRA fine-tuning on LLaMA-3.1-8B-Instruct (Dubey et al., 2024), Qwen-2.5-7B-Instruct (Yang et al., 2024) and Gemma-2-9B-it (Team et al., 2024). For all experiments, we set maximum target length of 512 tokens across all backbones. LoRA fine-tuning is conducted with an alpha of 8, and a dropout rate of 0, targeting all modules. All backbones are trained using their respective official chat templates.

More detailed settings across different backbones are listed in Table 14 and Table 15.

#### **G** Additional Experimental Results

# G.1 Analysis of Different Preference Optimization Algorithms

In addition to preference training based on DPO, we also investigated the performance of several other preference optimization algorithms, including SimPO (Meng et al., 2024), IPO (Azar et al., 2024), KTO (Ethayarajh et al., 2024), and ORPO (Hong et al., 2024). We tested these algorithms in both LoRA and full fine-tuning settings, maintaining the training parameters consistent with DPO, with adjustments made only to algorithm-specific parameters. For instance, the hyperparameter  $\beta$  was adjusted to suit each algorithm, and algorithm-specific parameters, such as Gemma-2 for SimPO, chosen weight and rejected weight for KTO, were tuned accordingly. A detailed list of the hyperparameters used for each algorithm can be found in

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$Q_{\mathcal{W}} \uparrow$	R - L↑
Qwen-2.5-7B-Instruct		19.84	2.47	28.12	23.52
	SFT	21.73	2.34	31.24	23.54
	CSO-DPO	28.78	1.92	34.39	26.16
FuLL.	CSO-SimPO	23.51	2.01	34.84	21.54
Tull	CSO-IPO	25.10	2.15	36.54	25.16
	CSO-KTO	25.56	1.56	38.44	26.65
	CSO-ORPO	22.20	1.24	35.91	24.58
	SFT	21.54	2.45	29.11	23.72
	CSO-DPO	23.16	2.09	32.26	24.17
LoRA	CSO-SimPO	25.91	2.02	30.45	23.32
LUKA	CSO-IPO	22.71	2.21	28.75	23.49
	CSO-KTO	22.91	2.17	30.82	23.61
	CSO-ORPO	22.49	2.12	28.42	23.35
	CSO-KTO	22.91	2.17	30.82	23.61

Table 10: The Results of different preference optimization algorithms on Qwen-2.5-7B-Instruct.

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$Q_{\mathcal{W}} \uparrow$	R - L↑
Gemm	a-2-9b-it	31.31	1.33	44.06	25.64
	SFT	32.52	1.29	46.45	25.25
	CSO-DPO	35.61	1.54	47.95	26.63
Full	CSO-SimPO	26.67	2.03	48.03	25.60
run	CSO-IPO	32.02	1.70	45.29	25.81
	CSO-KTO	39.73	0.81	48.87	27.84
	CSO-ORPO	34.80	1.24	48.28	27.52
	SFT	31.40	1.55	43.90	25.68
	CSO-DPO	35.77	1.23	52.34	26.61
LoRA	CSO-SimPO	34.95	1.19	51.62	26.22
LUKA	CSO-IPO	34.16	1.29	51.62	26.20
	CSO-KTO	35.89	0.99	48.53	26.45
	CSO-ORPO	32.35	1.30	48.14	25.85

Table 11: Result of different preference optimization algorithms on Gemma-2-9b-it.

the Table 14 and Table 15.

For all algorithms except KTO, we used the ESC-Pro dataset as the training set, which contains preference pairs directly derived from the original ESC-Pro data. In the case of KTO, the training set was derived by splitting each preference pair in the ESC-Pro dataset and removing duplicates. The experimental results are shown in Table 9, Table 10 and Table 11. From the results, it is evident that all tested preference optimization algorithms performed effectively after training on the ESC-Pro dataset, with some methods achieving higher performance than DPO. These results validate the efficacy and versatility of the ESC-Pro dataset for optimizing preference-based dialogue strategies.

#### **G.2** Results of Large-Scale Backbone

We conduct additional experiments on Qwen2.5-32B-Instruct using LoRA-based preference optimization. As shown in Table 12, CSO-LoRA consistently outperforms all baselines across all

	$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$Q_{\mathcal{W}} \uparrow$	R - L↑
Qwen2.5-32B-Instruct	37.37	1.40	41.97	24.37
Direct-Refine	34.87	1.68	40.69	23.14
Self-Refine	14.63	0.98	26.97	21.74
w/ Example	20.66	2.65	24.40	21.93
SFT-LoRA	37.69	1.60	42.09	24.40
CSO-LoRA	38.53	1.29	43.95	24.65

Table 12: Performance comparison on Qwen2.5-32B-Instruct using LoRA-based preference optimization.

four metrics, and long-term user satisfaction (R-L). Compared to standard supervised fine-tuning (SFT-LoRA), CSO-LoRA improves strategy quality by +0.84 and reduces strategy bias while further enhancing user-centered metrics. These results confirm that CSO remains effective even at larger scales, demonstrating strong scalability and robustness in enhancing both adaptability and emotional intelligence in LLM-based emotional support.

# H Case Study

Table 16 presents a dialogue example from the ESC-Pro dataset. We show one round of dialogue with the prefer and non-prefer responses, along with the contextual background of the conversation, including the scene and description from the seed data, as well as the previous conversation history. By comparing the responses, we can see that the prefer response corresponds to the highest-quality strategy and reply, while the non-prefer responses are based on other strategies.

In Table 17, we provide a case study to demonstrate that CSO can implicitly preserve emotional continuity over multiple turns by selecting context-appropriate strategies. This case shows that CSO can implicitly maintain emotional consistency across turns by adapting strategies to the seeker's changing emotions. The supporter begins with validation and empathy, then gradually introduces more proactive strategies like planning and reframing as the seeker becomes more hopeful. This smooth progression show that CSO responds not just to the current emotion, but also supports the emotional flow of the conversation as a whole.

#### I Human Evaluation Details

We conducted human evaluations and outsourced the evaluation tasks to anonymous crowd workers. We employed three evaluators for the human evaluation process, all of whom had prior annotation experience to ensure reliability. Before conducting evaluations, they received targeted training to thoroughly understand the criteria for each metric, promoting consistency and accuracy in their judgments. Evaluators were compensated based on the number of annotated samples they completed. To assess annotation consistency, we report the interannotator agreement using Cohen's Kappa scores, which are provided in the final column of Tables 1 and 3. Detailed guidelines and instructions were provided for each rating criterion to ensure the accuracy and consistency of the evaluations, as outlined in appendix C.2. For each evaluation task, we required three evaluators to assess 100 samples based on the corresponding metrics.

Criteria	Empathy	Information	Humanoid	Strategies
4 points	The system exhibits a high degree of anthropomorphism, going so far as to console users in a	There are many suggestions, and all of them are effective.	There is no apparent difference from human friends.	The strategies are numerous, well-tailored to the user's emotional state, and demonstrate high em-
	friendly manner and assist them in analyzing the underlying logic of emotions.			pathy and effectiveness in addressing the user's concerns.
3 points	Providing emotional comfort during conversations and assisting users in analyzing the underlying logical framework of their emotions.	There are more than five suggestions, but some of them are ineffective. There are fewer than five suggestions, but all of them are very effective.	1-2 traces can reveal that the AI assistant is a lan- guage model.	More than five strate- gies are provided, but some lack empathy or relevance. Alternatively, fewer than five strategies are shared, but they are highly empathetic and di- rectly address the user's core emotional needs.
2 points	The lack of understanding of user emotions or the absence of mechanisms to analyze user emotions are the main factors.	The suggestions are fewer than five, and some suggestions are effective, while others provide numerous suggestions, but none of them touch the root of the problem.	More than two traces can reveal that the AI assistant is a language model.	Fewer than five strategies are provided, and they are a mix of relevant and irrelevant approaches. Alternatively, a large number of strategies are given, but they fail to address the user's emotional root issues.
1 point	The lack of understanding of user emotions and the absence of mechanisms to analyze user emotions are the main factors.	Have suggestions but ineffective, as well as no suggestions.	Structured responses, or responses in the form of 'As a large language model' or robot-like replies.	Strategies are present but lack empathy or relevance. Some may appear dismissive or insufficiently supportive in the context of the user's concerns.
0 points	The disregard for user concerns, the absence of assistance in analyzing user issues, and even the imposition of negative effects on user emotions.	Suggestions were provided, but all of them were ineffective, and some even gave advice that could potentially harm the user.	The dialogue exhibits rigidity and lacks comprehension in terms of internalizing the content.	Strategies are counterproductive, exacerbating the user's distress or dismissing their concerns. Some suggestions may inadvertently harm the user's emotional well-being.

Table 13: Evaluation criteria of reward LLM.

	Epoch	Batch Size	Learning Rate	
LLaMA-3.1-8B-Instruct				
SFT	1 3	32	5.0e-7	
CSO		32	5.0e-7	
SFT	3	32	1.0e-6	
CSO	3	32	1.0e-6	
Qwen-2.5-7B-Instruct				
SFT	1 3	32	5.0e-7	
CSO		32	6.0e-7	
SFT	3	32	1.0e-6	
CSO	1	128	5.0e-7	
Gemma-2-9b-it				
SFT	1	32	5.0e-7	
CSO	1	32	5.0e-7	
SFT	3 3	32	5.0e-7	
CSO		8	6.0e-7	
	SFT CSO  SFT CSO  7B-Instruct SFT CSO  SFT CSO  9b-it SFT CSO SFT	SFT 1 CSO 3  SFT 3 CSO 3  7B-Instruct  SFT 1 CSO 3  SFT 3 CSO 1  9b-it  SFT 1 CSO 1  SFT 3	SFT	

Table 14: Overall hyper-parameter settings.

		beta	gemma	chosen weight	rejected weigh
LLaM	A-3.1-8B-Instruct				
FuLL	DPO	0.7	-	-	=
	SimPO	3.3	2.0	-	-
	IPO	1.0	-	=	=
	KTO	0.9	-	1.0	0.5
	ORPO	1.3	-	-	=
	DPO	0.01	-	-	-
	SimPO	2.0	0.5	=	-
LoRA	IPO	0.1	-	=	-
	KTO	0.01	-	1.0	1.0
	ORPO	1.0	-	=	=
Qwen-2	2.5-7B-Instruct				
	DPO	0.1	-	=	=
FuLL	SimPO	1.5	3.5	-	-
	IPO	0.5	-	-	-
	кто	0.07	-	1.0	0.5
	ORPO	0.7	-	-	-
	DPO	0.1	-	-	=-
	SimPO	3.5	3.0	=	-
LoRA	IPO	0.01	-	-	-
	KTO	0.05	-	1.0	0.5
	ORPO	2.0	-	-	-
Gemma	a-2-9b-it				
	DPO	0.5	-	=	-
	SimPO	3.5	2.0	-	-
FuLL	IPO	1.0	-	-	-
	кто	0.5	-	1.0	0.6
	ORPO	0.06	-		-
LoRA	DPO	0.1	-	-	-
	SimPO	1.5	2.0	-	-
	IPO	0.2	-	-	-
	кто	0.5	-	1.0	0.3
	ORPO	0.02	-	_	_

Table 15: Hyper-parameter settings of other preference optimization methods.

Conversation Ba	ackground		
Scene	Breakups or Divorce		
Description	After five years of marriage, my husband and I have decided to separate due to irreconcilable differences. I am feeling overwhelmed and uncertain as I navigate through this difficult time.		
Conversation H	istory		
Seeker	Hey, I really need someone to talk to right now. Can you lend me an ear?		
Supporter	(Emotional Validation) Of course! I'm here for you. What's been going on?		
Seeker	My husband and I have decided to separate. It's been really hard to accept this reality.		
Supporter	(Emotional Validation) I'm so sorry to hear that. Separating from a long-term partner can be incredibly challenging. It's normal to feel overwhelmed and unsure about what lies ahead. Can you share more about your feelings?		
Seeker	I feel like everything fell apart so quickly. We had dreams about raising a family together and now it's all shattered.		
Supporter	(Emotional Validation) It's really tough when our dreams and plans change unexpectedly. You're experiencing a whirlwind of emotions right now. It's completely understandable that you're feeling devastated and lost. Know that you're not alone in going through this difficult process.		
Seeker	Thank you. I can't help but blame myself for the breakdown of our marriage. I keep wondering if I did something wrong.		
Supporter	(Emotional Validation) It's natural to feel a sense of self-blame during such a challenging time. However, I want to emphasize that relationships involve the contributions of both partners, and it's rarely a reflection of just one person's actions. You are worthy of love and happiness, regardless of the outcome of this marriage.		
Seeker	I appreciate your kind words. It's just hard to let go of all the memories and years we spent together. I miss him so much.		
Supporter	(Empathetic Statements) It's completely understandable that you're missing your husband. Remember that healing takes time and it's okay to grieve the loss of your marriage. It's important to be patient with yourself and allow yourself to process the many emotions that may arise. Lean on your support system and engage in self-care activities that bring you comfort and joy.		
Seeker	I'll try my best. Right now, it feels impossible to imagine moving on. Do you have any advice on how to start the healing process?		
Responses			
Preferred	(Affirmation) While it may seem overwhelming now, know that healing and growth are possible. Focus on taking care of yourself both physically and emotionally. Engage in activities that bring you happiness and surround yourself with loved ones who uplift you. Consider seeking support from a therapist or counselor who can help guide you through this process. You are stronger than you realize, and I believe in your ability to navigate through this difficult time.		
Non-Preferred	(Collaborative Planning) Absolutely! Let's break it down together. What's one small step you think you could take this week to begin healing?		
Non-Preferred	(Avoid Judgment and Criticism) Absolutely! Let's break it down together. What's one small step you think you could take this week to begin healing?		
Non-Preferred	(Reframe Negative Thoughts) I get that moving on feels impossible right now. But think about it this way: healing is a journey, not a race. Each small step counts.		
Non-Preferred	(Share Information) Starting the healing process can be tough, but small steps can help. Journaling your feelings, talking to friends, or exploring new hobbies can be beneficial.		

Table 16: An example of preference dialogue in our ESC-Pro.

# Conversation Background Scene Depression and Low Mood I've been feeling emotionally exhausted and overwhelmed by everything happening in my life. It's hard to stay positive when it feels like things keep piling up and nothing changes. Role Emotion / Strategy Utterance

Role	Emotion / Strategy	Utterance
Seeker	Sadness	I've been feeling really overwhelmed lately. I don't know how to handle everything that's going on.
Supporter	Emotional Validation	It sounds like you're carrying a lot right now, and it's understandable that you're feeling overwhelmed. It's okay to feel that way.
Seeker	Frustration	Yeah, it's just so much, and I feel like I can't get a grip on it.
Supporter	Empathetic Statements	It's tough when everything seems like it's piling up. I can hear how much you're struggling with it.
Seeker	Hopelessness	I'm trying to stay positive, but it's hard.
Supporter	Affirmation	It's really impressive that you're trying to stay positive despite everything. That shows real strength.
Seeker	Overwhelm	I just don't know where to start.
Supporter	Collaborative Planning	Let's take it step by step. What's one small thing you could focus on right now?
Seeker	Норе	Maybe I can tackle some smaller tasks to get the ball rolling.
Supporter	Provide Different Perspectives	That's a great perspective. Starting small can help you feel more in control and build momentum.
Seeker	Defeat	I feel like I've been stuck in this cycle, and nothing changes.
Supporter	Reframe Negative Thoughts	It's easy to feel stuck when it seems like things aren't improving. But every small step forward counts, even if it doesn't feel like it.
Seeker	Acceptance	You're right. Maybe I need to be kinder to myself through this.
Supporter	Avoid Judgment and Criticism	Absolutely. Being kind to yourself is key. You're doing your best, and that's enough.
Seeker	Hopeful	I think I can manage it now, one step at a time.
Supporter	Emotional Validation	That's a wonderful mindset. Taking it one step at a time will help you regain your balance.

Table 17: Dialogue turns with emotional Context and support strategies