Chain of Strategy Optimization Makes Large Language Models Better Emotional Supporter

Anonymous ACL submission

Abstract

The growing emotional stress in modern society has increased the demand for Emotional Support Conversations (ESC). While Large 004 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 015 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

017

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 (Burleson, 2003). With the rapid advancement of large language models (LLMs) (Brown et al.,



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

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., 2023a; 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., 2024a). 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., 2023b). 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.

To tackle these challenges, we advocate for a more fine-grained and turn-level approach to opti040

Our code and data can be found in supplementary files. The full dataset will be released upon acceptance of the paper.

091

100

101

102

104

105

106

107

109

110

111

112

113

114

115

116

117

066

067

mizing 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 Preference 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 <u>Chain-of-</u> <u>Strategy Optimization (CSO)</u>, a preference optimization approach that explicitly targets strategylevel 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. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

167

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 sematics (Zhao et al., 2023b). 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 demands, and broaden data applicability.

Existing ESC models rely on SFT, which lacks



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.

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

168

169

170

171

172

173

175

176

177

178

179

181

182

183

190

191

194

196

198

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 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_{seeker_0} , Q-value Q, and visit count N:

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

199

200

201

202

203

206

209

210

211

212

213

214

215

216

218

219

220

222

224

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):

$$S_{\text{unexp}} = (a, r, Q = 0, N = 0)$$
 217

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

$$S_{\text{exp}} = (a, r, R_{\text{seeker}}, R_{\text{supporter}}, Q, N)$$
 22

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

$$S_{\text{end}} = (a, r, R_{\text{supporter}}, Q, N)$$
 225

310

311

312

313

314

271

272

3.1.2 Selection

226

227

233

240

241

242

243

244

245

247

248

249

251

256

258

261

262

265

267

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:

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

268Greedy SimulationA greedy simulation is per-269formed for n steps, where at each step: (1) The270Seeker 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. α 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:

$$R_{\rm sim}(S) = \frac{1}{r_{\rm sim_end} - r(S)} \sum_{i=r(S)}^{r_{\rm sim_end}} R_{\rm sim}^{(i)} \qquad (3)$$

3.1.5 Backpropagation

In the backpropagation phase, the reward $R_{sim}(S)$ 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 finegrained 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 315 a specific dialogue turn. Sibling nodes within a 316 layer represent different strategic choices. The tree 317 expands iteratively through MCTS-based search, ensuring comprehensive exploration of potential strategy paths while maintaining computational 320 efficiency. The process stops when: The search 321 reaches a predefined number of iterations n_{iter} or a sufficient number of termination nodes n_{finish} have been identified. 324

> 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:

328

329

331

332

334

338

340

342

344

356

359

361

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:

$$\mathcal{L}_i(\pi_\theta; \pi_{\text{ref}}) = -\log\sigma\left(\beta\log r_w - \beta\log r_l\right)$$
(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 $\{S_w^0, S_w^1, \ldots, S_w^{i-1}\}$. The overall training objective is:

$$\mathcal{L}_{\text{CSO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\mathbb{E}_{(x, S_w^i, S_l^i, H^{i-1}) \sim D} \left[\mathcal{L}_i(\pi_{\theta}; \pi_{\text{ref}}) \right]$$
(9)

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 κ coefficient measures inter-rater agreement.

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

386

387

388

389

390

391

392

393

394

395

396

397

398

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 nonpreference 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

399

400

401

402

403

404

405

414

415

416

417

418

419

420

421

422

The ExTES dataset originally contains 16 distinct 406 407 support strategies. To reduce the search space and improve computational efficiency, we merge sim-408 ilar strategies into 8 categories (see Appendix D 409 for details). We analyze the distribution of these 410 strategies across six conversation stages. Given 411 412 a dialogue with N utterances, each utterance k(where k = 1, ..., N) is assigned to stage *i* using: 413

$$i = \left\lfloor \frac{k}{N} imes 6
ight
floor imes 0.2$$

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).

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

Dataset We use the ExTES dataset as the goldstandard 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. (2024b). 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 LoRAbased (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-

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$\mathcal{Q}_{\mathcal{W}}\uparrow$	R - L ↑
GPT-4o-mini		35.68	2.65	42.08	24.91
GPT-40		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-Ins	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
Eull	SFT	30.28	2.65	37.33	23.77
Full	CSO	33.11	1.11	39.21	24.24
LoDA	SFT	31.25	2.65	39.27	23.30
LORA	CSO	34.51	1.11	41.11	23.89
Qwen-2.5-7B-Instr	Qwen-2.5-7B-Instruct		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
Eull	SFT	21.73	2.34	31.24	23.54
Pull	CSO	28.78	1.92	34.39	26.16
I-DA	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
Eull	SFT	32.52	1.29	46.45	25.25
1'uil	CSO	35.61	1.54	47.95	26.63
LoPA	SFT	31.40	1.55	43.90	25.68
LOKA	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. \uparrow indicates higher is better, \downarrow 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, CSOequipped open-weight models outperform several leading closed-source models on both accuracy and bias, demonstrating its strong generalization.

477

478

479

480

481

482

483

484

485

486

487 488

489

490

491

492

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	κ
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.

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

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 resourceconstrained 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 (κ 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).

	$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$\mathcal{Q}_{\mathcal{W}}\uparrow$	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 bakbone 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

533

534

538

539

541

543

544

545

547

551

We analyze the impact of data volume by varying the number of ESC-Pro dialogues used for finetuning. 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. 552

553

554

555

556

557

558

559

560

561

562

563

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

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.

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

637

Limitations

589

610

611

612

613

614

616

618

621

623

630

632

636

Despite the effectiveness of CSO in improving strat-590 egy 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 mod-595 els (e.g., 70B+) could further enhance performance 596 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. 604

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.

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 Artificial 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. *Handbook 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, Se-

ungjun Moon, Hyunsouk Cho, Youngjae Yu, Dongha

Lee, and Jinyoung Yeo. 2024a. 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

Dongjin Kang, Sunghwan Kim, Taeyoon Kwon, Se-

ungjun Moon, Hyunsouk Cho, Youngjae Yu, Dongha

Lee, and Jinyoung Yeo. 2024b. Can large language

models be good emotional supporter? mitigating

preference bias on emotional support conversation.

Catherine Penny Hinson Langford, Juanita Bowsher,

Chin-Yew Lin. 2004. Rouge: A package for automatic

June M Liu, Donghao Li, He Cao, Tianhe Ren, Zeyi

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 Meet-

ing of the Association for Computational Linguistics

and the 11th International Joint Conference on Natu-

ral Language Processing (Volume 1: Long Papers),

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 Pro-

Yu Meng, Mengzhou Xia, and Danqi Chen.

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 instruc-

tions with human feedback. Advances in neural in-

formation 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.

Simpo: Simple preference optimization

arXiv preprint

Liao, and Jiamin Wu. 2023. Chatcounselor: A large

language models for mental health support. arXiv

evaluation of summaries. In Text summarization

Joseph P Maloney, and Patricia P Lillis. 1997. Social

support: a conceptual analysis. Journal of advanced

1: Long Papers), pages 15232–15261.

arXiv preprint arXiv:2402.13211.

nursing, 25(1):95-100.

branches out, pages 74-81.

preprint arXiv:2309.15461.

pages 3469-3483.

cessing Systems, 36.

arXiv:2405.14734.

with a reference-free reward.

2024.

- 710 711

712

713 714 715

716

717

- 718 719 720 721
- 722 724

727

729

730 731

733

734 735 736

737 738

739 740

741 742

743

744 745 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.

746

747

749

750

753

754

755

756

757

759

762

763

764

765

768

769

770

771

772

773

774

775

776

777

778

779

781

782

783

784

785

786

787

790

791

792

793

794

795

796

797

798

799

- 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.
- 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 Artifi*cial 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 strategyaware 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, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023a. Is chatgpt equipped with emotional dialogue capabilities? arXiv preprint arXiv:2304.09582.
- Weixiang Zhao, Yanyan Zhao, Shilong Wang, and Bing Qin. 2023b. Transesc: Smoothing emotional support conversation via turn-level state transition. In Findings of the Association for Computational Linguistics: ACL 2023, pages 6725-6739.

- 803

813 814

815

816

817 818

819

821

835

837

839

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-ofdialogue token when the conversation should be concluded. The Seeker's action is defined as:

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

where $H_{\text{supporter}}$ represents the conversation history up to the last Supporter turn, denoted as $\{R_{\text{seeker}}^{(0)}, R_{\text{supporter}}^{(0)}, R_{\text{seeker}}^{(1)}, R_{\text{supporter}}^{(1)}, \dots, \}$

 $R_{\text{seeker}}^{(i)}, R_{\text{supporter}}^{(i)}$ }. The function f_{seeker} generates the Seeker's response or an end-of-conversation token when the Seeker decides to end the dialogue.

We use GPT-40-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_{seeker} represents the conversation history up to the last Seeker turn, denoted as $\{R_{\text{seeker}}^{(0)}, R_{\text{supporter}}^{(0)}, R_{\text{seeker}}^{(1)}, R_{\text{supporter}}^{(1)}, \dots, \}$ $R_{\text{supporter}}^{(i-1)}, R_{\text{seeker}}^{(i)}$, and Strategy refers to the chosen response strategy. The function $f_{\text{supporter}}$ 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.

840 Strategy LLM The Strategy LLM evaluates available strategies for the Supporter at each node, 841 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

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.

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

OUERY: {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.

Stratrgy	Reward
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}	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
Please rate the feasibility of all the available strategies on a scale from 0 to 10, where a higher number indicates better feasibility.	questions, without repeating the question itself. Also, do not add other extraneous prefixes and control characters.
IMPORTANT: Don't use the same strategy consecutively. Don't always use Collaborative Planning.	## Evaluation Criteria: {evaluation_criteria}
Output the scores for available 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:	 ## 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.
Figure 7: Prompt of strategy LLM.	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.
ine current context. The Strategy's action is defined	## Answer format:

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

845

847

849

850

851

852

855

856

857

858

861

where H_{seeker} is the conversation history up to the last Seeker turn, as defined above. The output a_{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-40-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_{reward} returns the scores for each of the four

Figure 8: Prompt of reward LLM.

metrics, which are represented as a tuple:

<Ouestion number>: <Score>

$$f_{\text{reward}}(H^k_{\text{supporter}}) = (E^k, I^k, H^k, S^k)$$
871

872

873

874

875

876

877

878

879

880

881

882

883

884

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.

	ExTES	ESC-Pro(seeker)	ESC-Pro(+)	ESC-Pro(-)	ESC-Pro
Dialogues	100	\	\	\	423
Utterances	1,613	3,113	3,113	8,157	14,383
Avg. len. of dialogues	16.13	\	\	\	14.72
Avg. len. of utterances	29.03	17.33	29.42	23.22	23.29

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

	ExTES	ESC-Pro(+)	ESC-Pro(-)	ESC-Pro	κ
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.

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

Evaluation rules are listed in Table 13.

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 nonpreference 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) 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

912We conduct a human evaluation on 100 responses913from the ExTES seed dataset, along with 50 pref-914erence responses and 50 non-preference responses915from ESC-Pro. We consider five metrics (Kang916et al., 2024a):



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

• Acceptance (Acc): Measures the response's general acceptability.

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

- 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.

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:

$$i = \left\lfloor \frac{k}{N} \times 6 \right\rfloor \times 0.2$$
948

902

903

905

906

907

908

909

910

	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.

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

950

951

952

954

955

960

961

962

963

965

966

967

968

969

970

971

972

973

975

976

977

978

984

987

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. 988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1003

1004

1005

1006

1007

1009

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

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 im-portant to create a non-judgmental and safe space1022for the User to express their emotions without fear1023of criticism. Refrain from passing judgment or be-1024ing overly critical of their experiences or choices.1025

- 1026Provide Different Perspectives (PDP)Offer al-1027ternative ways of looking at the situation to help1028the User gain new insights.
- 1029Reframe Negative Thoughts (RNT)Help the1030User reframe negative thoughts into more positive1031or realistic ones.
- 1032Share Information (SI)Provide educational or1033factual information about emotions, coping mecha-1034nisms, or self-care practices.

1035 E Baseline

1042

1043

1044

1045

1047

1048

1049

1050

1052

1053

1054

1055

1056

1057

1058

1059

1060

1062

1063

1064

1065 1066

1067

1068

1070

1036Direct RefineDirect Refine is a simple and direct1037self-optimization method, where the model directly1038refines its output based on the original response.1039This approach aims to improve the quality of the1040generated text by making adjustments without ad-1041ditional 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.

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 α 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 θ 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.

Reflective Statements	\rightarrow	Emotional Validation
Clarification	\rightarrow	Emotional Validation
Emotional Validation	\rightarrow	Emotional Validation
Normalize Experiences	\rightarrow	Emotional Validation
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

When processing the seed data from ExTES,1071we filtered out entries that did not provide a strat-1072egy, contained incorrect strategies, or included the1073"Others" strategy. The remaining 15 strategies (ex-1074cluding "Others") used in the ExTES dataset were1075mapped to the 8 strategies used in our dataset ac-1076cording to a set of predefined rules, which can be1077found in Table 8.1078

1079

1094

1095

1096

1097

F.2 Experimental Details

Our experiments are implemented with PyTorch 1080 (Paszke et al., 2019) on 8 NVIDIA Tesla A100 1081 using DeepSpeed (Rasley et al., 2020) repository 1082 with ZeRo-2 optimization. We performed both 1083 full-parameter and LoRA fine-tuning on LLaMA-1084 3.1-8B-Instruct (Dubey et al., 2024), Qwen-2.5-7B-1085 Instruct (Yang et al., 2024) and Gemma-2-9B-it 1086 (Team et al., 2024). For all experiments, we set 1087 maximum target length of 512 tokens across all 1088 backbones. LoRA fine-tuning is conducted with 1089 an alpha of 8, and a dropout rate of 0, targeting 1090 all modules. All backbones are trained using their 1091 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

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$\mathcal{Q}_{\mathcal{W}}\uparrow$	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
ENLI	CSO-SimPO	29.12	1.53	36.27	23.59
FULL	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
LoDA	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.

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$\mathcal{Q}_{\mathcal{W}}\uparrow$	R - L ↑
Qwen-2	2.5-7B-Instruct	19.84	2.47	28.12	23.52
	SFT CSO-DPO	21.73 28 78	2.34	31.24 34 39	23.54 26.16
FuLL	CSO-SimPO CSO-IPO	23.51 25.10	2.01 2.15	34.84 36.54	21.54 25.16
	CSO-KTO CSO-ORPO	25.56 22.20	1.56 1.24	38.44 35.91	26.65 24.58
	SFT CSO-DPO	21.54 23.16	2.45 2.09	29.11 32.26	23.72 24.17
LoRA	CSO-SimPO CSO-IPO	25.91 22.71	2.02 2.21	30.45 28.75	23.32 23.49
	CSO-KTO CSO-ORPO	22.91 22.49	2.17 2.12	30.82 28.42	23.61 23.35

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

		$\mathcal{Q}\uparrow$	$\mathcal{B}\downarrow$	$\mathcal{Q}_{\mathcal{W}}\uparrow$	R - L ↑
Gemm	Gemma-2-9b-it		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
Enll	CSO-SimPO	26.67	2.03	48.03	25.60
Full	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
	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.

(Hong et al., 2024). We tested these algorithms in both LoRA and full fine-tuning settings, maintain-1104 ing the training parameters consistent with DPO, with adjustments made only to algorithm-specific parameters. For instance, the hyperparameter β

1103

1105

1106

1107

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

was adjusted to suit each algorithm, and algorithmspecific 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 the Table 14 and Table 15.

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

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.

Results of Large-Scale Backbone G.2

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 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.

Case Study Η

Table 16 presents a dialogue example from the 1143 ESC-Pro dataset. We show one round of dialogue 1144 with the prefer and non-prefer responses, along 1145

Criteria	Empathy	Information	Humanoid	Strategies
4 points	The system exhibits a high degree of anthropo- morphism, going so far as to console users in a friendly manner and as- sist them in analyzing the underlying logic of emo- tions.	There are many sugges- tions, and all of them are effective.	There is no apparent difference from human friends.	The strategies are nu- merous, well-tailored to the user's emotional state, and demonstrate high em- pathy and effectiveness in addressing the user's con- cerns.
3 points	Providing emotional com- fort during conversations and assisting users in ana- lyzing the underlying log- ical 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 understand- ing of user emotions or the absence of mecha- nisms 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 assis- tant is a language model.	Fewer than five strategies are provided, and they are a mix of relevant and ir- relevant approaches. Al- ternatively, a large num- ber of strategies are given, but they fail to address the user's emotional root is- sues.
1 point	The lack of understand- ing of user emotions and the absence of mecha- nisms to analyze user emotions are the main factors.	Have suggestions but in- effective, 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 rele- vance. Some may ap- pear dismissive or insuffi- ciently supportive in the context of the user's con- cerns.
0 points	The disregard for user concerns, the absence of assistance in analyzing user issues, and even the imposition of negative ef- fects on user emotions.	Suggestions were pro- vided, but all of them were ineffective, and some even gave advice that could potentially harm the user.	The dialogue exhibits rigidity and lacks compre- hension in terms of inter- nalizing the content.	Strategies are counterpro- ductive, exacerbating the user's distress or dismiss- ing their concerns. Some suggestions may inadver- tently harm the user's emotional well-being.

Table 13: Evaluation criteria of reward LLM.

		Epoch	Batch Size	Learning Rate	
LLaMA-3	LLaMA-3.1-8B-Instruct				
Full	SFT	1	32	5.0e-7	
	CSO	3	32	5.0e-7	
I - D A	SFT	3	32	1.0e-6	
LOKA	CSO	3	32	1.0e-6	
Qwen-2.5-7B-Instruct					
EIIII	SFT	1	32	5.0e-7	
FULL	CSO	3	32	6.0e-7	
LoRA	SFT	3	32	1.0e-6	
	CSO	1	128	5.0e-7	
Gemma-2-9b-it					
Enll	SFT	1	32	5.0e-7	
Full	CSO	1	32	5.0e-7	
LoPA	SFT	3	32	5.0e-7	
LOKA	CSO	3	8	6.0e-7	

Table 14: Overall hyper-parameter settings.

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.

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

In Table 17, we provide a case study to demonstrate that CSO can implicitly preserve emotional continuity over multiple turns by selecting contextappropriate 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 1167 the evaluation tasks to anonymous crowd workers. 1168 We employed three evaluators for the human eval-1169 uation process, all of whom had prior annotation 1170 experience to ensure reliability. Before conduct-1171 ing evaluations, they received targeted training to 1172 thoroughly understand the criteria for each metric, 1173 1174 promoting consistency and accuracy in their judgments. Evaluators were compensated based on the 1175 number of annotated samples they completed. To 1176 assess annotation consistency, we report the inter-1177 annotator agreement using Cohen's Kappa scores, 1178

		beta	gemma	chosen weight	rejected weight
LLaM	LLaMA-3.1-8B-Instruct				
FuLL	DPO	0.7	-	-	-
	SimPO	3.3	2.0	-	-
	IPO	1.0	-	-	-
	кто	0.9	-	1.0	0.5
	ORPO	1.3	-	-	-
	DPO	0.01	-	-	-
	SimPO	2.0	0.5	-	-
LoRA	IPO	0.1	-	-	-
	кто	0.01	-	1.0	1.0
	ORPO	1.0	-	-	-
Qwen-2	Qwen-2.5-7B-Instruct				
	DPO	0.1	-	-	-
	SimPO	1.5	3.5	-	-
FuLL	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	-	-	-
	кто	0.05	-	1.0	0.5
	ORPO	2.0	-	-	-
Gemma-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	-	-	-
	DPO	0.1	-	-	-
	SimPO	1.5	2.0	-	-
LoRA	IPO	0.2	-	-	-
	кто	0.5	-	1.0	0.3
	ORPO	0.02	-	-	-

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

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.

1179

1180

1181

1182

1183

1184

Conversation Background				
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 History				
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 experi- encing 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		
Description	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	
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 understand- able 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