

# Dynamic Cognitive Planning for Cognitive-Functional Dialogue: A Case Study in Emotional Support Conversation

Jiaqi Liu, Yankun Yang, Jiakang Xu, Zhongqiang Du, Wenbin Jiang\*

School of Artificial Intelligence, Beijing Normal University, Beijing, China  
liujiaq1@mail.bnu.edu.cn, jiangwenbin@bnu.edu.cn

## Abstract

Cognitive-functional dialogues, such as those for persuasion, consultation, and question-answering, are prevalent throughout human social interaction. The core difference between these dialogues and casual chat lies in their objective: to guide a person's cognitive and psychological state toward a predetermined one. Existing conversational technologies perform poorly in handling such dialogues. The fundamental reason is that the transformation of human cognitive psychology follows specific patterns, yet existing technologies neither account for these patterns nor possess cognitive guidance planning based on them. This deficiency makes it difficult for dialogues to achieve their intended cognitive-functional goals effectively. To address this, we propose a dynamic cognitive planning method (DyCoP). By modeling the long-term evolution of a user's cognitive psychology during the dialogue process, this method dynamically generates dialogue guidance plans that align with the principles of cognitive-psychological evolution. This allows for the generation of appropriate dialogue responses based on prior user psychology and the immediate conversational context, thereby achieving cognitive-functional goals more efficiently and accurately. Simultaneously, we constructed an evaluation framework for cognitive-functional dialogues and constructed a richly annotated emotional support conversation dataset. Comprehensive automatic and human evaluations show that our proposed DyCoP method demonstrates significant advantages over existing baseline models.

## Introduction

Human society features a vast number of dialogue forms centered on cognitive intervention, such as persuasion in business negotiations, emotional counseling in psychotherapy, and knowledge clarification in educational settings. These can be collectively termed cognitive-functional dialogues. Unlike casual chats aimed at information exchange or emotional companionship, the core characteristic of cognitive-functional dialogues is their directionality: guiding the other party's cognitive psychology from an initial state to a preset target state through a series of logical linguistic interactions. While current dialogue systems have made progress (Chen et al. 2017; Quan et al. 2021; Zhang

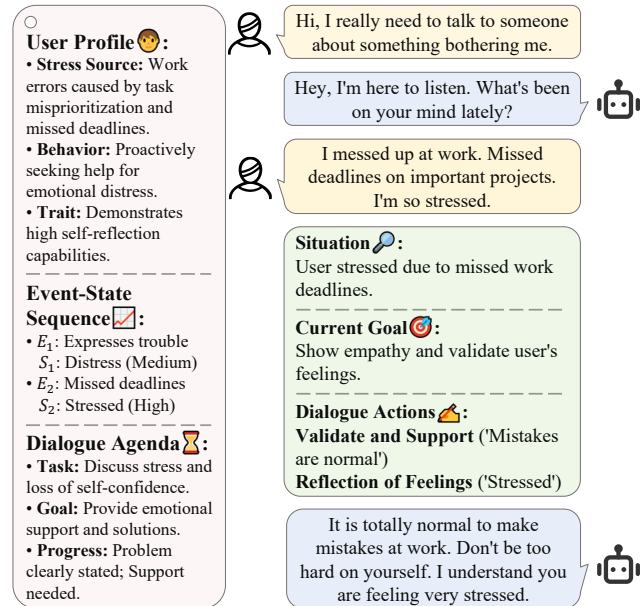


Figure 1: Example responses from our method. The figure shows how our method generates dialogue plans (green boxes on the right) based on a multi-dimensional understanding of the user (on the left) to generate responses (blue boxes on the right) that are both empathetic and policy-guided.

et al. 2020; Roller et al. 2020; Adiwardana et al. 2020), we observe that existing technologies fall short in cognitive-functional dialogues. This inadequacy is mainly manifested in two aspects: first, a reliance on surface-level linguistic patterns to generate responses, lacking long-term and dynamic planning for cognitive guidance logic; second, an inability to dynamically adjust strategies based on the other party's changing cognitive psychology, leading to inefficient cognitive guidance or deviation from the goal.

The root causes of the deficiencies in existing technologies for cognitive-functional dialogues can be attributed to both internal and external factors. Internally, the transformation of human cognitive psychology is not a random process but follows specific principles of cognitive-psychological

\*Corresponding Author

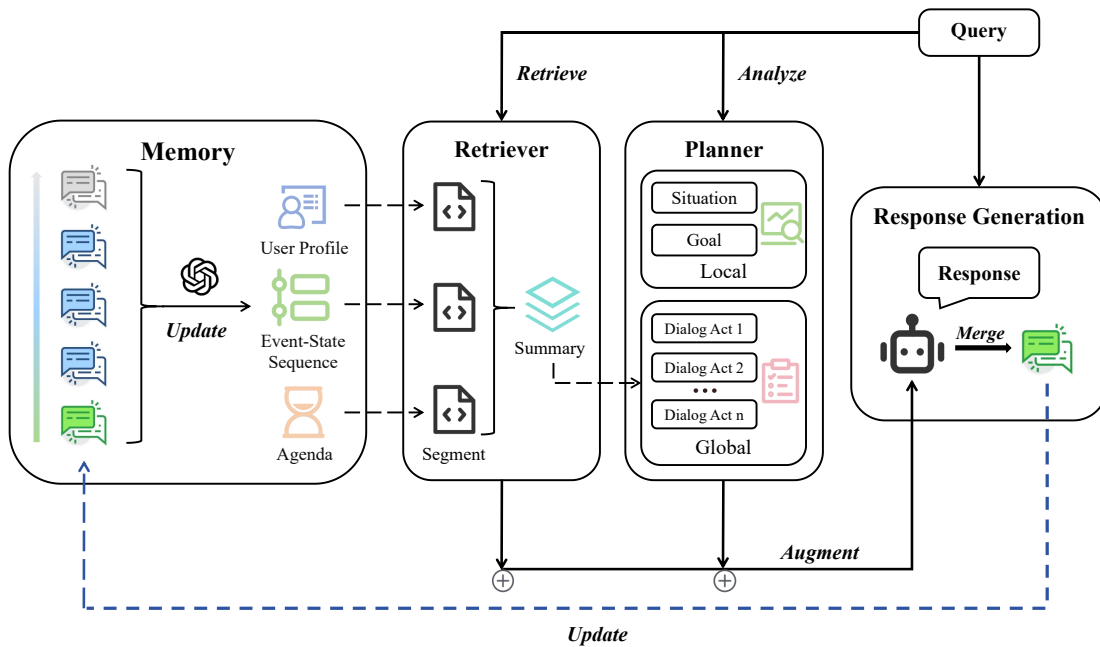


Figure 2: The overview architecture of DyCoP, including memory, retrieval, planning and generation.

evolution (Prochaska 2011). For instance, a student’s cognitive understanding in a Q&A scenario undergoes a psychological evolution from “phenomenal cognition” to “logical association” to “internalization and integration”. Similarly, an opponent’s attitude shift in a persuasion scenario involves a psychological evolution from “cognitive conflict” to “evidence weighing” to “stance adjustment”. However, current dialogue technologies have neither systematically modeled these evolutionary principles nor developed the capability for cognitive guidance planning based on them (Chen et al. 2017; Raheja and Tetreault 2019). Externally, the dialogue field lacks adequate supporting resources. On one hand, existing dialogue datasets primarily focus on chit-chat or task-based interactions, with few high-quality datasets annotating the evolutionary trajectory of user cognitive states. On the other hand, evaluation systems are still dominated by surface-level and local effectiveness metrics (Braggaar et al. 2024), lacking measures for deeper and global effects such as cognitive guidance efficiency and goal achievement, which hinders targeted technological optimization.

To address these issues, we creatively propose a dynamic cognitive planning method for cognitive-functional dialogues. This method analyzes user historical interaction data to model the user’s long-term cognitive psychological trajectory during the conversation process, thereby forming a long-term, personalized understanding of the user’s cognitive psychology and dynamically generates a dialogue guidance plan at each step that conforms to cognitive evolutionary laws, ensuring that each response serves the ultimate goal. We also constructed a comprehensive evaluation framework to measure the user’s deep-level cognitive-psychological transformation and the dialogue’s global ef-

fectiveness. Furthermore, we built a dataset for emotional support conversation (Rashkin et al. 2019), annotating the user’s cognitive evolution sequence, key events, and preset guidance goals to provide data support for related research. The primary contributions of our work are summarized as follows:

- We propose a dynamic cognitive planning method for cognitive-functional dialogue. By modeling the user’s cognitive-psychological evolution, it dynamically generates a cognitive guidance plan to effectively steer conversations toward preset functional goals.
- We also developed a comprehensive evaluation framework. It employs ten core metrics to holistically assess the efficiency and quality of cognitive guidance.
- Furthermore, we constructed a dataset of emotional support conversation with rich structured annotations, which provides an important resource for advancing research in this field.

## Cognitive-Functional Dialogue Task

### Task Definition

We use the term cognitive-functional dialogue to represent a class of dialogue forms whose core objective is to directionally guide the evolution of human cognitive-psychological states. Specifically, a cognitive-functional dialogue systematically intervenes in the conversational partner’s cognitive processes—including knowledge understanding, attitudes, and emotional states—through a series of logical linguistic interactions, guiding them from an initial cognitive state to a preset target cognitive state.

Cognitive-functional dialogues differ significantly from other dialogue types. First, unlike chit-chat, which aims

for emotional companionship or information exchange, it is strongly goal-oriented, with each conversational turn serving a long-term plan for cognitive guidance. Second, unlike task-oriented dialogues that aim to complete specific tasks like booking a flight or checking the weather, its core output is a change in cognitive psychology rather than the completion of a task result. It places greater emphasis on intervening in the internal cognitive processes of the dialogue partner.

## Evaluation Framework

The evaluation of cognitive-functional dialogues presents unique challenges, primarily due to the hidden nature of the evaluation object and the systematic nature of the evaluation dimensions.

Specifically, from the perspective of the evaluation object, its core objective—the evolution of cognitive states—is internal and implicit. It cannot be directly measured by task completion or surface text features, as in task-oriented dialogues. At the same time, from the perspective of evaluation dimensions, the dialogue’s goals rely on long-term strategic planning, demanding a holistic assessment. The value of the dialogue cannot be determined by the linguistic appropriateness or empathy of a single turn. Therefore, evaluation must transcend the local level to measure the overall contribution of each guidance step to the ultimate goal and the coherence of the strategy.

In summary, the evaluation of cognitive-functional dialogue requires an evaluation framework that captures the thorough evolution of cognitive psychology and encompasses the overall effectiveness of conversation guidance. We have designed a deep and comprehensive evaluation framework for cognitive-functional dialogue tasks, measuring both the efficiency of deep user cognitive psychological transformation and the overall effectiveness of the conversation. We propose a series of metrics designed to measure the core effectiveness of cognitive guidance: **Cohesion (Chr.)** for fluency of text; **Variety (Var.)** for response richness; **Empathy (Emp.)** for ability to empathize; **Utility (Uti.)** for suggestion quality; **Naturalness (Nat.)** for human-likeness; **Support (Sup.)** for overall supportive capabilities; **Preference (Pref.)** for subjective human preference; **Resolution (Res.)** for goal achievement; **Coherence (Coh.)** for logical consistency of content; and **Guidance (Gui.)** for proactivity and strategy. These metrics aim to capture deep cognitive and psychological changes and the global function of conversations. Detailed definitions and scoring criteria for each indicator are provided in the supplementary materials.

## Method

To achieve effective guidance in cognitive-functional dialogues, we innovatively propose the Dynamic Cognitive Planning method. This method uses an agent that integrates memory, retrieval, planning, and generation as its core framework, as shown in Figure 2.

The method consists of two interconnected core stages: Cognitive-Psychological Evolution Modeling and Dynamic Dialogue Plan Generation. The former is responsible for

continuously tracking and structurally representing the user’s internal state during the dialogue, providing the system with a profound understanding of the user. The latter, based on this understanding, formulates and adjusts dialogue strategies in real-time to ensure that the dialogue efficiently progresses toward the preset cognitive goal.

## Cognitive-Psychological Evolution Modeling

The objective of cognitive-psychological evolution modeling is to construct a dynamic, multidimensional representation of the user’s state, which not only records the surface-level dialogue history but also delves deeper to characterize the trajectory of the user’s internal psychological changes during the interaction, providing an interpretable and computable foundation in cognitive psychology for subsequent cognitive guidance planning. We achieve this by modeling two parallel dimensions:

**Psychological World Modeling** We describe the user’s inner mental states and their changes with a structured temporal cognitive-psychological representation. Adhering to cognitive psychology theories (Johnson-Laird 2010), we perform a fine-grained analysis to identify the key psychological state  $S_i$ , including its core emotion and intensity. Through this structured modeling, we can not only express deep cognitive-psychological states that far exceed single emotion labels but also capture the user’s long-term cognitive-psychological evolutionary trajectory.

**Physical World Modeling** We identify the external physical events that trigger the user’s psychological states, modeling them with a structured temporal event sequence. A Large Language Model performs Causal Event Attribution to analyze the user’s narrative and extract the specific event  $E_i$  causing the core emotion. This creates an objective and traceable link between external events and the user’s internal cognitive changes.

To fundamentally understand user context, we perform joint modeling to establish a causal link between external events (physical world) and internal states (psychological world), as an isolated understanding of the user’s actions or emotions is insufficient to support effective cognitive guidance. We construct a causal bridge binding a specific external event  $E_i$  to the resulting internal psychological state  $S_i$ . At any given time point  $i$ , we define an “event-state pair”  $U_i$  as the fundamental semantic unit of our joint model, formally represented as:

$$U_i = (E_i, S_i). \quad (1)$$

Connecting these pairs chronologically forms the user’s **Event-State Sequence**, which represents the complete evolutionary trajectory of the user’s cognitive psychology  $H_t$ :

$$\begin{aligned} H_t &= (U_1, U_2, \dots, U_t) \\ &= ((E_1, S_1), (E_2, S_2), \dots, (E_t, S_t)). \end{aligned} \quad (2)$$

This structured sequence precisely captures the causal logic from external events to internal states, providing a rich, machine-readable context for subsequent dialogue planning. The entire generation process of this cognitive-psychological evolution is automatically completed through the large model prompt method we designed.

Model	Chr.	Var.	Emp.	Uti.	Nat.	Sup.	Pref.	Res.	Coh.	Gui.	Total
Qwen	7.81	6.93	7.80	<u>6.83</u>	7.75	7.15	7.65	<u>7.29</u>	7.94	7.29	74.44
ESConv	6.27	5.20	5.63	4.45	5.43	5.00	5.29	5.06	6.31	5.31	53.95
ExTES	7.82	6.88	7.96	6.41	7.78	7.20	7.56	7.20	8.21	7.20	74.22
AUGESC	5.01	3.51	3.41	2.57	3.38	2.97	3.13	2.95	4.77	3.61	35.31
ServeForEmo	<b>8.00</b>	<u>7.30</u>	<b>8.12</b>	6.50	<b>8.06</b>	<u>7.38</u>	<u>7.68</u>	7.28	<b>8.36</b>	<u>7.32</u>	76.00
DyCoP	<u>7.94</u>	<b>7.35</b>	<u>8.11</u>	<b>7.18</b>	<u>8.01</u>	<b>7.61</b>	<b>7.93</b>	<b>7.66</b>	<u>8.29</u>	<b>7.69</b>	<b>77.77</b>

Table 1: Comparison of automatic evaluation results. All indicators except Total are scored on a 10-point scale, with higher scores being better. The best results are highlighted in **bold**. The second-best results are underlined. All automatic evaluation results are the average score of 3 runs.

## Dynamic Dialogue Plan Generation

Dynamic dialogue plan generation aims to generate and dynamically update a structured plan in real-time to guide subsequent response generation, based on the dialogue goal and the user’s cognitive-psychological evolution information. The plan expresses the long-term path from the user’s current cognitive state to the target cognitive state as a series of executable cognitive guidance steps, ensuring that the conversation proceeds along an optimal path that is consistent with the user’s cognitive-psychological development to achieve the dialogue goal.

The dialogue plan, also referred to as the cognitive guidance plan, is composed of a series of executable dialogue micro-strategies, also known as Dialogue Acts (DAs). Based on the emotional support conversation task, inspired by (Liu et al. 2021) and (Hill 2014), we have predefined a system of 7 core types of dialogue acts. The specific definitions are as follows:

- **Conversation Management:** Manages the macroscopic flow of a conversation, such as initiating greetings, changing topics, or ending the conversation, to ensure naturalness and coherence in interaction.
- **Questioning and Clarification:** Asks targeted questions to gather key information and clarify vague details, aiming for a deep and accurate understanding of the user’s situation and emotional state.
- **Restatement and Paraphrasing:** Verifies and confirms the system’s accurate understanding of factual content by restating the objective content of the user’s statements.
- **Reflection of Feelings:** Identifies and responds to the underlying emotions within the user’s utterances, aiming to confirm the perception of the emotional state and convey empathy.
- **Validation and Support:** Expresses clear understanding, acceptance, and affirmation of the user’s feelings, perspectives, or behaviors to build a safe conversational environment and enhance the user’s self-worth.
- **Providing Information:** Provides neutral, objective factual knowledge or universal viewpoints to help users expand their cognitive boundaries.
- **Guidance and Suggestion:** Propose specific action plans or problem-solving ideas to users to guide their cognition or behavior in a positive direction.

The types and definitions of nodes in the cognitive guidance plan are detailed in the supplementary material. The structure of the cognitive guidance plan is similar to the question decomposition semantic graph (Wolfson et al. 2020) and embodied intelligence task planning (Wu et al. 2023; Ahn et al. 2022), ensuring the logic and coherence of the guidance actions in a clear, structured manner.

The generation and application of this cognitive guidance plan follows a dynamic, real-time pipeline, which consists of the following key stages:

**Memory Retrieval and Compression** Upon receiving a new user query, the Retriever module is activated. Let the user query at turn  $t$  be  $Q_t$ , and the agent’s memory state from the previous turn be composed of the User Profile  $P_{t-1}$  (a long-term summary of the user), Event-State Sequence  $S_{t-1}$ , and Dialogue Agenda  $A_{t-1}$  (Ideal goals of the dialogue and current progress). The module first retrieves the most relevant raw memories  $M_{\text{recalled}}$  from the comprehensive memory storage based on the current query and memory state (Wang, Yang, and Wei 2023). This retrieval operation can be expressed as:

$$M_{\text{recalled}} = \mathcal{F}_{\text{recall}}(Q_t, P_{t-1}, S_{t-1}, A_{t-1}, H_{t-1}), \quad (3)$$

where  $H_{t-1}$  represents the recent dialogue history. The retrieved raw information is then fused and refined by a summarization model to generate a concise summary of recalled memories  $M_{\text{comp}}$ , which serves as the foundational context for subsequent steps.

$$M_{\text{comp}} = \mathcal{F}_{\text{compress}}(M_{\text{recalled}}). \quad (4)$$

**Multi-level Planning** The Planner module then receives the compressed memories  $M_{\text{comp}}$  and executes a multi-level planning process. First, the model analyzes this information along with the user’s latest input  $Q_t$  to understand the immediate context and define the goal for the current turn. This step generates a situation analysis  $C_t$  and a local goal  $G_{\text{local},t}$ . This process is represented by:

$$(C_t, G_{\text{local},t}) = \mathcal{F}_{\text{analyze}}(Q_t, M_{\text{comp}}, H_{t-1}). \quad (5)$$

After establishing the local goal, the Planner generates a specific plan  $\Pi_t$ . This planning step can be formulated as:

$$\Pi_t = \mathcal{F}_{\text{plan}}(Q_t, C_t, G_{\text{local},t}, M_{\text{comp}}). \quad (6)$$

The model is responsible for both selecting appropriate Dialogue Acts and populating their specific parameters within this plan.

**Plan-based Response Generation** Finally, the generated cognitive guidance plan  $\Pi_t$  is passed to the Response Generation module. This module transforms the structured, instruction-laden plan into a natural, fluent, and empathetic dialogue response  $R_t$ . This process is represented by:

$$R_t = \mathcal{F}_{\text{respond}}(\Pi_t, C_t, G_{\text{local},t}, Q_t, H_{t-1}). \quad (7)$$

After each turn, the agent’s memory  $P_t, S_t, A_t$  is updated based on the new dialogue segment  $(Q_t, R_t)$ . This structured process, from modeling and retrieval to multi-level planning, ensures that the generated response is grounded in a comprehensive understanding of the user’s cognitive state and guides the conversation toward the preset goal.

## Related Work

### Cognitive-Psychological Characterization

Modeling users’ internal states is fundamental to advanced human-computer interaction, with research focusing on emotion recognition and user profiling. Emotion recognition leverages deep learning to detect users’ immediate emotional states from various signals (Majumder et al. 2019; Xue and Bai 2024; Munikar, Shakya, and Shrestha 2019; Hazarika et al. 2018; Frye and Wilson 2022). In contrast, user profiling focuses on building long-term, static user representations (Eke et al. 2019; Wu et al. 2019; Zhang et al. 2018). Existing research focuses on capturing instantaneous states or static snapshots, whereas this study concentrates on capturing the dynamic evolution of users’ cognitive psychological processes.

### Dialogue Planning

Dialogue planning guides conversations toward specific goals, traditionally using statistical models like reinforcement learning for policy learning (Williams and Young 2007; Budzianowski and Vulić 2019). The advent of Large Language Models has transformed this paradigm, enabling more flexible and complex planning through techniques like Chain-of-Thought prompting (Wei et al. 2023; Ma, Jiang, and Huang 2025), reasoning-action frameworks (Yao et al. 2023b,a), and strategic generation (He et al. 2018; Park et al. 2023). These planning paradigms are primarily oriented toward transactional or informational tasks, while planning for cognitive-functional dialogue based on psychological principles remains a direction requiring deeper exploration.

### Global Evaluation of Dialogues

Evaluating generative dialogue systems is a recognized challenge, as traditional lexical overlap metrics correlate poorly with human judgment (Liu et al. 2017). This has led to more insightful methods, such as reference-free, learnable models for multi-dimensional quality assessment (Mehri and Eskenazi 2020). More recently, using Large Language Models as evaluators (LLM-as-a-judge) has shown high agreement with human experts (Zheng et al. 2023b; Zhou, Chen, and Yu 2024). Existing frameworks mainly measure general conversational properties. This study posits that evaluating cognitive-functional dialogue requires an assessment framework that directly measures the effectiveness of its core cognitive guidance function (Yalçın 2019).

## Experiments

To comprehensively evaluate the effectiveness of our proposed method (referred to as DyCoP) on cognitive-functional dialogue tasks, we designed a series of experiments. This chapter will detail the experimental setup, including the selected baseline models, the evaluation metric system, and the specific evaluation methods.

### Baseline Models

To verify the effectiveness of our method, we conduct a comprehensive comparison, detailed below:

**Qwen** (Yang et al. 2025): A foundational open-source LLM. In our experiments, we use the Qwen2.5-7B-Instruct version as the base model. We guide it to complete the cognitive-functional dialogue task through a task-specific prompt.

**ESConv** (Liu et al. 2021): Fine-tuned from Qwen using the ESConv dataset, which contains a large number of emotional support conversation from real-life scenarios.

**AUGESC** (Zheng et al. 2023a): Fine-tuned from Qwen using the AUGESC dataset, which is crafted by leveraging a fine-tuned dialogue LM to do dialogue completion task.

**ExTES** (Zheng et al. 2023c): Fine-tuned from Qwen using the ExTES dataset, which is generated by prompting a large language model to imitate and expand seed conversations.

**ServeForEmo** (Ye et al. 2024): Fine-tuned from Qwen using the ServeForEmo dataset, which was generated using a multi-role strategy-enhanced role-playing framework.

**DyCoP(ours)**: This is the dynamic cognitive planning approach proposed in this study that drives the system to produce high-quality responses that are both empathetic and goal-oriented.

### Evaluation Metrics

To comprehensively evaluate the model performance, we used the 10 core indicators defined above: Cohesion, Variety, Empathy, Utility, Naturalness, Support, Preference, Resolution, Coherence, and Guidance. Detailed scoring criteria are provided in the supplementary materials.

### Evaluation Method

To enable dynamic and scalable evaluation, we adopted the method of the ESC-Eval framework (Zhao et al. 2024), using a role-playing model as a user to interact with the dialogue systems being evaluated.

Specifically, we used the ESC-Role model proposed in that framework. ESC-Role is a specially trained language model capable of playing the role of a real person in emotional distress based on a preset Role Card. In our experiments, we configured the same pool of high-quality role cards for DyCoP and all baseline models. Subsequently, ESC-Role engaged in 5-turn dialogues with each model under test, generating a set of dialogue logs for each that reflect its true interaction capabilities. This process ensures that all models are tested under the same, realistic simulated scenarios.

Configuration	Chr.	Var.	Emp.	Ut.	Nat.	Sup.	Pref.	Res.	Coh.	Gui.	Total
Baseline	7.81	6.93	7.80	6.83	7.75	7.15	7.65	7.29	7.94	7.29	74.44
DyCoP w/o Memory	<b>7.94</b>	7.13	8.00	<b>6.99</b>	7.96	7.41	7.76	7.44	8.15	7.44	76.22
DyCoP w/o Retriever	<u>7.88</u>	7.02	7.93	6.89	7.75	7.36	7.53	7.40	8.10	7.46	75.31
DyCoP w/o Planner	<b>7.94</b>	<u>7.25</u>	<b>8.12</b>	6.97	<b>8.06</b>	<b>7.66</b>	<b>7.97</b>	<u>7.63</u>	<b>8.18</b>	<u>7.49</u>	<u>77.28</u>
DyCoP (Full)	<b>7.94</b>	<b>7.35</b>	<u>8.11</u>	<b>7.18</b>	<u>8.01</u>	<u>7.61</u>	<u>7.93</u>	<b>7.66</b>	<b>8.29</b>	<b>7.69</b>	<b>77.77</b>

Table 2: Ablation study based on the Qwen2.5-7B-Instruct model explores the impact of different DyCoP components on performance. All automatic evaluation results are the average score of 3 runs.

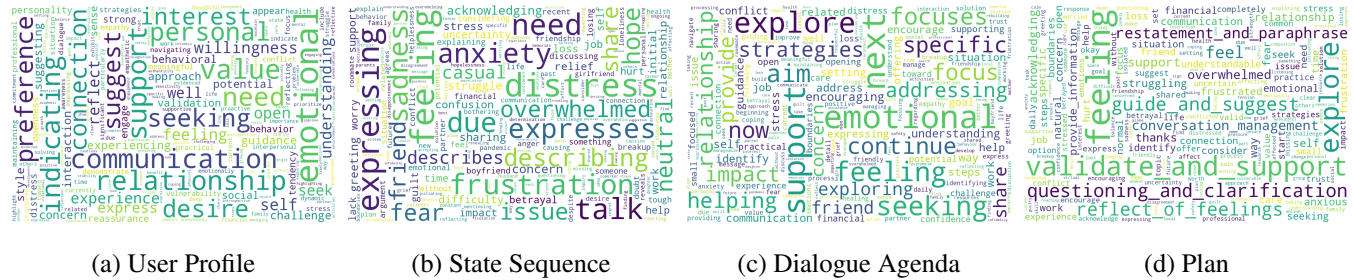


Figure 3: The word cloud of each component of dataset annotations.

After obtaining the conversation logs, we perform automatic and manual evaluations on them based on the evaluation metrics defined above.

**Automatic Evaluation** We utilized the powerful Gemini 2.5 Flash (Comanici et al. 2025) model as an automated judge. By designing specific evaluation prompts for each metric, we guided it to score the dialogue histories generated according to the criteria. This method allows for rapid, large-scale quantitative assessment of model performance.

**Human Evaluation** Considering the complexity and subjectivity of evaluating cognitive-functional dialogues, we recruited and trained professional evaluators for manual assessment. Evaluators scored the dialogues generated by all models independently without being informed of the model’s origin to minimize bias (Finch and Choi 2020).

## Dataset

Following the aforementioned evaluation method, we engaged the agent in multi-turn interactions with the ESC-Role user simulator to collect and build a richly annotated emotional support conversation dataset. This section provides a detailed description of the dataset’s composition, core characteristics, and processing pipeline. The statistics of the data set are shown in Table 3 and Table 4.

**Annotation Schema** For each turn in a dialogue, we not only saved the complete dialogue history but also recorded a series of internal state analyses produced by the agent during response generation. Consequently, each sample in the dataset is richly annotated, including the user profile, the user’s cognitive-psychological state sequence, the dialogue agenda, the plan, and the dialogue history. A detailed description of each component is available in the supplementary material.

Dialogue	Num. of Dialogues	331
	Num. of Turns	1966
	Avg. Turns per Dialogue	5.94
	Avg. Length per Dialogue	6727.14
Annotations (Avg. Length)	User Profile	815.39
	State Sequence	270.86
	Dialogue Agenda	213.89
	Plan	418.46

Table 3: Dialogue Data Statistics

Plan	Total Plans	1966
	Total Dialogue Act	9146
	Avg.Dialogue Act per Plan	4.65
Dialogue Act	Validation and Support	1963
	Questioning and Clarification	1762
	Guidance and Suggestion	1656
	Reflection of Feelings	1648
	Restatement and Paraphrase	1023
	Conversation Management	734
Providing Information	360	

Table 4: The Dialogue Act distribution of dataset.

**Data Processing** After obtaining the raw dialogue logs, we employed a multi-stage cleaning and correction pipeline to ensure data quality. This process began with the utilization of a large language model to automatically prune redundant segments from dialogues once the user’s core needs had been met. Subsequently, a combination of the model and human experts collaboratively proofread and corrected all structured annotations.

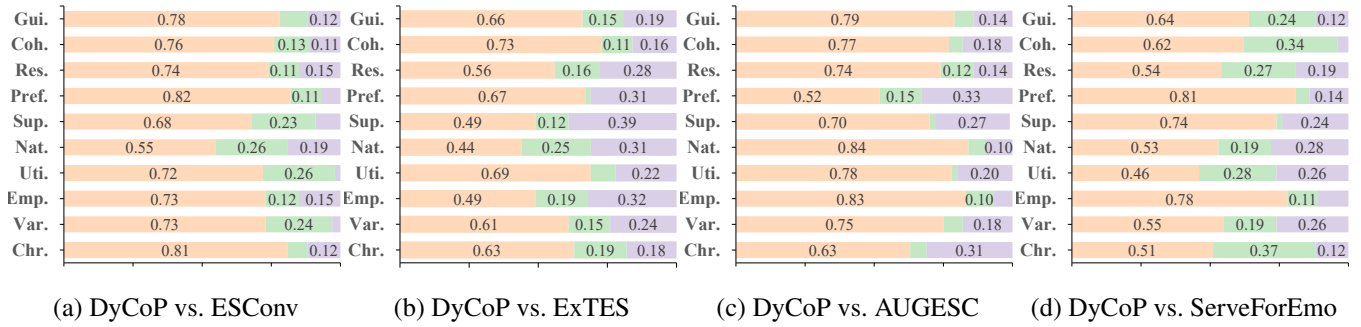


Figure 4: (a-d) Pairwise human evaluation results comparing DyCoP (A) against four baseline models (B). ■ indicates “A win”, ■ indicates “Tie”, and ■ indicates “B win”.

## Main Results

**Automatic Evaluation** The automatic evaluation results for all models are presented in Table 1. The results reveals that our proposed DyCoP achieved the best performance among all compared models. Its advantages are particularly prominent in core metrics that measure dialogue strategy and long-term goal achievement, such as Guidance, Resolution, and Utility. This strongly demonstrates that by modeling the user’s cognitive state and dynamically planning the dialogue path, our framework can more effectively guide the conversation toward preset cognitive-functional goals.

Notably, Qwen, which uses only prompt engineering, significantly outperformed several models that were fine-tuned on specific datasets like ESConv and AUGESC. This phenomenon indicates the limitations of conventional fine-tuning methods for cognitive-functional dialogue task. On the one hand, fine-tuning may lead to performance degradation due to issues like catastrophic forgetting or poor dataset quality. On the other hand, previous evaluation methods have been confined to superficial metrics such as BLEU and ROUGE. However, greater similarity to a dataset does not necessarily signify that a model possesses enhanced capabilities for cognitive guidance and emotional support.

In summary, the automatic evaluation results provide initial validation of our method’s superiority.

**Human Evaluation** To enhance the comprehensiveness and reliability of our evaluations, we also conduct human evaluations. Specifically, we engage three master students with psychology and computational linguistics backgrounds as annotators, human evaluators are instructed to compare the entire dialogues A and B generated by two models after interacting with the same virtual user (an ESC-Role configured with the same Role Card), and choose an option from “A wins”, “Tie”, and “B wins”. For the sake of fairness, we randomize the order of the dialogues to eliminate position bias. The evaluation is conducted on 50 randomly sampled dialogue pairs.

As illustrated, our method DyCoP demonstrates a clear and consistent superiority over all baselines. The results from the human evaluation are highly consistent with the automatic evaluation findings, further validating the effectiveness of our proposed framework. The consistent high perfor-

mance of our method in human preference underscores the value of DyCoP in achieving this goal.

**Ablation Study** To isolate the contribution of each component, we conducted an ablation study by systematically removing the Memory, Retriever, and Planner modules. As shown in Table 2, the full framework integrating all components achieved the best overall performance.

Removing any single component resulted in a noticeable degradation in performance, confirming the integral role of each. The most significant performance drop occurred upon removing the Retriever. Without the immediate contextual information it provides, the foundation for effective planning was undermined, critically impairing the agent’s performance. Ablating the Planner also led to a clear decline in performance, as the model’s ability to formulate structured, goal-oriented strategies was impaired, leading to weaker performance on key metrics like Guidance and Utility. Likewise, removing the Memory module, which prevents the agent from maintaining a long-term, evolving understanding of the user, also resulted in a discernible decrease in performance.

Collectively, these results demonstrate that the memory, retrieval, and planning stages are interconnected and indispensable. The integrity of this entire information processing pipeline is crucial for achieving high-performance cognitive guidance.

## Conclusion

In this paper, we innovatively propose a dynamic cognitive planning method (DyCoP) to address the challenges faced by current dialogue systems in handling cognitive-functional dialogues. This method models the evolution of user cognitive psychology in real time and dynamically generates a structured conversation guidance plan based on this model, thereby more accurately and efficiently guiding the conversation toward the pre-set goal. We propose a novel evaluation framework and have built a comprehensively annotated dataset for emotional support conversation. Extensive automatic and human evaluations provide clear evidence that our proposed method shows superior performance over existing baselines across metrics, including dialogue strategy and long-term objective fulfillment.

## Acknowledgments

This work was supported by the National Key Research and Development Program of China (No. 2024YFE0203000), the National Natural Science Foundation of China (No. 62437001), and the Fundamental Research Funds for the Central Universities.

## References

- Adiwardana, D.; Luong, M.-T.; So, D. R.; Hall, J.; Fiedel, N.; Thoppilan, R.; Yang, Z.; Kulshreshtha, A.; Nemade, G.; Lu, Y.; and Le, Q. V. 2020. Towards a Human-like Open-Domain Chatbot. arXiv:2001.09977.
- Ahn, M.; Brohan, A.; Brown, N.; Chebotar, Y.; Cortes, O.; David, B.; Finn, C.; Fu, C.; Gopalakrishnan, K.; Hausman, K.; Herzog, A.; Ho, D.; Hsu, J.; Ibarz, J.; Ichter, B.; Irpan, A.; Jang, E.; Ruano, R. J.; Jeffrey, K.; Jesmonth, S.; Joshi, N. J.; Julian, R.; Kalashnikov, D.; Kuang, Y.; Lee, K.-H.; Levine, S.; Lu, Y.; Luu, L.; Parada, C.; Pastor, P.; Quiambao, J.; Rao, K.; Rettinghouse, J.; Reyes, D.; Sermanet, P.; Sievers, N.; Tan, C.; Toshev, A.; Vanhoucke, V.; Xia, F.; Xiao, T.; Xu, P.; Xu, S.; Yan, M.; and Zeng, A. 2022. Do As I Can, Not As I Say: Grounding Language in Robotic Affordances. arXiv:2204.01691.
- Braggaar, A.; Liebrecht, C.; van Miltenburg, E.; and Kraemer, E. 2024. Evaluating Task-oriented Dialogue Systems: A Systematic Review of Measures, Constructs and their Operationalisations. arXiv:2312.13871.
- Budzianowski, P.; and Vulić, I. 2019. Hello, It's GPT-2 – How Can I Help You? Towards the Use of Pretrained Language Models for Task-Oriented Dialogue Systems. arXiv:1907.05774.
- Chen, H.; Liu, X.; Yin, D.; and Tang, J. 2017. A Survey on Dialogue Systems: Recent Advances and New Frontiers. *ACM SIGKDD Explorations Newsletter*, 19(2): 25–35.
- Comanici, G.; Bieber, E.; Schaekermann, M.; Pasupat, I.; Sachdeva, N.; Dhillon, I.; Blistein, M.; Ram, O.; et al. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. arXiv preprint arXiv:2507.06261.
- Eke, C. I.; Norman, A. A.; Shuib, L.; and Nweke, H. F. 2019. A Survey of User Profiling: State-of-the-Art, Challenges, and Solutions. *IEEE Access*, 7: 144907–144924.
- Finch, S. E.; and Choi, J. D. 2020. Towards Unified Dialogue System Evaluation: A Comprehensive Analysis of Current Evaluation Protocols. In Pietquin, O.; Muresan, S.; Chen, V.; Kennington, C.; Vandyke, D.; Dethlefs, N.; Inoue, K.; Ekstedt, E.; and Ultes, S., eds., *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 236–245. 1st virtual meeting: Association for Computational Linguistics.
- Frye, R. H.; and Wilson, D. C. 2022. Comparative Analysis of Transformers to Support Fine-Grained Emotion Detection in Short-Text Data. *The International FLAIRS Conference Proceedings*, 35.
- Hazarika, D.; Poria, S.; Zadeh, A.; Cambria, E.; Morency, L.-P.; and Zimmermann, R. 2018. Conversational Memory Network for Emotion Recognition in Dyadic Dialogue Videos. In Walker, M.; Ji, H.; and Stent, A., eds., *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2122–2132. New Orleans, Louisiana: Association for Computational Linguistics.
- He, H.; Chen, D.; Balakrishnan, A.; and Liang, P. 2018. Decoupling Strategy and Generation in Negotiation Dialogues. arXiv:1808.09637.
- Hill, C. E. 2014. *Helping Skills: Facilitating Exploration, Insight, and Action*. Washington, D.C. (Implied for APA): American Psychological Association, 4th edition.
- Johnson-Laird, P. 2010. Mental Models in Cognitive Science. *Cognitive Science*, 4: 71 – 115.
- Liu, C.-W.; Lowe, R.; Serban, I. V.; Noseworthy, M.; Charlin, L.; and Pineau, J. 2017. How NOT To Evaluate Your Dialogue System: An Empirical Study of Unsupervised Evaluation Metrics for Dialogue Response Generation. arXiv:1603.08023.
- Liu, S.; Zheng, C.; Demasi, O.; Sabour, S.; Li, Y.; Yu, Z.; Jiang, Y.; and Huang, M. 2021. Towards Emotional Support Dialog Systems. arXiv:2106.01144.
- Ma, X.; Jiang, W.; and Huang, H. 2025. Problem-Solving Logic Guided Curriculum In-Context Learning for LLMs Complex Reasoning. In *Findings of the Association for Computational Linguistics: ACL 2025*, 8394–8412. Association for Computational Linguistics.
- Majumder, N.; Poria, S.; Hazarika, D.; Mihalcea, R.; Gelbukh, A.; and Cambria, E. 2019. DialogueRNN: An Attentive RNN for Emotion Detection in Conversations. arXiv:1811.00405.
- Mehri, S.; and Eskenazi, M. 2020. USR: An Unsupervised and Reference Free Evaluation Metric for Dialog Generation. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 681–707. Online: Association for Computational Linguistics.
- Munikaar, M.; Shakya, S.; and Shrestha, A. 2019. Fine-grained Sentiment Classification using BERT. In *2019 Artificial Intelligence for Transforming Business and Society (AITB)*, volume 1, 1–5.
- Park, J. S.; O'Brien, J. C.; Cai, C. J.; Morris, M. R.; Liang, P.; and Bernstein, M. S. 2023. Generative Agents: Interactive Simulacra of Human Behavior. arXiv:2304.03442.
- Prochaska, J. O. 2011. Stages of change. *Journal of Clinical Psychology*, 67(8): 761–766.
- Quan, J.; Yang, M.; Gan, Q.; Xiong, D.; Liu, Y.; Dong, Y.; Ouyang, F.; Tian, J.; Deng, R.; Li, Y.; Yang, Y.; and Jiang, D. 2021. Integrating Pre-trained Model into Rule-based Dialogue Management. arXiv:2102.08553.
- Raheja, V.; and Tetreault, J. 2019. Dialogue Act Classification with Context-Aware Self-Attention. In Burstein, J.; Doran, C.; and Solorio, T., eds., *Proceedings of the 2019 Conference of the North American Chapter of the Association*

- for *Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 3727–3733. Minneapolis, Minnesota: Association for Computational Linguistics.
- Rashkin, H.; Smith, E. M.; Li, M.; and Boureau, Y.-L. 2019. Towards Empathetic Open-domain Conversation Models: a New Benchmark and Dataset. arXiv:1811.00207.
- Roller, S.; Dinan, E.; Goyal, N.; Ju, D.; Williamson, M.; Liu, Y.; Xu, J.; Ott, M.; Shuster, K.; Smith, E. M.; Boureau, Y.-L.; and Weston, J. 2020. Recipes for building an open-domain chatbot. arXiv:2004.13637.
- Wang, L.; Yang, N.; and Wei, F. 2023. Query2doc: Query Expansion with Large Language Models. arXiv:2303.07678.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Ichter, B.; Xia, F.; Chi, E.; Le, Q.; and Zhou, D. 2023. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. arXiv:2201.11903.
- Williams, J. D.; and Young, S. 2007. Partially observable Markov decision processes for spoken dialog systems. *Computer Speech Language*, 21(2): 393–422.
- Wolfson, T.; Geva, M.; Gupta, A.; Gardner, M.; Goldberg, Y.; Deutch, D.; and Berant, J. 2020. Break It Down: A Question Understanding Benchmark. arXiv:2001.11770.
- Wu, S.; Tang, Y.; Zhu, Y.; Wang, L.; Xie, X.; and Tan, T. 2019. Session-Based Recommendation with Graph Neural Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01): 346–353.
- Wu, Z.; Wang, Z.; Xu, X.; Lu, J.; and Yan, H. 2023. Embodied Task Planning with Large Language Models. arXiv:2307.01848.
- Xue, P.; and Bai, W. 2024. A Fine-Grained Sentiment Analysis Method Using Transformer for Weibo Comment Text. *Int. J. Inf. Technol. Syst. Approach*, 17(1): 1–24.
- Yalçın, N. 2019. Evaluating Empathy in Artificial Agents. In *2019 8th International Conference on Affective Computing and Intelligent Interaction (ACII)*, 1–7.
- Yang, Q. T. A.; Yang, B.; Zhang, B.; Hui, B.; Zheng, B.; Yu, B.; Li, C.; Liu, D.; Huang, F.; Wei, H.; Lin, H.; Yang, J.; Tu, J.; Zhang, J.; Yang, J.; Yang, J.; Zhou, J.; Lin, J.; Dang, K.; Lu, K.; Bao, K.; Yang, K.; Yu, L.; Li, M.; Xue, M.; Zhang, P.; Zhu, Q.; Men, R.; Lin, R.; Li, T.; Tang, T.; Xia, T.; Ren, X.; Ren, X.; Fan, Y.; Su, Y.; Zhang, Y.; Wan, Y.; Liu, Y.; Cui, Z.; Zhang, Z.; and Qiu, Z. 2025. Qwen2.5 Technical Report. arXiv:2412.15115.
- Yao, S.; Yu, D.; Zhao, J.; Shafran, I.; Griffiths, T. L.; Cao, Y.; and Narasimhan, K. 2023a. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. arXiv:2305.10601.
- Yao, S.; Zhao, J.; Yu, D.; Du, N.; Shafran, I.; Narasimhan, K.; and Cao, Y. 2023b. ReAct: Synergizing Reasoning and Acting in Language Models. arXiv:2210.03629.
- Ye, J.; Xiang, L.; Zhang, Y.; and Zong, C. 2024. SweetieChat: A Strategy-Enhanced Role-playing Framework for Diverse Scenarios Handling Emotional Support Agent. arXiv:2412.08389.
- Zhang, S.; Dinan, E.; Urbanek, J.; Szlam, A.; Kiela, D.; and Weston, J. 2018. Personalizing Dialogue Agents: I have a dog, do you have pets too? arXiv:1801.07243.
- Zhang, Y.; Sun, S.; Galley, M.; Chen, Y.-C.; Brockett, C.; Gao, X.; Gao, J.; Liu, J.; and Dolan, B. 2020. DialoGPT: Large-Scale Generative Pre-training for Conversational Response Generation. arXiv:1911.00536.
- Zhao, H.; Li, L.; Chen, S.; Kong, S.; Wang, J.; Huang, K.; Gu, T.; Wang, Y.; Jian, W.; Liang, D.; Li, Z.; Teng, Y.; Xiao, Y.; and Wang, Y. 2024. ESC-Eval: Evaluating Emotion Support Conversations in Large Language Models. arXiv:2406.14952.
- Zheng, C.; Sabour, S.; Wen, J.; Zhang, Z.; and Huang, M. 2023a. AugESC: Dialogue Augmentation with Large Language Models for Emotional Support Conversation. arXiv:2202.13047.
- Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E. P.; Zhang, H.; Gonzalez, J. E.; and Stoica, I. 2023b. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. arXiv:2306.05685.
- Zheng, Z.; Liao, L.; Deng, Y.; and Nie, L. 2023c. Building Emotional Support Chatbots in the Era of LLMs. arXiv:2308.11584.
- Zhou, R.; Chen, L.; and Yu, K. 2024. Is LLM a Reliable Reviewer? A Comprehensive Evaluation of LLM on Automatic Paper Reviewing Tasks. In *International Conference on Language Resources and Evaluation*.