

PFDial: A Structured Dialogue Instruction Fine-tuning Method Based on UML Flowcharts

Anonymous ACL submission

Abstract

Process-driven dialogue systems, which operate under strict predefined process constraints, are essential in customer service and equipment maintenance scenarios. Although Large Language Models (LLMs) have shown remarkable progress in dialogue and reasoning, they still struggle to solve these strictly constrained dialogue tasks. To address this challenge, we construct **Process Flow Dialogue (PFDial)** dataset, which contains 12,705 high-quality Chinese dialogue instructions derived from 440 flowcharts containing 5,055 process nodes. Based on PlantUML specification, each UML flowchart is converted into atomic dialogue units i.e., structured five-tuples. Experimental results demonstrate that a 7B model trained with merely 800 samples, and a 0.5B model trained on total data both can surpass 90% accuracy. Additionally, the 8B model can surpass GPT-4o up to 43.88% with an average of 11.00%. We further evaluate models' performance on challenging backward transitions in process flows and conduct an in-depth analysis of various dataset formats to reveal their impact on model performance in handling decision and sequential branches.

1 Introduction

Process-driven dialogue systems (Yi et al., 2024), as a special type of task-oriented dialogue systems, play a crucial role in various real-world applications, particularly in scenarios such as customer service, equipment maintenance, and medical consultation, where strict adherence to predefined process constraints is essential. In these contexts, dialogue systems must navigate through complex decision trees while maintaining precise control over the conversation flow, ensuring both compliance with established procedures and effective user interaction. These process flows typically contain two types of branches: sequential branches that follow a linear progression, decision branches that require conditional routing based on user input.

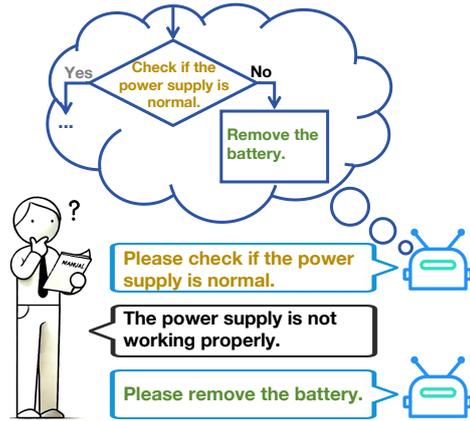


Figure 1: In this scenario, the system interacts with the user based on a flowchart that checks whether the power supply is functioning properly. If the power supply is faulty, the system guides the user to remove the battery. This interaction follows the decision-making process outlined in the flowchart.

The emergence of Large Language Models (LLMs) has brought unprecedented capabilities in natural language understanding and generation, demonstrating remarkable performance in both dialogue and reasoning tasks. These advances suggest potential solutions for process-driven dialogue systems (Yi et al., 2024; Zhang et al., 2023; Wu et al., 2020). However, our empirical evaluation reveals that even state-of-the-art (SOTA) LLMs like GPT-4o (OpenAI, 2023) struggle to consistently maintain process constraints while engaging in dialogue. Specifically, these models often deviate from predefined process constraints, make incorrect state transitions, or fail to properly handle complex decision branches, highlighting the need for more specialized solutions.

To address this challenge, we construct **Process Flow Dialogue (PFDial)** dataset, which contains 12,705 high-quality dialogue instructions derived from 440 flowcharts containing 5,055 process

nodes. Based on PlantUML specification, each UML flowchart is converted into atomic dialogue units, forming structured five-tuples (flowchart description, current state, user input, next state, robot output). This structured representation enables models to learn precise state transitions while maintaining natural dialogue capabilities. Through supervised fine-tuning (SFT) on PFDial, models can acquire strong controlled reasoning capabilities for process flows effectively following the prescribed state transitions and decision logic.

We conducted comprehensive experiments to address four key questions:

(1) How does our approach perform compared to SOTA LLMs? Our main experimental results demonstrate that models with varying parameter sizes can achieve excellent results after SFT on total data of PFDial. For instance, even a 0.5B model can achieve accuracy of 98.99% and 92.79% on in-domain and out-of-domain tests, respectively. More impressively, an 8B model achieves 97.02% accuracy on out-of-domain tests, surpassing GPT-4o by 11.00%, with over 43.88% improvement in decision branches.

(2) How does model performance scale with training data size? Our data scaling experimental results show that a 7B model can surpass 90% accuracy with only 800 training samples. As the training dataset increases, the overall performance of the model continues to improve. This underscores PFDial’s exceptional effectiveness in enhancing the model’s capability for controlled reasoning tasks with minimal data.

(3) How effective is our approach in handling backward transitions? We evaluated the model’s performance on more challenging backward transitions in decision branches using our constructed dataset, PFDial-H. This specialized benchmark focuses on cases where the next state returns to a previous point in the process flow. These transitions are particularly challenging due to their relative scarcity and the complex reasoning they require. Results on PFDial-H further validate our approach’s superiority in challenging controlled reasoning tasks.

(4) How do different state representation formats affect model performance? Through systematic analysis of three different dataset formats, we provide interpretable insights into how different formats impact model performance on decision and sequential branches, offering valuable guidance for future research.

Statistics	Train	ID Test	OOD Test
Flowcharts	440	80	80
State Nodes	5055	902	1262
Sequential Samples	9029	1628	2265
Decision Samples	3676	645	698
Dialogue Samples	12705	2273	2963
Avg. Length	277.16	270.57	326.10

Table 1: Statistics of the PFDial Dataset

Overall, our contributions can be summarized as follows:

- We have developed the Process Flow Dialogue (PFDial) dataset, which is derived from 440 flowcharts encompassing 5,055 process nodes. This dataset contains 12,705 high-quality dialogue instructions, serving as a valuable resource for training process-driven dialogue systems.
- The comprehensive experiments demonstrate that the PFDial dataset is highly effective. Even models with a smaller number of parameters (e.g., 0.5B, 1B, 1.5B) or those trained on relatively limited data (800 training samples) can achieve high accuracy. An 8B model achieves 97.02% accuracy on out-of-domain tests, surpassing GPT-4o by 11.00%, with over 43.88% improvement in decision branches.
- We demonstrate our model’s superior performance on complex backward transitions in decision branches using the PFDial-H benchmark, highlighting its capability in handling rare and intricate reasoning tasks. Additionally, we provide insights into the impact of different dataset formats on model performance, offering guidance for future research.

2 PFDial

In this section, we introduce PFDial, a Chinese dataset specifically designed for process-driven dialogue systems.

2.1 Dataset Overview

2.1.1 PFDial

The construction of PFDial follows a systematic pipeline, including flow chart collection, textual representation conversion, state transition Information extraction, prompt generation, and data validation. The dataset combines real-world scenarios

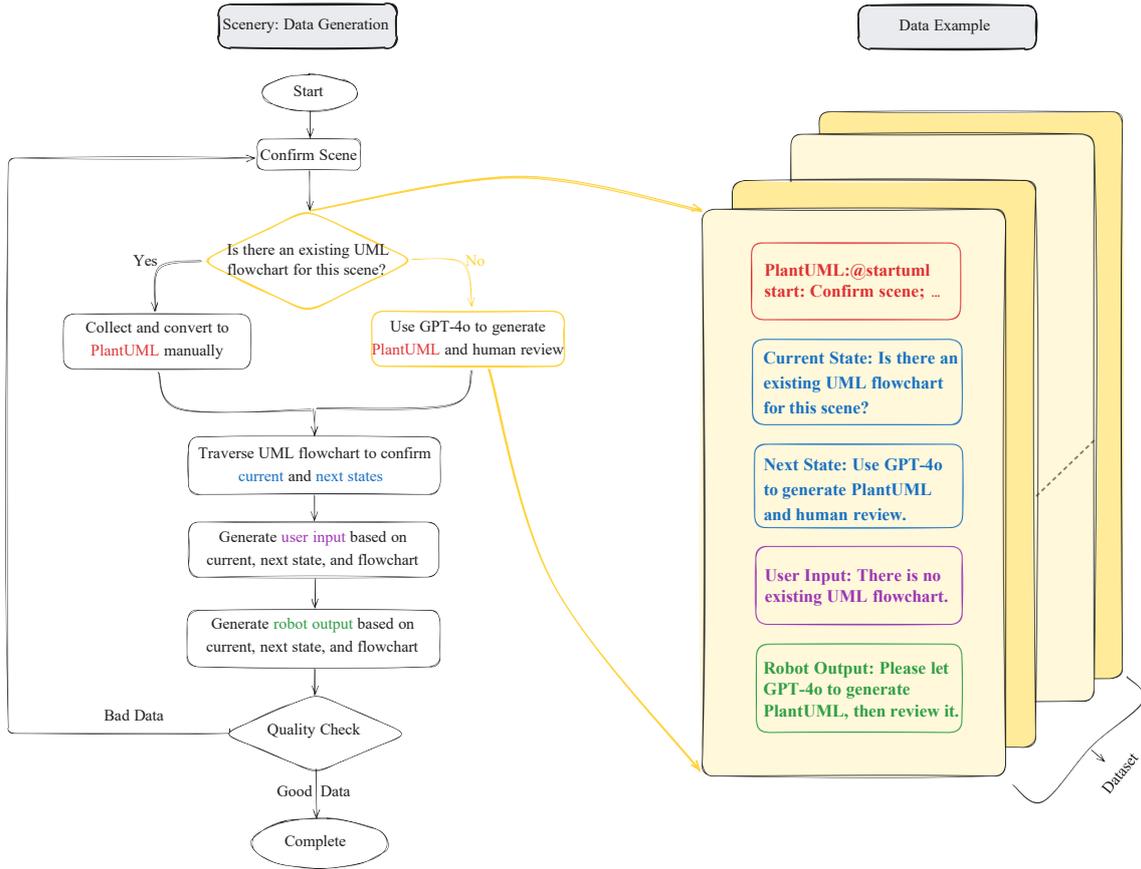


Figure 2: The left side illustrates the data construction process. The right side shows an example of the five-tuple dataset generated based on the leftside flowchart.

and synthetic data, achieving broad coverage of practical applications while maintaining high quality and diversity. Table 1 presents detailed statistics of the PFDial dataset.

2.1.2 PFDial-H

Considering the prevalence and importance of backward transitions in practical applications, we specifically constructed a supplementary dataset that incorporates backward transition mechanisms called PFDial-Hard (PFDial-H). By strategically adding backward transition nodes to existing flowcharts using GPT-4o, we implemented backward transition functionality for cases where conditions are not met. This improvement was applied to both out-of-domain test and training flowcharts, generating new training samples through the same prompt augmentation process. Table 5 in Appendix B.2 presents detailed statistics of the PFDial-H dataset.

2.2 Dataset Construction Process

2.2.1 Flow Chart Collection

Through extensive research, we identified and categorized 90 specific business scenarios, details of

which can be found in Appendix 16. Additionally, we designed a carefully constructed template dataset. After identifying the business scenarios, we collected the flowcharts manually or automatically, depending on whether pre-existing flowcharts were available. This process, combining automation with human review, significantly improved the efficiency and accuracy of UML flowchart generation.

2.2.2 Textual Representation Conversion

To efficiently convert flowcharts into machine-readable formats, we conducted a comparative study of several text-based representation schemes, including PlantUML¹, chart-mage², nomnoml³, and Mermaid⁴. After comprehensive evaluation, we selected PlantUML as the standard format. This decision was made for several reasons: First, PlantUML employs a code-like structured representation utilizing syntax features such as indentation,

¹<https://plantuml.com/>

²<https://chartmage.com/intro.html>

³<https://www.nomnoml.com/>

⁴<https://mermaid.js.org/>

194 branching, and loops, making flowchart descrip-
195 tions both intuitive to read and convenient for pro-
196 gram processing. Second, a preliminary experi-
197 ment using GPT-4o demonstrated that the Plant-
198 UML format exhibited superior accuracy com-
199 pared to other approaches. Details can be found in
200 Table 4 in Appendix B.1. During data processing,
201 we represented all flowcharts in PlantUML format
202 and generated standardized state nodes and their
203 transition relationships through parsing, providing
204 a unified representation for subsequent five-tuple
205 dataset generation.

2.2.3 State Transition Information Extraction

206 Based on the standardized PlantUML representa-
207 tion, we developed a specialized algorithm to ex-
208 tract complete state transition information, which
209 is shown in Algorithm 1 in Appendix A.1. During
210 this process, we got all existing paths and identi-
211 fied 5,055 distinct state nodes from the training
212 set. Each state transition pair (current state -> next
213 state) strictly corresponds to a specific path in the
214 flowchart, ensuring the accuracy and consistency
215 of the state information extraction.
216

2.2.4 Prompt Generation

217 To create quality training samples, we used the
218 GPT-4o model for bidirectional prompt augmen-
219 tation for each state transition. This involved gen-
220 erating user input and robot output based on the
221 current and next states, along with the flowchart.
222 Detailed prompts can be found in Appendix D.
223

2.2.5 Data Validation

224 To ensure the reliability of the dataset, we imple-
225 mented a rigorous multi-level validation process:
226 first, ensuring that all state nodes strictly corre-
227 spond to the original flowcharts; second, validating
228 the syntax correctness of the PlantUML; and finally,
229 checking the logical consistency between user in-
230 puts and state transitions. Any data that does not
231 meet these criteria will be regenerated.
232

233 This dataset construction methodology not only
234 ensures data quality and diversity but also pro-
235 vides a solid foundation for subsequent model train-
236 ing. Through systematic construction processes
237 and strict quality control, the PFDial dataset effec-
238 tively balances authenticity, standardization, and
239 scalability requirements.

3 Experiments

240 In this chapter, we present a series of comprehen-
241 sive experiments to address the four key questions
242 mentioned in Section 1. We conducted the main ex-
243 periment, data scaling experiment, backward tran-
244 sition studies, and format ablation studies respec-
245 tively.
246

3.1 Experimental Setup

247 **Base Models** We evaluate our method using two
248 series of base models with varying parameter sizes.
249 The first series includes Qwen2.5 models (Yang
250 et al., 2024) ranging from 0.5B to 7B param-
251 eters (Qwen2.5-0.5B, Qwen2.5-1.5B, Qwen2.5-3B,
252 and Qwen2.5-7B). The second series consists of
253 Llama3 models (Dubey et al., 2024) spanning from
254 1B to 8B parameters (Llama3.2-1B, Llama3.2-3B,
255 and Llama3.1-8B). This diverse selection of models
256 enables us to comprehensively analyze the impact
257 of model scale on performance.
258

259 **Baselines** For comparison, we select a compre-
260 hensive set of both open-source and proprietary
261 state-of-the-art (SOTA) LLMs that have demon-
262 strated strong performance across various NLP
263 tasks. These include proprietary models like GPT-
264 4o (OpenAI, 2023), GPT-3.5-turbo (Ouyang et al.,
265 2022a), Gemini-2.0-flash-exp (Anil et al., 2023),
266 and Claude-3.5-sonnet⁵, as well as open-source
267 models such as DeepSeek-v3 (DeepSeek-AI et al.,
268 2024), Llama3.1-8b-instruct (Dubey et al., 2024),
269 and Qwen2.5-7b-instruct (Yang et al., 2024). These
270 models represent the current frontier of language
271 model capabilities and serve as strong baselines for
272 evaluating our approach.

273 **Datasets** Our main experiments utilize the PF-
274 Dial dataset containing 12,705 training samples.
275 For evaluation, we maintain two test sets: our
276 out-of-domain test set comprising 698 decision
277 branches and 2,963 sequential branches, and an in-
278 domain test set containing 645 decision branches
279 and 1628 sequential branches. For backward tran-
280 sition experiments, we use the specially constructed
281 PFDial-H dataset. For format ablation studies, we
282 construct corresponding training and test sets in
283 three different state representation formats to sys-
284 tematically evaluate their impact on model perfor-
285 mance.

⁵<https://www.anthropic.com/news/claude-3-family>

Model	ID-test			OOD-test		
	Acc	Decision Acc	Sequential Acc	Acc	Decision Acc	Sequential Acc
Baselines						
LLaMA-3.1-8B-Instruct	—	—	—	13.20	2.16	15.00
Claude-3.5-Sonnet	—	—	—	62.74	22.06	69.40
Qwen2.5-7B-Instruct	—	—	—	65.87	37.88	71.34
Gemini-2.0-Flash-Exp	—	—	—	75.17	47.48	79.66
DeepSeek-v3	—	—	—	79.02	47.72	84.11
GPT-3.5-Turbo	—	—	—	79.76	39.57	86.29
GPT-4o	—	—	—	86.29	51.80	91.90
FineTuned on PFDial						
LLaMA-3.2-1B	98.90	98.59	98.96	93.57	87.05	94.62
LLaMA-3.2-3B	98.77	98.02	98.91	95.81	91.37	96.53
LLaMA-3.1-8B	99.03	98.31	99.17	97.29	96.88	97.35
Qwen2.5-0.5B	98.99	98.02	99.17	91.35	89.45	91.66
Qwen2.5-1.5B	98.90	97.46	99.17	94.00	88.97	94.82
Qwen2.5-3B	98.77	98.31	98.85	94.97	89.69	95.84
Qwen2.5-7B	98.94	98.02	99.11	96.51	90.65	97.47

Table 2: Results on PFDial: Decision Acc represents the accuracy of the decision branch, and Sequential Acc reflects the accuracy of the sequential branch. Acc is The overall accuracy.

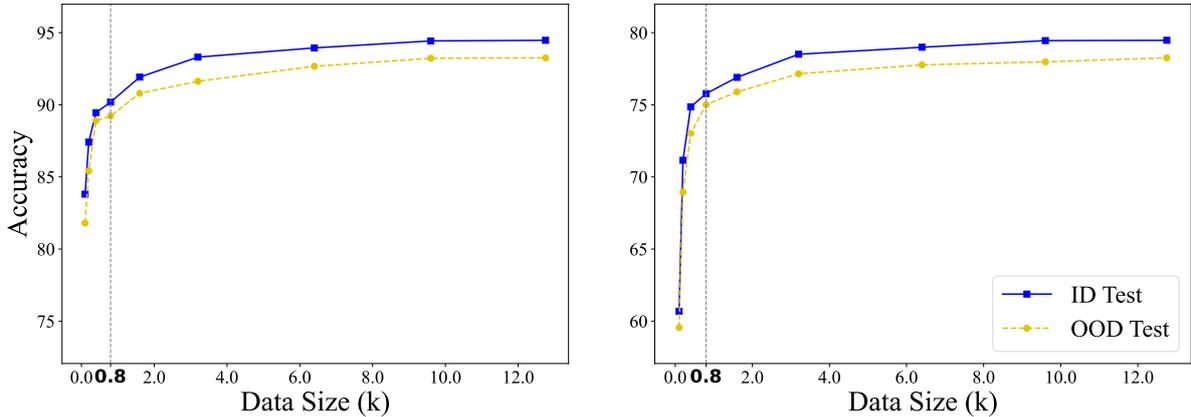


Figure 3: Results on Data Scaling: The left plot shows the accuracy of ID and OOD tests using data scaling strategy (a), while the right plot shows the accuracy using data scaling strategy (b).

Implementation Details We combine the PFDial dataset with the general dialogue dataset BELLE (BELLEGroup, 2023; Wen et al., 2023) in a 1:1 ratio for SFT. For the training process, we utilize the OpenRLHF(Hu et al., 2024) framework. For instance, the training of Qwen2.5-7B model is conducted on 8 H20 GPUs, with a total training time of approximately 1 hour. For detailed hyperparameter configurations, please refer to Table 6 in Appendix C.1.

Evaluation Metrics Model performance is evaluated by exact match accuracy between predicted and ground truth next states. For a prediction to be considered correct, it must exactly match the ground truth state transition. Specifically, our dataset contains two types of samples: sequential

samples and decision samples. Sequential samples refer to cases where there is only one unique next state given the current state. Decision samples refer to cases where there are at least two possible next states for the current state. Accordingly, our accuracy can be further refined into two types: sequential accuracy and decision accuracy.

3.2 Main Results

Table 2 presents our detailed experimental results on the main experiments. Our experiments demonstrate significant improvements over baseline models across all parameter scales. Even our smallest 0.5B parameter model achieves 91.35% accuracy on out-of-domain tests, with particularly strong performance on decision branches. Our 8B

Model	PFDial-H			OOD-test of PFDial		
	Acc	Backward Acc(Dist <5)	Backward Acc(Dist ≥5)	Acc	Decision Acc	Sequential Acc
Baselines						
LLaMA-3.1-8B-Instruct	15.00	13.64	15.52	13.20	2.16	15.00
Claude-3.5-Sonnet	57.50	59.09	56.90	62.74	22.06	69.40
Qwen2.5-7B-Instruct	31.25	31.82	31.03	65.87	37.88	71.34
Gemini-2.0-Flash-Exp	26.25	22.73	27.59	75.17	47.48	79.66
DeepSeek-v3	40.00	31.82	43.10	79.02	47.72	84.11
GPT-3.5-Turbo	51.25	54.55	50.00	79.76	39.57	86.29
GPT-4o	58.75	68.18	55.17	86.29	51.80	91.90
Secondary Fine-tuning						
Qwen2.5-0.5B	58.75	77.27	51.72	50.77	52.67	39.09
Qwen2.5-1.5B	47.50	45.45	48.28	87.90	88.82	82.25
Qwen2.5-3B	57.50	59.09	56.90	79.73	80.91	72.42
Qwen2.5-7B	66.25	90.91	56.90	76.04	76.82	71.22
Integrated Training						
Qwen2.5-0.5B	65.00	72.73	62.07	92.76	93.22	89.93
Qwen2.5-1.5B	70.00	72.73	68.97	93.87	94.82	88.01
Qwen2.5-3B	75.00	90.91	68.97	94.57	96.18	84.65
Qwen2.5-7B	76.25	86.36	72.41	96.31	96.96	92.33

Table 3: Results on PFDial and PFDial-H: Backward Acc represents the accuracy of backward transition.

model achieves sota performance with 97.29% accuracy on out-of-domain tests, surpassing GPT-4o by 11%. In decision branches, our 8B parameter model achieves 96.88% accuracy, surpassing GPT-4o by 43.88%.

3.3 Data Scaling Experiments

Experimental Setup To investigate data efficiency, we conducted experiments with varying training data sizes from 100 to 12,705 samples. We employed two data scaling strategies: **(a)** keeping fixed 12,000 general dialogue data samples while gradually increasing PFDial data, and **(b)** maintaining a 1:1 ratio mixing of general dialogue data and PFDial data.

Results and Analysis Figure 3 presents our experimental results on the data scaling experiments. The results demonstrate remarkable performance even with limited data: using only 800 samples, a 7B model can surpass 90% accuracy on OOD test. Performance continues to improve consistently with increased data volume, though we observe diminishing returns after approximately 3,000 samples. The comparable performance across both data scaling strategies validates the effectiveness of our PFDial dataset, demonstrating its robust data efficiency regardless of the mixing approach. For comprehensive results on decision accuracy and sequential accuracy across various data scaling configurations, please refer to Appendix C.2.

3.4 Backward Transition Studies

Experimental Setup We evaluated models’ capability in handling complex backward transitions on the PFDial-H test set, which provides a more rigorous assessment of models’ ability to strictly follow process constraints. The details of PFDial-H test set is shown in table 5 in Appendix B.2. For handling backward transitions, we compared two approaches: integrated training, which incorporates 440 PFDial-H training data samples to our PFDial training data during initial training, and secondary fine-tuning, which applies a secondary fine-tuning phase using PFDial-H data on the SFT-completed model with 440 PFDial-H training data samples.

Results and Analysis Detailed experimental results are shown in Table 3. Our results demonstrate that our method achieves optimal performance in handling backward branches, with the integrated training approach yielding superior results by maintaining robust performance on forward transitions while significantly improving accuracy in backward transition cases. Specifically, with integrated training, the Qwen2.5-7B model achieves 76.25% overall accuracy on PFDial-H, with 86.36% accuracy on backward transitions with distance less than 5, and 72.41% accuracy on transitions with distance greater than or equal to 5. Meanwhile, the model maintains a high accuracy of 96.31% on the OOD test.

In contrast, secondary fine-tuning not only fails

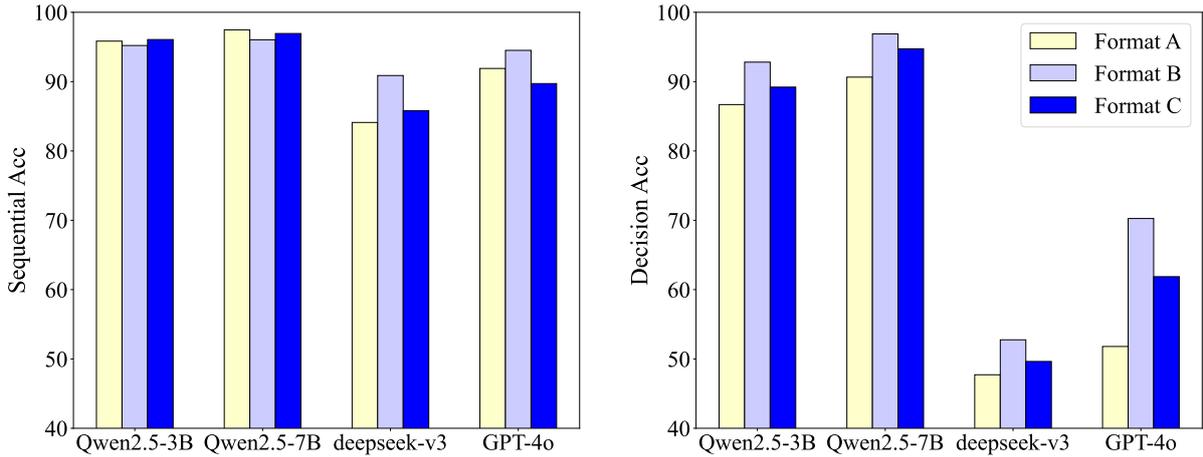


Figure 4: Results on Model Performance with Different Formats: The left plot shows the sequential accuracy for different models across three different formats (Format A, Format B, and Format C). The right plot shows the decision accuracy for the same models under the same formats.

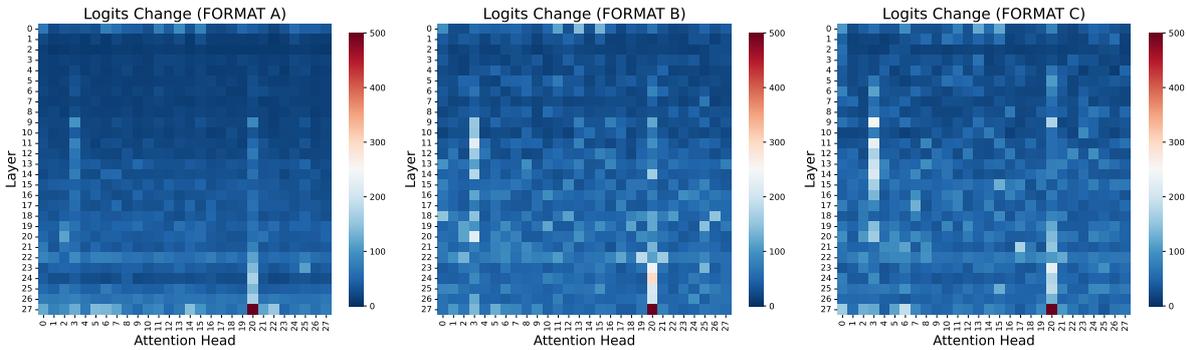


Figure 5: Logits changes of Three formats

376 to improve performance on backward transition cases but also reduces performance on the PFDial
 377 cases but also reduces performance on the PFDial dataset, with Qwen2.5-7B’s OOD test accuracy
 378 dropping from 96.31% to 76.04%. These results emphasize the importance of integrating backward
 379 transition samples during the initial training process rather than treating them as a post-hoc fine-
 380 tuning step.
 381
 382
 383

384 3.5 Format Ablation Studies

385 **Experimental Setup** We conducted experiments for three state representation formats: Format A:
 386 natural language description (default method), Format B: state codes (e.g., S1, C2), and Format C:
 387 combining codes and descriptions to explore the impact of different data formats on model performance
 388 and the underlying reasons. Specific cases of data in different formats and the visualization
 389 results can be seen in Appendix C.3. To ensure fairness, we fine-tuned the base model on data in
 390
 391
 392
 393
 394

395 all three formats and tested it with corresponding
 396 test sets containing the same content.

397 We then froze the attention heads of each layer and compared the model’s output logits with the
 398 original logits in these scenarios. By examining the logits changes, we assessed the impact of each
 399 attention head after fine-tuning with different formats. Finally, we visualized these attention heads
 400 to gain deeper insights into their roles.
 401
 402
 403

404 **Results and Analysis** The model performance with different formats are shown in Figure 4. The
 405 results indicate that Format B achieved the highest accuracy in most cases, particularly on Decision
 406 Accuracy, however, it performed slightly worse than other formats on Sequential Accuracy. The fol-
 407 lowing experimental results, to some extent, shed light on the reasons behind this phenomenon.
 408
 409
 410
 411

412 The comparison results of the three formats’s attention heads are shown in Figure 5. The models
 413 fine-tuned with FormatB and FormatC showed a
 414

more uniform and diverse distribution of attention head contributions. This can be explained that both of the latter formats introduced code state identifiers, requiring the language model to learn both the sequence of state transitions and the correspondence between code states and natural language states. Thus more attention heads with different functions were activated.

In particular, *head_layer27_head20* is crucial in all three formats. We visualized this attention head’s scores for all three formats in Figure 6, Figure 7, and Figure 8, respectively, focusing on significant differences. The attention scores for both the natural language and hybrid formats were concentrated at the intersection of the user’s current state and the corresponding state in the flowchart. Format C, which includes some state codes, showed a more moderated concentration. In contrast, Format B did not exhibit such clustering. We hypothesize that the introduction of state codes allows the model’s attention to generalize across different input parts, rather than being confined to specific segments. This enables models to better understand user inputs and facilitates learning of global logic, such as condition checking and state selection. This also reasonably explains previous results.

4 Related Work

4.1 Controllable Reasoning in LLMs

Controllable reasoning in LLMs has gained significant attention in recent years. Ouyang et al. (2022b) pioneered instruction-guided control via Reinforcement Learning from Human Feedback (RLHF), combining supervised fine-tuning (SFT), reward model training, and Proximal Policy Optimization (PPO) to align outputs with human preferences. While effective, this approach requires complex annotation and training (Li and Liang, 2021). In contrast, our approach simplifies the process by encoding reasoning logic into structured UML state flowcharts, guiding learning through SFT alone. This provides a clear, human-readable control mechanism, addressing the ‘reasoning opacity’ challenge (Liang et al., 2024b).

4.2 Graph-based Enhanced Reasoning

Previous research (Pan et al., 2024) has explored graph-based approaches to enhance the reasoning capabilities of LLMs. Several works (Liang et al., 2024a; Luo et al., 2024a,b) have used Knowledge Graph (KG) structural information to reduce hal-

lucinations by breaking reasoning into path extraction and inference steps. Similarly, Zhou et al. (2024) showed that graph-based training improves multi-hop reasoning accuracy. However, these methods mainly treat graphs as external knowledge sources rather than explicit control mechanisms.

4.3 LLM-based Dialogue Systems

Task-Oriented Dialogue (TOD) systems help users achieve specific goals through conversations (Yi et al., 2024). Current approaches fall into two main categories: Pipeline-based Approaches, which separate dialogue systems into modules with LLMs handling specific tasks (Comi et al., 2023; Parikh et al., 2023; Chen et al., 2019; Nguyen et al., 2023), and End-to-End Approaches, which use LLMs to generate responses based on the entire dialogue history (Hemanthage et al., 2023; Zhang et al., 2024, 2023; Wu et al., 2020; Algherairy and Ahmed, 2025). Pipeline-based Approaches offer better transparency but require extensive annotated data, while End-to-End Approaches are simpler but less controllable and demand higher model capabilities.

5 Conclusion

In this paper, we introduce the PFDial dataset, a novel resource designed to enhance process-driven dialogue systems. By utilizing structured dialogue instructions derived from UML flowcharts, PFDial provides a robust framework for training models to handle complex decision-making and sequential processes. Our experiments demonstrate that models fine-tuned on PFDial achieve high accuracy, even with limited training data, and outperform sota LLMs like GPT-4o on specific tasks.

We conducted an in-depth analysis of backward transitions using the PFDial-H dataset, highlighting the importance of integrated training approaches for maintaining strong performance across diverse dialogue scenarios. Additionally, we explored the impact of various data representation formats, finding that structured state codes significantly improve the accuracy of state transition predictions.

Overall, our work underscores the potential of structured datasets like PFDial to advance process-driven dialogue systems, offering new insights into the design and training of models for precise and controlled reasoning. Future research will focus on expanding the dataset to cover more scenarios and refining training methodologies to enhance model generalization and adaptability.

513 Limitations

514 Our research presents a comprehensive set of
515 experiments, yet it is not without limitations.
516 First, the Chinese-centric nature of our dataset
517 introduces potential cross-lingual generalization
518 constraints. Second, the scarcity of standardized
519 flowchart benchmarks in Chinese process speci-
520 fications increases the risk of domain-specific bi-
521 ases, despite our rigorous validation framework;
522 while the potential residual inconsistencies in flow-
523 to-text conversion may emerge from the inherent
524 subjectivity in interpreting semantic structures.

525 References

526 Atheer Algherairy and Moataz Ahmed. 2025. [Prompt-](#)
527 [ing large language models for user simulation in task-](#)
528 [oriented dialogue systems](#). *Comput. Speech Lang.*,
529 89:101697.

530 Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-
531 Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan
532 Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Mil-
533 lican, David Silver, Slav Petrov, Melvin Johnson,
534 Ioannis Antonoglou, Julian Schrittwieser, Amelia
535 Glaese, Jilin Chen, Emily Pitler, Timothy P. Lilli-
536 crap, Angeliki Lazaridou, Orhan Firat, James Molloy,
537 Michael Isard, Paul Ronald Barham, Tom Hennig-
538 an, Benjamin Lee, Fabio Viola, Malcolm Reynolds,
539 Yuanzhong Xu, Ryan Doherty, Eli Collins, Clemens
540 Meyer, Eliza Rutherford, Erica Moreira, Kareem
541 Ayoub, Megha Goel, George Tucker, Enrique Pi-
542 queras, Maxim Krikun, Iain Barr, Nikolay Savinov,
543 Ivo Danihelka, Becca Roelofs, Anaïs White, Anders
544 Andreassen, Tamara von Glehn, Lakshman Yagati,
545 Mehran Kazemi, Lucas Gonzalez, Misha Khalman,
546 Jakub Sygnowski, and et al. 2023. [Gemini: A fam-](#)
547 [ily of highly capable multimodal models](#). *CoRR*,
548 abs/2312.11805.

549 BELLEGroup. 2023. Belle: Be everyone’s large
550 language model engine. [https://github.com/](https://github.com/LianjiaTech/BELLE)
551 [LianjiaTech/BELLE](https://github.com/LianjiaTech/BELLE).

552 Qian Chen, Zhu Zhuo, and Wen Wang. 2019. [BERT](#)
553 [for joint intent classification and slot filling](#). *CoRR*,
554 abs/1902.10909.

555 Daniele Comi, Dimitrios Christofidellis, Pier Francesco
556 Piazza, and Matteo Manica. 2023. [Zero-shot-bert-](#)
557 [adapters: a zero-shot pipeline for unknown intent](#)
558 [detection](#). In *Findings of the Association for Comput-*
559 *ational Linguistics: EMNLP 2023, Singapore, De-*
560 *cember 6-10, 2023*, pages 650–663. Association for
561 Computational Linguistics.

562 DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingx-
563 uan Wang, Bochao Wu, Chengda Lu, Chenggang
564 Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan,
565 Damai Dai, Daya Guo, Dejian Yang, Deli Chen,
566 Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai,

Fuli Luo, Guangbo Hao, Guanting Chen, Guowei
567 Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng
568 Wang, Haowei Zhang, Honghui Ding, Huajian Xin,
569 Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang,
570 Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang,
571 Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie
572 Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu,
573 Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean
574 Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao,
575 Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang,
576 Mingchuan Zhang, Minghua Zhang, Minghui Tang,
577 Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang,
578 Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu
579 Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge,
580 Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin
581 Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao
582 Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,
583 Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu
584 Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou,
585 Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun,
586 W. L. Xiao, and Wangding Zeng. 2024. [Deepseek-v3](#)
587 [technical report](#). *CoRR*, abs/2412.19437.

588 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
589 Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
590 Akhil Mathur, Alan Schelten, Amy Yang, Angela
591 Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,
592 Archi Mitra, Archie Sravankumar, Artem Korenev,
593 Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien
594 Rodriguez, Austen Gregerson, Ava Spataru, Bap-
595 tiste Rozière, Bethany Biron, Binh Tang, Bobbie
596 Chern, Charlotte Caucheteux, Chaya Nayak, Chloe
597 Bi, Chris Marra, Chris McConnell, Christian Keller,
598 Christophe Touret, Chunyang Wu, Corinne Wong,
599 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-
600 lonsius, Daniel Song, Danielle Pintz, Danny Livshits,
601 David Esiobu, Dhruv Choudhary, Dhruv Mahajan,
602 Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,
603 Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,
604 Emily Dinan, Eric Michael Smith, Filip Radenovic,
605 Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Geor-
606 gia Lewis Anderson, Graeme Nail, Grégoire Mialon,
607 Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-
608 nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,
609 Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan
610 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan
611 Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,
612 Jeet Shah, Jelmer van der Linde, Jennifer Billock,
613 Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,
614 Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,
615 Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph
616 Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,
617 Kalyan Vasuden Alwala, Kartikeya Upasani, Kate
618 Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and
619 et al. 2024. [The llama 3 herd of models](#). *CoRR*,
620 abs/2407.21783.

621 Bhatthiya Hemanthage, Christian Dondrup, Phil Bartie,
622 and Oliver Lemon. 2023. [Simplemtd: A simple](#)
623 [language model for multimodal task-oriented dia-](#)
624 [logue with symbolic scene representation](#). *CoRR*,
625 abs/2307.04907.

626 Jian Hu, Xibin Wu, Weixun Wang, Xianyu, Dehao
627 Zhang, and Yu Cao. 2024. [Openrlhf: An easy-to-](#)
628

and Xuanjing Huang. 2024. *Transfertod: A generalizable chinese multi-domain task-oriented dialogue system with transfer capabilities*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 12750–12771. Association for Computational Linguistics.

Xiaoying Zhang, Baolin Peng, Kun Li, Jingyan Zhou, and Helen Meng. 2023. *SGP-TOD: building task bots effortlessly via schema-guided LLM prompting*. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, pages 13348–13369. Association for Computational Linguistics.

Jiaming Zhou, Abbas Ghaddar, Ge Zhang, Liheng Ma, Yaochen Hu, Soumyasundar Pal, Mark Coates, Bin Wang, Yingxue Zhang, and Jianye Hao. 2024. *Enhancing logical reasoning in large language models through graph-based synthetic data*. *CoRR*, abs/2409.12437.

A Algorithm

A.1 Parse PlantUML Code

This algorithm parses the given PlantUML code by line by line. It initializes an empty dictionary, *nodeDict*, to store state nodes. Then, it calls *PARSESEQUENTIAL* (Algorithm 2) to process the sequential flow. Finally, it returns all paths originating from the start node, representing the complete execution flow.

Algorithm 1 Parse PlantUML Code

```

1: function PARSEPLANTUML(code)
2:   Split code into lines
3:   Initialize an empty dictionary nodeDict to
   store nodes
4:   Call PARSESEQUENTIAL
5:   return all paths originating from the start
   node
6: end function

```

A.2 Parse Sequential Blocks

This algorithm parses the sequential execution flow by processing each line to identify states. Sequential states create and merge new nodes, while decision states call *PARSEDECISION* (Algorithm 3) for further processing.

A.3 Parse Decision Blocks

This algorithm handles decision structures by first parsing the "if" block using *PARSESEQUENTIAL* (Algorithm 2). If an "else if" block is encountered,

Algorithm 2 Parsing Sequential Blocks

```

1: function PARSESEQUENTIAL(startNode, lines,
   nodeDict)
2:   root ← startNode
3:   for each line in lines do
4:     Trim whitespace
5:     if the line represents a sequential state
   then
6:       Create a new node into nodeDict
7:       Merge new nodes
8:     else if the line represents a decision
   state then
9:       Call PARSEDECISION
10:      Connect new nodes
11:    else
12:      Continue processing
13:    end if
14:  end for
15:  return
16: end function

```

it recursively calls itself to process nested conditions. For an "else" block, it parses the sequence and connects the resulting nodes using *PARSESEQUENTIAL* (Algorithm 2). This ensures correct branching and flow control within decision blocks.

Algorithm 3 Parsing Decision Blocks

```

function PARSEDECISION(startNode, lines,
nodeDict)
  root ← startNode
  Call PARSESEQUENTIAL to parse "if" block
  if meet "else if" block then
    Recursively call PARSEDECISION
    Connect new nodes
    return
  else if meet "else" block then
    Parse the "else" block using PARSESE-
    QUENTIAL
    Connect new nodes
    return
  else
    return
  end if
end function

```

783
784
785
786
787

788

Format	PlantUML	ChartMage	NomNoml	Meamaid
Decision Block	if (D) then (C1) S1 else (C2) S2	D - C1 ->> S1 D - C2 ->> S2	[D] C1-> [Block1] [D] C2-> [Block2]	A{D} - C1 -> B[S1] A - C2 -> C[S2]
Input	Here is the flowchart code: [a UML flowchart with 21 paths, expressed in the corresponding format] Show all possible complete paths and count how many there are.			
Path Count	21	15	19	17

Table 4: Comparison of Flowchart Formats

B Dataset

B.1 Different Format to Represent Flowchart

In Table 4, we compares different formats to represent flowchart. We select a complex flowchart with 21 paths, and let GPT-4o to find all paths and count the number paths. Only with PlantUML, GPT-4o correctly output the right path number, which means flowchart with PlantUML is easy for models to understand. Moreover, syntax features, as is shown in Table 4, such as indentation branching, and loops, making PlantUML a good format to read and convenient for program processing.

B.2 PFDial-H Tests

Table 5 presents the details of PFDial-H tests, including the total number of dialogues in the test set, average dialogue length in turns, and the proportion of dialogues with backward transition distances divided by the threshold of 5.

C Experimental Details

C.1 Implementation Details

To ensure reproducibility of our experiments, we provide detailed hyper-parameter configurations used in our training process. All experiments were conducted using the AdamW optimizer with cosine annealing learning rate scheduling. The complete set of training hyper-parameters is presented in Table 6. We maintained consistent hyper-parameter settings across models of different scales to ensure fair comparison.

C.2 Supplementary Data for Data Scaling Experiment

While the section 3.3 presents the overall accuracy trend with increasing training samples, here we provide the complete experimental results. Tables 11 and 12 present the performance results for Strategy A and Strategy B respectively. Each table shows

Backward Distance < 5	22
Backward Distance ≥ 5	58
Dialogue Samples	80
Avg. Length	534.36

Table 5: PFDial-H Tests Data Overview

Hyperparameter	Value
Optimizer	AdamW
Learning Rate	5×10^{-6}
Learning Rate Scheduling	Cosine Annealing
Adam Beta1	0.9
Adam Beta2	0.95
Batch Size	128
Batch Size Per-Device	4
Training Epochs	5

Table 6: Training Hyper Parameters

overall accuracy, decision accuracy, and sequential accuracy metrics on both ID and OOD tests across different training sample sizes. The information of dataset with different training sample sizes is shown in Table 10.

C.3 Details in Format Ablation Study

In this section, we present specific cases of data in different formats and the visualization results for *head_layer27_head20*. As shown in Figure 6, 7, and 8, we visualize the attention patterns for Format A (Natural Language) in Table 13, Format B (State Code) in Table 14, and Format C (Hybrid) in Table 15 respectively.

D Prompt

D.1 Prompt for generating User Input

We use the prompt from Table 7 to generate user inputs. This prompt helps us create appropriate user input text based on the given state transitions.

User

These examples are four-tuples consisting of the PlantUML diagram, the current state, the next state, and the user input.

[several examples]

The user's input explains the change in state from the current state to the next state. For example, if the original state is A, the user might input "A has been completed." or, when a choice is required, the user selects an option based on the next state.

Now I have a four-tuple consisting of the PlantUML diagram, the current state, and the next state, without user input. Your task is to generate the user's input based on the rules provided.

This is the four-tuple need to be filled:

[four-tuple to be filled]

Table 7: The prompt for generating the user's input.

D.2 Prompt for generating Robot Output

We use the prompt from Table 8 to generate robot outputs. This prompt helps us create appropriate robot responses based on the current state, next state, and user input.

User

These examples are five-tuples consisting of the PlantUML diagram, the current state, the next state, user input, and the robot output.

[several examples]

The user's input explains the change in state from the current state to the next state. The robot output is related to next state. Robot acts as the server-provider. For example, if the current state is A, robot might output "Now process A." or, when a choice is required, robot lets user to make a choice. Now I have a five-tuple consisting of the PlantUML diagram, the current state next state, and the user input, without robot output. Your task is to generate the robot's output based on the rules provided.

This is the five-tuple need to be filled:

[five-tuple to be filled]

Table 8: The prompt for generating the robot's output.

D.3 Prompt for Adding Backward Transition

In the third step, we use the prompt from Table 9 to add backward transitions. This prompt guides us in adding logical loop structures to the flowchart.

User

This is a flowchart in PlantUML syntax and the result after adding a loop to itself:

[original PlantUML and revised PlantUML]

Your task is to follow this modification rule to add a loop to the plantuml that I will give you next. The following conditions must be met:

1. The added loop is logical
2. The conditional state of "repeat while" cannot be repeated with the any conditional state that already exists in the original PlantUML
3. Ensure that the syntax of PlantUML is correct
4. Add is([need to loop]) not([jump out of loop]) statements after repeat while as much as possible.

PlantUML to be modified:

[original PlantUML]

Table 9: Prompt for Adding Backward Transition.

843

844

845

846

847

848

849

850

851

Training Sample Size	100	200	400	800	1600	3200	6400	9600	12705
Flowcharts	3	5	12	29	59	113	208	308	440
State Nodes	49	79	149	304	612	1236	2473	3772	5055
Sequential Samples	62	134	264	561	1099	2185	4293	6717	9029
Decision Samples	38	66	136	239	501	1015	2107	2883	3676
Avg. Length	386.22	368.19	320.24	286.92	275.63	272.29	273.00	280.90	277.16

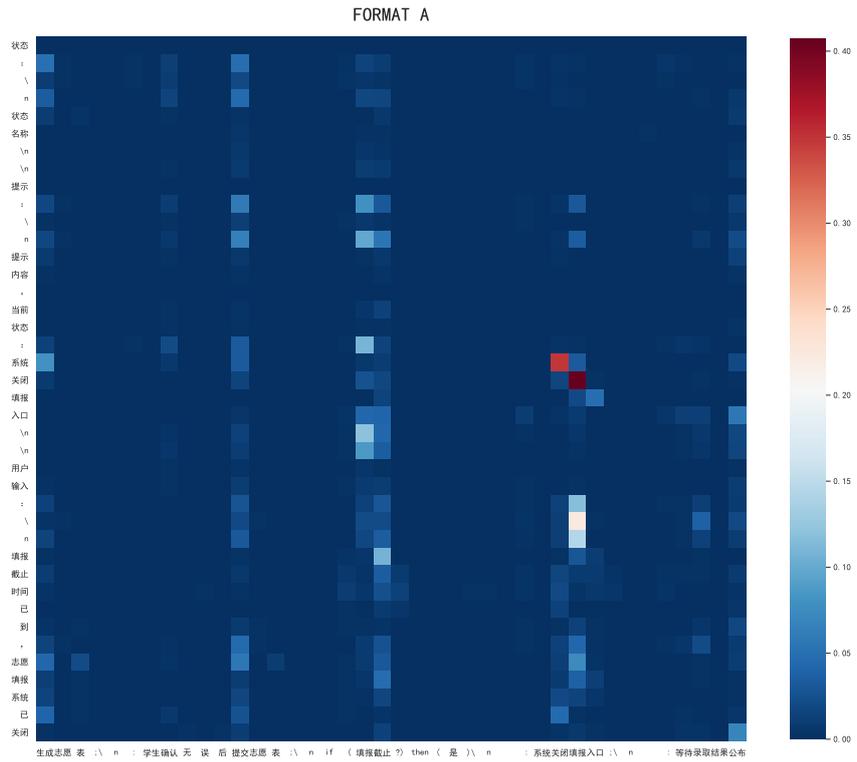
Table 10: Statistics of the PFDial Dataset with Different Training Sample Sizes

Training Sample Size	ID-test			OOD-test		
	Acc	Decision Acc	Sequential Acc	Acc	Decision Acc	Sequential Acc
100-all	77.61	54.52	81.87	73.59	39.81	79.08
200-all	84.82	83.33	85.10	80.83	62.35	83.83
400-all	88.91	90.68	88.59	87.80	76.50	89.64
800-all	90.37	92.37	89.99	88.44	78.90	89.99
1600-all	93.84	94.92	93.64	91.59	84.17	92.79
3200-all	96.61	97.46	96.46	93.26	89.69	93.84
6400-all	97.89	97.18	98.02	95.34	91.85	95.91
9600-all	98.86	98.31	98.96	96.45	90.17	97.47
all	98.94	98.02	99.11	96.51	90.65	97.47

Table 11: Performance with different training sample sizes across ID and OOD datasets after training on Qwen2.5-7B, with data scaling strategy (a)

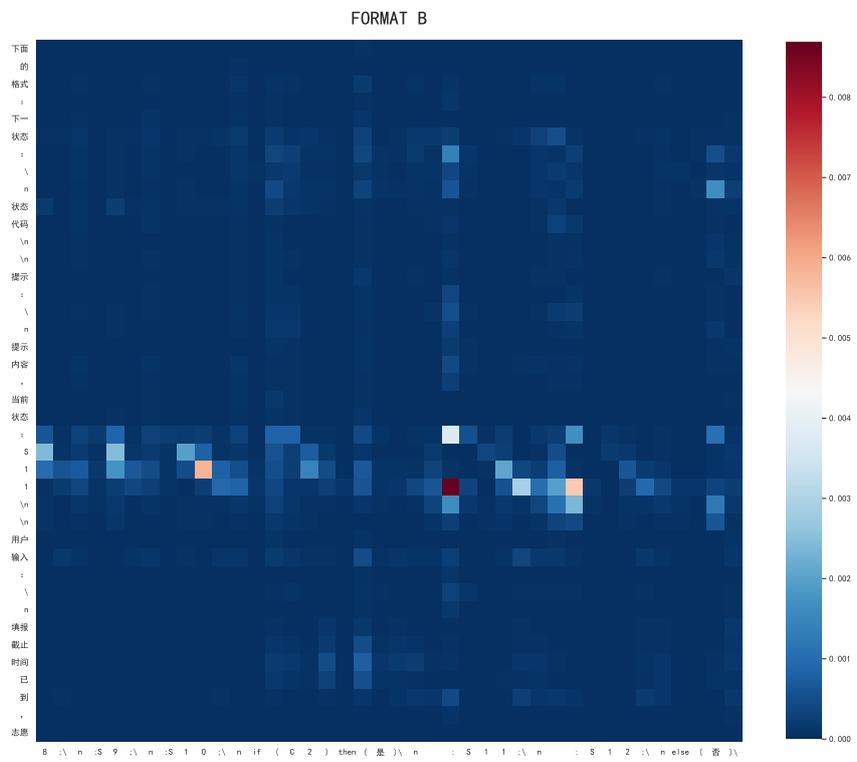
Training Sample Size	ID-test			OOD-test		
	Acc	Decision Acc	Sequential Acc	Acc	Decision Acc	Sequential Acc
100-100	61.37	23.45	68.37	59.12	13.91	66.46
200-200	82.31	84.46	81.92	77.88	65.47	79.90
400-400	89.71	89.55	89.73	86.03	75.30	87.77
800-800	91.55	92.66	91.35	90.01	81.06	91.47
1600-1600	93.80	93.50	93.85	91.79	85.13	92.87
3200-3200	97.01	98.02	96.82	94.30	90.65	94.90
6400-6400	97.98	98.31	97.92	95.54	90.65	96.34
9600-9600	98.90	97.46	99.17	95.95	89.45	97.00
all	98.94	98.02	99.11	96.51	90.65	97.47

Table 12: Performance with different training sample sizes across ID and OOD datasets after training on Qwen2.5-7B, with data scaling strategy (b)



(a) Format A

Figure 6: local attention score of $head_layer_{27_head_{20}}$ using Format A.



(b) Format B

Figure 7: local attention score of $head_layer_{27_head_{20}}$ using Format B.

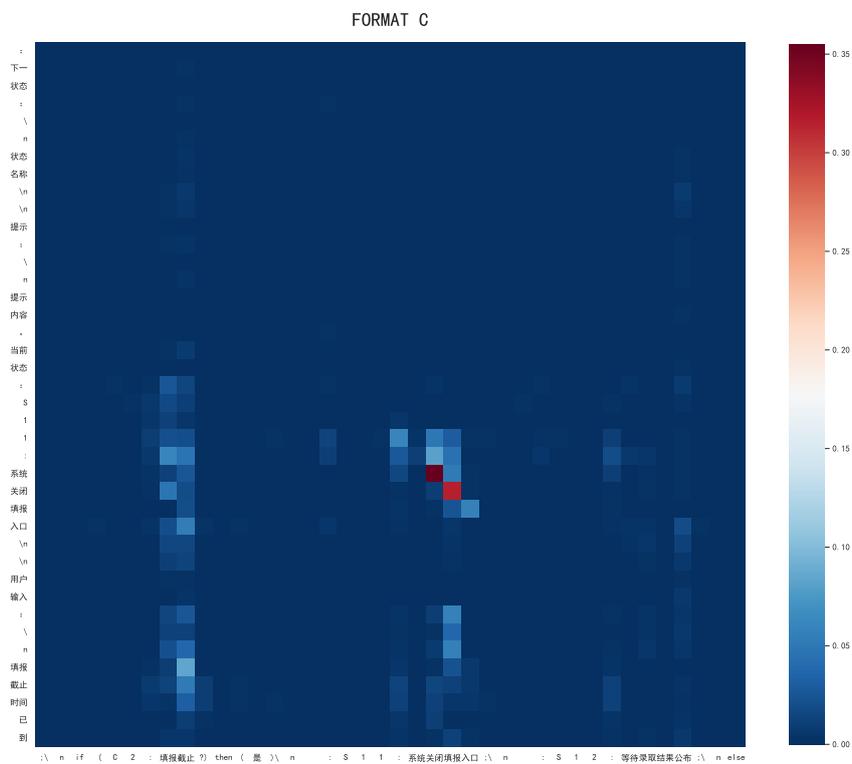


Figure 8: local attention score of $head_{layer27_head20}$ using Format C.

Format	Data Example
Format A	<p>PlantUML:</p> <pre>@startuml start :College entrance exam results announced; :Student obtains college entrance exam score report; :Student checks the exam results and determines the range of colleges and majors to apply for; if (Do you need to research colleges and majors in advance?) then (Yes) :Student conducts research on colleges and majors; else (No) :Skip this step; endif :Student logs into the application system; :System provides the application entry and instructions; :Student fills out the application, prioritizing choices; :After completing the application, the system generates an application form; :Student confirms the form and submits it; if (Application deadline?) then (Yes) :System closes the application entry; :Wait for the admission results to be announced; else (No) :Student can modify the application before the deadline; :Wait for the application deadline; endif :Admission results announced; if (Admitted?) then (Yes) :Student completes registration procedures according to the admission notice; else (No) :Student applies for supplementary applications or participates in the supplementary application process; endif :Enroll in the school; stop @enduml</pre> <p>Current state: System closes the application entry</p> <p>Next state: Wait for the admission results to be announced</p> <p>User input: The application deadline has passed, and the application system is now closed.</p> <p>Robot output: Please patiently wait for the announcement of the admission results.</p>

Table 13: Example Data Display - Format A

Format	Data Example
Format B	<p>PlantUML: @startuml start :S1; :S2; :S3; if (C1) then (Yes) :S4; else (No) :S5; endif :S6; :S7; :S8; :S9; :S10; if (C2) then (Yes) :S11; :S12; else (No) :S13; :S14; endif :S15; if (C3) then (Yes) :S16; else (No) :S17; endif :S18; stop @enduml</p> <p>Dictionary of the state codes: { "College entrance exam results announced": "S1", "Student obtains college entrance exam score report": "S2", "Student checks the exam results and determines the range of colleges and majors to apply for": "S3", "Do you need to research colleges and majors in advance?": "C1", "Student conducts research on colleges and majors": "S4", "Skip this step": "S5", "Student logs into the application system": "S6", "System provides the application entry and instructions": "S7", "Student fills out the application, prioritizing choices": "S8", "After completing the application, the system generates an application form": "S9", "Student confirms the form and submits it": "S10", "Application deadline?": "C2", "System closes the application entry": "S11", "Wait for the admission results to be announced": "S12", "Student can modify the application before the deadline": "S13", "Wait for the application deadline": "S14", "Admission results announced": "S15", "Admitted?": "C3", "Student completes registration procedures according to the admission notice": "S16", "Student applies for supplementary applications or participates in the supplementary application process": "S17", "Enroll in the school": "S18" }</p> <p>Current state: System closes the application entry</p> <p>Next state: Wait for the admission results to be announced</p> <p>User input: The application deadline has passed, and the application system is now closed.</p> <p>Robot output: Please patiently wait for the announcement of the admission results.</p>

Table 14: Example Data Display - Format B

Format	Data Example
Format C	<p>PlantUML:</p> <pre>@startuml start :S1: College entrance exam results announced; :S2: Student obtains college entrance exam score report; :S3: Student checks the exam results and determines the range of colleges and majors to apply for; if (C1: Do you need to research colleges and majors in advance?) then (Yes) :S4: Student conducts research on colleges and majors; else (No) :S5: Skip this step; endif :S6: Student logs into the application system; :S7: System provides the application entry and instructions; :S8: Student fills out the application, prioritizing choices; :S9: After completing the application, the system generates an application form; :S10: Student confirms the form and submits it; if (C2: Application deadline?) then (Yes) :S11: System closes the application entry; :S12: Wait for the admission results to be announced; else (No) :S13: Student can modify the application before the deadline; :S14: Wait for the application deadline; endif :S15: Admission results announced; if (C3: Admitted?) then (Yes) :S16: Student completes registration procedures according to the admission notice; else (No) :S17: Student applies for supplementary applications or participates in the supplementary application process; endif :S18: Enroll in the school; stop @enduml</pre> <p>Current state: System closes the application entry</p> <p>Next state: Wait for the admission results to be announced</p> <p>User input: The application deadline has passed, and the application system is now closed.</p> <p>Robot output: Please patiently wait for the announcement of the admission results.</p>

Table 15: Example Data Display - Format C

Category	Senarios
Lifestyle Services	Hairdressing, Phone Card, Car Wash, Agritainment, Printing & Copying, Lawyer, Yoga Studio, Music Classroom, Internship
Daily Convenience	Takeout, Scenic Spots, Furniture Cleaning, Wedding Photography, Movie, Cleaning Service, Self-service Car Wash, Second-hand Car Trading, Visa Application
Food & Education	Catering, Training Courses, Physiotherapy & Massage, Lighting Design, Gym, Express Delivery, Travel Group Purchase, Health Check-up, Group Tour
Entertainment & Transportation	Concert, Car Rental, Baking Studio, Digital Repair, Airplane Ticket, Water Delivery, Florist Shop, Photo Studio, Course Selection
Business & Professional Services	Commercial Photography, Hospital, Pet Boarding, Café, Dentist, Pet Grooming, Tea House, Outdoor Expansion, Electrician Inspection
Luxury & Specialized Services	Beauty Salon, Museum, Horticultural Design, Car Maintenance, Cruise, Photography Studio Rental, Piano Tuning, Basketball Court, Cargo Delivery into Cabin
Living-related Services	Laundry, House Rental, Education Consultation, Library, Cultural Exhibition, Health Consultation, Holiday Villa, Interior Design, Bank Account Operation
Maintenance & Care Services	Home Appliance Repair, Hotel, Leather Goods Care, Arcade, Furniture Installation, Medical Check-up, Car Insurance Claim, Bicycle Rental, Parts Inspection
Comprehensive Services	Moving, Resort, Language Translation, Medical Aesthetics, Driving School, Wedding Planning, Pet Hospital, Manicure, Document Approval
Shopping & Other Activities	Shopping, Train Ticket, Water & Electricity Repair, Ski Resort, Credit Card, College Entrance Examination Volunteer Filling, Cooking, DIY Handicraft, Content Creation

Table 16: Senarios