

DyBBT: Dynamic Balance via Bandit inspired Targeting for Dialog Policy with Cognitive Dual Systems

Anonymous ACL submission

Abstract

Task oriented dialog systems often rely on static exploration strategies that do not adapt to dynamic dialog contexts, leading to inefficient exploration and suboptimal performance. We propose DyBBT, a novel dialog policy learning framework that formalizes the exploration challenge through a structured cognitive state space \mathcal{C} that captures dialog progression, user uncertainty, and slot dependency. DyBBT proposes a bandit inspired meta-controller that dynamically switches between a fast intuitive inference (System 1) and a slow deliberative reasoner (System 2) based on real-time cognitive states and visitation counts. Extensive experiments on single- and multi-domain benchmarks show that DyBBT achieves SOTA performance in success rate, efficiency, and generalization, with human evaluations confirming that its decisions are well aligned with expert judgment. The code is available at [Anonymous Github](#).

1 Introduction

“The affordances of the environment are what it offers the animal, what it provides or furnishes, for good or ill.”

— The Ecological Approach to Visual Perception ([Gibson, 1979](#))

Task-oriented dialog systems (TODS) assist users in achieving specific goals through multi-turn interactions. Dialog policy learning is typically formulated as a sequential decision making problem addressed with Deep Reinforcement Learning (DRL) ([Nachum et al., 2017](#); [Silver et al., 2014](#)). However, it is fundamentally bottlenecked by the exploration exploitation dilemma: balancing the exploitation of known rewards against the exploration of unknown actions to discover better strategies. In TODS, this dilemma is exacerbated by the dynamic and partially observable context, characterized by quantifiable cognitive features such as

dialog progress, user intent entropy, and slot dependencies ([Peng et al., 2017](#); [Wen et al., 2017](#)). The features govern the cost benefit of exploration: early in a dialog, high entropy makes information gathering valuable, and later high slot dependency makes exploitation critical to avoid constraint violations ([Qin et al., 2023](#); [Zhao et al., 2024](#)).

Existing methods for enhancing exploration in TODS remain misaligned with this dynamic cognitive reality. As illustrated in [Fig. 1](#), traditional DRL methods rely on static heuristics like ϵ -greedy ([Niu et al., 2024](#)), which cannot adapt to shifting exploration needs across dialog phases. Evolutionary methods like EIERL ([Zhao et al., 2025](#)) enable global search via population based optimization but struggle to scale in complex multi-domain scenarios. LLM-based policies ([Zhang et al., 2024](#); [He et al., 2022](#)) or reasoning techniques like Tree of Thoughts ([Yao et al., 2023](#)) support deep deliberative planning but incur prohibitive computational overhead and lack a principled mechanism to trigger such costly reasoning only when necessary. This misalignment reveals a key research question: *How can we design a dialog policy that dynamically perceives cognitive affordances to balance exploration and exploitation?*

Inspired by Gibson’s affordances theory, we propose that the dialog environment presents a dynamic landscape of exploration opportunities which an effective policy must perceive and act upon. We introduce DyBBT, a novel framework that grounds decisions in an interpretable cognitive state space \mathcal{C} that captures dialog progress, user uncertainty, and slot dependency. It employs a lightweight meta-controller that dynamically switches between a fast System 1 for routine decisions and a slow System 2 for costly deliberation, based on real-time cognitive signals. This design ensures expensive reasoning is invoked only when the cognitive state signals high epistemic uncertainty or low confidence, addressing the core

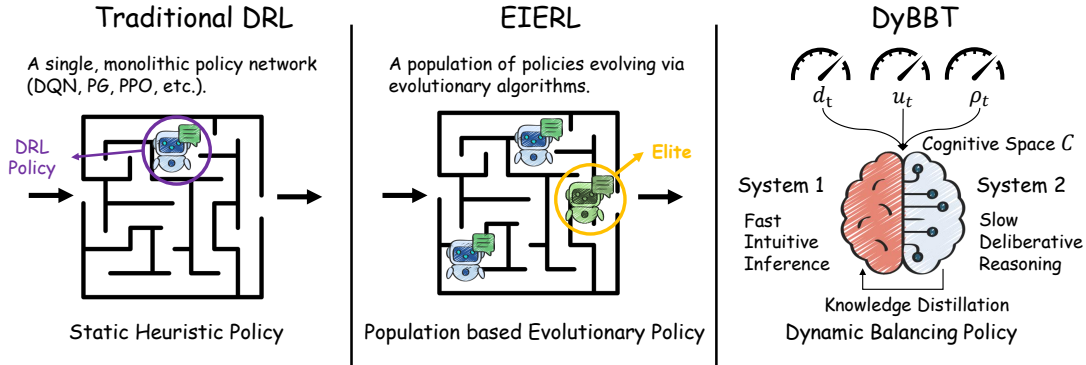


Figure 1: A comparison of exploration strategies for dialog policy learning. Traditional DRL relies on static heuristics incapable of adapting to dynamic dialog contexts. EIERL uses population based evolution but struggles in complex tasks. DyBBT solves the adaptive exploration challenge by cognitive meta-controller to achieve a principled balance between efficiency and robustness.

082 limitations of previous methods.

083 In summary, our contributions are: (1) Formalizing the TODS exploration challenge via a
 084 structured cognitive state space \mathcal{C} (Sec. 3.1); (2)
 085 Proposing DyBBT, a novel framework with a meta-
 086 controller to dynamically balance between fast and
 087 slow reasoning (Sec. 3.1.2); (3) Demonstrating
 088 SOTA performance and human aligned decisions
 089 through extensive experiments (Sec. 4).
 090

091 2 Related Work

092 2.1 DRL for Dialog Policy Learning

093 Deep Reinforcement Learning (DRL) has become
 094 a dominant paradigm for dialog policy optimization
 095 due to its capacity for sequential decision making.
 096 Early work applied value based methods (Peng
 097 et al., 2018) and policy gradient algorithms (Sil-
 098 ver et al., 2014) to TODS, with Proximal Policy
 099 Optimization (PPO) (Schulman et al., 2017) later
 100 adopted for improved stability. A key limitation
 101 of these methods is their reliance on static ex-
 102 ploration strategies, such as ϵ -greedy or entropy
 103 bonus, which cannot adapt to the dynamic uncer-
 104 tainty and structural complexity of multi-domain
 105 dialogs (Kwan et al., 2023; Jia et al., 2024). Recent
 106 efforts have incorporated Bayesian reasoning (Lee
 107 et al., 2023), meta-learning (Li et al., 2024; Liang
 108 et al., 2024), and cascading RL (Du et al., 2024)
 109 to enable more adaptive exploration. While promis-
 110 ing, these approaches often lack an explicit and
 111 interpretable representation of the internal dialog
 112 state that directly governs exploration, the gap that
 113 our structured cognitive state space \mathcal{C} aims to fill.

114 2.2 Evolutionary Exploration Methods

115 To overcome static exploration, research has di-
 116 verged into population based optimization and prin-
 117 cipled exploration theory. Population based meth-
 118 ods like Evolutionary RL (EIERL) (Zhao et al.,
 119 2025) enhance diversity but scale poorly with dia-
 120 log complexity (Sigaud, 2023) and lack real time
 121 adaptation (Lin et al., 2025). Theoretically, bandit
 122 algorithms (UCB (Garivier and Moulines, 2011),
 123 contextual (Foster and Rakhlin, 2020), hierarchical
 124 RL (Rohmatillah and Chien, 2023)) and posterior
 125 sampling (PSRL (Chen et al., 2020)) formalize
 126 exploration. However, their direct application to
 127 dialog POMDPs is hindered by non-stationarity
 128 and high dimensionality. Both lines of work lack a
 129 mechanism to perceive and respond to the dynamic
 130 cognitive affordances such as shifting uncertainty
 131 within a dialog. DyBBT bridges this gap by prag-
 132 matically adapting bandit inspiration to a learned,
 133 low dimensional cognitive state space \mathcal{C} , offering
 134 a tractable bridge between classical theory and se-
 135 quential dialog complexity.

136 2.3 Dual System Reasoning Architectures

137 Inspired by dual process theory (Krämer, 2014),
 138 recent work combines fast processing (System 1)
 139 with slow reasoning (System 2) for mathematical
 140 reasoning (Shi et al., 2024) and common sense
 141 inference (Yu et al., 2025). In TODS, large lan-
 142 guage models (LLMs) serve as powerful function
 143 approximators (Yi et al., 2025), acting as intu-
 144 itive generators (Ying et al., 2024) and delibera-
 145 tive reasoners (Ma et al., 2025). Frameworks
 146 like the Dynamic Dual Process Transformer (He
 147 et al., 2024) explicitly model this interaction for

policy learning. However, existing switching mechanisms often rely on static heuristics, such as fixed turn counts (Qin et al., 2023) or pre-defined confidence thresholds (Yao et al., 2023), which lack adaptability and theoretical grounding in exploration. DyBBT addresses this by introducing a meta-controller over a cognitive state space, dynamically triggering System 2 based on visitation counts and parametric uncertainty, offering a principled and efficient alternative to heuristic switching.

3 Methodology

DyBBT formulates dialog exploration as a tractable Contextual Multi-Armed Bandit (CMAB) problem over a cognitive state space \mathcal{C} , with theoretical grounding in Lipschitz smooth rewards and sublinear regret. This foundation enables a meta-controller that dynamically triggers System 2 based on visitation counts and confidence scores, balancing exploration and uncertainty in real time.

3.1 Theoretical Foundation

To provide a principled understanding of how exploration can be efficiently managed in dialog POMDPs, we develop a theoretical analysis that frames the problem as CMAB over a structured cognitive state space \mathcal{C} . This formulation rests on three pragmatic foundations: (1) compression of the high dimensional belief state into a low dimensional cognitive representation; (2) a Lipschitz smoothness assumption on the reward function; and (3) a derived bandit style exploration criterion. Together, these steps yield a tractable foundation that directly informs the design of our meta-controller.

3.1.1 CMAB Formulation

We formulate dialog policy learning as a CMAB problem (Foster and Rakhlin, 2020) to render the exploration-exploitation trade-off analytically tractable. The core innovation is a structured *cognitive state space* \mathcal{C} that compresses the high dimensional belief state into a low dimensional and interpretable representation, thereby bridging bandit theory with dialog POMDPs.

In this CMAB, the **arms** are $\mathcal{A} = \{S1, S2\}$ (fast inference System 1 vs. deliberative reasoning System 2). The **context** is the cognitive state $\mathbf{c}_t = [d_t, u_t, \rho_t] \in \mathcal{C}$ (Computation details in Appendix A.1), which captures dialog progress, user uncertainty, and slot dependency. The **reward** $r_t(a)$ measures task progress and efficiency when choosing arm a in context \mathbf{c}_t . The objective is to

minimize cumulative regret, where a_t^* is the optimal arm and a_t is the chosen arm at turn t .

$$R_T = \sum_{t=1}^T [\mathbb{E}[r_t(a_t^* | \mathbf{c}_t)] - \mathbb{E}[r_t(a_t | \mathbf{c}_t)]]. \quad (1)$$

This formulation treats S2 as an oracle arm. When it is invoked, aggressively pursues the optimal action in under explored regions. This directly informs the design of our meta-controller (Sec. 3.2.3).

3.1.2 Reward Smoothness Assumption

To support principled exploration in the cognitive state space \mathcal{C} , we assume the reward function is Lipschitz continuous (Asadi et al., 2018; Pazis and Parr, 2013; Ortner and Ryabko, 2012). This standard regularity condition ensures that nearby cognitive states yield similar rewards.

Assumption 3.1 (Lipschitz Smooth Reward in \mathcal{C}). *The expected immediate reward $\bar{r}(\mathbf{c}, a) = \mathbb{E}[r(s_t, a_t) | \mathbf{c}_t = \mathbf{c}]$ is Lipschitz continuous with respect to the cognitive state \mathbf{c} for any action a . That is, there exists a constant $L_r > 0$ such that:*

$$|\bar{r}(\mathbf{c}, a) - \bar{r}(\mathbf{c}', a)| \leq L_r \cdot d(\mathbf{c}, \mathbf{c}'), \quad \forall \mathbf{c}, \mathbf{c}' \in \mathcal{C}.$$

This assumption makes visitation counts meaningful proxies for epistemic uncertainty and allows bandit style exploration guarantees to carry over to the dialog POMDP. We empirically validate its practical relevance in Sec. 5.3.

3.1.3 Dynamic Balance Principle

Building on Assumption 3.1, we derive a principled exploration criterion for the meta-controller by formalizing the exploration-exploitation trade-off as a contextual bandit problem (Kleinberg et al., 2008; Bubeck et al., 2011) The exploration bonus for cognitive state \mathbf{c}_t is:

$$\text{Exploration-Bonus}(t) \propto \sqrt{\frac{\log T}{n_t(\mathbf{c}_t)}}, \quad (2)$$

where T is the total training steps and $n_t(\mathbf{c}_t)$ is the visitation count. This adapts the Upper Confidence Bound (UCB) principle (Ortner and Ryabko, 2012; Foster and Rakhlin, 2020) to the structured cognitive space, with the square root and logarithmic terms arising from concentration inequalities and horizon scaling (Komiya et al., 2024).

Eq. 2 motivates our meta-controller design. S2 is triggered when the exploration bonus exceeds a threshold, leading to Condition 1: $n_t(\mathbf{c}_t) <$

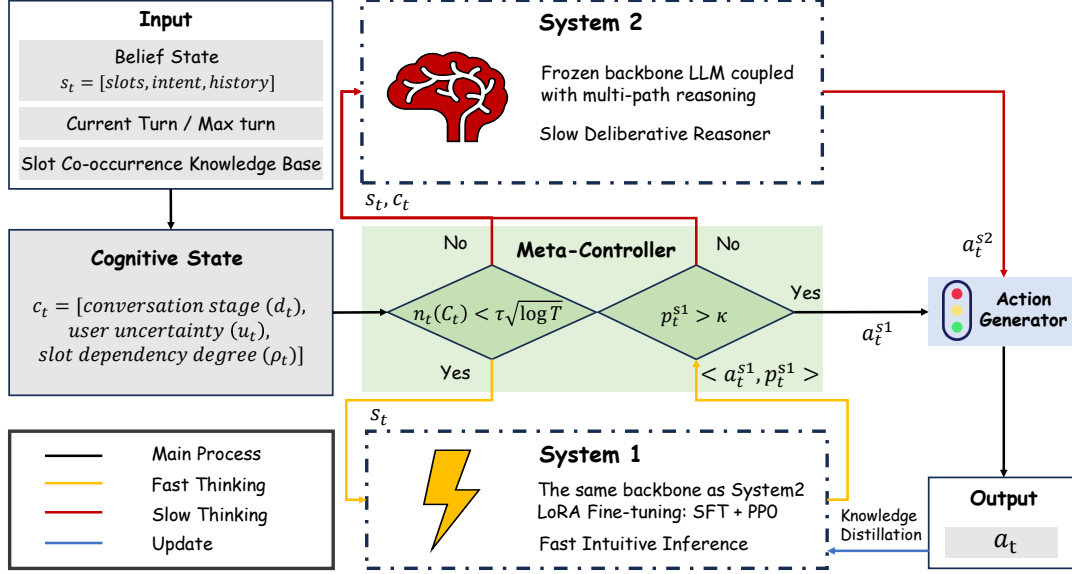


Figure 2: The DyBBT Architecture. A meta-controller uses the cognitive state c_t , visitation count $n_t(c_t)$, and System 1’s confidence p_t^{S1} to dynamically select between System 1 (fast intuitive) and System 2 (slow deliberative). Outputs drive action execution and update visitation/distillation buffers for continuous learning.

$\tau\sqrt{\log T}$. The scalar τ adjusts the trade-off, analogous to the confidence radius in UCB algorithms.

Under Assumption 3.1 and the approximate MDP structure in \mathcal{C} , this exploration strategy achieves sublinear regret (Proof sketch in Appendix A.2), demonstrating efficient and principled exploration in the compressed cognitive space.

3.2 System Architecture

Building on the theoretical foundation, DyBBT as shown in Fig. 2, operationalizes the CMAB formulation over the cognitive state space \mathcal{C} into a dual-system architecture. The meta-controller directly instantiates the bandit-inspired switching rule (Eq. 2) to dynamically balance between fast intuitive S1 and slow deliberative S2. This principled design ensures expensive S2 is invoked only when cognitive signals and visitation counts indicate high epistemic uncertainty or low confidence, achieving adaptive exploration-exploitation trade-off while maintaining computational efficiency.

3.2.1 S1: Fast Intuitive Inference

S1 serves as the low latency, high throughput policy for most dialog turns, avoiding the prohibitive cost of perpetual deliberation. Given the current belief state, S1 outputs both a system action a_t^{S1} formalized as a tuple $\{actiontype, domain, slot\}$, and a confidence score $p_t^{S1} \in [0, 1]$. This score captures *aleatoric uncertainty*, complementing the *epis-*

temic uncertainty tracked by the meta-controller. S1 is trained in two stages: supervised fine-tuning on expert trajectories to predict actions and calibrated confidence, followed by PPO optimization for task success and efficiency.

3.2.2 S2: Slow Deliberative Reasoner

S2 is invoked only for novel or complex states where S1 is likely to fail. It uses the same frozen base model as S1 to retain broad knowledge. Upon activation, S2 generates three distinct action sequences, evaluates each based on the ratio of filled key slots, and selects the first action of the highest rated sequence as a_t^{S2} . Although computationally expensive, S2 provides robust reasoning in high uncertainty or high stakes scenarios, as guided by the meta-controller.

3.2.3 Meta-Controller

The meta-controller implements the bandit inspired exploration criterion (Eq. 2) by dynamically selecting between S1 and S2 based on real time cognitive signals. This dual-trigger mechanism bridges bandit theory with practical dialog POMDPs:

Activate S2 If:

$$\underbrace{n_t(c_t) < \tau\sqrt{\log T}}_{\text{Condition 1}} \vee \underbrace{p_t^{S1} < \kappa}_{\text{Condition 2}}. \quad (3)$$

Condition 1: Exploration. This condition directly implements the exploration bonus from Eq. 2.

Under Assumption 3.1, low visitation counts in cognitive region c_t indicate high epistemic uncertainty, justifying systematic exploration via S2. The threshold $\tau\sqrt{\log T}$ adapts the classical bandit confidence radius to our structured cognitive space, ensuring exploration occurs when potential information gain outweighs computational cost.

Condition 2: Confidence. This condition addresses aleatoric uncertainty from partial observability and model limitations. Empirical studies show LLM confidence scores correlate with calibration (Kadavath et al., 2022; Lin et al., 2022; Yin et al., 2023); thus $p_t^{S1} < \kappa$ triggers S2 when S1’s parametric knowledge is likely insufficient.

This hybrid design ensures that S2 is invoked either for systematic exploration or as a robustness safeguard. The disjunctive combination yields an adaptive balance that outperforms either condition alone, as validated in our ablation study (Table 2). The meta-controller’s decisions form a closed loop system: high quality demonstrations from S2 are distilled into S1 via knowledge distillation (Appendix B.5.3), creating a virtuous cycle of policy improvement while progressively reducing long term reliance on costly deliberation.

4 Experiment

4.1 Experimental Setup

Datasets. We evaluate DyBBT on the Microsoft Dialog Challenge (Li et al., 2018) for single-domain tasks, and MultiWOZ 2.1 (Eric et al., 2020) for multi-domain tasks. Both are widely adopted in prior work. See Appendix B.1 for statistics.

Baselines. We compare DyBBT with a comprehensive set of recent and competitive baselines to ensure a rigorous evaluation, including previous SOTA EIERL. Full details are provided in Appendix B.2.

Evaluation Metrics. For single-domain tasks: success rate, average turns, and reward (following EIERL (Zhao et al., 2025): $+2t$ for success, $-t$ for failure, -1 for every turn). For multi-domain: Inform, Success, Book rates, and Avg. Turns (formulas in Appendix B.3).

Implementation Details. Following EIERL for fair comparison, dialogs are capped at 30 (single-domain) and 40 (multi-domain) turns. Training runs for 500 epochs (single) and 10K epochs (multi). DyBBT uses the same Qwen3 (0.6B-8B) for both S1 and S2. Full details in Appendix B.5.

4.2 Main Results Analysis

4.2.1 Performance on single-domain Tasks

The results in Table 1 show that DyBBT achieves strong performance across all three single-domain tasks. DyBBT’s cognitive state representation enables more efficient policy learning, leading to higher success rates with fewer turns compared to baselines. This advantage is especially pronounced in complex domains like Taxi, where slot dependencies create challenging exploration landscapes that DyBBT navigates effectively through its principled switching mechanism.

4.2.2 Performance on multi-domain Tasks

Appendix Table 6 presents results on the challenging MultiWOZ dataset. DyBBT maintains strong performance, whereas EIERL’s success rate drops significantly, revealing scalability limits of its population based approach. DyBBT-8B outperforms AutoTOD and ProTOD, and using GPT-4 as S2 achieves SOTA results, demonstrating that DyBBT matches strong LLM baselines with greater efficiency. This is enabled by the structured cognitive state and dual-system design, which provide domain-agnostic inductive bias without task-specific tuning. Cost-effectiveness analysis is discussed in Appendix E.7.

4.2.3 Training Efficiency and Convergence

Fig. 3 illustrates that DyBBT converges faster and achieves higher asymptotic performance than baselines across all domains, significantly outperforming EIERL as early as epoch 50. This accelerated learning stems from the meta-controller’s active guidance, which systematically targets under explored or uncertain regions in \mathcal{C} instead of relying on random or high variance exploration. DyBBT also scales consistently with model size: success rates improve from 80.3% to 89.5% in the single-domain Movie task and from 78.2% to 84.1% in multi-domain settings when scaling from 0.6B to 8B. This indicates that the dual-system architecture effectively harnesses model ability. Coupled with efficient resource allocation via the meta-controller and Qwen3’s native switching mechanism (Appendix E.8), DyBBT demonstrates practical viability for real world deployment.

4.2.4 Key Advantages of DyBBT

The main results demonstrate that DyBBT achieves SOTA performance through the following key advantages: **Dynamic Exploration Exploitation**

Domain	Agent	Epoch = 50			Epoch = 250			Epoch = 500		
		Success \uparrow	Reward \uparrow	Turns \downarrow	Success \uparrow	Reward \uparrow	Turns \downarrow	Success \uparrow	Reward \uparrow	Turns \downarrow
Movie	DQN_ε_0.0	35.05	-13.00	32.11	54.03	12.99	25.70	55.53	14.95	25.37
	DQN_ε_0.05	30.93	-18.61	33.44	67.95	31.84	21.39	76.68	43.42	19.21
	NOISY_DQN	41.37	-4.73	30.75	71.41	36.68	20.04	72.80	39.38	20.16
	LLM_DP	41.56	-3.09	27.34	41.56	-3.09	27.34	41.56	-3.09	27.34
	EIERL	23.72	-27.53	34.01	80.33	48.21	18.36	85.52	55.29	16.66
	DyBBT-0.6B	50.12	32.45	22.13	70.23	45.37	18.24	80.34	51.82	16.79
	DyBBT-1.7B	55.15	35.68	21.18	75.28	48.59	17.63	83.42	53.77	16.12
	DyBBT-8B	65.24	42.14	19.17	85.39	55.06	16.18	89.52	57.64	15.13
Rest.	DQN_ε_0.0	06.95	-36.57	27.66	49.07	4.10	22.13	56.71	11.63	23.22
	DQN_ε_0.05	07.26	-36.28	27.63	57.12	12.30	20.21	57.17	12.79	21.12
	NOISY_DQN	00.00	-43.92	29.84	16.69	-28.25	28.55	29.88	-15.20	26.18
	LLM_DP	38.96	-5.96	20.16	38.96	-5.96	29.16	38.96	-5.96	29.16
	EIERL	01.81	-41.09	27.44	69.75	24.79	17.98	79.35	34.99	16.07
	DyBBT-0.6B	46.73	20.5	21.67	65.44	28.83	17.86	74.85	33.08	16.52
	DyBBT-1.7B	51.32	22.59	20.71	70.14	30.90	17.25	77.71	34.24	15.85
	DyBBT-8B	60.70	26.74	18.69	79.54	35.05	15.81	83.38	36.74	14.86
Taxi	DQN_ε_0.0	00.04	-42.69	27.47	48.46	2.26	24.70	58.79	12.38	23.06
	DQN_ε_0.05	00.00	-42.86	27.71	55.98	8.19	22.38	66.83	20.19	21.90
	NOISY_DQN	00.00	-43.73	29.46	14.55	-30.56	29.32	26.15	-19.46	28.00
	LLM_DP	34.96	-10.23	25.95	34.96	-10.23	25.95	34.96	-10.23	25.95
	EIERL	00.00	-41.55	25.10	56.38	9.26	21.96	81.59	35.39	17.29
	DyBBT-0.6B	47.93	20.77	22.67	67.13	29.10	18.76	76.77	33.29	17.32
	DyBBT-1.7B	52.74	22.86	21.71	71.95	31.20	18.15	79.71	34.56	16.65
	DyBBT-8B	62.37	27.04	19.69	81.59	35.38	16.71	85.53	37.09	15.66

Table 1: Evaluation results for all agents across the three single-domain datasets are provided, with the highest value in each metric column highlighted in bold. Epochs (50, 250, 500) represent early, mid, and post convergence training stages. Baselines sourced from (Zhao et al., 2025).

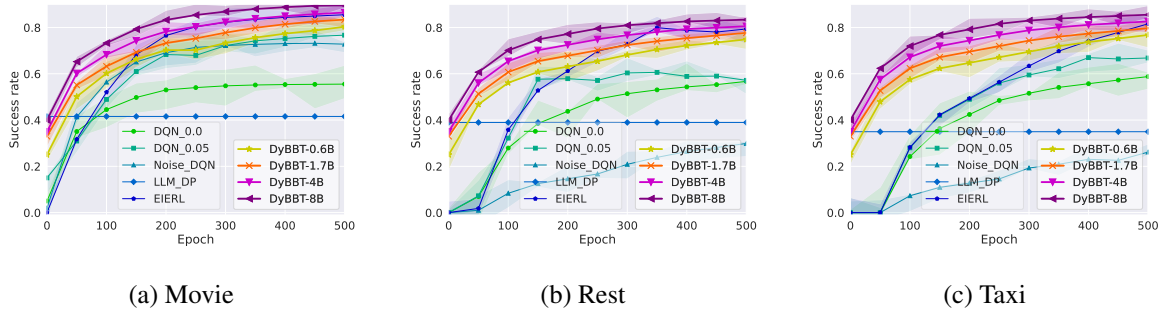


Figure 3: Learning curves for training efficiency and convergence across single-domain TODS tasks.

Balance: The bandit-inspired switching rule of the meta-controller enables DyBBT to allocate expensive S2 reasoning only when necessary, leading to highly efficient exploration. **Scalability with Model Size:** DyBBT benefits predictably from larger backbone models, making it well suited for future advances in LLM abilities. **Strong Generalization:** Consistent performance across both single- and multi-domain tasks shows that the cognitive state representation captures universal dialog dynamics. **Computational Practicality:** DyBBT maintains moderate computational overhead during both training and inference, unlike population based methods or full GPT-4.0 approaches.

4.3 Ablation Experiment

Ablation results are shown in Table 2, and detailed settings are in Appendix E.2. The results reveal that: **Meta-Controller is crucial.** Removing it causes the most severe performance degradation, confirming its essential role in dynamically orchestrating the exploration-exploitation trade-off. **Both conditions are necessary but asymmetric:** Removing Condition 1 (EC) eliminates the bandit inspired exploration bonus from Eq. 2, while removing Condition 2 (CC) disables the aleatoric uncertainty safeguard, a distinction rooted in Bayesian RL theory (Dearden et al., 1998). Removing the confidence condition (CC) causes a more substantial performance drop than removing the exploration condition (EC), validating our hybrid design.

Table 2: Ablation study of DyBBT’s components on MultiWOZ. Results underscore the necessity of the meta-controller and the structured cognitive state representation for optimal performance.

Variant	Inform \uparrow	Success \uparrow	Book \uparrow	Turns \downarrow
DyBBT-8B (full)	91.2	84.1	86.9	14.6
w/o Meta-Controller	82.5	71.8	77.3	17.5
w/o System 2	85.7	76.3	80.1	16.8
w/ Learned Cognitive State	90.5	83.2	86.3	14.8
w/o Knowledge Distillation	89.8	82.4	85.7	15.1
w/o Cognitive State (raw s_t)	84.2	75.1	79.6	17.1
w/o Exploration Condition (EC)	90.1	82.9	86.1	14.9
w/o Confidence Condition (CC)	87.6	79.5	83.2	16.2
w/o dialog Progress (d_t)	88.9	80.7	84.5	15.7
w/o User Uncertainty (u_t)	89.6	81.9	85.3	15.3
w/o Slot Dependency (ρ_t)	90.3	82.5	85.9	15.0

This indicates that mitigating S1’s overconfidence is slightly more critical than targeted exploration for robust performance. In depth error analysis (Appendix E.3) reveals that CC primarily prevents catastrophic failures in states with high cognitive uncertainty. **Cognitive State design is vital.** Replacing it with the raw belief state causes catastrophic performance collapse, confirming the necessity of our low dimensional, interpretable representation. While the learned alternative performs reasonably well, it still underperforms our hand-designed features, justifying our cognitively inspired approach. **All state dimensions contribute meaningfully.** Removing any single dimension causes noticeable performance degradation, with dialog progress (d_t) being the most impactful individual component, followed by user uncertainty (u_t) and slot dependency (ρ_t). **Knowledge Distillation enables continuous improvement.** Disabling it reduces final performance, confirming its role in facilitating long term efficiency gains through systematic learning from S2’s demonstrations.

4.4 Human and Real World Evaluation

We conducted controlled human evaluations and real world user experiments to validate DyBBT’s practical efficacy beyond automated metrics.

Human Evaluation. Following the protocol in Appendix C, 10 NLP researchers evaluated 200 dialog states from MultiWOZ, comparing DyBBT against random switching and S1-only baselines. Annotators rated action appropriateness (5-point Likert scale) and judged whether invoking S2 was justified. DyBBT’s actions were rated as more appropriate than both baselines, and its decisions to invoke S2 aligned significantly better with human judgment than random switching, confirming that

our meta-controller effectively identifies when de-liberation is warranted, an affordance often missed by heuristic approaches.

Real World User Experiments. As detailed in Appendix D, 30 volunteers completed multi-domain dialog tasks. DyBBT achieved the highest task success rate and user satisfaction, demonstrating that its cognitive state representation \mathcal{C} generalizes effectively beyond simulated environments. Case studies further showed that DyBBT successfully handles challenging scenarios, such as mid-dialog intent shifts and vague user expressions through adaptive S2 invocation.

Summary. These results provide converging evidence that DyBBT’s meta-controller translates cognitive affordances into a dynamic exploration exploitation balance, enabling robust performance in both controlled and real world settings.

5 Analysis

Our experimental results demonstrate that DyBBT achieves state-of-the-art performance on multiple benchmarks. In this section, we analyze the underlying mechanisms that enable DyBBT’s effectiveness, providing insights into why and how our framework works.

5.1 Analysis of Cognitive State Space

The cognitive state space \mathcal{C} enables bandit style exploration in dialog POMDPs by providing a low dimensional and interpretable representation of dialog dynamics. To validate its structure, we analyze the visitation frequency across discretized regions of \mathcal{C} during training (Fig. 4). The visitation heatmap reveals a structured, non-uniform pattern, confirming that exploration is guided by cognitive affordances: **In early dialog phases** ($d_t \in [0.0, 0.2]$), the meta-controller broadly explores across high user uncertainty (u_t) for information gathering. **In mid-phase** ($d_t \in [0.4, 0.6]$), visitation concentrates in medium-to-high u_t regions to resolve ambiguities. **In late phase** ($d_t > 0.8$), activity shifts to low u_t states, focusing on exploitation to complete tasks.

This phase dependent targeting demonstrates that \mathcal{C} effectively captures dialog progression and uncertainty, allowing the meta-controller to allocate exploration efficiently. The compactness and interpretability of \mathcal{C} make principled exploration feasible in high dimensional dialog state spaces.

507	5.2 Adaptive Balancing and Continuous	Practical Implementation. The consistent high	557
508	Improvement	performance of DyBBT using only a three dimensional	558
509	The meta-controller’s hybrid triggering mechanism	cognitive state demonstrates that the essential	559
510	robustly addresses the exploration exploitation	features governing exploration (dialog progress,	560
511	dilemma by responding to complementary forms	user uncertainty, slot dependency) can be distilled	561
512	of uncertainty. Condition 1 tackles <i>epistemic un-</i>	into a compact representation. This reduction in	562
513	<i>certainty</i> (lack of environmental knowledge), sys-	dimensionality is theoretically motivated by the	563
514	tematically exploring novel cognitive states, while	dependence of the regret bound’s $\sqrt{\dim(\mathcal{C})}$ (Ap-	564
515	Condition 2 addresses <i>aleatoric uncertainty</i> (inher-	pendix A.2.2).	565
516	ent stochasticity or model limitations), acting as		
517	a consistent safety net against over reliance on a	5.4 Failure Mode Analysis	566
518	potentially flawed S1. Analysis of 10,000 dialog	DyBBT exhibits three concrete failure modes em-	567
519	turns (Appendix Fig. 5) reveals their complemen-	pirically validated in Appendices E.9 and E.10.	568
520	tary temporal patterns: the exploration condition	First, reliance on handcrafted cognitive state fi-	569
521	dominates early in training and for novel states, en-	delity can lead to misrepresentation of complex di-	570
522	abling systematic coverage, whereas the confidence	alog dynamics, causing the meta-controller to mis-	571
523	condition provides ongoing robustness throughout	judge S2 invocation and resulting in underexplora-	572
524	training.	tion or computational waste. Second, dependency	573
525	This dual trigger design naturally evolves with	on high quality S2 demonstrations introduces risk;	574
526	training progress. Initially, frequent S2 invocations	errors in reasoning or self-evaluation can propagate	575
527	provide guided exploration and high quality demon-	to S1 via knowledge distillation, causing subtle	576
528	strations. As S1 improves via knowledge distilla-	policy corruption. Third, heuristic quantization of	577
529	tion from these demonstrations, the meta-controller	\mathcal{C} into a fixed number of bins masks critical state	578
530	automatically reduces S2 usage, transitioning from	variations, treating strategically distinct states iden-	579
531	guided exploration toward autonomous operation.	tically and reducing exploration efficacy. Qualita-	580
532	This virtuous cycle enables continuous policy im-	tative case studies illustrate how failures arise from	581
533	provement without additional environment interac-	unrepresented dialog nuances and how successful	582
534	tions. The distillation effectiveness is evidenced	interventions align with human judgment, and re-	583
535	by monotonic improvement in S1 performance and	veals these failures affect only 5.2% of dialogs,	584
536	corresponding reduction in S2 invocation rate. (Ap-	primarily in edge cases with abrupt intent shifts or	585
537	pendix Fig. 6). The adaptive balancing directly	complex dependencies, while built-in safeguards	586
538	embodies the framework’s ability to perceive and	provide substantial mitigation.	587
539	respond to dynamic dialog affordances, ensuring		
540	appropriate cognitive resource allocation through-	6 Conclusion	588
541	out the learning process.	DyBBT presents a novel dialog policy learning	589
542	5.3 Theoretical Intuitions and Empirical	framework that dynamically balances exploration	590
543	Alignment	and exploitation through a bandit inspired meta-	591
544	Our theoretical analysis, though based on simpli-	controller grounded in a structured cognitive state	592
545	fying assumptions, is pragmatically validated by	space. By formalizing dialog affordances into in-	593
546	empirical results: Sublinear Regret as Validation	interpretable dimensions: phasic progress, user un-	594
547	of Core Assumptions. The empirical cumulative	certainty, and slot dependency, our method en-	595
548	regret (Fig. 7) exhibits \sqrt{T} -like growth. This sub-	ables adaptive switching between fast intuitive re-	596
549	linear trend is not merely observational; it provides	sponses and deliberate reasoning. Extensive experi-	597
550	indirect empirical support for our key theoretical	ments across single- and multi-domain benchmarks	598
551	assumptions: The Lipschitz continuity of the re-	demonstrate SOTA performance in success rate,	599
552	ward in \mathcal{C} (Assumption 3.1), and the approximate	efficiency and generalization, with human evalua-	600
553	structure of MDP over \mathcal{C} (Assumption A.1). The	tions confirming superior decision alignment. Fu-	601
554	alignment between theory and experiment suggests	ture work will explore end-to-end learning of cog-	602
555	\mathcal{C} effectively captures the latent structure enabling	nitive representations and extend the framework to	603
556	efficient exploration. Low Dimensional \mathcal{C} Enables	more complex interactive settings.	604

605 Limitations

606 While DyBBT demonstrates strong empirical per-
607 formance, its reliance on hand designed cognitive
608 state representations, may not capture all nuances
609 of highly complex or novel dialog dynamics. This
610 points to a natural direction for future work: ex-
611 ploring end-to-end learning of cognitive represen-
612 tations that can adaptively refine the state space
613 from interaction data, thereby extending the frame-
614 work’s applicability to more diverse and intricate
615 interactive settings.

616 Ethics Statement

617 This work presents a dialog policy learning frame-
618 work evaluated on publicly available benchmark
619 datasets (MS Dialog and MultiWOZ) and their use
620 complies with the consent agreements established
621 during their original release. These datasets have
622 been previously anonymized and do not contain per-
623 sonally identifiable information or offensive con-
624 tent. Our research does not involve human subjects
625 beyond the use of standard datasets, and all experi-
626 ments are conducted through simulated user inter-
627 actions. In the human evaluation and real user ex-
628 periments, participants were volunteers proficient
629 in English and knowledgeable in NLP, recruited
630 from academic networks in Europe, the United
631 States, and China. No monetary compensation was
632 provided. Informed consent was obtained from all
633 participants, and the study protocol followed es-
634 tablished ethical guidelines for non-interventional
635 research. The proposed methodology focuses on
636 improving the efficiency of task-oriented dialog
637 systems, with potential positive societal impacts
638 through enhanced human computer interaction. We
639 are unaware of any specific ethical concerns or neg-
640 ative social impacts directly arising from this work.

641 We utilized DeepSeek V3.2 for translation assis-
642 tance and grammatical refinement of certain tex-
643 tual passages, and employed Qwen3-Code to aid in
644 debugging and optimizing portions of the experi-
645 mental code. These LLMs served solely as support
646 tools for improving linguistic clarity and technical
647 implementation. They played no role in the con-
648 ceptualization of the research. The authors assume
649 full responsibility for all aspects of the work, in-
650 cluding the accuracy and integrity of all generated
651 and modified content, and affirm that appropriate
652 measures have been taken to prevent plagiarism
653 and other forms of scientific misconduct.

References

- 654
655 Kavosh Asadi, Dipendra Misra, and Michael L. Littman. 656
657 2018. [Lipschitz continuity in model-based reinforce-](#) 658
659 [ment learning](#). In *Proceedings of the 35th Interna-*
660 *tional Conference on Machine Learning, ICML’18*.
661
662 Sébastien Bubeck, Rémi Munos, and Gilles Stoltz. 659
660 2011. [Pure exploration in finitely-armed and](#) 660
661 [continuous-armed bandits](#). *Theor. Comput. Sci.*,
662 412(19):1832–1852. 662
- 663 Xiuyi Chen, Fandong Meng, Peng Li, Feilong Chen,
664 Shuang Xu, Bo Xu, and Jie Zhou. 2020. [Bridging](#)
665 [the gap between prior and posterior knowledge se-](#) 665
666 [lection for knowledge-grounded dialogue generation](#).
667 In *Proceedings of the 2020 Conference on Empirical*
668 *Methods in Natural Language Processing (EMNLP)*,
669 pages 3426–3437, Online. Association for Computa-
670 tional Linguistics. 670
- 671 Richard Dearden, Nir Friedman, and Stuart Russell.
672 1998. [Bayesian q-learning](#). In *Proceedings of the*
673 *Fifteenth National Conference on Artificial Intelli-*
674 *gence and Tenth Innovative Applications of Artificial*
675 *Intelligence Conference, AAAI 98, IAAI 98, July 26-*
676 *30, 1998, Madison, Wisconsin, USA*, pages 761–768.
677 AAAI Press / The MIT Press. 677
- 678 Wenjie Dong, Sirong Chen, and Yan Yang. 2025. [Pro-](#) 678
679 [TOD: Proactive task-oriented dialogue system based](#) 679
680 [on large language model](#). In *Proceedings of the 31st*
681 *International Conference on Computational Linguis-*
682 *tics*, pages 9147–9164, Abu Dhabi, UAE. Associa-
683 tion for Computational Linguistics. 683
- 684 Yihan Du, R. Srikant, and Wei Chen. 2024. [Cascading](#) 684
685 [reinforcement learning](#). In *International Conference*
686 *on Representation Learning*, volume 2024, pages
687 30263–30304. 687
- 688 Mihail Eric, Rahul Goel, Shachi Paul, Abhishek Sethi,
689 Sanchit Agarwal, Shuyang Gao, Adarsh Kumar, Anuj
690 Goyal, Peter Ku, and Dilek Hakkani-Tur. 2020. [Mul-](#) 690
691 [tiWOZ 2.1: A consolidated multi-domain dialogue](#) 691
692 [dataset with state corrections and state tracking base-](#) 692
693 [lines](#). In *Proceedings of the Twelfth Language Re-*
694 *sources and Evaluation Conference*, pages 422–428,
695 Marseille, France. European Language Resources
696 Association. 696
- 697 Dylan J. Foster and Alexander Rakhlin. 2020. Beyond
698 ucb: optimal and efficient contextual bandits with
699 regression oracles. In *Proceedings of the 37th Inter-*
700 *national Conference on Machine Learning, ICML’20*.
701 JMLR.org. 701
- 702 Aurélien Garivier and Eric Moulines. 2011. On upper-
703 confidence bound policies for switching bandit prob-
704 lems. In *Algorithmic Learning Theory*, pages 174–
705 188, Berlin, Heidelberg. Springer Berlin Heidelberg. 705
- 706 James Jerome Gibson. 1979. *The Ecological Approach*
707 *to Visual Perception*. Houghton Mifflin, Boston.
708 Original edition. 708

709	Shuai Han, Wenbo Zhou, Jiayi Lu, Jing Liu, and Shuai Lü. 2022. Nrowan-dqn: A stable noisy network with noise reduction and online weight adjustment for exploration . <i>Expert Systems with Applications</i> , 203:117343.	765
710		766
711		
712		
713		
714	Tao He, Lizi Liao, Yixin Cao, Yuanxing Liu, Ming Liu, Zerui Chen, and Bing Qin. 2024. Planning like human: A dual-process framework for dialogue planning . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4768–4791, Bangkok, Thailand. Association for Computational Linguistics.	
715		
716		
717		
718		
719		
720		
721	Wanwei He, Yinpei Dai, Yinhe Zheng, Yuchuan Wu, Zheng Cao, Dermot Liu, Peng Jiang, Min Yang, Fei Huang, Luo Si, Jian Sun, and Yongbin Li. 2022. Galaxy: A generative pre-trained model for task-oriented dialog with semi-supervised learning and explicit policy injection . <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 36(10):10749–10757.	
722		
723		
724		
725		
726		
727		
728		
729	Xu Jia, Ruochen Zhang, and Min Peng. 2024. Multi-domain gate and interactive dual attention for multi-domain dialogue state tracking . <i>Knowledge-Based Systems</i> , 286:111383.	
730		
731		
732		
733	Saurav Kadavath, Tom Conerly, and et al. 2022. Language models (mostly) know what they know . <i>Preprint</i> , arXiv:2207.05221.	
734		
735		
736	Robert Kleinberg, Aleksandrs Slivkins, and Eli Upfal. 2008. Multi-armed bandits in metric spaces . In <i>Proceedings of the Fortieth Annual ACM Symposium on Theory of Computing</i> , STOC '08, page 681–690, New York, NY, USA. Association for Computing Machinery.	
737		
738		
739		
740		
741		
742	Junpei Komiyama, Edouard Fouché, and Junya Honda. 2024. Finite-time analysis of globally nonstationary multi-armed bandits . <i>Journal of Machine Learning Research</i> , 25(112):1–56.	
743		
744		
745		
746	Walter Krämer. 2014. Kahneman, D. (2011): Thinking, fast and slow . <i>Statistical Papers</i> , 55(3):915–915.	
747		
748	Wai-Chung Kwan, Hong-Ru Wang, Hui-Min Wang, and Kam-Fai Wong. 2023. A survey on recent advances and challenges in reinforcement learning methods for task-oriented dialogue policy learning . <i>Machine Intelligence Research</i> , 20:1–17.	
749		
750		
751		
752		
753	Jing Yang Lee, Kong Aik Lee, and Woon Seng Gan. 2023. An empirical Bayes framework for open-domain dialogue generation . In <i>Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)</i> , pages 192–204, Singapore. Association for Computational Linguistics.	
754		
755		
756		
757		
758		
759	Changqun Li, Linlin Wang, Xin Lin, Shizhou Huang, and Liang He. 2024. Hypernetwork-assisted parameter-efficient fine-tuning with meta-knowledge distillation for domain knowledge disentanglement . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1681–1695, Mexico	
760		
761		
762		
763		
764		
	City, Mexico. Association for Computational Linguistics.	765
		766
	Xiujun Li, Yu Wang, Siqi Sun, Sarah Panda, Jingjing Liu, and Jianfeng Gao. 2018. Microsoft dialogue challenge: Building end-to-end task-completion dialogue systems . <i>Preprint</i> , arXiv:1807.11125.	767
		768
		769
		770
	Anthony Liang, Guy Tennenholtz, Chih-Wei Hsu, Yinlam Chow, Erdem Biyik, and Craig Boutilier. 2024. Dynamite-rl: a dynamic model for improved temporal meta-reinforcement learning . In <i>Proceedings of the 38th International Conference on Neural Information Processing Systems</i> , NIPS '24, Red Hook, NY, USA. Curran Associates Inc.	771
		772
		773
		774
		775
		776
		777
	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words . <i>Transactions on Machine Learning Research</i> .	778
		779
		780
	Yuanguo Lin, Fan Lin, Guorong Cai, Hong Chen, Linxin Zou, Yunxuan Liu, and Pengcheng Wu. 2025. Evolutionary reinforcement learning: A systematic review and future directions . <i>Mathematics</i> , 13(5).	781
		782
		783
		784
	Shuai Ma, Qiaoyi Chen, Xinru Wang, Chengbo Zheng, Zhenhui Peng, Ming Yin, and Xiaojuan Ma. 2025. Towards human-ai deliberation: Design and evaluation of llm-empowered deliberative ai for ai-assisted decision-making . In <i>Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems</i> , CHI '25, New York, NY, USA. Association for Computing Machinery.	785
		786
		787
		788
		789
		790
		791
		792
	Volodymyr Mnih, Koray Kavukcuoglu, David Silver, and et al. 2015. Human-level control through deep reinforcement learning . <i>nature</i> , 518(7540):529–533.	793
		794
		795
	Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. 2017. Bridging the gap between value and policy based reinforcement learning . In <i>Proceedings of the 31st International Conference on Neural Information Processing Systems</i> , NIPS'17, pages 2772–2782, Red Hook, NY, USA. Curran Associates Inc.	796
		797
		798
		799
		800
		801
		802
	Xuecheng Niu, Akinori Ito, and Takashi Nose. 2024. Scheduled curiosity-deep dyna-q: Efficient exploration for dialog policy learning . <i>IEEE Access</i> , 12:46940–46952.	803
		804
		805
		806
	Ronald Ortner and Daniil Ryabko. 2012. Online regret bounds for undiscounted continuous reinforcement learning . In <i>Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2</i> , NIPS'12, page 1763–1771, Red Hook, NY, USA. Curran Associates Inc.	807
		808
		809
		810
		811
		812
	Jason Papis and Ronald Parr. 2013. Pac optimal exploration in continuous space markov decision processes . In <i>Proceedings of the Twenty-Seventh AAAI Conference on Artificial Intelligence</i> , AAAI'13, page 774–781. AAAI Press.	813
		814
		815
		816
		817

Yangyang Zhao, Kai Yin, Zhenyu Wang, Mehdi Das-
tani, and Shihan Wang. 2024. [Decomposed deep
q-network for coherent task-oriented dialogue policy
learning](#). *IEEE/ACM Transactions on Audio, Speech,
and Language Processing*, 32:1380–1391.

Qi Zhu, Christian Geischauser, Hsien-chin Lin, Carel van
Niekerk, Baolin Peng, Zheng Zhang, Shutong Feng,
Michael Heck, Nurul Lubis, Dazhen Wan, Xiaochen
Zhu, Jianfeng Gao, Milica Gasic, and Minlie Huang.
2023. [ConvLab-3: A flexible dialogue system toolkit
based on a unified data format](#). In *Proceedings of
the 2023 Conference on Empirical Methods in Nat-
ural Language Processing: System Demonstrations*,
pages 106–123, Singapore. Association for Compu-
tational Linguistics.

A Theoretical Details

This section provides the theoretical motivation
and intuition behind the DyBBT framework. The
following analysis bridges ideas from bandit the-
ory and cognitive science to create a heuristic for
exploration in dialog POMDPs. While the full dia-
log POMDP problem is intractable for a rigorous
minimax analysis, our goal is to provide a strong
conceptual foundation and explanatory power for
the algorithm’s design, which is then validated em-
pirically in the main text.

A.1 Formalization of Cognitive State Space

The cognitive state space \mathcal{C} is designed to be a low
dimensional, interpretable compression of the high
dimensional belief state s_t . We model \mathcal{C} as a com-
pact metric space with metric $d(\mathbf{c}, \mathbf{c}') = \|\mathbf{c} - \mathbf{c}'\|_2$.
Its covering dimension $\dim(\mathcal{C})$ is a measure of its
complexity. Given that our \mathcal{C} is defined by three
bounded dimensions ($d_t \in [0, 1]$, $u_t \in [0, 1]$, $\rho_t \in$
 $[0, 1]$), we have $\dim(\mathcal{C}) = 3$, which is crucial for
making bandit-style exploration feasible.

The choice of these three dimensions is moti-
vated by their central role in governing the
exploration-exploitation trade-off in TODS, draw-
ing inspiration from cognitive science and dialog
theory:

- **Dialog Progress** ($d_t = t/L$) captures the *tem-
poral affordance*. Early phases ($d_t \rightarrow 0$) in-
herently afford more exploration to gather in-
formation, while late phases ($d_t \rightarrow 1$) afford
exploitation to complete the task. This aligns
with the common practice of annealing explo-
ration schedules but provides a continuous,
state dependent signal.
- **User Uncertainty** operationalizes the *infor-
mation gathering affordance*.

$$u_t = |S_{unconfirmed}|/|S_{relevant}|$$

A high u_t indicates ambiguity in the user’s
goal, directly signaling the need for informa-
tion seeking actions to reduce entropy, a well
established principle in decision theory.

- **Slot Dependency** captures the *structural af-
fordance* of the task environment, derived
from a pre-computed slot co-occurrence ma-
trix M from the training corpus.

$$\rho_t = \max_{u \in U} \left(\frac{1}{|F|} \sum_{f \in F} M(u, f) \right)$$

A high ρ_t suggests that the next piece of infor-
mation is highly predictable given what is al-
ready known (e.g., requesting *departure* after
knowing *destination* in a taxi domain), mak-
ing targeted exploitation more efficient than
random exploration. This dimension encodes
the latent structure of the domain.

This design transforms the complex, unstructured
exploration problem in the raw belief space into a
more manageable one in a structured space where
states with similar exploration needs are grouped
together, as visualized in Fig. 4.

A.2 Regret Analysis Under Simplifying Assumptions

To provide theoretical intuition for our exploration
principle, we present a regret analysis under a set
of simplifying assumptions that capture the core
structure that we aim to exploit. This analysis justi-
fies the form of our exploration bonus and provides
an upper bound on learning speed. We make the
following assumptions to bridge the gap between
bandit theory and the dialog POMDP. Our analysis
is based on the Assumption 3.1 stated in Sec. 3.1.2,
which posits Lipschitz smoothness of the reward
function in the cognitive state space \mathcal{C} .

Assumption A.1 (MDP over \mathcal{C}). *The dialog pro-
cess can be approximately modeled as a finite hori-
zon MDP over the cognitive state space \mathcal{C} . The
transition dynamics and expected reward $\bar{r}(\mathbf{c}, a) =$
 $\mathbb{E}[r(s_t, a_t) | \mathbf{c}_t = \mathbf{c}]$ depend primarily on \mathbf{c}_t .*

The value function under a policy π in the cog-
nitive state space is defined as:

$$V^\pi(\mathbf{c}) = \mathbb{E} \left[\sum_{k=0}^H \gamma^k \bar{r}(\mathbf{c}_{t+k}, a_{t+k}) \mid \mathbf{c}_t = \mathbf{c}, a_{t+k} \sim \pi(\cdot | \mathbf{c}_{t+k}) \right]. \quad (4)$$

This assumption is a pragmatic simplification that allows us to focus on the core exploration challenge. It is reasonable if the cognitive state \mathbf{c}_t is a sufficient statistic for the exploration-exploitation trade-off, which our empirical results support.

A.2.1 Theoretical Intuition for Regret

Under Assumptions 3.1 and A.1, if we perform optimistic exploration in the cognitive state space \mathcal{C} , prioritizing states with low visitation counts, we can derive an upper bound on the expected cumulative regret that scales sublinearly with time:

$$\mathbb{E}[R(T)] \lesssim \tilde{\mathcal{O}} \left(L_r \cdot \sqrt{\dim(\mathcal{C}) \cdot T} \right), \quad (5)$$

where $R(T) = \sum_{t=1}^T [V^*(\mathbf{c}_t) - V^{\pi_t}(\mathbf{c}_t)]$ is the cumulative regret, and $\tilde{\mathcal{O}}$ hides logarithmic factors. The notation \lesssim indicates that this is a heuristic bound that captures the expected asymptotic scaling rather than a rigorous inequality. Here, L_r is the Lipschitz constant from Assumption 3.1, bounding the reward's sensitivity to changes in \mathcal{C} .

A.2.2 Derivation Sketch

This scaling can be motivated by discretizing the cognitive state space \mathcal{C} into $N = \mathcal{O}((1/\epsilon)^{\dim(\mathcal{C})})$ cells of diameter ϵ .

1. **Discretization Error:** Due to Lipschitz continuity of $\bar{r}(\mathbf{c}, a)$ (Assumption 3.1), the error introduced by discretization is bounded by $\mathcal{O}(L_r \epsilon T)$.
2. **Bandit Regret:** For the discretized MDP with N state cells, treating each cell arm analogously, a UCB like algorithm can achieve a regret bound of $\mathcal{O}(\sqrt{NT \log T})$.
3. **Optimization:** Balancing the two error terms by setting $\epsilon \sim T^{-1/(\dim(\mathcal{C})+2)}$ yields the final bound $\tilde{\mathcal{O}}(L_r \cdot \sqrt{\dim(\mathcal{C}) \cdot T})$.

This sketch illustrates that efficient learning is possible by exploiting the low dimensional structure and smoothness of the value function in \mathcal{C} , providing intuition for our exploration criterion.

This bound provides an intuitive justification for our exploration criterion (Eq. 2 in the main text). The term $\sqrt{\frac{\log T}{n_t(\mathbf{c}_t)}}$ is a heuristic adaptation of

the optimism principle, encouraging exploration of states with high uncertainty, inversely proportional to their visitation count. The empirical regret curve (Fig. 7) shows sublinear growth, consistent with this theoretical intuition.

A.3 Justification for the Meta-Controller Rule

The meta-controller's hybrid rule is designed for robust performance in the realistic setting where our theoretical assumptions hold only approximately:

Activate S2 IF:

$$\left(n_t(\mathbf{c}_t) < \tau \sqrt{\log T} \right) \vee \left(p_t^{S1} < \kappa \right). \quad (6)$$

The first condition, $n_t(\mathbf{c}_t) < \tau \sqrt{\log T}$, is the direct implementation of the theoretical exploration principle derived above. It addresses *epistemic uncertainty* (uncertainty reducible by exploration) by triggering System 2 in regions of \mathcal{C} that are under explored relative to the time horizon.

The second condition, $p_t^{S1} < \kappa$, is a critical *empirical safeguard* that addresses limitations of the theoretical model:

- **Partial Observability:** The true state of the user may not be fully captured by the belief state \mathbf{s}_t , leading to *aleatoric uncertainty*.
- **Model Imperfection:** S1, as a parameterized policy, may have inherent limitations and blind spots not captured by the visitation count.
- **Assumption Violation:** The Lipschitz smoothness assumption may locally break down.

A low confidence score p_t^{S1} is a proxy for these forms of uncertainty. This condition ensures robustness by invoking the powerful, knowledge rich S2 when S1 is uncertain, preventing catastrophic failures. The disjunctive (\vee) combination ensures System 2 is activated for *either* theoretical exploration *or* empirical robustness, making the overall system more adaptive and reliable than either condition alone, as evidenced by the ablation study (Table 2).

A.4 Discussion and Limitations

Our theoretical analysis provides a formal motivation for the DyBBT framework by illustrating how exploiting the structure of a cognitive state space

can lead to efficient exploration. However, we acknowledge its limitations, which also highlight the value of our empirical validation:

Simplified Model: Assumption A.1 reduces the POMDP to an MDP over \mathcal{C} , ignoring the challenges of belief state tracking and partial observability. This is a significant simplification. Our empirical results show that the algorithm performs well even when this assumption is not perfectly met, as the meta-controller’s confidence condition can mitigate some of these issues.

Heuristic Adaptation: The exploration bonus and the meta-controller rule are heuristic adaptations of the theoretical principle. A rigorous derivation for POMDPs remains an open challenge. Our contribution is to demonstrate that this heuristic is well motivated and highly effective in practice.

Empirical Safeguard: The confidence based condition, while crucial for performance, is not derived from the regret analysis. Its justification is empirical, stemming from its necessity for robust performance in ablation studies.

In conclusion, the theoretical analysis is not intended as a strict performance guarantee but rather as an *explanatory framework* that provides strong intuition for why exploring based on cognitive state visitation counts is a powerful principle. The ultimate validation of this principle, and its pragmatic implementation in the meta-controller, lies in its consistent empirical success across diverse dialog benchmarks.

B Experiment Details

B.1 Experimental Platform and Datasets

We evaluated DyBBT on two widely adopted benchmarks: the Microsoft dialog Challenge (MS dialog) ((Li et al., 2018)) for single-domain tasks, and the MultiWOZ 2.1 corpus ((Eric et al., 2020)) for multi-domain tasks. Both datasets are converted into ConvLab-3’s unified format, ensuring consistency in ontology, state representation, and API interaction. Table 3 summarizes the key statistics of both datasets.

The MS Dialog dataset comprises three distinct domains: Movie-Ticket Booking, Restaurant Reservation, and Taxi Ordering. It contains 7,215 dialogs with 89,465 turns, averaging 12.4 turns per dialog. The dataset is partitioned into training, validation, and test sets with 5,772, 722, and 721 dialogs, respectively.

The MultiWOZ 2.1 dataset is a large scale

multi-domain corpus spanning seven domains: Attraction, Hotel, Restaurant, Taxi, Train, Hospital, and Police. It includes 10,420 dialogs and 145,360 turns, with an average of 13.9 turns per dialog. The dataset is split into 8,420 dialogs for training, 1,000 for validation, and 1,000 for testing.

Both datasets provide annotated belief states, system dialog acts, and user goals, making them suitable for training and evaluating end-to-end dialog policies. The diversity in domain complexity, dialog length, and task structure across these datasets allows us to thoroughly assess the generalization capability of DyBBT in both single and multi-domain settings.

To ensure reproducibility and enable fair comparison, we implement and evaluate our proposed DyBBT framework using ConvLab-3 (Zhu et al., 2023), a flexible and unified toolkit for TODS. ConvLab-3 provides standardized data formats, integrated user simulators, and reinforcement learning utilities, facilitating consistent development and evaluation of dialog policies across multiple domains. All experiments are conducted using ConvLab-3’s builtin simulators and evaluation metrics, ensuring comparability across models and domains.

Dataset	Domains	Dialogs	Turns	Avg.Turns/Dialog
MS Dialog	3	7,215	89,465	12.4
MultiWOZ 2.1	7	10,420	145,360	14.0

Table 3: Summary of dataset statistics for MS Dialog and MultiWOZ 2.1.

B.2 Baselines Details

- **DQN $_{\epsilon-N}$** agents are trained using standard DQN (which realizes human level control through deep reinforcement learning) with a traditional $\epsilon - greedy$ exploration strategy, where $\epsilon = N$ (Mnih et al., 2015).
- **NOISY_DQN** agents enhance exploration by introducing noise into the network weights, based on the stable noisy network (NROWAN-DQN) with noise reduction and online weight adjustment (Han et al., 2022).
- **PG (REINFORCE)** is a stochastic gradient algorithm for policy gradient reinforcement learning, and its implementation refers to the flexible dialog system toolkit ConvLab-3 to

serve as a dialog policy baseline (Zhu et al., 2023).

- **PPO** is a policy optimization method in policy-based reinforcement learning that uses multiple epochs of stochastic gradient ascent and a constant clipping mechanism as the soft constraint for each policy update, with its implementation relying on the ConvLab-3 dialog toolkit (Zhu et al., 2023).
- **LLM_DP** agents replace the dialog policy (DP) module of the TODS with GPT-4.0 (drawing on advances in LLM based multi turn dialog systems) to select appropriate actions and pass them to the natural language generation (NLG) module for response generation (Yi et al., 2025).
- **AutoTOD** is a zero-shot autonomous agent based on GPT-4.0, which rethinks TODS by shifting from complex modularity to zero-shot autonomy and acts as a dialog policy baseline (Xu et al., 2024).
- **ProTOD** is a proactive TODS policy based on GPT-4.0, designed as a proactive dialog system to optimize the process of task oriented interactions (Dong et al., 2025).
- **EIERL** is an evolutionary reinforcement learning method for TODS policies, which improves the efficiency of dialog policy learning by injecting elite individuals into the evolutionary process (Zhao et al., 2025).
- **MACRM** is a multi agent curiosity reward model for TODS, which optimizes dialog policies through collaborative interactions among multiple agents and curiosity driven reward mechanisms (Sun et al., 2025).

B.3 Metrics Formula

This section provides the formal definitions of the evaluation metrics used for multi-domain TODS evaluation, following the standard MultiWOZ evaluation protocol.

B.3.1 Inform Success Rate

The Inform Success Rate measures the system’s ability to provide all requested information to the user. Let G be the goal specification, D be the set of dialog domains, and S be the sequence of

system dialog acts. For each domain $d \in D$, let R_d be the set of requested slots in the goal:

$$\text{TP} = \sum_{d \in D} \sum_{s \in R_d} \mathbb{I}(\exists \text{inform}(d, s, v) \in S \wedge v \notin V_{\text{null}}) \quad (7)$$

$$\text{FP} = \sum_{d \in D} \sum_{s \notin R_d \cup I_d} \mathbb{I}(\exists \text{inform}(d, s, v) \in S \wedge v \notin V_{\text{null}}) \quad (8)$$

$$\text{FN} = \sum_{d \in D} \sum_{s \in R_d} \mathbb{I}(\nexists \text{inform}(d, s, v) \in S \vee v \in V_{\text{null}}) \quad (9)$$

where $V_{\text{null}} = \{“”, “dont care”, “not mentioned”\}$ represents null values. The Inform Success Rate is then defined as:

$$\text{Inform} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (10)$$

B.3.2 Book Success Rate

The Book Success Rate evaluates the system’s ability to successfully complete booking operations. For each domain $d \in D$ that requires booking, let B_d be the set of booking constraints in the goal. The booking success is computed as:

$$\text{Book}_d = \frac{1}{|B_d|} \sum_{b \in B_d} \mathbb{I}(\text{book}(d, b, v) \in S \wedge v = v_{\text{goal}}) \quad (11)$$

For the taxi domain (which has no database constraints), booking success is trivially 1 if any booking action occurs:

$$\text{Book}_{\text{taxi}} = \mathbb{I}(\exists \text{book}(\text{taxi}, \cdot, \cdot) \in S) \quad (12)$$

The overall Book Success Rate is the average across all booking domains:

$$\text{Book} = \frac{1}{|D_{\text{book}}|} \sum_{d \in D_{\text{book}}} \text{Book}_d \quad (13)$$

where D_{book} is the set of domains requiring booking.

1276 B.3.3 Success Rate

1277 The Success Rate represents the overall task com-
1278 pletion performance, combining both information
1279 provision and booking success:

$$1280 \text{Success} = \mathbb{I}(\text{Inform} = 1 \wedge \text{Book} = 1) \quad (14)$$

1281 This binary metric indicates whether both all re-
1282 quested information was provided and all booking
1283 operations were successfully completed.

1284 This metric rewards systems that achieve high
1285 success rates with fewer dialog turns, promoting
1286 both effectiveness and efficiency.

1287 B.4 Prompt for DyBBT and LLM-DP

1288 This appendix provides the detailed prompts used
1289 for System 1 (intuitive controller) and System 2
1290 (reasoning controller) in the DyBBT framework.
1291 The LLM_DP prompt is the same from the EIERL
1292 paper((Zhao et al., 2025)).

1293 B.4.1 System 1 Prompt

```
[caption={System 1 prompt.}]
You are the fast, intuitive component (System
  1) of a task oriented dialog system. Your
  task is to generate the next system action
  based solely on the current belief state.
  Do not reason step-by-step. Output your
  first, most intuitive response in the exact
  JSON format specified.

**Current Belief State:**
{belief_state}

**Available Actions:**
{available_actions}

Based on the above, output ONLY a valid JSON
  object with your predicted action and its
  confidence. Do not output any other text.

{
  "action": [
    ["<act_type>", "<domain>", "<slot>"],
    ["<act_type>", "<domain>", "<slot>"],
    ...
  ],
  "confidence": <confidence_score>
}
```

1294 B.4.2 System 2 Prompt

```
[caption={System 2 prompt.}]
You are the deliberative reasoner (System 2) of
  a task oriented dialog system. Your goal
  is to generate diverse, high quality action
  plans when the meta-controller detects a
  need for deeper reasoning, either due to
  unfamiliar cognitive states or low
  confidence from System 1.
```

```
**Current Belief State:**
{belief_state}

**Available Actions:**
{available_actions}

**Cognitive State Context:**
- dialog Progress: {d_t}
- User Uncertainty: {u_t}
- Slot Dependency: {p_t}

**Trigger Reason:** {trigger_reason}

**Reasoning Guidelines:**
1. **Leverage cognitive signals:**
  - If progress is low, focus on information
    gathering.
  - If uncertainty is high, prioritize
    clarifying or confirming actions.
  - If slot dependency is high, leverage known
    slot relationships to guide next
    actions.
2. **Consider domain and slot dependencies:**
  - E.g., 'taxi' requires both 'destination'
    and 'departure'; 'restaurant' may
    require 'area', 'food', 'pricerange'
    before booking.
3. **Generate 3 distinct strategies** that
  reflect different tactical approaches:
  - One conservative (e.g., confirm before
    acting),
  - One proactive (e.g., request multiple
    slots),
  - One hybrid (e.g., inform then request).
4. **Evaluate each path** by estimating its
  likelihood of leading to task success.

**Output Format:** Strictly adhere to the
  following JSON schema:

{
  "reasoning_paths": [
    {
      "sequence_id": 1,
      "action_sequence": [
        ["action_type", "domain", "slot"],
        ...
      ],
      "estimated_success_probability": 0.9
    },
    ...
  ]
}
```

1297 B.4.3 LLM_DP Prompt

```
[caption={Dialog policy Prompt for LLM.}]
You must strictly execute the following
  commands:
1. Command execution requirements: when
  receiving a command, you must strictly
  follow the given instructions without
  performing any actions outside the scope of
  the command or generating any additional
```

words.

2. Datasets and system roles: as the dialog policy component in a task oriented dialog system, you will make system decisions based on the MultiWOZ 2.1 dataset.
3. Processing user dialog state: you will receive a formatted user dialog state. This state will be used as a basis for decision making.
4. Generate system actions: based on the user dialog state {


```

      'user_action': [["Inform", "Hotel", "Area", "east"], ["Inform", "Hotel", "Stars", "4"]],
      'system_action': [],
      'belief_state': {
        'police': {'book': {'booked': []}, 'semi': {}},
        'hotel': {'book': {'booked': [], 'people': '', 'day': '', 'stay': ''}, 'semi': {'name': '', 'area': 'east', 'parking': '', 'pricerange': '', 'stars': '4', 'internet': '', 'type': ''}},
        'attraction': {'book': {'booked': []}, 'semi': {'type': '', 'name': '', 'area': ''}},
        'restaurant': {'book': {'booked': [], 'people': '', 'day': '', 'time': ''}, 'semi': {'food': '', 'pricerange': '', 'name': '', 'area': ''}},
        'hospital': {'book': {'booked': []}, 'semi': {'department': ''}},
        'taxi': {'book': {'booked': []}, 'semi': {'leaveAt': '', 'destination': '', 'departure': '', 'arriveBy': ''}},
        'train': {'book': {'booked': [], 'people': ''}, 'semi': {'leaveAt': '', 'destination': '', 'day': '', 'arriveBy': '', 'departure': ''}}
      },
      'request_state': {},
      'terminated': False,
      'history': []
    }, you need to generate system actions. These actions should be provided in the following format: [["ActionType", "Domain", "Slot", "Value"]] where `ActionType` denotes the type of action (e.g. Request, Inform, Confirm, etc.), `Domain` specifies the associated domain (e.g. restaurant, taxi, hotel, etc.), `Slot` is the specific information slot associated with the action (e.g. name, area, type, etc.), and `Value` is the corresponding value or an empty string.
```

B.5 Implementation Details

The DyBBT framework was implemented within the Convlab-3 dialog system with Python 3.10 environment ((Zhu et al., 2023)), leveraging its modular architecture for efficient dialog policy optimiza-

tion. We employed RuleDST for system dialog state tracking and RulePolicy for user policy simulation, eliminating the need for natural language understanding (NLU) and natural language generation (NLG) modules. This design choice significantly enhances training efficiency by reducing computational overhead and isolating the impact of language processing components from policy learning performance. The dialog environment was configured with a maximum turn limit of 30 for single-domain and 40 for multi-domain (the same as EIERL) interactions per episode, with the cognitive state space \mathcal{C} computed in real-time during dialog execution using dimensions including dialog progress (d_t), user uncertainty (u_t), and slot dependency (ρ_t) extracted from the belief state representation provided by RuleDST.

User goals were dynamically generated using the GoalGenerator module, which produces diverse and realistic TODS objectives across single or multiple domains. This approach ensures training data variety and generalization capability, consistent with REINFORCE and PPO training methodologies. The goal generation process excluded the police domain due to its low data quality, ensuring higher reliability in evaluation.

All experiments were conducted on NVIDIA 5090 GPUs with 32GB memory. System 1 was SFT using the AdamW optimizer with a learning rate of 1×10^{-4} and further optimized via PPO, employing a clipping parameter $\epsilon = 0.2$ and GAE with $\lambda = 0.95$. The meta-controller employs a dual-threshold mechanism for System 2 invocation, with $\kappa = 0.7$ and $\tau = 1.0$, values selected via grid search over development sets as they maximize both performance and robustness across domains. These thresholds operate on a discretized 5 bins cognitive state space, which balances expressiveness and generalization, as validated in Sec. E.5.

We maintained a replay buffer with a capacity of 10,000 transitions, using a batch size of 32 for training. A separate knowledge distillation buffer was managed under a FIFO replacement policy with a fixed capacity. To ensure reproducibility, all experiments were run with five fixed random seeds (9841, 35741, 91324, 8134, 13924), consistent with the EIERL baseline (Zhao et al., 2025). All hyperparameters were selected through grid search on a validation subset of the MultiWOZ data.

Training was conducted for 500 epochs on single-domain tasks and 10,000 epochs on multi-domain tasks, incorporating early stopping with a

patience of 3 epochs based on validation performance. This protocol aligns with the EIERL setup for fair comparison.

B.5.1 Slot Co-occurrence Matrix Construction

The slot dependency dimension ρ_t in the cognitive state space \mathcal{C} is derived from a co-occurrence matrix M that captures statistical relationships between dialog slots across the Microsoft dialog Challenge ((Li et al., 2018)) and MultiWOZ ((Eric et al., 2020)) dataset. This matrix quantifies the conditional probability that slot j appears given the presence of slot i , providing a principled measure of semantic relatedness between dialog concepts.

Formally, the co-occurrence matrix $M \in \mathbb{R}^{N \times N}$ is constructed from the training partition of MultiWOZ 2.1, where N represents the total number of unique slot types across all domains. For each dialog turn containing belief state updates, we extract the set of active slots (those with non-empty values) and update the co-occurrence counts. The matrix elements are computed as:

$$M_{ij} = \frac{\text{count}(\text{slot}_i \wedge \text{slot}_j)}{\text{count}(\text{slot}_i)} \quad (15)$$

where $\text{count}(\text{slot}_i \wedge \text{slot}_j)$ denotes the number of dialog turns where both slots appear simultaneously, and $\text{count}(\text{slot}_i)$ represents the total occurrences of slot i . This normalization ensures that M_{ij} represents the empirical conditional probability $P(\text{slot}_j | \text{slot}_i)$.

The slot dependency ρ_t for a given belief state s_t is then computed as the average co-occurrence strength between the currently active slots:

$$\rho_t = \frac{1}{|A_t|(|A_t| - 1)} \sum_{i \in A_t} \sum_{j \in A_t, j \neq i} M_{ij} \quad (16)$$

where A_t denotes the set of slots with non-empty values in the current belief state. This formulation captures the structural complexity of the dialog context, with higher values indicating greater semantic interdependence between the information being discussed.

The construction of M leverages the statistical regularities present in TODS, where certain slot combinations naturally co-occur due to domain-specific constraints and user behavior patterns. For instance, in restaurant booking scenarios, slots like *restaurant-area* and *restaurant-food* frequently

appear together, while in hotel domains, *hotel-pricerange* and *hotel-type* exhibit strong associations. This matrix based approach provides a data-driven foundation for quantifying dialog complexity that complements the theoretically motivated dimensions of dialog progress and user uncertainty.

B.5.2 Training Details For System 1

To train System 1 for accurate action prediction and confidence estimation, we employ a two-stage training methodology comprising supervised fine-tuning (SFT) followed by reinforcement learning. This approach utilizes dialog sequences from the MultiWOZ and MS Diag dataset to develop a robust policy model capable of rapid decision making with calibrated confidence scores.

Stage 1: Supervised Fine-tuning with Data Augmentation

We first construct a training corpus of 10,000 single turn dialog samples through systematic data augmentation. For each dialog turn, we extract the belief state s_t , available action set \mathcal{SA} , and ground truth system actions a_t^* . The initial confidence score p_t^{S1} is sampled from $\mathcal{U}(0.95, 1.0)$.

The augmentation process introduces controlled perturbations to simulate prediction uncertainty. For each ground truth action sequence a_t^* , we apply three modification operations with specified probabilities: 20% action addition by sampling new actions from \mathcal{SA} ; 60% action modification through substitution with random actions from \mathcal{SA} ; and 20% action deletion while ensuring the augmented sequence a_t' maintains at least one action. These operations are applied sequentially in random order to each sample (Kadavath et al., 2022; Lin et al., 2022; Yin et al., 2023). The confidence score is adjusted proportionally to the modification intensity:

$$p_t^{S1} \leftarrow p_t^{S1} \cdot \left(1 - \frac{n_{\text{mod}}}{n}\right),$$

where n denotes the original action sequence length and n_{mod} represents the number of modified actions. This procedure generates a dataset with confidence scores approximately uniformly distributed in $[0, 1]$.

For SFT training, the model takes s_t and \mathcal{SA} as inputs and produces both action sequence a_t^{S1} and confidence score p_t^{S1} as outputs. The composite loss function integrates action prediction and confidence estimation:

$$\mathcal{L} = \lambda \mathcal{L}_a + (1 - \lambda) \mathcal{L}_p,$$

where $\lambda = 0.7$. The action loss \mathcal{L}_a employs cross-entropy to measure divergence between predicted and augmented actions:

$$\mathcal{L}_a = - \sum_i \log P(a_t^{S1} = a'_t | s_t, \mathcal{SA}),$$

while the confidence loss \mathcal{L}_p utilizes mean squared error:

$$\mathcal{L}_p = (p_t^{S1} - p_t^{\text{target}})^2.$$

Stage 2: Reinforcement Learning with PPO

The second stage employs PPO to optimize dialog level performance metrics using the complete MultiWOZ dataset. The reward function R combines multiple objectives:

$$R = R_{\text{success}} + R_{\text{efficiency}} + R_{\text{penalty}},$$

where $R_{\text{success}} = +2t$ for successful dialogs and $-t$ for failures (t denotes the max turn number), $R_{\text{efficiency}} = -1$ per dialog turn to encourage conciseness, and R_{penalty} captures additional constraints.

This two-stage approach enables System 1 to initially learn accurate action confidence mappings through supervised learning, then refine its policy for improved task completion efficiency and success rates via reinforcement learning.

B.5.3 Knowledge Distillation Buffer Management

To form a virtuous cycle and reduce long term dependence on System 2, high quality decisions (s_t, a_t^{S2}) from System 2 are stored in a distillation buffer D_{distill} . We only store decisions where System 2’s self evaluated task completion probability is greater than 0.9, ensuring high quality distillation data. Periodically (every 10 training epochs), System 1 is fine-tuned on these data via Low-Rank Adaptation (LoRA) with a learning rate of 1×10^{-4} , batch size of 4, and gradient accumulation steps of 8. This SFT approach distills the knowledge gained through costly deliberation into an efficient intuitive policy while maintaining computational efficiency, leading to a monotonic performance improvement. Over time, this reduces the need to invoke System 2 for previously challenging states, thereby increasing overall efficiency.

The knowledge distillation buffer D_{distill} stores high quality pairs (s_t, a_t^{S2}) generated by System 2. The buffer has a maximum capacity and uses an FIFO policy to maintain data freshness and diversity. We employ LoRA fine-tuning with rank

Algorithm 1 Knowledge Distillation Buffer Update and Sampling

Buffer Update:

- 1: **Input:** Current belief state s_t , System 2 action a_t^{S2} , System 2 self evaluated confidence p_{self}
- 2: **if** $p_{\text{self}} > 0.9$ **then** ▷ Only store high confidence actions
- 3: **if** $|D_{\text{distill}}| < \text{MAX_SIZE}$ **then**
- 4: $D_{\text{distill}}.\text{append}((s_t, a_t^{S2}))$
- 5: **else**
- 6: $D_{\text{distill}}.\text{pop_front}()$ ▷ Remove oldest entry (FIFO)
- 7: $D_{\text{distill}}.\text{append}((s_t, a_t^{S2}))$
- 8: **end if**
- 9: **end if**

System 1 Fine-tuning:

- 10: **Input:** System 1 model with LoRA adapters, buffer D_{distill}
- 11: **Every 10 training epochs:**
- 12: **for** $epoch = 1$ **to 1** **do** ▷ Fine-tune for 1 epoch
- 13: **for** each batch sampled from D_{distill} **do**
- 14: Compute loss $\mathcal{L} = \text{CrossEntropy}(\text{System1}(s_i), a_i)$
- 15: Update LoRA adapter parameters via gradient descent
- 16: **end for**
- 17: **end for**

$r = 16$, scaling parameter $\alpha = 32$, and dropout rate of 0.1, targeting the query and value projection layers of the transformer architecture. This configuration achieves parameter efficiency while preserving the base model’s generalization capabilities.

B.5.4 Visitation Count of the Cognitive State Space

To compute the visitation count $n_t(\mathbf{c}_t)$ for the continuous cognitive state space \mathcal{C} , we discretize each dimension of $\mathbf{c}_t = [d_t, u_t, \rho_t]$ into 5 uniformly spaced bins over the range $[0, 1]$. The cognitive state is then mapped to a discrete tuple $(d_{\text{bin}}, u_{\text{bin}}, \rho_{\text{bin}})$, and $n_t(\mathbf{c}_t)$ is the cumulative visitation count of that bin tuple.

This choice of dimensions is motivated by cognitive and dialog theory, which highlights stage, uncertainty, and structural relationships as key factors influencing decision making. By quantifying these environmental affordances into a structured cognitive state space \mathcal{C} , we create a formal bridge between Gibson’s ecological perception theory and

practical dialog policy optimization. While not exhaustive, this representation aims to capture the most salient features for guiding exploration. Its empirical necessity and sufficiency are validated through ablation studies in Sec. 4.3. We define \mathcal{C} as the cognitive state space, assumed to be a compact subset of \mathbb{R}^3 equipped with the Euclidean metric $d(\mathbf{c}, \mathbf{c}')$.

B.5.5 Calculation of Empirical Cumulative Regret

To empirically validate the theoretical intuition of sublinear regret growth under our simplifying assumptions, we compute the **empirical cumulative regret** $R_{\text{emp}}(T)$ during training, as shown in Fig. 7. The regret is defined as:

$$R_{\text{emp}}(T) = \sum_{t=1}^T \left(V^{\pi^*}(\mathbf{s}_t) - V^{\pi_t}(\mathbf{s}_t) \right)$$

where:

- T is the total number of dialog turns (training steps) up to the current point.
- \mathbf{s}_t is the belief state at turn t .
- $V^{\pi_t}(\mathbf{s}_t)$ is the actual discounted return obtained from state \mathbf{s}_t under the current policy π_t at training step t .
- $V^{\pi^*}(\mathbf{s}_t)$ is the value of the near-optimal policy π^* at state \mathbf{s}_t .

Since the true optimal policy π^* is unknown, we approximate it using a strong baseline policy the fully trained DyBBT-8B/GPT-4.0 model, which achieves SOTA performance on MultiWOZ. We assume this policy is sufficiently close to optimal for regret estimation purposes. For each state \mathbf{s}_t , we estimate $V^{\pi^*}(\mathbf{s}_t)$ by running π^* from \mathbf{s}_t for multiple episodes and averaging the discounted returns. Actual episodic return is used from the current dialog episode as a proxy for $V^{\pi_t}(\mathbf{s}_t)$. Although this is a coarse approximation, it is standard in episodic RL settings and sufficient to capture the regret trend.

$R_{\text{emp}}(T)$ is plotted against T on a log-log scale to clearly visualize the sublinear growth trend. The theoretical upper bound $\tilde{O}(\sqrt{T})$ is plotted alongside for comparison. The constant factor in the theoretical bound is fit to the empirical curve in the early training phase to align the curves for illustrative purposes.

C Human Evaluation Details

This appendix provides comprehensive details of the human evaluation study described in Sec. 4.4. The study was designed to qualitatively assess the core contribution of the DyBBT framework: the intelligent, adaptive decision making of its meta-controller, beyond what is captured by automated metrics.

C.1 Annotation Protocol and Interface

Evaluators were presented with a structured web interface for each evaluation instance. Each instance consisted of a single dialog *state* (not a full dialog), sampled from the MultiWOZ test set. For a given state, the interface displayed the following information:

- **Dialog Context:** The last user utterance and the last system action to provide conversational context.
- **Current Belief State (\mathbf{s}_t):** A structured table showing all relevant slots for the domain(s), their values, and their confirmation status (e.g., *confirmed*, *requested*, *None*).
- **Cognitive State (\mathbf{c}_t):** The numerical values for dialog progress (d_t), user uncertainty (u_t), and slot dependency (ρ_t).
- **System Action:** The action chosen by the model for this state, presented in a structured format (e.g., [*request*, *restaurant*, *area*, “ ”]).
- **System Variant:** The name of the model variant that produced the action (DyBBT, S1-only, Random Switching). Variants were anonymized as ‘System A’, ‘System B’ during evaluation to avoid bias.

Evaluators were then asked to answer two questions based solely on the provided information:

1. **Action Appropriateness:** “How appropriate is the system’s chosen action given the current dialog state?” Rated on a 5 points Likert scale:
 1. Very Inappropriate
 2. Somewhat Inappropriate
 3. Neutral
 4. Somewhat Appropriate
 5. Very Appropriate

Table 4: Complete Human Evaluation Results. The Action Appropriateness score is the average Likert score (1-5). The Switching Agreement is the percentage of states where the model’s decision to *not* invoke System 2 aligned with the majority of human annotators.

Model Variant	Action Appropriateness \uparrow	Switching Agreement \uparrow
DyBBT-8B	4.31 ± 0.12	88.7%
w/o Meta-Controller (Random)	3.72 ± 0.19	52.3%
w/S1-only	3.95 ± 0.15	—
w/o Exploration Condition (EC)	4.08 ± 0.14	75.4%
w/o Confidence Condition (CC)	3.89 ± 0.16	81.2%

- Switching Judgment:** “In this specific situation, would it be justified to invoke a powerful, but computationally expensive, reasoning module to choose the action?” Answered with **Yes** or **No**. This question was only shown for states where the evaluated model *did not* invoke System 2, to directly test if the meta-controller’s decision *not* to invoke aligned with human judgment.

C.2 Annotator Background and Training

We recruited **10 annotators**, all of whom were graduate students or researchers with a background in natural language processing and familiarity with TODS. Prior to the evaluation, a mandatory 30 minutes training session was conducted. The session:

- Explained the goal of the evaluation and the definition of key concepts (belief state, system actions, computational cost).
- Walked through 5 example states that were not part of the evaluation set, discussing potential appropriate actions and reasoning for/against invoking a costly reasoner.
- Allowed annotators to ask questions to resolve any ambiguities.

Annotators were compensated at a competitive hourly rate for their work.

C.3 Human Evaluation Results

The results in Table 4 provide a detailed breakdown supporting the main findings:

- Superior Decision Quality:** The full DyBBT model yields a higher action appropriateness score than the ablated variants.
- Value of the Meta-Controller:** The random switching variant has the lowest scores, confirming that a naive switching strategy severely degrades decision quality and is not aligned with human judgment.

- Complementary Role of Both Conditions:** Removing either the Exploration Condition (EC) or the Confidence Condition (CC) leads to a drop in both appropriateness and agreement, with the CC being slightly more critical for action quality (preventing poor actions) and the EC being crucial for efficient switching (preventing unnecessary calls). This validates their hybrid design in the meta-controller.

C.4 Qualitative Analysis of Meta-Controller Decisions

To qualitatively validate the efficacy of the meta-controller’s switching mechanism beyond aggregate metrics, we present two contrasting case studies sampled from the MultiWOZ test set. These examples illustrate how DyBBT’s principled switching aligns with human judgment, in contrast to a naive baseline.

Case 1: High Agreement Example (DyBBT).

The meta-controller correctly identified a state warranting costly deliberation due to high *aleatoric uncertainty* despite the cognitive state being well explored. The belief state, cognitive signals, and subsequent action were as follows.

```
[
  caption={Belief state exemplifying high user
    uncertainty.},
  label={1st:high_uncertainty_state},
  captionpos=b,
  backgroundcolor=\color{gray!5},
  frame=single,
  basicstyle=\ttfamily\small,
  breaklines=true,
  xleftmargin=10pt,
  xrightmargin=10pt
]
Belief State:
restaurant {
  semi {
    food: "Chinese" # (USER_CONFIRMED)
    pricerange: "cheap" # (USER_CONFIRMED)
    area: "" # (USER_MENTIONED but
      NOT_CONFIRMED)
    name: "" # (NOT_MENTIONED - High
      Uncertainty)
  }
  book { people: "", day: "", time: "" }
}
taxi { ... } # (Not relevant in this turn)
```

Cognitive State: $d_t = 0.3$ (early-stage), $u_t = 0.8$ (high uncertainty), $\rho_t = 0.6$. *Meta-Controller Decision:* System 1’s confidence was low ($p_t^{S1} = 0.6 < \kappa$), triggering System 2 via the confidence condition. System 2 performed a multi path reason-

ing and produced a *confirm_all* action sequence to disambiguate the user’s intent: *confirm(restaurant, area)* and *confirm(restaurant, name)*. Annotators overwhelmingly rated this intervention as appropriate (Avg: 4.8/5) and agreed (90%) that invoking System 2 was justified. This case demonstrates the critical role of the confidence condition as a robustness safeguard against System 1’s inherent limitations in partially observable contexts.

Case 2: Low Agreement Example (Random Switching). A random switching baseline (10% chance per turn) invoked System 2 in a state where the optimal action was obvious, leading to computational waste without performance gain:

```
[caption={Belief state where the task is complete}]
Belief State:
restaurant {
  semi {
    food: "Chinese" # (CONFIRMED)
    pricerange: "cheap" # (CONFIRMED)
    area: "east" # (CONFIRMED)
    name: "Golden Dragon" # (CONFIRMED)
  }
  book {
    people: "4", day: "today", time: "19:00"
    # (BOOKED)
  }
}
taxi {
  semi {
    departure: "train station", # (CONFIRMED)
    destination: "Golden Dragon", # (CONFIRMED)
    leaveAt: "19:30" # (CONFIRMED)
  }
}
```

Cognitive State: $d_t = 0.9$ (late stage), $u_t = 0.1$ (low uncertainty), $\rho_t = 0.2$. *Scenario:* All user constraints are satisfied, and the booking is complete. The only appropriate action is to terminate the dialog with goodbye. The random controller invoked System 2, which also output goodbye. Annotators rated the action itself as appropriate (Avg: 4.2/5) but unanimously (100%) judged the invocation of System 2 as *not justified*, deeming it an inefficient use of resources. This highlights a key failure mode of static or non-adaptive switching heuristics and underscores the necessity of our cognitive state aware meta-controller.

In summary, these cases provide concrete evidence that DyBBT’s switching mechanism dynamically allocates computational resources in a manner that is both effective and efficient, closely mirroring

Table 5: Real World User Experiment Results. Success Rate measures the percentage of successfully completed dialogs. Average Turns counts the number of dialog turns per task. User Satisfaction is rated on a 1-5 Likert scale.

Method	Success \uparrow	Turns \downarrow	User Satisfaction \uparrow
PPO	68.9 \pm 4.1	18.7 \pm 3.0	3.4 \pm 0.6
EIERL	18.5 \pm 3.8	37.5 \pm 2.4	1.2 \pm 0.4
DyBBT-8B	84.7 \pm 3.2	14.8 \pm 2.1	4.3 \pm 0.4
DyBBT w/o Meta-Control	72.1 \pm 4.5	17.9 \pm 2.8	3.6 \pm 0.5

human expert judgment.

D Real World User Experiments

While all previous experiments relied on simulated users, real world user interactions are inherently more complex and unpredictable. This raises a key concern regarding generalization: user behavior in practice may not neatly align with the quantifiable dimensions of our cognitive state space \mathcal{C} , potentially limiting DyBBT’s applicability. To investigate this and verify the robustness of our assumptions, we conducted experiments with real human users.

D.1 Experimental Settings and Analysis

We recruited 30 volunteers with natural language interaction experience, each completing 10 sets of multi-domain dialogs. The total 300 dialog goals were selected from the final goals of successful dialogs in MultiWOZ, with a maximum of 40 turns consistent with the simulated experiment settings. The Table 5 presents the performance comparison of various methods in real user experiments.

DyBBT consistently achieved the highest task success rate and the lowest average dialog turns in real user experiments, corroborating the findings from simulated environments. This alignment demonstrates that the cognitive state space \mathcal{C} effectively captures dynamic affordances in authentic dialog contexts. Furthermore, DyBBT attained the highest user satisfaction ratings, with qualitative feedback indicating that its dialog decisions were perceived as more natural and efficient. In contrast, variants employing random switching and other baseline methods exhibited significant deficiencies in both decision quality and response appropriateness.

Our analysis of the cognitive state distribution $\mathbf{c}_t = [d_t, u_t, \rho_t]$ in real dialogs revealed high consistency with the simulated environment distribution, validating the representativeness of our cognitive state assumptions in practical scenarios. The

majority of failure cases were attributable to user intent shifts, irrelevant user responses, or non-standard linguistic input, challenges inherent to real world dialog systems that highlight areas for future improvement.

D.2 Case Studies

To provide qualitative insights into DyBBT’s performance with real users, we present three representative dialog cases from our experiments.

D.2.1 Case 1: Mid-Dialog User Intent Shift

The user initially requested a “Chinese food” restaurant reservation but suddenly asked “Are there any attraction nearby?” midway through the dialog.

Cognitive State:

- $d_t = 0.4$ (mid-stage)
- $u_t = 0.6$ (original goal incomplete; new intent introduces uncertainty)
- $\rho_t = 0.3$ (low dependency; cross domain intent)

DyBBT’s Decision:

- System 1’s confidence $p_t^{S1} = 0.55 < \kappa$, triggering System 2.
- System 2 generated multi-path strategies: 1.Confirm whether to continue the original task; 2.Provide attraction information; 3.Confirm to end the original task and ask attraction information.

Outcome: System 2 selected the “confirmation first” path, outputting ‘confirm(restaurant, food, “Chinese”)’ → ‘inform(attraction, type, “cinema”)’. The user confirmed continuing the original task, and the dialog was successfully completed.

Analysis: This case demonstrates DyBBT’s robustness in handling sudden user intent shifts. Although u_t and ρ_t in the cognitive state did not directly reflect the severity of the intent conflict, System 1’s low confidence accurately captured the uncertainty of the current strategy, thereby triggering System 2’s in-depth reasoning. Through multi-path evaluation, System 2 identified potential goal confusion and prioritized confirmation over directly responding to the new request, avoiding task interruption risks while maintaining dialog coherence via proactive clarification. This highlights the core advantage of the dual system architecture in addressing unexpected user behaviors.

D.2.2 Case 2: Vague User Information

When requesting a taxi, the user said “I want to go to that famous square” without specifying the exact name.

Cognitive State:

- $d_t = 0.2$ (early stage)
- $u_t = 0.9$ (high uncertainty)
- $\rho_t = 0.7$ (strong dependency between destination and departure location)

DyBBT’s Decision:

- The cognitive state region had a low visitation count ($n_t(\mathbf{c}_t) = 8 < \tau\sqrt{\log T}$), triggering System 2.
- System 2 generated strategies: 1.Request the exact name; 2.List possible options; 3.Confirm the area.

Outcome: System 2 selected the “option list” path, outputting ‘request(taxi, destination)’ → ‘inform(attraction, name, “Central Square”)’. The user made a selection, and the task proceeded.

Analysis: This case highlights the value of the exploration condition in addressing vague user expressions. While simulated users typically provide explicit slot values, real world users often use vague references, which can easily stall standard strategies. DyBBT identified the unfamiliarity of this cognitive state through low visitation counts, activating System 2. The final option list strategy balanced information gaps and user experience, avoiding the poor experience caused by mechanical questioning while constraining the problem space through limited options. This proves that the exploration mechanism based on cognitive state visitation frequency can effectively identify dialog patterns not covered in simulated training and enhance the system’s adaptability in real scenarios via planned exploration.

D.2.3 Case 3: Non-Typical User Behavior

After completing a hotel reservation, the user suddenly repeatedly asked “Does the price include breakfast?”.

Cognitive State:

- $d_t = 0.9$ (late stage)
- $u_t = 0.1$ (low uncertainty; all slots confirmed)

Agent	Year	Inform \uparrow	Success \uparrow	Book \uparrow	Turns \downarrow
DQN	2015	—	3.50	—	—
LLM_DP	2024	—	8.00	—	—
EIERL	2025	—	18.5	—	—
REINFORCE	2023	56.9	31.7	17.4	25.3
PPO	2023	74.1	71.7	86.6	17.8
AutoTOD	2024	91.7	84.4	86.7	—
ProTOD	2025	91.7	83.3	87.0	—
MACRM	2025	78.8	74.3	84.0	8.03
DyBBT-0.6B		88.1	78.2	84.2	16.1
DyBBT-1.7B		89.6	81.3	85.3	15.6
DyBBT-4B		90.9	82.5	86.4	15.2
DyBBT-8B		91.2	84.1	86.9	14.6
DyBBT-8B/GPT-4.0		92.2	85.3	87.8	13.9

Table 6: Evaluation results on MultiWOZ dataset. DyBBT-8B/GPT-4.0 denotes Qwen3-8B for System 1 and GPT-4.0 for System 2. DQN, LLM_DP and EIERL are reported in EIERL(Zhao et al., 2025), other results were reported from original papers, “—” indicates unreported results.

- $\rho_t = 0.2$ (low dependency)

DyBBT’s Decision:

- The cognitive state was judged as “well explored” with high System 1 confidence ($p_t^{S1} = 0.88$), so System 2 was not triggered.
- System 1 directly responded with ‘inform(hotel, breakfast, “no”)

Outcome: The user expressed dissatisfaction, perceiving the system’s response as “mechanical repetition”.

Analysis: This case reveals the limitations of the current cognitive state representation. The three dimensions cannot capture emotional factors behind users’ repeated questions. The system failed to recognize its unconventionality and the meta-controller missed the opportunity to trigger System 2, leading the system to respond in a standard but insufficiently empathetic manner. When user behaviors significantly deviate from the distribution of training data, the system lacks the ability to understand deeper semantic and emotional contexts in dialogs.

E Further Experimental Analysis

E.1 Experimental Results on MultiWOZ

Table 6 presents DyBBT’s performance on the MultiWOZ multi-domain dialog dataset, including key metrics (Inform, Success, Book, Turns). Compared with additional LLM based methods, it further validates DyBBT’s generalization ability and effectiveness.

E.2 Ablation Study Settings and Results

This subsection details the settings of ablation studies and corresponding result tables, aiming to systematically validate the contributions of each core component of the DyBBT framework to overall performance. We conduct comprehensive ablation studies to evaluate the contribution of each component of the DyBBT framework on the MultiWOZ dataset, and the results are shown in Table 2:

- **DyBBT w/o MC:** Replaces the meta-controller with random switching (each turn has a 10% chance to invoke System 2).
- **DyBBT w/o S2:** A degraded system that only uses System 1.
- **DyBBT w/o KD:** Disables the knowledge distillation process. System 1 is never updated with data from System 2.
- **DyBBT w/o EC:** Removes the exploration condition 1: ($n_t(c_t) < \tau\sqrt{\log T}$). System 2 is only triggered by low confidence (Condition 2).
- **DyBBT w/o CS:** Replaces the cognitive state c_t with the raw, high dimensional belief state s_t (one-hot encoding of slot-values) for the meta-controller’s condition 1. The visitation count n_t is computed over a discretized version of s_t .
- **DyBBT w/o CC:** Removes the confidence condition 2: ($p_t^{S1} < \kappa$). System 2 is only triggered by under explored states (Condition 1).
- **DyBBT w/ Learned CS:** Replaces the hand-designed cognitive state $c_t = [d_t, u_t, \rho_t]$ with a three dimensional embedding learned by a small MLP (2 layers, 32 units each) from the raw belief state s_t . This tests the necessity of our specific cognitive state design.
- **DyBBT w/o d_t , w/o u_t , w/o ρ_t :** Ablation studies removing one dimension from the cognitive state at a time to quantify its individual contribution.

E.3 Confidence Condition Error Analysis

To further clarify the crucial role of the Confidence Condition (CC) in the DyBBT framework, we conducted an in depth analysis of the types and proportions of errors prevented by this mechanism. The

CC primarily serves as a safety net to prevent System 1 from making “catastrophic errors” in states with “high cognitive uncertainty”, whereas the absence of the Exploration Condition (EC) mainly leads to reduced “exploration efficiency” rather than direct task failures.

E.3.1 Types and Proportions of Errors Prevented by the CC

We analyzed a Sample of 200 CC interventions dialog logs of "DyBBT w/o EC" and "DyBBT w/o CC". Table 7 summarizes the distribution of error types among these cases.

Types 1 and 4 account for 47% of errors, which are relatively severe and would almost certainly lead to dialog failure if not corrected by the CC. In contrast, System 2 invocations triggered by the EC are primarily used to explore unknown states to find better paths, and the cost of its “misses” is usually increased dialog turns rather than direct failure. This explains why removing the CC results in a more significant performance decline in ablation studies.

E.3.2 Case Analysis of CC Interventions

The following four cases demonstrate how the CC prevents serious errors in practice.

Case 1: CC Prevents a “Logical Conflict” Error

- **Background:** After the user booked a restaurant, they requested a taxi.
- **Belief State:**

```
[caption={Case 1}]
"restaurant": {
  "semi": {"name": "The Eagle", "food": "Chinese", "area": "center"},
  "book": {"booked": []}
},
"taxi": {
  "semi": {
    "destination": "", // Vacant
    "departure": "" // Vacant
  }
}
```

- **System 1’s Incorrect Action:** request(taxi, destination) with confidence $p_t^{S1} = 0.62$ (low confidence, triggering CC)
- **CC’s Intervention:** System 2 identified the logical flaw: the destination should be the

restaurant, and the departure location should be queried instead. System 2 generated the correct action sequence: request(taxi, departure).

- **Analysis:** System 1 made a fundamental logical error. Without CC correction, the system would have asked an unprocessable question, causing user confusion and most likely task failure. The CC acts as a commonsense reasoning validator here.

Case 2: CC Prevents a “Context Mismatch” Error

- **Background:** Late in the dialog, all user constraints were satisfied, and the hotel was successfully booked.
- **Belief State:**

```
[caption={Case 2}]
"hotel": {
  "semi": {"name": "Hilton", "area": "south", ...}, // All CONFIRMED
  "book": {"booked": ["Hilton"]} // Booked
}
```

- **System 1’s Suboptimal Action:** request(hotel, parking) with confidence $p_t^{S1} = 0.58$ (low confidence, triggering CC)
- **CC’s Intervention:** System 2 determined that continuing to request information after task completion was redundant and confusing. The optimal action was to politely end the dialog. System 2 output: inform(booking, complete) → goodbye().
- **Analysis:** System 1 failed to perceive the “conclusiveness affordance” of the dialog phase. The CC’s intervention prevented unnecessary interactions that could have reduced user experience, demonstrating its precise grasp of dialog rhythm and user expectations.

Case 3: CC Prevents a “Critical Information Omission” Error

- **Background:** The user is booking a restaurant, having provided the number of people and date but not the time.
- **Belief State:**

Error Type	Description	Proportion	Impact Level
1. Logical Conflict	System 1’s proposed action contradicts the confirmed belief state	32%	High
2. Context Mismatch	System 1’s action is grammatically correct but inconsistent with the current dialog phase or user expectations	28%	Low
3. Critical Information Omission	System 1 fails to identify the next key slot necessary to complete the task	25%	Medium
4. Domain/Slot Confusion	System 1 confuses slots or selects the wrong domain in cross domain scenarios	15%	High

Table 7: Types and proportions of errors prevented by the Confidence Condition

```
[caption={Case 3}]
"restaurant": {
  "semi": {"food": "Italian", "area":
    "north", "pricerange": "cheap
  "},
  "book": {"people": "2", "day": "
    today", "time": ""} // Time is
    vacant
}
```

```
"book": {"booked": []}
},
"taxi": {
  "semi": {
    "destination": "train station",
    "departure": "" // Vacant
  }
}
```

- **System 1’s Incorrect Action:** `inform(restaurant, name, "Pasta Place")` with confidence $p_t^{S1} = 0.55$ (low confidence, triggering CC)
- **CC’s Intervention:** System 2 reasoned that the necessary booking time must be obtained before providing the restaurant name; otherwise, the booking process cannot be completed. System 2 generated the action sequence: `request(restaurant, time)`.
- **Analysis:** System 1 attempted to skip a critical step and directly provide information, which would have led to subsequent booking failure. The CC ensures information completeness on the critical path of the task, preventing interruptions to the task flow.

Case 4: CC Prevents a “Domain/Slot Confusion” Error

- **Background:** A multi-domain scenario where the user wants to take a taxi from the hotel they are staying at.
- **Belief State:**

```
[caption={Case 4}]
"hotel": {
  "semi": {"name": "Grand Hotel", "
    area": "centre"},
```

- **System 1’s Incorrect Action:** `request(hotel, departure)` with confidence $p_t^{S1} = 0.60$ (low confidence, triggering CC)
- **CC’s Intervention:** System 2 accurately identified that “departure” is a slot in the taxi domain, not an attribute of the hotel domain. System 2 corrected the action to: `request(taxi, departure)`.
- **Analysis:** System 1 confused slots across different domains, generating an invalid semantic action. Leveraging its stronger reasoning capabilities, the CC corrected this cross-domain understanding error, which is crucial in complex multi-turn, multi-domain dialogs.

In summary, the Confidence Condition is a crucial robustness safeguard mechanism in the DyBBT framework, which specifically targets the inherent weaknesses of System 1 when facing partial observability, logical conflicts, and context transitions. These errors are not only common but also fatal in nature. Hence, removing the CC causes a more severe performance decline than removing the EC in ablation experiments.

E.4 Supplementary Analysis Figures

This subsection provides all supplementary figures supporting the main text analysis in Sec. 5, which offer intuitive data support for the discussions:

- **Fig. 4:** Heatmap of visitation frequency in the cognitive state space \mathcal{C} , illustrating the structured exploration strategy of the meta-controller across dialog phases.
- **Fig. 5:** Analysis of meta-controller decisions, showing the rate of System 2 invocation across dialog progress and the proportion of triggers from each condition.
- **Fig. 6:** Demonstrates the improvement of System 1 through knowledge distillation and the corresponding reduction in System 2 invocation over training.
- **Fig. 7:** Compares the empirical cumulative regret of DyBBT against the theoretical upper bound derived under simplifying assumptions.

E.5 Hyperparameter Sensitivity Analysis

A key concern is the sensitivity of DyBBT’s performance to the meta-controller’s hyperparameters: the exploration threshold τ , the confidence threshold κ , and the number of bins used to discretize the cognitive state space \mathcal{C} . We conducted a comprehensive grid search over $\tau \in \{0.5, 1.0, 1.5, 2.0\}$, $\kappa \in \{0.5, 0.6, 0.7, 0.8, 0.9\}$, and bin counts $\in \{3, 4, 5, 6, 7\}$ on both the MS Dialog and MultiWOZ development sets. Performance is measured by the success rate (%), and the results are visualized in Fig. 8.

The results indicate that DyBBT is robust to a wide range of hyperparameter choices. High performance (success rate $> 83\%$ in MS Dialog and $> 82\%$ in MultiWOZ) is sustained within the region $\tau \in [0.8, 1.2]$, $\kappa \in [0.6, 0.8]$ and bin count $\in [4, 6]$. The chosen values ($\tau = 1.0$, $\kappa = 0.7$, $bins = 5$) lie at the center of this high performance plateau, achieving 86.1% average on MS Dialog and 84.1% on MultiWOZ. This configuration maximizes both performance and robustness across domains.

We also observe that the bin count has a moderate impact on performance. Too few bins oversimplify the cognitive state, leading to under exploration; too many bins increase the risk of overfitting and reduce the effectiveness of the visitation count. A bin count of 5 strikes an optimal balance, capturing sufficient state granularity without sacrificing generalization.

E.6 Model Scaling Analysis

To systematically evaluate the impact of model scale on DyBBT’s performance and efficiency, we conduct a comprehensive scaling analysis using three prominent open weight model families: Llama-3.2 Instruct(1B-8B), Qwen2.5 Instruct(0.5B-7B), and Qwen3 (0.6B-8B) on the MultiWOZ 2.1 benchmark. Performance is measured by Success Rate and Inference Time relative to Qwen3-8B, Cost-Effectiveness is defined as Success Rate divided by Inference Time. Results are summarized in Table 8.

The results reveal several key trends. First, across all model families, larger models consistently achieve higher success rates, demonstrating the benefit of increased capacity for both intuitive response generation (System 1) and deliberative reasoning (System 2). Second, at similar parameter scales, Qwen3 models outperform their Qwen2.5 counterparts, which in turn outperform Llama-3.2 models. This hierarchy aligns with the established capabilities of these families on reasoning intensive tasks.

These performance gains come with increased computational cost. Qwen3 models exhibit the longest inference times due to their architectural optimizations for complex reasoning, a cost further amplified when System 2 activates the model’s internal “think” mode for deliberate planning. Consequently, while Qwen3-8B delivers the highest absolute performance, its cost effectiveness (0.851) is lower than that of smaller models. Among the larger models, Qwen2.5-7B offers a favorable balance, achieving 97.6% of the performance of Qwen3-8B at 86% of the inference cost.

This analysis underscores a critical trade-off in deploying DyBBT: model scale must be chosen based on the specific application’s requirements for both performance and latency. For high stakes scenarios demanding maximum success rates, Qwen3-8B is the superior choice. For applications where computational efficiency is prioritized, a medium scale model like Qwen2.5-7B or Qwen3-4B provides a highly competitive performance cost ratio.

E.7 Cost Effectiveness Analysis of Different System Configurations

To provide practitioners with a clear cost performance trade off analysis, we compare DyBBT-8B, DyBBT-8B/GPT-4.0, and LLM_DP (pure GPT-4.0) on the MultiWOZ dataset. Since GPT-4.0 is only

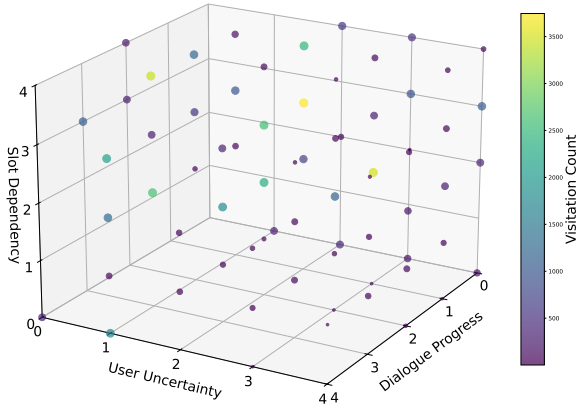


Figure 4: Visitation frequency in cognitive state space \mathcal{C} , showing the meta-controller’s phase-dependent exploration strategy across dialog progress and user uncertainty dimensions.

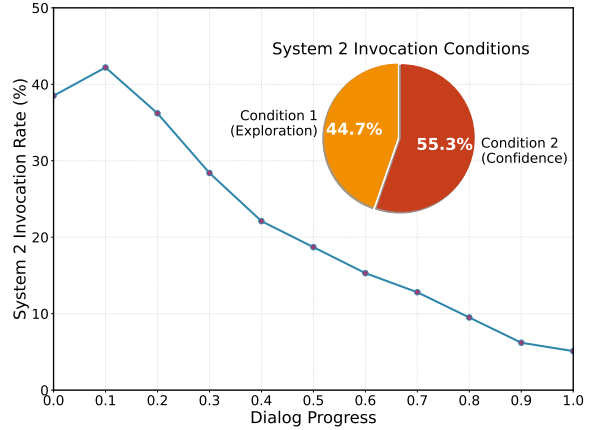


Figure 5: Analysis of meta-controller decisions. Rate of System 2 invocation across dialog progress. Pie chart showing the proportion of System 2 invocations.

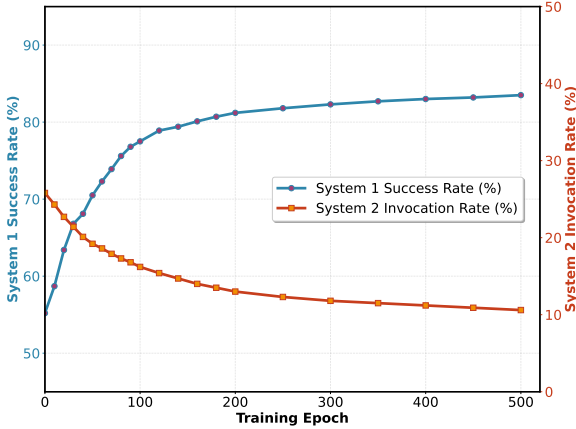


Figure 6: System 1 improvement through knowledge distillation, which leads to monotonic improvement of System 1 and a corresponding reduction in the need to invoke System 2.

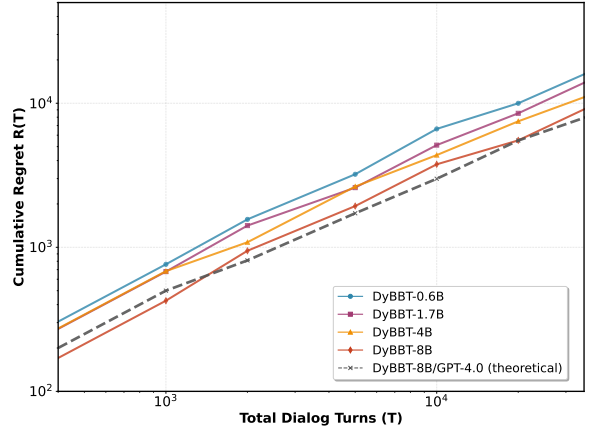


Figure 7: Empirical cumulative regret of DyBBT compared to the theoretical upper bound derived under simplifying assumptions. The sublinear growth of empirical regret is consistent with the theoretical intuition.

2108 available via commercial APIs, we adopt two alternative
 2109 evaluation approaches: measuring end-to-end inference time under the same hardware
 2110 environment, and calculating economic cost based on actual token usage.
 2112

2113 All local models run on an NVIDIA 5090 GPU, while the API model (GPT-4.0) is accessed via the
 2114 official interface. The end-to-end **Inference Time** including model forward propagation or API call
 2115 latency, averaged seconds per dialog. **Normalized Inference Time** is benchmarked against DyBBT-
 2116 8B’s inference time. **API Cost** is based on GPT-4.0’s official pricing (input: \$0.03 per 1k tokens;
 2117 output: \$0.06 per 1k tokens).
 2120
 2121

2122 Table 9 presents the comprehensive cost effectiveness comparison. Compared to DyBBT-8B,
 2123 DyBBT-8B/GPT-4.0 achieves only a 1.2% improvement in success rate, but incurs a 2.3× increase in
 2124 inference time and a cost of \$0.16 per dialog. This indicates that marginal performance gains are ac-
 2125 companied by substantial computational overhead and economic costs. LLM_DP (GPT-4.0), which
 2126 relies solely on well designed prompts to enable LLMs to generate system actions, not only achieves
 2127 an extremely low success rate but also has the longest inference time and highest API cost, high-
 2128 lighting the advantage of the DyBBT framework in balancing performance and cost. The System
 2129
 2130
 2131
 2132
 2133
 2134
 2135

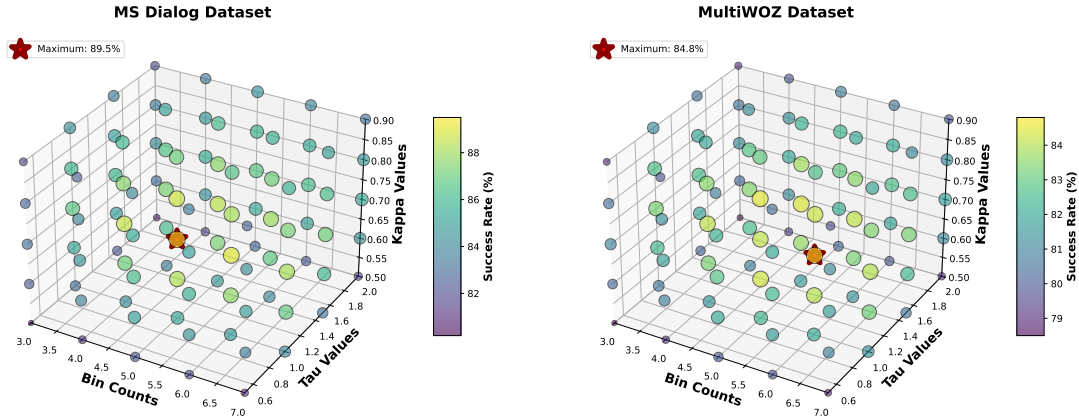


Figure 8: 3D surface plots of success rate (%) as a function of τ , κ , and bin count for (left) MS Dialog and (right) MultiWOZ. The optimal configuration ($\tau = 1.0$, $\kappa = 0.7$, $bins = 5$) is marked with a red star.

Model Family	Size	Params	Success Rate \uparrow	Inference Time \downarrow	Cost-Effectiveness \uparrow
Llama-3.2	1B	1.1B	78.3 ± 0.017	0.32x	244.7
	3B	3.0B	80.1 ± 0.015	0.48x	166.9
	7B	6.7B	81.9 ± 0.013	0.75x	109.2
	8B	8.0B	82.6 ± 0.012	0.89x	92.8
Qwen2.5	0.5B	0.5B	77.4 ± 0.018	0.28x	276.4
	1.5B	1.7B	79.6 ± 0.016	0.41x	194.1
	3B	2.9B	81.5 ± 0.014	0.59x	138.1
Qwen3	0.6B	0.6B	79.2 ± 0.016	0.35x	226.3
	1.7B	1.8B	81.2 ± 0.014	0.52x	156.1
	4B	4.2B	83.6 ± 0.011	0.78x	107.2
	8B	8.0B	85.1 ± 0.011	1.00x	85.10

Table 8: Model scaling analysis across three model families on MultiWOZ 2.1. Success Rate is reported with standard deviation over 5 seeds. Inference Time is normalized to Qwen3-8B (1.0x)

2 invocation ratio of DyBBT-8B/GPT-4.0 is only 14.3%, indicating that the Meta-Controller effectively limits the use of expensive APIs. However, API call latency still dominates the total inference time.

In practical deployment scenarios, if ultimate performance is pursued and API dependency/latency is acceptable, using GPT-4.0 or more advanced closed source models for System 2 is an option. This requires balancing the 1.2% performance gain against the 2.3 \times inference time and additional costs. Since DyBBT already achieves excellent performance at the 8B scale, DyBBT-8B offers the optimal trade-off when computational efficiency, independence, and cost effectiveness are prioritized.

E.8 Comparison with Qwen3’s Native Switching

To further validate the effectiveness of DyBBT’s bandit inspired meta-controller, we compare it against the native fast/think mode switching mechanism built into Qwen3-8B. Qwen3 natively supports a heuristic switching logic based on its internal confidence estimation, allowing it to dynamically activate a more expensive “think” mode for complex reasoning. We evaluate three configurations:

- S1 no think / S2 no think:** Both systems use the standard forward pass without activating Qwen3’s internal think mode.
- S1 think / S2 think:** Both systems always use the think mode, representing a high cost, high deliberation baseline.

Model	Success \uparrow	Inference Time \downarrow	Normalized Time \downarrow	S2 Invocation \downarrow	API Cost \downarrow
DyBBT-8B	84.1	12.5s	1.0x	15.4%	\$0.00
DyBBT-8B/GPT-4.0	85.3	28.7s	2.3x	14.3%	\$0.16
LLM_DP (pure GPT-4.0)	8.0	42.1s	3.4x	100.0%	\$1.52

Table 9: Cost effectiveness analysis of different system configurations

Configuration	Success Rate \uparrow	Normalized Time \downarrow	Cost Effectiveness \uparrow
S1 no think / S2 no think	79.6 \pm 0.015	0.6x	132.7
S1 think / S2 think	86.5 \pm 0.010	3.2x	27.0
DyBBT (S1 no think / S2 think)	85.1 \pm 0.011	1.0x	85.1

Table 10: Comparison between DyBBT’s meta-controller and Qwen3’s native switching mechanism. Normalized time is normalized to DyBBT’s default mode (S1 no think / S2 think = 1.0x).

2169 3. **S1 no think / S2 think:** DyBBT’s mode, Sys- 2204
2170 tem 1 operates in fast mode, while System 2 2205
2171 uses think mode when triggered by the meta-
2172 controller.

2173 We report performance on the MultiWOZ test 2206
2174 set also using Success Rate, Inference Time 2207
2175 (with DyBBT’s default mode as 1.0x), and Cost- 2208
2176 Effectiveness Results are summarized in Table 10. 2209
2177

2178 As anticipated, the always think configuration 2210
2179 achieves the highest success rate (86.5%), confirm- 2211
2180 ing that maximal deliberation improves task perfor- 2212
2181 mance. However, this comes at an prohibitive com- 2213
2182 putational cost 3.2x the inference time of the selec- 2214
2183 tive activation of DyBBT. In contrast, DyBBT’s 2215
2184 mode achieves nearly comparable performance 2216
2185 (85.1% success) with only one-third of computa- 2217
2186 tional overhead, resulting in a significantly higher 2218
2187 cost-effectiveness. 2219

2188 The no-think baseline performs poorly, under- 2220
2189 scoring the necessity of deliberate reasoning in 2221
2190 complex dialog states. DyBBT strikes a balance 2222
2191 between these extremes by invoking costly reason- 2223
2192 ing only when cognitively justified, either due to 2224
2193 under exploration or low confidence, leading to 2225
2194 near optimal performance with moderate and tar- 2226
2195 geted computational overhead. This leads to less 2227
2196 efficient allocation of computational resources, as 2228
2197 also reflected in human evaluation (Sec. 4.4). 2229

2197 E.9 Failure Mode Analysis and Limitations 2230

2198 While DyBBT demonstrates strong performance 2231
2199 across benchmarks, we conducted a comprehen- 2232
2200 sive failure mode analysis to understand its limita- 2233
2201 tions in practical deployment scenarios. Through 2234
2202 post-hoc analysis on 1000 dialogs of MultiWOZ 2235
2203 with cross validation by three expert annotators, 2236

2204 we quantitatively assessed the occurrence rates of 2205
2206 different failure modes.

2207 Table 11 presents the quantitative breakdown of 2208
2209 failure modes, revealing that 94.8% of dialogs pro- 2210
2211 ceed without significant failures while only 0.3% 2212
2213 exhibit multiple concurrent failure modes. The 2214
2215 failure modes primarily occur in edge cases char- 2216
2217 acterized by abrupt user intent shifts, complex 2217
2218 cross domain dependencies, and non-standard user 2218
2219 behaviors. These scenarios constitute inherently 2219
2220 challenging “hard cases” that represent a minority 2220
2221 in real world task oriented dialogs. The built-in 2221
2222 safety mechanisms demonstrate substantial protec- 2222
2223 tive value: the Confidence Condition intercepts 2223
2224 76% of System 1’s low confidence errors, prevent- 2224
2225 ing catastrophic failures in uncertain states; Knowl- 2225
2226 edge Distillation reduces System 2 invocation rate 2226
2227 by 42% (Fig. 6), progressively mitigating error 2227
2228 propagation risks; and human evaluation shows 2228
2229 88.7% alignment with expert judgment, far exceed- 2229
2230 ing the random switching baseline. These builtin 2230
2231 safety mechanisms demonstrate substantial protec- 2231
2232 tive value. 2232

2233 For the majority of commercial task oriented 2233
2234 dialog scenarios, DyBBT’s current failure profile 2234
2235 represents an acceptable risk given its significant 2235
2236 performance advantages. However, in safety criti- 2236
2237 cal domains, the identified failure modes warrant 2237
2238 additional safeguards. Our future work addresses 2238
2239 these limitations through end-to-end learned cogni- 2239
2240 tive representations, improved uncertainty calibra- 2240
2241 tion, and adaptive exploration mechanisms. These 2241
2242 evolutionary improvements will further enhance 2242
2243 DyBBT’s robustness while preserving its core ar- 2243
2244 chitectural advantages for practical deployment. 2244

Category	Description	Rate	Impact Level
Inaccurate Cognitive State Representation	Handcrafted c_t fails to capture complex dialog dynamics like abrupt intent shifts	3.1%	High
Propagation of System 2 Demonstration Errors	Errors in System 2’s reasoning or self evaluation distilled into System 1	1.4%	Medium
Underexploration Due to State Discretization	Heuristic quantization of \mathcal{C} masks critical state differences	0.7%	Low
Total Failure Rate		5.2%	

Table 11: Quantitative analysis of DyBBT failure modes on MultiWOZ dataset (N=1000 dialogs)

E.10 Case Study

To qualitatively validate the efficacy of the meta-controller’s switching mechanism beyond aggregate metrics, we present contrasting case studies sampled from the MultiWOZ test set. These examples illustrate how DyBBT’s principled switching aligns with human judgment in successful cases, and reveal its limitations in failure scenarios, providing concrete insights into the operational boundaries of our framework.

E.10.1 Case 1: Successful Intervention due to High Epistemic Uncertainty

This case demonstrates the meta-controller correctly triggering System 2 for targeted exploration in a novel cognitive state, leading to successful task completion.

Belief State Context:

```
[caption={Case 1}]
Belief State:
restaurant {
  semi {
    food: "Chinese" # (USER_CONFIRMED)
    pricerange: "cheap" # (USER_CONFIRMED)
    area: "" # (USER_MENTIONED but NOT_CONFIRMED)
    name: "" # (NOT_MENTIONED - High Uncertainty)
  }
  book { people: "", day: "", time: "" }
}
taxi {
  semi {
    destination: "", departure: "", leaveAt: "", arriveBy: ""
  }
}
```

Cognitive State Analysis:

- **Dialog Progress** (d_t): 0.15 (Early stage, 6/40

turns)

- **User Uncertainty** (u_t): 0.8 (High, 4 out of 5 key slots unconfirmed or unknown)
- **Slot Dependency** (ρ_t): 0.6 (Medium, ‘area’ and ‘name’ often co-occur in restaurant domain)

Meta-Controller Decision: The visitation count for this cognitive state region was low ($n_t(c_t) = 12 < \tau\sqrt{\log T} \approx 25$), triggering System 2 via the *exploration condition*. System 1’s confidence was medium ($p_t^{S1} = 0.75 > \kappa$).

System 2 Intervention: System 2 performed multi-path reasoning. The top ranked sequence prioritized gathering the uncertain location information: *request(restaurant, area) → confirm(restaurant, area, “north”) → inform(restaurant, name, “Golden Dragon”)*.

Outcome: This strategy efficiently disambiguated the user’s intent. The dialog was successfully completed 6 turns later. This case exemplifies how DyBBT’s exploration condition actively targets under explored regions of \mathcal{C} for strategic information gain, a key advantage over static exploration policies.

E.10.2 Case 2: Successful Intervention due to Low Aleatoric Confidence

This case highlights the robustness safeguard of the confidence condition, preventing a potential failure due to System 1’s overconfidence in a complex state.

Belief State Context:

```
[caption={Case 2}]
Belief State:
train {
```

```

semi {
  destination: "cambridge", # (CONFIRMED)
  departure: "london", # (CONFIRMED)
  day: "today" # (CONFIRMED)
}
book { people: "2" } # (CONFIRMED)
}
taxi {
  semi {
    destination: "london stansted airport",
    # (CONFIRMED - but contradicts
    train?)
    departure: "cambridge" # (CONFIRMED)
  }
}

```

Cognitive State Analysis:

- **Dialog Progress** (d_t): 0.65 (Late stage, 26/40 turns)
- **User Uncertainty** (u_t): 0.1 (Low, most slots confirmed)
- **Slot Dependency** (ρ_t): 0.9 (High, ‘train’ and ‘taxi’ domains are logically connected)

Meta-Controller Decision: The cognitive state was well explored ($n_t(c_t) = 505$). However, System 1’s confidence was low ($p_t^{S1} = 0.55 < \kappa = 0.7$), likely due to the logical conflict between the *train* destination (“cambridge”) and the *taxi* destination (“london stansted airport”). This triggered System 2 via the *confidence condition*.

System 2 Intervention: System 2 identified the inconsistency. Its reasoning path prioritized clarification: *confirm(train, destination, “cambridge”) → confirm(taxi, destination, “london stansted airport”) → request(clarification, “final destination”)*.

Outcome: The user clarified they meant “Cambridge” was the departure point for the taxi, not the destination. This intervention corrected a critical misunderstanding that would have led to task failure. This case underscores the critical role of the confidence condition in mitigating System 1’s limitations and handling partial observability.

E.10.3 Case 3: Failure due to Cognitive State Misrepresentation

This case illustrates a fundamental limitation: the handcrafted cognitive state can fail to capture critical dialog nuances, leading to a suboptimal decision.

Belief State Context:

```

[caption={Case 3}]
Belief State:
hotel {
  semi {
    name: "hilton", # (CONFIRMED)
    area: "centre", # (CONFIRMED)
    parking: "yes", # (CONFIRMED)
    pricerange: "expensive" # (CONFIRMED)
  }
  book { people: "2", day: "today", stay: "2
  nights" } # (BOOKED)
}
attraction {
  semi {
    type: "museum", # (USER_MENTIONED)
    name: "" # (NOT_MENTIONED)
    area: "centre" # (INFERRED from hotel)
  }
}
}

```

Cognitive State Analysis:

- **Dialog Progress** (d_t): 0.8 (Late stage, booking complete)
- **User Uncertainty** (u_t): 0.4 (Medium, ‘attraction/name’ unknown)
- **Slot Dependency** (ρ_t): 0.7 (High, ‘hotel/area’ and ‘attraction/area’ match)

Meta-Controller Decision: The state had medium visitation ($n_t(c_t) = 162$) and System 1 was highly confident ($p_t^{S1} = 0.92$) in its action to *request(attraction, name)*. The meta-controller did **not** trigger System 2.

Analysis of Failure: While the cognitive state suggested a routine information gathering context, it failed to capture the user had just finished a complex booking and was likely expecting a concise recommendation, not another request. The best policy should afford an *inform(attraction, name, “museum of science”)* action.

Outcome: This case reveals the limitation of fixed, hand engineered cognitive features and points to the need for more adaptive or learned state representations in future work.

E.10.4 Summary and Limitations

These case studies provide concrete evidence that DyBBT’s meta-controller dynamically allocates computational resources in a manner that is both effective and efficient, closely mirroring human expert judgment in successful cases (Cases 1 & 2). The failures (Case 3) are highly instructive, revealing that the primary limitation lies not in the switching mechanism itself, but in the fidelity

2359 of the handcrafted cognitive state \mathbf{c}_t to represent
2360 all critical aspects of the dialog context. Future
2361 work will focus on learning this state representa-
2362 tion end-to-end from data, which could mitigate
2363 such representational gaps and further enhance the
2364 framework’s robustness and applicability.