# DYBBT: DYNAMIC BALANCE VIA BANDIT INSPIRED TARGETING FOR DIALOG POLICY WITH COGNITIVE DUAL-SYSTEMS

#### **Anonymous authors**

Paper under double-blind review

#### 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 state-of-theart 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 https://anonymous.4open.science/r/DyBBT-C6B7.

#### 1 Introduction

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

— James J. Gibson, The Ecological Approach to Visual Perception (1979)

Task oriented dialog system (TODS) assist users in achieving specific goals, like booking flights or reserving restaurants, via multi-turn natural language interactions. Dialog policy typically formulated as a sequential decision making problem addressed with Deep Reinforcement Learning (DRL) Nachum et al. (2017); Silver et al. (2014), is bottlenecked by the exploration-exploitation dilemma: balancing exploitation of known rewards against exploration of unknown actions to discover better strategies. Unlike in standard RL, this dilemma in TODS is fundamentally exacerbated by its intrinsic cognitive structure, a dynamic partially observable context characterized by quantifiable features such as the progress ratio of filled goal slots, the entropy of user intent over possible values, and the conditional dependency between unfilled slots based on domain ontology Peng et al. (2017); Wen et al. (2017). These features directly govern the cost benefit analysis of exploration: early in a dialog, high entropy makes information gathering actions valuable; late in a dialog, high slot dependency makes exploitation critical to avoid constraint violations Qin et al. (2023); Zhao et al. (2024).

Exploration in TODS is fundamentally challenging due to its dynamic, partially observable nature Lee et al. (2023), characterized by three key cognitive properties that unfold in distinct dialog phases. Early dialog stages afford information gathering, as user goals are often ambiguous and multiple slots remain unfilled Kwan et al. (2023); Mid-stages afford clarification and confirmation as slots begin to fill and dependencies emerge Jia et al. (2024); and late stages afford task completion, where actions must adhere to strict slot-value dependencies, for example, a taxi cannot be booked without both "departure" and "destination" Niu et al. (2024). This dynamic "affordance landscape" demands adaptive exploration: static strategies cause inefficiencies, premature exploitation fails tasks, while aimless exploration wastes turns.

Current methods for enhancing exploration in TODS, while powerful, are fundamentally misaligned with this dynamic cognitive reality. As illustrated in Figure 1, traditional DRL methods rely on static heuristics such as  $\epsilon$ -greedy Niu et al. (2024), which cannot adapt to shifting exploration needs between dialog phases. Evolutionary methods like EIERL Zhao et al. (2025) enable global search

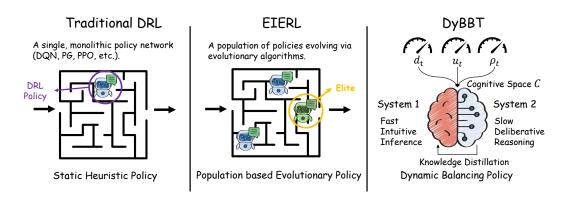


Figure 1: Traditional DRL methods (left) employ a static exploration strategy with a single policy. EIERL (middle) uses population based evolutionary optimization with elite injection but struggles to scale to complex multi domain tasks. DyBBT (right) introduces a cognitive meta-controller dynamically balances fast intuitive responses and slow deliberative reasoning for adaptive policy selection.

via population based optimization and elite injection to accelerate evolution, yet struggle in complex multi-domain scenarios due to poor scalability and unflexible updates. LLM based policies Zhang et al. (2024); He et al. (2022) or reasoning techniques such as Tree of Thoughts (ToT) 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 to design a dialog policy that dynamically perceives cognitive affordances to balance exploration and exploitation?* 

To solve the above challenges, we propose DyBBT, a novel framework that grounds decisions in an interpretable cognitive state space  $\mathcal{C}$  that captures dialog progress  $d_t$ , user uncertainty  $u_t$ , and slot dependency  $p_t$ , as shown in Figure 1. DyBBT introduces 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 and visitation counts. This design ensures that expensive reasoning is invoked only when the cognitive state signals under exploration or high uncertainty, addressing the core limitations (RQ) of previous methods. By formalizing dialog affordances and embedding them into a bandit inspired switching mechanism, DyBBT achieves a principled and efficient balance between exploration and exploitation.

In summary, our work makes the following contributions: (1) Formalization of TODS exploration challenge via a structured cognitive state space  $\mathcal{C}$  (Section 3.1). (2) Proposal of DyBBT, a novel framework with bandit inspired meta-controller to dynamically balance between fast System 1 and deliberate System 2 reasoning (Section 3.2). (3) Demonstration of state-of-the-art (SOTA) performance and human aligned decisions through extensive experiments (Section 4).

#### 2 RELATED WORK

#### 2.1 DIALOG POLICY LEARNING WITH DEEP REINFORCEMENT LEARNING

Deep Reinforcement Learning (DRL) has become a dominant paradigm for dialog policy optimization due to its capacity for sequential decision making. Early work applied value based methods Peng et al. (2018) and Policy Gradient Silver et al. (2014) to TODS, Proximal Policy Optimization (PPO) Schulman et al. (2017) was later adopted for improved stability and has become a common baseline. A key limitation of these methods is their reliance on static exploration strategies, such as  $\epsilon$ -greedy or entropy bonus. These heuristics cannot adapt to the dynamic uncertainty and structural complexity of multi-domain dialogs Kwan et al. (2023); Jia et al. (2024). Recent efforts have incorporated Bayesian reasoning Lee et al. (2023), meta-learning Li et al. (2024); Liang et al. (2025), Cascading RL Du et al. (2024), CB-RL Thoma et al. (2025) to allow more adaptive exploration. While promising, they often lack an explicit and interpretable representation of the internal cognitive dialog state that directly governs exploration, a gap our cognitive state space  $\mathcal C$  aims to fill.

#### 2.2 EVOLUTIONARY AND POPULATION BASED METHODS FOR EXPLORATION

Evolutionary Reinforcement Learning (ERL) combines population based global search with gradient-based optimization to enhance exploration diversity. Methods such as EIERL Zhao et al. (2025) inject elite policies to accelerate evolution, enabling escape from local optima. However, ERL scales poorly with dialog complexity due to exponential growth in population size Sigaud (2023). Moreover, these methods often rely on fixed schedules for policy replacement, lacking dynamic adaptation to real-time dialog progression and cognitive state changes Bai et al. (2023). In contrast, DyBBT replaces expensive population evolution with a single, efficient dual-system architecture guided by a structured cognitive state space, enabling fine grained, context aware exploration without the scalability limitations of population based approaches.

#### 2.3 CLASSICAL AND MODERN EXPLORATION THEORIES

The exploration-exploitation trade-off is a cornerstone of sequential decision theory. Bandit algorithms, such as Upper Confidence Bound (UCB) Garivier & Moulines (2011), provide theoretical guarantees for stationary settings, and their principles have been extended to contextual bandits Foster & Rakhlin (2020) and hierarchical RL Rohmatillah & Chien (2023). However, directly applying these theories to dialog Partially Observable Markov Decision Processes (POMDPs) faces significant challenges due to non-stationarity, partial observability, and the high dimensional nature of the state space. Our work draws inspiration from the optimism principle of UCB but makes a pragmatic *heuristic adaptation* to a learned cognitive state space  $\mathcal{C}$ . This approach preserves the interpretability and theoretical intuition of bandit algorithms while specifically addressing the complexities of sequential dialog environments. Compared to methods like PSRL Chen et al. (2020) that require maintaining a posterior over the entire MDP, our method focuses exploration on a compact cognitive space, offering a computationally efficient alternative better suited to dialog POMDPs.

#### 2.4 Dual-System Architectures and LLMs for dialog Policies

Krämer (2014), combine fast, intuitive processing (System 1) with slow, deliberative reasoning (System 2), have been applied to mathematical reasoning Shi et al. (2024) and common sense inference Yu et al. (2025). In dialog systems, large language models (LLMs) serve as powerful function approximators Yi et al. (2024), acting as intuitive generators Ying et al. (2024) and deliberative reasoners Ma et al. (2025). Recent work, such as the Dynamic Dual-Process Transformer He et al. (2024), explicitly models the interaction for dialog 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 guided by a bandit inspired principle, dynamically triggering System 2 based on cognitive state visitation counts and parametric uncertainty, offering a principled and efficient alternative to heuristic switching.

#### 3 METHODOLOGY

To address the limitations of *static exploration*, *high computational cost*, and *lack of cognitive grounding*, we propose DyBBT: a framework that grounds adaptive exploration in a structured *cognitive state space*  $\mathcal{C}$ , orchestrated by a lightweight *meta-controller* over a dual-system architecture. Its core innovation is using  $\mathcal{C}$  as the *decision basis*: the meta-controller dynamically triggers costly System 2 based on real-time cognitive signals, enabling efficient, context aware exploration.

#### 3.1 Theoretical Foundation

#### 3.1.1 COGNITIVE STATE SPACE: A BRIDGE FOR EXPLORATION

We define a low dimensional interpretable representation of the dialog context *cognitive state*  $\mathbf{c}_t = [d_t, u_t, \rho_t]$ , to bridge bandit inspired exploration principles with the sequential POMDPs setting. Based on Gibson's affordance theory,  $\mathbf{c}_t$  captures action possibilities that directly govern the exploration-exploitation trade-off, where  $d_t$ ,  $u_t$ , and  $\rho_t$  quantify dialog progress, user uncertainty, and slot dependency, respectively (see Appendix A.1 for computation).

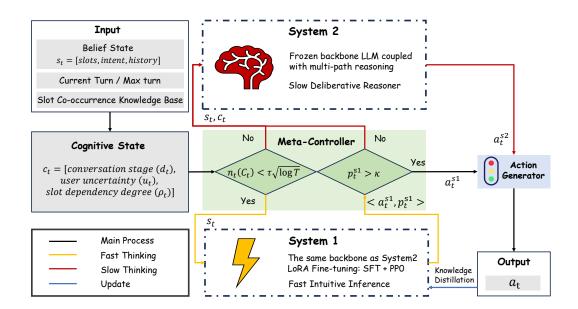


Figure 2: The DyBBT Architecture. A meta-controller uses the cognitive state  $\mathbf{c}_t$ , visitation count  $n_t(\mathbf{c}_t)$ , and System 1's confidence  $p_t^{S1}$  to dynamically select between System 1 (fast, intuitive) and System 2 (slow, deliberative). Outputs drive action execution and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

#### 3.1.2 REWARD SMOOTHNESS: A PRAGMATIC ASSUMPTION FOR STRUCTURED TASKS

To enable theoretical analysis, we assume that the reward function is Lipschitz continuous (Asadi et al., 2018) in the cognitive state space C. We provide empirical validation in Section 5.4.

**Assumption 1** (Lipschitz Smooth Reward in 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:

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

#### 3.1.3 DYNAMIC BALANCE PRINCIPLE: A BANDIT INSPIRED HEURISTIC

To address the intractability of optimal exploration in dialog POMDPs, we propose a bandit (Komiyama et al., 2024) inspired heuristic: the exploration bonus for cognitive state  $\mathbf{c}_t$  is inversely proportional to the square root of its visitation count:

Exploration-Bonus
$$(t) \propto \sqrt{\frac{\log T}{n_t(\mathbf{c}_t)}}$$
. (1)

This yields sublinear regret:  $\mathbb{E}[R(T)] \lesssim \widetilde{\mathcal{O}}\left(L_r \cdot \sqrt{\dim(\mathcal{C}) \cdot T}\right)$ , where  $L_r$  is the Lipschitz constant. See Appendix A.2 for the derivation and assumptions.

#### 3.2 System Architecture

The proposed DyBBT framework operationalizes its theoretical principles into a dynamic dual-system architecture (Fig. 2), where the meta-controller leverages the cognitive state  $\mathbf{c}_t$  to orchestrate a cost aware trade-off: invoking the fast, learned System 1 for routine decisions, or the powerful System 2 for high uncertainty or under explored states. This design ensures that expensive deliberation is reserved for contexts where the bandit inspired criterion (Eq. 2) signals its necessity, either due to low visitation  $(n_t(\mathbf{c}_t))$  or low confidence  $(p_t^{S1})$ . The result is a closed loop system that translates cognitive perception into adaptive action selection, without redundant or static heuristics.

#### 3.2.1 System 1 (S1): The Fast Intuitive Inference

To provide a low latency, high throughput baseline policy for the majority of dialog turns, mitigating the prohibitive cost of always using a deliberative reasoner, S1 embodies the fast and intuitive system. It is implemented as a base pre-trained Transformer model augmented with LoRA modules Hu et al. (2022). It is first trained via Supervised Fine-Tuning (SFT) on expert trajectories and subsequently optimized using PPO to maximize the expected task reward. At inference time, it takes the current belief state  $\mathbf{s}_t$ , formatted via a concise prompt, and generates an action  $a_t^{S1}$  along with a calibrated confidence score  $p_t^{S1} \in [0,1]$ , which reflects its intrinsic uncertainty.

#### 3.2.2 System 2 (S2): The Slow Deliberative Reasoner

To handle novel or complex situations where fast policy (S1) is likely to fail, thus addressing the suboptimal performance of static DRL policies in under explored regions, S2 represents the slow and analytical system. It utilizes the same base model as S1, but remains frozen to preserve its broad knowledge and reasoning capabilities. The prompt instructs S2 to generate Top-3 distinct action sequences. Each sequence's quality is evaluated using the ratio of filled key slots. We extract the first action from the highest quality sequence as output  $a_t^{S2}$ . This system is computationally expensive but is designed to handle novel or high stakes situations identified by the meta-controller.

#### 3.2.3 META-CONTROLLER: DYNAMIC ORCHESTRATION VIA COGNITIVE AFFORDANCES

The meta-controller bridges the gap between adaptive exploration and computational efficiency by dynamically invoking S2 only when the cognitive state  $\mathbf{c}_t$  signals a clear affordance for deliberation. It operates on  $\mathbf{c}_t = [d_t, u_t, \rho_t]$ , its derived visitation count  $n_t(\mathbf{c}_t)$ , and S1's confidence  $p_t^{S1}$ , ensuring that decisions are grounded in real time cognitive perception. This hybrid rule combines two practical triggers: (1) exploration in under visit cognitive states, motivated by bandit principles, and (2) uncertainty mitigation when System 1 lacks confidence. The disjunctive design ensures S2 activation for systematic exploration or robust decision making.

Activate S2 IF: 
$$\left(n_t(\mathbf{c}_t) < \tau \sqrt{\log T}\right) \vee \left(p_t^{S1} < \kappa\right)$$
. (2)

Condition 1:  $n_t(\mathbf{c}_t) < \tau \sqrt{\log T}$ , targets under explored regions of  $\mathcal{C}$ , triggering S2 to perform theoretically motivated exploration where reward potential is high (Assumption 1). This replaces static global heuristics with a local adaptive mechanism. Condition 2:  $p_t^{S1} < \kappa$ , acts as a robustness safeguard, activating S2 when S1 exhibits high parametric or aleatoric uncertainty, often due to partial observability or distributional shift. This prevents catastrophic failures in high stakes states, addressing the brittleness of single policy or fixed schedule methods.

This disjunctive design ensures S2 is invoked only when its cost is justified: to explore promising new territory (curiosity driven) or to mitigate imminent failure (uncertainty driven), achieving both efficiency and robustness. The selected actions are executed via a lightweight *action generator* ensuring semantic validity. To reduce long term S2 reliance, high quality S2 decisions are distilled into S1 (see Appendix B.5.2 & Algorithm 1).

#### 4 EXPERIMENT

#### 4.1 EXPERIMENTAL SETUP

**Datasets.** We conduct experiments on two of the most prominent TODS benchmark datasets which are also used in baselines. The Microsoft Dialog Challenge platform Li et al. (2018); Zhao et al. (2024); Niu et al. (2024) for single domain, while the MultiWOZ2.1 dataset Budzianowski et al. (2018) for multi domains. Statistics in Appendix B.1.

**Baselines.** We compare DyBBT against four kinds of comprehensive suite of strong and recent baselines to ensure a rigorous evaluation, and details in Appendix B.2. **DRL Series**: DQN\_ $\epsilon$ \_N (agents are trained using standard DQN 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 Han et al. (2022)), PG (REINFORCE, a stochastic gradient algorithm for policy

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).

Domain	Agent	Epoch = 50			Epoch = 250			Epoch = 500		
Domain	rigent	Success↑	Reward↑	Turns↓	Success↑	Reward↑	Turns↓	Success↑	Reward↑	Turns↓
	DQN_ε_0.0	35.05	-13.00	32.11	54.03	12.99	25.70	55.53	14.95	25.37
	$DQN_{\epsilon}_{0.05}$	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
Movie	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-4B	60.21	38.91	20.14	80.35	51.83	17.15	86.47	55.71	15.64
	DyBBT-8B	65.24	42.14	19.17	85.39	55.06	16.18	89.52	57.64	15.13
	DQN_ε_0.0	06.95	-36.57	27.66	49.07	4.10	22.13	56.71	11.63	23.22
	$DQN_{\epsilon}_{0.05}$	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
Rest.	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-4B	56.03	24.68	19.67	74.86	32.98	16.78	80.55	35.49	15.37
	DyBBT-8B	60.70	26.74	18.69	79.54	35.05	15.81	83.38	36.74	14.86
	DQN_ $\epsilon$ _0.0	00.04	-42.69	27.47	48.46	2.26	24.70	58.79	12.38	23.06
	$DQN_{\epsilon}_{0.05}$	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
Taxi	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-4B	57.57	24.95	20.67	76.78	33.29	17.68	82.62	35.83	16.17
	DyBBT-8B	62.37	27.04	19.69	81.59	35.38	16.71	85.53	37.09	15.66

gradient reinforcement learning Zhu et al. (2023)), PPO (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 to perform each policy update. Zhu et al. (2023)). **LLM based DP**: LLM\_DP (agents use the DP module with GPT-4.0 Yi et al. (2024)), AutoTOD (a zero-Shot autonomous agent with GPT-4.0 Xu et al. (2024), ProTOD (proactive dialog policy based on GPT-4.0 Dong et al. (2025)). **ERL**: EIERL(evolutionary reinforcement learning injected by elite individuals Zhao et al. (2025)). **Multi Agent Collaborative**: MACRM (a multi agent curiosity reward mode for dialog policy Sun et al. (2025))

**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

#### 4.2.1 Performance on Single Domain Tasks

The evaluation results on single domain dialog tasks are presented in Table 1. DyBBT demonstrates strong performance across all three domains. The results reveal that DyBBT's cognitive enables more efficient policy learning: by dynamically allocating computational resources based on real-time cognitive signals, DyBBT achieves higher task success with significantly fewer dialog turns compared to methods relying on static exploration heuristics, population level evolution or GPT-4 based policy. This efficiency gain is particularly pronounced in complex domains like Taxi, where slot dependencies create challenging exploration landscapes that DyBBT navigates more effectively through its principled switching mechanism.

#### 4.2.2 Performance on Multi Domain Task

Results on the challenging MultiWOZ dataset are provided in Table 4 (Appendix D.1). While EIERL's success rate drops significantly in this complex multi domain setting, highlighting the

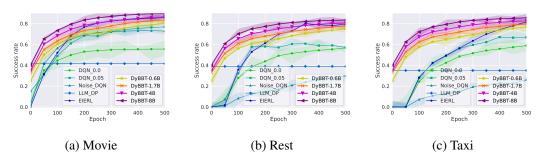


Figure 3: Learning curves for training efficiency and convergence across single-domain TODS tasks.

scalability limits of its population based approach, DyBBT maintains strong performance. DyBBT-8B performs slightly better than AutoTOD/ProTOD (AutoTOD, ProTOD), and using GPT-4 as S2 yields SOTA results, showing that DyBBT matches strong LLM baselines while being more efficient. This is enabled by the structured cognitive state and dual-system design, which provide a domain agnostic inductive bias without requiring task specific tuning.

#### 4.2.3 Training Efficiency and Convergence

Figure 3 illustrates the learning curves of DyBBT compared to baselines. DyBBT converges faster and achieves higher asymptotic performance across all domains, outperforming EIERL significantly at epoch 50. This accelerated learning stems from the meta-controller's active guidance of exploration from the outset, which systematically targets under explored or uncertain regions in  $\mathcal C$  rather than relying on random exploration or high-variance evolutionary mechanisms.

Furthermore, DyBBT exhibits consistent scaling with model size, for instance, success rates improve from 80.34% to 89.52% in the single domain Movie task and from 78.2% to 84.1% in the multi domain setting when scaling from 0.6B to 8B parameters. This trend indicates that the dual-system architecture effectively harnesses the increased representational capacity of larger backbone models. When coupled with the meta-controller's efficient resource allocation with Qwen3's native switching mechanism in balancing performance and computational cost (Appendix D.6). DyBBT underscores its practical viability for real-world deployment.

#### 4.2.4 Summary of Strengths

The main results demonstrate that DyBBT achieves state-of-the-art performance through: **Dynamic Exploration-Exploitation Balance:** The meta-controller's bandit inspired switching rule allows DyBBT to dynamically 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 advancements in LLM capabilities. **Strong Generalization:** Consistent performance across both single and multi domain tasks shows that the cognitive state representation captures universal dialog dynamics. **Computational Practicality:** Unlike population based methods (EIERL) or full GPT-4.0 approaches, DyBBT maintains moderate computational overhead during both training and inference.

#### 4.3 ABLATION EXPERIMENT

Ablation settings and results (Table 5) are detailed in Appendix D.2, revealing 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 the confidence condition (CC) causes a more substantial performance drop than removing the exploration condition (EC), validating our hybrid design. This indicates that mitigating S1's overconfidence is slightly more critical than targeted exploration for robust performance. **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 con-**

**tribute 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  $(\rho_t)$ . **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 EVALUATION

Automated metrics cannot fully capture the quality of decision making, such as action appropriateness or timing of reasoning. We conduct a human evaluation (details in Appendix C) focusing on the meta-controller's switching decisions. Ten NLP researchers evaluated 200 dialog states from Multi-WOZ, comparing DyBBT against random switching and S1 only baselines. Annotators assessed action appropriateness (5 point Likert scale) and whether invoking S2 was justified (binary judgment). The results show that DyBBT's actions are more appropriate than both baselines. Its decisions to invoke S2 align substantially better with human judgment than random switching, providing qualitative evidence that our meta-controller effectively identifies when deliberation is warranted a key affordance often missed by heuristic approaches.

Experimental results demonstrate that DyBBT's meta-controller effectively translates the perception of local cognitive affordances into a dynamic exploration-exploitation balance. This leads to better task success and efficiency, proving that a cognitively grounded adaptive switching mechanism is the key to high performing dialog policies.

#### 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 THE EMERGENT STRUCTURE OF COGNITIVE STATE SPACE

The cognitive state space  $\mathcal C$  serves as the foundational bridge that enables the transfer of bandit exploration principles to the complex dialog POMDP. To empirically validate its utility, we analyze the *visitation frequency* of different regions within the discretized  $\mathcal C$  over training (Fig. 4; detailed computation in Appendix B.5.3). The heatmap reveals a highly structured, non-uniform occupancy pattern, directly validating our core hypothesis. The meta-controller's exploration is not random but strategically focused: in the **early dialog phase**  $(d_t \in [0.0, 0.2])$ , it broadly explores across user uncertainty  $(u_t)$  for information gathering. In the **mid-phase**  $(d_t \in [0.4, 0.6])$ , visitation concentrates in regions of **medium-to-high**  $u_t$ , targeting ambiguity resolution. In the **late phase**  $(d_t > 0.8)$ , activity focuses on states with **low**  $u_t$ , exploiting known information to complete tasks.

This phase dependent targeting demonstrates that  $\mathcal C$  successfully captures the dialog's dynamic "affordances". The meta-controller learns to allocate its exploration budget to the most relevant regions of  $\mathcal C$  for the current dialog stage, enabling highly efficient and context aware exploration. The effectiveness of  $\mathcal C$  stems from its ability to distill the high dimensional belief state into a low dimensional, actionable representation, making principled exploration computationally feasible.

#### 5.2 Adaptive Balancing Through Dual Triggers

The meta-controller's hybrid triggering mechanism provides a robust solution to the exploration-exploitation dilemma by responding to different types of uncertainty:

**Epistemic vs.** Aleatoric Uncertainty Distinction: Two trigger conditions address fundamentally different types of uncertainty. The exploration condition  $(n_t(\mathbf{c}_t) < \tau \sqrt{\log T})$  targets *epistemic uncertainty*, lack of knowledge about the environment that can be reduced through exploration. The confidence condition  $(p_t^{S1} < \kappa)$  addresses *aleatoric uncertainty*, inherent stochasticity or model limitations irreducible via exploration alone. **Complementary Trigger Patterns:** Analyzing 10,000 dialog turns reveals complementary triggering patterns (Fig. 5 in Appendix D.3). The exploration condition dominates in early training phases and for novel state regions, enabling systematic coverage of

the state space. The confidence condition acts as a consistent safety net throughout training, preventing overreliance on a potentially flawed System 1. This complementary design ensures robustness across diverse dialog scenarios. **Progressive Adaptation:** The triggering rate evolves naturally with training progress. Initially, frequent System 2 invocations offer guided exploration and high quality demos. As training progresses and System 1 improves through distillation, the meta-controller automatically reduces System 2 usage, transitioning from guided exploration to autonomous operation. This adaptive balancing is key to DyBBT's computational efficiency and crucially, it is the core manifestation of DyBBT's ability to perceive and respond to the dynamic "affordances" of the dialog environment, ensuring the right cognitive system is invoked at the right time.

#### 5.3 KNOWLEDGE DISTILLATION AS IMPLICIT POLICY IMPROVEMENT

The knowledge distillation process creates a virtuous cycle that enables continuous policy improvement without additional environment interactions. The effectiveness of distillation is evidenced by the monotonic improvement of System 1 and corresponding reduction in System 2 invocation rate (Fig. 6 in Appendix D.3), demonstrating successful knowledge transfer.

#### 5.4 THEORETICAL INTUITIONS AND EMPIRICAL ALIGNMENT

Our theoretical analysis, though based on simplifying assumptions, is pragmatically validated by empirical results: **Sublinear Regret as Validation of Core Assumptions.** The empirical cumulative regret (Fig. 7) exhibits  $\sqrt{T}$ -like growth. This sublinear trend is not merely observational; it provides indirect empirical support for our key theoretical assumptions: The Lipschitz continuity of the reward in  $\mathcal C$  (Assumption 1), and the approximate structure of MDP over  $\mathcal C$  (Assumption 2). The alignment between theory and experiment suggests  $\mathcal C$  effectively captures the latent structure enabling efficient exploration. **Low Dimensional**  $\mathcal C$  **Enables Practical Implementation.** The consistent high performance of DyBBT using only a three dimensional cognitive state demonstrates that the essential features governing exploration (dialog progress, user uncertainty, slot dependency) can be distilled into a compact representation. This reduction in dimensionality is theoretically motivated by the dependence of the regret bound's  $\sqrt{\dim(\mathcal C)}$  (Appendix A.2.2).

#### 5.5 FAILURE MODE ANALYSIS AND LIMITATIONS

Despite its strong performance, DyBBT exhibits three key failure modes that constrain its robustness. It is **over reliant on cognitive state fidelity**: the handcrafted  $\mathbf{c}_t$  can misrepresent complex dialog dynamics such as abruptting intent shifts, leading the meta-controller to misjudge when to invoke System 2, either under exploring or wasting computation. It **depends on high quality System 2 demonstrations**: errors in System 2's reasoning or selfevaluation, even if confident, can be distilled into System 1, corrupting its policy in a subtle cascading failure. It is **sensitive to discretization**: the heuristic quantization of 5 bins of  $\mathcal C$  can mask critical state variations, which underexplore by treating strategically distinct states as identical. These limitations reveal a tension between DyBBT's elegant, theory driven design and the messy reality of dialog environments. Qualitative examples that illustrate these failure modes are provided in D.7. Future work will explore handcrafted and learned hybrid state representations, more robust uncertainty calibration for System 2, and adaptive or continuous exploration bonuses to mitigate these issues.

#### 6 CONCLUSION

DyBBT introduces a principled, cognitively grounded framework for dialog policy learning that dynamically balances exploration and exploitation through a bandit inspired meta-controller operating over a structured cognitive state space. By formalizing dialog affordances, phasic progression, user uncertainty, and slot dependency, our approach enables adaptive, context aware switching between fast intuitive responses and deliberate reasoning. Extensive experiments demonstrate state-of-the-art performance across single and multi domain benchmarks, with human evaluations confirming superior decision quality and alignment with expert judgment. DyBBT offers a scalable, efficient, and interpretable alternative to static or population based methods, bridging cognitive theory with practical dialog optimization. Future work will focus on end-to-end learning of cognitive representations and extending the framework to more complex, interactive settings.

#### **ETHICS STATEMENT**

This work presents a dialog policy learning framework evaluated on publicly available benchmark datasets (MS Dialog and MultiWOZ). Our research does not involve human subjects beyond the use of standard datasets, and all experiments are conducted through simulated user interactions. The proposed methodology focuses on improving the efficiency of task-oriented dialog systems, with potential positive societal impacts through enhanced human-computer interaction. We are unaware of any specific ethical concerns or negative social impacts directly arising from this work.

#### REPRODUCIBILITY STATEMENT

To ensure reproducibility, we have made our code and datasets publicly available at https://anonymous.4open.science/r/DyBBT-C6B7. The appendix provides comprehensive implementation details, including: hyperparameters (Section B.5), dataset statistics (Section B.1), cognitive state computation (Section A.1), and full experimental configurations. All baselines are implemented using standard toolkits (ConvLab-3) with referenced parameter settings. The prompts for System 1 and System 2 are detailed in Section B.4, and the evaluation metrics are formally defined in Section B.3.

#### LLM USE STATEMENT

We utilized DeepSeek V3.1 for translation assistance and grammatical refinement of certain textual passages, and employed Qwen3-Code to aid in debugging and optimizing portions of the experimental code. These LLMs served solely as support tools for improving linguistic clarity and technical implementation. They played no role in the conceptualization of the research, the formulation of methodologies, the analysis of results, or the derivation of scientific conclusions. Consequently, their use does not qualify them as contributors under the authorship criteria. The authors assume full responsibility for all aspects of the work, including the accuracy and integrity of all generated and modified content, and affirm that appropriate measures have been taken to prevent plagiarism and other forms of scientific misconduct.

#### REFERENCES

- Kavosh Asadi, Dipendra Misra, and Michael L. Littman. Lipschitz continuity in model-based reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, ICML'18, 2018.
- Hui Bai, Ran Cheng, and Yaochu Jin. Evolutionary reinforcement learning: A survey. *Intelligent Computing*, 2:0025, 2023. doi: 10.34133/icomputing.0025.
- Paweł Budzianowski, Tsung-Hsien Wen, Bo-Hsiang Tseng, and et al. MultiWOZ a large-scale multi-domain Wizard-of-Oz dataset for task-oriented dialogue modelling. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 5016–5026, Brussels, Belgium, October-November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1547.
- Xiuyi Chen, Fandong Meng, Peng Li, Feilong Chen, Shuang Xu, Bo Xu, and Jie Zhou. Bridging the gap between prior and posterior knowledge selection for knowledge-grounded dialogue generation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 3426–3437, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.275.
- Wenjie Dong, Sirong Chen, and Yan Yang. ProTOD: Proactive task-oriented dialogue system based on large language model. In *Proceedings of the 31st International Conference on Computational Linguistics*, pp. 9147–9164, Abu Dhabi, UAE, January 2025. Association for Computational Linguistics.

- Yihan Du, R. Srikant, and Wei Chen. Cascading reinforcement learning. In B. Kim, Y. Yue, S. Chaudhuri, K. Fragkiadaki, M. Khan, and Y. Sun (eds.), *International Conference on Representation Learning*, volume 2024, pp. 30263–30304, 2024.
  - Dylan J. Foster and Alexander Rakhlin. Beyond ucb: optimal and efficient contextual bandits with regression oracles. In *Proceedings of the 37th International Conference on Machine Learning*, ICML'20. JMLR.org, 2020.
  - Aurélien Garivier and Eric Moulines. On upper-confidence bound policies for switching bandit problems. In *International conference on algorithmic learning theory*, pp. 174–188. Springer, 2011.
  - Shuai Han, Wenbo Zhou, Jiayi Lu, and et al. NROWAN-DQN: A stable noisy network with noise reduction and online weight adjustment for exploration. *Expert Syst. Appl.*, 203:117343, 2022.
  - Tao He, Lizi Liao, Yixin Cao, and et al. Planning like human: A dual-process framework for dialogue planning. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4768–4791, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.262.
  - Wanwei He, Yinpei Dai, Yinhe Zheng, Yuchuan Wu, Zheng Cao, Dermot Liu, Peng Jiang, Min Yang, Fei Huang, Luo Si, et al. Galaxy: A generative pre-trained model for task-oriented dialog with semi-supervised learning and explicit policy injection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
  - Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
  - Xu Jia, Ruochen Zhang, and Min Peng. Multi-domain gate and interactive dual attention for multi-domain dialogue state tracking. *Knowledge-Based Systems*, 286:111383, 2024. ISSN 0950-7051. doi: https://doi.org/10.1016/j.knosys.2024.111383.
  - Junpei Komiyama, Edouard Fouché, and Junya Honda. Finite-time analysis of globally nonstationary multi-armed bandits. *Journal of Machine Learning Research*, 25(112):1–56, 2024.
  - Walter Krämer. Kahneman, D. (2011): Thinking, fast and slow. *Statistical Papers*, 55(3):915–915, 2014. ISSN 1613-9798. doi: 10.1007/s00362-013-0533-y.
  - Wai-Chung Kwan, Hong-Ru Wang, Hui-Min Wang, and Kam-Fai Wong. A survey on recent advances and challenges in reinforcement learning methods for task-oriented dialogue policy learning. *Machine Intelligence Research*, 20(3):318–334, June 2023. ISSN 2731-5398. doi: 10.1007/s11633-022-1347-y.
  - Jing Yang Lee, Kong Aik Lee, and Woon Seng Gan. An empirical Bayes framework for open-domain dialogue generation. In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pp. 192–204, Singapore, December 2023. Association for Computational Linguistics.
  - Changqun Li, Linlin Wang, Xin Lin, and et al. Hypernetwork-assisted parameter-efficient fine-tuning with meta-knowledge distillation for domain knowledge disentanglement. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pp. 1681–1695, 2024.
  - Xiujun Li, Yu Wang, Siqi Sun, and et al. Microsoft dialogue challenge: Building end-to-end task-completion dialogue systems. *arXiv preprint arXiv:1807.11125*, 2018.
- Anthony Liang, Guy Tennenholtz, Chih-Wei Hsu, Yinlam Chow, Erdem Biyik, and Craig Boutilier. Dynamite-rl: a dynamic model for improved temporal meta-reinforcement learning. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, NIPS '24, Red Hook, NY, USA, 2025. Curran Associates Inc. ISBN 9798331314385.

Shuai Ma, Qiaoyi Chen, Xinru Wang, Chengbo Zheng, Zhenhui Peng, Ming Yin, and Xiaojuan Ma. Towards human-ai deliberation: Design and evaluation of llm-empowered deliberative ai for ai-assisted decision-making. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*, CHI '25, New York, NY, USA, 2025. Association for Computing Machinery. ISBN 9798400713941. doi: 10.1145/3706598.3713423.

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, and et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and et al. Bridging the gap between value and policy based reinforcement learning. *Advances in neural information processing systems*, 30, 2017.
- Xuecheng Niu, Akinori Ito, and Takashi Nose. Scheduled curiosity-deep dyna-q: Efficient exploration for dialog policy learning. *IEEE Access*, 12:46940–46952, 2024.
- Baolin Peng, Xiujun Li, Lihong Li, Jianfeng Gao, Asli Celikyilmaz, Sungjin Lee, and Kam-Fai Wong. Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2231–2240, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1237.
- Baolin Peng, Xiujun Li, Jianfeng Gao, and et al. Deep Dyna-Q: Integrating planning for task-completion dialogue policy learning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2182–2192, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1203.
- Libo Qin, Wenbo Pan, Qiguang Chen, and et al. End-to-end task-oriented dialogue: A survey of tasks, methods, and future directions. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 5925–5941, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.363.
- Mahdin Rohmatillah and Jen-Tzung Chien. Hierarchical reinforcement learning with guidance for multi-domain dialogue policy. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 31:748–761, January 2023. ISSN 2329-9290. doi: 10.1109/TASLP.2023.3235202.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- Wenhao Shi, Zhiqiang Hu, Yi Bin, and et al. Math-LLaVA: Bootstrapping mathematical reasoning for multimodal large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 4663–4680, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.268.
- Olivier Sigaud. Combining evolution and deep reinforcement learning for policy search: A survey. *ACM Trans. Evol. Learn. Optim.*, 3(3), September 2023. doi: 10.1145/3569096.
- David Silver, Guy Lever, Nicolas Heess, and et al. Deterministic policy gradient algorithms. In *International conference on machine learning*, pp. 387–395. Pmlr, 2014.
- Jingtao Sun, Jiayin Kou, Wenyan Hou, and Yujei Bai. A multi-agent curiosity reward model for task-oriented dialogue systems. *Pattern Recognition*, 157:110884, 2025. ISSN 0031-3203. doi: https://doi.org/10.1016/j.patcog.2024.110884.
- Vinzenz Thoma, Barna Pasztor, Andreas Krause, Giorgia Ramponi, and Yifan Hu. Contextual bilevel reinforcement learning for incentive alignment. In *Proceedings of the 38th International Conference on Neural Information Processing Systems*, NIPS '24, Red Hook, NY, USA, 2025. Curran Associates Inc. ISBN 9798331314385.
- Tsung-Hsien Wen, David Vandyke, Nikola Mrkšić, and et al. A network-based end-to-end trainable task-oriented dialogue system. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pp. 438–449, Valencia, Spain, April 2017. Association for Computational Linguistics.

- Heng-Da Xu, Xian-Ling Mao, Puhai Yang, Fanshu Sun, and Heyan Huang. Rethinking task-oriented dialogue systems: From complex modularity to zero-shot autonomous agent. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2748–2763, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.152.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- Zihao Yi, Jiarui Ouyang, Yuwen Liu, and et al. A survey on recent advances in llm-based multi-turn dialogue systems. *CoRR*, abs/2402.18013, 2024. doi: 10.48550/ARXIV.2402.18013.
- Jiahao Ying, Yixin Cao, Kai Xiong, and et al. Intuitive or dependent? investigating LLMs' behavior style to conflicting prompts. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4221–4246, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.232.
- Jianxing Yu, Shiqi Wang, Han Yin, Qi Chen, Wei Liu, Yanghui Rao, and Qinliang Su. Diversified generation of commonsense reasoning questions. *Expert Systems with Applications*, 263:125776, 2025. ISSN 0957-4174. doi: https://doi.org/10.1016/j.eswa.2024.125776.
- Ming Zhang, Caishuang Huang, Yilong Wu, Shichun Liu, Huiyuan Zheng, Yurui Dong, Yujiong Shen, Shihan Dou, Jun Zhao, Junjie Ye, Qi Zhang, Tao Gui, and Xuanjing Huang. Transfer TOD: A generalizable Chinese multi-domain task-oriented dialogue system with transfer capabilities. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pp. 12750–12771, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.emnlp-main.710.
- Yangyang Zhao, Kai Yin, Zhenyu Wang, and et al. Decomposed deep q-network for coherent task-oriented dialogue policy learning. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:1380–1391, 2024. doi: 10.1109/TASLP.2024.3357038.
- Yangyang Zhao, Ben Niu, Libo Qin, and Shihan Wang. An efficient task-oriented dialogue policy: Evolutionary reinforcement learning injected by elite individuals. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3429–3442, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.171.
- Qi Zhu, Christian Geishauser, Hsien-Chin Lin, and et al. Convlab-3: A flexible dialogue system toolkit based on a unified data format. In *EMNLP*, pp. 106–123, 2023.

#### A THEORETICAL DETAILS

This section provides the theoretical motivation and intuition behind the DyBBT framework. The following analysis bridges ideas from bandit theory and cognitive science to create a heuristic for exploration in dialog POMDPs. While the full dialog 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 empirically 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  $\mathbf{s}_t$ . We model  $\mathcal{C}$  as a compact 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  $\dim(\mathcal{C}) \leq 3$ , which is crucial for making bandit-style exploration feasible.

The choice of these three dimensions is motivated by their central role in governing the explorationexploitation trade-off in TODS, drawing inspiration from cognitive science and dialog theory:

- 702 704 705
- 706 708 709 710 711 712 713
- 714 715 716 717

- 719 720 721
- 722 723 724
- 725 726 727 728
- 729 730 731
- 732 733 734
- 735 736 737
- 738 739 740 741
- 742 743 744
- 745 746 747

748

749

- 750 751
- 752 753
- 754 755

- Dialog Progress  $(d_t = t/L)$  captures the temporal affordance. Early phases  $(d_t \to 0)$ inherently afford more exploration to gather information, while late phases  $(d_t \rightarrow 1)$  afford exploitation to complete the task. This aligns with the common practice of annealing exploration schedules but provides a continuous, state dependent signal.
- User Uncertainty ( $u_t = |S_{unconfirmed}|/|S_{relevant}|$ ) operationalizes the information gathering affordance. A high  $u_t$  indicates ambiguity in the user's goal, directly signaling the need for information seeking actions to reduce entropy, a well established principle in decision theory.
- Slot Dependency ( $\rho_t = \max_{u \in U} (\frac{1}{|F|} \sum_{f \in F} M(u, f))$ ) captures the structural affordance of the task environment, derived from a pre-computed slot co-occurrence matrix M from the training corpus. A high  $\rho_t$  suggests that the next piece of information is highly predictable given what is already known (e.g., requesting departure after knowing destination in a taxi domain), making 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 Figure 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 justifies 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 1 stated in Section 3.1.2, which posits Lipschitz smoothness of the reward function in the cognitive state space  $\mathcal{C}$ .

**Assumption 2** (MDP over C). The dialog process can be approximately modeled as a finite horizon MDP over the cognitive state space 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  $\pi$  in the cognitive state space is defined as:

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

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.2.1THEORETICAL INTUITION FOR REGRET

Under Assumptions 1 and 2, 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:

$$\mathbb{E}[R(T)] \lesssim \widetilde{\mathcal{O}}\left(L_r \cdot \sqrt{\dim(\mathcal{C}) \cdot T}\right),\tag{3}$$

where  $R(T) = \sum_{t=1}^{T} [V^*(\mathbf{c}_t) - V^{\pi_t}(\mathbf{c}_t)]$  is the cumulative regret, and  $\widetilde{\mathcal{O}}$  hides logarithmic factors. The notation  $\leq$  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 1, bounding the reward's sensitivity to changes in 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  $\epsilon$ .

- 758
- 759 760
- 761 762 763
- 764 765 766
- 767 768 769 770
- 771 772 773
- 774 775 776
- 777 778
- 779 780 781
- 782 783
- 784 785
- 786
- 787 788
- 789 790 791
- 793 794
- 796

798

799 800 801

802 803 804

809

- 1. **Discretization Error:** Due to Lipschitz continuity of  $\bar{r}(\mathbf{c}, a)$  (Assumption 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  $\mathcal{O}(L_r \cdot \sqrt{\dim(\mathcal{C}) \cdot T})$ .

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

This bound provides an intuitive justification for our exploration criterion (Eq. 1 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 (Figure 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 System 2 IF: 
$$\left(n_t(\mathbf{c}_t) < \tau \sqrt{\log T}\right) \vee \left(p_t^{S1} < \kappa\right)$$
.

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 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  $s_t$ , leading to aleatoric uncertainty.
- Model Imperfection: System 1, 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 System 2 when System 1 is uncertain, preventing catastrophic failures. The disjunctive (V) 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 5).

## 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 2 reduces the POMDP to an MDP over 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 (Budzianowski et al. (2018)) for multi domain tasks. Both datasets are converted into ConvLab-3's unified format, ensuring consistency in ontology, state representation, and API interaction. Table 2 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.

Table 2: Summary of dataset statistics for MS Dialog and MultiWOZ 2.1.

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

#### **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., 2024).
- **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 base-line (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:

$$TP = \sum_{d \in D} \sum_{s \in R_d} \mathbb{I}\left(\exists \operatorname{inform}(d, s, v) \in S \land v \notin V_{\operatorname{null}}\right) \tag{4}$$

$$FP = \sum_{d \in D} \sum_{s \notin R_d \cup I_d} \mathbb{I}\left(\exists \operatorname{inform}(d, s, v) \in S \land v \notin V_{\operatorname{null}}\right)$$

$$\tag{5}$$

$$FN = \sum_{d \in D} \sum_{s \in R_d} \mathbb{I} \left( \nexists \operatorname{inform}(d, s, v) \in S \lor v \in V_{\text{null}} \right)$$
 (6)

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

$$Inform = \frac{TP}{TP + FN} \tag{7}$$

#### 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:

$$Book_d = \frac{1}{|B_d|} \sum_{b \in B_d} \mathbb{I}\left(book(d, b, v) \in S \land v = v_{goal}\right)$$
(8)

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

$$Book_{taxi} = \mathbb{I}\left(\exists book(taxi,\cdot,\cdot) \in S\right) \tag{9}$$

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

919 920 921

918

$$Book = \frac{1}{|D_{book}|} \sum_{d \in D_{book}} Book_d$$
 (10)

922 923

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

#### B.3.3 SUCCESS RATE

928 929 The Success Rate represents the overall task completion performance, combining both information provision and booking success:

930

$$Success = \mathbb{I} (Inform = 1 \land Book = 1)$$
 (11)

931 932 933

This binary metric indicates whether both all requested information was provided and all booking operations were successfully completed.

934 935

This metric rewards systems that achieve high success rates with fewer dialog turns, promoting both

936 937

938

939 940

This appendix provides the detailed prompts used for System 1 (intuitive controller) and System 2 (reasoning controller) in the DyBBT framework. The LLM\_DP prompt is the same from the EIERL

942 943 944

945

954 955 956

957 958 959

960 961

962

963

964

965 966 967

968 969 970

971

# effectiveness and efficiency.

# B.4 PROMPT FOR DYBBT AND LLM-DP

paper(Zhao et al. (2025)).

#### B.4.1 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>}
```

#### B.4.2 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}

```
972
973
        **Cognitive State Context:**
974
        - dialog Progress: {d_t}
975
        - User Uncertainty: {u_t}
        - Slot Dependency: {p_t}
976
977
        **Trigger Reason:** {trigger_reason}
978
979
        **Reasoning Guidelines:**
980
        1. **Leverage cognitive signals**:
           - If progress is low, focus on information gathering.
981
           - If uncertainty is high, prioritize clarifying or confirming
982
               actions.
983
           - If slot dependency is high, leverage known slot relationships to
984
               guide next actions.
985
        2. **Consider domain and slot dependencies**:
986
           - E.g., 'taxi' requires both 'destination' and 'departure'; 'restaurant' may require 'area', 'food', 'pricerange' before
987
988
               booking.
989
        3. **Generate 3 distinct strategies** that reflect different tactical
990
            approaches:
991
           - One conservative (e.g., confirm before acting),
992
           - One proactive (e.g., request multiple slots),
993
           - One hybrid (e.g., inform then request).
994
        4. **Evaluate each path** by estimating its likelihood of leading to
995
            task success.
996
997
        **Output Format: ** Strictly adhere to the following JSON schema:
998
999
          "reasoning_paths": [
1000
            {
1001
               "sequence_id": 1
1002
               "action_sequence": [
1003
                 ["action_type", "domain", "slot"],
1004
                 . . .
              ],
1005
               "estimated_success_probability": 0.9
1006
            },
1007
1008
          ]
1009
        }
1010
```

#### B.4.3 LLM\_DP PROMPT

1027

1028

1029

1031

1032

1033

1034

1035

1036

1039 1040

1041 1042 1043

1044

1045

1046

1047

1048 1049

1050

1051

1052 1053 1054

1055 1056

1057

1058

1059

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1074

1075

1076

1077

1078

1079

```
'belief_state': {
            'police': {'book': {'booked': []}, 'semi': {}},
            'hotel': {'book': {'booked': [], 'people': '', 'day': '',
                       'semi': {'name': '', 'area': 'east', 'parking': ''
1030
                          'pricerange': '', 'stars': '4', 'internet': ''
                          'type': ''}},
            'attraction': {'book': {'booked': []}, 'semi': {'type': '', 'name': '', 'area': ''}},
                'name': ''
            restaurant': {'book': {'booked': [], 'people': '', 'day': '',
                'time': ''},
                            'semi': {'food': '', 'pricerange': '', 'name': '',
                               'area': ''}},
            'hospital': {'book': {'booked': []}, 'semi': {'department': ''}},
            'taxi': {'book': {'booked': []},
                      'semi': {'leaveAt': '', '
'', 'arriveBy': ''}},
                                              'destination': '', '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 environment (Zhu et al. (2023)), leveraging its modular architecture for efficient dialog policy optimization. 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  $(\rho_t)$  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  $1 \times 10^{-4}$  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 Section D.4.

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  $\rho_t$  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 (Budzianowski et al. (2018)) 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:

$$M_{ij} = \frac{\operatorname{count}(\operatorname{slot}_i \wedge \operatorname{slot}_j)}{\operatorname{count}(\operatorname{slot}_i)}$$
(12)

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

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

$$\rho_t = \frac{1}{|A_t|(|A_t| - 1)} \sum_{i \in A_t} \sum_{j \in A_t, j \neq i} M_{ij}$$
(13)

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 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  $D_{\text{distill}}$ . 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  $1 \times 10^{-4}$ , 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  $D_{\text{distill}}$  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 r=16, scaling parameter  $\alpha=32$ , and dropout rate of 0.1,

1135

1136 1137

1156 1157

1158 1159

1160

1161

1162

1163

1164

1165

1167

1169

1170

1171 1172

1173

1174

1175

11761177

1182

1183 1184

1186

1187

targeting the query and value projection layers of the transformer architecture. This configuration achieves parameter efficiency while preserving the base model's generalization capabilities.

#### Algorithm 1 Knowledge Distillation Buffer Update and Sampling

```
1138
               Buffer Update:
1139
           1: Input: Current belief state s_t, System 2 action a_t^{S2}, System 2 self evaluated confidence p_{\text{self}}
1140
                                                                                   > Only store high confidence actions
           2: if p_{\text{self}} > 0.9 then
1141
                   if |D_{
m distill}| < {
m MAX\_SIZE} then
1142
                        D_{\text{distill}}.\text{append}((s_t, a_t^{S2}))
           4:
1143
           5:
1144
                        D_{\text{distill}}.pop\_front()
                                                                                           ▶ Remove oldest entry (FIFO)
           6:
                        D_{	ext{distill}}.append((s_t, a_t^{S2}))
1145
           7:
1146
           8:
                   end if
           9: end if
1147
               System 1 Fine-tuning:
1148
          10: Input: System 1 model with LoRA adapters, buffer D_{\text{distill}}
1149
          11: Every 10 training epochs:
1150
          12: for epoch = 1 to 1 do
                                                                                                   ⊳ Fine-tune for 1 epoch
                   for each batch sampled from D_{
m distill} do
          13:
1152
          14:
                        Compute loss \mathcal{L} = \text{CrossEntropy}(\text{System1}(s_i), a_i)
1153
          15:
                        Update LoRA adapter parameters via gradient descent
1154
                   end for
          16:
1155
          17: end for
```

#### B.5.3 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 Section 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.5.4 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**  $R_{\rm emp}(T)$  during training, as shown in Figure 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  $\pi_t$  at training step t.
- $V^{\pi^*}(\mathbf{s}_t)$  is the value of the near-optimal policy  $\pi^*$  at state  $\mathbf{s}_t$ .

Since the true optimal policy  $\pi^*$  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  $\pi^*$  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.

 $R_{\rm emp}(T)$  is plotted against T on a log-log scale to clearly visualize the sublinear growth trend. The theoretical upper bound  $\widetilde{\mathcal{O}}(\sqrt{T})$  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 Section 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 ( $s_t$ ): A structured table showing all relevant slots for the domain(s), their values, and their confirmation status (e.g., confirmed, requested, None).
- Cognitive State ( $c_t$ ): The numerical values for dialog progress ( $d_t$ ), user uncertainty ( $u_t$ ), and slot dependency ( $\rho_t$ ).
- **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
- 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.

Table 3: 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 ↑	Switching Agreement ↑
DyBBT-8B	$\textbf{4.31} \pm \textbf{0.12}$	88.7%
w/o Meta-Controller (Random)	$3.72 \pm 0.19$	52.3%
w/ S1-only	$3.95 \pm 0.15$	<del>_</del>
w/o Exploration Condition (EC)	$4.08 \pm 0.14$	75.4%
w/o Confidence Condition (CC)	$3.89 \pm 0.16$	81.2%

#### 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 3 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.

```
Belief State:
restaurant {
    semi {
        food: "Chinese"  # (USER_CONFIRMED)
        pricerange: "cheap"  # (USER_CONFIRMED)
        area: ""  # (USER_MENTIONED but NOT_CONFIRMED)
        name: ""  # (NOT_MENTIONED - High Uncertainty)
```

1298

1299 1300

1301 1302 1303

1304

1305

1306

1307

1308

1309

1310

1311

1312 1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1326

1327

1328

1330

1331 1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1344 1345

1347

1348

1349

```
}
book { people: "", day: "", time: "" }
}
taxi { ... } # (Not relevant in this turn)
```

Listing 1: Belief state exemplifying high user uncertainty.

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 reasoning 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:

```
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)
    }
}
```

Listing 2: Belief state where the task is complete.

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 allocits computational resources in a manner that is both effective and efficient, closely mirroring human expert judgment.

#### D FURTHER EXPERIMENTAL ANALYSIS

#### D.1 EXPERIMENTAL RESULTS ON MULTIWOZ

Table 4 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.

Table 4: 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 EIERLZhao et al. (2025), other results were reported from original papers, "—" indicates unreported results.

Agent	Year	<b>Inform</b> ↑	<b>Success</b> ↑	Book↑	Turns↓
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
ĎyBBT-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

#### D.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 5:

- **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:  $(n_t(\mathbf{c}_t) < \tau \sqrt{\log T})$ . 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  $\mathbf{c}_t = [d_t, u_t, \rho_t]$  with a three dimensional embedding learned by a small MLP (2 layers, 32 units each) from the raw belief state  $\mathbf{s}_t$ . This tests the necessity of our specific cognitive state design.
- **DyBBT w/o**  $d_t$ , **w/o**  $u_t$ , **w/o**  $\rho_t$ : Ablation studies removing one dimension from the cognitive state at a time to quantify its individual contribution.

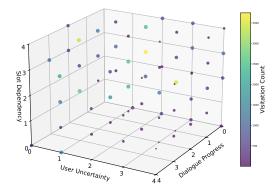
#### D.3 SUPPLEMENTARY ANALYSIS FIGURES AND TABLES

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

- **Figure 4**: Heatmap of visitation frequency in the cognitive state space C, illustrating the structured exploration strategy of the meta-controller across dialog phases.
- **Figure 5**: Analysis of meta-controller decisions, showing the rate of System 2 invocation across dialog progress and the proportion of triggers from each condition.
- **Figure 6**: Demonstrates the improvement of System 1 through knowledge distillation and the corresponding reduction in System 2 invocation over training.

Table 5: 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 <sup>†</sup>	Success <sup>↑</sup>	Book↑	Turns↓
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 $(\rho_t)$	90.3	82.5	85.9	15.0



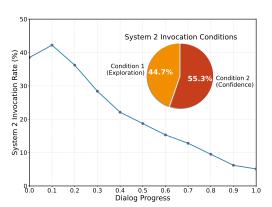


Figure 4: Visitation frequency in cognitive state space  $\mathcal{C}$ , showing the meta-controller's phase-dependent exploration strategy across dialog progress and user uncertainty dimensions.

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

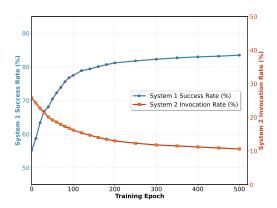
 Figure 7: Compares the empirical cumulative regret of DyBBT against the theoretical upper bound derived under simplifying assumptions.

#### D.4 HYPERPARAMETER SENSITIVITY ANALYSIS

A key concern is the sensitivity of DyBBT's performance to the meta-controller's hyperparameters: the exploration threshold  $\tau$ , the confidence threshold  $\kappa$ , 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 Figure 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



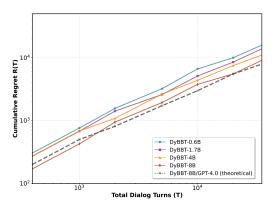


Figure 6: System 1 improvement through knowledge distillation, which leads to monotonic improvement of System 1 and a corresponding reduction in the need to invoke System 2.

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.

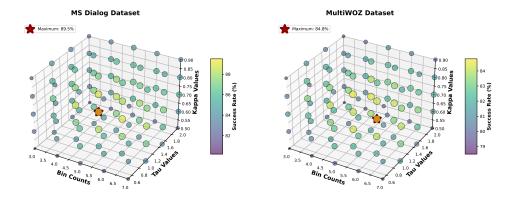


Figure 8: 3D surface plots of success rate (%) as a function of  $\tau$ ,  $\kappa$ , and bin count for (left) MS Dialog and (right) MultiWOZ. The optimal configuration ( $\tau = 1.0$ ,  $\kappa = 0.7$ , bins = 5) is marked with a red star.

reduce the effectiveness of the visitation count. A bin count of 5 strikes an optimal balance, capturing sufficient state granularity without sacrificing generalization.

#### D.5 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 6.

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

Table 6: 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)

<b>Model Family</b>	Size	Params	Success Rate ↑	Inference Time $\downarrow$	<b>Cost-Effectiveness</b> ↑
	1B	1.1B	$78.3 \pm 0.017$	0.32x	244.7
Llama-3.2	3B	3.0B	$80.1 \pm 0.015$	0.48x	166.9
Liaina-3.2	7B	6.7B	$81.9 \pm 0.013$	0.75x	109.2
	8B	8.0B	$82.6 \pm 0.012$	0.89x	92.8
	0.5B	0.5B	$77.4 \pm 0.018$	0.28x	276.4
Owen 2.5	1.5B	1.7B	$79.6 \pm 0.016$	0.41x	194.1
Qwen2.5	3B	2.9B	$81.5 \pm 0.014$	0.59x	138.1
	7B	6.6B	$83.1 \pm 0.012$	0.86x	96.6
	0.6B	0.6B	$79.2 \pm 0.016$	0.35x	226.3
Overage 2	1.7B	1.8B	$81.2 \pm 0.014$	0.52x	156.1
Qwen3	4B	4.2B	$83.6 \pm 0.011$	0.78x	107.2
	8B	8.0B	$85.1 \pm 0.011$	1.00x	85.10

Table 7: Comparison between DyBBT's meta-controller and Qwen3's native switching mechanism. Inference time is normalized to DyBBT's default mode (S1 no think / S2 think = 1.0x).

Configuration	Success Rate ↑	Inference Time ↓	<b>Cost-Effectiveness</b> ↑
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

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.

#### D.6 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 Owen3's internal think mode.
- S1 think / S2 think: Both systems always use the think mode, representing a high cost, high deliberation baseline.
- 3. **S1 no think / S2 think**: DyBBT's mode, System 1 operates in fast mode, while System 2 uses think mode when triggered by the meta-controller.

We report performance on the MultiWOZ test set also using Success Rate, Inference Time (with DyBBT's default mode as 1.0x), and Cost-Effectiveness Results are summarized in Table 7.

As anticipated, the always think configuration achieves the highest success rate (86.5%), confirming that maximal deliberation improves task performance. However, this comes at an prohibitive computational cost 3.2× the inference time of the selective activation of DyBBT. In contrast, DyBBT's

mode achieves nearly comparable performance (85.1% success) with only one-third of computational overhead, resulting in a significantly higher cost-effectiveness.

The no-think baseline performs poorly, underscoring the necessity of deliberate reasoning in complex dialog states. DyBBT strikes a balance between these extremes by invoking costly reasoning only when cognitively justified, either due to under exploration or low confidence, leading to near optimal performance with moderate and targeted computational overhead. This leads to less efficient allocation of computational resources, as also reflected in human evaluation (Section 4.4).

#### D.7 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.

#### D.7.1 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:**

```
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<sub>t</sub>): 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  $(n_t(\mathbf{c}_t) = 12 < \tau \sqrt{\log T} \approx 25)$ , 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:  $request(restaurant, area) \rightarrow confirm(restaurant, area, "north") \rightarrow inform(restaurant, name, "Golden Dragon").$ 

**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.

#### D.7.2 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:**

1620

1621 1622

1623 1624

1640

1641 1642

1643

1644 1645

1646

1647

1648

1649

1650 1651

1652

1653 1654

1655

1657

1658 1659

1660

1663

```
Belief State:
1626
       train {
            semi {
1628
                destination: "cambridge", # (CONFIRMED)
1629
                departure: "london",
                                           # (CONFIRMED)
1630
                day: "today"
                                            # (CONFIRMED)
1631
            book { people: "2" }
                                            # (CONFIRMED)
1632
       }
1633
       taxi {
1634
            semi {
1635
                destination: "london stansted airport", # (CONFIRMED - but
1636
                    contradicts train?)
                departure: "cambridge"
                                            # (CONFIRMED)
1637
            }
       }
1639
```

#### **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 ( $\rho_t$ ): 0.9 (High, 'train' and 'taxi' domains are logically connected)

**Meta-Controller Decision**: The cognitive state was well explored  $(n_t(\mathbf{c}_t) = 505)$ . 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:  $confirm(train, destination, "cambridge") \rightarrow confirm(taxi, destination, "london stansted airport") \rightarrow request(clarification, "final destination").$ 

**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.

#### D.7.3 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:**

```
1664
       Belief State:
       hotel {
1666
            semi {
1667
                name: "hilton"
                                           # (CONFIRMED)
                area: "centre"
1668
                                           # (CONFIRMED)
                parking: "yes",
                                           # (CONFIRMED)
                pricerange: "expensive" # (CONFIRMED)
1670
1671
            book { people: "2", day: "today", stay: "2 nights" } # (BOOKED)
1672
       }
1673
       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 ( $\rho_t$ ): 0.7 (High, 'hotel/area' and 'attraction/area' match)

**Meta-Controller Decision**: The state had medium visitation  $(n_t(\mathbf{c}_t) = 162)$  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.

#### D.7.4 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 of the handcrafted cognitive state  $\mathbf{c}_t$  to represent all critical aspects of the dialog context. Future work will focus on learning this state representation end-to-end from data, which could mitigate such representational gaps and further enhance the framework's robustness and applicability.