DUAL-SCALE WORLD MODELS FOR LLM AGENTS TOWARDS HARD-EXPLORATION PROBLEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

LLM-based agents have seen promising advances, yet they are still limited in "hard-exploration" tasks requiring *learning new knowledge through exploration*. We present GLoW, a novel approach leveraging dual-scale world models, maintaining a trajectory frontier of high-value discoveries at the global scale, while learning from local trial-and-error in exploration through a Multi-path Advantage Reflection mechanism which infers advantage-based progress signals to guide exploration. To evaluate our framework for hard-exploration, we tackle the Jericho benchmark suite of text-based games, where GLoW achieves a new state-of-theart performance for LLM-based approaches. Compared to state-of-the-art RL-based methods, our approach achieves comparable performance while requiring 100-800× fewer environment interactions.¹

1 Introduction

While LLM agents (Yao et al., 2023; Sumers et al., 2024; Wang et al., 2024) excel at leveraging vast pre-trained knowledge in tasks such as robotic planning, software engineering, and web automation (Ahn et al., 2022; Yang et al., 2024; 2025), they are reportedly limited in *hard-exploration problems* (Sutton & Barto, 2018; Ecoffet et al., 2019). Hard exploration problems are typically characterized by large state–action spaces, deceptive local optima, and sparse rewards. These factors often trap naive exploration in local optima, such that exploration fails to reach deeper states with rewards. For LLM agents, such problems pose two central challenges: (1) Global learning, for maintaining long-term knowledge of valuable discoveries during exploration, (2) Local trial-and-error, for quickly refining exploration policies from sparse environmental feedback. Current LLM agent approaches such as ReAct (Yao et al., 2023) or Reflexion (Shinn et al., 2023) support local trial-and-error, but lack mechanisms for long-term knowledge accumulation. Consequently, LLM agents fall short on hard-exploration tasks that humans can often solve effectively (Cui et al., 2025; Phan et al., 2025).

In this work, we introduce Global-Local World Models (GLoW), a framework for LLM agents that enables effective exploration in hard-exploration problems, by maintaining structured world models at two complementary scales for global and local learning. Our approach builds on Go-Explore (Ecoffet et al., 2019) algorithm, which achieves breakthroughs on hard-exploration problems by enhancing the exploration capabilities of RL and LLM-based agents (Lu et al., 2025). The key idea of Go-Explore is to store discovered states into a *state archive*: Then, based on this archive, Go-Explore decomposes hard-exploration into alternating between: (1) a *selection* phase, choosing a *promising* state from the archive to return to, and (2) an *exploration* phase, to continue discovering new states from the selected state. In its original implementation, Go-Explore used hand-crafted heuristics for selection, and random action sampling for exploration, while later work, such as IGE (Lu et al., 2025) improved selection to leverage LLM inference.

In this work, our core insight is that both selection and exploration require structured learning from past exploration experiences, but at different scales: we first enrich beyond an archive of isolated states, by additionally maintaining a trajectory frontier, which keeps the full temporal context of how high value states were reached and why progress stalled, into a **global world model** for richer structured learning. This allows an LLM-based analysis across the frontier to infer high-value regions

¹Code will be open sourced after blind review

as well as bottleneck states with high future potential, enabling principled state selection in GLoW, beyond heuristic or LLM-internalized notions of interestingness. At the local scale, to guide exploration actions from the state, we draw insights that advantage-based rewards better capture progress signals than Q-values (Kazemnejad et al., 2025; Setlur et al., 2025): Our Multi-path Advantage Reflection mechanism explores multiple trajectories from the same starting state and leverages LLM reasoning to infer *intermediate advantages at key state-action pairs*. Through these advantage signals, the **local world model** enables controlled exploration under sparse environmental feedback.

To evaluate the capability of LLM agents in hard-exploration problems, we study the Jericho benchmark suite of text-based games (Hausknecht et al., 2019), where SOTA has been RL-based solutions (Hausknecht et al., 2019; Ammanabrolu & Hausknecht, 2020; Guo et al., 2020) with ε -greedy or softmax exploration or MCTS-based exploration (Jang et al., 2021; Shi et al., 2025). However, they suffer from poor sample efficiency, relying on extensive trial-and-error which requires **hundreds of thousands** of environment interactions. Meanwhile, existing LLM agents were insufficient to address the challenge of learning from exploration in Jericho games, showing limited performance compared to humans (Cui et al., 2025; Phan et al., 2025).

Through extensive experiments, we show that GLoW improves the performance of LLM-based agents while achieving orders of magnitude improvement in sample efficiency compared to RL baselines. Our contributions are summarized as follows:

- We propose GLoW, a novel LLM agent framework for hard-exploration problems through global-local world models,.
- We conduct comprehensive comparisons with existing agent approaches (RL, MCTS, LLM) and ablation studies to validate components of our method.
- We achieve a new state of the art for LLM-based approaches on Jericho, achieving comparable performance with RL-based SOTA, while reducing environment interactions required by 100-800x.

2 BACKGROUND

 Jericho Benchmark The Jericho benchmark (Hausknecht et al., 2019) remains an unsolved hard-exploration problem, where the text-based game environments provide two fundamental challenges (Ammanabrolu & Riedl, 2021): (1) partial observability, requiring agents to construct models of the world from local textual descriptions, and (2) combinatorial state-action spaces. For example in Zork1, the game vocabulary has 697 words and up to five-word commands, resulting in $O(697^5) = 1.64 \times 10^{14}$ possible actions per step, though only a tiny fraction are grammatically coherent and contextually relevant. As a result, RL approaches, with simple exploration strategies, incur hundreds of thousands interactions to offset sample inefficiencies in exploration. This makes Jericho an ideal testbed for evaluating whether agents learn by exploring, rather than brute-force discovery.

Methods for Hard-Exploration Problems Go-Explore (Ecoffet et al., 2019) achieved breakthroughs in hard-exploration problems by maintaining an archive of discovered states as global knowledge to (1) select promising states and (2) explore from the state. Algorithm 1 illustrates this high-level view, using example strategies from the original algorithm: selection returns the next state s_{next} , based on novelty driven heuristics (e.g., less visited states), and explore generates actions (e.g., random action sampling in the original implementation), returning trajectory τ . Appendix B shows adaptations improve upon these heuristics. XTX (Tuyls et al., 2022) adapts imitation learning for selection and DQN for explore, and IGE uses LLM inference for both. Beyond Go-Explore family, MCTS-based methods like MC-LAVE (Jang et al., 2021) and MC-DML (Shi et al., 2025) leverage tree search with language-driven exploration and LLM priors respectively, though requiring 400,000+ interactions.

3 Method

In this section, we describe the dual-scale learning paradigm of GLoW in detail.

Algorithm 1 Go-Explore Algorithm

```
1: procedure GO-EXPLORE-FAMILY(s_0, n_{iter})
2:
          \mathcal{A} \leftarrow \{(s_0, \text{score}_0)\}
                                                                                                                                        ▶ Archive
3:
         for i=1 to n_{iter} do
4:
              s_{next} \sim select(\mathcal{A}) \propto \frac{1}{\text{visits}(s)^{\alpha}}
                                                                                                                                        ▶ Novelty
5:
              \tau \leftarrow explore(s_{next}) \propto \text{RandomActions}(s_{next})
                                                                                                                                   No learning
6:
7:
         end for
8: end procedure
```

3.1 GLOBAL WORLD MODEL FOR STATE SELECTION

The global world model extracts value signals from accumulated exploration trajectories. Unlike traditional state-based archives, we maintain trajectories in a value-ranked frontier. The global world model additionally maintains LLM-generated trajectory analysis.

Value-Ranked Trajectory Frontier As the source of value information, the global world model maintains a trajectory frontier $\mathcal{F}=\{\tau_1,\tau_2,...,\tau_k\}$, containing the k highest-value trajectories discovered during exploration, ranked by a value function $v:\mathcal{T}\to\mathbb{R}$. Each trajectory $\tau_i=(s_0^i,a_1^i,r_1^i,s_1^i,...,a_T^i,r_T^i,s_T^i)$ represents a complete episode generated by the exploration policy π_{explore} defined by the LLM agent, where $s_t\in\mathcal{S}$ are states, $a_t\in\mathcal{A}$ are actions, and $r_t\in\mathbb{R}$ are rewards. For the trajectory value function v, we use the maximum cumulative reward achieved during the episode, $v(\tau_i)=\max_{t\in[1,T]}\sum_{j=1}^t r_j^i$. This is an effective choice for Jericho's sparse reward structure, and the possibility of encountering negative rewards or terminal failures. In contrast to state-only representations, which lose the context of action and observation sequences, preserving complete trajectories enables accurate credit assignment and value estimation in sparse-reward environments where success depends on precise action sequences.

The frontier evolves progressively through iterative exploration. When exploration from selected states (detailed in Section 3.2) produces trajectory τ_{new} with value $v(\tau_{\text{new}})$, the frontier is updated:

$$\mathcal{F}_{t+1} = \text{top-}k(\mathcal{F}_t \cup \{\tau_{\text{new}}\}, v) \tag{1}$$

This sliding window mechanism ensures the frontier maintains diverse high-value strategies, while allowing newly discovered superior trajectories to replace outdated ones. For any state s_i , we can derive the achieved value $v(s_i) = \max_{\tau \in \mathcal{F}, s_i \in \tau} v(\tau)$, representing the maximum value reached from state s_i across all frontier trajectories. By tracking complete trajectories, the frontier serves as both an estimator of achieved values and a repository of successful action sequences.

Motivation: Decomposing value for *select* **and** *explore* Inspired by UCB's value decomposition which balances exploitation with exploration bonus as:

$$\bar{V}(s) + c\sqrt{\frac{\log(N)}{n_s}}$$

where V(s) is the empirical mean value and the second term is the exploration bonus based on visit count n_s , we annotate two types of values v and v', corresponding to each term, by analyzing patterns across all frontier trajectories \mathcal{F} , to extract a set of critical global states with value annotations:

$$W_{\text{global}} = g_{\text{LLM}}(\mathcal{F}) = \{ (s_1, v_1, v_1'), (s_2, v_2, v_2'), \dots, (s_k, v_k, v_k') \}$$
(2)

Here, each (s_i, v_i, v_i') represents a key state identified from frontier analysis by a prompted LLM g_{LLM} , v_i denotes the achieved value from s_i , while v_i' reflects LLM's estimate of future value potential. Importantly, this potential value v_i' cannot be derived from trajectory scores alone, requiring LLM's reasoning about why trajectories fail and what progress could be achieved by resolving current bottlenecks. For instance, a state where multiple trajectories fail might have low achieved value, but have high potential value when: (1) multiple high-value trajectories converge but fail to progress further, suggesting unexplored regions beyond, (2) partial solution patterns indicate missing components, or (3) environmental hints suggest valuable areas remain undiscovered. This implements a semantic form of optimism under uncertainty (Auer, 2003; Brafman & Tennenholtz, 2003) where UCB uses statistical bonuses while we derive optimistic values from LLM analysis of bottlenecks. See Appendix E.1 for a full example of W_{qlobal} generated for Zork1.

Balancing Exploitation and Exploration in State Selection We maintain a state archive $\mathcal{A} = \{(s_i, \operatorname{score}(s_i))\}$ containing discovered states with their achieved scores. Given W_{global} , we select the next exploration state s_{next} by balancing achieved and potential values via LLM:

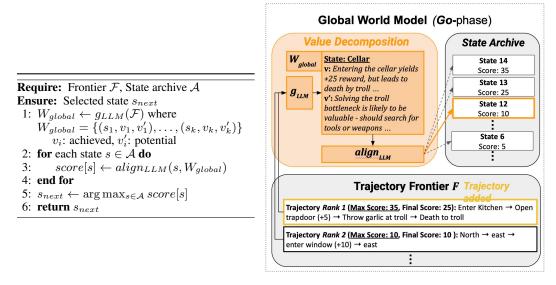


Figure 1: (a) Select procedure in GLoW, (b) Illustration of selection with Global World Model

where $\operatorname{align}_{\operatorname{LLM}}$ evaluates how well each archived state s aligns with the high-value patterns identified in $W_{\operatorname{global}}$. Since $W_{\operatorname{global}}$ contains both achieved and potential values for key frontier states, this alignment naturally balances exploitation (favoring states similar to proven high-reward regions), with exploration (prioritizing states near identified bottlenecks with high potential). Fig. 1 illustrates selection (l.4 in Alg. 1) in GLoW with the Global World Model where a new trajectory (highlighted in gold) has been added to the frontier. Once a state is chosen, we replay the stored sequence of actions to return to the state, which becomes the starting point of the next exploration phase, described in the following section.

3.2 LOCAL WORLD MODEL FOR EXPLORATION

In addition to the selection of states which align with exploration goals with high potential value, exploration can be enhanced by learning which actions are likely to lead to further progress, which is the objective of the local world model.

Motivation: From Q-values to Advantages for Exploration Existing LLM learning methods like self-reflection can be viewed as estimating state-action values (Q-values) from single trajectories. However, Q-value estimation from sparse rewards is notoriously high-variance (Sutton et al., 1999; Schulman et al., 2017), and we observe the same challenge in LLM-based learning: inferences from entire trajectories with sparse feedback are prone to incorrect causal attribution.

Drawing from RL theory, advantage functions A(s,a) = Q(s,a) - V(s) reduce variance by comparing actions to a baseline rather than estimating absolute values. Recent work on process reward models (PRMs) further demonstrates that advantage-based rewards are more suited for exploration, by better capturing progress signals than Q-values, which tend to exploit known strategies (Setlur et al., 2025; Kazemnejad et al., 2025).

Multi-path Advantage Reflection (MAR) Inspired by TRPO (Schulman et al., 2015), computing robust advantage in sparse-reward setting over multiple rollouts from the same state, we propose Multi-path Advantage Reflection to compare multiple trajectories from the same starting state, to produce pseudo-dense advantage signals from sparse environmental feedback. This effectively densifies the reward signal by inferring intermediate advantages at key state-action pairs, providing rich guidance for exploration where environmental rewards are insufficient.

Given a state s selected by the global world model, we perform iterative exploration by sampling n trajectories sequentially: after each trajectory τ_i , we perform MAR to extract learnings that in-

form the next trajectory τ_{i+1} , in the form of world representation W_{local} . This creates a sequence $\{\tau_1, \tau_2, ..., \tau_n\}$ where each trajectory benefits from insights gained from previous attempts.

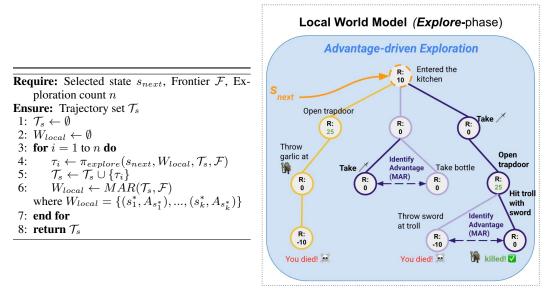


Figure 2: (a) Explore procedure in GLoW, (b) Illustration of exploration with Local World Model

where \mathcal{F} provides global best trajectories as a stable value baseline, $\mathcal{T}_s = \{\tau_1, ..., \tau_n\}$ contains the trajectories sampled during the current exploration phase from state s. States $s_1^*, ..., s_k^*$ are k critical states (typically 2-4) where MAR identifies valuable advantage information can be extracted, either from divergent outcomes revealing good/bad actions, or from consistent patterns confirming reliable strategies. MAR focuses on these few decision points rather than annotating entire trajectories, enabling focused identification of which state-action pairs provide advantages.²

Semantic Advantage Representation Unlike scalar advantages A(s,a), MAR produces W_{local} containing rich semantic advantages which encodes not just which actions are beneficial, but why they work and under what conditions, and captures progress signals which are not expressed by sparse rewards. See Appendix E.2 for a full example of W_{local} generated for Zork1.

Exploration Policy The local world model enhances the *explore* procedure in Alg. 1(l.5), by guiding a policy defined by an LLM agent, as:

$$\pi_{\text{explore}}(a|s_t, h_t) = \text{Agent}_{\text{LLM}}(s_t, h_t, W_{\text{local}}, T_s, \mathcal{F})$$
(3)

where h_t is the current trajectory history, T_s contains previous trajectories in the same exploration phase, and the policy leverages both learned advantages from W_{local} and successful strategies from frontier F. Fig. 2 illustrates exploration (l.5 in Alg. 1) in GLoW with the Local World Model. Consider a trajectory (gold) that reached the cellar but failed at the troll bottleneck without the sword. After analysis by the global world model (Fig. 1), which identifies high v' at the cellar state, this state becomes s_{next} (orange root, Fig. 2). The local world model drives multiple exploration attempts (purple paths), where MAR identifies advantages for "taking sword" despite no immediate reward. This advantage learning guides successful exploration through the troll bottleneck (rightmost path). To address Jericho's exponential action space, we implement a hybrid approach combining free generation with soft constraints. While previous works use either constrained selection from valid actions in RL agents (Hausknecht et al., 2019; Ammanabrolu & Hausknecht, 2020; Tuyls et al., 2022) or pure free-form generation in LLM agents such as ReAct, we provide the valid actions to the LLM as a soft constraint, while still allowing free-form generation. This avoids failure modes of both approaches, where constrained selection can harm action diversity, while pure generation can produce many invalid actions. As we show in Section 4.2, this hybrid approach, which we use consistently across both GLoW as well as all LLM baselines, significantly improves the base LLM performance with only a lightweight prompt and no few-shot examples.

 $^{^2}$ In Appendix A, we provide theoretical analysis showing how MAR reduces variance through both multi-trajectory comparison and the stable baseline based on \mathcal{F} .

4 RESULTS

We evaluate GLoW on the Jericho benchmark suite, We present baselines (Sec. 4.1), setup (Sec. 4.2), main results demonstrating the effectiveness of GLoW (Sec. 4.3), and ablation studies (Sec. 4.4) isolating each module contribution. Lastly, we provide detailed analysis of exploration dynamics in (Sec. 4.5).

4.1 BASELINES

We perform comprehensive comparison against baselines spanning RL-based, MCTS-based, and LLM-based approaches. Furthermore, we compare with specialized methods for hard-exploration problems in each type of baseline. All methods assume access to valid actions from Jericho.

RL-Based Methods: DRRN (He et al., 2016) is a value-based RL approach for choice-based games, learning Q-values for valid actions using GRU encoders and decoders trained via TD loss. **KG-A2C** (Ammanabrolu & Hausknecht, 2020) is a on-policy RL agent that adapts Advantage Actor Critic (A2C) (Mnih et al., 2016), augmented by a dynamic knowledge graph as a state representation that is learned during exploration. Similar to DRRN, **RC-DQN** (Guo et al., 2020) is a DQN-based agent (Mnih et al., 2015), but leverages object-centric neural reading comprehension architectures (Seo et al., 2017) for computing Q-values from observations. **eXploit-TheneXplore** (XTX) (Tuyls et al., 2022) is the current state-of-the-art method in Jericho, implementing Go-Explore with imitation learning on promising trajectories for state selection, and DQN with intrinsic curiosity reward for exploration. RL-based methods rely on million-scale interaction data to learn, leveraging parallel environments for training, with the exception of RC-DQN which leverages 100,000 interactions.

MCTS-Based Methods: Monte Carlo Tree Search is widely adopted for large sequential decision-making problems (Browne et al., 2012; Silver et al., 2016), which explores effectively by combining random sampling and tree search. MC-LAVE (Jang et al., 2021) combines MCTS with language-driven exploration, concentrating search effort on promising actions identified based on value estimates from semantically similar past actions. MC-DML (Shi et al., 2025) enhances MCTS by incorporating LLMs as action priors in the PUCT algorithm (Silver et al., 2016), which balances exploration and exploitation during tree search. The LLM is equipped with a cross-trial memory mechanism, allowing it learn from past experiences such as death in Zork1. Both methods require around 400,000 environment interactions to build comprehensive search trees.

LLM-Based Methods: ReAct (Yao et al., 2023) is the widely adopted standard LLM agent approach interleaving reasoning and acting. **Reflexion** (Shinn et al., 2023) is a multi-episode approach building on ReAct, incorporating self-reflection on each episode to guide future episodes. **In-context Reinforcement Learning** (ICRL) (Song et al., 2025) is another multi-episode approach leveraging in-context reinforcement learning, using cumulative history of past trajectories and rewards as context for future episodes. **Intelligent Go-Explore** (**IGE**) (Lu et al., 2025) implements Go-Explore with LLMs, leveraging LLM-based state selection from a state archive, combined with ReAct-based exploration. As LLM-based baseline methods were not originally applied on Jericho, we re-implement them for Jericho using the action generation approach with valid action soft-constraint described in Sec. 3.2. All LLM-based approaches use 1,000 interactions to balance performance and API cost.³

4.2 EXPERIMENTAL SETUP

Implementation Details Each method is evaluated over 3 random seeds, reporting mean and standard deviation of maximum achieved scores. ReAct performs 20 independent 50-step episodes. Reflexion performs 20 trials of 50-step episodes, incorporating sliding-window memories from up to 10 previous attempts. Likewise, ICRL includes a sliding window of 10 previous trajectories as in-context examples. IGE and GLoW adaptively alternate between state selection and 50-step exploration episodes within the total 1,000 step budget. We use temperature 0.5 for all methods except IGE, which uses 0.3 following Lu et al. (2025). For GLoW hyperparameters, n=3 exploration trajectories and k=5 trajectory frontier size is used.

³We provide details of LLM API usage and cost in Appendix C.1.

Evaluation We evaluate on 10 games from the Jericho benchmark (Hausknecht et al., 2019), spanning different difficulty levels. Following the benchmark's categorization, we test on *possible games* (Pentari, Detective, Temple, Ztuu) featuring moderate puzzles and frequent rewards, *difficult games* (Zork1, Zork3, Deephome, Ludicorp) requiring more complex inventory management, puzzle-solving and navigation, and *extreme games* (Enchanter) involving non-standard actions and spell mechanics. We use the standard Jericho interface providing textual observations and access to valid actions at each step. Unlike some prior work, we do not augment observations with explicit "look" or "inventory" commands, instead allowing agents to learn these through play.

| Games | | RL-l | based | | MCTS | -based | | | LLM-base | d | |
|-----------|-----------|-----------|---------|------------|-----------|-------------|-----------|------------|------------|-----------|-------------------|
| Games | DRRN | KG-A2C | RC-DQN | XTX | MC-LAVE | MC-DML | ReAct | Reflexion | ICRL | IGE | GLoW (Ours |
| Steps | 1,000,000 | 1,600,000 | 100,000 | 800,000 | ~400,000 | ~400,000 | 1000 | 1000 | 1000 | 1000 | 1000 |
| Enchanter | 20 | 12.1 | 20 | 52.0 | - | 20±0.0 | 46.7±9.4 | 48.3±9.4 | 43.3±8.5 | 50.0±7.1 | 61.7 ±20.1 |
| Zork1 | 32.6 | 40.2±0.4 | 38.8 | 103.4±10.9 | 45.2 | 48.66±1.89 | 48.3±4.7 | 48.0±5.0 | 51.7±4.7 | 44.3±0.5 | 73.0±4.5 |
| Zork3 | 0.5 | 0.0 | 2.83 | 4.2±0.1 | _ | 3±0.0 | 3.0±0.0 | 2.7±0.5 | 3.0±0.0 | 3.7±0.9 | 4.3±0.9 |
| Deephome | 1 | 20±2.1 | 1 | 77.7±2.1 | 35 | 67±1.41 | 11.0±4.2 | 22.0±1.6 | 24.0±5.7 | 71.3±4.9 | 75.0±8.7 |
| Ludicorp | 13.8 | 19.8±1.0 | 17 | 78.8 | 22.8 | 19.67±1.7 | 19.7±0.9 | 21.7±1.2 | 32.0±7.1 | 28.3±11.3 | 73.7±11.0 |
| Balances | 10 | 10 | 10 | 24 | 10 | 10±0.0 | 10±0.0 | 10±0.0 | 11.7±2.4 | 10.0±0.0 | 16.7±2.4 |
| Pentari | 27.2 | 44±0.9 | 43.8 | 49.6 | <u>68</u> | 70±0.0 | 30.0±0.0 | 30.0±0.0 | 26.7±4.7 | 30.0±0.0 | 30.0±0.0 |
| Detective | 197.8 | 338±3.4 | 291.3 | 312.2 | 330 | 346.67±9.43 | 113.3±4.7 | 166.7±20.5 | 233.3±47.8 | 316.7±4.7 | 310.0±8.2 |
| Temple | 7.4 | 8 | 8 | - | 8±0.0 | 8±0.0 | 8.7±0.9 | 8.7±0.9 | 8±0.0 | 13.7±0.9 | 13.0±0.0 |
| Ztuu | 21.6 | 5±0.0 | _ | _ | 7 | 23.67±1.9 | 18.7±2.4 | 18.3±2.6 | 16.7±4.1 | 15.0±9.1 | 29.3±4.0 |

Table 1: Comparison of RL-based, MCTS-based, and LLM-based methods on Jericho benchmark games. We report mean \pm standard error over 3 runs. **Bold** indicates best overall performance, and <u>underline</u> indicates second-best. Steps shows total environment interactions. The color of game name indicates original game difficulty categories from Hausknecht et al. (2019): **extreme**, **difficult**, and possible. GLoW achieves state-of-the-art among LLM-based approaches in 7/10 games, and is overall best among all compared approaches in 3/10, second-best in 5/10.

4.3 MAIN RESULTS

We report our main results in Table 1. GLoW achieves a new state-of-the-art performance among LLM approaches across 7 out of 10 games. On Zork1, a canonical game of the Jericho suite, our method reaches a score of 73.0, a significant improvement over the next best LLM method (ICRL at 51.7), and surpassing all compared approaches (with the exception of XTX), including RL and MCTS baselines that use orders of magnitude more interactions. We observe the same strong improvements over the closest LLM method in Ludicorp (73.7 vs. 32.0 for ICRL), Enchanter (61.7 vs. 50.0 for IGE), Ztuu (29.3 vs. 18.7 for ReAct), and Balances (16.7 vs. 11.7 for ICRL).

Notably, our implementation of baselines with hybrid action generation approach shows surprisingly strong performance, whereas prior works reported near-zero scores for LLM agents on Jericho (Shi et al., 2025; Cui et al., 2025; Phan et al., 2025). Our implementation enables ReAct, Reflexion and ICRL to reach 48.3, 48.0, 51.7 on Zork1, respectively, and similarly on par with RL baselines such as KG-A2C and RC-DQN across the board. While this reveals the sample efficiency of LLM agents, these baselines still fall far short of more advanced exploration methods such as XTX and MC-DML, demonstrating the necessity of effective exploration for LLM agents.

Next we compare GLoW against advanced exploration approaches. First, comparing with IGE which is the most directly comparable to ours as an LLM-based Go-Explore method, GLoW substantially outperforms with better performance on 8 out of 10 games. GLoW also achieves competitive performance with state-of-the-art RL and MCTS methods, XTX and MC-DML. We nearly match the overall state-of-the-art XTX, which uses 800× more interactions, on both Deephome (75.0 vs. 77.7) and Ludicorp (73.7 vs. 78.8), and notably surpass it on Enchanter (61.7 vs. 52.0). It also outperforms MC-DML, which employs extensive MCTS-based exploration around 400× more interactions, on most games including Zork1 (73.0 vs. 48.66), Deephome (75.0 vs. 67.0), and Ludicorp (73.7 vs. 19.67). These results demonstrate that our dual-scale approach combining global world models for value-based state selection, with advantage learning for exploration, enables significant performance gains in LLM agents, competitive with sample-intensive RL approaches.

4.4 ABLATION STUDY

To validate the contribution of each component of GLoW, we perform systematic ablations and report the results in Table 2.

Effectiveness of Local World Model We first analyze the efficacy of our local world model by ablating MAR. We replace MAR by Reflexion, which performs the same multi-path exploration but does not leverage our proposed advantage learning, instead performing single-trajectory reflection on the latest trajectory. The results show that the performance drops significantly across most games, demonstrating that MAR's advantage-based formulation more effectively leverages multi-trajectory information than Reflexion, improving exploration under sparse rewards.

Effectiveness of Global World Model Next, we analyze the effectiveness of the global world model, which consists of the frontier of high-value trajectories, and the LLM-based value analysis and alignment state selection. We first ablate the LLM-based value analysis W_{global} , leveraging the raw frontier trajectories for state selection. The negative performance impact shows that, using LLM to reason across the frontier trajectories to infer potential value is indeed effective. Next, we ablate the trajectory frontier $\mathcal F$ altogether, such that it is not used for state selection or leveraged by the exploration policy. This causes further decrease in performance, confirming the contribution of the trajectory frontier in both phases.

Synergy of LWM and GWM Finally, we ablate all the above components together. The resultant model is similar to IGE, with multi-path Reflexion for exploration. The results show that simply adding multi-path reflection does not lead to a clear improvement over IGE, indicating that the overall performance of GLoW comes from the complementary synergy of its components.

| Ablation Variants | Zork1 | Zork3 | Enchanter | Deephome | Ludicorp | Balances |
|--|-----------|---------------------------|-----------------------------|-----------------------------|----------------------------|----------------------------|
| GLoW (Full) | 73.0±4.5 | 4.3±0.9 | 61.7±20.1 | 75.0±8.7 | 73.7±11.0 | 16.7±2.4 |
| X [Local WM] Multi-path Advantage Reflection (MAR) | 70.0±13.6 | 4.3±0.5 | 51.7±9.4 | 56.7±21.7 | 54.7±22.4 | 11.7±2.4 |
| X [Global WM] State selection with W_{qlobal} | 62.0±15.6 | 4.3 ± 0.9 | $60.0{\scriptstyle\pm10.8}$ | 61.3 ± 26.0 | 63.3 ± 14.7 | 13.3 ± 2.4 |
| X [Global WM] Trajectory frontier F √ | 61.7±1.9 | $4.0{\scriptstyle\pm0.8}$ | 53.3 ± 10.3 | 57.7 ± 23.3 | 63.3 ± 12.3 | 11.7 ± 2.4 |
| ✗ All above | 51.3±5.2 | $4.3{\scriptstyle\pm0.9}$ | $51.7{\scriptstyle\pm9.4}$ | $56.0{\scriptstyle\pm21.2}$ | $22.0{\scriptstyle\pm0.8}$ | $10.0{\scriptstyle\pm0.0}$ |
| Standard IGE | 44.3±0.5 | 3.7±0.9 | 50.0±7.1 | 71.3±4.9 | 28.3±11.3 | 10.0±0.0 |

Table 2: Ablation study on GLoW components. We evaluate the contribution of: (1) Local world model through Multi-path Advantage Reflection, (2) Global world model for state selection, (3) trajectory frontier \mathcal{F} .

4.5 ANALYSIS

Controlling global vs local focus with n exploration parameter. We study the tradeoff between local learning depth and global exploration coverage by varying n, the number of explorations per selected state. Larger n enables MAR to learn from more trajectories, while smaller n increases state selection frequency, helping escape local minima. With budget B=1000 and steps s=50, minimum state selections is $m = \lfloor B/(s \cdot n) \rfloor - 1$. With n=1, MAR is turned off. With n>1, MAR analyzes n-1 local trajectories plus the global frontier trajectories.

Table (3) shows that extreme values of n generally yield suboptimal performance. When n=1, effectively disabling MAR, performance drops significantly on certain games like Ludicorp (34.0 vs 73.7 with n=3). Conversely, Deephome shows consistent improvement with increasing n, suggesting it particularly benefits from deeper local exploration. The results demonstrate that moderate increases in n improve performance across several games, consistent with our theoretical analysis (Appendix A) that MAR benefits from variance reduction through multi-trajectory advantage calculation. However, setting n=5 begins to degrade performance, as excessive commitment to individual exploration phases reduces minimum state selection frequency to just 3, increasing susceptibility to local optima. These findings indicate that balancing global and local learning is crucial. We select n=3 as our default parameter, as it achieves the best overall performance by providing sufficient trajectories for robust advantage estimation while maintaining adequate state selection frequency to escape local minima.

Table 3: Controlling the focus on global (less explorations per state but more frequent state selection) vs local learning (more explorations per state). The results demonstrate n=3 exploration from promising states strikes a good balance between the two.

| | Max. | Min. | | | S | core | | |
|------------------------|-----------------------------------|--------------------|------------------------------------|-----------------------------|----------------------------|-------------------------------------|-----------------------------|----------------------------|
| Explorations per State | Steps per Exploration Phase | State Selection | Zork1 | Zork3 | Enchanter | Deephome | Ludicorp | Balances |
| 1 (no MAR) | 50 × 1 | 19 | 59.0±5.7 | 3.7±0.9 | 58.3±9.4 | 59.7±22.6 | 34.0±15.6 | 13.3±4.7 |
| 2 (MAR w/ 1) | 50×2 | 9 | $67.3{\scriptstyle\pm8.7}$ | 3.7 ± 1.2 | 55.0 ± 7.1 | 43.3 ± 26.6 | 66.0 ± 3.7 | $11.7{\scriptstyle\pm2.4}$ |
| 3 (MAR w/ 2) | 50×3 | 5 | $\textbf{73.0} {\pm} \textbf{4.5}$ | $4.3{\scriptstyle\pm0.9}$ | 61.7 ± 20.1 | 75.0 ± 8.7 | $73.7{\scriptstyle\pm11.0}$ | $16.7{\scriptstyle\pm2.4}$ |
| 4 (MAR w/ 3) | 50×4 | 4 | 63.0 ± 6.5 | $4.3{\scriptstyle \pm 0.9}$ | 66.7 ± 10.3 | 73.7 ± 4.5 | 62.0 ± 12.4 | 16.7 ± 2.4 |
| 5 (MAR w/ 4) | 50×5 | 3 | $59.3{\scriptstyle\pm13.8}$ | $4.0{\scriptstyle \pm 0.8}$ | $46.7{\scriptstyle\pm6.2}$ | $\textbf{76.3}{\scriptstyle\pm6.8}$ | $53.3{\scriptstyle\pm7.0}$ | $15.0{\scriptstyle\pm0.0}$ |

5 RELATED WORKS

Go-Explore-based Methods Go-Explore (Ecoffet et al., 2019) enables effective exploration in sparse-reward environments by decomposing exploration into state selection and exploration IGE (Lu et al., 2025) adapts Go-Explore for LLMs, using LLM-based "promisingness" for state selection and ReAct for exploration. However, IGE's limited exploration and ill-defined selection criteria limit its effectiveness in complex environments like Jericho. Our work addresses these limitations through principled value decomposition for selection, and multi-path advantage learning for exploration.

Agents for Text-based Games RL approaches to Jericho include DRRN (He et al., 2016), KG-A2C (Ammanabrolu & Hausknecht, 2020), and RC-DQN Guo et al. (2020), and the aforementioned XTX, where all are sample-intensive, relying on hundreds of thousands of interactions. MCTS-based methods like MC-LAVE Jang et al. (2021) and MC-DML Shi et al. (2025) leverage tree search but still rely on a similar scale of interactions. We show that LLM agents can achieve comparable performance to RL methods, while requiring orders of magnitude fewer interactions through structured exploration and learning mechanisms.

Learning in LLM Agents Recent works have studied how LLMs can learn from experience. Reflexion (Shinn et al., 2023) enables learning through self-reflection on failed attempts, while in-context reinforcement learning (ICRL) (Song et al., 2025) leverages previous trajectories' history as context. However, these approaches struggle with sparse rewards due to noisy learning signals. Our MAR mechanism addresses this challenge through multi-path advantage-based learning, providing more robust learning signals.

World Models for LLM Agents While traditional world models in model-based RL focused on transition dynamics (Ha & Schmidhuber, 2018; Hafner et al., 2024), recent works show that in the context of LLMs, world models are usefully expanded as mechanisms for extracting task-sufficient state representations (Tang et al., 2024; Li et al., 2024). Our dual-scale world models build on these insights, to learn both value patterns across global discoveries, and local advantage signals for exploration.

6 Conclusion

We introduce GLoW, a dual-scale world model framework to tackle hard-exploration problems. GLoW leverages a global world model that enables principled decomposition of state values, and a local world model that integrates trajectories from the same state as controlled exploration feedbacks. Our approach achieves state-of-the-art performance among LLM methods on the challenging Jericho benchmark, while matching RL-based methods that require $800\times$ more environment interactions. By learning global value patterns across discoveries, and local progress signals from multi-path exploration, GLoW overcomes a key limitation of LLM agents in hard-exploration tasks, demonstrating a sample efficient yet high performance results.

REPRODUCIBILITY STATEMENT

To ensure reproducibility of our results, we provide comprehensive implementation details in the paper. Algorithm 2 provides the complete pseudocode for GLoW, and hyperparameters are detailed in Section 4.2 (n=3 exploration trajectories, temperature=0.5, k=5 frontier size, 1000 environment steps). All prompts used for the global world model (Appendix D.1), LLM-based state selection (Appendix D.2), MAR (Appendix D.3), and exploration policy (Appendix D.4) are provided in full. Experiments use GPT-4.1-mini-2025-04-14 as the LLM backbone, reporting results averaged over 3 random seeds with standard deviations. We implement all LLM baselines using the same action generation approach (Section 3.2) for fair comparison. The Jericho benchmark is publicly available, and we use the standard evaluation protocol from Hausknecht et al. (2019). Code implementation will be publicly released upon publication.

REFERENCES

- Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, Alex Herzog, Daniel Ho, Jasmine Hsu, Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil J Joshi, Ryan Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang-Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Sichun Xu, Mengyuan Yan, and Andy Zeng. Do as i can, not as i say: Grounding language in robotic affordances, 2022. URL https://arxiv.org/abs/2204.01691.
- Prithviraj Ammanabrolu and Matthew Hausknecht. Graph constrained reinforcement learning for natural language action spaces. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=Blx6w0EtwH.
- Prithviraj Ammanabrolu and Mark Riedl. Modeling worlds in text. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. URL https://openreview.net/forum?id=7FHnnENUG0.
- Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.*, 3(null):397–422, March 2003. ISSN 1532-4435.
- Ronen I. Brafman and Moshe Tennenholtz. R-max a general polynomial time algorithm for near-optimal reinforcement learning. *J. Mach. Learn. Res.*, 3(null):213–231, March 2003. ISSN 1532-4435. doi: 10.1162/153244303765208377. URL https://doi.org/10.1162/153244303765208377.
- Cameron B. Browne, Edward Powley, Daniel Whitehouse, Simon M. Lucas, Peter I. Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1):1–43, 2012. doi: 10.1109/TCIAIG.2012.2186810.
- Christopher Cui, Xingdi Yuan, Ziang Xiao, Prithviraj Ammanabrolu, and Marc-Alexandre Côté. Tales: Text-adventure learning environment suite. *arXiv preprint arXiv:2504.14128*, 2025. URL https://arxiv.org/abs/2504.14128.
- Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. Go-explore: a new approach for hard-exploration problems. *CoRR*, abs/1901.10995, 2019. URL http://arxiv.org/abs/1901.10995.
- Xiaoxiao Guo, Mo Yu, Yupeng Gao, Chuang Gan, Murray Campbell, and Shiyu Chang. Interactive fiction game playing as multi-paragraph reading comprehension with reinforcement learning. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 7755–7765, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main. 624. URL https://aclanthology.org/2020.emnlp-main.624/.

- David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/2de5d16682c3c35007e4e92982f1a2ba-Paper.pdf.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models, 2024. URL https://arxiv.org/abs/2301.04104.
- Matthew Hausknecht, Prithviraj Ammanabrolu, Marc-Alexandre Côté, and Xingdi Yuan. Interactive fiction games: A colossal adventure. In AAAI 2020, October 2019. URL https://www.microsoft.com/en-us/research/publication/interactive-fiction-games-a-colossal-adventure/. ArXiv.
- Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. Deep reinforcement learning with a natural language action space. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pp. 1621–1630, Berlin, Germany, August 2016. Association for Computational Linguistics. doi: 10.18653/v1/P16-1153. URL https://aclanthology.org/P16-1153/.
- Youngsoo Jang, Seokin Seo, Jongmin Lee, and Kee-Eung Kim. Monte-carlo planning and learning with language action value estimates. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=7_G8JySGecm.
- Amirhossein Kazemnejad, Milad Aghajohari, Eva Portelance, Alessandro Sordoni, Siva Reddy, Aaron Courville, and Nicolas Le Roux. VinePPO: Refining credit assignment in RL training of LLMs. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=Myx2kJFzAn.
- Zichao Li, Yanshuai Cao, and Jackie CK Cheung. Do LLMs build world representations? probing through the lens of state abstraction. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=lzfzjYuWgY.
- Cong Lu, Shengran Hu, and Jeff Clune. Intelligent go-explore: Standing on the shoulders of giant foundation models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=apErWGzCAA.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, February 2015. ISSN 00280836. URL http://dx.doi.org/10.1038/nature14236.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In Maria Florina Balcan and Kilian Q. Weinberger (eds.), *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pp. 1928–1937, New York, New York, USA, 20–22 Jun 2016. PMLR. URL https://proceedings.mlr.press/v48/mniha16.html.
- Long Phan, Mantas Mazeika, Andy Zou, and Dan Hendrycks. Textquests: How good are llms at text-based video games?, 2025. URL https://arxiv.org/abs/2507.23701.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 1889–1897, Lille, France, 07–09 Jul 2015. PMLR. URL https://proceedings.mlr.press/v37/schulman15.html.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL https://arxiv.org/abs/1707.06347.
 - Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. Bidirectional attention flow for machine comprehension. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=HJ0UKP9ge.
 - Amrith Setlur, Chirag Nagpal, Adam Fisch, Xinyang Geng, Jacob Eisenstein, Rishabh Agarwal, Alekh Agarwal, Jonathan Berant, and Aviral Kumar. Rewarding progress: Scaling automated process verifiers for LLM reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=A6Y7AqlzLW.
 - Zijing Shi, Meng Fang, and Ling Chen. Monte carlo planning with large language model for text-based game agents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=r1KcapkzCt.
 - Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 8634–8652. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/1b44b878bb782e6954cd888628510e90-Paper-Conference.pdf.
 - David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–503, 2016. URL http://www.nature.com/nature/journal/v529/n7587/full/nature16961.html.
 - Kefan Song, Amir Moeini, Peng Wang, Lei Gong, Rohan Chandra, Yanjun Qi, and Shangtong Zhang. Reward is enough: Llms are in-context reinforcement learners, 2025. URL https://arxiv.org/abs/2506.06303.
 - Theodore Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas Griffiths. Cognitive architectures for language agents. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=1i6ZCvflQJ. Survey Certification.
 - Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL http://incompleteideas.net/book/the-book-2nd.html.
 - Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), Advances in Neural Information Processing Systems, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
 - Hao Tang, Darren Yan Key, and Kevin Ellis. Worldcoder, a model-based LLM agent: Building world models by writing code and interacting with the environment. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=QGJSXMhVaL.
 - Chen Feng Tsai, Xiaochen Zhou, Sierra S. Liu, Jing Li, Mo Yu, and Hongyuan Mei. Can large language models play text games well? current state-of-the-art and open questions, 2025. URL https://arxiv.org/abs/2304.02868.
 - Jens Tuyls, Shunyu Yao, Sham M. Kakade, and Karthik R Narasimhan. Multi-stage episodic control for strategic exploration in text games. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=Ek7PSN7Y77z.

Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Wayne Xin Zhao, Zhewei Wei, and Jirong Wen. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6), March 2024. ISSN 2095-2236. doi: 10.1007/s11704-024-40231-1. URL http://dx.doi.org/10.1007/s11704-024-40231-1.

John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=mXpq6ut8J3.

Ke Yang, Yao Liu, Sapana Chaudhary, Rasool Fakoor, Pratik Chaudhari, George Karypis, and Huzefa Rangwala. Agentoccam: A simple yet strong baseline for LLM-based web agents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=oWdzUpOlkX.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. ReAct: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.

A THEORETICAL ANALYSIS OF MULTI-PATH ADVANTAGE REFLECTION

A.1 VARIANCE REDUCTION IN MAR

Proposition 1. Let MAR analyze n trajectories $\{\tau_1, ..., \tau_n\}$ starting from state s. For any critical state s^* identified by MAR, let $\hat{A}_{single}(s^*)$ denote an advantage estimate from analyzing a single trajectory through s^* , and $\hat{A}_{MAR}(s^*)$ denote MAR's advantage estimate from comparing $m \leq n$ trajectories that pass through s^* . Under the assumption of bounded variance across trajectories:

$$\operatorname{Var}[\hat{A}_{\operatorname{MAR}}(s^*)] \leq \frac{\operatorname{Var}[\hat{A}_{\operatorname{single}}(s^*)]}{m}$$

Proof. For a trajectory j passing through state s^* and taking action a_j , let $R_j(s^*, a_j)$ denote the random variable representing the sum of future rewards from s^* onward. This provides an unbiased estimate of the true $Q(s^*, a_j)$.

The single-trajectory advantage estimate for action a is:

$$\hat{A}_{\text{single}}(s^*, a) = R_j(s^*, a) - \hat{V}(s^*)$$

where $\hat{V}(s^*)$ is an estimate of the state value. This estimate has high variance because it relies on a single sample: $Var[\hat{A}_{single}(s^*,a)] = Var[R_i(s^*,a)]$ when $\hat{V}(s^*)$ is held constant.

Now consider MAR's approach. From the m trajectories passing through s^* , let m_a denote the number of trajectories taking action a. MAR computes an improved Q-value estimate by averaging outcomes:

$$\hat{Q}_{MAR}(s^*, a) = \frac{1}{m_a} \sum_{j: a_j = a} R_j(s^*, a)$$

Using basic properties of variance for independent random variables with equal variance σ_a^2 :

$$\mathrm{Var}[\hat{Q}_{\mathrm{MAR}}(s^*,a)] = \mathrm{Var}\left[\frac{1}{m_a}\sum_{j:a_j=a}R_j\right] = \frac{1}{m_a^2}\cdot m_a\cdot \sigma_a^2 = \frac{\sigma_a^2}{m_a}$$

This shows variance reduction by factor m_a for the Q-estimate. For the baseline in the advantage calculation, MAR combines the global frontier trajectories \mathcal{F} , with local trajectories through s^* . The advantage estimate is:

$$\hat{A}_{\text{MAR}}(s^*, a) = \hat{Q}_{\text{MAR}}(s^*, a) - \hat{V}_{\text{MAR}}(s^*)$$

MAR's hybrid baseline incorporating frontier trajectories serve a similar role to target networks in DQN (Mnih et al., 2015). While target networks update parameters periodically to provide stable targets, our frontier baseline updates only when superior trajectories are discovered, providing stable value estimates that reduce learning instability. Under the well-founded assumption that this stable baseline has low variance relative to the Q-estimate, the variance of the advantage estimate is dominated by the Q-component:

$$\operatorname{Var}[\hat{A}_{\operatorname{MAR}}(s^*, a)] \approx \operatorname{Var}[\hat{Q}_{\operatorname{MAR}}(s^*, a)] = \frac{\sigma_a^2}{m_a} \leq \frac{\sigma_a^2}{1} = \operatorname{Var}[\hat{A}_{\operatorname{single}}(s^*, a)]$$

More generally, for any action with $m_a \geq 1$ samples, we achieve variance reduction by a factor of m_a . This confirms that MAR reduces variance at each critical state, with greater reduction for actions sampled more frequently. \square

Remark. The proven variance reduction factor of 1/m represents a conservative lower bound for three reasons beyond the statistical averaging captured in the proof:

First, MAR strategically identifies critical states $s_1^*, ..., s_k^*$ where advantage information is most valuable, rather than analyzing entire trajectories. This focused analysis avoids diluting the signal with irrelevant state transitions.

Second, the LLM provides semantic reasoning at these critical states, identifying causal patterns (e.g., lamp necessity for combat in darkness), generalizing across similar states, and leveraging prior knowledge, which are capabilities beyond pure statistical averaging.

Third, our sequential sampling with intermediate MAR reflection means each τ_{i+1} benefits from analysis of $\{\tau_1,...,\tau_i\}$, allowing later trajectories to avoid known failure modes and actively reduce uncertainty about critical decisions.

These enhancements explain why MAR succeeds with small m (typically 2-4 trajectories) in practice.

B ALGORITHMS

We provide the detailed overview of Go-Explore-based algorithms in Alg. 3, and the full algorithm of GLoW in Alg. 2.

C CONTAMINATION CHECK

Table 4: Data contamination analysis: LLM accuracy (%) on navigation questions without seeing gameplay.

| Game | # Questions | Accuracy (%) |
|-----------|-------------|--------------|
| Zork1 | 230 | 10.9 |
| Zork3 | 194 | 8.2 |
| Enchanter | 239 | 9.2 |
| Detective | 66 | 9.1 |
| Balances | 54 | 1.9 |
| Library | 26 | 15.4 |
| Pentari | 70 | 1.4 |
| Deephome | 288 | 17.0 |
| Temple | 92 | 12.0 |
| Ludicorp | 320 | 19.7 |
| Ztuu | 71 | 9.9 |

To assess whether large language models have prior knowledge of Jericho games, we conducted a data contamination analysis following the methodology of Tsai et al. (2025). We evaluate contamination by testing whether models can navigate between locations without being shown any gameplay. Specifically, we: (1) collect a walkthrough trajectory by executing up to 300 steps from

Algorithm 2 GLoW: Global-Local World Models

756

792 793 794

796

797

798

799

800

801

802

803

804

805

806

807

808

809

```
757
               1: procedure GLoW(s_0, n_{iter}, n_{explore}, k)
758
               2:
                          \mathcal{F} \leftarrow \emptyset
                                                                                                                                                     ▶ Initialize frontier
759
               3:
                          \mathcal{A} \leftarrow \{(s_0,0)\}
                                                                                                                                              760
               4:
                          for i=1 to n_{iter} do
                               s_{\text{next}} \leftarrow \text{SELECTSTATE}(\mathcal{F}, \mathcal{A})
                5:
761
                               \mathcal{T} \leftarrow \text{Explore}(s_{\text{next}}, \mathcal{F}, n_{explore})
                6:
762
               7:
                               UPDATEARCHIVE(\mathcal{T}, \mathcal{F}, \mathcal{A}, k)
                          end for
               8:
764
               9:
                          return \arg \max_{\tau \in \mathcal{F}} v(\tau)
765
              10: end procedure
              11:
766
              12: procedure SELECTSTATE(\mathcal{F}, \mathcal{A})
767
              13:
                          W_{\text{global}} \leftarrow g_{\text{LLM}}(\mathcal{F})
768
              14:
                          s_{\text{next}} \leftarrow \arg\max_{s \in \mathcal{A}} \operatorname{align}_{\text{LLM}}(s, W_{\text{global}})
                                                                                                                   ▶ Select state based on decomposed value
769
              15:
                          return s_{\text{next}}
770
              16: end procedure
              17:
771
              18: procedure EXPLORE(s, \mathcal{F}, n)
772
              19:
                          \mathcal{T} \leftarrow \emptyset
                                                                                                 ▶ Initialize trajectory set for current exploration phase
773
              20:
                          W_{\text{local}} \leftarrow \emptyset
774
              21:
                          for j = 1 to n do
775
              22:
                               \tau_j \leftarrow \pi_{\text{explore}}(s, W_{\text{local}}, \mathcal{T}, \mathcal{F})

ightharpoonup Rollout full trajectory from s
                               \mathcal{T} \leftarrow \mathcal{T} \cup \{\tau_j\}
776
              23:
              24:
                               W_{\text{local}} \leftarrow \text{MAR}(\mathcal{T}, \mathcal{F})
777
              25:
                          end for
778
                          return \mathcal{T}
              26:
779
              27: end procedure
              28:
              29: procedure MAR(\mathcal{T}, \mathcal{F})
781
              30:
                          W_{\text{local}} \leftarrow f_{\text{LLM}}(\mathcal{T}, \mathcal{F})
                                                                                                                 782
                          return W_{local}
              31:
783
              32: end procedure
784
              33:
785
              34: procedure UPDATEARCHIVES(\mathcal{T}, \mathcal{F}, \mathcal{A}, k)
                          for \tau \in \mathcal{T} do
              35:
786
                               \mathcal{F} \leftarrow \mathsf{top-}k(\mathcal{F} \cup \{\tau\}, v)
              36:

    □ Update the trajectory frontier

787
              37:
                               for s' \in \tau do
788
              38:
                                     \mathcal{A} \leftarrow \mathcal{A} \cup \{(s', score(s'))\}
                                                                                                                                       ▶ Add states to state archive
789
              39:
                               end for
790
              40:
                          end for
              41: end procedure
791
```

each game's built-in Jericho walkthrough actions, (2) build a graph of locations and transitions from this walkthrough, (3) generate navigation questions asking for paths between observed locations, and (4) query the model with these questions without providing any context. Navigation questions take the form: "In [GAME], what steps would you take to go to [LOCATION B] from [LOCATION A]?" We evaluate responses using strict pattern matching with word boundaries, requiring the exact sequence of navigation commands to appear consecutively in the model's response.

Table 4 shows results of contamination checks for GPT-4.1-mini across 11 Jericho games. We observe minimal contamination, with all games showing below 20% accuracy. Most games (8 out of 11) show less than 10% accuracy, consistent with random guessing or generic text adventure knowledge. The slightly higher accuracies for Ludicorp (19.7%), Deephome (17.0%), and Library (15.4%) likely reflect the model providing common navigation commands (e.g., "go south") that occasionally match by chance. Even famous games like Zork1 (10.9%) show accuracy near chance level, while less-known games like Balances (1.9%) and Pentari (1.4%) show essentially no prior knowledge. These low accuracy rates, combined with the model's generic responses that lack game-specific details, indicate that our experimental results reflect genuine exploration and reasoning capabilities rather than memorized solutions.

Algorithm 3 Go-Explore-based Algorithms

1: **procedure** GO-EXPLORE-FAMILY (s_0, n_{iter}) 2: $\mathcal{A} \leftarrow \{(s_0, \text{score}_0)\}$ ▶ Archive 3: for i = 1 to n_{iter} do — Go Phase (State Selection) — 4: **Go-Explore A:** $s_{next} \sim \text{Uniform}(\mathcal{A})$ ▶ Random sampling Go-Explore B: $s_{next} \sim P(s) \propto \frac{1}{\text{visits}(s)^{\alpha}}$ 5: Novelty **Go-Explore C:** $s_{next} \sim P(s) \propto \text{domain}(s)$ 6: Domain heuristics Domain heuris 7: **XTX:** $s_{next} \leftarrow \text{ImitationLearning}(\mathcal{T})$ ▶ Imitation learning

```
8:
              IGE: s_{next} \leftarrow \text{LLM.SelectPromising}(A)
                                                                                                     ▶ Ill-defined promising-ness
9:
              GLoW: s_{next} = \operatorname{align}_{LLM}(s, W_{global})
                                                                                  ▶ Principled value decomposition (Sec. 3.1)
10:
            — Explore Phase —
11:
              Go-Explore: \tau \leftarrow \text{RandomActions}(s_{next})
                                                                                                                       No learning
12:
              XTX: \tau \leftarrow \text{DQN}(s_{next})
                                                                                                    DQN with curiosity reward
13:
              IGE: \tau \leftarrow \text{ReAct}(s_{next})

    Standard LLM agent

              GLoW: For j = 1 to n:
                                                                ▷ LLM agent with advantage-driven exploration (Sec. 3.2)
14:
15:
                  \tau_j \leftarrow \pi(s_{next}, W_{local})
                  W_{\text{local}} \leftarrow \text{MAR}(W_{\text{local}}, \tau_j, \mathcal{F})
16:
17:
            — Archive Update —
18:
              for each state s' in \tau do
                  if IsNotRedundant(s', A) then
```

C.1 LLM API COST

24: end procedure

end for

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835 836 837

838

839

840

841

842 843 844

845 846

847

848

849 850

851 852

853

854

855 856

857

858

859

861

862

863

19:

20:

21:

22:

23:

We use gpt-4.1-mini-2025-04-14 for all LLM components (\$0.40/\$1.60 per million input/output tokens). Per-run costs of all LLM-based approaches with 1,000 environment steps range from \$4 to \$6, with negligible differences across approaches, maintaining practicality for research iteration.

▷ Domain-specific novelty

D **PROMPTS**

We present the full prompts used in GLoW. Our prompts rely solely on simple instructions and structured output formats without requiring few-shot exemplars, enabling the method to generalize across diverse game scenarios.

D.1 FRONTIER TRAJECTORY ANALYSIS

 $\mathcal{A} \leftarrow \mathcal{A} \cup \{s'\}$

end if

end for

```
Analysis (g_{LLM}) Prompt
Analyze these successful game trajectories to identify patterns and
strategy:
{For each trajectory in \mathcal{F}:}
Trajectory N (Peak: X, Final: Y):
  [score] action -> observation (reward: +/-N if non-zero)
  [score] action -> observation
Based on these trajectories, provide a strategic analysis:
   FRONTIER & EXPLORATION STATUS:
```

```
864
           - What areas/locations have been successfully reached?
865
           - What remains unexplored or inaccessible?
866
867
         2. GAME CHECKPOINTS & PROGRESS:
868
           - What are the key milestones/checkpoints identified?
869
           - What items or abilities unlock new areas?
           - What phase of the game are we in?
870
871
         3. BOTTLENECKS & CHALLENGES:
872
            Where do trajectories commonly get stuck?
873
           - What obstacles block further progress?
874
           - What resources or knowledge are we missing?
875
         4. REWARD STRUCTURE:
876
           - When and how are points earned?
877
           - What actions yield the highest rewards?
878
           - Are there patterns to the scoring?
879
         5. NEXT INVESTIGATION GOALS:
880
           - What specific objectives should we pursue?
881
           - Which unexplored areas are most promising?
882
           - What items or states do we need to reach?
883
884
         Provide a concise strategic summary focusing on actionable
         insights.
885
886
```

D.2 STATE SELECTION

887 888

```
890
891
          State Selection (align<sub>LLM</sub>) Prompt
892
893
          === STRATEGIC GAME ANALYSIS ===
          \{Analysis of frontier trajectories W_{	t global}\}
894
895
896
          Based on the above analysis, select the state from the archive
897
          - Best aligns with the identified investigation goals
898
          - Can help overcome identified bottlenecks
899
          - Explores promising frontiers
900
          - Has potential for high rewards based on patterns
901
902
          Current state archive:
903
          0: [Score: X, Steps: Y, Visits:
  Observation: {state observation}
904
905
            Inventory: {state inventory}
906
907
          1: [Score: X, Steps: Y, Visits:
            Observation: {state observation}
908
            Inventory: {state inventory}
909
910
          . . .
911
912
          Choose state index (0-N).
          Respond in JSON format:
913
914
            "thought": "Your reasoning about which state best aligns with the
915
          strategic goals",
916
            "index": <number>
917
```

919 920

921

922

923

924

925 926

927

928

929

930

931

933

934

935 936

937

938 939

940 941

942

943

944

945

946

947

948

949 950

951

952

953

954

955

956 957

958

959

960

961

962

963

964

965966967

968 969

970

971

D.3 MULTI-PATH ADVANTAGE REFLECTION (MAR)

The MAR prompt generates $W_{\rm local}$ as described in Section 3.2, identifying critical decision points and their associated advantages from multiple exploration trajectories. The prompt incorporates three inputs: (1) the global trajectory frontier containing highest-value trajectories that serve as value baselines, (2) local exploration attempts from the current phase showing different outcomes from the same starting state, and (3) previous $W_{\rm local}$ outputs when available, enabling cumulative learning within the exploration phase.

By comparing outcomes across these trajectory sources, MAR produces $W_{\text{local}} = \{(s_i^*, A_{s_i^*})\}_{i=1}^k$, identifying where specific actions provide clear advantages. This semantic representation captures causal relationships (e.g., "taking the lamp enables combat in darkness") rather than strictly scalar values, enabling the exploration policy to leverage both statistical patterns from trajectory comparison and LLM reasoning about game mechanics at critical states.

```
W_{local} Generation Prompt (MAR)
Review these exploration attempts and identify KEY STATE
ADVANTAGES:
{Previous W_{local} from earlier iterations, if any}
{Global frontier trajectories \mathcal{F}}
{Local exploration trajectories from state s}
Analyze all trajectories and identify ADVANTAGES at KEY STATES:
For each important location/state observed across ALL attempts,
list:
- STATE: [description of state/location]
- ADVANTAGES discovered:
  • [specific action] → [specific benefit/outcome] (score impact if
clear)
  • [what to avoid] → [consequence] (score impact if clear)
  • [optimal sequence] → [why it's better]
Example format:
STATE: At the house entrance with lamp
ADVANTAGES:
   "go east" → finds sword (enabled +15 points later)
  • "open mailbox first" → gets crucial map (+5 immediate)
  • avoid "go upstairs" early \rightarrow wastes moves in empty attic (-7
overall)
Focus on:
1. States that appear across multiple attempts (to see different
outcomes)
2. Critical decision points where scores diverged significantly
3. Action sequences that consistently led to success or failure
   Items or information that enabled later progress
Provide 2-4 KEY STATES with their discovered advantages.
Be specific about actions, items, and locations from the actual
game.
```

D.4 EXPLORATION POLICY IN GLOW

The prompt for the exploration policy in GLoW is shown below. We instruct the LLM to use a JSON output format with reasoning and action fields.

```
You are exploring a text adventure game. Your goal is to make progress and increase your score.

Generate actions that explore new possibilities and make progress.

Respond in JSON format:
{
    "thought": "Your reasoning about what to try",
    "action": "the exact command to execute"
}
```

```
User Prompt (subsequent steps)

Step {N}:
Observation: {current observation}
Score: {current score}
Valid actions: {comma-separated list of valid actions}

What is your next move?
```

E QUALITATIVE EXAMPLES

E.1 GLOBAL WORLD MODEL

```
1008
         Global World Model Example (Zork1)
1009
1010
         Strategic Analysis of Game Trajectories
1011
1012
         ______
1013
1014
         1. FRONTIER & EXPLORATION STATUS
1015
         Successfully Reached Areas:
1016
         - Starting point: Mailbox and House exterior (north, south, east,
1017
         west of house)
1018
         - Forest Path and Clearing with grating
1019
         - Up a Tree (bird's nest with egg and leaflet)
         - Behind House (window entry)
1020
         - Kitchen and Living Room inside the house
1021
         - Attic (Trajectory 1)
1022
         - Cellar and Troll Room (Trajectory 1 & 2)
1023
1024
         Unexplored or Inaccessible Areas:
         - Upstairs beyond the dark staircase (pitch black, dangerous)
1025
```

| - West door with gothic lettering (nailed shut) - Crawlway south in cellar and forbidding hole west in Troll Room - Steep metal ramp in cellar (unclimbable) - Chasm area (Trajectory 1 ends in death) - Any further passages beyond troll or trap door |
|--|
| - Crawlway south in cellar and forbidding hole west in Troll Room - Steep metal ramp in cellar (unclimbable) - Chasm area (Trajectory 1 ends in death) - Any further passages beyond troll or trap door |
| - Crawlway south in cellar and forbidding hole west in Troll Room - Steep metal ramp in cellar (unclimbable) - Chasm area (Trajectory 1 ends in death) - Any further passages beyond troll or trap door |
| - Steep metal ramp in cellar (unclimbable) - Chasm area (Trajectory 1 ends in death) - Any further passages beyond troll or trap door |
| - Chasm area (Trajectory 1 ends in death) - Any further passages beyond troll or trap door |
| - Any further passages beyond troll or trap door |
| 2. GAME CHECKPOINTS & PROGRESS Key Milestones: Obtaining jewel-encrusted egg and leaflet from mailbox/tree Using egg and leaflet in grating to unlock cellar access Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) Opening trap door and descending into cellar Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: Egg + Leaflet → Grating → Cellar entrance Lantern (turned on) → Safe navigation in dark cellar Sword (elvish, glowing) → Combat with troll Rope and knife (potential utility, not fully leveraged yet) Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters 3. BOTTLENECKS & CHALLENGES Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death Navigating dark areas without proper light (pitch black upstairs Attempting risky actions (jumping over chasm leads to death) Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) Hostile NPC (troll) Environmental hazards (darkness, chasm) Limited inventory or missing key items Missing Resources/Knowledge: Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| **Rey Milestones: - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar **Items/Abilities** Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) **Current Phase:* Early—to—mid game — exploration of house and cellar, initial combat encounters - Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death Navigating dark areas without proper light (pitch black upstairs Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) **Obstacles Blocking Progress:* - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items **Missing Resources/Knowledge:* - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| **Rey Milestones: - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar **Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) **Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters - Combat with troll (Trajectories 1 & 2 end in death or near death - Navigating dark areas without proper light (pitch black upstairs - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) **Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items **Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| **Rey Milestones: - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar **Items/Abilities** Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) **Current Phase:* Early-to-mid game - exploration of house and cellar, initial combat encounters - Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death Navigating dark areas without proper light (pitch black upstairs Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) **Obstacles Blocking Progress:* - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items **Missing Resources/Knowledge:* - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Wey Milestones: - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Obtaining jewel-encrusted egg and leaflet from mailbox/tree - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Using egg and leaflet in grating to unlock cellar access - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Entering house via window and acquiring items (food, garlic, water, lantern, sword, rope, knife) - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| water, lantern, sword, rope, knife) Opening trap door and descending into cellar Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: Egg + Leaflet → Grating → Cellar entrance Lantern (turned on) → Safe navigation in dark cellar Sword (elvish, glowing) → Combat with troll Rope and knife (potential utility, not fully leveraged yet) Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Opening trap door and descending into cellar - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Defeating or engaging troll in cellar Items/Abilities Unlocking New Areas: - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| <pre>Items/Abilities Unlocking New Areas: Egg + Leaflet → Grating → Cellar entrance Lantern (turned on) → Safe navigation in dark cellar Sword (elvish, glowing) → Combat with troll Rope and knife (potential utility, not fully leveraged yet) Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters </pre> |
| - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Egg + Leaflet → Grating → Cellar entrance - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Lantern (turned on) → Safe navigation in dark cellar - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| - Sword (elvish, glowing) → Combat with troll - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters - Combat with croll (Trajectories 1 & 2 end in death or near deat - Navigating dark areas without proper light (pitch black upstair - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| - Rope and knife (potential utility, not fully leveraged yet) - Food, garlic, water (likely survival or puzzle items) Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| Current Phase: Early-to-mid game - exploration of house and cellar, initial combat encounters |
| Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near deat - Navigating dark areas without proper light (pitch black upstair - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death Navigating dark areas without proper light (pitch black upstain Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| 3. BOTTLENECKS & CHALLENGES Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death nea |
| Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death - Navigating dark areas without proper light (pitch black upstairs - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death - Navigating dark areas without proper light (pitch black upstains: Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Common Stuck Points: - Combat with troll (Trajectories 1 & 2 end in death or near death - Navigating dark areas without proper light (pitch black upstain - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| <pre>- Combat with troll (Trajectories 1 & 2 end in death or near death - Navigating dark areas without proper light (pitch black upstairs - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs</pre> |
| Navigating dark areas without proper light (pitch black upstairs - Attempting risky actions (jumping over chasm leads to death) Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: Physical barriers (nailed door, barred trap door) Hostile NPC (troll) Environmental hazards (darkness, chasm) Limited inventory or missing key items Missing Resources/Knowledge: Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door Safe traversal methods for chasm or dark upstairs |
| - Attempting risky actions (jumping over chasm leads to death) - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| - Locked or barred doors (gothic door, trap door closing behind player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| <pre>player) Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs</pre> |
| Obstacles Blocking Progress: - Physical barriers (nailed door, barred trap door) - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Physical barriers (nailed door, barred trap door) Hostile NPC (troll) Environmental hazards (darkness, chasm) Limited inventory or missing key items Missing Resources/Knowledge: Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door Safe traversal methods for chasm or dark upstairs |
| Physical barriers (nailed door, barred trap door) Hostile NPC (troll) Environmental hazards (darkness, chasm) Limited inventory or missing key items Missing Resources/Knowledge: Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door Safe traversal methods for chasm or dark upstairs |
| - Hostile NPC (troll) - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| - Environmental hazards (darkness, chasm) - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| - Limited inventory or missing key items Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Missing Resources/Knowledge: - Effective combat strategy or stronger weaponry to defeat troll safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs |
| Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door Safe traversal methods for chasm or dark upstairs |
| Effective combat strategy or stronger weaponry to defeat troll safely Means to reopen or bypass barred trap door Safe traversal methods for chasm or dark upstairs |
| <pre>safely - Means to reopen or bypass barred trap door - Safe traversal methods for chasm or dark upstairs</pre> |
| Means to reopen or bypass barred trap doorSafe traversal methods for chasm or dark upstairs |
| - Safe traversal methods for chasm or dark upstairs |
| |
| |
| J 111, 1, 1 00000 200000 |
| |
| ======================================= |
| |
| A DEMADD CODITION |
| 4. REWARD STRUCTURE |
| When Points Are Earned: |
| |
| |
| - +5 for taking the egg |
| - +5 for taking the egg - +10 for acquiring food and water items in kitchen |
| - +5 for taking the egg |
| - +5 for taking the egg - +10 for acquiring food and water items in kitchen |

1080 Highest Reward Actions: 1081 - Descending into cellar (+25) 1082 - Collecting key items early (+5 to +10) 1083 - Progressing past major checkpoints 1084 1085 Scoring Patterns: - Early exploration yields small incremental rewards 1086 - Major area transitions yield large rewards 1087 - Death penalizes score, encouraging cautious play 1088 1089 _____ 1090 5. NEXT INVESTIGATION GOALS 1091 1092 Specific Objectives: 1093 - Develop reliable strategy to defeat or bypass troll without dying 1094 - Investigate use of rope and knife for puzzles or combat 1095 - Explore crawlway south and forbidding hole west in troll room - Find safe method to cross or bypass chasm 1096 - Attempt to open or circumvent nailed gothic door 1097 1098 Most Promising Unexplored Areas: 1099 - Crawlway south and forbidding hole west in cellar/troll room 1100 - Upstairs dark staircase (with lantern or other light source) - Beyond barred trap door if it can be reopened 1101 - Areas beyond chasm once safe crossing method found 1102 1103 1104 E.2 LOCAL WORLD MODEL 1105 1106 Local World Model Example (Zork1) 1107 1108 Based on detailed exploration attempts, here are KEY STATES with 1109 discovered advantages: 1110 1111 _____ 1112 STATE: The Troll Room 1113 Small room with passages east/south, hole west; troll blocks all 1114 exits; sword glows 1115 1116 **ADVANTAGES:** - "hit troll with sword" → kills troll, clears all exits for 1117 exploration (score: 40) 1118 - "take axe" after troll death → obtains valuable weapon for later 1119 puzzles 1120 - avoid engaging without glowing sword → likely failure 1121 - optimal: kill troll first, then take axe (mandatory to proceed) 1122 _____ 1123 1124 STATE: Maintenance Room (Flood Control Dam #3) 1125 Room with colored buttons, tool chests, wrench, screwdriver, 1126 toothpaste, leaking pipe 1127 ADVANTAGES: 1128 - "take wrench, screwdriver, tube" → essential tools for 1129 environment interaction 1130 - "push blue button" o triggers leak, raises water level to access 1131 new areas - "push red button" → toggles lights, affects water level

- avoid throwing lantern → breaks critical light source

1132

```
1134
         - optimal: collect tools → manage buttons → control water without
1135
         drowning
1136
1137
         ______
1138
1139
         STATE: Temple / Torch Room / Dome Room / Altar
         Large temple with inscriptions; dome with railing; rope for
1140
         descent; ivory torch; brass bell; gold coffin
1141
1142
         ADVANTAGES:
1143
         - "take ivory torch" → stable light for deeper cave exploration
         - "take bell" → key item for spirit/wraith interaction
1144
         - "ring bell at Entrance to Hades" → paralyzes wraiths, enables
1145
         passage
1146
         - "blow out candles" → enables safe descent or passage
1147
         - optimal: acquire torch → bell → sceptre → manipulate altar →
1148
         control spirits
1149
1150
1151
         STATE: East-West Passage / Chasm Area
1152
         Narrow passage with stairs; chasm with paths; multiple routes
1153
         (north/east/west/up/down)
1154
         ADVANTAGES:
1155
         - "east" then "north" \rightarrow leads to Reservoir South and further areas
1156
         - "tie rope to railing" → enables safe descent into lower levels
1157
         - avoid getting stuck in loops → wastes moves
1158
         - optimal: explore chasm edges → use rope for vertical → access
1159
         Dome/Torch
1160
         ______
1161
1162
         Cross-Cutting Insights:
1163
         - Inventory Management: Strategic dropping/picking essential for
1164
         critical artifacts
         - Light Preservation: Maintaining lantern/torch crucial for dark
1165
         exploration
1166
         - Combat Readiness: Glowing sword indicates combat opportunity
1167
         (essential for progress)
1168
```

F LLM USAGE

We utilized Claude for minor grammar and language edits in paper writing.