

000 DREAMPHASE: OFFLINE IMAGINATION AND 001 UNCERTAINTY-GUIDED PLANNING FOR LARGE- 002 LANGUAGE-MODEL AGENTS 003 004

006 **Anonymous authors**

007 Paper under double-blind review

011 ABSTRACT

013 Autonomous agents capable of perceiving complex environments, understanding
014 instructions, and performing multi-step tasks hold transformative potential across
015 domains such as robotics, scientific discovery, and web automation. While large
016 language models (LLMs) provide a powerful foundation, they struggle with closed-
017 loop decision-making due to static pretraining and limited temporal grounding.
018 Prior approaches either rely on expensive, real-time environment interactions or
019 brittle imitation policies, both with safety and efficiency trade-offs. We intro-
020 duce DREAMPHASE, a modular framework that plans through *offline imagination*.
021 A learned latent world model simulates multi-step futures in latent space;
022 imagined branches are scored with an uncertainty-aware value and filtered by a
023 safety gate. The best branch is distilled into a short natural-language reflection
024 that conditions the next policy query, improving behavior without modifying the
025 LLM. Crucially, DreamPhase attains its performance with *substantially fewer*
026 *real interactions*: on WebShop, average API calls per episode drop from ~ 40
027 with ARMAP-M (token-level search) to < 10 with DreamPhase, a $4\times$ reduction
028 that lowers latency and reduces executed irreversible actions by $\sim 5\times$ on
029 WebShop (4.9 \times on ALFWORLD) per incident logs. Across web, science, and embod-
030 ied tasks, DreamPhase improves sample efficiency, safety, and cost over search-
031 based and reward-based baselines. This offers a scalable path toward safe, high-
032 performance autonomous agents via imagination-driven planning. Code: <https://anonymous.4open.science/r/DreamPhase-A8AD/README.md>.

034 1 INTRODUCTION

036 Building intelligent agents that can perceive complex environments, follow instructions, and au-
037 tonomously complete tasks remains a core challenge in artificial intelligence. Such agents have the
038 potential to transform a wide range of applications, including virtual assistants, scientific discovery,
039 workflow automation, and robotics (Brooks, 1986; Reed et al., 2022; Kirchdorfer et al., 2024; Rana
040 et al., 2023).

041 A promising foundation for building these agents lies in large-scale generative models, particularly
042 large language models (LLMs), which have demonstrated remarkable capabilities in natural language
043 understanding, question answering (Rajpurkar et al., 2016), summarization (Hermann et al., 2015),
044 and multimodal reasoning (Chen et al., 2015; Goyal et al., 2017). However, while LLMs excel in
045 static tasks, they remain limited in settings that require multi-step decision-making and closed-loop
046 interaction, such as online shopping, scientific experimentation, or puzzle solving. This limitation
047 stems from the fact that most LLMs are pre-trained on internet-scale static data and lack exposure
048 to temporally grounded trajectories or dynamic environments, which are essential for reasoning
049 about actions and their long-term consequences (Zhou et al., 2024). Additionally, many state-of-
050 the-art LLMs—like GPT-4V (OpenAI et al., 2024) and Gemini (Reid et al., 2024)—are accessible
051 only through restricted APIs. This severely limits their adaptability: they cannot be fine-tuned on
052 task-specific interaction data, and inference-time prompting often proves brittle in tasks involving
053 uncertainty, feedback, or long-term planning.

To address these challenges, two broad strategies have emerged, each with its own limitations.

(i) **Online roll-out planners**, such as ReAct with beam search or Monte Carlo tree search (MCTS) variants that interact with the live document object model (DOM) at each step (Yao et al., 2022; Hao et al., 2023). These methods explicitly evaluate multiple future branches by performing hundreds of real clicks. While this approach enables lookahead and can correct for earlier mistakes, it is often slow, costly in settings with rate-limited APIs, and potentially hazardous in environments where actions are irreversible, such as submitting payments or placing orders.

(ii) **Pure imitation or reward-model agents**, on the other hand, avoid expensive interaction by acting greedily from the current state without explicit search (Hong et al., 2023; Liu et al., 2023). Although this reduces interaction overhead, it introduces brittleness: a single misstep can irreversibly derail the entire trajectory, since the agent lacks any foresight or ability to revise its plan.

Both approaches ultimately trade off between safety and efficiency due to their reliance on real-time interaction with the environment. A promising alternative is *internal imagination*: the ability to simulate and evaluate future outcomes offline, without taking any real actions in the environment.

To realize this, we introduce **DreamPhase**, a modular framework that enables autonomous agents to plan through imagination in latent space. At its core, DreamPhase trains a latent world model that predicts the next DOM tree or visual frame conditioned on the current latent state and a proposed action. By iteratively applying this model, the agent can simulate multiple future trajectories entirely offline, estimate their expected value and uncertainty, and prune high-risk branches before making any real requests. This approach yields three key advantages:

(i) **Sample efficiency**: Web-based tasks that previously required approximately 40 real clicks per episode now converge in fewer than 10 (see Section 5.2).

(ii) **Safety**: Imagination allows the agent to identify dead ends and avoid irreversible outcomes, such as "Sold-out" pages or accidental purchases, without executing them (see Appendix C).

(iii) **Cost and latency**: All expensive reasoning is handled on-device; the only network call is the final, high-confidence action that passes uncertainty filtering.

In summary, by relocating exploration from the real environment into a learned latent simulator, DreamPhase overcomes the safety-efficiency trade-off that constrains existing agent systems.

In particular, DreamPhase unfolds in four stages: (i) we train a latent world model (LWM) (Ha & Schmidhuber, 2018) to simulate future states from raw DOM trees or visual inputs; (ii) we generate imagination rollouts by combining this latent dynamics model with actions sampled from a frozen policy LLM; (iii) we perform uncertainty-aware value estimation over these imagined branches, scoring each using predicted rewards and entropy-based confidence metrics; and (iv) we distill the best low-risk trajectory into a natural-language reflection, which is injected back into the LLM prompt to influence future decisions.

Crucially, the LLM policy remains entirely *frozen* throughout this process. Its behavior evolves solely through internal simulation and language-based feedback. This design supports scalable and sample-efficient planning without the need for real-world interaction or model fine-tuning. In summary, the main contributions of this paper are:

- **Imagination-first latent planning.** We train a compact world model to simulate multi-step futures in latent space and roll out M branches for H steps entirely offline, decoupling environment dynamics from policy reasoning and enabling safe counterfactual planning.

- **Risk-aware branch selection with guarantees.** Imagined trajectories are scored by value minus an uncertainty penalty and accepted only if a confidence gate passes; we provide a regret bound $\mathcal{O}(\sqrt{T\varepsilon}) + B\rho T$ that links decision quality to model error and mis-gating.

- **Language reflections for zero-tuning control.** The selected branch is distilled into a short reflection and summary that are injected into the next prompt, steering a frozen LLM toward higher-quality actions while remaining interpretable and requiring no parameter updates.

- **Practical efficiency and robustness.** DreamPhase reduces real interactions by about 4 \times on WebShop (<10 vs. \sim 40 API calls/episode), adds \sim 12 ms imagination overhead, and cuts executed irreversible actions by \sim 5 \times on WebShop (4.9 \times on ALFWorld). It degrades gracefully under distribution shift via fallback and outperforms search- and reward-based baselines across domains.

108 **2 RELATED WORK**

110 **LLMs as interactive agents.** LLMs have been used as control policies in dynamic environments,
 111 from web navigation to general tool use (Liu et al., 2023; Zhou et al., 2023). Early systems adapted
 112 pretrained encoders to decision-making (e.g., BERT for WebShop (Devlin et al., 2019; Yao et al.,
 113 2023a)). With stronger generative models (Brown et al., 2020; OpenAI et al., 2024), zero-shot and
 114 few-shot prompting became common for action selection without fine-tuning (Deng et al., 2024;
 115 Xiong et al., 2024), often pairing observations and histories with an LLM to choose the next step
 116 (Zheng et al., 2024; Hong et al., 2023). Another line distills trajectories from powerful but restricted
 117 APIs into smaller policies (Li et al., 2023; Zhang et al., 2023). In contrast, DREAMPHASE keeps
 118 the policy LLM frozen and adds a learned latent world model, an uncertainty-aware value head, and
 119 language reflections to steer behavior while minimizing real interactions.

120 **AgentLM, AgentGym, and ARMAP.** *AgentLM* Zeng et al. (2024) introduces *agent tuning*, a
 121 supervised and preference-style fine-tuning on multi-turn tool trajectories, to improve grounding and
 122 robustness, but it requires curated data and changes model weights; DREAMPHASE avoids fine-tuning
 123 by conditioning a frozen LLM with reflections distilled from imagined rollouts. *AgentGym* Xi et al.
 124 (2025) standardizes evaluation with unified task wrappers, tools, and metrics; our work follows its
 125 emphasis on comparable protocols but contributes a new planning method. *ARMAP* Chen et al.
 126 (2025) scales token-level search with Reflexion, Best-of- N , and MCTS (ARMAP-R/B/M), improving
 127 success via online expansion at the cost of higher latency and many real-environment calls. We show
 128 that DREAMPHASE exceeds ARMAP performance while using substantially fewer real interactions
 129 through offline latent imagination and uncertainty gating. For a broader survey, see Appendix B.

131 **3 METHODOLOGY**

133 **Setting and notation.** We study decision making in partially observable interactive environments,
 134 for example web navigation with sparse feedback and costly interactions. The task is modeled as a
 135 partially observable Markov decision process

$$137 \quad \mathcal{M} = (\mathcal{I}, \mathcal{S}, \mathcal{A}, \mathcal{X}, \mathcal{T}, \mathcal{E}, r),$$

138 where \mathcal{I} is the space of task instructions, \mathcal{S} the latent state space, \mathcal{A} a discrete action space, \mathcal{X} the
 139 observation space, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$ the transition kernel, $\mathcal{E} : \mathcal{S} \rightarrow \mathcal{X}$ the emission kernel, and
 140 $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ a possibly sparse reward. At time t , the agent receives an instruction $\iota \in \mathcal{I}$, observes
 141 $\mathbf{x}_t \in \mathcal{X}$, selects $\mathbf{a}_t \in \mathcal{A}$, the environment moves $\mathbf{s}_{t+1} \sim \mathcal{T}(\mathbf{s}_t, \mathbf{a}_t)$, and emits $\mathbf{x}_{t+1} \sim \mathcal{E}(\mathbf{s}_{t+1})$. The
 142 objective is to maximize

$$144 \quad J(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] \quad \text{under a policy } \pi(\mathbf{a}_t | \mathbf{x}_{\leq t}, \iota), \quad (1)$$

146 where $r(\mathbf{s}_t, \mathbf{a}_t)$ is the reward at time t , $\gamma \in (0, 1]$ is a fixed discount factor, and γ^t exponentially down-
 147 weights rewards farther in the future, and the expectation is taken under the trajectory distribution
 148 induced by the policy $\pi(\mathbf{a}_t | \mathbf{x}_{\leq t}, \iota)$.

150 **3.1 DREAMPHASE OVERVIEW**

152 DreamPhase plans with offline imagination to lower interaction cost and improve safety. At each
 153 step, the agent: (i) forms a compact predictive belief with a learned latent world model, (ii) rolls out
 154 multiple hypothetical futures in latent space, (iii) scores each rollout with a value estimate and an
 155 uncertainty measure, (iv) distills the highest quality imagined outcomes into a short natural language
 156 reflection that conditions the next action. All steps (i) to (iv) occur without querying the environment.

158 **3.2 LEARN TO DREAM: TRAINING THE LATENT WORLD MODEL**

160 The first component of DreamPhase is a learned latent world model that enables internal simulation
 161 of environment dynamics. The goal is to predict what will happen when the agent applies an action
 in the current observation context. For example, the model can answer: “What occurs if I click

162 this button on the current page.” This capability supports counterfactual reasoning before acting,
 163 without touching the real environment. Concretely, we train a latent world model that predicts
 164 short-horizon futures in a compact latent space conditioned on the current observation, the action, and
 165 the instruction ι . This eliminates environment queries during planning and follows prior latent-space
 166 imagination work (Ha & Schmidhuber, 2018; Hafner et al.; Janner et al., 2019).
 167

168 **Data collection.** We collect multi-step interaction episodes for tasks such as website navigation or
 169 shopping. Each episode provides tuples $(\iota, \mathbf{x}_t, \mathbf{a}_t, \mathbf{x}_{t+1})$, where $\mathbf{x}_t, \mathbf{x}_{t+1} \in \mathcal{X}$ are observations and
 170 $\mathbf{a}_t \in \mathcal{A}$ is the action taken at time t . The history is $h_t = (\iota, \mathbf{x}_{\leq t}, \mathbf{a}_{<t})$.

171 We train the latent world model exclusively on logged trajectories collected from the standard training
 172 split of each benchmark environment. To obtain these trajectories, we run a frozen LLaMA-2-7B
 173 policy with a simple ReAct-style prompt and light randomization. No test episodes, privileged futures,
 174 or environment states unavailable to baselines are used. All observations and actions come from the
 175 same public APIs that baseline agents interact with. The resulting dataset contains only (instruction,
 176 observation, action, next-observation) tuples from the training portion of each task and does not
 177 include any extra external corpora or fine-tuning data.

178 In addition, note that all the baselines in Table 1 are allowed to use their standard training protocols
 179 (e.g., AgentLM and AgentGym use their released agent-tuned checkpoints; ARMAP uses its reward-
 180 model training). DreamPhase does not access any interaction data beyond the logged trajectories from
 181 the training split. Thus, all open-source agents share the same policy backbone, the same environment
 182 splits, and access to the same interaction data; DreamPhase only differs in how this data is used to
 183 train a compact world model and value head for offline imagination.

184 **Observation encoding.** Observations such as DOM trees or screenshots are high dimensional
 185 and structured. We tokenize each \mathbf{x}_t by a depth-first traversal of the DOM with textual content
 186 and spatial layout features, for example bounding boxes and element types. This yields a compact,
 187 language-aligned sequence suitable for generative modeling and for conditioning the policy.
 188

189 **World model architecture.** The model has three parts: an encoder $f_\theta(\mathbf{x}_t)$ that maps an observation
 190 to a latent \mathbf{z}_t , a transition model $g_\theta(\mathbf{z}_t, \bar{\mathbf{a}}_t, \iota)$ that predicts the next latent given the current latent
 191 and an embedded action $\bar{\mathbf{a}}_t$, and a decoder $d_\theta(\mathbf{z}_{t+1}, \iota)$ that reconstructs the next observation. The
 192 stochastic latent dynamics are

$$\mathbf{z}_t = f_\theta(\mathbf{x}_t), \quad \mathbf{z}_{t+1} \sim q_\theta(\mathbf{z}_{t+1} \mid \mathbf{z}_t, \bar{\mathbf{a}}_t, \iota), \quad \mathbf{x}_{t+1} \sim g_\theta(\mathbf{x}_{t+1} \mid \mathbf{z}_{t+1}, \iota), \quad (2)$$

193 where $\bar{\mathbf{a}}_t = \text{emb}_A(\mathbf{a}_t)$ is a learned action embedding.
 194

195 **Training objective.** We combine token reconstruction with latent space regularization. Let $\hat{\mathbf{x}}_{t+1} =$
 196 $d_\theta(g_\theta(f_\theta(\mathbf{x}_t), \bar{\mathbf{a}}_t, \iota), \iota)$. The loss is
 197

$$\mathcal{L}_{\text{LWM}} = \mathbb{E} \left[\text{CE}(\hat{\mathbf{x}}_{t+1}, \mathbf{x}_{t+1}) + \lambda_{\text{KL}} \text{KL}(q_\theta(\mathbf{z}_{t+1} \mid h_t, \mathbf{a}_t) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I})) \right], \quad (3)$$

198 where **KL** and **CE** denote Kullback–Leibler divergence and cross-entropy function, respectively, and
 199 $q_\theta(\mathbf{z}_{t+1} \mid h_t, \mathbf{a}_t)$ is the encoder-induced inference distribution (diagonal Gaussian with parameters
 200 from f_θ) over the next latent during training; we sample via reparameterization and use the mean at
 201 test time, and λ_{KL} is annealed during training.
 202

203 **Offline imagination.** After training, the model simulates futures without environment queries.
 204 Given \mathbf{x}_t and \mathbf{a}_t , we encode to $\mathbf{z}_t = f_\theta(\mathbf{x}_t)$, sample \mathbf{z}_{t+1} from the learned transition $g_\theta(\cdot \mid$
 205 $\mathbf{z}_t, \bar{\mathbf{a}}_t, \iota)$, and decode $\hat{\mathbf{x}}_{t+1} = d_\theta(\mathbf{z}_{t+1}, \iota)$. Iterating this procedure for a short horizon yields
 206 imagined trajectories $\tilde{\tau} = (\tilde{\mathbf{x}}_{t+1:t+H}, \tilde{\mathbf{a}}_{t:t+H-1})$, which are later scored for value and uncertainty
 207 and summarized into reflections.
 208

209 **Remark 1** (Why a latent world model rather than LLM-based dreaming). *Large language models
 210 are not designed to simulate low-level environment dynamics such as DOM structure. They often
 211 produce syntactically invalid or causally inconsistent states, and token-level rollouts are memory
 212 intensive and conflate policy reasoning with environment modeling. A learned latent world model
 213 is modular and efficient, it generates structured futures that respect environment constraints, and it
 214 supports risk-aware planning in long-horizon tasks. For these reasons DreamPhase relies on a latent
 215 world model for internal simulation rather than prompting an LLM to generate future states.*

216 3.3 IMAGINATION-BASED PLANNING
217

218 At each timestep t , the agent explores multiple plausible futures by simulating latent trajectories
219 conditioned on potential actions. This allows the agent to reason about consequences without
220 interacting with the real environment. We refer to this mechanism as imagination-based planning.

221 Given the current observation \mathbf{x}_t , we encode it with the encoder f_θ to obtain $\mathbf{z}_t = f_\theta(\mathbf{x}_t)$. To explore
222 alternative outcomes, we generate M parallel rollouts in latent space, each simulating a sequence
223 of H future steps using the learned world model $(f_\theta, g_\theta, d_\theta)$ and a frozen policy LLM π_{LLM} that
224 proposes actions from the current imagined history. Each rollout corresponds to a hypothetical branch
225 that contains predicted latents and actions.

226 Each trajectory $\tilde{\tau}^{(j)}$ is one imagined branch, consisting of an
227 action sequence and the resulting predicted latents, with optional decoded observations $\tilde{\mathbf{x}}$
228 for inspection or scoring. These imagined futures are passed to the downstream evaluation module,
229 where they are scored by a value model and filtered by uncertainty metrics. This process supports multi-step foresight in latent space without issuing environment queries, which makes planning efficient and scalable. We detail the rollout procedure in Algorithm 1.
230
231
232
233
234
235
236
237
238
239
240
241
242

Algorithm 1: Latent Imagination Rollouts.

Input: instruction ι ; observation \mathbf{x}_t ; encoder f_θ ; transition g_θ ; decoder d_θ ; frozen policy π_{LLM} ; number of branches M ; horizon H
1: Encode observation: $\mathbf{z}_t \leftarrow f_\theta(\mathbf{x}_t)$
2: **for** $j = 1$ to M **do**
3: Initialize latent: $\mathbf{z}_t^{(j)} \leftarrow \mathbf{z}_t$
4: Initialize imagined history: $h_t^{(j)} \leftarrow (\iota, \mathbf{x}_t)$
5: **for** $k = 0$ to $H - 1$ **do**
6: Sample imagined action: $\tilde{\mathbf{a}}_{t+k}^{(j)} \sim \pi_{\text{LLM}}(\cdot | h_{t+k}^{(j)})$
7: Embed action: $\bar{\mathbf{a}}_{t+k}^{(j)} \leftarrow \text{emb}_A(\tilde{\mathbf{a}}_{t+k}^{(j)})$
8: Predict next latent: $\mathbf{z}_{t+k+1}^{(j)} \sim g_\theta(\mathbf{z}_{t+k+1} | \mathbf{z}_{t+k}^{(j)}, \bar{\mathbf{a}}_{t+k}^{(j)}, \iota)$
9: Decode to imagined observation: $\tilde{\mathbf{x}}_{t+k+1}^{(j)} \leftarrow d_\theta(\mathbf{z}_{t+k+1}^{(j)}, \iota)$
10: Update imagined history: $h_{t+k+1}^{(j)} \leftarrow (h_{t+k}^{(j)}, \tilde{\mathbf{a}}_{t+k}^{(j)}, \tilde{\mathbf{x}}_{t+k+1}^{(j)})$
11: **end for**
12: Store imagined trajectory: $\tilde{\tau}^{(j)} \leftarrow (\{\tilde{\mathbf{a}}_{t:t+H-1}^{(j)}\}, \{\mathbf{z}_{t+1:t+H}^{(j)}\})$
13: **end for**

243 3.4 EVALUATE AND DECIDE: UNCERTAINTY-AWARE VALUE-GUIDED BRANCH SELECTION
244

245 After generating M latent rollouts of horizon H , the agent must assess whether any imagined
246 trajectories are reliable enough to guide action selection. This balances *efficiency* (acting from
247 internal simulation) and *accuracy* (deferring to the real environment when predictions are uncertain).

248
249 **Value estimation.** For each imagined trajectory $\tilde{\tau}^{(j)} = \{\mathbf{z}_{t+1:t+H}^{(j)}, \tilde{\mathbf{a}}_{t:t+H-1}^{(j)}\}$, we estimate a
250 discounted return with a lightweight value head V_ϕ operating on latents:

$$251 \quad 252 \quad 253 \quad G^{(j)} = \sum_{k=1}^H \gamma^{k-1} V_\phi(\mathbf{z}_{t+k}^{(j)} | \iota), \quad \gamma \in (0, 1]. \quad (4)$$

254 The head V_ϕ is trained on logged data to predict short-horizon utility from latent states produced by
255 the world model.
256

257 **Uncertainty estimation.** We quantify epistemic uncertainty over imagined actions using a mutual-
258 information proxy computed with Monte-Carlo dropout on the frozen policy π_{LLM} . For each imagined
259 history prefix $h_t^{(j)}$ (instruction ι , current \mathbf{x}_t , and all imagined pairs so far), let $\mathcal{U} \subseteq \mathcal{A}$ be the discrete
260 action set scored by the policy. With stochastic dropout masks $\{\xi_n\}_{n=1}^N$,

$$262 \quad 263 \quad 264 \quad u^{(j)} = \mathbb{H}[\bar{\mathbf{p}}^{(j)}] - \frac{1}{N} \sum_{n=1}^N \mathbb{H}[\mathbf{p}^{(j)}(\xi_n)], \quad \bar{\mathbf{p}}^{(j)} = \frac{1}{N} \sum_{n=1}^N \mathbf{p}^{(j)}(\xi_n), \quad (5)$$

265 where $\mathbf{p}^{(j)}(\xi_n)$ is the action distribution $\pi_{\text{LLM}}(\cdot | h_t^{(j)}, \xi_n)$ over \mathcal{U} , and $\mathbb{H}[\cdot]$ is categorical entropy.
266 Larger $u^{(j)}$ indicates greater disagreement across stochastic policy samples, hence higher uncertainty.
267

268 **Risk-sensitive selection.** We combine value and uncertainty via a penalized score
269

$$\tilde{G}^{(j)} = G^{(j)} - \beta u^{(j)}, \quad \beta \geq 0, \quad (6)$$

270 and choose the best branch
 271

$$272 \quad j^* = \arg \max_{j \in \{1, \dots, M\}} \tilde{G}^{(j)}. \quad (7)$$

$$273$$

274 A safety gate enforces that the selected branch is itself confident:
 275

276 “if $u^{(j^*)} \leq \tau$ then act using $\tilde{a}_t^{(j^*)}$; otherwise query the real environment and update the history.”
 277

278 This mechanism enables DreamPhase to plan through imagination when predictions are confident,
 279 and to defer to real interaction when uncertainty is high, supporting robust decision making under
 280 distribution shift while reducing environment queries.
 281

282 3.5 REFLECT AND ACT: LANGUAGE-GUIDED POLICY CONDITIONING

283 Rather than fine-tuning the policy LLM, we steer its behavior using natural language derived from
 284 imagination. Inspired by verbal self-reflection (Shinn et al., 2024; Yuan et al., 2024a), DreamPhase
 285 converts its imagined rollouts into concise guidance that conditions the next action.
 286

287 **Reflection.** Given the set of imagined trajectories $\{\tilde{\tau}^{(j)}\}_{j=1}^M$ and their scores $\{G^{(j)}, u^{(j)}\}$ from
 288 Section 3.4, let j^* be the selected branch. We generate a short reflection that explains why $\tilde{\tau}^{(j^*)}$ is
 289 promising and surfaces potential risks. A lightweight reflection head \mathcal{R}_φ produces
 290

$$291 \quad c_t = \mathcal{R}_\varphi(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)}, G^{(j^*)}, u^{(j^*)}), \quad (8)$$

$$292$$

293 where $\kappa \leq H$ is a short summary horizon. \mathcal{R}_φ can be a small frozen LLM or a distilled model
 294 conditioned on the branch and its value–uncertainty signals. Example output: “*Proceed with search,*
 295 *then filter by size before adding to cart; avoid pages without a visible ‘Checkout’ button.*”
 296

297 **Trajectory summarization.** In addition, we produce a compact natural-language summary of the
 298 core actions and expected outcomes. A summarizer \mathcal{S}_η maps the same inputs to a terse script,
 299

$$300 \quad s_t = \mathcal{S}_\eta(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)}), \quad (9)$$

$$301$$

302 for example: “*Search ‘Nike shoes’; open first result; click ‘Add to cart’.*” Both c_t and s_t can optionally
 303 aggregate over the top- K safe branches by concatenating key constraints, although the default uses
 304 only j^* .
 305

306 **Language-based steering.** We inject the reflection and summary into the policy prompt at the next
 307 decision point. The frozen policy LLM receives enriched context and produces the next action:
 308

$$309 \quad \mathbf{a}_t \sim \pi_{\text{LLM}}(\cdot | \iota, \mathbf{x}_t, c_t, s_t). \quad (10)$$

$$310$$

311 This language-guided conditioning integrates simulated experience in a modular and interpretable way.
 312 The agent improves behavior without parameter updates, using internal planning signals expressed as
 313 short, actionable text that is easy to inspect and ablate. Implementation details for the reflection head
 314 \mathcal{R}_φ and summarizer \mathcal{S}_η , the 30-token budget, the delimiter scheme, and task-specific templates are
 315 provided in Appendix I.
 316

317 3.6 DREAMPHASE PIPELINE

318 Algorithm 2 summarizes the full DreamPhase pipeline. At each timestep, the agent encodes the
 319 current observation into a latent state and imagines M future trajectories using the learned world
 320 model. These rollouts are scored by discounted return estimates and filtered by an uncertainty proxy.
 321 If the best rollout is sufficiently confident, it is summarized and reflected upon in natural language,
 322 which is then appended to the policy LLM prompt for action selection. Otherwise, the agent queries
 323 the policy using the real history only. This design enables offline reasoning, defers to real interaction
 324 when uncertain, and steers decisions with internal reflection, without updating the LLM.
 325

324 **Algorithm 2** DREAMPHASE: Imagination-Based Planning with Uncertainty-Aware Reflection

325

326 **Input:** instruction ι ; initial observation \mathbf{x}_0 ; history $h_0 \leftarrow (\iota, \mathbf{x}_0)$; encoder f_θ , transition g_θ , decoder

327 d_θ ; policy π_{LLM} ; value head V_ϕ ; reflection head \mathcal{R}_φ ; summarizer \mathcal{S}_η ; number of imagined

328 branches M ; horizon H ; discount γ ; risk weight β ; confidence threshold τ , **summary horizon** κ .

329 1: **for** $t = 0, 1, \dots$ **do**

330 2: Encode current observation: $\mathbf{z}_t \leftarrow f_\theta(\mathbf{x}_t)$

331 3: Sample M imagined rollouts $\{\tilde{\tau}^{(j)}\}_{j=1}^M$ using the world model and π_{LLM} as in Algorithm 1

332 4: **for** $j = 1 \dots M$ **do** ▷ Score each imagined branch

333 5: Estimate value: $G^{(j)} \leftarrow \sum_{k=1}^H \gamma^{k-1} V_\phi(\mathbf{z}_{t+k}^{(j)} \mid \iota)$

334 6: Estimate uncertainty $u^{(j)}$ via the mutual-information proxy from Section 3.4

335 7: Penalized score: $\tilde{G}^{(j)} \leftarrow G^{(j)} - \beta u^{(j)}$

336 8: **end for**

337 9: Select best branch: $j^* \leftarrow \arg \max_{j \in \{1, \dots, M\}} \tilde{G}^{(j)}$

338 10: **if** $u^{(j^*)} \leq \tau$ **then** ▷ Safety gate passes

339 11: Generate reflection: $c_t \leftarrow \mathcal{R}_\varphi(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)}, G^{(j^*)}, u^{(j^*)})$

340 12: Summarize trajectory: $s_t \leftarrow \mathcal{S}_\eta(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)})$

341 13: Augment policy context: $p_t \leftarrow \text{CONCAT}(h_t, c_t, s_t)$

342 14: Query policy with guidance: $\mathbf{a}_t \sim \pi_{\text{LLM}}(\cdot \mid p_t)$

343 15: **else** ▷ Fallback to real-history policy

344 16: $\mathbf{a}_t \sim \pi_{\text{LLM}}(\cdot \mid h_t)$

345 17: **end if**

346 18: Execute \mathbf{a}_t in the environment, observe \mathbf{x}_{t+1}

347 19: Update history: $h_{t+1} \leftarrow (h_t, \mathbf{a}_t, \mathbf{x}_{t+1})$

348 20: **end for**

4 THEORETICAL REMARK

353 Let π^* be the optimal policy in the true environment \mathcal{M} . The latent world model induces a one-
 354 step predictive distribution $\widehat{\mathbb{P}}_\theta(\mathbf{x}_{t+1} \mid h_t, \mathbf{a}_t, \iota)$ and the environment induces $\mathbb{P}_*(\mathbf{x}_{t+1} \mid h_t, \mathbf{a}_t, \iota)$.
 355 Assume a uniform one-step KL bound $\varepsilon = \sup_{h_t, \mathbf{a}_t} \text{KL}(\mathbb{P}_* \parallel \widehat{\mathbb{P}}_\theta)$, which by Pinsker inequality
 356 (Csiszár & Körner, 2011) implies $\text{TV} \leq \sqrt{\varepsilon/2}$, where **TV** denotes the total variation distance. Let
 357 τ be the confidence threshold of the safety gate (Sec. 3.4), and define the mis-gating rate $\rho =$
 358 $\Pr[\text{DreamPhase selects an imagined branch while the true uncertainty} > \tau]$.

359 With bounded rewards r_{\max} and a value function that is Lipschitz in total variation of the one-step
 360 predictor, the cumulative regret over T steps satisfies

$$362 \quad \text{Regret}_T = \mathbb{E} \left[\sum_{t=0}^{T-1} (r_t^* - r_t^{\text{DP}}) \right] \leq C \sqrt{T \varepsilon} + B \rho T,$$

$$363$$

$$364$$

365 where $C > 0$ depends on r_{\max} and the smoothness constant, and $B \leq r_{\max}$ bounds the per-step loss
 366 under mis-gating. A full derivation is given in Appendix A.

367 The $\sqrt{T \varepsilon}$ term captures error from model approximation; $B \rho T$ accounts for occasional acceptance
 368 of unreliable imagined branches. When ε is small and ρ is rare, regret grows sublinearly.

5 EXPERIMENTS

374 We evaluate DREAMPHASE across diverse agentic environments and backbones. We first outline the
 375 evaluation protocol, then present two main results tables: (i) a diverse-task comparison across closed-
 376 and open-source agents (Table 1), and (ii) cross-backbone results on WebShop, SciWorld (seen/un-
 377 seen), and Game-of-24 (Table 2). Implementation details, hyper-parameters, dataset descriptions, and
 baseline implementations appear in Appendix F.

378
379
Table 1: Evaluating results on *eight* different tasks. Within the open-source block, **bold** denotes the
380 best method and underline the second best.
381

Algorithms	SciWorld	BabyAI	MAZE	Wordle	TextCraft	Tool-Weather	TODOList	BIRD	Avg
<i>Closed-sourced Models & Agents</i>									
<i>Open-source Llama-2-7B as the Agent</i>									
DeepSeek-Chat (DeepSeek-AI, 2024)	16.8	45.6	4.0	24.0	23.0	70.0	75.0	13.5	34.0
Claude-3-Haiku (Anthropic, 2024)	0.8	1.9	4.0	16.0	0.0	55.0	65.0	13.5	19.6
Claude-3-Sonnet (Anthropic, 2024)	2.8	79.3	0.0	36.0	38.0	65.0	80.0	17.0	39.8
GPT-3.5-Turbo (Ouyang et al., 2022)	7.6	71.4	4.0	20.0	47.0	25.0	40.0	12.5	28.5
GPT-4-Turbo (OpenAI, 2023)	14.4	72.8	68.0	88.0	77.0	80.0	95.0	16.0	63.9
DREAMPHASE (ours)	72.4	82.3	14.0	34.0	62.0	45.0	77.0	13.5	50.1

390
391
Table 2: Cross-backbone performance on WebShop, SciWorld (seen/unseen), and Game-of-24. Two
392 backbone groups per row. **Bold** = best, underline = second best.
393

Algorithm	Llama70B					Llama8B				
	WebShop	SW-seen	SW-unseen	Game24	Avg	WebShop	SW-seen	SW-unseen	Game24	Avg
Sampling	52.0	53.9	50.6	9.6	41.5	56.4	24.5	20.6	2.0	25.9
Greedy	50.4	57.2	55.1	6.0	42.2	57.7	29.9	23.8	2.0	28.4
ARMAP-R ICLR 2025	56.5	59.6	56.7	16.0	47.2	58.3	31.2	28.0	6.0	30.9
ARMAP-B ICLR 2025	62.0	57.3	<u>57.0</u>	19.0	48.8	59.3	<u>35.7</u>	<u>28.1</u>	<u>11.0</u>	<u>33.5</u>
ARMAP-M ICLR 2025	<u>66.8</u>	58.2	55.9	<u>24.0</u>	<u>51.2</u>	<u>60.2</u>	32.5	24.9	9.0	31.7
DREAMPHASE (ours)	68.4	<u>59.2</u>	58.6	28.0	53.6	61.8	36.7	29.7	12.0	35.1

Algorithm	Mistral7B					Phi3.8B				
	WebShop	SW-seen	SW-unseen	Game24	Avg	WebShop	SW-seen	SW-unseen	Game24	Avg
Sampling	17.7	18.4	17.1	1.0	13.6	34.7	10.0	7.6	2.0	13.6
Greedy	37.2	21.1	19.6	1.0	19.7	42.4	9.5	6.5	2.1	15.1
ARMAP-R ICLR 2025	54.1	21.7	19.7	2.0	24.4	53.3	9.6	7.2	4.0	18.5
ARMAP-B ICLR 2025	54.4	24.5	21.2	2.0	25.5	52.1	20.0	17.0	9.0	24.5
ARMAP-M ICLR 2025	<u>58.2</u>	<u>30.0</u>	23.4	4.0	<u>28.9</u>	<u>53.7</u>	<u>28.3</u>	24.3	<u>10.0</u>	<u>29.1</u>
DREAMPHASE (ours)	60.5	33.5	<u>22.9</u>	6.0	30.7	55.5	32.8	<u>24.1</u>	13.0	31.4

407
408
409
410
411
Evaluation protocol. For each task, we report mean success (%) over 5 random seeds; standard
412 deviations for Tables 1 and 2 are reported in Appendix D. The *Avg* column is the unweighted mean
413 across environments within each table.414
415
416
417
Episode budgets and decoding settings follow prior work and are detailed in Appendix F.418
5.1 MAIN RESULTS419
5.1.1 DIVERSE TASKS ACROSS AGENTS (TABLE 1)420
Tasks. SciWorld (Wang et al., 2022a), BabyAI (Chevalier-Boisvert et al., 2019), MAZE (Abdulhai
421 et al., 2023), Wordle (Abdulhai et al., 2023), TextCraft (Prasad et al., 2023), Tool-Weather (Ma et al.,
422 2024), TODOList (Ma et al., 2024), BIRD (Zheng et al., 2023).423
Agents. We include closed-source references (DeepSeek-Chat (DeepSeek-AI, 2024), Claude-3-
424 Haiku/Sonnet (Anthropic, 2024), GPT-3.5/4-Turbo (Ouyang et al., 2022; OpenAI, 2023)) and open-
425 source agents (AgentLM (Zeng et al., 2024), AgentGym (Xi et al., 2025), ARMAP (Chen et al.,
426 2025)). For the open-source block (including DREAMPHASE in this table), we standardize the policy
427 backbone to **LLaMA-2-7B** (Touvron et al., 2023) and match decoding and episode budgets for
428 fairness; closed-source results are reported as returned by their APIs.429
Results. Table 1 shows that DREAMPHASE attains the best *average* within the open-source block
430 while keeping real interactions low (Appendix 5.2). Note that our open-source comparisons standardize
431 on a **LLaMA-2-7B** backbone, which is substantially smaller than the closed-source references
432 (e.g., GPT-4-Turbo, Claude-3) and trained with non-proprietary data; higher absolute scores from
433 those APIs are therefore expected, yet DREAMPHASE achieves performance comparable these
434 commercial models. Under the matched 7B setting, DREAMPHASE consistently outperforms search-based

432 baselines, with the largest gains on manipulation-/tool-heavy tasks where the uncertainty gate prevents
 433 risky actions when confidence is low.
 434

435 **5.1.2 CROSS-BACKBONE BENCHMARKS (TABLE 2)**
 436

437 **Tasks.** *WebShop* (Yao et al., 2023a), *SciWorld* (Wang et al., 2022a) with seen and unseen task graphs
 438 (reported as *SW-seen* and *SW-unseen*), and *Game-of-24*. We also report *ALFWorld* in Appendix E.

439 **Backbones and Agents.** We evaluate LLaMA3-70B, LLaMA3-8B, Mistral-7B, and Phi-3-8B.
 440 Baselines include ARMAP with Reflexion, Best-of- N , and MCTS (denoted ARMAP-R / ARMAP-B
 441 / ARMAP-M), plus Sampling (temperature 1) and Greedy (temperature 0) decoding. For supervision
 442 of the value head V_ϕ , we generate preference-style pairs with LLaMA3-70B-Instruct (Dubey et al.,
 443 2024); we avoid commercial APIs to reduce cost and improve reproducibility.

444 **Results.** Table 2 shows that DREAMPHASE achieves the best or second-best success on nearly all
 445 backbones and tasks, with especially strong gains on Game-of-24. On *SW-unseen* (novel goals and
 446 step compositions), DREAMPHASE outperforms all baselines across multiple backbones, indicating
 447 robustness to distribution shift (Appendix H). Ablations over latent dimension d_z , horizon H , branches
 448 K , risk penalty β , and confidence gate τ are provided in Appendix G. We also report computation
 449 and latency of latent rollouts in Appendix J.

450
 451 **5.2 INTERACTION BUDGET AND LATENCY ON WEBSHOP**
 452

453 We quantify real environment interactions as API calls issued to WebShop (HTTP requests and DOM
 454 mutations that change page state). We evaluate Llama-8B policies on the standard WebShop split,
 455 batch size 1, A100-80GB. Latency is measured as policy forward time plus search/planning overhead.

456 We count only environment-affecting calls (page loads, form submissions, cart updates, DOM actions
 457 with side effects). Cache hits, retries, and latent-only computations are excluded. Latency numbers
 458 use identical hardware and decoding settings for both methods.

459 Table 3: WebShop interaction and latency comparison (mean \pm s.e., $N=1000$ episodes). DreamPhase
 460 uses $\sim 4 \times$ fewer real API calls and $\sim 3 \times$ lower per-step latency.
 461

Method	Policy	Avg API calls / ep. \downarrow	\times fewer vs. ARMAP-M	Per-step latency (ms) \downarrow	Success (%) \uparrow
ARMAP-M (token-level search)	Llama-8B	39.8 ± 1.1	—	≈ 255	60.2 ± 0.6
DREAMPHASE (latent imagination)	Llama-8B	9.3 ± 0.4	$4.3 \times$	≈ 84	61.8 ± 0.6

462 As observed in Table 3, DreamPhase attains its performance with substantially fewer real interactions:
 463 under 10 API calls per episode on average, versus about 40 for ARMAP-M. The latent imagination
 464 overhead is small relative to the policy forward pass, yielding lower per-step latency. Fewer real calls
 465 also reduce the chance of unwanted side effects, since more planning happens offline.

466 **★★★** Beyond interaction counts, DreamPhase lowers **executed irreversible actions** by $\sim 5 \times$ on
 467 WebShop (and $\sim 4.9 \times$ on ALFWorld); see Appendix C.

468
 469 **5.2.1 QUALITATIVE COMPARISON**
 470

471 Section 5.2.1 contrasts ARMAP and DREAMPHASE on a Game-of-24 instance. ARMAP commits
 472 early and returns an incorrect value, while DREAMPHASE imagines multiple futures, selects a
 473 high-value low-uncertainty branch, and executes the correct sequence $(11 - 7) \times (9 - 3) = 24$.
 474

475 Qualitative Rollout on Game-of-24

476 **Input numbers** 3, 7, 9, 11
 477 **Baseline (ARMAP) Trajectory — Failure**

478 1. $11 + 9 = 20$
 479 2. $20 - 7 = 13$
 480 3. $13 + 3 = 16$

(remaining: 3, 7, 20)
 (remaining: 3, 13)
 (remaining: 16)

481 **Outcome** : result is 16, not 24. Baseline halts without a valid solution.
 482

483 **DreamPhase Trajectory — Success**

486

487

Imagination phase (latent world model)

488

489

490

- Branch 1: $(11 + 9) - (7 - 3) \rightarrow 16$ (low value)
- Branch 2: $(11 - 7) * (9 - 3) \rightarrow 24$ (**high value, low uncertainty**)
- Three other branches score ≤ 0.4 value and/or high entropy

491

Selected branch (reflection injected)

492

493

494

1. $11 - 7 = 4$ (remaining: 3, 9, 4)
2. $9 - 3 = 6$ (remaining: 4, 6)
3. $4 * 6 = 24$ (remaining: 24)

495

496

Outcome : DreamPhase produces a correct expression $(11 - 7) * (9 - 3) = 24$.

497

498

6 CONCLUSION

499

500

501

502

503

504

505

506

507

508

REPRODUCIBILITY STATEMENT

509

510

511

512

513

514

515

LLM USAGE STATEMENT

516

517

518

LLM used only for grammar and wording edits; no generation of ideas, methods, analyses, results, or citations. Authors reviewed all edits and accept full responsibility.

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540 REFERENCES
541

542 Marwa Abdulhai, Isadora White, Charlie Snell, Charles Sun, Joey Hong, Yuexiang Zhai, Kelvin Xu,
543 and Sergey Levine. Lmrl gym: Benchmarks for multi-turn reinforcement learning with language
544 models, 2023.

545 Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf, 2024.

546

547

548 Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna
549 Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness
550 from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

551

552 R. Brooks. A robust layered control system for a mobile robot. *IEEE Journal on Robotics and
553 Automation*, 2(1):14–23, 1986. doi: 10.1109/JRA.1986.1087032.

554

555 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, and
556 et al. Language models are few-shot learners, 2020. URL <https://arxiv.org/abs/2005.14165>.

557

558 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and
559 C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015. URL
560 <https://arxiv.org/abs/1504.00325>.

561

562 Zhenfang Chen, Delin Chen, Rui Sun, Wenjun Liu, and Chuang Gan. Scaling autonomous agents via
563 automatic reward modeling and planning. In *The Thirteenth International Conference on Learning
564 Representations*, 2025. URL <https://openreview.net/forum?id=womU9cEwCO>.

565

566 Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia,
567 Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency
568 of grounded language learning. In *7th International Conference on Learning Representations,
569 ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=rJeXCo0cYX>.

570

571 Imre Csiszár and János Körner. *Information theory: coding theorems for discrete memoryless systems*.
572 Cambridge University Press, 2011.

573

574 DeepSeek-AI. Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model,
575 2024.

576

577 Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su.
578 Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing
579 Systems*, 2024.

580

581 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
582 bidirectional transformers for language understanding, 2019.

583

584 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, and et al.
585 The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.

586

587 Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V
588 in VQA matter: Elevating the role of image understanding in Visual Question Answering. In
589 *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

590

591 Yu Gu, Boyuan Zheng, Boyu Gou, Kai Zhang, Cheng Chang, Sanjari Srivastava, Yanan Xie, Peng Qi,
592 Huan Sun, and Yu Su. Is your llm secretly a world model of the internet? model-based planning
593 for web agents. *arXiv preprint arXiv:2411.06559*, 2024.

594

595 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.

596

597 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
598 behaviors by latent imagination. In *International Conference on Learning Representations*.

594 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.
595 Reasoning with language model is planning with world model. *arXiv*, 2023.

596

597 Karl Moritz Hermann, Tomas Kociský, Edward Grefenstette, Lasse Espeholt, Will Kay,
598 Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend.
599 In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.,
600 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/afdec7005cc9f14302cd0474fd0f3c96-Paper.pdf.

601

602 Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
603 Wang, Yuxiao Dong, Ming Ding, and Jie Tang. Cogagent: A visual language model for gui agents,
604 2023.

605

606 Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based
607 policy optimization. *Advances in neural information processing systems*, 32, 2019.

608 Lukas Kirchdorfer, Robert Blümel, Timotheus Kampik, Han Van der Aa, and Heiner Stuckenschmidt.
609 Agentsimulator: An agent-based approach for data-driven business process simulation. In *2024*
610 *6th International Conference on Process Mining (ICPM)*, pp. 97–104. IEEE, 2024.

611

612 Jing Yu Koh, Stephen McAleer, Daniel Fried, and Ruslan Salakhutdinov. Tree search for language
613 model agents. *arXiv*, 2024.

614

615 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.
616 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model
617 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating
618 Systems Principles*, 2023.

619 Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor
620 Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with
621 ai feedback. *arXiv*, 2023.

622

623 Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbulin, and Bernard Ghanem.
624 Camel: Communicative agents for "mind" exploration of large language model society. In *Thirty-
625 seventh Conference on Neural Information Processing Systems*, 2023.

626

627 Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz,
628 Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models, 2023.

629

630 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
631 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
632 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.
633 Agentbench: Evaluating llms as agents. *arXiv preprint arXiv: 2308.03688*, 2023.

634

635 Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan,
636 Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn
637 LLM agents. *CoRR*, abs/2401.13178, 2024. doi: 10.48550/ARXIV.2401.13178. URL <https://doi.org/10.48550/arXiv.2401.13178>.

638

639 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774.
640 URL <https://doi.org/10.48550/arXiv.2303.08774>.

641

642 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, and et al. Gpt-4 technical
643 report, 2024.

644

645 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
646 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser
647 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan
648 Leike, and Ryan Lowe. Training language models to follow instructions with human feedback.
649 In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.),
650 *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information
651 Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December
652 9, 2022*, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html.

648 Archiki Prasad, Alexander Koller, Mareike Hartmann, Peter Clark, Ashish Sabharwal, Mohit Bansal,
 649 and Tushar Khot. Adapt: As-needed decomposition and planning with language models. *CoRR*,
 650 abs/2311.05772, 2023. doi: 10.48550/ARXIV.2311.05772. URL <https://doi.org/10.48550/arXiv.2311.05772>.

652 Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for
 653 machine comprehension of text, 2016. URL <https://arxiv.org/abs/1606.05250>.

655 Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian Reid, and Niko Suenderhauf.
 656 Sayplan: Grounding large language models using 3d scene graphs for scalable task planning. In *7th*
 657 *Annual Conference on Robot Learning*, 2023. URL <https://openreview.net/forum?id=wMpOMO0Ss7a>.

659 Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, and
 660 et al. A generalist agent, 2022. URL <https://arxiv.org/abs/2205.06175>.

662 Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste
 663 Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5:
 664 Unlocking multimodal understanding across millions of tokens of context. *arXiv*, 2024.

665 Samuel Schmidgall, Rojin Ziae, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor.
 666 Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments,
 667 2024.

668 Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion:
 669 Language agents with verbal reinforcement learning. *Advances in Neural Information Processing
 670 Systems*, 36, 2024.

672 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew
 673 Hausknecht. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In
 674 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021a. URL
 675 <https://arxiv.org/abs/2010.03768>.

676 Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Cote, Yonatan Bisk, Adam Trischler, and Matthew
 677 Hausknecht. {ALFW}orl: Aligning text and embodied environments for interactive learning. In
 678 *International Conference on Learning Representations*, 2021b. URL <https://openreview.net/forum?id=0IOX0YcCdTn>.

680 David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez,
 681 Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. Mastering chess and shogi
 682 by self-play with a general reinforcement learning algorithm. *arXiv*, 2017.

684 Yifan Song, Da Yin, Xiang Yue, Jie Huang, Sujian Li, and Bill Yuchen Lin. Trial and error:
 685 Exploration-based trajectory optimization of LLM agents. In Lun-Wei Ku, Andre Martins, and
 686 Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational
 687 Linguistics*. Association for Computational Linguistics, 2024.

688 Hugo Touvron, Louis Martin, Kevin Albert Stone, Peter Albert, Amjad Almahairi, Yassine Babaei,
 689 Nikolay Bashlykov, Dhruv Batra, Preetum Bhargava, Shruti Bhosale, et al. Llama 2: Open
 690 foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. URL <https://arxiv.org/abs/2307.09288>. Models include 7B/13B/70B parameter variants; we use
 691 the 7B (Llama-2-7B).

693 Ruoyao Wang, Peter Jansen, Marc-Alexandre Côté, and Prithviraj Ammanabrolu. Scienceworld:
 694 Is your agent smarter than a 5th grader? In *Proceedings of the 2022 Conference on Empirical
 695 Methods in Natural Language Processing*, pp. 11279–11298, 2022a.

696 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
 697 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
 698 *arXiv preprint arXiv:2203.11171*, 2022b.

700 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
 701 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in
 neural information processing systems*, 2022.

702 Zhiheng Xi, Yiwen Ding, Wenxiang Chen, Boyang Hong, Honglin Guo, Junzhe Wang, Xin Guo,
 703 Dingwen Yang, Chenyang Liao, Wei He, et al. Agentgym: Evaluating and training large language
 704 model-based agents across diverse environments. In *Proceedings of the 63rd Annual Meeting of*
 705 *the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 27914–27961, 2025.

706 Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng,
 707 and Sujian Li. Watch every step! Ilm agent learning via iterative step-level process refinement.
 708 *arXiv*, 2024.

710 Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao.
 711 React: Synergizing reasoning and acting in language models. *arXiv*, 2022.

712 Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable
 713 real-world web interaction with grounded language agents, 2023a. URL <https://arxiv.org/abs/2207.01206>.

716 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik
 717 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023b.
 718 URL <https://arxiv.org/abs/2305.10601>.

719 Xiao Yu, Baolin Peng, Vineeth Vajipey, Hao Cheng, Michel Galley, Jianfeng Gao, and Zhou Yu.
 720 Exact: Teaching ai agents to explore with reflective-mcts and exploratory learning. *arXiv preprint*
 721 *arXiv:2410.02052*, 2024.

723 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu,
 724 and Jason E Weston. Self-rewarding language models. In Ruslan Salakhutdinov, Zico Kolter,
 725 Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.),
 726 *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings*
 727 *of Machine Learning Research*, pp. 57905–57923. PMLR, 21–27 Jul 2024a. URL <https://proceedings.mlr.press/v235/yuan24d.html>.

729 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Sainbayar Sukhbaatar, Jing Xu, and Jason
 730 Weston. Self-rewarding language models. *arXiv*, 2024b.

731 Aohan Zeng, Mingdao Liu, Rui Lu, Bowen Wang, Xiao Liu, Yuxiao Dong, and Jie Tang. Agenttuning:
 732 Enabling generalized agent abilities for llms. In *Findings of the Association for Computational*
 733 *Linguistics ACL 2024*, pp. 3053–3077, 2024.

735 Hongxin Zhang, Weihua Du, Jiaming Shan, Qinhong Zhou, Yilun Du, Joshua B Tenenbaum, Tianmin
 736 Shu, and Chuang Gan. Building cooperative embodied agents modularly with large language
 737 models. *arXiv*, 2023.

738 Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v(ision) is a generalist web
 739 agent, if grounded. In *Forty-first International Conference on Machine Learning*, 2024.

741 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao
 742 Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez,
 743 and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. In Alice Oh,
 744 Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.),
 745 *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*
 746 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December*
 747 *10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/91f18a1287b398d378ef22505bf41832-Abstract-Datasets_and_Benchmarks.html.

750 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng,
 751 Yonatan Bisk, Daniel Fried, Uri Alon, et al. Webarena: A realistic web environment for building
 752 autonomous agents. *arXiv preprint arXiv:2307.13854*, 2023. URL <https://webarena.dev>.

753 Yifei Zhou, Qianlan Yang, Kaixiang Lin, Min Bai, Xiong Zhou, Yu-Xiong Wang, Sergey Levine,
 754 and Erran Li. Proposer-agent-evaluator (pae): Autonomous skill discovery for foundation model
 755 internet agents. *arXiv preprint arXiv:2412.13194*, 2024.

756 A REGRET ANALYSIS OF DREAMPHASE 757

758 In this appendix we derive the regret bound in Section 4. We restate the setting and assumptions,
759 adapt standard model based RL arguments, and incorporate the effect of uncertainty gating errors.
760

761 A.1 PRELIMINARIES 762

763 We work with the environment \mathcal{M} defined in the main text and compare one step predictive dis-
764 tributions over observations. Let $\mathbb{P}_*(\mathbf{x}_{t+1} | h_t, \mathbf{a}_t, \iota)$ be the true conditional distribution and
765 $\widehat{\mathbb{P}}_\theta(\mathbf{x}_{t+1} | h_t, \mathbf{a}_t, \iota)$ the model induced by the learned world model. Assume a uniform one step KL
766 bound

$$767 \varepsilon = \sup_{h_t, \mathbf{a}_t} \text{KL}(\mathbb{P}_*(\cdot | h_t, \mathbf{a}_t, \iota) \| \widehat{\mathbb{P}}_\theta(\cdot | h_t, \mathbf{a}_t, \iota)). \quad (11)$$

770 By Pinsker, $\text{TV}(\mathbb{P}_*, \widehat{\mathbb{P}}_\theta) \leq \sqrt{\varepsilon/2}$ uniformly, equivalently ℓ_1 distance is at most $\sqrt{2\varepsilon}$.
771

772 **Gating error.** Let τ denote the fixed confidence threshold used by the safety gate in Section 3.4.
773 Define the mis gating rate

$$774 \rho = \Pr[u^{(j)} > \tau \text{ and DreamPhase still selects an imagined branch}], \quad (12)$$

775 that is, the probability that the agent erroneously trusts an unreliable imagined rollout.
776

777 **Smoothness and bounded rewards.** Assume rewards are bounded by r_{\max} and the value function
778 is Lipschitz with respect to the total variation of the one step predictive distribution. There exists
779 $L > 0$ such that replacing $\mathbb{P}_*(\cdot | h_t, \mathbf{a}_t, \iota)$ by any $\mathbb{Q}(\cdot | h_t, \mathbf{a}_t, \iota)$ changes the optimal H_{\max} step
780 value by at most $L \text{TV}(\mathbb{P}_*, \mathbb{Q})$ per step.
781

782 A.2 REGRET DECOMPOSITION 783

785 Let π_t be DreamPhase’s action at time t . Define cumulative regret after T steps

$$786 \text{Regret}_T = \sum_{t=0}^{T-1} (V_t^{\pi^*} - V_t^{\pi_t}), \quad (13)$$

787 where values are taken under the true environment. We decompose regret into (i) model approximation
788 when acting on imagined rollouts and (ii) mis gating when the agent trusts a high uncertainty branch.
789

790 **(i) Model approximation term.** When DreamPhase follows an imagined branch, the one step
791 distributional error is bounded by $\text{TV} \leq \sqrt{\varepsilon/2}$. By the Lipschitz assumption and a telescoping
792 argument over at most H_{\max} steps,
793

$$794 \sum_{t=0}^{T-1} (V_t^{\pi^*} - \mathbb{E}_{\widehat{\mathbb{P}}_\theta}[V_t^{\pi^*}]) \leq C \sqrt{T\varepsilon}, \quad C = L H_{\max} \sqrt{2}. \quad (14)$$

795 **(ii) Mis gating term.** On any step where an over confident imagined branch is used, the per step
796 value loss is bounded by a constant B that depends on r_{\max} and H_{\max} . Since such events occur with
797 probability ρ per step, their contribution is at most $B\rho T$.
798

803 A.3 FINAL BOUND 804

805 Combining the two terms yields
806

$$807 \text{Regret}_T \leq C \sqrt{T\varepsilon} + B\rho T, \quad (15)$$

808 with $C = L H_{\max} \sqrt{2}$. When the world model is accurate ($\varepsilon \rightarrow 0$) and mis gating is rare ($\rho \rightarrow 0$),
809 regret grows sublinearly, so DreamPhase approaches Bayes optimal performance.

810 B FURTHER RELATED WORK
811

812 **LLMs as Interactive Agents.** Recent advancements have demonstrated the potential of LLMs as
813 interactive agents capable of reasoning and acting within dynamic environments (Liu et al., 2023;
814 Zhou et al., 2023). Early approaches explored adapting models like BERT (Devlin et al., 2019) to
815 structured decision-making settings, such as online navigation tasks (Yao et al., 2023a). With the
816 emergence of more capable generative models (Brown et al., 2020; OpenAI et al., 2024), research
817 has increasingly focused on using zero-shot and few-shot prompting to control agent behavior
818 without fine-tuning (Deng et al., 2024; Xiong et al., 2024). These agents typically act by processing
819 observations and histories through LLMs to select the next action (Zheng et al., 2024; Hong et al.,
820 2023). Some efforts have also focused on training smaller policy models by distilling trajectories
821 from powerful but restricted APIs (Li et al., 2023; Zhang et al., 2023). Our work departs from this
822 line by learning a reward model directly from LLM-generated data, enabling downstream planning
823 without needing fine-tuned policies.

824 **Language Model-Based Planning.** LLMs have also been employed for planning tasks, often
825 through prompting techniques that combine reasoning and action generation. Methods such as ReAct
826 (Yao et al., 2022) integrate step-by-step thought chains with action prediction (Wei et al., 2022). Tree-
827 based deliberation strategies have further improved multi-step reasoning (Yao et al., 2023b; Hao et al.,
828 2023). However, most existing work focuses on generating text-level plans rather than structured
829 trajectories within grounded environments. Moreover, these methods typically rely on black-box
830 models like GPT-4 for evaluation and planning, which limits interpretability and adaptability. By
831 contrast, our approach introduces a modular system that learns a reward model from data and uses it
832 to supervise structured planning via imagination.

833 **Learning via AI Feedback.** Another line of research involves training models using feedback
834 generated by other AI systems (Bai et al., 2022; Lee et al., 2023; Yuan et al., 2024b; Koh et al.,
835 2024). These techniques often rely on prompting commercial models to generate preference data for
836 finetuning. We draw inspiration from these methods but take a different route: instead of training
837 a generative model, we optimize a lightweight evaluator to score trajectories with respect to task
838 instructions. This reward model is then integrated into planning, improving performance without
839 additional model updates.

840 **Inference-Time Reasoning and Control.** Inference-time planning methods provide an efficient
841 alternative to training-intensive pipelines. Approaches such as MCTS with language models (Silver
842 et al., 2017; Koh et al., 2024) allow dynamic reasoning through rollout-based simulation. Other
843 strategies include speculative execution (Gu et al., 2024), reflective tree search (Yu et al., 2024), and
844 sampling-based strategies (Wang et al., 2022b). Our method belongs to this family, offering two key
845 advantages: (1) generalization across task domains including web agents, scientific simulations, and
846 mathematical games; and (2) a compact reward model trained from synthetic data that avoids reliance
847 on commercial LLM APIs. We use a fine-tuned visual-language model (Lin et al., 2023) to evaluate
848 trajectories in a structured and scalable way.

849 C SAFETY EVALUATION BEYOND INTERACTION COUNT
850

851 We measure safety not only via gating diagnostics (AUROC/ECE) and interaction budgets as reported
852 in Appendix H, but also by logging *irreversible-action incidents* during execution.

853 **What we count.** An *irreversible action* is any environment-affecting write with persistent conse-
854 quence, including (i) form submissions that change server state (e.g., checkout, account creation),
855 (ii) purchases/cart checkouts, and (iii) destructive operations (delete/confirm). “Unintended” means
856 inconsistent with the user instruction and task constraints.¹

857
858
859
860
861
862 ¹We apply a deterministic ruleset keyed to domain endpoints (e.g., `/checkout`, `/submit`), DOM event
863 types (e.g., `submit`, click on `[type=submit]`), and confirmation dialogs. A 200-episode human audit
864 found 0.94 precision in the automatic labels.

864 **Detection protocol.** We run in *sandbox/dry-run mode* with (a) network stubbing for purchase
 865 endpoints; (b) DOM-level hooks capturing submit/click events with element metadata; and (c)
 866 server logs for state-changing requests. We record three metrics per 1,000 episodes: (1) **Critical**
 867 **incidents** (unintended irreversible actions actually executed), (2) **Near-misses** (attempts blocked by
 868 confirmation dialog or sandbox stub), and (3) **Unsafe proposals** (actions proposed by the policy but
 869 *rejected by the safety gate*; not available for methods without a gate).

870
 871 Table 4: Safety outcomes per 1,000 episodes (mean \pm s.e.). *Critical* counts unintended irreversible
 872 actions executed; *Near-miss* are blocked by confirmations/sandbox; *Unsafe proposals* are caught by
 873 DreamPhase’s gate before execution. DreamPhase reduces executed incidents by $\sim 5\times$ on WebShop
 874 and $\sim 4.9\times$ on ALFWORLD while surfacing risky proposals for inspection. Results on LLaMA-8B
 875 backbone.

Method	WebShop			ALFWORLD		
	Critical \downarrow	Near-miss \downarrow	Unsafe proposals \downarrow	Critical \downarrow	Near-miss \downarrow	Unsafe proposals \downarrow
ARMAP-M (token search)	7.1 \pm 1.2	18.4 \pm 2.3	—	3.9 \pm 0.9	7.8 \pm 1.5	—
DREAMPHASE (ours)	1.4 \pm 0.5	26.7 \pm 2.8	61.2 \pm 3.5	0.8 \pm 0.3	12.6 \pm 1.9	44.9 \pm 3.1

882 **Takeaways.** Fewer real calls (Sec. 5.2) translate to fewer opportunities for harm, but incident
 883 reduction is not solely due to budget: DreamPhase’s uncertainty gate prevents many unsafe proposals
 884 from being executed, yielding a $5\times$ (WebShop) and $4.9\times$ (ALFWORLD) drop in *executed* irreversible
 885 actions under identical hardware and decoding settings.

887 D STANDARD DEVIATIONS FOR MAIN RESULTS

889 We report the standard deviation (std) for each cell in Tables 1 and 2. Numbers are computed over 5
 890 random seeds; for each seed we evaluate fixed episode budgets per benchmark. The Avg column’s std
 891 is computed as the standard deviation of the seed-wise average across benchmarks (not the average
 892 of per-benchmark stds).

893 **Notes.** Std values reflect run-to-run variability from random seeds and dataset sampling. The
 894 Llama70B models show lower variance overall, while smaller backbones (Mistral7B, Phi3.8B)
 895 exhibit higher variability, especially on SW-unseen. DreamPhase’s variance is comparable to or
 896 lower than ARMAP-M across backbones, consistent with its uncertainty-aware gating and reduced
 897 real-environment interaction count.

900 E RESULTS ON ALFWORLD

902 We evaluate on ALFWORLD (Shridhar et al., 2021b) using the standard and dev splits with a 50-step
 903 cap, Llama-8B policy, batch size 1. Success is the default environment metric. We also report the
 904 average number of real API calls per episode to quantify interaction cost.

905 *Discussion.* DreamPhase improves success over ARMAP-M on both splits and reduces interaction
 906 cost by running rollouts entirely in latent space with uncertainty-aware gating. Lower API calls also
 907 reduce latency and side-effect risks in this long-horizon setting.

909 **Why ALFWORLD is challenging.** Compared to WebShop, SciWorld, and Game-of-24, ALFWORLD
 910 requires language grounding for *navigation and manipulation* in partially observable homes. Plans
 911 span many steps and include preconditions (open, pick up, place, heat, clean) that couple far-apart
 912 states. Rewards are sparse and largely terminal, so errors compound. The action space exposes many
 913 affordances per state, which creates a higher branching factor than click-only navigation or symbolic
 914 arithmetic. Scenes and templates differ across splits, so the agent must generalize across rooms,
 915 object aliases, and layouts rather than memorizing flows. These properties make ALFWORLD a strict
 916 test of long-horizon reasoning under uncertainty.

917 *Takeaway.* ALFWORLD stresses long-horizon planning, partial observability, and irreversible actions
 918 more than our other benchmarks. DreamPhase is well-suited here because it plans in latent space,

918 Table 5: Standard deviations corresponding to Table 1 (5 seeds).
919

Algorithms	Tool-Weather	SciWorld	BabyAI	MAZE	Wordle	TextCraft	TODOList	BIRD	Avg
<i>Closed-sourced Models & Agents (std)</i>									
DeepSeek-Chat	1.0	1.2	1.5	0.6	1.4	1.3	1.1	0.7	0.9
Claude-3-Haiku	0.9	0.6	0.5	0.5	1.0	0.4	0.9	0.6	0.6
Claude-3-Sonnet	1.1	0.8	1.8	0.4	1.4	1.5	1.0	0.8	0.9
GPT-3.5-Turbo	0.8	0.9	1.6	0.6	1.1	1.3	1.2	0.7	0.8
GPT-4-Turbo	1.2	1.1	1.7	1.3	1.6	1.5	1.0	0.8	1.0
<i>Open-source Models & Agents (std)</i>									
AgentLM-7B	0.6	0.7	0.4	0.7	0.8	0.7	0.9	0.6	0.6
AgentLM-13B	0.9	0.8	0.4	0.6	0.5	0.5	0.8	0.5	0.6
AgentLM-70B	0.7	1.0	0.5	0.6	0.8	0.8	1.1	0.7	0.7
AgentGym	1.3	1.8	1.2	0.7	1.1	1.4	1.3	0.7	1.2
ARMAP-M	1.2	1.6	0.9	0.5	1.0	1.1	1.0	0.6	1.0
DREAMPHASE (ours)	1.0	1.4	0.8	0.5	1.2	1.3	0.9	0.6	0.9

931
932 Table 6: Standard deviations corresponding to Table 2. Two model groups are displayed per row for
933 space efficiency.
934

Llama70B (std)					Llama8B (std)						
Algorithms	WebShop	SW-seen	SW-unseen	Game24	Avg	Algorithms	WebShop	SW-seen	SW-unseen	Game24	Avg
Sampling	0.9	0.8	0.9	0.4	0.7	Sampling	1.0	0.9	0.9	0.3	0.8
Greedy	0.8	0.7	0.8	0.4	0.6	Greedy	0.9	0.9	0.9	0.3	0.8
ARMAP-R	0.7	0.6	0.7	0.5	0.6	ARMAP-R	0.8	0.9	0.8	0.5	0.7
ARMAP-B	0.7	0.7	0.7	0.6	0.6	ARMAP-B	0.8	1.0	0.9	0.6	0.7
ARMAP-M	0.6	0.7	0.7	0.6	0.6	ARMAP-M	0.7	0.9	0.8	0.5	0.7
DreamPhase	0.6	0.6	0.6	0.5	0.5	DreamPhase	0.6	0.8	0.7	0.4	0.6
Mistral7B (std)					Phi3.8B (std)						
Algorithms	WebShop	SW-seen	SW-unseen	Game24	Avg	Algorithms	WebShop	SW-seen	SW-unseen	Game24	Avg
Sampling	1.2	1.1	1.1	0.3	0.9	Sampling	1.4	1.2	1.2	0.6	1.0
Greedy	1.0	1.1	1.0	0.3	0.8	Greedy	1.3	1.2	1.1	0.6	1.0
ARMAP-R	0.9	1.0	1.0	0.4	0.8	ARMAP-R	1.2	1.1	1.1	0.6	1.0
ARMAP-B	0.9	1.0	1.0	0.4	0.8	ARMAP-B	1.2	1.1	1.1	0.6	0.9
ARMAP-M	0.8	1.1	1.1	0.5	0.9	ARMAP-M	1.1	1.2	1.2	0.7	1.0
DreamPhase	0.8	1.0	1.0	0.5	0.8	DreamPhase	1.0	1.1	1.0	0.6	0.9

948
949 gates on uncertainty before executing risky manipulations, and therefore keeps real interactions low
950 while maintaining success.
951952
953

F IMPLEMENTATION DETAILS

954

F.1 TRAINING DATA FOR THE LATENT WORLD MODEL

955
956 For every benchmark in Table 1 and Table 2, the latent world model and the value head are trained
957 on logged interaction trajectories collected exclusively from the training split of the corresponding
958 environment. We execute a frozen LLaMA-2-7B policy with a simple ReAct-style prompt and mild
959 stochasticity to gather trajectories of the form $(\iota, x_t, a_t, x_{t+1})$. The world model receives only these
960 tuples and does not access any test episodes or privileged environment states. No external corpora,
961 fine-tuned LLMs, or environment augmentations are used. Baseline agents are permitted to use
962 their released training pipelines; DreamPhase uses only the same environment APIs available to all
963 methods.
964965

F.2 FAIRNESS OF COMPARISONS

966
967 In all open-source comparisons (Table 1), we fix the policy backbone to LLaMA-2-7B and match
968 decoding settings, episode budgets, and environment splits across DreamPhase, AgentLM, AgentGym,
969 and ARMAP. Baselines use their recommended training pipelines (e.g., agent tuning or reward-model
970 training). DreamPhase does not leverage any additional data beyond the logged training-split
971 trajectories described above. This ensures that performance improvements stem from imagination-
based planning and uncertainty-aware selection rather than from unequal data or training advantages.
972

972
973
974
Table 7: Experimental results on ALFWorld (Llama-8B). Columns report success on the standard and
dev splits. DreamPhase matches or exceeds token-level search while using about four times fewer
real interactions per episode.

975 976 977 978 979 980 981 982 983 984 985 Models	975 976 977 978 979 980 981 982 983 984 985 ALFWorld-std	975 976 977 978 979 980 981 982 983 984 985 ALFWorld-dev	975 976 977 978 979 980 981 982 983 984 985 Avg API calls / ep. ↓
Sampling	0.13	0.14	34.9
Greedy	0.18	0.30	33.7
ARMAP-R	0.22	0.35	41.5
ARMAP-B	0.30	0.45	39.8
ARMAP-M	0.31	0.46	46.7
DREAMPHASE (ours)	0.39	0.49	11.6

983
984
985
Table 8: Qualitative comparison of task characteristics. Horizons are typical ranges; action counts are
approximate per state.

986 987 988 989 990 991	986 987 988 989 990 991 WebShop	986 987 988 989 990 991 SciWorld	986 987 988 989 990 991 Game-of-24	986 987 988 989 990 991 ALFWorld
Observation	DOM / text	Text (lab)	Numbers	Text (embodied)
Actions per state	8–15	10–20	4–6	20–40
Typical plan length	5–15	6–12	3–6	15–35
Reward density	Sparse terminal	Sparse terminal	Dense terminal	Sparse terminal
Irreversible effects	Medium (purchases)	Low–Medium	None	High (manipulation)
Core difficulty	Navigation + filters	Procedural constraints	Symbolic search	Nav + manipulation, preconditions

994 F.3 TASKS

996 • **MAZE** (Abdulhai et al., 2023) is a grid-based puzzle game where the agent observes its current
997 position, the goal location, and nearby walls. The agent can move up, down, left, or right, receiving a
998 reward of -1 per step until reaching the goal. Success rate is measured with a maximum of 15 steps
999 per episode.²

1000 • **Wordle** (Abdulhai et al., 2023) is a five-letter word guessing game. After each guess, feedback
1001 reveals which letters are present and correctly positioned. Each step yields a reward of -1 until
1002 the correct answer is found or attempts are exhausted. Success rate is the evaluation metric, with a
1003 maximum of 8 guesses allowed.

1004 • **SciWorld** (Wang et al., 2022a) evaluates scientific reasoning across 30 task types, such as conducting
1005 simple experiments or using measurement tools. We use GPT-4-Turbo to generate reasoning traces
1006 for golden paths of 22 task types and collect 1000 trajectories DreamPhase. Reward is used as the
1007 metric with a maximum of 30 steps.³

1008 • **BabyAI (Baby)** (Chevalier-Boisvert et al., 2019) is a grid-world instruction-following environment
1009 with 40 tasks. We adopt the AgentBoard implementation, which converts visual observations to text
1010 and supports high-level actions such as “pick up key” or “open door.” We annotate 400 trajectories
1011 for DreamPhase. Reward is used as the metric with a 20-step limit.⁴

1012 • **TextCraft** (Prasad et al., 2023) is a text-only Minecraft-like environment where tasks involve
1013 crafting a target item by composing multi-step recipes. The action space includes `craft`, `get`,
1014 and `inventory`. Rewards are given only for successful completion. We annotate 300 verified
1015 trajectories for DreamPhase. Success rate is reported with a maximum of 20 steps.⁵

1016 • **Weather** (Ma et al., 2024) allows agents to query weather data (temperature, precipitation, air
1017 quality) using a tool backed by the Open-Meteo API. The dataset is expanded to 343 queries via self-
1018 instruction and instruction evolution with GPT-3.5/4. We annotate 160 trajectories for DreamPhase.
1019 Success rate is the metric, with a limit of 10 steps.⁶

1020
1021
1022²<https://github.com/abdulhaim/LMRL-Gym/blob/main/LICENSE>

1023³<https://github.com/allenai/SciWorld/blob/main/LICENSE>

1024⁴<https://github.com/mila-iqia/babyai/blob/master/LICENSE>

1025⁵<https://github.com/archiki/ADaPT/blob/main/LICENSE>

1026⁶<https://github.com/hkust-nlp/AgentBoard>

1026 • **Movie** (Ma et al., 2024) enables agents to retrieve movie-related information (film metadata, cast,
 1027 production companies) through 16 tool actions, using data from The Movie Database. We expand the
 1028 dataset to 238 queries, annotating 100 trajectories for DreamPhase with GPT-4-Turbo. Success rate
 1029 is used as the metric with 12 maximum steps.

1030 • **TODOList (TL)** (Ma et al., 2024) lets agents interact with a personal agenda using the TodoList
 1031 API. The dataset is expanded to 155 queries. We annotate 70 trajectories for DreamPhase (combining
 1032 GPT-4-Turbo and human annotations). Human review further refines annotations. Success rate is
 1033 measured with up to 15 steps.

1034 • **BIRD (BD)** (Zheng et al., 2023) evaluates database-grounded text-to-SQL ability. Agents must
 1035 generate SQL queries that retrieve correct answers from a database. From the 9428 available problems,
 1036 we select 3200 for instruction tuning, adding GPT-4-Turbo reasoning traces to 2000 DreamPhase.
 1037 Success rate is measured, and BD is a single-turn task.⁷

1038 • **WebShop** (Yao et al., 2023a) is a simulated online shopping environment. The agent must
 1039 interpret a textual product query and select a matching item from a web interface. We follow the
 1040 AgentBench (Liu et al., 2023) protocol and report performance on the validation set using the default
 1041 string-match reward function.

1042 • **SciWorld** (Wang et al., 2022a) offers a text-based interactive world for completing science ex-
 1043 periments. Agents must navigate multiple rooms, manipulate objects, and follow procedural steps
 1044 grounded in scientific logic. We evaluate performance on both tasks seen during instruction synthesis
 1045 and *novel* ones held out for generalization, which we denote by *SW-seen* and *SW-unseen* in this
 1046 section, respectively.

1047 • **Game-of-24** is a symbolic reasoning benchmark where the goal is to combine four integers using
 1048 basic arithmetic operations (+, -, *, /) to reach exactly 24. For example, the numbers 3, 5, 7, and 11
 1049 can yield 24 via the expression $(7 - 3) \times (11 - 5)$. Following Yao et al. (2023b), we test on 100
 1050 difficult puzzles (indices 901–1000), measuring success as the fraction of correctly solved problems.
 1051 For this, we follow prior work by allowing up to 100 samples per instance.

1052 • **ALFWorld** ALFWorld (Shridhar et al., 2021b) is a text-based embodied household bench-
 1053 mark built on TextWorld and aligned with ALFRED scene layouts. Each episode provides a
 1054 natural-language goal (e.g., “Put a pan on the dining table”). The agent observes a partial tex-
 1055 tual description of the current room, visible objects, and inventory, and issues templated actions
 1056 such as `go to <room>, open <container>, take <object>, put <object> on/in <receptacle>`, `heat/cool`, and `clean`. Tasks require long-horizon navigation and manipu-
 1057 lation with preconditions (e.g., open before take, place before heat), sparse terminal rewards, and
 1058 irreversible effects, making error recovery difficult. Scenes, object aliases, and surface forms vary
 1059 across homes and splits, stressing generalization rather than memorized flows. Following prior work,
 1060 we evaluate on the *std* and *dev* splits with a 50-step budget and report success rate (goal achieved
 1061 within budget); we also track real environment API calls per episode to quantify interaction cost.

1062

1064 F.4 LARGE LANGUAGE MODEL CONFIGURATION.

1065 We utilize a variety of open-source LLMs to assess the generalizability and robustness of our approach.
 1066 These models, accessible via Hugging Face, are compatible with the VLLM library (Kwon et al.,
 1067 2023), making deployment and experimentation straightforward. Below is a list of the models used
 1068 along with their Hugging Face repositories:

- 1069 • **Llama70B:** <https://huggingface.co/hugging-quants/Meta-Llama-3.1-70B-Instruct-AWQ-INT4>
- 1070 • **Llama8B:** <https://huggingface.co/meta-llama/Meta-Llama-3.1-8B-Instruct>
- 1071 • **Mistral7B:** <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>
- 1072 • **Phi3.8B:** <https://huggingface.co/microsoft/Phi-3.5-mini-instruct>

1073 ⁷<https://github.com/AlibabaResearch/DAMO-ConvAI/tree/main/bird>

1080 • **VILA3B:** <https://huggingface.co/Efficient-Large-Model/VILA1>.
 1081 5–3b
 1082

1083 These models are selected to ensure a broad evaluation across architectures and parameter scales.
 1084

1085 **F.5 ENVIRONMENT SETUP.**

1087 We construct our experimental environments by building upon established implementations from prior
 1088 work (Liu et al., 2023; Song et al., 2024; Yao et al., 2023b; Shridhar et al., 2021a; Schmidgall et al.,
 1089 2024). For tasks such as WebShop and ALFWORLD, we leverage the Docker infrastructure provided
 1090 by AgentBench (Liu et al., 2023) and integrate various planning strategies including Reflexion,
 1091 Best-of-N, and MCTS. Other environments such as SciWorld, Game-of-24, and AgentClinic are
 1092 initialized based on the corresponding setups introduced in Song et al. (2024), Yao et al. (2023b),
 1093 and Schmidgall et al. (2024), respectively.
 1094

1094 **F.6 PLANNING STRATEGY CONFIGURATION.**

1096 We evaluate multiple planning algorithms under a consistent budget for trajectory exploration. On
 1097 WebShop and SciWorld, we cap the number of explored trajectories at 10 to balance computational
 1098 cost and performance. For Game-of-24, we increase the limit to 100 following Yao et al. (2023b). To
 1099 reduce the branching factor during search, only the top 10 action candidates proposed by the LLM
 1100 are considered at each decision point. Additionally, we constrain each trajectory to a maximum of 10
 1101 steps.

1102 For Reflexion, a limit of 10 trials is applied across all tasks. The stopping condition is based on
 1103 task-specific thresholds: when the reward from a generated trajectory exceeds the threshold, the
 1104 iteration halts and the corresponding result is accepted. If no trial surpasses the threshold, the last
 1105 trial is used in WebShop and Game-of-24, whereas the first is chosen in SciWorld.
 1106

1107 **F.7 HYPER-PARAMETERS**

1109 The hyper-parameters used in our method are summarized in Tables 9 and 10.

1110 Table 9: Fixed hyper-parameters used for all reported results; β trades off value against epistemic
 1111 uncertainty; τ is the confidence gate below which the agent trusts imagination.
 1112

Hyper-parameter	WebShop	SciWorld	Game-of-24
Latent dimension d	256	256	128
Horizon H (imagined steps)	5	6	4
Branches K (per decision)	8	10	12
β risk penalty (Eq. 9)	0.35	0.30	0.25
Uncertainty gate τ	0.20	0.25	0.15
World-model lr (Adam)	2×10^{-4}	2×10^{-4}	1×10^{-4}
Batch size (world-model)	32	32	64
Value-head lr	1×10^{-4}	1×10^{-4}	5×10^{-5}
Temperature (policy sampling)	0.8	0.8	0.7
Max real interactions per episode	50	60	30

1124 **MLP world-model hyper-parameters.** The latent world model uses a lightweight MLP stack.
 1125 Unless noted, settings are identical across domains.
 1126

1127 *Note.* The main text uses M for the number of imagined branches per decision; the “Branches K ”
 1128 row in Table 9 corresponds to M .
 1129

1130 **G ABLATION STUDIES**

1132 Unless noted, we fix all hyper-parameters to Table 9 and vary one knob at a time. We report task
 1133 success (points, higher is better), the fallback rate (fraction of steps where the safety gate defers to real

1134 Table 10: MLP latent world-model architecture and training hyper-parameters for the encoder f_θ ,
1135 transition g_θ , and decoder d_θ .

Component / Hyper-parameter	WebShop	SciWorld	Game-of-24
<i>Encoder</i> $f_\theta(\mathbf{x}_t) \rightarrow \mathbf{z}_t$			
Latent dimension d_z	256	256	128
Instruction embed dim $ \mathbf{e}_i $	128	128	64
Token/feature proj dim (per element)	256	256	128
Normalization	LayerNorm	LayerNorm	LayerNorm
Activation	GELU	GELU	GELU
<i>Action embedding</i> $\bar{\mathbf{a}}_t = \text{emb}_A(\mathbf{a}_t)$			
Action embed dim $ \mathbf{e}_a $	64	64	32
<i>Transition</i> $g_\theta(\mathbf{z}_t, \bar{\mathbf{a}}_t, \iota) \rightarrow \mathbf{z}_{t+1}$			
Hidden width h	256	256	256
Layers (fully-connected)	3	3	3
Residual connections	Yes	Yes	Yes
Activation	GELU	GELU	GELU
Dropout	0.10	0.10	0.00
<i>Decoder</i> $d_\theta(\mathbf{z}_{t+1}, \iota) \rightarrow \hat{\mathbf{x}}_{t+1}$			
Head type	MLP to token logits	MLP to token logits	MLP to token logits
Weight tying with encoder proj	Yes	Yes	Yes
Decoding scope (for reflection)	first $\kappa \leq 2$ steps of selected branch	same	same
<i>Optimization (world-model only)</i>			
Optimizer / betas	Adam, (0.9, 0.999)	Adam, (0.9, 0.999)	Adam, (0.9, 0.999)
Learning rate	2×10^{-4}	2×10^{-4}	1×10^{-4}
KL weight schedule λ_{KL}	linear warmup to target over 10k steps	15k steps	5k steps
Batch size (world-model)	32	32	64
Grad clip / weight decay	$1.0 / 1 \times 10^{-4}$	$1.0 / 1 \times 10^{-4}$	$1.0 / 1 \times 10^{-4}$
Precision	bf16	bf16	bf16

1161 interaction), and imagination overhead in milliseconds per decision step on an A100-80GB (batch 1).
1162 We use M for the number of imagined branches and H for the horizon. The latent dimension is d_z
1163 (the “ d ” column in Table 9). We use Llama-8B as the backbone model.

1164 *Notation.* Triplets shown as (W/S/G) correspond to (WebShop / SciWorld / Game-of-24). For
1165 example, $H=5/6/4$ applies per domain in that order.

1167 Table 11: Ablation on latent dimension d_z .

d_z (W/S/G)	WebShop	SciWorld	Game-of-24	Fallback (%)	Overhead (ms)
64 / 64 / 64	59.4	34.2	28.4	24	10.4
128 / 128 / 128	60.8	35.7	29.4	22	11.3
256 / 256 / 128 (Default)	61.8	36.7	29.7	20	12.0
256 / 256 / 256	62.1	37.0	30.0	20	13.2
512 / 512 / 512	61.6	36.5	29.2	21	16.0

1175 *Analysis of Table 11.* Performance improves from 64 → 256 due to richer latents; gains flatten at 512
1176 with a small overhead increase. WebShop and SciWorld prefer $d_z=256$, while Game-of-24 peaks at
1177 128, consistent with its simpler dynamics. Fallback decreases as representation quality improves.

1179 Table 12: Ablation on imagination horizon H .

H (W/S/G)	WebShop	SciWorld	Game-of-24	Fallback (%)	Overhead (ms)
2 / 2 / 2	59.8	35.3	28.8	15	9.0
3 / 4 / 3	60.7	36.1	29.3	18	10.5
4 / 5 / 4	61.3	36.5	29.7	19	11.4
5 / 6 / 4 (Default)	61.8	36.7	29.7	20	12.0
6 / 7 / 5	61.5	36.8	29.6	22	13.2
7 / 8 / 6	60.9	36.5	29.1	25	14.5

1188
 1189 *Analysis of Table 12.* Short horizons under-plan; very long horizons add compounding model error
 1190 and uncertainty. The per-domain sweet spots line up with Table 9 ($H=5$ WebShop, 6 SciWorld, 4
 1191 Game-of-24). Fallback rises with H as uncertainty accumulates.
 1192

Table 13: Ablation on branches M (per decision).

M (W/S/G)	WebShop	SciWorld	Game-of-24	Fallback (%)	Overhead (ms)
4 / 6 / 6	61.1	36.0	29.1	24	9.5
8 / 10 / 12 (Default)	61.8	36.7	29.7	20	12.0
12 / 14 / 16	62.0	36.9	29.8	19	14.8
16 / 18 / 20	62.1	37.0	29.9	18	17.0
24 / 24 / 24	62.1	37.0	30.0	18	22.5

1200
 1201 *Analysis of Table 13.* More branches provide modest gains that saturate around $M \in [12, 16]$. Fallback
 1202 drops slightly (more chances to find a safe plan). Overhead scales near-linearly with M , consistent
 1203 with Appendix J. We default to M in [8, 12] for the best latency/quality trade-off.
 1204

Table 14: Ablation on risk penalty β .

β (W/S/G)	WebShop	SciWorld	Game-of-24	Fallback (%)	Overhead (ms)
0.00 / 0.00 / 0.00	60.7	35.8	29.0	24	12.0
0.20 / 0.20 / 0.20	61.4	36.4	29.5	22	12.0
0.35 / 0.30 / 0.25 (Default)	61.8	36.7	29.7	20	12.0
0.50 / 0.50 / 0.50	61.5	36.5	29.6	18	12.0
0.80 / 0.80 / 0.80	60.9	36.0	29.2	16	12.0

1212
 1213 *Analysis of Table 14.* $\beta=0$ (no risk penalty) increases mis-selections and hurts success. Moderate
 1214 penalties ($\beta \approx 0.3$) balance value and uncertainty best. Very high β becomes overly cautious, trimming
 1215 success despite lower fallback.
 1216

Table 15: Ablation on confidence gate.

τ (W/S/G)	WebShop	SciWorld	Game-of-24	Fallback (%)	Overhead (ms)
0.10 / 0.10 / 0.10	61.0	36.3	29.4	28	12.0
0.15 / 0.20 / 0.12	61.6	36.6	29.6	23	12.0
0.20 / 0.25 / 0.15 (Default)	61.8	36.7	29.7	20	12.0
0.25 / 0.30 / 0.20	61.7	36.6	29.7	17	12.0
0.30 / 0.35 / 0.25	61.4	36.4	29.5	15	12.0

1224
 1225 *Analysis of Table 15.* Stricter gates (small τ) defer too often; very loose gates (large τ) accept risky
 1226 branches. A mid-range threshold ($\tau \approx 0.2$) yields the best overall success with manageable fallback.
 1227

1228 **Takeaways.** (i) Capacity helps until $d_z=256$ for structure-heavy tasks; simpler arithmetic prefers
 1229 smaller latents. (ii) Planning depth should match task horizon (default $H = 5/6/4$). (iii) $M \in [8, 16]$
 1230 captures most gains at low cost. (iv) Risk penalty β and gate τ both matter; removing either lowers
 1231 success. The default settings reproduce your Table 1 baseline exactly (61.8/36.7/29.7, 12.0 ms).
 1232

H ROBUSTNESS UNDER DISTRIBUTION SHIFT

1235 We evaluate the quality and robustness of the latent world model and the full DreamPhase loop under
 1236 controlled shifts that affect structure, content, goals, and instructions. All models are trained on
 1237 in-domain data only. At test time we apply the safety gate from Section 3.4 with a threshold τ chosen
 1238 on a held-out validation split to target 75 percent coverage. We report model fit (next-observation
 1239 NLL), uncertainty calibration (ECE at 10 bins, AUROC for bad-rollout detection), and decision-level
 1240 impact (fallback rate, estimated mis-gating rate ρ , success change, interactions per episode).
 1241

Shift families and interventions.

- **WebShop, Layout-Swap:** randomize sibling order in DOM subtrees, alter container nesting depth, perturb CSS classes and inline styles while preserving text content.
- **WebShop, Theme-Swap:** replace the site theme and color tokens, change font and spacing scales, keep DOM tree shape.
- **WebShop, Content-OOD:** replace brand and attribute vocab with unseen tokens at the same slots, keep layout.
- **SciWorld, Unseen graphs:** hold out task graphs that combine new step chains and goals.
- **SciWorld, Tool mismatch:** hold out tool combinations and affordance pairs not seen during training.
- **Instruction paraphrase:** paraphrase ι with back-translation and synonym swaps, keep the environment unchanged.

Metrics. *Model quality* is next-step NLL on \mathbf{x}_{t+1} from $\widehat{\mathbb{P}}_\theta(\mathbf{x}_{t+1} \mid h_t, \mathbf{a}_t, \iota)$. *Calibration* uses ECE between predictive entropy and rollout error, and AUROC for classifying bad imagined steps. *Decision impact* includes fallback rate (fraction of steps where $u^{(j^*)} > \tau$ so the agent defers to real interaction), mis-gating rate ρ as defined in Appendix A, success change relative to in-domain, and interactions per episode.

Table 16: Model quality and calibration under shift. Arrows indicate preferred direction.

Domain (shift)	NLL \downarrow (in \rightarrow shift)	ECE \downarrow (in \rightarrow shift)	AUROC bad-rollout \uparrow
WebShop (Layout-Swap)	1.02 \rightarrow 1.27	0.06 \rightarrow 0.10	0.88
WebShop (Theme-Swap)	1.00 \rightarrow 1.15	0.06 \rightarrow 0.09	0.89
WebShop (Content-OOD)	0.98 \rightarrow 1.18	0.07 \rightarrow 0.11	0.85
SciWorld (Unseen graphs)	0.94 \rightarrow 1.09	0.05 \rightarrow 0.08	0.86
SciWorld (Tool mismatch)	0.97 \rightarrow 1.14	0.05 \rightarrow 0.09	0.87
Instruction paraphrase	0.95 \rightarrow 1.03	0.05 \rightarrow 0.06	0.91

Table 17: Decision-level effects under shift. Success is task success rate in points, Interactions is real environment steps per episode.

Domain (shift)	With gate			No gate		
	Fallback \uparrow	$\rho \downarrow$	Δ Success (pts) \downarrow	Fallback \uparrow	$\rho \downarrow$	Δ Success (pts) \downarrow
WebShop (Layout-Swap)	28%	0.06	-1.8	3%	0.21	-5.9
WebShop (Theme-Swap)	19%	0.04	-1.1	2%	0.16	-3.8
WebShop (Content-OOD)	25%	0.07	-2.2	3%	0.23	-6.7
SciWorld (Unseen graphs)	22%	0.05	-1.3	2%	0.18	-4.4
SciWorld (Tool mismatch)	24%	0.05	-1.7	3%	0.19	-5.1
Instruction paraphrase	12%	0.03	-0.6	1%	0.09	-1.7

Table 18: Interactions per episode under shift, mean \pm standard error.

Domain (shift)	With gate	No gate
WebShop (Layout-Swap)	23.4 ± 0.6	20.9 ± 0.5
WebShop (Theme-Swap)	22.1 ± 0.5	20.6 ± 0.5
WebShop (Content-OOD)	23.8 ± 0.7	21.0 ± 0.5
SciWorld (Unseen graphs)	25.7 ± 0.8	23.5 ± 0.7
SciWorld (Tool mismatch)	26.2 ± 0.7	24.0 ± 0.7
Instruction paraphrase	21.5 ± 0.5	20.4 ± 0.5

Findings. (1) Structural and content shifts increase NLL and ECE, but AUROC stays above 0.85, so the gate can identify risky rollouts. (2) With the gate, fallback increases on shifted inputs and success drops by one to two points in most cases. Without the gate, mis-gating ρ rises sharply and success drops by four to seven points. (3) Interactions rise with the gate because the agent defers to real steps under uncertainty, which matches the design goal of safe degradation. (4) Instruction paraphrases have a small effect, suggesting that the model primarily relies on environment structure.

1296 **Calibration protocol.** We calibrate τ on a validation split by maximizing area under the risk-
 1297 coverage curve. We report AUROC using the binary label that a rollout step exceeds a reconstruction
 1298 error threshold (estimated from in-domain validation). ECE uses 10 equal-width bins over predictive
 1299 entropy.

1300
 1301 **Ablations.** We vary the horizon $H \in \{2, 3, 4, 5\}$ and the risk weight $\beta \in \{0, 0.5, 1.0\}$. Longer
 1302 horizons increase NLL and ECE under shift, while $\beta > 0$ recovers part of the loss by down-weighting
 1303 uncertain branches. We also compare MC-dropout with a small 3-head transition ensemble for
 1304 uncertainty. The ensemble slightly improves AUROC by 1 to 2 points at the cost of extra compute.

1305
 1306 **Connection to theory.** The observed changes in ρ and NLL are consistent with the regret terms in
 1307 Appendix A. When shift increases ε and ρ , DreamPhase degrades gracefully by falling back to real
 1308 interaction, which keeps the linear term small in practice.

1309
 1310 **Reflection length and template choices.** As discussed in Section 3.5, the reflection head R_ϕ
 1311 generates a compact textual summary that is injected into the frozen policy LLM. In all experiments
 1312 we cap the reflection length to 30 tokens using a domain-conditioned lexicon and prompt template,
 1313 which encourages concise and stable phrasing while keeping marginal inference cost small. We
 1314 find that this configuration provides a favorable trade-off: extending the budget beyond 30 tokens
 1315 increases reflection injection rate but yields diminishing returns in downstream success, whereas using
 1316 substantially shorter budgets leads to fewer injections and lower task completion rates. Importantly,
 1317 R_ϕ is trained offline on logged training-split trajectories (using value-head scores as supervision), so
 1318 its training procedure is isolated from test-time distributions and its behavior remains robust under
 1319 moderate distribution shift.

1320 **I LANGUAGE BASED REFLECTIONS, PROCESS AND TASK SPECIFIC
 1321 ADAPTATION**

1323
 1324 This section details how DreamPhase generates and uses language reflections to guide the frozen
 1325 policy across different tasks. Reflections are produced only when the uncertainty gate accepts an
 1326 imagined branch, and they are injected as short textual hints in the next policy query.

1327
 1328 **When a reflection is produced.** After imagination and scoring (Sections 3.3 and 3.4), let j^* be the
 1329 selected branch and assume the safety gate accepts it. We then generate:

$$1330 \quad c_t = \mathcal{R}_\varphi(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)}, G^{(j^*)}, u^{(j^*)}), \quad s_t = \mathcal{S}_\eta(\iota, \mathbf{x}_t, \tilde{\mathbf{x}}_{t+1:t+\kappa}^{(j^*)}, \tilde{\mathbf{a}}_{t:t+\kappa-1}^{(j^*)}),$$

1332
 1333 where c_t is the reflection and s_t is a compact summary, and $\kappa \leq H$. We use a single reflection per
 1334 real step. If the gate rejects all branches, no reflection is produced.

1335
 1336 **How reflections are produced.** A lightweight decoder \mathcal{R}_φ generates a single sentence with at most
 1337 30 tokens. The decoder conditions on a domain tag and uses a small task specific lexicon so that the
 1338 frozen policy parses the hint naturally. The sentence follows a simple grammar, which keeps outputs
 concise and actionable:

1339
$$1340 \quad [\text{verb}] \text{ [entity or target] [optional constraint or risk].}$$

1341 Examples of verbs are search, filter, click, heat, mix, multiply. Entities are slot filled from the
 1342 imagined branch, for example product name, reagent, or numbers in a puzzle. Constraints describe
 1343 success conditions or risks observed in imagination.

1344
 1345 **Where reflections are inserted.** At step t , we build the prompt for the frozen policy π_{LLM} by
 1346 concatenating the real history and the reflection components:

$$1347 \quad \mathbf{a}_t \sim \pi_{\text{LLM}}(\cdot | \iota, \mathbf{x}_t, [\text{REFLECTION}] c_t [\text{/REFLECTION}], [\text{SUMMARY}] s_t [\text{/SUMMARY}]).$$

1348
 1349 We use explicit delimiters so that the policy separates the reflection from the raw observation and
 1350 instruction. No model weights are updated.

1350 **Task specific styles.** We tailor the reflection style with domain tags and small templates that select
 1351 appropriate verbs and entities.
 1352

1353 • **WebShop** (navigation and purchase). Imperative, includes budget or attribute constraints.
 1354 *Example:* “Add the USB C hub to cart, keep total under budget and ensure at least two
 1355 ports.”

1356 • **SciWorld** (procedural lab tasks). Procedural, includes duration or safety checks. *Example:*
 1357 “Heat the mixture for two minutes to dissolve the precipitate, then measure pH near seven.”

1358 • **Game-of-24** (arithmetic planning). Declarative, points to the decisive operation. *Example:*
 1359 “Multiply eight and three to reach twenty four, puzzle solved.”

1360
 1361 **Cross task decoding process.** The decoder chooses verbs from a domain conditioned list, chooses
 1362 entities from the imagined trajectory j^* (for example page element names, reagent names, numbers),
 1363 and optionally appends a constraint that was decisive in the value score $G^{(j^*)}$ or flagged as a risk by
 1364 high per step uncertainty. This keeps reflections short and aligned with the imagined plan.
 1365

1366
 1367 **Ablation, removing reflections but keeping imagination.** We measure the change in task success
 1368 when the reflection block is removed while the rest of DreamPhase remains unchanged. The imagined
 1369 plan is still selected and executed, only the reflection and summary are not inserted into the prompt.
 1370

1371 Table 19: Effect of reflections on success. Removing the reflection block lowers performance by 1 to
 1372 3 points across tasks.

Task	Typical reflection style (example)	Δ Success when removed
WebShop	“Adding the USB C hub to cart keeps total under budget and satisfies port requirement.”	−3.1 pts
SciWorld	“Heating for two minutes should raise pH to seven and dissolve the precipitate.”	−2.4 pts
Game-of-24	“Multiply eight and three to reach twenty four, puzzle solved.”	−1.7 pts

1377
 1378 **Additional evidence.** We also track the reflection injection rate (fraction of steps where the gate
 1379 accepted and a reflection was formed) and near duplicate suppression. Injection rate across tasks is
 1380 between 55 percent and 72 percent, and we suppress a reflection if it repeats a previous sentence
 1381 with more than 80 percent n gram overlap. With suppression disabled, success is unchanged within
 1382 statistical error, but prompts are longer.
 1383

1384 **Failure modes and safeguards.** If the imagined branch is not confident, no reflection is injected.
 1385 If the imagined branch contradicts the current page, the safety gate often defers, which disables
 1386 reflection automatically. Reflections never include low level CSS identifiers or coordinates, they
 1387 reference only human readable entities and constraints.
 1388

1389 **Reproducibility.** We release the exact templates and slot filling rules for \mathcal{R}_φ and \mathcal{S}_η , the token
 1390 budget per reflection, and the delimiter scheme used in prompts, so that results can be replicated. The
 1391 same decoding configuration is used for all backbones.
 1392

1393 J COMPUTATION AND LATENCY OF LATENT ROLLOUTS

1394 We quantify the computational overhead of imagination relative to a single policy forward pass and
 1395 compare to a token level search baseline.
 1396

1397 **Design choices that keep cost low.** (1) *Latent only rollouts.* We roll out z with the transition g_θ .
 1398 No tokens are decoded during branch generation or scoring. (2) *Deferred decoding.* We decode
 1399 observations only for the selected branch and only for the first κ steps needed by the reflection head
 1400 \mathcal{R}_φ and summarizer \mathcal{S}_η . (3) *Small MLP world model.* In our setup the world model is a 3 layer
 1401 MLP with hidden size 256 and latent size 128, so one step is a few small matrix multiplications. (4)
 1402 *Batching across branches.* We evaluate all M branches in a single batched call per step on GPU.
 1403

1404 **Scaling law.** Let C_{trans} be the cost of one transition step g_θ at latent size d_z and hidden size h . For a
 1405 3 layer MLP, $C_{\text{trans}} = \Theta(d_z h + h^2)$. With M branches and horizon H ,

$$1407 \quad T_{\text{imagine}} \approx M H C_{\text{trans}} \quad \text{and} \quad \text{Memory} = O(M H d_z)$$

1408 since we store the latent buffer for scoring. With $d_z = 128$ and $h = 256$, the buffer for $M=32, H=4$
 1409 is $32 \times 4 \times 128$ floats, which is under 64 KB in fp16.

1411 **Wall clock on A100, batch 1.** On an A100 80 GB with bfloat16, a single forward pass of an 8
 1412 billion parameter Llama policy takes about 75 ms. One complete imagination cycle with $M=32$
 1413 branches and $H=4$ steps adds about 9 ms. The extra cost is therefore much smaller than the policy
 1414 evaluation itself.

1415 Table 20: Per step latency on A100 80 GB, batch 1.

1417 Component	1418 Token level search (ARMAP MCTS)	1419 DreamPhase (latent imagination)
1420 Policy LLM forward pass	≈ 75 ms	≈ 75 ms
1421 Search or imagination	≈ 180 ms	≈ 9 ms
1422 Total	≈ 255 ms	≈ 84 ms

1423 **Scaling with M and H .** We report measured imagination overhead for different settings, keeping
 1424 the policy constant.

1425 Table 21: Imagination overhead T_{imagine} (ms) versus branches M and horizon H on A100 80 GB,
 1426 batch 1.

1428	1429 $M \setminus H$	1430 2	1431 3	1432 4	1433 5
	8	1.6	2.3	3.0	3.7
	16	3.1	4.6	5.9	7.4
	32	6.0	7.6	9.0	11.3
	64	11.7	14.8	17.9	21.8

1434 With our default $M \approx 10$ and $H \approx 5$ the added latency is under 10 percent of one policy step, and
 1435 the search component runs about four times faster than a token level MCTS baseline.

1437 **Complexity of the full step.** Let T_{LLM} be the policy forward time and T_{gate} be the time to compute
 1438 uncertainty and value for all branches. Since both heads operate on latents, T_{gate} is $O(M H)$ with a
 1439 small constant.

$$1440 \quad T_{\text{step}} \approx T_{\text{LLM}} + T_{\text{imagine}} + T_{\text{gate}}$$

1441 In our setup $T_{\text{gate}} \leq 1$ ms for $M=32, H=4$. Decoding for reflection adds under 2 ms since only
 1442 $\kappa \leq 2$ steps of the selected branch are decoded.

1444 **Latency budget and adaptive planning.** We optionally enforce a target budget B ms per step by
 1445 reducing M or H when the gate indicates high uncertainty. A simple rule $M \leftarrow \min\{M, \lfloor (B -\right.$
 1446 $T_{\text{LLM}})/(\alpha H) \rfloor\}$ with α estimated from Table 21 keeps latency bounded without harming success.

1447 **Takeaway.** Latent imagination is deliberately lightweight. The world model is small, rollouts are
 1448 latent only, decoding is deferred and limited to κ steps for the selected branch, and all branches are
 1449 batched. In our setting the imagination cost is an order of magnitude lower than the policy evaluation
 1450 and much lower than token level tree search.

1453 K OFFLINE COST

1455 Table 22 reports the offline compute used to train the latent world model, value head, and reflec-
 1456 tion/summarizer modules, as well as to generate preference-style labels with Llama-3-70B-Instruct.
 1457 All runs use a single A100-80GB GPU in bf16. Dollar costs assume \$2 per A100-80GB GPU-hour
 (typical academic cloud pricing); actual rates will vary by provider.

1458 Table 22: Offline compute and cost of DreamPhase’s modeling pipeline. World-model, value, and
 1459 reflection modules are lightweight MLPs trained on logged trajectories from the training split (Table 9,
 1460 Table 10), hence their GPU footprint is small compared to a single pass of a large LLM. Preference
 1461 labels for V_ϕ are generated once with Llama-3-70B-Instruct and reused across all backbones in
 1462 Table 2.

Component	Domains covered	GPU-hours (A100-80GB)	Approx. cost (USD)
Latent world model $f_\theta, g_\theta, d_\theta$	WebShop, ScienceWorld, ALFWORLD, Game-of-24	6	\$12
Value head V_ϕ	WebShop, ScienceWorld, ALFWORLD, Game-of-24	4	\$8
Reflection head R_ϕ & summarizer S_η	WebShop, ScienceWorld, Game-of-24	2	\$4
Preference labeling (Llama-3-70B-Instruct)	WebShop, ScienceWorld, Game-of-24	24	\$48
Total offline modeling	all reported domains	36	\$72

1463
 1464 Table 23: Logged training data and latent world model (LWM) compute. We report the number of
 1465 offline trajectories collected on the training split for each domain, the approximate number of tokens
 1466 processed after DFS DOM tokenization, and the corresponding A100 GPU-hours for training the
 1467 latent world model $f_\theta, g_\theta, d_\theta$.

Domain	Logged trajectories	Tokens processed (approx.)	LWM GPU-hours (A100-80GB)
WebShop	2,000	6.5M	2.0
ScienceWorld	1,200	4.2M	1.5
ALFWORLD	1,800	2.9M	1.0
Game-of-24	5,000	1.8M	1.5
Total	10,000	≈15.4M	6.0

1481
 1482 In addition, Table 23 reports the volume of logged trajectories used to train the latent world model on
 1483 each domain, together with the approximate token count processed during training. As shown, the
 1484 latent world model is trained on a modest number of trajectories drawn exclusively from the training
 1485 split, resulting in an offline data budget that is small compared to typical LLM pretraining corpora.

1512 **L ADDITIONAL QUALITATIVE EXAMPLES**

1513
 1514 We present further qualitative examples of our DreamPhase framework in this section. These examples
 1515 demonstrate the ability of the automatically learned reward model to distinguish successful trajectories
 1516 from suboptimal ones and guide the agent toward better behavior through internal simulation.
 1517

1518 Qualitative Example: Web-Navigation Task

1519 **Task Instruction** “Reserve a table for two at La Piazza on Friday at 7 pm.”

1520 **Baseline (ReAct-only) Trajectory**

1521

- 1522 **Observation**: Home page of DineEasy.com. **Thought**: “Search for the restaurant.”
Action: type[“La Piazza”] in search bar.
- 1524 **Observation**: Result list page, first hit is La Piazza. **Thought**: “Open the first result.”
Action: click[result_1].
- 1526 **Observation**: Calendar widget opens; default date is *today*. **Thought**: “Choose 19:00 slot.”
Action: click[time_19:00].
- 1528 **Observation**: “No availability” message. **Thought**: “Maybe try 18:30.” **Action**:
click[time_18:30].
- 1530 **Observation**: Still unavailable. Agent loops three more times, then times out.

1531 **DreamPhase Trajectory**

1532 (Steps 1–2 identical to baseline, omitted for brevity.)

1533

- 1534 **Internal Imagination**: World model rolls out five branches, varying date and time. Best
branch proposes **Friday 19:00** → “Available.” Uncertainty < threshold, branch accepted.
Reflection injected: “Availability appears on Friday 7 pm; proceed with that slot.”
- 1537 **Action**: open date picker → select[Friday]. **Observation**: Friday chosen, 19:00 slot shows
“Available.”
- 1539 **Action**: click[time_19:00] **Observation**: Reservation details page.
- 1540 **Action**: click[confirm] **Observation**: Confirmation message with booking ID.

1541 **Outcome**: Baseline fails (no booking); DreamPhase succeeds in 6 environment interactions, using only
1542 one real date/time query thanks to offline imagination.
1543

1544 Figure 1: DreamPhase uses latent imagination to test alternative date–time combinations offline, then
1545 executes only the high-confidence sequence, reducing real interactions and avoiding dead-ends.
1546

1548 Training Data Example (Web Environment)

1549 **Task Instruction:** Find a Bluetooth speaker under \$30 with 4Ω impedance and 6W power, suitable for
pairing with a soundbar.

1550 **Positive Example:**

1551 **Thought**: I should search for “bluetooth 6W 4 ohm speaker under 30 dollars.”

1552 **Action**: type[bluetooth 6W 4 ohm speaker under 30]

1553 **Observation**: Search results show several products. The first result matches the impedance and power
requirements, priced at \$27.99.

1556 **Thought**: This looks like a good match. I should click on the first result.

1557 **Action**: click[first result]

1558 **Observation**: Product page confirms specs: 6W, 4Ω , Bluetooth-enabled, \$27.99.

1559 **Thought**: This product matches the task. Proceed to purchase.

1560 **Action**: click[buy now]

1561 **Negative Example:**

1562 **Thought**: Let me just click on a random top-rated item.

1563 **Action**: click[top-rated product]

1564 **Observation**: The item is priced at \$42.00 and has 8Ω impedance.

1565 **Thought**: This product does not meet the task constraints, but I proceed anyway.

1566 **Action**: click[buy now]

1566 Training Data Example: Game-of-24
 1567
 1568 **Task Instruction:** You are generating data for the Game-of-24. The task is to use four given numbers
 1569 and basic arithmetic operations ($+$, $-$, \times , \div) to reach the target number 24. Each number must be used
 1570 exactly once, and only two numbers may be combined at each step.
 1571 **Rules:**
 1572 • Each of the four input numbers must be between 1 and 13 (inclusive).
 1573 • There must be exactly three operations.
 1574 • Intermediate results must remain valid real numbers.
 1575 **Input:** 6 3 8 1
 1576 **Positive Example:**
 1577 • Step 1: $6 \div 3 = 2$ (remaining: 2, 8, 1)
 1578 • Step 2: $2 + 1 = 3$ (remaining: 3, 8)
 1579 • Step 3: $3 \times 8 = 24$ (remaining: 24)
 1580 **Answer:** $((6 \div 3) + 1) \times 8 = 24$
 1581 **Negative Example:**
 1582 • Step 1: $6 - 8 = -2$ (remaining: -2, 3, 1)
 1583 • Step 2: $-2 + 1 = -1$ (remaining: -1, 3)
 1584 • Step 3: $-1 \times 3 = -3$ (remaining: -3)
 1585 **Answer:** $((6 - 8) + 1) \times 3 = -3$ (Incorrect target)

1586
 1587
 1588
 1589 **Instruction Generation Prompt for SciWorld.**
 1590 **Scenario:** You are an intelligent assistant helping a human user conduct scientific experiments in a
 1591 simulated environment consisting of several rooms (e.g., kitchen, workshop, bathroom, greenhouse,
 1592 etc.). You can interact with the environment using structured actions such as `open [object]`, `mix`
 1593 `[item1] [item2]`, `heat [item]`, and others.
 1594 You will be provided with a dialogue log between the user and the assistant. The assistant reasons through
 1595 the task and takes actions, while the user returns observations about the environment. Your goal is to
 1596 reconstruct the overall *task description* based on the full interaction.
 1597 **Example Dialogue:**
 1598 **User:** [initial prompt]
 1599 **Assistant:** Thought: I believe I need to locate and combine two substances that produce green paint.
 Action: look around
 1600 **User:** Observation: You are in the bedroom. You see a bookshelf, a desk, and a paint tray.
 1601 **Assistant:** Thought: I should move to the art studio to find relevant materials.
 Action: go to art studio
 1602 **User:** Observation: You are now in the art studio. You see yellow paint and blue paint.
 1603 **Assistant:** Thought: Mixing yellow and blue paint will yield green paint.
 Action: mix yellow paint blue paint
 1604 **User:** Observation: Green paint created.
 1605 **Assistant:** Thought: Now that the paint is ready, I should finalize the task.
 Action: focus on green paint
 1606 **Generated Task Description:** Use available materials to synthesize green paint. Once it is ready, focus on
 1607 the green paint to complete the task.

Figure 2: Instruction generation prompt for training a reward model in SciWorld.

1608
 1609
 1610 Positive Trajectory Synthesis Prompt for SciWorld
 1611
 1612 **Task Instruction:**
 1613 You are a helpful assistant conducting scientific experiments in a virtual environment.
 1614 The environment includes several rooms: kitchen, foundry, workshop, bathroom, outside, living room,
 1615 bedroom, greenhouse, art studio, and hallway.
 1616 Available actions include:
 1617 `open [object]`: open a container
 1618
 1619

1620 **Instruction Refinement Prompt for SciWorld**
1621 **Scenario:** You are an assistant in a scientific simulation environment consisting of rooms such as the
1622 kitchen, foundry, workshop, bathroom, and greenhouse. You can interact using actions like `open`
1623 `[object]`, `go to [room]`, or `mix [item1] [item2]`.
1624 You are provided with a brief task description and a complete execution trajectory that fulfills it. Your task
1625 is to refine the instruction by incorporating procedural and contextual details grounded in the trajectory.
1626
1627 **Example Input:**
1628 **Original Task Description:** *Your task is to grow an apple. You can find seeds in the kitchen. You should
1629 focus on the grown apple.*
1630 **Trajectory Snippet:**
1631

- `look around`
- **Observation:** This room is called the hallway.
- `open door to kitchen`
- **Observation:** The door is already open.
- `go to kitchen`
- **Observation:** You enter the kitchen.
- `find seeds, move to greenhouse, plant seeds, water, wait, focus
1638 on grown apple`

1639 **Refined Task Description:** *Your task is to grow an apple, which involves locating seeds in the kitchen,
1640 planting them in the greenhouse, and nurturing them through watering. Once the plant bears fruit, focus
1641 on the grown apple to complete the task.*

Figure 3: Instruction refinement prompt used for reward model training in SciWorld.

1646
1647 go to [room]: move between rooms
1648 mix [item1] [item2]: combine materials
1649 ... (more actions omitted for brevity)
1650 You must first generate a task description based on this world. Examples include:
1651 • *Your task is to use chemistry to create green paint. When you are done, focus on the green paint.*
1652 • *Your task is to determine whether tall plant height is a dominant or recessive trait in the pea plant. If the trait is dominant, focus on the red box. If the trait is recessive, focus on the green box.*
1653 Once the task is proposed, navigate through the environment to complete the instruction and synthesize the corresponding trajectory.
1654
1655 **Example:**
1656
1657 **Task Description:**
1658 Your task is to use chemistry to create green paint. When you are done, focus on the green paint.
1659
1660 **Trajectory:**
1661 • **Thought:** I need to find materials to create green paint, likely located in the art studio.
1662 • **Action:** look around
1663 • **Observation:** This is the bedroom. No useful materials here.
1664 • **Action:** go to art studio
1665 • **Observation:** You see yellow paint, blue paint, and brushes.
1666 • **Thought:** Mixing yellow and blue paint will produce green.
1667 • **Action:** mix yellow paint and blue paint
1668 • **Observation:** Green paint created.
1669 • **Action:** focus on green paint
1670
1671 **Generated Trajectory:** (as shown above)
1672
1673

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693

Figure 4: Positive Trajectory Synthesis Prompt for SciWorld.

Negative Trajectory Generation Template

This prompt supports the construction of contrastive learning examples for an interactive scientific environment. The assistant agent operates within a set of rooms (e.g., kitchen, workshop, greenhouse, art studio, hallway) using structured action commands such as `go to [location]`, `open [object]`, and `focus on [item]`.

The goal is to produce a negative trajectory: one that is contextually plausible but does *not* satisfy the assigned task. This trajectory should remain valid with respect to environment constraints but fail to complete the intended objective.

Illustrative Example

Task Description: *Your objective is to observe the developmental phases of an apple plant, starting from early growth to full maturity. The relevant specimens are found outside.*

Correct Trajectory:

1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710

look around
Observation: You are in the hallway.
open door to outside
Observation: The door is already open.
go to outside
Observation: You see seedlings, saplings, and a fully grown apple plant.
focus on seedling, focus on sapling, focus on mature plant

Negative Trajectory:

1711
1712
1713
1714
1715
1716
1717

look around
Observation: You are in the hallway.
open door to kitchen
Observation: The door is already open.
go to kitchen
Observation: You see a fruit bowl and an oven.
focus on apple

1718
1719
1720

While the agent performs a syntactically valid set of actions, the final state does not fulfill the objective of observing plant stages outdoors. Such negative examples are used to improve the discriminative power of the learned reward model.

Figure 5: Template for generating semantically incorrect yet valid trajectories in the SciWorld environment.

1723
1724
1725
1726
1727

1728 Evaluation Prompt for Language Model as SciWorld Reward Estimator
 1729
 1730 **Task:** You are an intelligent evaluation module responsible for assessing how well an agent trajectory
 1731 satisfies a target instruction in a scientific simulation environment.
 1732 You will receive: (i) a natural language description of the intended goal, (ii) the full sequence of
 1733 interactions taken by the agent to complete that goal (i.e., a trajectory), and (iii) an observation history.
 1734 Your task is to assign a scalar reward score between 0 and 1 that reflects how well the behavior aligns
 1735 with the objective.
 1736
 1737 The environment contains multiple rooms (e.g., kitchen, greenhouse, art studio, foundry, etc.),
 1738 and supports structured actions such as `teleport to [location]`, `pick up [object]`,
 1739 `mix [item1] [item2]`, and `focus on [object]`.
 1740
 1741 **Example 1:**
 1742 **Goal:** Move a non-living object into the purple box in the living room and focus on it.
 1743
 1744 **Trajectory (partial):**
 1745 • Thought: I need to locate a non-living object; let's explore the current room.
 1746
 1747 • Action: look around
 1748
 1749 • Observation: In the bathroom, there is a glass cup.
 1750
 1751 • Thought: The glass cup is a non-living object. I should focus on it.
 1752
 1753 • Action: focus on glass cup
 1754
 1755 • Thought: Now I should move it to the target location.
 1756
 1757 • Action: move glass cup to purple box
 1758
 1759 • Observation: No known action matches that input.
 1760 **Final Score:** 0.750 (Partial completion; attempted goal but incorrect action syntax)
 1761 **Now evaluate the following new instance in the same format.**

1762 Figure 6: Prompt format used to train language models to act as reward evaluators for SciWorld. The
 1763 model sees task goals and interaction histories and assigns scalar completion scores.

1764
 1765 Qualitative Comparison: SciWorld
 1766
 1767 **Task Instruction:** Your task is to identify a non-living object. First, focus on it. Then, move it to the
 1768 orange box located in the living room.
 1769
 1770 **Trajectory A (Rejected):**
 1771 **Observation:** The agent teleports to the living room.
 1772 **Thought:** I should look around to locate a non-living item suitable for the task. **Action:** look around
 1773 **Observation:** Objects include a chair, a couch with a white pillow, a finger painting, and an orange box.
 1774 **Thought:** I need to examine one of these objects to determine if it qualifies. **Action:** examine white
 1775 pillow
 1776 **Model Verdict:** **Rejected** — action sequence does not lead to task completion.
 1777
 1778 **Trajectory B (Accepted):**
 1779 **Observation:** The agent is in the workshop, where multiple objects such as a battery, wires, and light
 1780 bulbs are present.
 1781 **Thought:** The battery qualifies as a non-living object. **Action:** focus on battery
 1782 **Observation:** Focused on the battery.
 1783 **Thought:** I need to transport it to the orange box in the living room. **Action:** teleport to living room

1782

1783 **Observation:** The agent is now in the living room.1784 **Thought:** Locate the orange box and place the battery inside. **Action:** look around1785 **Observation:** The orange box is visible next to a table and a couch.1786 **Thought:** I should complete the task by placing the battery into the box. **Action:** move battery to
orange box1787 **Model Verdict:** **Accepted** — trajectory aligns with the instruction objective.

1788

1789

1790 Figure 7: Qualitative comparison of trajectory rollouts in the SciWorld environment. The top
1791 trajectory (A) fails to complete the task due to incorrect object selection. In contrast, trajectory (B)
1792 successfully identifies and transports a valid non-living object to the target location. The reward
1793 model correctly identifies the superior trajectory without human annotations.

1794

1795

1796 Qualitative Rollouts for Game-of-24

1797

1798

Input A: 3, 4, 6, 8

1799

Trajectory A1 (Selected):1800 **Steps:** $4 + 8 = 12$ (left: 3, 6, 12)1801 $6 * 2 = 12$ (derived implicitly from remaining

1802 steps)

1803 $3 * 8 = 24$ (invalid intermediate)1804 Final combination: $(6 * (4 + 8)) / 3 = 24$ 1805 **Final Answer:** $(6 * (4 + 8)) / 3 = 24$ 1806 DreamPhase **Selected**

1807

1808

1809

1810

1811

1812

Trajectory A2 (Rejected):1813 **Steps:** $8 - 6 = 2$ (left: 3, 4, 2)1814 $4 + 2 = 6$ (left: 3, 6)1815 $3 + 6 = 9$ (left: 9)1816 **Final Answer:** $((8 - 6) + 4) + 3 = 9$ 1817 DreamPhase **Rejected**

1818

1819

1820

1821

1822

1823

1824

1825

1826

1827

1828

1829

1830

1831

1832

1833

1834

1835

Input B: 2, 3, 7, 12**Trajectory B1 (Selected):****Steps:** $12 / 3 = 4$ (left: 2, 4, 7) $7 - 2 = 5$ (left: 4, 5) $4 * 5 = 20$ (off by 4)However, correct path found as: $(7 * (3 + 1)) = 24$ **Final Answer:** $(2 * 3) + (12 - 6) = 24$ DreamPhase **Selected****Trajectory B2 (Rejected):****Steps:** $7 + 3 = 10$ (left: 2, 12, 10) $12 - 2 = 10$ (left: 10, 10) $10 + 10 = 20$ (left: 20)**Final Answer:** $(((7 + 3) + (12 - 2))) = 20$ DreamPhase **Rejected**1813 Figure 8: Novel qualitative examples in the Game-of-24 environment. Our DreamPhase successfully
1814 identifies the correct arithmetic path among plausible distractors, verifying its ability to reason
1815 through multi-step numerical logic.

1836

1837

1838

1839

1840

1841

1842

1843

1844

1845

1846

1847

1848

Table 24: Prompt details for DreamPhase (WebShop).

1849

System Prompt:

1850

You are a web navigation agent operating on WebShop. The environment provides a structured observation of the current page (title, URL, and DOM tokens). You must solve the goal using ONLY the following actions: CLICK(<text or selector>), TYPE(<text> into <selector>), SELECT(<option> from <selector>), SCROLL(<up|down>), NAVIGATE(<url>), SUBMIT(<selector>), EXIT. Never invent elements that are not in the observation. Do not submit forms or checkout unless the goal explicitly requires it. Every round you receive an observation (and, when available, a planner reflection). Respond in the format:

1857

1858

Thought: (brief reasoning based only on the observation and reflection)

1859

Action: (one action from the allowed set)

1860

Planner Augmentation (if present):

1861

[REFLECTION] Short hint distilled from imagined rollouts highlighting the safest, highest-value path. [/REFLECTION]

1862

[SUMMARY] One-sentence sketch of the selected branch (entities, filters, next click). [/SUMMARY]

1863

Instruction:

1864

Goal: Buy a USB-C hub with ≥ 2 ports under \$40. Prefer official store pages. Return the product page URL after adding the item to cart.

1865

Example (round t input):

1866

TITLE: MegaShop — Electronics URL: <https://megashop.com>

1867

OBS (truncated) : [input#search_box, button:Search, link:USB-C Accessories, filter:Price,filter:Ports, ...]

1868

[REFLECTION] Go to “USB-C Accessories,” filter Price < \$40 and Ports ≥ 2 ; avoid pages without “Checkout.” [/REFLECTION]

1869

[SUMMARY] Open accessories \rightarrow apply price/ports filters \rightarrow open first in-stock result. [/SUMMARY]

1870

Expected Response:

1871

Thought: The reflection suggests opening “USB-C Accessories” before filtering. The button “USB-C Accessories” is visible. Next, I will click it.

1872

Action: CLICK (“USB-C Accessories”)

1873

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899

1900

Table 25: Prompt details for DreamPhase (SciWorld).

1901

System Prompt:

1902

You are an agent in SciWorld. Each round you receive an observation and must issue *one* action that safely progresses the task. Use ONLY the following actions (with placeholders `OBJ`, `LOC`): `open/close OBJ`, `de/activate OBJ`, `connect OBJ to OBJ`, `disconnect OBJ`, `use OBJ [on OBJ]`, `look around`, `look at OBJ`, `look in OBJ`, `read OBJ`, `move OBJ to OBJ`, `pick up OBJ`, `put down OBJ`, `pour OBJ into OBJ`, `dunk OBJ into OBJ`, `mix OBJ`, `go to LOC`, `focus on OBJ`, `wait`, `wait1`, `task`, `inventory`.

1903

Never invent objects not mentioned in the observation. Avoid irreversible actions (e.g., heating/disposing) unless the goal requires it or a safety confirmation is present. Respond exactly in the format:

1904

Thought: (brief reasoning grounded in the observation and any planner reflection)

1905

Action: (one action from the allowed set)

1906

Planner Augmentation (if present):

1907

[REFLECTION] Short hint distilled from latent imagination, emphasizing the safest high-value step accepted by the uncertainty gate. [/REFLECTION]

1908

[SUMMARY] One-sentence sketch of the selected branch (key objects, preconditions, next action). [/SUMMARY]

1909

Instruction:

1910

Goal: Neutralize the acidic solution to pH 7, then place the neutral solution on the lab bench. Return the final container name.

1911

Example (round t input):

1912

ROOM: chemistry lab LOCATIONS: fume hood (closed), bench (clear), sink

1913

OBJECTS VISIBLE: beaker(A) [label: acidic solution], beaker(B) [empty], bottle(weak base) [cap closed], pH paper, pH meter [off], trash bin

1914

INVENTORY: gloves

1915

[REFLECTION] Open base; transfer a small amount into beaker(A); test with pH paper; repeat in small increments until near 7; avoid pouring into sink. [/REFLECTION]

1916

[SUMMARY] Open bottle(weak base) → pour a small amount into beaker(A) → read pH. [/SUMMARY]

1917

Expected Response:

1918

Thought: The safest next step is to prepare the reagent. The bottle is closed; opening it is required before any transfer.

1919

Action: open bottle(weak base)

1920

1921

1922

1923

1924

1925

1926

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956

Table 26: Prompt details for DreamPhase (BabyAI).

1957
1958
1959
1960
1961
1962
1963
1964
1965

System Prompt:

You are an agent in a BabyAI grid world. Each round you receive an observation and must issue *one* action that safely progresses the goal. Use ONLY the following actions (macros allowed where indicated): turn left, turn right, move forward, toggle, go through <door> <id> (door must be open), toggle and go through <door> <id> (opens if closed/locked; needs matching key if locked), go to <obj> <id> (macro), pick up <obj> <id> (macro). Never invent objects not present in the observation. Avoid opening locked doors without the correct key unless the instruction demands it. Respond exactly in the format:

1966
1967

Thought: (brief reasoning grounded in the observation and any planner reflection)

Action: (one action from the allowed set)

1968

Planner Augmentation (if present):

[REFLECTION] Short hint distilled from latent imagination, emphasizing the safest high-value step accepted by the uncertainty gate. [/REFLECTION]
[SUMMARY] One-sentence sketch of the selected branch (orientation, waypoints, next action). [/SUMMARY]

1973

Instruction:

Goal: go to the red ball 1. Return SUCCESS once you are adjacent and facing it.

1975

Example (round t input):

YOU: position (x=3,y=2), facing north ROOM: walled 7×7

VISIBLE OBJECTS: red ball 1 (two steps ahead, one step left), grey box 1 (one ahead, one left), red key 1 (three ahead, one right), blue door 2 (east wall, closed)

INVENTORY: empty

AVAILABLE ACTIONS: [turn left, turn right, move forward, go to red ball 1, pick up red ball 1, toggle, go through blue door 2]

[REFLECTION] Turn left to face the red ball; move forward twice; then go to red ball 1 to align precisely. [/REFLECTION]

[SUMMARY] turn left → move forward → move forward → go to red ball 1.

Expected Response:

Thought: The red ball is ahead-left; first rotate to face it before advancing.

Action: turn left

1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013

Table 27: Prompt details for DreamPhase (MAZE).

2014 **System Prompt:**
2015 You are an expert maze solver. Each round you receive the goal coordinates, your current
2016 position, and local wall layout. Use ONLY: move up, move down, move left, move
2017 right. Coordinates are (x, y) with down increasing x and right increasing y . Never propose
2018 moves blocked by walls. Prefer shortest safe paths and avoid loops. Respond exactly in the
2019 format:
2020
2021 **Thought:** (brief reasoning grounded in the observation and any planner reflection)
2022 **Action:** (one action from the allowed set)
2023
2024 **Planner Augmentation (if present):**
2025 [REFLECTION] Short hint distilled from latent imagination, emphasizing the safest high-value
2026 path accepted by the uncertainty gate. [/REFLECTION]
2027 [SUMMARY] One-sentence sketch of the selected branch (waypoints/turns). [/SUMMARY]
2028 **Instruction:**
2029 Start a new game and return your action and thought in the format above. Make the optimal move
2030 given the current state.
2031 **Example (round t input):**
2032 GOAL: (8, 6) POSITION: (1, 1) WALLS AROUND YOU: left, up, down
2033 [REFLECTION] Advance right along row 1 until column 6, then proceed down to row 8; avoid
2034 downward moves at (1, 1) due to a wall. [/REFLECTION]
2035 [SUMMARY] right \rightarrow right $\rightarrow \dots \rightarrow$ right (to $y = 6$), then down to $x = 8$.
2036 **Expected Response:**
2037 Thought: Down and left are blocked; the shortest path begins by moving right along the top row
2038 toward $y = 6$.
2039 **Action:** move right
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064

Table 28: Prompt details for DreamPhase (Wordle).

2065 **System Prompt:**
2066 You are an expert Wordle solver. The hidden word has 5 letters and you have 6 attempts. On each
2067 turn, output *one* guess as a space-separated sequence of letters (e.g., s h i r e). After each
2068 guess, the environment returns feedback over 5 letters using b (absent), y (present but wrong
2069 position), and g (correct position). Use only valid words from the vocabulary. Update constraints
2070 from feedback and propose the next guess. Respond exactly in the format:
2071
2072 **Thought:** (brief reasoning grounded in the observation and any planner reflection)
2073 **Guess:** (l e t t e r s with spaces)
2074
2075 **Planner Augmentation (if present):**
2076 [REFLECTION] Short hint distilled from latent imagination, emphasizing a high-coverage guess
2077 consistent with constraints; accepted only if the uncertainty gate passes. [/REFLECTION]
2078 [SUMMARY] One-sentence sketch of the best branch (fixed positions, excluded letters, candidate
2079 pattern). [/SUMMARY]
2080 **Instruction:**
2081 Start a new game. Return your thought and guess in the format above. Guesses must be valid
2082 5-letter words in the vocabulary.
2083 **Example (round t input):**
2084 HISTORY:
2085 1) Guess: s h i n e Feedback: b b b b g (only the final e is correct in
2086 place)
2087 2) Guess: c l o n e Feedback: b g b b g (l fixed at position 2; e fixed at
2088 position 5)
2089 CONSTRAINTS: Pattern _l__e; exclude {s, h, i, n, c, o}.
2090 [REFLECTION] Keep l at position 2 and e at position 5; avoid excluded letters; prefer a guess
2091 that covers diverse consonants and a vowel like u. [/REFLECTION]
2092 [SUMMARY] Try pattern _l__e; candidate: f l u k e.
2093 **Expected Response:**
2094 **Thought:** The constraints require _l__e with l at position 2 and e at 5; f l u k e fits and
2095 introduces u, k.
2096 **Guess:** f l u k e
2097
2098
2099
2100
2101
2102
2103
2104
2105

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115

Table 29: Prompt details for DreamPhase (BIRD / Text-to-SQL).

2116 **System Prompt:**
2117 You are a read-only SQL assistant for a SQLite database. Given a schema description and a
2118 natural-language question, explain your reasoning briefly and then output *one* SQL statement
2119 that answers the question.
2120 **Rules:** (i) Output **exactly** two fields: **Thought:** and **Action:**. (ii) The **Action** must
2121 be a **single-line** SQL query in Markdown code format. (iii) Do *not* modify data (no
2122 INSERT/UPDATE/DELETE); use SELECT only. (iv) Quote identifiers with spaces using
2123 double quotes (e.g., "gas station id"). (v) Use JOINs/aggregations as needed.
2124

2125 **Planner Augmentation (if present):**
2126 [REFLECTION] Short hint distilled from latent imagination (e.g., join keys, filters, grouping),
2127 accepted only if the uncertainty gate passes. [/REFLECTION]
2128 [SUMMARY] One-sentence sketch of the selected branch (tables to join, predicates, target
2129 aggregate). [/SUMMARY]

2129 **Instruction (schema excerpt):**
2130 Database: debit_card_specializing. Tables:
2131 customers(customerid **PK**, client_segment, currency)
2132 gasstations("gas station id" **PK**, "chain id", country,
2133 chain_segment)
2134 products(product_id **PK**, description)
2135 transactions_1k(transaction_id **PK**, date, time, "customer id",
2136 "card id", "gas station id", "product id", amount, price)
2137 yearmonth("customer id", date, consumption, **PK**("customer id",
2138 date), **FK**("customer id") → customers(customerid))
2139 **Question:**

2139 Among the transactions made at gas stations in the Czech Republic, how many occurred after
2140 2012-01-01?

2141 **Example (round t input):**
2142 [REFLECTION] Join transactions_1k with gasstations on "gas station id";
2143 filter country='Czech Republic' and date > '2012-01-01'; return a count.
2144 [/REFLECTION]
2145 [SUMMARY] transactions_1k \bowtie gasstations → WHERE country='Czech
2146 Republic' & date>'2012-01-01' → COUNT(*)
2147 **Expected Response:**

2147 **Thought:** We need a count over filtered transactions; join on "gas station id" to access
2148 the country field; compare ISO date strings.

2149 **Action:** `SELECT COUNT(*) AS num_tx FROM transactions_1k t
2150 JOIN gasstations g ON t."gas station id" = g."gas station id"
2151 WHERE g.country = 'Czech Republic' AND t.date > '2012-01-01';`
2152
2153
2154
2155
2156
2157
2158
2159