# AGENT2WORLD: A UNIFIED LLM-BASED MULTI-AGENT FRAMEWORK FOR SYMBOLIC WORLD-MODEL GENERATION

**Anonymous authors**Paper under double-blind review

# **ABSTRACT**

Symbolic world models, which formally represent environment dynamics and constraints, are essential for model-based planning. While leveraging large language models (LLMs) to automatically generate these models from natural language has shown promise, existing approaches predominantly rely on scripted workflows that follow predetermined execution paths regardless of intermediate outcomes, often leading to inefficient computations and suboptimal solutions. In this paper, we propose AGENT2WORLD, a novel paradigm that employs autonomous tool-augmented LLM-based agents to generate symbolic world models adaptively. We further introduce AGENT2WORLD<sub>Multi</sub>, a unified multi-agent framework with specialized agents: (i) A Deep Researcher agent performs knowledge synthesis by web searching to address specification gaps; (ii) A Model Developer agent implements executable world models; and (iii) a specialized Testing Team conducts evaluation-driven refinement via systematic unit testing and simulation-based validation. AGENT2WORLD<sub>Multi</sub> demonstrates superior performance across three benchmarks spanning both Planning Domain Definition Language(PDDL) and executable code representations, achieving consistent stateof-the-art results through a single unified framework. By enabling proactive, knowledge-grounded world-model generation, this work opens new possibilities for AI systems that can reliably understand and formalize complex environments.

# 1 Introduction

In recent years, researchers have explored *symbolic world models*, a formal representation of an environment's dynamics and constraints, which is widely used in model-based planning (Guan et al., 2023; LeCun, 2022; Craik, 1967). The task of *symbolic world-model generation* involves automatically synthesizing these models from natural language descriptions, eliminating the need for domain experts to manually design and specify complex rules and dynamics. Large language models (LLMs) (Guo et al., 2025; Zhao et al., 2023; Bai et al., 2023) have made this automation increasingly possible by combining two key capabilities: commonsense knowledge about how the world works, and code generation abilities that formalize this knowledge into executable representations.

As illustrated in Figure 1, prior work in this domain largely follows two paradigms: (i) direct generation of symbolic world models, and (ii) scripted workflows that couple generation with iterative verification and repair. Early exemplars of the latter include Guan et al. (2023) and Hu et al. (2025a), using LLMs to produce Planning Domain Definition Language(PDDL)-based world models. Furthermore, GIF-MCTS (Dainese et al., 2024) combines a code executor with trajectories collected in a real environment to furnish feedback, driving a generate–fix-improve loop that progressively refines the generated MuJoCo-style code world models. While scripted workflows achieve better results than direct generation, they suffer from fundamental limitations: (i) Passive and rigid execution: these methods reactively respond to validation failures through predetermined repair sequences, leading to unnecessary computations when simpler solutions exist or inadequate exploration when complex problems require adaptive strategies; (ii) Knowledge isolation: Most approaches rely solely on the LLM's internal knowledge, lacking mechanisms to access external information when specifications are incomplete or ambiguous. While Guan et al. (2023) leverage human feedback as external knowledge, the labor-intensive nature limits their scalability for large-

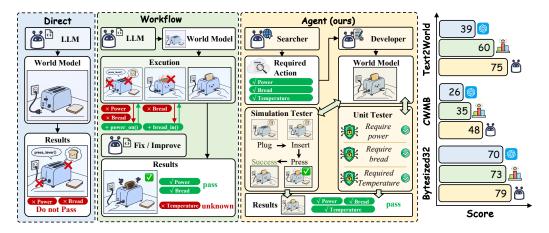


Figure 1: Comparison of AGENT2WORLD and previous world-model generation paradigms.

scale applications. (iii) *Representation fragmentation*: existing approaches typically target either PDDL-based symbolic representations or executable code exclusively, requiring separate systems and workflows for different output formats.

We propose AGENT2WORLD, a new paradigm for symbolic world-model generation that leverages tool-augmented, autonomous LLM-based agents, which plan and call tools adaptively. The key advantages lie in: (i) *Proactive and adaptive execution*. Rather than passively fixing errors through rigid sequences, AGENT2WORLD proactively gathers information and dynamically adjusts strategies based on intermediate feedback, automatically deciding when to terminate and which tools to use for maximum efficiency. (ii) *Scalable external knowledge integration*. Unlike knowledge-isolated approaches or labor-intensive human-in-the-loop methods, AGENT2WORLD incorporate web search as first-class tools to automatically fill specification gaps and enforce commonsense regularities, minimizing LLM hallucination (Huang et al., 2025). (iii) *Unified cross-representation framework*. While prior systems suffer from representation fragmentation, AGENT2WORLD seamlessly handles both PDDL-based and code-based models through lightweight tool adapters, enabling the same agentic framework to work across different symbolic representations.

Specifically, AGENT2WORLD<sub>Single</sub> employs a single ReAct-style agent (Yao et al., 2023) to invoke all available tools (code sandbox, web search, etc.). However, due to the extensive array of available tools and the token-heavy nature of world models (Wang et al., 2023), the single-agent approach faces context length limitations and tool coordination challenges. Furthermore, we propose AGENT2WORLD<sub>Multi</sub> to address these issues, which features role specialization and tool partitioning. Specifically, AGENT2WORLD<sub>Multi</sub> is structured as a three-stage pipeline as shown in Figure 2: (i) Knowledge Synthesis: A Deep Researcher agent proactively identifies knowledge gaps in ambiguous specifications and systematically gathers authoritative information via web search; (ii) World Model Generation: a Model Developer agent generates concrete implementations with iterative code execution and refinement; (iii) Evaluation-Driven Refinement: specialized testing agents evaluate and diagnose the generated models through both simulation for holistic behavioral fidelity and unit testing for component-level correctness, providing feedback to guide iterative refinement.

Our experiments on three benchmarks validate our approach with insights into the fundamental challenges of world model generation: (i) on *Text2World* (Hu et al., 2025a) (PDDL-based world models), AGENT2WORLD<sub>Multi</sub> achieves both high semantic fidelity and syntactic correctness (75.4 macroaveraged F1 with 93.1% executability) (ii) on *CWMB* (Dainese et al., 2024) (code-based world models), it establishes new state-of-the-art with 54.4% predictive accuracy and 48.1% normalized return. Notably, while predictive accuracy improvements are modest, the normalized return gains are dramatic (+13.2% over GIF-MCTS), revealing that our *Evaluation-Driven Refinement* (§3.3) addresses the critical gap between accurate next-state prediction and effective model-based planning; (iii) on *ByteSized32* (Wang et al., 2023) (reasoning-heavy text games), it demonstrates superior performance across technical validity and physical reality alignment (+28.4% alignment score). This validates our hypothesis that external *Knowledge Synthesis* (§3.1) is essential for grounding world models in reliable commonsense knowledge, especially for tasks requiring complex reasoning about everyday physical constraints.

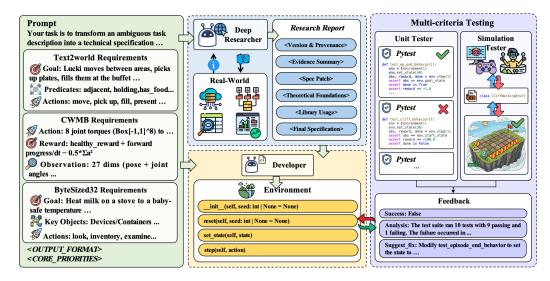


Figure 2: Overall Pipeline of AGENT2WORLD.

# 2 Preliminary

# 2.1 PROBLEM DEFINITION

We investigate the problem of *symbolic world-model generation* from natural language. Given a textual description x, the objective is to synthesize an *executable* program  $\mathcal{M}_F$  that faithfully captures the dynamics and constraints of the environment. Such a program may take various forms, for instance, a specification in the Planning Domain Definition Language (PDDL) (Hu et al., 2025a; McDermott et al., 1998) or an implementation in Python (Dainese et al., 2024; Wang et al., 2023). Formally, an environment is defined by a set of predicates  $\mathcal{P}$ , a set of actions  $\mathcal{A}$ , and a transition function  $T: S \times \mathcal{A} \to S$ , where S denotes the set of possible states. Semantically,  $\mathcal{M}_F$  encodes these components to represent the environment in an executable manner. We therefore define the task as a mapping  $P(x) = \mathcal{M}_F$  where  $\mathcal{M}_F = \langle \mathcal{P}, \mathcal{A}, T \rangle$ , where P is a synthesis procedure that generates the world-model program from natural language input x.

# 2.2 Autonomous Agent

We consider an *autonomous agent* that *interleaves reasoning with tool use* (Qin et al., 2023; Yao et al., 2023). Concretely, the agent firstly produces reasoning traces and issues tool calls in an alternating loop. Let  $\mathcal{T} = \{t_1, \ldots, t_m\}$  denote a set of callable tools. At discrete time k, the agent maintains a history  $h_k = (x, o_{\leq k}, a_{< k}, r_{< k})$  consisting of the task description x, observations o, past actions a, and internal reasoning traces r. A *policy*  $\pi_\theta$  maps histories to the next reasoning trace and action:  $(r_k, a_k) \sim \pi_\theta(\cdot \mid h_k)$ ,  $a_k \in \mathcal{T}$ . In practice,  $\pi_\theta$  is implemented by an LLM.

# 3 METHODOLOGY

As is shown in Figure 2, AGENT2WORLD<sub>Multi</sub> unfolds in three stages: (i) **Knowledge Synthesis** (§3.1): As outlined in Section 1, a key challenge in symbolic world-model construction arises from incomplete descriptions. For example, commonsense knowledge may be missing both from LLMs and from the given specifications. To address this limitation, we employ a *deep researcher agent* that interacts with external resources such as the internet or structured databases, thereby enriching the specification and producing an intermediate representation. (ii) **World Model Generation** (§3.2): At this stage, a *developer agent* equipped with a code execution tool constructs the symbolic world model. The process is iteratively refined based on execution feedback, ensuring both correctness and executability. (iii) **Evaluation-Driven Refinement** (§3.3): We enhance the semantic fidelity by designing two complementary *test agents*: one that generates unit tests to validate functional

behavior, and another that simulates downstream usage to evaluate performance through trajectory-based testing. We also provide a pseudo code in Algorithm 1.

#### 3.1 STAGE I: KNOWLEDGE SYNTHESIS

We introduce a *Deep Researcher* agent designed to gather background knowledge and fill in missing details that are not explicitly provided in the world model description. By leveraging external information sources, this agent not only compensates for potential knowledge gaps inherent in large language models but also enhances the factual reliability of world model descriptions. Equipped with web search and retrieval tools, it iteratively retrieves the knowledge required for world model construction from the internet and ultimately outputs a structured intermediate representation with the missing information completed.

#### 3.2 STAGE II: WORLD MODEL GENERATION

After obtaining the comprehensive world-model description from the previous stage, the *Model Developer* takes this as input and generates a concrete implementation of the world model in the required formalism (e.g., PDDL or executable code). To support iterative refinement, the *Model Developer* is equipped with a sandboxed code-execution tool, enabling it to test and debug implementations in multiple rounds until the code is functional and consistent with the specification.

### 3.3 STAGE III: EVALUATION-DRIVEN REFINEMENT

A key component of our approach is the refinement of a bug-free, code-based world model. Unlike prior works that rely on annotated gold trajectories (Dainese et al., 2024) or human feedback (Guan et al., 2023), our method is fully autonomous and does not require manual labels. More specifically, we introduce a two-agent *Testing Team* to evaluate and diagnose the generated models: (i) The *Simulation Tester* evaluates the world model in a play-testing manner by attempting to perform tasks, explore actions, and issue queries within the environment. Specifically, it interacts with the environment in a ReAct-style (Yao et al., 2023) loop to collect trajectories for subsequent behavior and reward analysis, which uncovers execution-time failures such as unreachable goals, missing preconditions, or inconsistent state updates. (ii) The *Unit Tester* complements play-testing with systematic, programmatic verification. It automatically generates Pytest-style unit tests targeting the predicates, actions, and invariants specified in the world-model descriptions.

Together, these agents produce a detailed *test report* that assesses the quality of the generated world model and provides fine-grained diagnostic signals on correctness, coverage, logical consistency, and compliance with physical requirements. This report is fed back to the *Model Developer*; if inconsistencies or failures are detected, the Model Developer revises the implementation, triggering another evaluation round by both testers. This loop continues until all checks are satisfied or a predefined convergence criterion is reached.

# 4 EXPERIMENTS

In this section, we first describe the baselines (§ 4.1) and implementation details (§ 4.2), and then present experiments on three benchmarks: (i) *Text2World* (Hu et al., 2025a) (§ 4.3): A PDDL-centric benchmark for text-to-symbolic world modelling. (ii) *Code World Models Benchmark (CWMB)* (Dainese et al., 2024) (§ 4.4): A code-based world-model benchmark comprising MuJoCo-style environments, designed to assess *predictive correctness* and *downstream control utility* under both discrete and continuous action settings. (iii) *ByteSized32* (Wang et al., 2023) (§ 4.5): A suite of reasoning-heavy text games requiring executable Python environments. A side-by-side comparison and a detailed metric explanation are shown in Appendix C.

#### 4.1 BASELINES

We compare AGENT2WORLD<sub>Multi</sub> against the following methods:

(i) **Direct Generation** (**Direct**): Single-shot generation of the symbolic world model without tool use, external retrieval, or feedback. (ii) **Agent2World**<sub>single</sub>: A single agent closes the loop by in-

voking code execution/validators/web search tools for self-repair and information synthesis, without multi-agent specialization. (iii) Text2World (EC=k) (Hu et al., 2025a): directly using large language models to generate PDDL-based world model and iteratively repairing with planner/validator signals, where EC denotes the error-correction budget. (iv) WorldCoder (Tang et al., 2024): A plan-code-execute-repair search that scores and iteratively improves candidate programs using simulator/planner signals to select runnable hypotheses. (v) GIF-MCTS (Dainese et al., 2024): A macro-action MCTS that orchestrates Generate/Improve/Fix steps, guided by unit tests and trajectory-based feedback for code world-model synthesis. (vi) ByteSized32 baseline (Wang et al., 2023). The reference pipeline introduced by Wang et al. (2023). In order to avoid metric leakage, we do not use the official checker's evaluation signals.

#### 4.2 IMPLEMENTATION DETAILS

We employ the OpenAI GPT-4.1-mini model via the official API, setting the decoding temperature to 0 and top\_p to 1 for deterministic reproducibility. All agents operate within a ReAct (Yao et al., 2023) framework, following a "think → act (tool) → observe" loop for a maximum of 10 turns. The *Deep Researcher* agent utilizes the Serper API for web searching. We blocked some websites to ensure experimental integrity and prevent information leakage <sup>1</sup> Regarding the configuration of refinement turns, we set Text2World and ByteSized32 to 2 iterations and CWMB to 3 iterations based on the complexity of environments. For automated evaluation on the ByteSized32 benchmark, we leverage GPT-40 (Hurst et al., 2024) as the LLM evaluator. All experiments are conducted on a CPU server without GPU acceleration. The prompt examples could be found at Appendix F.

# 4.3 Text2World

Table 1: Benchmark results on Text2World (Hu et al., 2025a). Following the reporting convention in Text2World, all metrics are presented as percentage scores (%).

Methods	Executability (†)	) Similarity (↑) Component-wise F1 (↑)			F1 <sub>AVG</sub> (†)		
		Z	$F1_{PRED}$	$F1_{PARAM}$	$F1_{\text{PRECOND}}$	$F1_{\text{EFF}}$	AvG ( )
Text2World <sub>EC=3</sub>	78.2	81.1	73.4	64.5	49.3	53.3	60.1
Direct Generation	45.5	82.8	45.0	40.3	33.9	34.9	38.5
AGENT2WORLD <sub>Single</sub>	79.2	82.5	76.5	75.8	60.1	66.0	69.6
AGENT2WORLD <sub>Multi</sub>	93.1	81.0	87.2	82.3	63.7	68.2	75.4

We evaluate the Planning Domain Definition Language (PDDL)-based world model generation of AGENT2WORLD on Text2World Hu et al. (2025a), which comprises 103 PDDL domains paired with natural language descriptions. The evaluation metrics are: (i) **Executability**: whether the generated PDDL can be parsed and validated; (ii) **Structural Similarity**: the normalized Levenshtein similarity; (iii) **Component-wise F1**: the macro-averaged F1 of predicates (F1<sub>PRED</sub>) and action components, including parameters (F1<sub>PARAM</sub>), preconditions (F1<sub>PRECOND</sub>), and effects (F1<sub>EFF</sub>).

**Results.** We can draw several conclusions from Table 1: (i) Direct Generation attains the highest Similarity (82.8) yet performs poorly on executability (45.5) and all component-wise F1s, underscoring that surface-level textual overlap is a weak proxy for runnable, semantically correct PDDL. (ii) While agent-based methods achieve executability comparable to the reference solution (e.g., AGENT2WORLD<sub>Single</sub> with 79.2 vs. Text2World<sub>EC=3</sub> with 78.2), they exhibit substantial gaps in F1 scores (AGENT2WORLD<sub>Single</sub>: 69.6 vs. Text2World<sub>EC=3</sub>: 60.1). This suggests that while integrating validators for iterative correction can significantly improve syntactic validity, the semantic utility of the generated world models remains limited without comprehensive knowledge synthesis. (iii) AGENT2WORLD<sub>Multi</sub> achieves both the highest executability (+14.9 points over Text2World<sub>EC=3</sub>) and superior F1 performance (+15.3 points), demonstrating the synergistic benefits of multi-agent specialization. These patterns align with our design philosophy: knowledge synthesis combined with evaluation-driven refinement steers the model to recover the correct predicate inventory and logical gating constraints, producing domains that are both syntactically valid and semantically solvable, even when surface-level representations diverge from reference implementations.

<sup>&</sup>lt;sup>1</sup>For example, the original Text2World and ByteSized32 huggingface pages, CWMB source code, OpenAI-Gym code repository are blocked.

Table 2: Benchmark results on CWMB. † We adopted the official implementation of GIF-MCTS (Dainese et al., 2024) and their reimplementation of WorldCoder (Tang et al., 2024).

Method	Discrete Actio	on Space	<b>Continuous Action Space</b>		Overall		
1/20/2004	Accuracy (†)	$\mathcal{R}\left(\uparrow ight)$	Accuracy (†)	ccuracy $(\uparrow)$ $\mathcal{R}(\uparrow)$		$\mathcal{R}\left(\uparrow ight)$	
WorldCoder <sup>†</sup>	0.9024	0.5399	0.3303	0.2097	0.5210	0.3197	
GIF-MCTS <sup>†</sup>	0.9136	0.6842	0.2748	0.1811	0.4877	0.3488	
Direct Generation	0.7321	0.4527	0.3038	0.1666	0.4466	0.2620	
AGENT2WORLDSingle	0.7897	0.5418	0.1917	0.2420	0.3911	0.3419	
AGENT2WORLD <sub>Multi</sub>	0.9174	0.8333	0.3575	0.3050	0.5441	0.4811	

The CWMB (Dainese et al., 2024) evaluates the ability of generated executable code to serve as faithful world models across 18 MuJoCo-style environments. It measures both the predictive accuracy of next-state dynamics and the normalized return ( $\mathcal{R}$ ) when the model is used by a planner, where  $\mathcal{R}$  reflects the gap between a random policy and an oracle planner with the true environment. This setup ensures CWMB jointly assesses correctness of the simulation code and its practical utility for downstream control.

**Results.** Table 2 reveals several key findings. (i) All methods demonstrate superior performance in discrete spaces compared to continuous settings, reflecting the inherent difficulty of modeling continuous dynamics. (ii) Workflow-based approaches consistently outperform both *Direct Generation* and AGENT2WORLD<sub>Single</sub>, indicating that LLMs' native world model generation capabilities are limited and require expert-designed iterative refinement to achieve competitive performance. (iii) AGENT2WORLD<sub>Multi</sub> establishes new state-of-the-art results, surpassing the previous best method GIF-MCTS by  $+0.132~\mathcal{R}$  points in overall normalized return. Notably, while other methods achieve comparable predictive accuracy (e.g., 0.917 vs 0.914 on discrete settings), our simulation-based testing framework significantly enhances the downstream utility of generated world models, demonstrating that accurate next-state prediction alone is insufficient for effective model-based planning.

# 4.5 BYTESIZED32

Table 3: Technical Validity and Physical Reality Alignment scores on ByteSized32

Method		Technical Validit	y (†)	Physical Reality Alignment (†)
	Game Init.	Possible Actions	Runnable Game	Alignment Score
ByteSized32 <sub>0-shot</sub>	0.9792	0.9375	0.7292	0.0600
ByteSized32 <sub>1-shot</sub>	0.9792	0.8958	0.7500	0.1748
Direct Generation	0.9271	0.8854	0.7604	0.0000
AGENT2WORLD <sub>Single</sub>	0.9792	0.9375	0.7708	0.1920
AGENT2WORLD <sub>Multi</sub>	0.9896	0.9583	0.8958	0.4768

Table 4: Specification Compliance and Winnability scores on ByteSized32

Method	Specific	Specification Compliance $(\uparrow)$				
	Critical Objects	Critical Actions	Distractors	Winnable Game		
ByteSized32 <sub>0-shot</sub>	0.9375	0.9375	0.8750	0.0625		
ByteSized32 <sub>1-shot</sub>	1.0000	0.9375	0.9375	0.1354		
Direct Generation	1.0000	0.9375	0.9375	0.1354		
AGENT2WORLDSingle	1.0000	1.0000	0.8438	0.1354		
AGENT2WORLD <sub>Multi</sub>	1.0000	0.9688	0.8750	0.1458		

The ByteSized32 (Wang et al., 2023) benchmark consists of 32 reasoning-heavy text games, each implemented as an executable Python environment. Models are required to generate runnable game

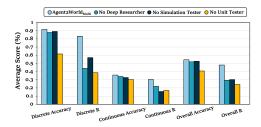




Figure 3: Ablation Study on CWMB.

Figure 4: Pair-wise Win-Tie-Loss analysis.

code that captures task-specific dynamics, objects, and rules, allowing direct interaction and evaluation. The benchmark evaluates four dimensions: **Technical Validity** (whether the code runs), **Specification Compliance** (whether all required elements are present), **Winnability** (whether the task can be completed), and **Physical Reality Alignment** (whether the environment dynamics are consistent with commonsense constraints). This setting emphasizes both logical fidelity and practical executability, making it a stringent testbed for language models as world-model generators.

**Results.** Several conclusions could be drawn from Table 3 and 4: (i) The official reference pipeline outperforms direct generation with in-context learning and shows comparable performance to AGENT2WORLD<sub>Single</sub> on certain metrics. (ii) The AGENT2WORLD<sub>Single</sub> baseline shows moderate gains in game solvability, yet its alignment with physical reality is slightly weaker. (iii) AGENT2WORLD<sub>Multi</sub> outperforms both baselines across almost all dimensions, especially +0.2848 physical alignment score, which stems from *Deep Researcher* agent synthesizing the commonsense knowledge required for reasoning-heavy games.

# 5 ANALYSIS

# 5.1 ABLATION STUDY

Ablations from Figure 3 and Table 8 clarify where the lift comes from: (i) Removing the *Unit Tester* causes the most significant performance drop, with accuracy declining by 0.3008 and reward  $\mathcal{R}$  by 0.4470 in discrete action spaces. (ii) The *Deep Researcher* primarily impacts reward  $\mathcal{R}$  quality, showing a substantial decrease of 0.3926 for discrete spaces when removed. (iii) Although the removal of *Simulation Tester* results in the smallest overall performance drops, the reward  $\mathcal{R}$  decreases by 0.2615 and 0.1473 for discrete space and continuous space, respectively. These results collectively validate our design choices and highlight the complementary nature of the three components.

# 5.2 Pair-wise Evaluation

To quantify the effect of AGENT2WORLD<sub>Multi</sub> and the refinement procedure, we perform instance-level pairwise comparisons, recording a Win–Tie–Loss (WTL) outcome according to the benchmarks' primary metric: (i) F1<sub>AVG</sub> for Text2World; (ii)  $\mathcal{R}$  for CWMB; and (iii) the mean of all official metrics for ByteSized32. As shown in Figure 4, the left panel contrasts the final-turn model with its first-turn counterpart. Refinement yields consistent gains on CWMB and ByteSized32 (68.8% and 93% wins, respectively; no losses), largely preserves performance on Text2World while delivering occasional improvements (14% wins vs. 7% losses). The right panel compares AGENT2WORLD<sub>Multi</sub> against previous state-of-the-art systems. AGENT2WORLD<sub>Multi</sub> attains clear advantages across all three benchmarks, most notably on ByteSized32 (87% wins) and CWMB (66.7% wins).

# 5.3 FEEDBACK ITERATION

To understand the dynamics of performance improvement through iterative feedback, we analyze how model performance evolves as the number of testing team feedback iterations increases. Figure 5 illustrates the relationship between feedback iteration count and model performance across the evaluation benchmarks. The results reveal several key patterns: (i) Text2World shows rapid initial

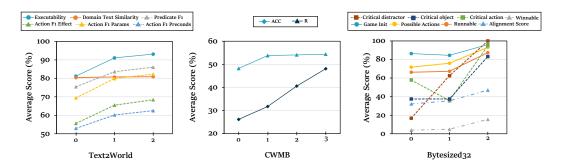


Figure 5: Performance evolution of iterative refinement (§ 3.3.)

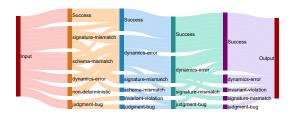




Figure 6: Error distribution on CWMB of evaluationdriven refinement. Due to page limit, the error distribution of other benchmarks is presented in Appendix E.

Figure 7: Performance comparisons and specialized efficiency of multi-agent and single-agent architectures.

improvements. Notably, execution-based metrics improve substantially while similarity measures remain stable, suggesting that refinement enhances functional correctness rather than surface-level similarities. (ii) CWMB demonstrates sustained improvement across iterations, reflecting the compound complexity of physics simulation where numerical accuracy and dynamics must be jointly optimized. (iii) ByteSized32 exhibits the most dramatic gains, with several metrics showing step-function improvements that reflect the discrete nature of game logic debugging.

#### 5.4 MANUAL ERROR ANALYSIS

We conducted a manual error analysis to examine the evolution of error patterns throughout the refinement process of AGENT2WORLD<sub>Multi</sub>. Taken CWMB in Figure 6 as an example, the initial turn predominantly exhibits superficial errors like *signature-arity mismatches* and *representation mismatches*, stemming from inadequate adherence to world model specifications. Throughout the iterative refinement process, these surface-level inconsistencies are systematically eliminated, with the error landscape shifting toward more fundamental *dynamics mismatches* in later iterations. This pattern demonstrates remarkable consistency across all benchmarks: refinement consistently shifts the error distribution from form-oriented problems (*syntax*, *arity*) to substance-oriented challenges (*dynamics*, *state transitions*) as shown in Figure 8 and 9. The systematic progression from surface to substance reflects the hierarchical nature of world model correctness and validates our multi-turn refinement architecture. We also provide the detailed proportion of each error type in Appendix E.

# 5.5 Multi-Agent vs. Single-Agent Architecture Analysis

To quantify the benefits of multi-agent specialization, we compare AGENT2WORLD<sub>Multi</sub> against AGENT2WORLD<sub>Single</sub> using each benchmark's primary metric as in Section 5.2, alongside specialization efficiency calculated as  $\frac{\text{Performance Improvement (\%)}}{\text{Token Cost Increase (\%)}}$ . Figure 7 reveals an efficiency pattern: 4.6% for Text2World, 10.5% for CWMB, and 30.2% for ByteSized32. This trend correlates with the *marginal cost structure* of each benchmark, where ByteSized32's higher efficiency stems from its token-intensive world models, where specialization enables more targeted improvements per unit of computational investment. In contrast, Text2World's compact PDDL representations limit the absolute token overhead of specialization, resulting in lower efficiency despite meaningful perfor-

mance gains. The results suggest that *multi-agent architectures yield disproportionate returns when* world models are computationally expensive, as specialized agents can more effectively amortize their coordination costs across complex, token-heavy implementations.

# 6 RELATED WORK

432

433

434

435 436

437 438 439

440

441

442

443

444

445

446

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474 475 476

477 478

479

480

481

482

483

484

485

World Models. World models are widely applied in reinforcement learning, robotics, and autonomous systems for planning, etc (Hao et al., 2023; Ha & Schmidhuber, 2018). Generally, there are two types of world models: (i) neural world models, which employ neural networks to approximate dynamics (Ha & Schmidhuber, 2018; Hafner et al., 2019), and (ii) symbolic world models, which are represented using formal languages such as the Planning Domain Definition Language (PDDL) or code-based implementations. In this paper, we investigate symbolic world-model generation, which involves transforming world model descriptions into formal representations that can subsequently be utilized for applications such as model-based planning (Guan et al., 2023; Dainese et al., 2024), dataset construction (Hu et al., 2025b), and so on. Most prior work follows a draftrepair workflow where the model first proposes an initial implementation, then gradually refines under closed-loop diagnosis from some kind of feedback. The feedback mechanisms can vary significantly across different approaches. For instance, Guan et al. (2023) employs human feedback to provide corrective signals and perform iterative modifications. Other works, such as Hu et al. (2025a) and Tang et al. (2024), utilize executors and validators to generate feedback. GIF-MCTS (Dainese et al., 2024) leverages gold experiences as feedback signals. Compared to these scripted workflows, the AGENT2WORLD paradigm introduced in this paper can more flexibly adjust subsequent strategies based on feedback signals. Furthermore, benefiting from the scalability of AGENT2WORLD, we incorporate additional external feedback mechanisms such as web search to supplement the intrinsic knowledge limitations of the underlying models. A side-by-side of AGENT2WORLD and existing methods can be found in Appendix B.1.

Large Language Model-based Agent. In recent years, benefiting from the rapid advancement of large language models (LLMs), LLM-based agents have emerged as particularly powerful systems that accept natural language user intentions as input and achieve goal states through planning and sequential decision-making (Hu et al., 2024; Yao et al., 2023; Schick et al., 2023). These autonomous agents have demonstrated remarkable effectiveness across diverse applications, ranging from web navigation (Yao et al., 2022; Nakano et al., 2021; Wang et al., 2025) and software development (Qian et al., 2023; Hong et al., 2024) to scientific research (Lu et al., 2024; Chen et al., 2025) and robotic planning (Huang et al., 2022). Prominent examples of such systems include ReAct (Yao et al., 2023), which synergizes reasoning and acting in language models by interleaving thought, action, and observation steps; Existing research has explored how world models can assist LLM-based agents in planning, such as RAP (Hao et al., 2023), which uses Monte Carlo Tree Search with world models for improved reasoning, and Guan et al. (2023), which leverages pre-trained LLMs to construct world models for model-based task planning. These approaches primarily focus on *utilizing* existing world models rather than generating them. Similarly, recent work has investigated how world models can enhance training of LLM-based agents, as demonstrated by AgentGen (Hu et al., 2025b) and Kimi-K2 (Team et al., 2025). To our best knowledge, our work represents the first systematic investigation into using autonomous agents for world model generation, bridging the gap between agent-based problem solving and symbolic world modeling.

# 7 Conclusion

We introduced AGENT2WORLD<sub>Multi</sub>, a unified multi-agent framework that employs autonomous LLM-based agents to generate symbolic world models across both PDDL and executable code representations. The framework operates through three specialized stages: knowledge synthesis via web search, world model development with iterative refinement, and evaluation-driven testing through unit tests and simulation. Experimental results demonstrate consistent state-of-the-art performance across three world-model generation benchmarks of different types. By enabling fully autonomous world model generation without human feedback or manual annotations, this work opens new possibilities for AI systems that can reliably understand and formalize complex environments from natural language.

# **ETHICS STATEMENT**

All authors have read and will adhere to the ICLR Code of Ethics. This work does not involve human subjects, personal data, demographic attributes, or user studies; IRB approval was therefore not required. Our experiments use public, non-sensitive benchmarks: *Text2World*, *Code World Models Benchmark (CWMB)*, and *ByteSized32*. We complied with dataset licenses and did not attempt to deanonymize or enrich any data with personal information. Because AGENT2WORLD uses web search as a tool (§ 3.1), we enforced safeguards to reduce legal and research-integrity risks: we retrieved only publicly accessible pages, implemented a denylist to avoid solution leakage from benchmark source repositories or discussion pages as discussed in Section 4.2; we have no conflicts of interest or undisclosed sponsorship related to this work.

# REPRODUCIBILITY STATEMENT

Implementation details needed to re-create agents, tools, and evaluation are specified in Section 4.2 and Appendix B.2; algorithmic workflow and role specialization are detailed in  $\S 3.2-\S 3.3$  (with pseudo code in Alg. 1); prompts are provided in Appendix F. For datasets and metrics, benchmark compositions and metric definitions are summarized in  $\S C$  and Appendix C.2; ablation settings and additional figures/tables appear in Appendix D. As discussed in Section 4.2, to facilitate exact runs, we fix decoding parameters (temperature = 0, top\_p = 1) and cap agent turns, specify external services (web search API) and denylisted domains to prevent leakage, and report hardware assumptions (CPU-only). We also provide the source code of our methods in supplementary materials.

# REFERENCES

- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Qiguang Chen, Mingda Yang, Libo Qin, Jinhao Liu, Zheng Yan, Jiannan Guan, Dengyun Peng, Yiyan Ji, Hanjing Li, Mengkang Hu, et al. Ai4research: A survey of artificial intelligence for scientific research. *arXiv* preprint arXiv:2507.01903, 2025.
- Kenneth James Williams Craik. The nature of explanation, volume 445. CUP Archive, 1967.
- Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. Generating code world models with large language models guided by monte carlo tree search. *arXiv preprint arXiv:2405.15383*, 2024. URL https://arxiv.org/abs/2405.15383. Introduces the Code World Models Benchmark (CWMB).
- Lin Guan, Karthik Valmeekam, Sarath Sreedharan, and Subbarao Kambhampati. Leveraging pretrained large language models to construct and utilize world models for model-based task planning. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pp. 79081–79094, 2023. URL https://arxiv.org/abs/2305.14909.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- David Ha and Jürgen Schmidhuber. World models. 2018. doi: 10.5281/ZENODO.1207631. URL https://zenodo.org/record/1207631.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*, 2023.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for a multi-agent collaborative framework. International Conference on Learning Representations, ICLR, 2024.
- Mengkang Hu, Tianxing Chen, Qiguang Chen, Yao Mu, Wenqi Shao, and Ping Luo. Hiagent: Hierarchical working memory management for solving long-horizon agent tasks with large language model. *arXiv preprint arXiv:2408.09559*, 2024.
- Mengkang Hu, Tianxing Chen, Yude Zou, Yuheng Lei, Qiguang Chen, Ming Li, Yao Mu, Hongyuan Zhang, Wenqi Shao, and Ping Luo. Text2world: Benchmarking large language models for symbolic world model generation, 2025a. URL https://arxiv.org/abs/2502.13092.
- Mengkang Hu, Pu Zhao, Can Xu, Qingfeng Sun, Jianguang Lou, Qingwei Lin, Ping Luo, and Saravan Rajmohan. Agentgen: Enhancing planning abilities for large language model based agent via environment and task generation. *arXiv preprint arXiv:2408.00764*, 2025b. URL https://arxiv.org/abs/2408.00764. Accepted by KDD 2025 (Research Track).
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55, 2025.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, pp. 9118–9147. PMLR, 2022.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.

- Adam Ishay and Joohyung Lee. Llm+al: Bridging large language models and action languages for complex reasoning about actions, 2025. URL https://arxiv.org/abs/2501.00830.
- Yann LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. *Open Review*, 62(1):1–62, 2022.
  - Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.
  - Drew McDermott, Malik Ghallab, Adele Howe, Craig Knoblock, Ashwin Ram, Manuela Veloso, Daniel Weld, and David Wilkins. Pddl the planning domain definition language (version 1.2). Technical Report CVC TR-98-003 / DCS TR-1165, Yale Center for Computational Vision and Control, 1998. URL https://www.cs.cmu.edu/~mmv/planning/readings/98aips-PDDL.pdf.
  - Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
  - James Oswald, Kavitha Srinivas, Harsha Kokel, Junkyu Lee, Michael Katz, and Shirin Sohrabi. Large language models as planning domain generators, 2024. URL https://arxiv.org/abs/2405.06650.
  - Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *arXiv* preprint arXiv:2307.07924, 2023.
  - Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong, Xiangru Tang, Bill Qian, et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. *arXiv preprint arXiv:2307.16789*, 2023.
  - Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
  - Pavel Smirnov, Frank Joublin, Antonello Ceravola, and Michael Gienger. Generating consistent pddl domains with large language models, 2024. URL https://arxiv.org/abs/2404.07751.
  - Hao Tang, Darren Key, and Kevin Ellis. Worldcoder: A model-based llm agent for building world models by writing code and interacting with the environment. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 37, pp. 70148–70212, 2024.
  - Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv* preprint arXiv:2507.20534, 2025.
  - Peng Wang, Ruihan Tao, Qiguang Chen, Mengkang Hu, and Libo Qin. X-webagentbench: A multilingual interactive web benchmark for evaluating global agentic system. *arXiv* preprint *arXiv*:2505.15372, 2025.
  - Ruoyao Wang, Graham Todd, Xingdi Yuan, Ziang Xiao, Marc-Alexandre Côté, and Peter Jansen. Bytesized32: A corpus and challenge task for generating task-specific world models expressed as text games. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023. doi: 10.18653/v1/2023.emnlp-main.830. URL https://aclanthology.org/2023.emnlp-main.830/.
  - Shunyu Yao, Howard Chen, John Yang, and Karthik Narasimhan. Webshop: Towards scalable real-world web interaction with grounded language agents. *Advances in Neural Information Processing Systems*, 35:20744–20757, 2022.

Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023. URL https://arxiv.org/abs/2210.03629.

Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv preprint arXiv:2303.18223, 1(2), 2023.

# **Appendix**

# A THE USE OF LARGE LANGUAGE MODELS

We used a large language model (LLM) strictly as a general-purpose writing assistant for surface-level editing, such as grammar correction, wording polish, and minor style consistency. The LLM was not used for conceptual ideation, literature review, algorithm or model design, data collection or labeling, experiment setup, result analysis, drafting of technical content.

# B More Details on Methodology

#### **B.1** METHOD COMPARISONS

Table 5: Comparison of AGENT2WORLD and related approaches. **Feedback** stands for whether the method uses execution/checker/test signals during generation. **External Knowledge** stands for whether explicit web/external knowledge retrieval is a first-class component. **Type** represents environment types, where ● and ○ stand for discrete and continuous environments, respectively.

Method	Representation	Core Paradigm	Feedback	External Knowledge	<b>Environment Type</b>
Text2World (Hu et al., 2025a)	PDDL	Static Workflow	Y	N	•
Guan et al. (2023)	PDDL	Static Workflow	Y	Human	•
Oswald et al. (2024)	PDDL	Static Workflow	Y	N	•
Smirnov et al. (2024)	PDDL	Static Workflow	Y	N	•
AgentGen (Hu et al., 2025b)	PDDL	Static Workflow	Y	N	•
WorldCoder (Tang et al., 2024)	Code	Static Workflow	Y	N	• 0
GIF-MCTS (Dainese et al., 2024)	Code	Static Workflow	Y	N	• 0
ByteSized32 Wang et al. (2023)	Code	Static Workflow	Y	N	•
LLM+AL (Ishay & Lee, 2025)	Action Language	Static Workflow	Y	N	•
Direct Generation	PDDL & Code	Direct	N	N	• 0
AGENT2WORLD <sub>Single</sub>	PDDL & Code	Adaptive Agent	Y	N	• 0
AGENT2WORLD <sub>Multi</sub>	PDDL & Code	Adaptive Agent	Y	Internet	• 0

As is shown in Table 5, we compare the related methods.

# **B.2** Per-agent Tool Configuration

Table 6: Per-agent configuration.

Agent	Tools
Deep Researcher	browser_search; browser_open
Model Developer	file_tool; sandbox; run_code
Simulation Tester	play_env; file_tool
Unit Tester	run_code; run_bash; file_tool

Detailed per-agent tool configuration is presented in Table 6.

# B.3 PSEUDO CODE OF AGENT2WORLD<sub>Multi</sub>

We formalize the process of AGENT2WORLD<sub>Multi</sub> in Algorithm 1.

```
702
            Algorithm 1: The execution pipeline of AGENT2WORLD<sub>Multi</sub>
703
            Input: T, N
704
            Output: e
705
             N_r \leftarrow predefined integers;
706
             R \leftarrow \emptyset;
707
             E \leftarrow \emptyset;
708
            Q \leftarrow \text{ExtractQuestions}(T);
709
            for r \leftarrow 1 to N_r do
710
                  q \leftarrow \text{ResearchAgent}(\text{select}, \{Q, E, R\});
                  if q = \emptyset then
711
                   break
712
                  L \leftarrow WebSearch(q);
713
                  E \leftarrow E \cup \text{ResearchAgent}(\text{summarize}, \{L\});
714
                  R \leftarrow \text{ResearchAgent}(\text{update}, \{T, E, R\});
715
            F_t \leftarrow R;
716
            C_{\text{last}} \leftarrow \emptyset;
717
            \quad \textbf{for } n \leftarrow 1 \textbf{ to } N \textbf{ do}
718
                  C_d \leftarrow \text{DevelopAgent}(T, F_t);
719
                  if C_d \neq \emptyset then
720
                       C_{\text{last}} \leftarrow C_d;
721
                       p_{\text{code}} \leftarrow \text{FileTool}(\text{save}, C_d);
722
                  else
723
                   continue
                  C_t \leftarrow \text{UnitTestAgent}(C_d, T, R);
724
                  p_{\text{test}} \leftarrow \text{FileTool}(\text{save}, C_t);
725
                  U_t \leftarrow \texttt{CodeTool}(\texttt{run\_tests}, \{p_{\texttt{code}}, p_{\texttt{test}}\});
726
                  S_t^{\star} \leftarrow \text{PlayEnv}(C_d);
727
                  S_t \leftarrow \text{SimulationTestAgent}(S_t^{\star}, T);
728
                  if U_t.pass \wedge S_t.pass then
729
                        e \leftarrow C_{\text{last}};
730
                       return e;
731
                  F_t \leftarrow \texttt{MergeFeedback}(U_t, S_t);
732
            e \leftarrow C_{\text{last}};
733
             return e
734
```

# C MORE DETAILS ON BENCHMARKS

# C.1 SIDE-BY-SIDE COMPARISON

Table 7: Overview of Text2World (Hu et al., 2025a), Code World Models Benchmark (CWMB) (Dainese et al., 2024), and ByteSized32 (Wang et al., 2023). "Type" denotes the target representation (PDDL vs. executable code). Metrics are shown at the family level. A detailed explanation of each metrics is presented in Appendix C.2.

Benchmark	#Environments	Туре	Metrics (core)
Text2World	103	PDDL	Executability; Domain similarity; F1 scores
CWMB	18	Code (Python)	Accuracy; Normalized return $\mathcal{R}$ (discrete/continuous)
ByteSized32	32	Code (Python)	Technical validity; Specification compliance; Winnability; Physical reality alignment

A side-by-side comparison of the evaluated benchmarks in this paper is presented in Table 7

# C.2 METRIC EXPLANATION

# Text2World

735 736

737 738

739 740

741

742

743

752753754

**Executability.** Name: Exec. Range: [0,1] (higher is better). Whether the generated {domain, problem} can be successfully parsed and validated by standard PDDL validators; reported as the fraction (percentage) over all test cases. Fine-grained metrics below are computed only when Exec=1.

**Domain similarity.** *Name: Sim. Range:* [0, 1] (higher is better). Textual/structural similarity between the generated and gold PDDL measured by a *normalized Levenshtein ratio*.

Let X and Y be the character sequences of the two files with lengths |X| and |Y|, and let Lev(X,Y) denote their Levenshtein distance, then

$$Sim(X,Y) = 1 - \frac{Lev(X,Y)}{max\{|X|,|Y|\}} \in [0,1].$$
 (1)

**F1 scores.** Range: [0,1] (higher is better). When Exec=1, we parse both generated and gold PDDL into structured representations and report macro-averaged F1 for the following components: **Predicates** (F1<sub>PRED</sub>), **Parameters** (F1<sub>PARAM</sub>), **Preconditions** (F1<sub>PRECOND</sub>), and **Effects** (F1<sub>EFF</sub>). We use the standard definition of  $F_1$ , where P and R denote precision and recall, respectively:

$$F_1 = \frac{2 P R}{P + R}$$

# **CWMB**

**Prediction Accuracy.** Symbol: Acc<sub>pred</sub>. Range: [0,1] (higher is better). Definition: We use the same accuracy metric as in the evaluation phase of GIF–MCTS (Sec. 4). Given a validation set  $D = \{(s_i, a_i, r_i, s_i', d_i)\}_{i=1}^N$  and CWM predictions  $(\hat{s}_i', \hat{r}_i, \hat{d}_i) = \text{CWM.step}(s_i, a_i)$ , the accuracy uniformly weights next state, reward, and termination:

$$Acc_{pred} = \frac{1}{N} \sum_{i=1}^{N} \left[ \frac{1}{3} \mathbf{1} (\hat{s}'_{i} = s'_{i}) + \frac{1}{3} \mathbf{1} (\hat{r}_{i} = r_{i}) + \frac{1}{3} \mathbf{1} (\hat{d}_{i} = d_{i}) \right].$$
 (2)

**Normalized Return.** *Symbol:*  $\mathcal{R}$ . *Range:* unbounded (higher is better;  $\mathcal{R} > 0$  means better than random;  $\mathcal{R} \to 1$  approaches the oracle). *Definition:* 

$$\mathcal{R} = \frac{R(\pi_{\text{CWM}}) - R(\pi_{\text{rand}})}{R(\pi_{\text{true}}) - R(\pi_{\text{rand}})},$$
(3)

where  $R(\pi)$  denotes the return. *Protocol:* as in the original setup, we use vanilla MCTS for discrete action spaces and CEM for continuous action spaces;  $R(\cdot)$  is averaged across a fixed number of episodes per environment (10 in the original), and  $R(\pi_{\rm rand})$  uses the environment's random policy baseline.

# ByteSized32

**Technical Validity.** Range: [0,1]. Measured in the order of API calls, such that failure of an earlier function implies failure of subsequent tests. Game initialization is evaluated once at the beginning of the game, whereas GENERATEPOSSIBLEACTIONS() and STEP() are evaluated at *every step*. We check:

- Game initialization: the game/world initializes without errors;
- *Valid actions generation*: the routine that enumerates valid actions for the current state returns without errors (verified via a bounded path crawl);
- Runnable game: a bounded-depth crawl of trajectories executes without errors.

**Specification Compliance.** Range: [0,1]. An LLM acts as the judge for *true/false* compliance against the task specification. The prompt provides the task spec {GAME\_SPEC}, the game code {GAME\_CODE}, and an evaluation question {EVAL\_QUESTION}; the LLM is instructed to first output Yes/No and then a brief rationale. To reduce variance, we use a fixed prompt template and perform multiple independent runs with majority vote/mean aggregation. We report three submeasures: *Task-critical objects*, *Task-critical actions*, and *Distractors*.

**Physical Reality Alignment.** Range: [0,1]. Automatic evaluation proceeds in two stages:

(1) Trajectory generation: perform a breadth-first crawl using the action strings returned by GENERATEPOSSIBLEACTIONS () at each step; actions are grouped by verb (first token) and expanded in a bounded manner. If an error occurs, the error message is recorded as the observation and the search continues.

(2) Sampling and judgment: group paths by the last action verb, draw a fixed-size subsample approximately balanced across groups, and submit each path—together with the task description {GAME\_TASK}—to an LLM for a binary judgment (yes/no; errors are treated as failures). The final score is the fraction judged aligned.

Winnability. Range: [0,1]. A text-game agent (LLM agent) attempts to reach a terminal win within horizon H; we report the fraction of tasks deemed winnable. Given the limited agreement between automatic and human assessments for this metric, we prioritize human evaluation in the main results and use the automatic estimate as auxiliary reference.

# D More Details on Ablation Study

Table 8: Ablation Study of AGENT2WORLD on CWMB (Dainese et al., 2024).

Method	Discrete A	ction Space	Continuous	Action Space	Overall	
Withing	Accuracy (†)		$\mathcal{R}\left(\uparrow\right)$ Accuracy $\left(\uparrow\right)$ $\mathcal{R}\left(\uparrow\right)$		Accuracy (†)	$\mathcal{R}\left(\uparrow ight)$
AGENT2WORLD	0.9174	0.8333	0.3575	0.3050	0.5441	0.4811
No Deep Researcher No Simulation Tester No Unit Tester	$\begin{array}{c} 0.8794_{-0.0380} \\ 0.8920_{-0.0254} \\ 0.6166_{-0.3008} \end{array}$	$\begin{array}{c} 0.4407_{-0.3926} \\ 0.5718_{-0.2615} \\ 0.3863_{-0.4470} \end{array}$	$\begin{array}{c} 0.3404_{-0.0171} \\ 0.3288_{-0.0287} \\ 0.3025_{-0.0550} \end{array}$	$\begin{array}{c} 0.2201_{-0.0849} \\ 0.1577_{-0.1473} \\ 0.1704_{-0.1346} \end{array}$	0.5201 <sub>-0.0240</sub> 0.5275 <sub>-0.0166</sub> 0.4072 <sub>-0.1369</sub>	$0.2936_{-0.1875}$ $0.3039_{-0.1772}$ $0.2423_{-0.2388}$

Detailed experimental results of the ablation study are presented in Table 8.

# E More Details on Error Analysis

# E.1 ERROR ANALYSIS ON TEXT2WORLD AND BYTESIZED32

We visualize the error patterns during evaluation-driven refinement on Text2World and ByteSized32 in Figure 8 and Figure 9.

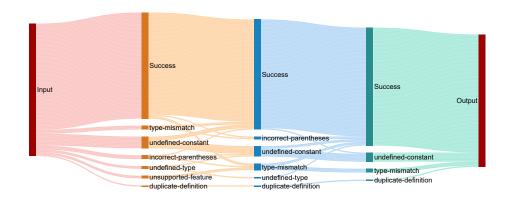


Figure 8: Error distribution of AGENT2WORLD<sub>Multi</sub> on Text2World.

# E.2 DISTRIBUTION OF ERROR TYPES

A detailed proportion of error types on Text2World, CWMB, ByteSized32 are presented in Table 9, Table 10 and Table 11, respectively.

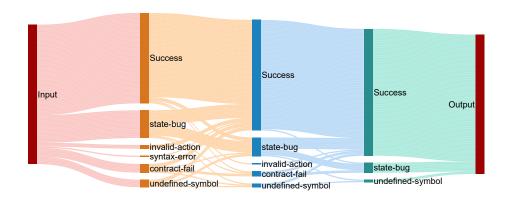


Figure 9: Error distribution of AGENT2WORLD<sub>Multi</sub> on ByteSized32.

Table 9: Distribution of Syntax Errors in Text2World Across Turns

Error Type	Explanation	Turn 0 (%)	Turn 1 (%)	Turn 2 (%)
undefined-constant	Reference to undeclared constants in predicates or actions.	8.91	7.92	6.93
type-mismatch	Parameter type conflict with declared type constraints.	2.97	4.95	2.97
incorrect-parentheses	Invalid or mismatched parentheses.	2.97	1.98	0.00
undefined-type	Undeclared parent type in hierarchical type definitions.	1.98	0.99	0.00
unsupported-feature	Parser-incompatible features (e.g.,either types).	1.98	0.00	0.00
duplicate-definition	Multiple declarations of identical domain elements.	0.99	0.99	0.99

Table 10: Distribution of Syntax Errors in CWMB Across Turns

Error Type	Explanation		Turn 1 (%)	Turn 2 (%)	Turn 3 (%)
signature-mismatch	Arity/types do not match the declared signature.	27.78	11.11	11.11	5.56
schema-mismatch	Value type/shape/dtype violates the expected schema.	22.22	5.56	0.00	0.00
dynamics-error	State or reward deviates from expected dynamics.	11.11	44.44	33.33	11.11
non-deterministic	Results are inconsistent under fixed conditions.	11.11	0.00	0.00	0.00
judgment-bug	Environment setup inconsistent with the description.	11.11	5.56	11.11	5.56
invariant-violation	Internal invariants are broken (e.g., illegal config).	0.00	5.56	0.00	5.56

Table 11: Distribution of Syntax Errors in ByteSized32 Across Turns

Error Type	Explanation	Turn 0 (%)	Turn 1 (%)	Turn 2 (%)
state-bug	State inconsistency across steps.	19.82	13.51	7.21
contract-fail	Task/API contract not satisfied.	6.31	3.60	0.00
undefined-symbol	Reference to undeclared name/type/domain/constant.	5.41	2.70	1.80
invalid-action	Unknown or unsupported action not safely handled.	2.70	0.90	0.00
syntax-error	Load/parse failure (e.g., null bytes, syntax/indentation errors).	0.90	0.00	0.00

919 920

921

922

923

924

925

926

927

928

929

930

931

932

933 934

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956 957

958

959

960

961

962

963

964

965

966 967

968

969 970

# F PROMPT EXAMPLES

# **Deep Researcher**

You are a world-class Systems Analyst and Technical Specification Writer, specializing in creating reinforcement learning environments compliant with the Gymnasium API. Your mission is to transform an ambiguous task description into a precise, actionable, and verifiable technical specification.

<Environment Name>

#### CliffWalking-v0

</Environment Name>

# <TASK DESCRIPTION>

Cliff walking involves crossing a gridworld from start to goal while avoiding falling off a cliff. ## Description ...

# </TASK DESCRIPTION>

#### <Workflow>

Please strictly follow the following six-step process:

# Deconstruction and Analysis (Use Version Locking)

- Identify all ambiguities, gaps, and conflicts in the task description.
- Lock the exact environment version and all key library versions (record name, version, and source link).
- Categorize gaps by type: missing value/unit/range boundary/time- sensitive/ambiguous reference/unclosed list/conflict/no provenance.

# Planning and Investigation (Authoritative Search + Evidence Log)

- For each high/mid-level project, write 1–2 focused queries that include: synonyms/abbreviations, site filters, authoritative domains (e.g., site:numpy.org, site:mujoco.readthedocs.io, site:doi.org), and recency windows (e.g., after 2024-01-01 or "last 2 years").
- Execute the query using browser\_search and open >= 2 trusted results with browser\_open.
- If the top source disagrees, open >= 1 additional authoritative sources and triangulate.
- Create an evidence log entry for each opened page: Title | Organization/Author
   | Version/Submission | URL (+ archived URL) | Publication Date | Access Date
   (Asia/Singapore) | 3 Key Facts | Confidence (High/Medium/Low).

### Synthesis and Citation (Conflict Resolution)

- Integrate the findings into a concise evidence summary with citations.
- When sources conflict, explain the differences and justify the chosen resolution (related to version locking).

# Refinement and Improvement (Specification Patch)

Generate a structured "diff": action/observation space; rewards; termination/truncation; timing (dt/frame\_skip); seeding and certainty; numerical tolerances; dependencies; interface flags.

# • Formalization and Finalization (Ready-to-Use Specification)

 Write the final specification according to the <Output Format>, including the public API, core logic, usage scenarios, and a verification plan aligned with metrics and statistical validation.

# Review and Self-Correction (Compliance Check)

Verify conformance to the <Output Constraints> (OUT-PUT\_CONSTRAINTS>), version consistency, SI units, ISO dates, and the inclusion of any code.

# </Workflow>

### <OUTPUT\_CONSTRAINTS>

980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1012
1013
1017 1018
1019
1020
1021
1022
1023
1024
1025

973 974

975 976

977

978

979

- Strictly adhere to the structure defined in <PLANNING\_STRUCTURE>.
- Do **NOT** output runnable code definitions (classes, functions). Only may include short illustrative snippets or pseudo-code.
- All claims about industry standards or common practices MUST be supported by citations.
- Use ISO-8601 dates (e.g., 2025-09-02).
- Use SI units for physical and mathematical quantities.
- Data-leakage rule: Do not access, copy, quote, or derive from raw source code in the OpenAI/Gym/Gymnasium repositories or similar code repositories. Do not include any repository code in the output. Prefer official documentation, standards, papers, or reputable secondary sources. If the only available evidence is a code repository, summarize behavior without copying code and mark it as an inference with risks.

# <PLANNING\_STRUCTURE>

- Your output must begin with this planning and analysis section.
- Ambiguity Analysis
  - List each ambiguity/vagueness/conflict and mark Impact: High / Medium / Low.
  - Cover at least: missing numeric value, missing unit, missing boundary/range, time-sensitive items, unclear references, open lists ("etc."/"e.g."), conflicts, and missing citation.

# Investigation Plan

- For each High/Medium item, provide one atomic question.
- For each question, provide 1-2 executable queries including: synonyms/abbreviations, a site filter to authoritative domains, and a time window (e.g., after: 2024-01-01 or "past 2 years").
- State the minimum evidence policy: High/Medium → ≥ 2 credible sources; if disagreement → add ≥ 1 more for triangulation.

# <Formula requirements>

- For any formula, define all symbols, units, and applicability constraints.
- Cite the source of the formula immediately after its definition.
- Provide the complete formula rather than a descriptive explanation.

# <OUTPUT\_FORMAT>

- Please provide the final specification document structured as follows. This is the primary deliverable. Do **NOT** include code.
- <Version & Provenance>
- <Evidence Summary>
- <Spec Patch>
- <Theoretical Foundations>
- <Final Specification>
- <Assumptions & Risks>
- <Third-Party Library Usage>

# **Model Developer**

# <TASK DESCRIPTION>

Cliff walking involves crossing a gridworld from start to goal while avoiding falling off a cliff. ## Description ...

# </TASK DESCRIPTION>

026	<research report=""></research>
027	# CliffWalking-v0 Environment Specification
028	Version & Provenance
029 030	Environment: CliffWalking-v0
031	• Gymnasium version: v0.26.3 (2022-09-15) and latest documentation snapshot (2025-
)32	01-01)
33	• Scope: This specification normalizes the CliffWalking environment as implemented
34	in Gymnasium v0.26.3, aligning with the classic gridworld cliff walking task.
35	Accessed date: 2024-06-01 (Asia/Singapore)
6 7	Evidence Summary
	• The environment is a 4x12 grid with 48 discrete states indexed by row-major flatten-
	ing: state = $row * 12 + col$ .
	• Valid states exclude cliff cells ([3,110]) and the goal ([3,11]). The player can occupy all cells in the first 3 rows plus the bottom-left cell [3,0].
	• The start state is 36 (row=3, col=0), and the goal state is 47 (row=3, col=11).
	• The action space is Discrete(4) with actions: 0 (up), 1 (right), 2 (down), 3 (left).
	• Each step yields a reward of -1. Stepping into the cliff yields a reward of -100 and
	resets the player to the start state; the episode continues.
	• The episode terminates only when the player reaches the goal state.
	• The set_state method is not explicitly documented but is inferred to accept any valid
	non-terminal, non-cliff state and sets the environment to not done.
	• No explicit error handling for invalid actions is documented; valid inputs are as-
	sumed.
	• The environment is deterministic.
	Spec Patch
	• action.space = Discrete(4), values {0,1,2,3}, shape (1,)
	• observation.space = Discrete(48), integer in [0,47], representing flattened grid index:
	row * 12 + col
	• valid states = all grid cells except cliff cells ([3,110]) and goal ([3,11])
	• start_state = 36 (row=3, col=0)
	• goal_state = 47 (row=3, col=11)
	• reward.step = -1 per step
	• reward.cliff = -100 on stepping into cliff
	• episode termination = True if state == goal_state; else False
	<ul> <li>stepping into cliff resets player to start_state, episode continues</li> </ul>
	• set_state(state: int) sets environment state to given valid state, done = False
	<ul> <li>no maximum episode length or truncation condition</li> </ul>
	<ul> <li>no explicit error handling for invalid actions; assume valid input</li> </ul>
	• environment deterministic transitions
	Theoretical Foundations
	MDP Definition:
	- State space $S = \{0, \dots, 47\} \setminus \{\text{cliff states, goal state}\}.$
	- Action space $A = \{0, 1, 2, 3\}$ .
	- Transition function $T(s, a, s')$ (deterministic):
	* If $s'$ is a cliff cell, next state = start_state. * Else next state = $s + \delta(a)$ (with grid boundary checks).
	- Libe next state — o   0 (a) ( with gird boulluly clicks).

```
1080
                        - Reward function R(s, a, s'):
1081
                            * R = -100 if s' is a cliff cell.
1082
                            * R = -1 otherwise.
1083
                        - Episode ends when s' = \text{goal\_state}.
1084
1085

    Symbol Table:

1086
                        - s: state (int), flattened grid index, 0 \le s \le 47, s not in cliff or goal.
1087
                        - a: action (int), in \{0 : \text{up}, 1 : \text{right}, 2 : \text{down}, 3 : \text{left}\}.
1088
                        - r: reward (float), -100 or -1.
1089
                        - done: boolean, True if s = goal\_state.
1090
                        - start_state: 36 (int).
1091
                        - goal_state: 47 (int).
1092
                        - cliff_cells: set of ints corresponding to grid positions [3, 1..10].
1093
1094
                        - nrows: 4 (int).
1095
                        - ncols: 12 (int).
1096
                   • Assumptions:
                        - No shaping rewards.
1098
                        - Deterministic environment.
1099
                        - No truncation or max step limit.
1100
1101
            Final Specification
1102
                   • Environment Name
1103
                        - CliffWalking-v0
1104
                   · Public API
1105
1106
                        - __init__()
1107
                        - set state(state: int)
1108
                        - step(action: int) \rightarrow (observation: int, reward: float, done: bool)
1109
                   • Core Logic Description
1110
                        - init :
1111
                            * Initialize grid size (4 rows \times 12 columns).
1112
                            * Define start_state = 36 (row=3, col=0).
1113
                            * Define goal state = 47 (row=3, col=11).
1114
                            * Define cliff cells as positions [3, 1..10].
1115
                            * Initialize current state to start state.
1116
                            * Initialize done flag to False.
1117
                        - set_state(state: int):
1118
                            * Validate that state is a valid observation (not cliff or goal).
1119
                            * Set current state to state.
1120
                            * Set done flag to False.
1121
1122
                        - step(action: int) \rightarrow (obs, reward, done):
1123
                            * Validate action is in \{0, 1, 2, 3\}.
1124
                            * Compute next position based on action with boundary checks:
1125
                                \cdot 0: move up (row -1)
1126
                                • 1: move right (col +1)
1127
                                \cdot 2: move down (row +1)
1128
                                \cdot 3: move left (col -1)
1129
                            * If next position is outside grid, remain in current position.
1130
                            * If next position is a cliff cell:
1131
                                \cdot reward = -100
1132
                                \cdot next state = start state
                               \cdot done = False
1133
```

1134	* Else if next position is goal_state:
1135	• Else it liext position is goal_state. • reward = $-1$
1136	$\cdot \text{ done} = \text{True}$
1137	• next state = goal_state
1138	* Else:
1139	$\cdot$ reward = $-1$
1140	$\cdot \text{ done} = \text{False}$
1141	• next state = computed position
1142	* Update current state to next state and return (next state, reward, done).
1143	• Usage Scenarios
1144 1145	
1146	- Scenario 1: Normal move into cliff
1147	* Initial state: 36 ([3,0])
1148	* Action: 1 (right)  Next state: 27 ([2,1]) is aliff. A recet to 26
1149	* Next state: 37 ([3,1]) is cliff $\rightarrow$ reset to 36 * Reward: $-100$
1150	* Reward. – 100  * done: False
1151	- Scenario 2: Move to goal
1152	
1153	* Current state: 46 ([3,10])  * Action: 1 (right)
1154	* Next state: 47 ([3,11]) goal
1155	* Reward: -1
1156	* done: True
1157	- Scenario 3: Move up from start
1158	* Current state: 36 ([3,0])
1159	* Action: 0 (up)
1160	* Next state: 24 ([2,0])
1161	* Reward: -1
1162	* done: False
1163	Assumptions & Risks
1164	
1165 1166	<ul> <li>No truncation or maximum episode length specified; risk of infinite episodes if goal not reached.</li> </ul>
1167	• set_state accepts only valid non-terminal, non-cliff states; no explicit error handling
1168	for invalid states.
1169	<ul> <li>No explicit error handling for invalid actions; assumed caller responsibility.</li> </ul>
1170 1171	<ul> <li>Reward and reset behavior on cliff stepping inferred from documentation; no episode termination on cliff.</li> </ul>
1172	Observation space includes all grid cells but player never occupies cliff or goal states.
1173	Potential future Gymnasium version changes may alter behavior.
1174 1175	Third-Party Library Usage
1176	Gymnasium API conventions inform environment interface: discrete action and ob-
1177	servation spaces, step and reset semantics.
1178	No new runtime dependencies introduced.
1179 1180	RNG seeding not specified; environment deterministic.
1181 1182	<ul> <li>No source code copied; specification based on official Gymnasium documentation and Sutton &amp; Barto RL book.</li> </ul>
1183	Numpy indexing conventions used for grid flattening.
1184	
1185	Amount in Police
1186	<workflow></workflow>
1187	1. <b>Deconstruct Specification:</b> Carefully review the <task describe="">,<research report=""></research></task>

188	
189	to fully understand the environment's specification, including state/action spaces, dynamics,
190	reward function, and termination conditions.
191	2. <b>Physics Engine Selection:</b> Evaluate if the task requires physics simulation. If so, choose
192	an appropriate physics engine for the specific task requirements.
193	3. <b>Model Design:</b> If using a physics engine, design the model structure and embed it as needed in the Python file.
194	4. <b>Plan Class Structure:</b> Outline the 'Environment' class, including its internal state
195	variables, helper methods, and the public interface ('init', 'reset', 'set_state', 'step').
196	5. <b>Implement Complete Code:</b> Write the full implementation of the 'Environment' class. 6.
197	Self-Correction Review: Meticulously check that the generated code fully complies with the
198	<task description="">, the <research report="">, and all <implementationrequirements>.</implementationrequirements></research></task>
199	7. <b>Finalize Output:</b> Present the complete, reviewed, and runnable single-file code in the
200	specified final format.
201	
202	<implementationrequirements></implementationrequirements>
203	1. Interface (single file):
204	
205	• Implement a complete, self-contained Python class Environment with:
206	init(self, seed: int   None = None)
207 208	<ul> <li>reset(self, seed: int   None = None) → ndarray (reinitialize the episode and return the initial observation in canonical shape)</li> </ul>
209	- set_state(self, state) (must accept ndarray or list/tuple in canonical shape)
210	<ul> <li>step(self, action) → tuple[ndarray, float, bool] (returns: observation, reward,</li> </ul>
211	done)
212	Requirements:
213	- Single-file constraint: all code, including any model definitions, must be
214	contained in one Python file.
215 216	<ul> <li>For physics-based environments, embed model definitions as string con- stants within the class.</li> </ul>
217	<ul> <li>Explicitly define state, action, and observation spaces (types, shapes, ranges,</li> </ul>
218	formats).
219	- Provide reproducibility (seeding) via the constructor and/or a seed(int)
220	method.
221	<ul> <li>Be robust to common representations:</li> </ul>
222	* set_state: accept list/tuple/ndarray of the same logical content.
223	* step: accept int / NumPy integer scalar / 0-D or 1-D len-1 ndarray (convert
224	to canonical form; raise clear TypeError/ValueError on invalid inputs).
225	No dependence on external RL frameworks; no Gym inheritance.
226	- No external file dependencies (model definitions must be embedded).
227	Maintain internal state consistency; allow reconstruction from observations  where applicable.
228	where applicable.  - Clean, readable code suitable for RL experimentation.
229	•
230	2. Determinism & validation:
231	<ul> <li>Provide reproducibility via seed (constructor and/or seed(int) method).</li> </ul>
232	Normalize inputs: accept equivalent representations (e.g., NumPy scalar/int/len-
233	1 array) and convert to a canonical form.
234 235	<ul> <li>Validate inputs; raise clear one-line errors (ValueError/TypeError) on invalid shapes or ranges.</li> </ul>
236	
237	3. Dynamics (MCTS/control oriented):
238	• For physics-based tasks, prefer suitable physics simulation methods with em-
239	bedded model definitions over custom physics implementations.
240	• Choose and document an integration scheme (e.g., implicit integrator, explicit
241	Euler) consistent with the research report.

	ullet Use a stable time step $dt$ ; clamp to safety bounds; keep all values finite (no
	NaN/Inf).  • Keep per-step computation efficient and allocation-light.
	4. Dependencies & style:
	•
	<ul> <li>No Gym inheritance or external RL frameworks unless explicitly allowed.</li> <li>Allowed: third-party libraries as needed (e.g., NumPy, physics engines, SciPy, Numba, JAX, PyTorch, etc.).</li> </ul>
	<ul> <li>For robotics/physics tasks, physics engines with embedded model definitions are recommended over custom implementations.</li> </ul>
	Clean, readable code suitable for RL experimentation.
	<ul> <li>All dependencies must be importable standard libraries or commonly available packages.</li> </ul>
	plementationRequirements>
	tput Format>
ent	al> <code_file_path> The entrypoint file path of the generated code. </code_file_path> rypoint_code> "'python # Your complete, runnable single-file implementation here. "'trypoint_code>
	ttput Format>
Jnit '	Tester
TA:	SK DESCRIPTION> Cliff walking involves crossing a gridworld from start to goal
	e avoiding falling off a cliff.
	escription
	SK DESCRIPTION> leArtifact path="environment.py"> {code}
	ecutionPolicy>
	Do not modify the student's source file.
	• Create exactly one pytest file at "tests/test_env.py" using file_tool("save").
	• Import the module from "environment.py" via importlib (spec_from_file_location + module_from_spec).
	• Run tests with code_tool("run", "pytest -q"); capture exit_code, duration, and std-out/stderr tail.
	ecutionPolicy> tPlan>
	• Sanity: class Environment can be imported and instantiated, e.g.,
	Environment (seed=0).
	• Contract:
	<ol> <li>set_state accepts list/tuple/ndarray of the same logical content (convert to canonical).</li> </ol>
	2. step (action) returns a 3-tuple: (observation, reward, done) with expected
	types/shapes.
	3. Determinism: with the same seed and same initial state, the first step with the
	same action yields identical outputs.
	4. Action space validation: actions within bounds are accepted, out-of-bounds actions are handled greenfully.
	tions are handled gracefully.  5. Observation space validation: observations match declared space bounds and
	shapes.
	6. State space consistency: internal state dimensions match expected environment
	specifications.
	specifications.

1290	
1297	
1298	<reportingguidelines></reportingguidelines>
1299	• Summarize pytest results in 2–4 sentences; mention the first failing nodeid/assert if
1300	any.
1301	• Provide a brief contract coverage assessment and the most probable root cause for
1302	failures.
1303	• If failing, add 1–3 concise actionable fixes (no long logs).
1304	
1305 1306	<b><outputformat></outputformat></b> Return exactly one <final> block containing a single JSON object that</final>
1307	matches PytestReport: {
1308	"success": truelfalse, "analysis": "<2–4 sentence summary/diagnosis>",
1309	"suggest_fix": " 1–3 bullets with minimal actionable changes>"
1310	No extra text outside <final>. No additional code fences.</final>
1311	<pre><final> { "success": false, "code_result": "", "analysis": "", "suggest_fix": "" } </final></pre>
1312	
1313	
1314	Simulation Tester
1315	Your task is to interact with the environment code and then analyze the feedback from the
1316	interaction and propose modifications
1317	<task description=""></task>
1318	Cliff walking involves crossing a gridworld from start to goal while avoiding falling off a cliff.
1319	## Description
1320	<task descrif="" tion=""> <codeartifact path="environment.py"></codeartifact></task>
1321 1322	{code}
1323	
1324	
1325	<executionpolicy></executionpolicy>
1326	- Use the play_env tool exactly once on "environment.py" - If the tool throws or cannot run, perform diagnosis from static review only; still produce output in the required format.
1327	<pre> </pre> <pre> <pre> <pre> <pre> <pre> <pre> <pre> </pre> <pre> <pre< th=""></pre<></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre></pre>
1328	<rubric></rubric>
1329	Success (boolean) must be decided from the available signals with graceful degradation:
1330	• Primary (step-level signals present):
1331	<ul> <li>success = true iff the run finished without exceptions AND there is NO misclas-</li> </ul>
1332	sified_transition with (valid == false OR state_matches == false).
1333 1334	- If only observation deltas are available, use obs_matches instead of
1335	state_matches.
1336	- When numeric deltas are provided, treat matches = true if max_abs_error $\leq 10^{-3}$ or rel error $\leq 10^{-3}$ .
1337	
1338	• Secondary (no per-step signals):
1339	- If success_rate exists: success = true iff no exceptions AND success_rate ≥
1340	0.95.
1341	<ul> <li>Else: success = true iff no exceptions AND no invariant/contract violations you can substantiate from code and logs.</li> </ul>
1342	• Reward/termination:
1343	
1344	<ul> <li>If reward_matches == true AND done_matches == true, explicitly state they match and DO NOT propose changes to reward or termination logic.</li> </ul>
1345	Action space consistency (discrete & continuous):
1346 1347	
1348	<ul> <li>If GT exposes Box (low, high, shape): align predicted bounds and expose them (e.g., env.action_space or a getter). Never place clipping inside</li> </ul>
1349	the integrator; clamp only at action ingestion or at observation.

	- If GT exposes Discrete (n): actions must be integer indices in $[0, n-1]$ ;
	expose n (e.g., gym.spaces.Discrete(n)); if indices map to continuous
	commands/torques, list the mapping table and align it with GT; never float-clip
	discrete actions.
	<ul> <li>If action-space info is missing, skip these checks (do not speculate).</li> </ul>
	• Internal vs observation:
	- If clipping or angle normalization is found inside the integrator step (e.g., in
	_rk4_step), this likely causes trajectory drift; propose moving them to the observation path (e.g., _get_observation) unless GT specifies otherwise.
	<ul> <li>If latent state is unavailable but observations exist, compare observations instead and state this explicitly.</li> </ul>
	• Integrator & timestep:
	- Mismatches in integrator method (e.g., RK4 vs Euler) or dt can cause state di-
	vergence even when reward/done match; acknowledge and, if state mismatches
	persist, propose aligning method/dt to GT.
	Batched/multiple transitions:
	- If multiple transitions are reported, aggregate sensibly (e.g., mean success_rate
	or fraction matched $\geq 0.95$ ) before deciding success.
	<procedure></procedure>
	1) Static review: scan for action bounds, clipping/normalization inside integrator, integrator/dt
	choice, and how observation is formed.
	2) Execute: call play_env once.
	3) Diagnose: reconcile play_env signals with code; if reward/done matched, explicitly say so.
	If state mismatched, point to ONE OR TWO most likely roots.
	4) Suggest: 1–3 smallest patches that directly address the identified root causes.
	<outputformat></outputformat>
	Return exactly one <final> block containing a single JSON object that matches PlayReport:</final>
	{ "success": truelfalse, "analysis": "<2–4 sentences summarizing what happened and why;
	mention matches/mismatches explicitly>", "suggest_fix": "- bullet 1\n- bullet 2\n- bullet 3
	(optional)" } No extra text outside <final>. No additional code fences.</final>
4	
C	G EXAMPLES OF DEEP RESEARCHER REPORT
ſ	Deep Researcher's Report
	Vancian & Duamanana
	<version &="" provenance=""></version>
	• Environment: Pusher-v4
	• MuJoCo version: $\geq 2.1.3$ and $< 3.0.0$
	• Source release date: 2024-01-01 (approximate for v4)
	• Accessed date: 2025-06-01 (Asia/Singapore)
	<pre><spec patch=""></spec></pre>
	• action.space = Box(low=-2.0, high=2.0, shape=(7,),
	dtype=np.float32)
	• observation.space: Box with low $= -\infty$ , high $= \infty$ , shape (23,), dtype
	nn float6/
	np.float64
	• reward.weights = {dist: $1.0$ , ctrl: $0.1$ , near: $0.5$ }
	<ul> <li>reward.weights = {dist: 1.0, ctrl: 0.1, near: 0.5}</li> <li>reward.formula: reward = reward_dist + 0.1 reward_ctrl +</li> </ul>
	• reward.weights = {dist: 1.0, ctrl: 0.1, near: 0.5}

```
1404
1405
                   • reward_ctrl = -\|\mathbf{a}\|_2^2 (negative squared Euclidean norm)
1406
                   • reward_dist = -\|\mathbf{o} - \mathbf{g}\|_2 (negative Euclidean distance)
1407
                   • reward_near = -\|\mathbf{f} - \mathbf{o}\|_2 (negative Euclidean distance)
1408
                   • episode.max_steps = 100
1409
1410
                   • episode.termination: never (terminated always False)
1411
                   • episode.truncation: True if step count > max steps or any observation value is
1412
                     non-finite
1413
                   • dt = 0.05\,\mathrm{s} (frame_skip= 5, frame_duration= 0.01)
1414
                   • initial_state.joint_positions = \mathbf{0}_7
1415
1416
                   • initial_state.joint_velocities: each joint \sim U(-0.005, 0.005)
1417
                   • initial_state.object_position: x \sim U(-0.3, 0), y \sim U(-0.2, 0.2); resample until
1418
                     ||(x,y)||_2 > 0.17
1419
                   • initial_state.object_velocity = 0<sub>3</sub>
1420
                   • initial_state.goal_position = [0.45, -0.05, -0.323]
                   • initial_state.goal_velocity = 0_3
1422
1423
                   • step_function:
1424
                        - clips action to action-space bounds
1425
                        - applies action torques to joints
1426

    advances MuJoCo simulation by dt

1427
                        - updates observation vector from simulation state
1428
                        - computes reward components and total reward
1429
                        - checks for finite observation values
1430
                        - increments step count
1431
1432
                                                 (observation, reward, terminated=False,
1433
                           truncated=done_flag, info) with reward components in info
1434
                   • seeding: via reset (seed=...) method only
1435
                   • info dict keys: "reward_dist", "reward_ctrl", "reward_near"
1436
                   • fingertip position: MuJoCo site "tips_arm"
1437
                   • object & goal positions: MuJoCo bodies "object" and "goal"
1438
1439
                   • internal state: corresponds to MuJoCo simulation state
1440
            <Theoretical Foundations>
1441
                   • MDP tuple: (S, A, P, R)
1442
                        - S: continuous state space \mathbb{R}^{23} (23-dimensional real vector)
1443
                        - A: continuous action space \mathbb{R}^7 with elementwise bounds [-2, 2] (torques)
1444
1445
                        - P: transition probability induced by MuJoCo physics with dt = 0.05 s
1446
                        - R: reward function as defined below
1447

    Reward function

1448
                        - Definition:
1449
                              r(s,a) = -\|P_{\text{object}} - P_{\text{goal}}\|_{2} - 0.1 \|a\|_{2}^{2} - 0.5 \|P_{\text{fingertip}} - P_{\text{object}}\|_{2}
1450
1451
                        - Where:
1452
                            * P_{\text{object}}, P_{\text{goal}}, P_{\text{fingertip}} \in \mathbb{R}^3 are positions in meters
1453
                            * \|\cdot\|_2 is the Euclidean norm
1454
                            * a \in \mathbb{R}^7 is the action torque vector in N \cdot m
1455
1456

    Episode ending

1457
```

1458	
1459	<ul> <li>Truncation after 100 steps or if any observation value is non-finite</li> </ul>
1460	<ul> <li>No termination condition (i.e., terminated=False always)</li> </ul>
1461	Fingertip position
1462	<ul> <li>Computed via forward kinematics internally by MuJoCo</li> </ul>
1463	
1464	• Symbol Table
1465	- a: Action torque vector; Unit $N \cdot m$ ; Range: each element in $[-2, 2]$
1466	- $P_{\text{fingertip}}$ : Fingertip 3D position; Unit m; Range: real values, unbounded
1467	<ul> <li>P<sub>object</sub>: Object 3D position; Unit m; Range: real values, unbounded</li> </ul>
1468	- $P_{\text{goal}}$ : Goal 3D position; Unit m; Value: fixed at $[0.45, -0.05, -0.323]$
1469	- r: Reward; Unit: unitless; Range: real values (sum of negative components)
1470	- dt: Simulation timestep; Unit s; Value: 0.05
1471	- step_count: Current timestep count; Unit: integer; Range: 0 to 100
1472	Final Specification
1473 1474	- Environment Name: Pusher-v4
1475	- Class Name: Environment
1476	- Class Name. Environment - Public API
1477	
1478	*init(self)
1479	<ul> <li>Initializes the MuJoCo simulation environment with Pusher-v4 model</li> <li>Sets initial internal variables including step count</li> </ul>
1480	Defines action and observation spaces as specified
1481	Sets fixed goal position
1482	* set_state(self, state: np.ndarray) -> None
1483	• Input: state vector of shape (23,) matching observation space
1484	• Sets the internal MuJoCo simulation state to correspond to the given ob-
1485	servation
1486	· Resets step count to zero; assumes state is valid and episode not done
1487	* step(self, action: np.ndarray) ->
1488	<pre>Tuple[np.ndarray, float, bool, dict]</pre>
1489	· Input: action vector of shape $(7,)$ clipped to $[-2, 2]$
1490	$\cdot$ Applies action torques to the simulation joints and advances by $dt =$
1491	$0.05\mathrm{s}$
1492	· Updates observation vector; computes reward components and total re-
1493	ward
1494	• Truncates if any observation value is non-finite or if step count $\geq 100$
1495	Returns: observation (23,) float64, reward float, terminated False, truncated bool, info with keys "reward_dist",
1496	"reward_ctrl", "reward_near"
1497	- Core Logic
1498	* Simulates a 7-DOF robotic arm pushing an object to a goal
1499	* State includes joint angles/velocities, fingertip, object, and goal positions
1500	* Actions are joint torques clipped to [-2, 2]
1501	* Physics advances by 0.05 s per step
1502	* Reward encourages moving object to goal, penalizes large actions, and en-
1503	courages fingertip proximity to object
1504	* Episode truncates after 100 steps or on non-finite state
1505	* Fingertip and object positions come from MuJoCo sites/bodies
1506	* Seeding via reset (seed=) (not shown here)
1507 1508	- Usage Scenarios & Examples
1506	* Scenario 1: Initial step with zero action
1510	· Observation: initial state vector with zero joint positions, small velocity

noise, specified object and goal

```
1512
1513
                           · Action: np.zeros(7, dtype=np.float32)
1514
                           · Reward: negative sum of distances (control penalty = 0)
1515
                           · Terminated: False; Truncated: False
1516
                           · Info: {"reward_dist": value, "reward_ctrl":
                                                                                            0.0,
1517
                              "reward_near": value}
1518
                         * Scenario 2: Non-zero action moves fingertip closer
1519
                           · Observation: updated state after physics step
1520
                                        np.array([1.0, -0.5, 0.0, 0.2, 0.0, -1.0,
                             0.5], dtype=np.float32)
1521
                           · Reward: computed per formula
1522
                           · Terminated: False; Truncated: False (unless step count reached 100)
                           · Info: {"reward_dist": value, "reward_ctrl": negative
1524
                             squared norm of action, "reward near": value}
1525
1526
           <Assumptions & Risks>
1527

    Assumes MuJoCo physics engine is used as backend with version constraints.

    Assumes internal MuJoCo state can be set and queried to implement set_state

                   and step.
1530

    Fingertip position is assumed to be accessible via MuJoCo site "tips_arm".

1531
                 • Reward weights are fixed as per v4; no user-configurable weights.
1532
1533
                 • Episode termination is only by truncation; no natural termination.
1534
                 • Seeding is only via reset (seed=...); no separate seed method.
1535

    Risks include dependency on MuJoCo simulation correctness and numerical stability.

1536
1537

    No explicit handling of invalid actions beyond clipping.

1538

    No rendering or visualization included.

1539
           <Library Usage>
1540

    Uses Python standard library and NumPy only.

1541

    NumPy used for array operations, clipping (np.clip), norm calculations

                   (np.linalg.norm), and finite checks (np.isfinite).
1543
1544

    No external dependencies beyond MuJoCo simulation backend (assumed).

1545
                 • All units are SI: torque in N·m, angles in radians, positions in meters, time in seconds.
1546

    No source code copied; specification based on official documentation and API stan-

1547
                   dards.
1548
```

# </final>

1549 1550 1551

1552 1553

1554 1555

1556

1557

1558

1560

1561

1562

# H CASE STUDY

#### H.1 TEXT2WORLD

**Task Describe.**Prepare and deliver sandwiches to a set of children: assemble sandwiches in the kitchen (optionally gluten-free), place them on a tray, move the tray to the child's location, and serve subject to allergy constraints. Gluten-allergic children must receive a gluten-free sandwich; non-allergic children may receive any sandwich. Serving requires the sandwich on the tray and the tray at the child's location; making a sandwich switches it from "notexist" to "prepared." The goal is that all children become served.

#### Prev. SOTA.

```
1566
           bread-portion
1567
           content-portion
1568
           sandwich
1569
           tray
           place
1570
1571
1572
         (:constants
           kitchen - place
1573
1574
1575
         (:predicates
1576
           (at_kitchen_bread ?b - bread-portion)
1577
           (at_kitchen_content ?c - content-portion)
           (at_kitchen_sandwich ?s - sandwich)
1578
           (no_gluten_bread ?b - bread-portion)
1579
           (no_gluten_content ?c - content-portion)
1580
           (ontray ?s - sandwich ?t - tray)
1581
           (no_gluten_sandwich ?s - sandwich)
           (allergic_gluten ?c - child)
1582
           (not_allergic_gluten ?c - child)
1583
           (served ?c - child)
1584
           (waiting ?c - child ?p - place)
1585
           (at ?t - tray ?p - place)
           (notexist ?s - sandwich)
1586
1587
1588
         ;; Action: make_sandwich_no_gluten
1589
         (:action make_sandwich_no_gluten
1590
           :parameters (?s - sandwich ?b - bread-portion ?c - content-portion)
           :precondition (and
             (at_kitchen_bread ?b)
1592
             (at_kitchen_content ?c)
1593
             (no_gluten_bread ?b)
1594
             (no_gluten_content ?c)
1595
             (notexist ?s)
1596
           :effect (and
1597
             (not (at_kitchen_bread ?b))
1598
             (not (at_kitchen_content ?c))
             (notexist ?s)
1599
             (not (notexist ?s))
1600
             (at_kitchen_sandwich ?s)
1601
             (no_gluten_sandwich ?s)
1602
1603
         )
1604
         ;; Action: make_sandwich
1605
         (:action make_sandwich
1606
           :parameters (?s - sandwich ?b - bread-portion ?c - content-portion)
1607
           :precondition (and
             (at_kitchen_bread ?b)
             (at_kitchen_content ?c)
1609
             (notexist ?s)
1610
1611
           :effect (and
             (not (at_kitchen_bread ?b))
1612
             (not (at_kitchen_content ?c))
1613
             (notexist ?s)
1614
             (not (notexist ?s))
1615
             (at_kitchen_sandwich ?s)
1616
             ;; no effect on gluten status, sandwich may or may not be

→ gluten-free

1617
1618
         )
1619
```

```
1620
         ;; Action: put_on_tray
1621
         (:action put_on_tray
1622
           :parameters (?s - sandwich ?t - tray)
1623
           :precondition (and
             (at_kitchen_sandwich ?s)
1624
             (at ?t kitchen)
1625
1626
           :effect (and
             (not (at_kitchen_sandwich ?s))
1627
             (ontray ?s ?t)
1628
1629
1630
1631
         ;; Action: serve_sandwich_no_gluten
         (:action serve_sandwich_no_gluten
1632
           :parameters (?s - sandwich ?c - child ?t - tray ?p - place)
1633
           :precondition (and
1634
             (ontray ?s ?t)
1635
             (no_gluten_sandwich ?s)
             (allergic_gluten ?c)
1636
             (waiting ?c ?p)
1637
             (at ?t ?p)
1638
1639
           :effect (and
1640
             (not (ontray ?s ?t))
             (served ?c)
1641
             (not (waiting ?c ?p))
1642
1643
         )
1644
         ;; Action: serve_sandwich
         (:action serve_sandwich
1646
           :parameters (?s - sandwich ?c - child ?t - tray ?p - place)
1647
           :precondition (and
1648
             (ontray ?s ?t)
1649
             (not (no_gluten_sandwich ?s))
             (not_allergic_gluten ?c)
1650
             (waiting ?c ?p)
1651
             (at ?t ?p)
1652
1653
           :effect (and
             (not (ontray ?s ?t))
1654
             (served ?c)
1655
             (not (waiting ?c ?p))
1656
1657
         )
1658
         ;; Action: move_tray
1659
         (:action move_tray
1660
           :parameters (?t - tray ?p1 - place ?p2 - place)
1661
           :precondition (at ?t ?p1)
           :effect (and
1662
             (not (at ?t ?p1))
1663
             (at ?t ?p2)
1664
1665
         )
1666
1667
       Agent2World.
1668
1669
       (define (domain child-snack)
1670
         (:requirements :typing :negative-preconditions :strips)
1671
         (:types
           child
1672
           bread-portion
1673
           content-portion
```

```
1674
           sandwich
1675
           tray
1676
           place
1677
         (:constants
1678
           kitchen - place
1679
1680
         (:predicates
           (at_kitchen_bread ?b - bread-portion)
1681
           (at_kitchen_content ?c - content-portion)
1682
           (at_kitchen_sandwich ?s - sandwich)
1683
           (no_gluten_bread ?b - bread-portion)
1684
           (no_gluten_content ?c - content-portion)
1685
           (ontray ?s - sandwich ?t - tray)
           (no_gluten_sandwich ?s - sandwich)
1686
           (allergic_gluten ?c - child)
1687
           (not_allergic_gluten ?c - child)
1688
           (served ?c - child)
1689
           (waiting ?c - child ?p - place)
           (at ?t - tray ?p - place)
1690
           (notexist ?s - sandwich)
1691
1692
1693
         (:action make_sandwich_no_gluten
1694
           :parameters (?s - sandwich ?b - bread-portion ?c - content-portion)
           :precondition (and
1695
             (at_kitchen_bread ?b)
1696
             (at_kitchen_content ?c)
1697
             (no_gluten_bread ?b)
1698
             (no_gluten_content ?c)
             (notexist ?s)
1700
           :effect (and
1701
             (not (notexist ?s))
1702
             (at_kitchen_sandwich ?s)
1703
             (no_gluten_sandwich ?s)
             (not (at_kitchen_bread ?b))
1704
             (not (at_kitchen_content ?c))
1705
1706
         )
1707
         (:action make_sandwich
1708
           :parameters (?s - sandwich ?b - bread-portion ?c - content-portion)
1709
           :precondition (and
1710
             (at_kitchen_bread ?b)
1711
             (at_kitchen_content ?c)
             (notexist ?s)
1712
1713
           :effect (and
1714
             (not (notexist ?s))
1715
             (at_kitchen_sandwich ?s)
             (not (at_kitchen_bread ?b))
1716
             (not (at_kitchen_content ?c))
1717
1718
         )
1719
         (:action put_on_tray
1720
           :parameters (?s - sandwich ?t - tray)
1721
           :precondition (and
1722
             (at_kitchen_sandwich ?s)
1723
             (at ?t kitchen)
1724
           :effect (and
1725
             (ontray ?s ?t)
1726
             (not (at_kitchen_sandwich ?s))
1727
```

```
1728
1729
1730
         (:action serve_sandwich_no_gluten
           :parameters (?s - sandwich ?c - child ?t - tray ?p - place)
1731
           :precondition (and
1732
              (ontray ?s ?t)
1733
              (no_gluten_sandwich ?s)
1734
              (allergic_gluten ?c)
              (waiting ?c ?p)
1735
              (at ?t ?p)
1736
1737
           :effect (and
1738
              (served ?c)
1739
              (not (ontray ?s ?t))
1740
         )
1741
1742
         (:action serve_sandwich
1743
           :parameters (?s - sandwich ?c - child ?t - tray ?p - place)
            :precondition (and
1744
              (ontray ?s ?t)
1745
              (waiting ?c ?p)
1746
              (not_allergic_gluten ?c)
1747
              (at ?t ?p)
           )
1748
           :effect (and
1749
              (served ?c)
1750
              (not (ontray ?s ?t))
1751
         )
1752
1753
         (:action move_tray
1754
           :parameters (?t - tray ?p1 - place ?p2 - place)
1755
           :precondition (at ?t ?p1)
1756
           :effect (and
              (at ?t ?p2)
1757
              (not (at ?t ?p1))
1758
1759
         )
1760
       )
1761
```

**Analysis.** Compared to the baseline domain, our *Child-Snack* formulation introduces three taskaligned modifications that improve state consistency, compositionality, and plan feasibility. (i) Creation-valid effects. During sandwich construction we flip the existence status from "nonexistent" to "prepared," and record gluten-free status when applicable, thereby avoiding contradictory postconditions at creation time; this yields deterministic successor states and reduces backtracking caused by ill-defined truth values. (ii) Serve-focused effects. During serving we only transfer the item off the tray and mark the child as served, leaving the waiting label untouched; this separation of concerns prevents nonessential side-effects, preserves modular composability with downstream routines (e.g., queueing or follow-up allocation), and promotes goal-monotonic progress on the served objective. (iii) **Permissive-serving preconditions.** For non-allergic children we do not exclude gluten-free items, weakening preconditions to accept any admissible sandwich; this enlarges the feasible search space and prevents avoidable dead-ends when only gluten-free inventory remains, while safety for allergic children is still enforced via a dedicated gluten-free serving action. Collectively, these choices align with the ground-truth specification, produce cleaner state transitions, and yield empirically favorable search dynamics—smaller inconsistent-state frontiers and fewer spurious deletions—resulting in a more robust make—put-on-tray—move-tray—serve pipeline for the objective of "serving each child an acceptable sandwich."

# H.2 CWMB

1762

1763

1764

1765

1766

1767

1768

1769

1770

1771

1772

1773

1774

1775

1776

177717781779

1780

1781

**Task Describe.** Control a 3D Ant (one free-body torso with four 2-DoF legs; nine bodies, eight hinge joints) to move forward along the +x axis by applying torques to the eight joints at each

step. The action space is  $Box([-1,1]^8)$  (joint torques). Observations list positions then velocities (27-D by default; 29-D if current x,y are included), and optionally +84 contact-force terms when use\_contact\_forces=True or version < v4. The reward is  $r=r_{\rm healthy}+r_{\rm forward}-{\rm ctrl}_{\rm cost}$  (and -contact\_cost if contact forces are used), where  $r_{\rm forward}\approx \Delta x/\Delta t$  is positive for motion in +x and  $\Delta t={\rm frame\_skip}\times 0.01=0.05$  by default. Episodes start from an upright, slightly noisy state, truncate at 1000 steps, and (by default) terminate early if the agent becomes unhealthy (non-finite state or torso  $z\notin [0.2,1.0]$ ).

# Prev SOTA vs Agent2World.

1782

1783

1784

1785

1786

1787

```
1790
1791
      import numpy as np
      from typing import Optional, Tuple
1792
      # [D-0] Toggle: False = Ours (physics-like integration), True = Baseline
1793
       1794
      BASELINE = False
1795
1796
      class Environment:
          def ___init___(
1797
              self,
1798
              # ---- Common params ----
1799
              frame_skip: int = 5,
1800
              max_episode_steps: int = 1000,
              healthy_reward: float = 1.0,
1801
              terminate_when_unhealthy: bool = True,
1802
              healthy_z_range: Tuple[float, float] = (0.2, 1.0),
1803
              ctrl_cost_weight: float = 0.5,
1804
              # ---- Ours-only params ---
              contact_cost_weight: float = 5e-4,
              contact_force_range: Tuple[float, float] = (-1.0, 1.0),
1806
              exclude_current_positions_from_observation: bool = True,
1807
              use_contact_forces: bool = False,
1808
              reset_noise_scale_ours: float = 0.1,
1809
              forward_reward_weight: float = 1.0,
               # ---- Baseline-only params ----
1810
              reset_noise_scale_baseline: float = 0.01,
1811
              seed: Optional[int] = None,
1812
          ):
1813
              # Core configuration
              self.frame_skip = frame_skip
1814
              self.dt = 0.01 * frame_skip
1815
              self.max_episode_steps = max_episode_steps
1816
              self.healthy_reward = healthy_reward
1817
              self.terminate_when_unhealthy = terminate_when_unhealthy
1818
              self.healthy_z_range = healthy_z_range
              self.ctrl_cost_weight = ctrl_cost_weight
1819
              self.forward_reward_weight = forward_reward_weight # [D-5] Only
1820
               \hookrightarrow used by Ours.
1821
              self.np_random = np.random.RandomState(seed)
1822
               # [D-1] Observation schema differs:
1823
                   Baseline: fixed 27 = 13 positions (z, quat4, joint8) + 14
                  velocities (flat vector).
1825
                   Ours:
                           positions (15) + velocities (14) (+ optional
1826
                  torso_xy, contact forces),
                             with option to exclude torso x, y from observation.
1827
              self.exclude_current_positions_from_observation = (
1828
                   False if BASELINE else
1829
                   → bool(exclude_current_positions_from_observation)
1830
1831
               # [D-2] Contacts: Baseline has no contact forces/cost; Ours can
1832
               \rightarrow include 84-dim contact forces + cost.
1833
               self.use_contact_forces = False if BASELINE else
1834
               → bool(use_contact_forces)
               self.contact_cost_weight = 0.0 if BASELINE else
               → float(contact_cost_weight)
```

```
1836
               self.contact_force_range = contact_force_range
1837
1838
               self.step\_count = 0
1839
               if BASELINE:
1840
                   # Baseline state: flat (27,) observation vector
1841
                   self.obs\_shape = (27,)
1842
                   self.state = np.zeros(self.obs_shape, dtype=np.float64)
                   self.x_position = 0.0 # [D-4] Progress tracked separately
1843

→ (not in observation)

                   self.y_position = 0.0
1845
                   self.reset_noise_scale = float(reset_noise_scale_baseline)
1846

    □ D - 4 ]

1847
                   self.contact_forces = None
                   self.observation_dim = 27 # [D-1]
               else:
1849
                   # Ours state: split positions(15) / velocities(14)
1850
                   self.pos_dim = 15
                                       # torso_pos(3), torso_quat(4),
1851

    joint_angles(8)

                   self.vel_dim = 14 # torso_lin_vel(3), torso_ang_vel(3),
1852
                       joint_vel(8)
1853
                   self.positions = np.zeros(self.pos_dim, dtype=np.float64)
1854
                   self.velocities = np.zeros(self.vel_dim, dtype=np.float64)
1855
                   self.reset_noise_scale = float(reset_noise_scale_ours)
1856
                   \hookrightarrow [D-4]
                   self.last_x_position = 0.0 \# [D-4] (used when restoring
1857

    state)

1858
                   self.contact_forces = (
                       np.zeros(84, dtype=np.float64) if self.use_contact_forces
1860

→ else None

                   # Compute observation length for Ours
1862
                   base_pos_len = self.pos_dim
1863
                   if self.exclude_current_positions_from_observation:
1864
                       base_pos_len -= 2 # drop torso x,y
                   self.obs_pos_len = base_pos_len
1865
                   self.obs_vel_len = self.vel_dim
1866
                   self.obs_contact_len = 84 if self.use_contact_forces else 0
1867
                   self.obs_torso_xy_len = 0 if
1868
                   → self.exclude_current_positions_from_observation else 2
1869
                   self.observation_dim = (
                       self.obs_pos_len + self.obs_vel_len +
1870
                        ⇔ self.obs_torso_xy_len + self.obs_contact_len
1871
1872
1873
               self.reset (seed)
1874
           # ------ Helpers (shared) ------
1875
           def _is_healthy(self, z_value: Optional[float] = None) -> bool:
1876
               # [D-3] Health z-source differs:
1877
                   Baseline: use state[0] (z) from flat observation.
1878
                   Ours:
                             use positions[2] (torso z) from split state.
               if BASELINE:
1879
                   st = self.state
1880
                   if not np.all(np.isfinite(st)):
1881
                       return False
1882
                   z = st[0] if z_value is None else z_value
                   return self.healthy_z_range[0] <= z <=
1883

    self.healthy_z_range[1]

1884
               else:
1885
                   if not np.all(np.isfinite(self.positions)) or not
1886
                   → np.all(np.isfinite(self.velocities)):
1887
                       return False
                   z = self.positions[2]
1888
                   return self.healthy_z_range[0] <= z <=</pre>
1889

    self.healthy_z_range[1]
```

```
1890
1891
           def _get_observation(self) -> np.ndarray:
1892
               # [D-1] Observation layout differs (see __init__ comment).
1893
              if BASELINE:
                  return self.state.copy()
1894
               if self.exclude_current_positions_from_observation:
1895
                  pos_obs = self.positions[2:].copy() # exclude torso x,y
1896
                   torso_xy = np.array([], dtype=np.float64)
1897
               else:
                  pos_obs = self.positions.copy()
1898
                   torso_xy = self.positions[0:2].copy()
1899
              vel_obs = self.velocities.copy()
1900
              obs = np.concatenate((pos_obs, vel_obs))
1901
               if not self.exclude_current_positions_from_observation:
                  obs = np.concatenate((obs, torso_xy))
1902
               if self.use_contact_forces:
1903
                  obs = np.concatenate((obs, self.contact_forces))
1904
              return obs
1905
           @staticmethod
1906
          def _quat_multiply(q1: np.ndarray, q2: np.ndarray) -> np.ndarray:
1907
              w1, x1, y1, z1 = q1
1908
               w2, x2, y2, z2 = q2
1909
              return np.array(
1910
                   Γ
                       w1 * w2 - x1 * x2 - y1 * y2 - z1 * z2,
1911
                       w1 * x2 + x1 * w2 + y1 * z2 - z1 * y2,
1912
                       w1 * y2 - x1 * z2 + y1 * w2 + z1 * x2,
1913
                       w1 * z2 + x1 * y2 - y1 * x2 + z1 * w2,
1914
                   dtype=np.float64,
1915
1916
1917
           # ----- Public API -----
1918
          def set_state(self, state: np.ndarray) -> None:
1919
              state = np.asarray(state, dtype=np.float64)
               if BASELINE:
1920
                   # [D-8] Baseline expects a flat 27-dim observation
1921
                   \hookrightarrow (pos13+vel14).
1922
                   if state.shape != (27,):
                       raise ValueError(f"set_state input must have shape (27,),
1923

    got {state.shape}")

1924
                   if not np.all(np.isfinite(state)):
1925
                       raise ValueError("set_state input contains non-finite
1926

    values")

1927
                   self.state = state.copy()
                   self.step\_count = 0
1928
                   self.x.position = 0.0 # [D-4] progress variable is external
1929

→ to obs

1930
                   self.y_position = 0.0
1931
               else:
1932
                   # [D-8] Ours expects the current obs layout length and

→ reconstructs split state.

1933
                   expected_len = self.observation_dim
1934
                   if state.ndim != 1 or state.shape[0] != expected_len:
1935
                       raise ValueError (f"State must be 1D array of length
1936
                       pos_len = self.obs_pos_len
1937
                   vel_len = self.obs_vel_len
1938
                   pos_part = state[:pos_len]
1939
                   vel_part = state[pos_len : pos_len + vel_len]
1940
                   if self.exclude_current_positions_from_observation:
1941
                       full_positions = np.zeros(self.pos_dim, dtype=np.float64)
                       full_positions[2:] = pos_part
1942
                       full_positions[0] = 0.0
1943
                       full_positions[1] = 0.0
```

```
1944
                   else:
1945
                       full_positions = pos_part.copy()
1946
                   self.positions = full_positions
1947
                   self.velocities = vel_part.copy()
                   self.step\_count = 0
1948
                   self.last_x_position = self.positions[0] # [D-4]
1949
1950
          def reset(self, seed: Optional[int] = None) -> np.ndarray:
               if seed is not None:
1951
                   self.np_random.seed(seed)
1952
               self.step\_count = 0
1953
1954
               if BASELINE:
1955
                   \# [D-2][D-4] Baseline init: 13 pos (z, quat, joints) + 14

→ vel; small noise.

1956
                   pos = np.zeros(13, dtype=np.float64)
1957
                   pos[0] = 0.75 # z
1958
                   pos[1:5] = np.array([1.0, 0.0, 0.0, 0.0]) # quaternion
1959
                   \hookrightarrow (w,x,y,z)
                   noise_pos = self.np_random.uniform(-self.reset_noise_scale,
1960

    self.reset_noise_scale, size=13)

1961
                   pos = pos + noise_pos
1962
                   vel = self.np_random.normal(0, self.reset_noise_scale,
1963
                   \hookrightarrow size=14)
1964
                   self.state = np.concatenate([pos, vel])
                   self.x_position = 0.0 # [D-4] progress variable
1965
                   self.y_position = 0.0
1966
               else:
1967
                   # Ours init: positions(15) / velocities(14) with larger noise
1968
                   \hookrightarrow and full torso pose.
                   base_positions = np.zeros(15, dtype=np.float64)
                   base\_positions[2] = 0.75 \# torso z
1970
                   base_positions[3] = 1.0
                                             # quat.w
1971
                   noise_pos = self.np_random.uniform(-self.reset_noise_scale,
1972

    self.reset_noise_scale, size=15)
1973
                   self.positions = base_positions + noise_pos
                   self.velocities = self.np_random.normal(loc=0.0,
1974
                   1975
                   self.last_x_position = self.positions[0] # [D-4]
1976
                   if self.use_contact_forces and self.contact_forces is not
1977
                      None:
                       self.contact_forces[:] = 0.0
1978
               return self._get_observation()
1979
1980
          def step(self, action: np.ndarray):
1981
               action = np.asarray(action, dtype=np.float64)
1982
               if action.shape != (8,):
                   raise ValueError(f"Action must be of shape (8,), got
1983
                      {action.shape}")
1984
1985
               # [D-6] Action bound handling differs:
                   Baseline: out-of-bounds raises; Ours: clip to [-1, 1].
               if BASELINE:
1987
                   if np.any(action < -1.0) or np.any(action > 1.0):
1988
                       raise ValueError("Action values must be in [-1, 1] for
1989
                        → Baseline")
1990
               else:
                   action = np.clip(action, -1.0, 1.0)
1991
1992
               # ---- Inline divergence (single function, two branches) ----
1993
               if BASELINE:
1994
                   # Forward progress proxy from hip joints
1995
                   prev_x_pos = self.x_position # [D-4] external tracker (not
                   \hookrightarrow in obs)
1996
                   forward_force = float(np.sum(action[[0, 2, 4, 6]]))
1997
```

```
1998
                    self.x_position += forward_force * self.dt * 0.1 # arbitrary
1999

→ scale

2000
2001
                    # Stochastic updates (no true physics)
                    z = float(np.clip(self.state[0] +
2002
                    \rightarrow self.np_random.uniform(-0.01, 0.01), 0.0, 2.0))
2003
                    # [D-7] Orientation update: Baseline = add noise to
2004
                       quaternion then renormalize.
2005
                    orientation = self.state[1:5] + self.np_random.uniform(-0.01,
                    \leftrightarrow 0.01, size=4)
2006
                   norm = float(np.linalg.norm(orientation))
2007
                    orientation = orientation / norm if norm > 0 else
2008
                    \rightarrow np.array([1.0, 0.0, 0.0, 0.0], dtype=np.float64)
2009
                    # Joint angles: integrate action + small noise; wrap to [-pi,
2010
                    → pi]
                    joint_angles = self.state[5:13] + action * self.dt +
2011
                    \rightarrow self.np_random.uniform(-0.005, 0.005, size=8)
2012
                    joint_angles = (joint_angles + np.pi) % (2 * np.pi) - np.pi
2013
                    # Velocities: torso (noise) + joints (action + noise)
2014
                   torso_vel = self.np_random.normal(0, 0.01, size=6)
                    joint_vel = action + self.np_random.normal(0, 0.01, size=8)
2015
                    velocities = np.concatenate([torso_vel, joint_vel])
2016
                    # Compose new flat state
2017
                   pos = np.empty(13, dtype=np.float64)
2018
                   pos[0] = z
                   pos[1:5] = orientation
2019
                   pos[5:13] = joint\_angles
2020
                    self.state = np.concatenate([pos, velocities])
2021
2022
                   healthy = self.\_is\_healthy(z\_value=z) # [D-3]
                    forward_delta = self.x_position - prev_x_pos
2023
                    contact\_cost = 0.0 \# [D-2] no contacts in Baseline
2024
                   weight = 1.0
                                         # [D-5] forward reward weight fixed to
2025
                    \hookrightarrow
                       1.0
2026
               else:
2027
                    # Physics-like integration
                    old_x = float(self.positions[0]) # [D-4] directly from torso
2028
                    \hookrightarrow x in positions
2029
                    # Joint dynamics: dv = (u - damp*v)*dt; dq = v*dt
2030
                    joint_damping = 0.1
2031
                    joint_vel_prev = self.velocities[6:]
2032
                    joint_acc = action - joint_damping * joint_vel_prev # [D-1]
                    \hookrightarrow acts on split state
2033
                    joint_vel_new = joint_vel_prev + self.dt * joint_acc
2034
                    joint_ang_new = self.positions[7:] + self.dt * joint_vel_new
2035
                    # Torso linear & angular velocity damping
2036
                   lin_damping = 0.1
                    ang_damping = 0.1
2037
                   torso_lin_vel_new = self.velocities[0:3] \star (1 - lin_damping \star
2038

    self.dt)

2039
                   torso_ang_vel_new = self.velocities[3:6] * (1 - ang_damping *
2040

    self.dt)

                    # Integrate torso position
2041
                   torso_pos_new = self.positions[0:3] + self.dt *
2042

→ torso_lin_vel_new

2043
                    # [D-7] Orientation update: Ours = quaternion integration
2044

→ from angular velocity.

2045
                    q = self.positions[3:7]
                   omega = torso_ang_vel_new
2046
                    omega_quat = np.array([0.0, omega[0], omega[1], omega[2]),
2047

    dtype=np.float64)

2048
                    q_dot = 0.5 * self._quat_multiply(omega_quat, q)
2049
                    q_new = q + self.dt * q_dot
                   norm = float(np.linalg.norm(q_new))
2050
                    q_new = q_new / norm if norm > 0 else np.array([1.0, 0.0,
2051
                    \rightarrow 0.0, 0.0], dtype=np.float64)
```

2089

2090

2091

2092

2093

2094

2095

2096

2097

2098

2099

2100 2101

2102

2103

2104

2105

```
2052
                   # Write back split state
2053
                  self.positions[0:3] = torso_pos_new
2054
                  self.positions[3:7] = q_new
                  self.positions[7:] = joint_ang_new
2055
                  self.velocities[0:3] = torso_lin_vel_new
2056
                  self.velocities[3:6] = torso_ang_vel_new
2057
                  self.velocities[6:] = joint_vel_new
2058
                   # Contacts (optional)
                  if self.use_contact_forces and self.contact_forces is not
2059
                      None:
2060
                       self.contact_forces.fill(0.0) # [D-2]
2061
                  new_x = float(self.positions[0])
2062
                  healthy = self._is_healthy()
                                                # [D-3]
2063
                  forward_delta = new_x - old_x
                   # Contact cost (if enabled)
2064
                  if self.use_contact_forces and self.contact_forces is not
2065
2066
                       clipped = np.clip(self.contact_forces,
2067

→ self.contact_force_range[0],
                          self.contact_force_range[1])
2068
                       contact_cost = self.contact_cost_weight *
2069

    float(np.sum(np.square(clipped)))  # [D-2]

2070
                  else:
2071
                      contact_cost = 0.0
2072
                  weight = self.forward_reward_weight # [D-5]
2073
               # ---- Shared reward & termination (single exit) ---
2074
              forward_reward = weight * (forward_delta / self.dt) # [D-5]
2075
              ctrl_cost = self.ctrl_cost_weight *
2076

    float (np.sum (np.square (action)))

              reward = (self.healthy_reward if healthy else 0.0) +
2077
                 forward_reward - ctrl_cost - contact_cost
2078
2079
              self.step_count += 1
2080
              done = (self.terminate_when_unhealthy and not healthy) or
2081
               return self._get_observation(), reward, done
2082
```

**Analysis.** On the **Ant-v4** forward-locomotion task, AGENT2WORLD surpasses the *Baseline* with higher success, smoother gait, and lower energy per meter under identical horizons and z-health checks. (i) **State & sensing.** The *Baseline* exposes a flat 27-D observation, while we adopt a taskaligned layout that separates positions/velocities and can hide global (x, y) by default ([D-1]). We additionally support contact forces for foot–ground cues ([D-2]). Health uses torso z from split state rather than the flat vector slot ([D-3]). State restoration matches each layout: the Baseline ingests a 27-D vector, whereas ours reconstructs split buffers from the current observation setting ([D-8]). (ii) **Dynamics & orientation.** The *Baseline* updates orientation by quaternion noise plus renormalization, and treats actions as noisy joint velocities; we integrate damped joint accelerations and update attitude via  $\dot{\mathbf{q}} = \frac{1}{2} \omega_{q} \otimes \mathbf{q}$  with renormalization ([D-7]). This physically consistent pipeline—enabled by the split state design ([D-1])—low-passes high-frequency actuation, reduces roll/pitch jitter, and yields more phase-coordinated gaits. (iii) Control semantics & reward. The Baseline hard-errors on out-of-range actions and uses a fixed forward-reward weight; forward progress is tracked by an external x variable and reset noise is smaller. Ours clips actions to [-1,1] ([D-6]), uses a tunable forward-reward weight ([D-5]), measures progress directly from torso x in the state and employs a different reset scale ([D-4]); an optional contact-cost term can be included when contact signals are enabled ([D-2]). Together these choices stabilize training signals and improve sample efficiency.

Summary of diffs. [D-1] Observation schema: Baseline uses a flat 27-D vector; Ours uses split positions+velocities with optional hidden (x, y) and optional contact forces; [D-2] Contacts: Baseline has no contact forces/cost; Ours optionally exposes 84-D contact forces and a contact-cost term; [D-3] Health source: Baseline takes z from the flat vector slot; Ours uses torso z from split positions; [D-4] Progress & reset: Baseline tracks forward x as an external variable and uses smaller reset noise; Ours reads torso x from state and uses a different reset scale; [D-5] Forward-reward

weight: Baseline fixed to 1.0; Ours is tunable; [D-6] Action bounds: Baseline errors on out-of-range actions; Ours clips to [-1,1]; [D-7] Orientation update: Baseline adds noise then renormalizes quaternion; Ours integrates  $\dot{\mathbf{q}} = \frac{1}{2} \, \omega_q \otimes \mathbf{q}$  then renormalizes; [D-8] State setting: Baseline ingests a flat 27-D state; Ours reconstructs split buffers from the current observation layout.

#### H.3 BYTESIZED32

2106

2107

2108

2109

2110 2111

2112

2113

2114

2115

2116

2117

2118

2119

2120

2121

2122

**Task Description.** We build a lightweight, text-interactive micro-simulation of pea growth in a small garden. The world contains a Pea, a FlowerPot, a Jug, and a Sink; water is represented as scalar levels in the Jug and FlowerPot and as an internal level in the Pea. The agent can look/examine, take/put objects, switch the sink on/off, fill the jug from the sink (effective only when the sink is on), and pour water from the jug into the flower pot. After each action, a tick advances processes: the sink supplies water if on; the pot passively transfers its water to the pea; and the pea consumes water and progresses from seed  $\rightarrow$  sprout  $\rightarrow$  young plant  $\rightarrow$  mature  $\rightarrow$  reproducing when sufficiently hydrated for several consecutive ticks. Episodes start with an unplanted pea and an empty pot; the goal is to plant the pea and water it repeatedly until it reaches the reproducing stage.

# Prev SOTA vs Agent2World.

```
2123
      import random
2124
       # [D0] Toggle: False = Ours, True = Baseline
2125
      BASELINE = False
2126
2127
                      --- Core object model (minimal API) ----
2128
      class GameObject:
2129
          def __init__(self, name):
               self.name, self.parent, self.contains = name, None, []
2130
              self.props = {"isContainer": False, "isMoveable": True}
          def get(self, k, d=None):
2132
              return self.props.get(k, d)
2133
          def add(self, obj):
               obj.removeSelf(); self.contains.append(obj);
2134
               obj.parent = self
2135
          def remove(self, obj):
2136
               self.contains.remove(obj); obj.parent = None
2137
          def removeSelf(self):
               if self.parent: self.parent.remove(self)
2138
          def allContained(self):
2139
               out = []
2140
               for o in self.contains: out += [o] + o.allContained()
2141
               return out
          def tick(self): pass
2142
2143
      class Container(GameObject):
2144
          def __init__(self, name): super().__init__(name);
2145
               self.props["isContainer"] = True
2146
           def place(self, obj):
               if not obj.get("isMoveable"):
2147
               return ("Can't move that object.", False)
2148
               self.add(obj); return ("OK.", True)
2149
          def take(self, obj):
2150
               if obj not in self.contains:
                   return ("Object not here.", None, False)
2151
               if not obj.get("isMoveable"):
2152
                   return ("Can't move that object.", None, False)
2153
               obj.removeSelf(); return ("OK.", obj, True)
2154
      class Device(Container):
2155
          def __init__(self, name): super().__init__(name);
2156
               self.props.update({"isDevice": True, "isOn": False})
2157
          def turnOn(self):
2158
               if self.props["isOn"]:
2159
                   return (f"{self.name} is already on.", False)
               self.props["isOn"] = True
```

```
2160
               return (f"{self.name} turned on.", True)
2161
           def turnOff(self):
2162
               if not self.props["isOn"]:
2163
                   return (f"{self.name} is already off.", False)
               self.props["isOn"] = False
2164
               return (f"{self.name} turned off.", True)
2165
2166
      class World(Container):
           def __init__(self): super().__init__("world")
2167
2168
       class Agent (Container):
2169
           def __init__(self): super().__init__("entity")
2170
2171
       # ----- Task objects -----
       class Pea(GameObject):
2172
           STAGES = ["seed", "sprout", "young plant", "mature plant", "reproducing"]
2173
           MAX_WATER, CONSUME, NEED, TICKS = 100, 5, 30, 3
2174
           def __init__(self):
               super().__init__("pea"); self.props["isMoveable"]=True
self.stage, self.water, self.hydrated = 0, 0, 0
2175
2176
           @property
2177
           def stage_name(self):
2178
               return self.STAGES[self.stage]
2179
           def addWater(self, n):
2180
               self.water = min(self.water + n, self.MAX_WATER)
           def tick(self):
2181
               \# [D3] Growth rule: Baseline = simple (>=2 -> +stage, else -1 if
2182
               \rightarrow >0); Ours = threshold + accumulation.
2183
               if BASELINE:
2184
                   if self.stage < len(self.STAGES)-1:
                        if self.water >= 2: self.water -= 2; self.stage += 1
2185
                        elif self.water > 0: self.water -= 1
2186
                    return
2187
               self.water = max(self.water - self.CONSUME, 0)
2188
               if self.water >= self.NEED:
2189
                   self.hydrated += 1
                   if self.hydrated >= self.TICKS and self.stage <</pre>
2190
                    → len(self.STAGES)-1:
2191
                       self.stage += 1; self.hydrated = 0
2192
2193
       class FlowerPot(Container):
           MAX_WATER = 100
2194
           def __init__(self):
2195
               super().__init__("flower pot")
2196
               self.water = 0
2197
           def addWater(self, n):
               add = min(self.MAX_WATER - self.water, n);
2198
               self.water += add; return add
2199
           def consume(self, n):
2200
               use = min(self.water, n); self.water -= use; return use
2201
           def tick(self):
2202
               # [D2] Passive transfer: Baseline = none; Ours = transfer
               → pot.water to pea on each tick.
2203
               if BASELINE: return
2204
               pea = next((o for o in self.contains if isinstance(o, Pea)),
2205
               → None)
2206
               if pea and self.water > 0:
                   x = self.consume(min(self.water, Pea.MAX_WATER))
2207
                   pea.addWater(x)
2208
2209
       class Jug(Container):
2210
           MAX_WATER = 100
           def __init__(self):
2211
               super().__init__("jug")
2212
               self.water = 0
2213
           def fill(self, n):
```

```
2214
               add = min(self.MAX_WATER - self.water, n)
2215
               self.water += add
2216
               return add
2217
           def pour(self, n):
               out = min(self.water, n
2218
               self.water -= out
2219
               return out
2220
      class Sink (Device):
2221
           MAX_WATER = 1000
2222
           def __init__(self): super().__init__("sink");
2223
               self.water = self.MAX_WATER
2224
           def tick(self):
2225
               self.water = self.MAX_WATER if self.props["isOn"] else 0
2226
       # ----- Minimal game scaffold (only actions we need) -----
2227
       class TextGame:
2228
           MAX\_STEPS = 50
2229
           def __init__(self, seed=0):
               random.seed(seed)
2230
               self.world, self.agent = World(), Agent();
2231

    self.world.add(self.agent)

2232
               self.pea, self.pot, self.jug, self.sink = Pea(), FlowerPot(),
2233
               \hookrightarrow Jug(), Sink()
               # [D7] Movability: Baseline pins the sink as immovable (ours
2234

→ keeps defaults).

2235
               if BASELINE: self.sink.props["isMoveable"] = False
2236
               for o in (self.pot, self.jug, self.sink, self.pea):
2237
                   self.world.add(o)
2238
               self.score = self.steps = 0
               self.over = self.won = False
2239
2240
           # API of interest (matching both variants); unchanged helpers omitted
2241
           \hookrightarrow for brevity.
2242
           def _obj(self, name):
2243
               for o in [self.world] + self.world.allContained():
                   if o.name == name: return o
2244
               for o in self.agent.contains:
2245
                   if o.name == name: return o
2246
               return None
2247
           def calculateScore(self):
2248
               # [D5] Reward: Baseline = stage*10; Ours = stage*20 + water bonus
2249
               \hookrightarrow (<=20).
2250
               if BASELINE:
2251
                   self.score = self.pea.stage*10
2252
               else:
                   self.score = self.pea.stage*20 +
2253
                    → int(self.pea.water/Pea.MAX_WATER*20)
2254
               if self.pea.stage_name == "reproducing":
2255
                   self.won = self.over = True
2256
               if self.steps >= self.MAX_STEPS and not self.won:
                   self.over = True
2257
2258
           # ----- actions -----
2259
           def take(self, name):
2260
               o = self.\_obj(name);
               if not o: return f"No {name}."
2261
               if not o.get("isMoveable"): return f"Can't take {name}."
2262
               if o.parent != self.world: return f"{name} not here."
2263
                _, got, ok = self.world.take(o);
2264
               if ok: self.agent.add(got);
               return "OK." if ok else "Fail."
2265
2266
           def put(self, obj, cont):
2267
               o, c = self._obj(obj), self._obj(cont)
```

```
2268
               if not o or o.parent != self.agent:
2269
                   return f"No {obj} in inventory."
2270
               if not c or not c.get("isContainer"):
2271
                   return f"{cont} not a container."
               # [D6] Placement constraint: Ours restricts pea -> pot only;
2272
                \hookrightarrow Baseline has no special rule.
2273
               if (not BASELINE) and isinstance(o, Pea) and not isinstance(c,
2274
                   FlowerPot):
2275
                   return "Pea must go into flower pot."
               _, ok = c.place(o); return "OK." if ok else "Fail."
2276
2277
           def turn_on(self, dev):
2278
               d = self._obj(dev);
2279
               if not d or not d.get("isDevice"): return f"No device {dev}."
               msg,_ = d.turnOn(); return msg
           def turn_off(self, dev):
2281
               d = self._obj(dev);
2282
               if not d or not d.get("isDevice"): return f"No device {dev}."
2283
               msg,_ = d.turnOff(); return msg
2284
           def fill_from_sink(self):
2285
               # [D1] Fill gating: Baseline ignores sink.on; Ours requires
2286

    sink.on == True.

2287
               if (not BASELINE) and (not self.sink.props["isOn"]): return "Sink
2288

    is off."

               need = self.jug.MAX_WATER - self.jug.water
2289
               if need <= 0: return "Jug already full."
2290
               self.jug.fill(need) # treat sink as infinite when allowed
2291
               return "Jug filled."
2292
           def pour_to_pot(self):
               if self.jug.water <= 0: return "Jug empty."</pre>
2294
               # [D8] Pour semantics: Baseline feeds pea directly; Ours fills
2295
               → pot; pea drinks via [D2].
2296
               poured = self.jug.pour(10)
2297
               if BASELINE and (self.pea in self.pot.contains):
                   self.pea.addWater(3); return "Poured; pea absorbs water."
2298
               added = self.pot.addWater(poured)
2299
               if added < poured: self.jug.fill(poured - added)</pre>
2300
               return "Poured into pot.'
2301
           # ----- driver -----
2302
           def step(self, cmd):
2303
               self.steps += 1
2304
               # [D4] Update order: Baseline ticks BEFORE action; Ours ticks
2305
               if BASELINE:
2306
                   for o in [self.world] + self.world.allContained(): o.tick()
2307
2308
               parts = cmd.lower().strip().split()
2309
               out = "Unknown."
               try:
2310
                   if parts[:1] == ["take"]: out = self.take(" ".join(parts[1:]))
2311
                   elif parts[:1] == ["put"] and "in" in parts:
2312
                        i = parts.index("in");
2313
                        out = self.put(" ".join(parts[1:i]), "
2314
                        \hookrightarrow ".join(parts[i+1:]))
                   elif parts[:2] == ["turn", "on"]:
2315
                        out = self.turn_on(" ".join(parts[2:]))
2316
                   elif parts[:2] == ["turn", "off"]:
2317
                        out = self.turn_off(" ".join(parts[2:]))
2318
                   elif parts[:3] == ["fill", "jug", "from"]:
2319
                        out = self.fill_from_sink()
                   elif parts[:4] == ["pour", "water", "from", "jug"] and "in" in
2320
                    → parts:
2321
                        out = self.pour_to_pot()
```

```
2322
                   # distractor (spec only)
2323
                   elif parts[:1] == ["use"]: out = "Nothing happens."
2324
               except Exception as e:
                   out = f"Error: {e}"
2325
2326
               if not BASELINE:
2327
                   for o in [self.world] + self.world.allContained(): o.tick()
2328
               old = self.score; self.calculateScore();
               reward = self.score - old
               return out, self.score, reward, self.over, self.won
2330
```

**Analysis.** Under identical initialization and evaluation (capacity limits and preconditions enforced), AGENT2WORLD outperforms a Baseline on the pea-growing (water-transfer) task, yielding higher success, shorter trajectories, and fewer invalid actions.(i) Action space and dynamics. We expose a precondition-aware interface and decouple water flow from uptake: fill is effective only when the sink is on ([D1]); pour increases the pot's water and the pea hydrates asynchronously via tick ([D2], [D8]). We advance environment dynamics after the action to preserve causal credit assignment ([D4]). By contrast, the Baseline exposes fill irrespective of sink state, credits hydration at pour time, and updates before acting.(ii) Physical consistency and constraints. We enforce finite capacities with overflow returned to the jug and constrain placement so the pea can only be planted in the flower pot ([D6]). These constraints prune degenerate branches without removing valid solutions. The Baseline omits the planting constraint and hydrates synchronously, which increases misleading transitions. (Regarding movability, the Baseline pins the sink as immovable while Ours keeps defaults; this ablation affects search but not preconditions, [D7].)(iii) Growth model and **reward.** Plant physiology follows thresholded, accumulated growth with per-tick water consumption ([D3]). Reward shaping combines stage progress with a bounded water bonus, and immediate rewards are score deltas ([D5]). The Baseline uses a stage-only score without water shaping, weakening the learning signal.

**Summary of diffs:** [D1] preconditioned fill; [D2] passive pot→pea transfer; [D3] threshold+consumption growth; [D4] post-action ticking; [D5] shaped reward (stage+water); [D6] pea→pot placement constraint; [D7] sink movability ablation; [D8] pour affects pot first (not the pea).