

STORI: A BENCHMARK AND TAXONOMY FOR STOCHASTIC ENVIRONMENTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Reinforcement learning (RL) techniques have achieved impressive performance on simulated benchmarks such as Atari100k, yet recent advances remain largely confined to simulation and show limited transfer to real-world domains. A central obstacle is environmental stochasticity, as real systems involve noisy observations, unpredictable dynamics, and non-stationary conditions that undermine the stability of current methods. Existing benchmarks rarely capture these uncertainties and favor simplified settings where algorithms can be tuned to succeed. The absence of a well-defined taxonomy of stochasticity further complicates evaluation, as robustness to one type of stochastic perturbation, such as sticky actions, does not guarantee robustness to other forms of uncertainty. To address this critical gap, we introduce STORI (STOchastic-ataRI), a benchmark that systematically incorporates diverse stochastic effects and enables rigorous evaluation of RL techniques under different forms of uncertainty. We propose a comprehensive five-type taxonomy of environmental stochasticity and demonstrate systematic vulnerabilities in state-of-the-art model-based RL algorithms through targeted evaluation of DreamerV3 and STORM. Our findings reveal that world models dramatically underestimate environmental variance, struggle with action corruption, and exhibit unreliable dynamics under partial observability. We release the code and benchmark publicly at <https://anonymous.4open.science/r/stori-353D>, providing a unified framework for developing more robust RL systems.

1 INTRODUCTION

Reinforcement learning (RL) techniques have achieved impressive performance on simulated benchmarks such as Atari100k, yet recent advances remain largely confined to simulation and show limited transfer to real-world domains. A central obstacle is environmental stochasticity, as real systems involve noisy observations, unpredictable dynamics, and non-stationary conditions that undermine the stability of current methods (Antonoglou et al., 2022; Paster et al., 2022). This challenge is especially acute for model-based RL, which must build world models to capture environment dynamics, a task that becomes significantly more complex when the environment exhibits multiple forms of uncertainty.

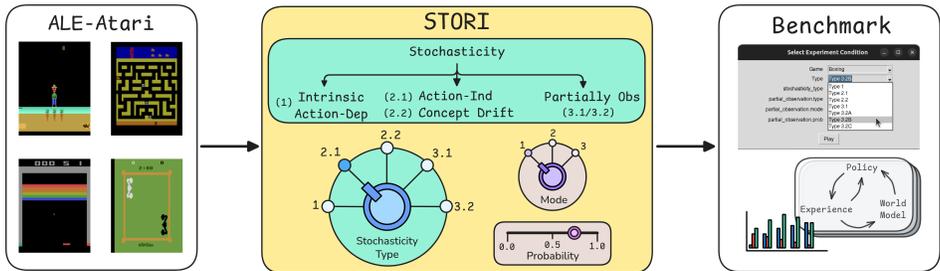
However, we lack a comprehensive stochastic environment benchmark that enables systematic development of RL methods robust to environmental stochasticity. Most widely used benchmarks, such as Atari games in the Arcade Learning Environment (ALE) (Bellemare et al., 2013), are deterministic or nearly deterministic (Paster et al., 2022). Although several approaches have introduced limited stochasticity, including sticky actions (Machado et al., 2018), no-ops (Mnih et al., 2015), human starts (Nair et al., 2015), and random frame skips (Brockman et al., 2016)—these modifications remain narrow in scope. To develop truly robust RL agents, we need benchmarks that systematically incorporate diverse forms of environmental uncertainty with granular control over both types and intensities of stochastic effects.

In this paper, we introduce STORI (STOchastic-ataRI), a benchmark that systematically incorporates diverse stochastic effects and enables rigorous evaluation of RL techniques under different forms of uncertainty. Alongside, we propose an updated taxonomy of stochasticity in RL environments, providing a unified framework for analyzing and comparing approaches. We leverage STORI to systematically investigate the reliability of world models under diverse forms of stochasticity and

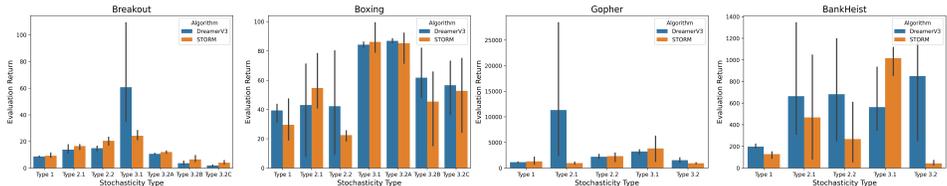
perform targeted evaluation probes to examine how learned dynamics respond to different stochastic challenges.

Our key contributions include:

- **A comprehensive stochasticity taxonomy** with five distinct types: action-dependent noise, action-independent randomness, concept drift, representation learning challenges, and missing state information
- **STORI benchmark implementation** that systematically incorporates these stochasticity types into four Atari environments with granular parameter control
- **Systematic evaluation** of state-of-the-art model-based RL algorithms (DreamerV3 and STORM) revealing fundamental vulnerabilities to environmental uncertainty
- **Targeted failure mode analysis** demonstrating that world models systematically underestimate variance, struggle with action corruption, and show unreliable dynamics under partial observability
- **Open-source framework** enabling researchers to develop and evaluate stochasticity-aware RL algorithms. We release the code and benchmark publicly at <https://anonymous.4open.science/r/stori-353D>



(a) STORI benchmark framework



(b) Performance across stochasticity types

Figure 1: STORI benchmark and results. (a) Framework for systematic stochasticity evaluation. (b) DreamerV3 and STORM performance degradation under uncertainty (Types: 1=action corruption, 2.1=random events, 2.2=concept drift, 3.1=default, 3.2=missing information).

2 RELATED WORKS

Stochastic Environment Benchmarks Recent efforts have addressed limitations of deterministic RL benchmarks by incorporating stochastic perturbations. Robust-Gymnasium (Gu et al., 2025) provides a modular framework for robust evaluation across sixty robotics and control tasks, introducing observation-disruptors, action-disruptors, and environment-disruptors for systematic robustness assessment. Similarly, Zouitine et al. (2024) introduced a benchmark extending Gymnasium-MuJoCo with six tasks capturing environmental shifts. STORI shares similar motivations but offers complementary contributions. First, it uses Atari games as a canonical testbed for discrete, high-dimensional decision-making. Second, STORI explicitly includes temporal non-stationarities such as concept drift, a critical yet underexplored aspect of real-world uncertainty. Beyond comprehensive benchmarks, several works target specific perturbation types. Zhang et al. (2020) examined adversarial state perturbations, Han et al. (2024) studied multi-agent adversarial attacks, and Park et al. (2025) analyzed offline goal-conditioned RL under perturbations. These reveal vulnerabilities

to specific noise sources but lack unified robustness suites. Sandbox frameworks like MiniHack (Samvelyan et al., 2021) allow custom stochastic environments.

Stochasticity Taxonomy Early RL literature formalized uncertainty through Partially Observable Markov Decision Processes (POMDPs) (Sutton & Barto, 2018), addressing partial observability and stochastic transitions. Recent works like Vamplew et al. (2022) propose broader stochastic MDP classifications for policy evaluation and learning. Beyond formal stochasticity, Liu et al. (2023) introduces environment heterogeneity concerning spatial layout and dynamics variations. Lu et al. (2018) addresses temporal non-stationarity or concept drift, where target concepts change due to shifting hidden contexts. STORI’s taxonomy integrates existing perspectives into a unified, practical framework explicitly designed for benchmarking, allowing systematic instantiation of stochasticity classes as configurable perturbations.

World Model Benchmarks Recent world modeling advances include DreamerV3 (Hafner et al., 2024), IRIS (Micheli et al., 2023), STORM (Zhang et al., 2023), TransDreamer (Chen et al., 2024), TWM (Robine et al., 2023), and DIAMOND (Alonso et al., 2024). These have been evaluated using Atari 100k (Ye et al., 2021), BSuite (Osband et al., 2020), Crafter (Hafner, 2021), and DMLab (Beattie et al., 2016). STORI introduces a stochasticity-driven framework to analyze world model performance under environmental uncertainties.

3 ENVIRONMENT STOCHASTICITY AND OUR MOTIVATION

3.1 ENVIRONMENT STOCHASTICITY

RL environments can be categorized by their predictability. Deterministic environments have fully predictable outcomes, while stochastic environments introduce uncertainty requiring agents to handle outcome variability. Following Kumar & Varaiya (1986), stochasticity includes: **Stationary Stochastic (Objective)** with consistent statistical distributions (coin tosses); **Stationary Stochastic (Subjective)** based on personal beliefs (expert forecasts); **Non-Stationary Stochastic** with time-varying dynamics (traffic patterns); and **Illusory or Complex Uncertainty** where probabilities are unreliable (nuclear accidents).

3.2 MATHEMATICAL TAXONOMY OF STOCHASTICITY

We formalize five key types using transition function $P(s'|s, a)$ where s is the current state, a is the action, and s' is the next state.

3.2.1 TYPE 1: INTRINSIC ACTION-DEPENDENT STOCHASTICITY

Unreliable action channels corrupt intended action a into executed action \tilde{a} .

$$P_{AD}(s'|s, a) = \sum_{\tilde{a} \in \mathcal{A}} C(\tilde{a}|a)P(s'|s, \tilde{a}) \quad (1)$$

For sticky actions: $C(\tilde{a}|a) = (1 - \alpha)\mathbb{I}\{\tilde{a} = a\} + \alpha \cdot u(\tilde{a})$

3.2.2 TYPE 2.1: INTRINSIC ACTION-INDEPENDENT STOCHASTICITY (RANDOM)

Exogenous random events independent of agent actions, with random variable $\xi \sim F(\xi)$:

$$P_{AI}(s'|s, a) = \int P(s'|s, a, \xi)dF(\xi) \quad (2)$$

3.2.3 TYPE 2.2: INTRINSIC ACTION-INDEPENDENT STOCHASTICITY (CONCEPT DRIFT)

Time-dependent dynamics $P_t(s'|s, a)$ with drift magnitude:

$$\text{Drift}(t, t + \Delta t) = \mathcal{D}(P_t(\cdot|s, a)||P_{t+\Delta t}(\cdot|s, a)) \quad (3)$$

3.2.4 TYPE 3.1: AGENT-INDUCED STOCHASTICITY (REPRESENTATION LEARNING)

Rich observations requiring representation learning where $I(S; O) \approx H(S)$. Belief states update via:

$$b_{t+1}(s') \propto O(o_{t+1}|s') \sum_{s \in \mathcal{S}} P(s'|s, a_t)b_t(s) \quad (4)$$

State aliasing is resolvable through better feature extraction. *Example:* Standard Atari RGB frames contain all game information but require learning to extract from pixels.

3.2.5 TYPE 3.2: AGENT-INDUCED STOCHASTICITY (MISSING STATE VARIABLES)

Critical state variables are completely omitted where $I(S; O) \ll H(S)$. Creates fundamental state aliasing $O(o|s_i) = O(o|s_j) = 1$ for $s_i \neq s_j$ that persists regardless of representational capacity. Same belief update as Type 3.1 but fundamentally limited.

Requires history tracking: $h_t = \{o_1, a_1, o_2, a_2, \dots, o_t\}$

Examples: Breakout with invisible ball regions; Boxing with hidden opponents; partial screen occlusion.

3.3 CHALLENGES FOR WORLD MODEL LEARNING AND MODEL BASED RL IN STOCHASTIC ENVIRONMENTS

In this section, we analyze potential challenges for MBRL in different types of stochastic environments. A world model, denoted \tilde{P}_θ , aims to learn the true transition dynamics from data. Each form of stochasticity introduces a distinct challenge that can cause a mismatch between \tilde{P}_θ and the true dynamics.

Challenge from Type 1 Stochasticity The world model must learn not only the environment’s response to actions, $P(s'|s, \tilde{a})$, but also the action channel itself, $C(\tilde{a}|s, a)$. If the model fails to account for the action channel, its predictions will be systematically biased. The prediction error is the divergence between the model’s direct prediction and the true, action-corrupted dynamics: $\text{Error} = \mathcal{D}(P_{AD}(s'|s, a) || \tilde{P}_\theta(s'|s, a))$, where P_{AD} represents the true dynamics and \tilde{P}_θ represents the model prediction. The model’s ability to control the environment is limited by the **action channel capacity**, which can be measured by the mutual information $I(A; \tilde{A}|S)$. A low-capacity channel is fundamentally difficult to model and exploit.

Challenge from Type 2.1 Stochasticity This introduces irreducible **aleatoric uncertainty** into the environment. A deterministic world model will fail completely. A probabilistic world model must accurately capture the variance of the outcomes. The world model must learn a distribution over next states. The core challenge is to match the variance of this distribution to the true environmental variance, which is inherent and cannot be reduced with more data. The prediction error is tied to the model’s ability to capture this spread: $\text{Aleatoric Error} = |\text{Var}_{s' \sim P_{AI}}[s'] - \text{Var}_{s' \sim \tilde{P}_\theta}[s']|$. The world model must avoid being overconfident in its predictions and instead represent the full range of possible outcomes.

Challenge from Type 2.2 Stochasticity Concept drift causes the world model’s learned parameters θ to become outdated. A model trained on data from time t will be inaccurate at time $t + \Delta t$. The prediction error grows over time as the environment drifts away from the data the model was trained on. The accumulated error is a function of the drift magnitude: $\text{Prediction Error}(t + \Delta t) \propto \mathcal{D}(P_t(\cdot|s, a) || \tilde{P}_\theta(\cdot|s, a))$, where \tilde{P}_θ was trained on data from distributions around time t . This forces the model to either continuously adapt its parameters or have a mechanism to detect and react to the drift.

Challenge from Type 3.1 & 3.2 Stochasticity The world model cannot operate on the true state s and must instead learn a latent state representation z_t from a history of observations $o_{1:t}$. The primary challenge is **state aliasing**. The uncertainty a world model faces is not just the environment’s true stochasticity, but also the aliasing-induced variance. The total variance in outcomes given an observation o is decomposed as: $\text{Var}[s'|o, a] = \mathbb{E}_{s \sim b(s|o)}[\text{Var}[s'|s, a]] + \text{Var}_{s \sim b(s|o)}[\mathbb{E}[s'|s, a]]$, where the first term represents true aleatoric uncertainty and the second term represents aliasing-induced uncertainty. The world model’s latent dynamics, $\tilde{P}_\theta(z'|z, a)$, must implicitly handle the second term, which is purely an artifact of perception. In Type 3.2 environments, where entire state variables are missing, this aliasing uncertainty can become overwhelmingly large, making it nearly impossible to form an accurate belief state and rendering long-term prediction unreliable.

4 STORI - A BENCHMARK OF STOCHASTIC ENVIRONMENTS FOR RL

In this section, we describe in details the benchmark environments we built for different types of stochasticity based on Atari-Arcade learning environment. Atari games such as Breakout, Boxing, Gopher and BankHeist were modified to allow fine-grain control of these stochasticity. The taxonomy for stochasticity in STORI is an extension of the summary of classification of stochasticity according to Antonoglou et al. (2022). Table 1 presents the taxonomy of stochasticity, listing

each type, its corresponding subtype, and the associated ID. Each type is explained in the following sections.

Table 1: Environment stochasticity taxonomy

ID	Type	Sub Type
0	Deterministic	NA
1	Action Dependent	NA
2.1	Action Independent	Random
2.2	Action Independent	Concept Drift
3.1	Partially Observed	Representation
3.2	Partially Observed	Missing State

ATARI - ARCADE LEARNING ENVIRONMENT

The Arcade Learning Environment (ALE) provides a foundational framework for applying RL to Atari 2600 games (Bellemare et al., 2013). Built on the Stella emulator and integrated with Gymnasium (Brockman et al., 2016), ALE supports over a hundred games with extensive configurability including observation types (RGB, grayscale, RAM), action spaces, and stochasticity parameters like sticky actions (Machado et al., 2018). The Atari 100K benchmark (Ye et al., 2021) evaluates sample efficiency by assessing agents after only 100,000 environment steps (approximately two hours of gameplay).

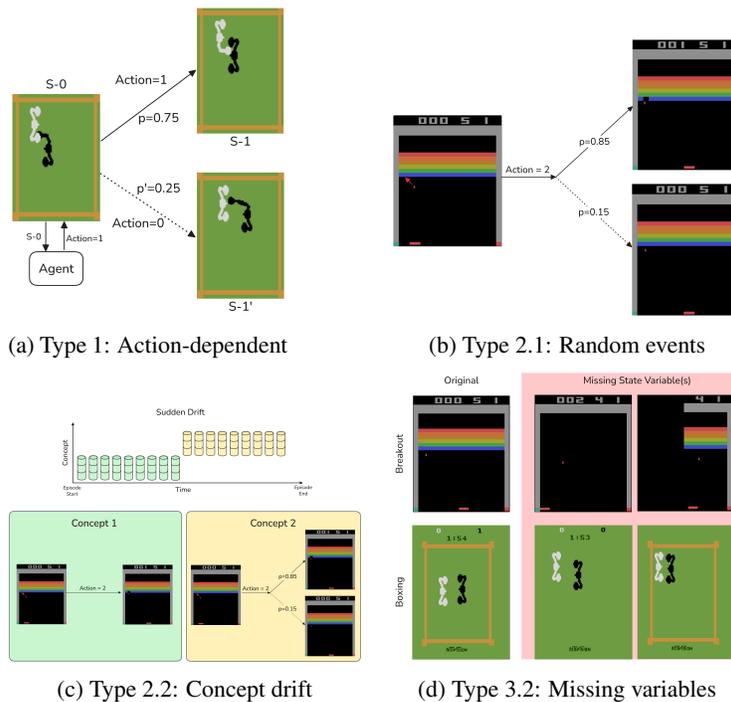


Figure 2: Stochasticity types in STORI benchmark.

TYPE 0: DETERMINISTIC ENVIRONMENT

Deterministic environments are those in which the next state is fully determined by the current state and the action taken. The state is completely observable and there is no randomness in the state transitions or rewards, meaning that the outcome of any action is predictable.

In the case of Atari, we consider the ground-truth labels of various state variables obtained directly from the RAM for each observation, following the approach of Anand et al. (2020). No additional

270 stochasticity parameters are introduced, meaning that the environment is fully deterministic and
 271 corresponds to Type 0 in our taxonomy. Example for Breakout can be seen in figure 5.
 272

273 TYPE 1: INTRINSIC ACTION DEPENDENT STOCHASTIC ENVIRONMENT

274 In environments with action-dependent intrinsic stochasticity, the environment may replace the
 275 agent’s chosen action with different action, by default, with a random one. For instance, in
 276 *sticky_action* (Machado et al., 2018) scenarios, the environment can repeat the previous action
 277 with some probability. This results in varied outcomes even from the same state, with the stochastic
 278 effects limited to the state variables that can be influenced by the agent’s actions. An example of
 279 action-dependent intrinsic stochasticity with Atari Boxing can be seen in the 2a.
 280

281 TYPE 2.1: INTRINSIC ACTION INDEPENDENT STOCHASTIC ENVIRONMENT - RANDOM

282 In action-independent random stochastic environments, randomness arises independently of the
 283 agent’s choices and affects state variables outside the agent’s direct control. This stochasticity, often
 284 due to external factors or inherent environmental noise, means that even with complete knowledge
 285 of the environment and carefully chosen actions, the next state cannot be predicted with certainty.

286 The figure 2b illustrates an example of how this type of stochasticity can be modeled in Atari Break-
 287 out, where the paddle is moved to the right while the ball is on a trajectory to hit a block. In this
 288 case, there is a 0.15 probability that the ball bounces back without destroying the block or yielding
 289 any reward. Notably, this stochastic behavior is independent of the action of moving the paddle to
 290 the right.
 291

292 TYPE 2.2: INTRINSIC ACTION INDEPENDENT STOCHASTIC ENVIRONMENT - CONCEPT DRIFT

293 Environments with intrinsic action-independent concept drift can change over time independently
 294 of the agent’s actions, a phenomenon known as concept drift (Lu et al., 2018), which can generally
 295 be categorized into three types according to how the drift unfolds over time. In *sudden drift*, the
 296 environment undergoes abrupt changes, forcing the agent to quickly adapt to new dynamics. In
 297 *gradual or incremental drift*, the transition to new dynamics occurs slowly over time, requiring the
 298 agent to continuously adjust its policy. Finally, in *recurring drift*, previously observed dynamics
 299 reappear in a cyclical or context-dependent manner, making long-term adaptation more challenging.
 300 Learning in such environments demands flexibility and the ability to detect and respond to changes.

301 In the case of Atari, most games have intrinsic incremental drift. As the agent levels up in the game,
 302 the difficulty of the game also increases. With a more fine-grain control over concept drift, other
 303 types of drift can also be achieved in Atari games as shown in figure 2c.
 304

305 TYPE 3.1: PARTIALLY OBSERVED ENVIRONMENT (REPRESENTATION LEARNING)

306 In partially observed environments, the agent does not have access to the full state information.
 307 When the state variables are represented differently from the true underlying environment, the agent
 308 must infer hidden information or learn an appropriate representation. This increases the complexity
 309 of decision making, since the agent must rely on approximate observations.

310 A typical example is the **Default Atari setting**, where the agent perceives only the screen image
 311 produced by the emulator after each action. These images are designed to approximate the true
 312 state, but they do not capture it fully. To enrich the observation, many implementations use a 4-
 313 frame skip with aggregation, which allows the agent to infer additional information, such as motion
 314 or rate of change over time, that is not apparent from a single frame.
 315

316 TYPE 3.2: PARTIALLY OBSERVED ENVIRONMENT (MISSING STATE VARIABLE(S))

317 An important subclass of partially observed environments arises when information about certain
 318 state variables is missing, leaving the agent unable to observe critical aspects of the environment.
 319 This lack of information demands strategies that can manage uncertainty and make robust decisions
 320 despite gaps in perception. Such environments are common in real-world scenarios where sensors
 321 are limited, noisy, or unreliable.

322 Figure 2d illustrates type 3.2 environments using Atari Breakout and Boxing. In Breakout, examples
 323 include invisible blocks or a partially hidden screen, while in Boxing, examples include a hidden
 boxing ring or concealed clock and score information.

5 EXPERIMENTS AND RESULTS

We evaluated DreamerV3 (Hafner et al., 2024) and STORM (Zhang et al., 2023) using STORI. DreamerV3 features a learned world model with actor-critic architecture, achieving robustness through fixed hyperparameters, normalization, and effective scaling. STORM employs a transformer backbone with stochastic variational modeling for strong sequence modeling and robustness. We focus on MBRL methods as model-free approaches require substantially more computational resources for meaningful results.

We selected four Atari 100K environments: *Breakout*, *Boxing*, *Gopher* (Agent-Optimal (Lim et al., 2025)), and *BankHeist* (Human-Optimal (Lim et al., 2025)). Breakout, Boxing, and Gopher are classified as Agent-Optimal tasks, where agents achieve strong performance under default conditions, facilitating the study of performance degradation under stochasticity. BankHeist, a Human-Optimal task, provides a preliminary understanding of how added stochasticity interacts with environments where agents perform substantially worse than humans. The set also offers diversity in action-space sizes: Breakout (4 actions), Gopher (8 actions), and Boxing/BankHeist (18 actions each), ensuring that robustness observations are not biased by a particular action-space complexity.

For each environment, stochasticity types were assigned mostly at random, except for Type 2.2 (concept drift), where Concept 2 was intentionally selected from other stochasticity types to enable controlled comparison between “pure” stochasticity and its drifting variant. Multiple parameter settings for each stochasticity type and game were provided, allowing users of the benchmark to select perturbation regimes relevant to their research focus. Implementation details are in Appendix A.1, and experiment-specific stochasticity settings are provided in Appendix B.1. Each algorithm was trained for 100K steps across baseline and modified environments for three seeds and evaluated on 100 episodes, with mean return reported. Note that Type 3.1 is the default ALE environment setting.

5.1 PERFORMANCE OF DREAMERV3 AND STORM IN DIFFERENT STOCHASTIC ENVIRONMENTS

5.1.1 BREAKOUT

Stochasticity introduction caused marked performance decline versus default Type 3.1 environment (Figure 1b), aligning with theoretical predictions. DreamerV3 initially outperformed STORM (60.71 ± 41.89 vs 24.17 ± 3.55) but STORM showed greater robustness across stochasticity types.

Table 2: Variance underestimation by world models.

Model	Type	Diff.
DreamerV3	3.1	1.25
DreamerV3	2.1	300.34
STORM	3.1	1.32
STORM	2.1	325.21

Breakout’s small action space (4 actions) creates high sensitivity to perturbations as incorrect LEFT/RIGHT actions immediately cause failure. Unlike Boxing, Breakout offers no recovery margin, amplifying uncertainty’s long-term impact.

Type 2.1 environments caused severe struggles, with performance dropping to 15% of baseline, confirming that irreducible aleatoric uncertainty challenges deterministic world models. Type 2.2 concept drift (default \rightarrow Type 3.2A after 300 steps) showed better performance than standalone Type 3.2A, suggesting adaptive mechanisms can leverage temporal structure.

ture.

5.1.2 BOXING

Boxing showed less pronounced performance decline. STORM initially led (86.18 ± 11.29 vs 84.22 ± 1.68) but DreamerV3 outperformed across several stochasticity types. Boxing’s resilience stems from: (1) larger action space (18 actions) providing redundancy with functionally similar actions, and (2) recovery mechanisms through retreating/repositioning.

Type 3.2A (hidden score/clock) counterintuitively improved DreamerV3 performance (86.90 ± 1.33 vs 84.22 ± 1.68), suggesting non-essential information removal simplifies representation learning.

For Type 3.2B (75% right-half occlusion), agents showed adaptive behavior as they confined opponents to visible areas, transforming partial observability into strategic constraints.

Type 2.2 concept drift performed worse than standalone Type 3.2C, except DreamerV3’s third seed learned to maximize early-episode scores before opponent invisibility, demonstrating strategic adaptation to predictable timing.

5.1.3 GOPHER AND BANKHEIST

Gopher showed high variability with Type 2.1 producing anomalous DreamerV3 performance ($11,333.53 \pm 14,761.12$) due to beneficial reward cancellation dynamics—suggesting implementation-specific edge case exploitation.

BankHeist exhibited divergent Type 3.2 outcomes: DreamerV3 achieved 849.80 ± 514.15 while STORM fell to 43.10 ± 22.85 . DreamerV3 consistently adopted an unconventional policy, remaining near city gates and triggering inter-city transitions to loot nearby banks while avoiding stochastic modifications—effectively exploiting structural features to reduce tasks to near-deterministic sub-problems rather than demonstrating genuine robustness. STORM explored broadly within cities, exposing itself to full stochastic impact. This highlights that high returns may reflect reward-maximizing shortcuts exploiting environment dynamics rather than genuine uncertainty resilience.

5.2 ANALYSIS OF ERROR TYPES AND WORLD MODEL FAILURES

To understand the specific failure modes of model-based RL under different stochasticity types, we conducted targeted analyses for some error categories defined in Section 3.3.

5.2.1 ERRORS CAUSED BY TYPE 2.1 STOCHASTICITY: ALEATORIC UNCERTAINTY

We ran a controlled probe of a single repeated action (action 3) for 1000 steps in both Type 3.1 (default setting, partially observed) and Type 2.1 (action-independent stochasticity) BankHeist environments, collecting states from the environment and predictions from the world models from DreamerV3 and STORM (Table 2). The resulting variance differences include:

Environment variance difference:

$$\text{Var}_{\text{env}}(\text{Type 2.1}) - \text{Var}_{\text{env}}(\text{Type 3.1})$$

→ DreamerV3: 299.097, STORM: 323.904. This confirms that Type 2.1 environments exhibit significantly higher true variance due to stochasticity.

Model variance difference:

$$\text{Var}_{\text{model}}(\text{Type 2.1}) - \text{Var}_{\text{model}}(\text{Type 3.1})$$

→ DreamerV3: 0.00465, STORM: 0.01216. Both models predict nearly identical variance between environments despite the true variance increasing substantially.

Both DreamerV3 and STORM significantly underestimate the increased stochasticity present in Type 2.1 environments in BankHeist. While environment variance increases by approximately 300, the models’ predicted variances remain nearly constant. This mismatch highlights a lack of variance calibration under action-independent stochastic conditions, revealing a limitation in the world models’ ability to capture environment uncertainty accurately.

5.2.2 ERRORS CAUSED BY TYPE 2.2 STOCHASTICITY: CONCEPT DRIFT ANALYSIS

For concept drift stochasticity (*BankHeist* Type 2.2), we measured the degradation ratio of model performance before and after the drift point.

The degradation ratio is defined as:

$$\text{Degradation ratio (dyn_loss)} = \frac{\text{dyn_loss}_{\text{Concept 2}}}{\text{dyn_loss}_{\text{Concept 1}}}$$

Both DreamerV3 and STORM exhibit substantial increases in dynamics loss after the concept change, indicating that world model accuracy deteriorates significantly when the environment transitions to Concept 2.

STORM – BankHeist (Type 2.2):

- Concept 1: $\text{dyn_loss} = 4.51 \pm 2.42$
- Concept 2: $\text{dyn_loss} = 22.71 \pm 4.20$
- Degradation ratio = 5.04

DreamerV3 – BankHeist (Type 2.2):

- Concept 1: $\text{dyn_loss} = 5.43 \pm 15.86$
- Concept 2: $\text{dyn_loss} = 37.67 \pm 36.37$
- Degradation ratio = 6.94

These results highlight that both model-based agents experience a sharp decline in world model performance following the concept drift, emphasizing the challenge of adapting to non-stationary dynamics.

5.2.3 ERRORS CAUSED BY TYPE 3 STOCHASTICITY: STATE ALIASING EFFECTS

To investigate how well world models handle missing information, we designed a controlled experiment using BankHeist Type 3.2, where city blocks are randomly hidden in 75% of observations. The key question: does a model’s prediction accuracy depend on whether it can initially see the environment clearly?

Experimental design: We created six test scenarios and compared two starting conditions for each:

- **Clear-start:** Model begins with city blocks visible, takes an action, observes the result
- **Obscured-start:** Model begins with city blocks hidden, takes the same action, sees the same result

We then measured how “surprised” each model was by computing the difference in prediction error: $\Delta\text{NLL} = \text{NLL}(\text{obscured-start}) - \text{NLL}(\text{clear-start})$ as in figure 3.

Key findings: As in table 3, DreamerV3 shows positive ΔNLL values (1.15), meaning it makes significantly worse predictions when starting from obscured observations. In contrast, STORM shows negative ΔNLL values (-3.32), indicating it actually performs slightly better when starting from limited information.

This reveals that DreamerV3’s world model relies heavily on having complete initial observations to make accurate predictions. When city blocks are initially hidden, DreamerV3 struggles to maintain accurate beliefs about the environment state, requiring larger “corrections” to its internal model after seeing the action’s outcome.

Critical insight: High task performance does not guarantee robust world model dynamics. Despite achieving strong returns in partially observable environments, DreamerV3’s world model is more brittle when dealing with missing information compared to STORM.

6 LIMITATIONS AND FUTURE WORK

Table 3: Model prediction errors under partial observability.

Model	ΔNLL	ΔKL
DreamerV3	1.15 ± 2.46	-0.18 ± 2.83
STORM	-3.32 ± 2.86	0.18 ± 0.62

Limited coverage of model-free baselines. Our benchmark currently focuses on DreamerV3 and STORM, two state-of-the-art model-based algorithms. This choice was driven by practical constraints of the evaluation setting: we study robustness in a 100K-interaction regime across multiple stochasticity types, magnitudes, and four Atari games, resulting in substantial computational cost. Many widely used model-free agents

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

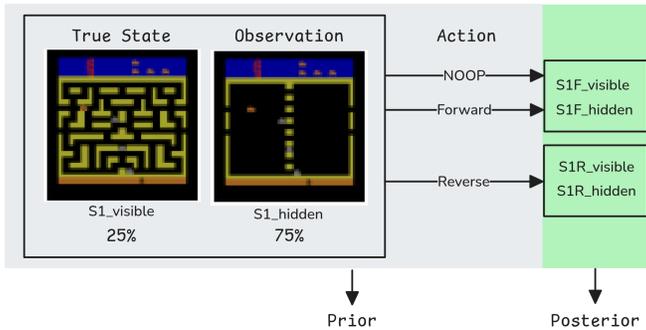


Figure 3: Partial observability probe showing prediction errors when models start with clear vs. obscured observations.

(e.g., PPO, Rainbow, IQN) are not competitive under 100K steps, as they typically require millions of frames for reasonable performance. Including them would therefore provide limited insight into robustness. Although fast-learning model-free methods such as BBF (“Bigger, Better, Faster”) exist, the public JAX implementation did not integrate cleanly with our PyTorch-based evaluation pipeline, leading to compatibility overheads. We acknowledge the value of including sample-efficient model-free algorithms and are actively extending our framework to benchmark such agents in future versions.

Cross-type comparison constraints. Our results focus on within-type rather than cross-type robustness comparisons. This reflects the fact that cross-type comparisons are meaningful only when reward semantics remain consistent across stochasticity types of a given environment. For `Boxing` and `BankHeist`, all stochasticity types preserve reward consistency and are directly comparable. For `Breakout` and `Gopher`, all types except Type 2.1 satisfy this criterion. In Type 2.1 for these games, perturbations (e.g., stochastically canceled hits) fundamentally alter the achievable returns and distort the optimal reward distribution. Such shifts make the environments structurally incomparable to other types; reporting cross-type results in these settings would be misleading.

Future directions. A promising extension of this work is to evaluate agents under *compositional* uncertainty, where multiple stochasticity types are activated simultaneously. This would more faithfully represent real-world uncertainty, where observation noise, transition randomness, and action corruption may interact in nontrivial ways. Another direction is to benchmark uncertainty-aware agents, including Bayesian world models, ensemble-based approaches, adversarially trained policies, and domain-randomization methods. Studying such techniques within our stochasticity taxonomy may yield deeper insights into how agents can model, anticipate, and adapt to diverse forms of uncertainty.

7 CONCLUSIONS

We introduced STORI, a systematic benchmark with a five-type taxonomy for evaluating RL algorithms under environmental stochasticity: action-dependent noise, action-independent randomness, concept drift, representation learning challenges, and missing state information.

Evaluation of DreamerV3 and STORM revealed systematic vulnerabilities in model-based approaches. Both algorithms struggle with action corruption, underestimate environmental variance by 300x, degrade 5-7x after concept drift, and show inconsistent reliability under partial observability. Strong task performance does not guarantee robust world model dynamics. STORI provides a foundation for building more robust RL systems capable of handling real-world uncertainty.

8 REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our findings. All data, code, and implementation details necessary to replicate our experiments will be made available to the research community. Careful documentation accompanies the released resources to facilitate independent verification and reuse. The authors affirm that the results reported in this paper can be fully reproduced using the provided materials.

9 ETHICS STATEMENT

This work was conducted in accordance with established ethical standards for scientific research. All methods, analyses, and interpretations were carried out with a commitment to transparency, integrity, and responsible reporting. The authors confirm that no part of this research involved practices that could compromise fairness, safety, or the ethical treatment of data, participants, or systems.

REFERENCES

- Eloi Alonso, Adam Jelley, Vincent Micheli, Anssi Kanervisto, Amos Storkey, Tim Pearce, and François Fleuret. Diffusion for world modeling: Visual details matter in atari. In *Thirty-eighth Conference on Neural Information Processing Systems*, 2024. URL <https://arxiv.org/abs/2405.12399>.
- Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon Hjelm. Unsupervised state representation learning in atari, 2020. URL <https://arxiv.org/abs/1906.08226>.
- Ioannis Antonoglou, Julian Schrittwieser, Sherjil Ozair, Thomas K Hubert, and David Silver. Planning in stochastic environments with a learned model. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=X6D9bAHhBQ1>.
- Charles Beattie, Joel Z. Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler, Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, Julian Schrittwieser, Keith Anderson, Sarah York, Max Cant, Adam Cain, Adrian Bolton, Stephen Gaffney, Helen King, Demis Hassabis, Shane Legg, and Stig Petersen. Deepmind lab, 2016. URL <https://arxiv.org/abs/1612.03801>.
- M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, jun 2013.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016. URL <https://arxiv.org/abs/1606.01540>.
- Chang Chen, Yi-Fu Wu, Jaesik Yoon, and Sungjin Ahn. Transdreamer: Reinforcement learning with transformer world models, 2024. URL <https://arxiv.org/abs/2202.09481>.
- Shangding Gu, Laixi Shi, Muning Wen, Ming Jin, Eric Mazumdar, Yuejie Chi, Adam Wierman, and Costas Spanos. Robust gymnasium: A unified modular benchmark for robust reinforcement learning. In Y. Yue, A. Garg, N. Peng, F. Sha, and R. Yu (eds.), *International Conference on Representation Learning*, volume 2025, pp. 102014–102041, 2025. URL https://proceedings.iclr.cc/paper_files/paper/2025/file/fcc22e5b7d5d2155d994da22d045f0a6-Paper-Conference.pdf.
- Danijar Hafner. Benchmarking the spectrum of agent capabilities. *arXiv preprint arXiv:2109.06780*, 2021.
- Danijar Hafner, Jurgis Pasukonis, Jimmy Ba, and Timothy Lillicrap. Mastering diverse domains through world models, 2024. URL <https://arxiv.org/abs/2301.04104>.
- Songyang Han, Sanbao Su, Sihong He, Shuo Han, Haizhao Yang, and Fei Miao. What is the solution for state-adversarial multi-agent reinforcement learning? *Transactions on Machine Learning Research (TMLR)*, 2024.

- 594 P. R. Kumar and Pravin Varaiya. *Stochastic systems: estimation, identification and adaptive control*.
595 Prentice-Hall, Inc., USA, 1986. ISBN 013846684X.
596
- 597 Jing Yu Lim, Zarif Ikram, Samson Yu, Haozhe Ma, Tze-Yun Leong, and Dianbo Liu. Jedi: Latent
598 end-to-end diffusion mitigates agent-human performance asymmetry in model-based reinforce-
599 ment learning, 2025. URL <https://arxiv.org/abs/2505.19698>.
- 600 Dianbo Liu, Vedant Shah, Oussama Boussif, Cristian Meo, Anirudh Goyal, Tianmin Shu,
601 Michael Curtis Mozer, Nicolas Heess, and Yoshua Bengio. Stateful active facilitator: Co-
602 ordination and environmental heterogeneity in cooperative multi-agent reinforcement learning.
603 In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=B4mazQLLW0_.
- 604
605
- 606 Jie Lu, Anjin Liu, Fan Dong, Feng Gu, Joao Gama, and Guangquan Zhang. Learning under con-
607 cept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2018.
608 ISSN 2326-3865. doi: 10.1109/tkde.2018.2876857. URL [http://dx.doi.org/10.1109/](http://dx.doi.org/10.1109/TKDE.2018.2876857)
609 [TKDE.2018.2876857](http://dx.doi.org/10.1109/TKDE.2018.2876857).
- 610 Marlos C. Machado, Marc G. Bellemare, Erik Talvitie, Joel Veness, Matthew J. Hausknecht, and
611 Michael Bowling. Revisiting the arcade learning environment: Evaluation protocols and open
612 problems for general agents. *Journal of Artificial Intelligence Research*, 61:523–562, 2018.
613
- 614 Vincent Micheli, Eloi Alonso, and François Fleuret. Transformers are sample-efficient world
615 models. In *The Eleventh International Conference on Learning Representations*, 2023. URL
616 <https://openreview.net/forum?id=vhFulAcb0xb>.
- 617 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G.
618 Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fijdeland, Georg Ostrovski, Stig Pe-
619 tersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmashan Kumaran, Daan
620 Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement
621 learning. *Nature*, 2015. doi: 10.1038/nature14236. URL [https://doi.org/10.1038/](https://doi.org/10.1038/nature14236)
622 [nature14236](https://doi.org/10.1038/nature14236).
- 623
- 624 Arun Nair, Praveen Srinivasan, Sam Blackwell, Cagdas Alcicek, Rory Fearon, Alessandro De
625 Maria, Vedavyas Panneershelvam, Mustafa Suleyman, Charles Beattie, Stig Petersen, Shane
626 Legg, Volodymyr Mnih, Koray Kavukcuoglu, and David Silver. Massively parallel methods for
627 deep reinforcement learning, 2015. URL <https://arxiv.org/abs/1507.04296>.
- 628 Ian Osband, Yotam Doron, Matteo Hessel, John Aslanides, Eren Sezener, Andre Saraiva, Katrina
629 McKinney, Tor Lattimore, Csaba Szepesvári, Satinder Singh, Benjamin Van Roy, Richard Sutton,
630 David Silver, and Hado van Hasselt. Behaviour suite for reinforcement learning. In *International*
631 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?](https://openreview.net/forum?id=rygf-kSYwH)
632 [id=rygf-kSYwH](https://openreview.net/forum?id=rygf-kSYwH).
- 633
- 634 Seohong Park, Kevin Frans, Benjamin Eysenbach, and Sergey Levine. Ogbench: Benchmarking
635 offline goal-conditioned rl, 2025. URL <https://arxiv.org/abs/2410.20092>.
- 636 Keiran Paster, Sheila A. McIlraith, and Jimmy Ba. You can’t count on luck: why decision trans-
637 formers and rvs fail in stochastic environments. In *Proceedings of the 36th International Conference*
638 *on Neural Information Processing Systems*, NIPS ’22, Red Hook, NY, USA, 2022. Curran Asso-
639 ciates Inc. ISBN 9781713871088.
- 640
- 641 Jan Robine, Marc Höftmann, Tobias Uelwer, and Stefan Harmeling. Transformer-based world mod-
642 els are happy with 100k interactions. In *The Eleventh International Conference on Learning*
643 *Representations*, 2023. URL <https://openreview.net/forum?id=TdBaDGCpjly>.
- 644
- 645 Mikayel Samvelyan, Robert Kirk, Vitaly Kurin, Jack Parker-Holder, Minqi Jiang, Eric Hambro,
646 Fabio Petroni, Heinrich Kuttler, Edward Grefenstette, and Tim Rocktäschel. Minihack the planet:
647 A sandbox for open-ended reinforcement learning research. In *Thirty-fifth Conference on Neural*
Information Processing Systems Datasets and Benchmarks Track (Round 1), 2021. URL <https://openreview.net/forum?id=skFwlyefkWJ>.

648 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
649 Book, Cambridge, MA, USA, 2018. ISBN 0262039249.

650

651 Peter Vamplew, Cameron Foale, and Richard Dazeley. The impact of environmental stochasticity
652 on value-based multiobjective reinforcement learning. *Neural Computing and Applications*, 34
653 (3):1783–1799, 2022. ISSN 1433-3058. doi: 10.1007/s00521-021-05859-1. URL <https://doi.org/10.1007/s00521-021-05859-1>.

654

655 Weirui Ye, Shaohuai Liu, Thanard Kurutach, Pieter Abbeel, and Yang Gao. Mastering atari games
656 with limited data, 2021. URL <https://arxiv.org/abs/2111.00210>.

657

658 Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-
659 Jui Hsieh. Robust deep reinforcement learning against adversarial perturbations on state ob-
660 servations. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.),
661 *Advances in Neural Information Processing Systems*, volume 33, pp. 21024–21037. Curran
662 Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper_files/
663 paper/2020/file/f0eb6568ea114ba6e293f903c34d7488-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f0eb6568ea114ba6e293f903c34d7488-Paper.pdf).

664 Weipu Zhang, Gang Wang, Jian Sun, Yetian Yuan, and Gao Huang. STORM: Efficient stochastic
665 transformer based world models for reinforcement learning. In *Thirty-seventh Conference on
666 Neural Information Processing Systems*, 2023. URL [https://openreview.net/forum?
667 id=WxnrX42rnS](https://openreview.net/forum?id=WxnrX42rnS).

668 Adil Zouitine, David Bertoin, Pierre Clavier, Matthieu Geist, and Emmanuel Rachelson. Rrls :
669 Robust reinforcement learning suite, 2024. URL <https://arxiv.org/abs/2406.08406>.

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

A APPENDIX

A.1 STORI IMPLEMENTATION

The STORI framework is built around a sophisticated wrapper-based architecture that introduces various types of uncertainty and partial observability into deterministic Atari environments with a granular control over the modifications.

A.1.1 CORE ARCHITECTURE AND WRAPPER SYSTEM

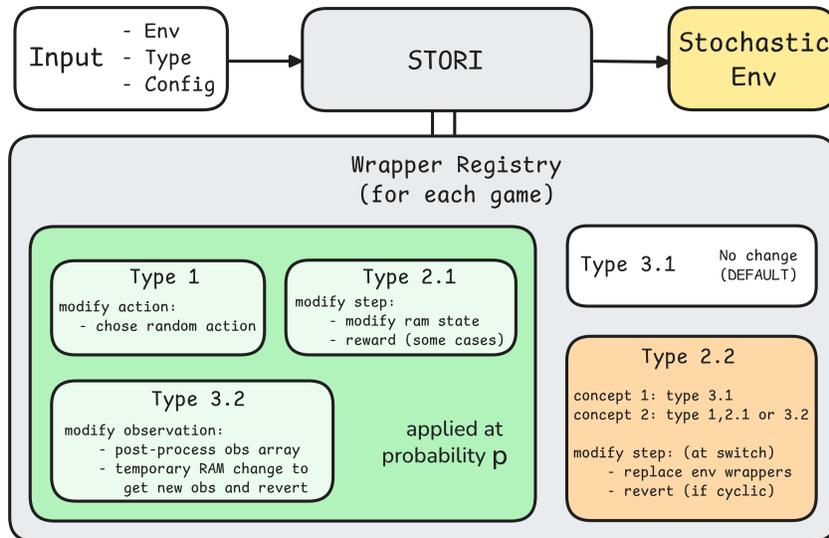


Figure 4: STORI Implementation Overview

The implementation uses a hierarchical wrapper system built on top of the Atari Learning Environment (ALE). The main ‘StochasticEnv’ class serves as the entry point, which applies different types of wrappers from ‘wrapper_registry’ of specified environment. The system supports five distinct types, each introducing different forms of stochasticity. The system is highly configurable through a dictionary-based configuration system. Users can specify probabilities for different stochasticity effects, choose between different modes of operation, and configure temporal parameters for concept drift. The wrapper registry system allows for easy extension and customization of stochasticity types for new games or research requirements.

A.1.2 STOCHASTICITY WRAPPERS

- Type 0: This type returns the RAM state of the game (a 1-D numpy array) with state labels as the observation. This implementation is an extension of Atari Annotated RAM Interface (Anand et al., 2020).
- Type 1: The ‘ActionDependentStochasticityWrapper’ randomly replaces the agent’s intended action with a random action from the action space with a specified probability.
- Type 2.1: The ‘ActionIndependentRandomStochasticityWrapper’ implements environment specific random events that occur independently of the agent’s actions. These effects are applied probabilistically and create unpredictable environmental changes to which the agent must adapt. Read more about game-specific modifications in section D.
- Type 2.2: This introduces temporal concept drift where the environment dynamics change over time. The ‘ActionIndependentConceptDriftWrapper’ supports both sudden and cyclic modes between 2 concepts. The concept 1 is the default environment (type 3.1) and concept 2 can be any other environment stochasticity types out of 1, 2.1 and 3.2. In sudden mode,

the environment switches to concept 2 after a fixed number of steps. In cyclic mode, it alternates between the concept 1 and 2 every specified number of steps, creating a challenging environment where the agent must continuously adapt to changing dynamics.

- Types 3.1: This stochasticity type returns the default ALE environment without any modifications.
- Types 3.2: The ‘PartialObservationWrapper’ introduces partial observability by modifying the agent’s observations. The system supports multiple observation modification techniques including cropping (removing portions of the screen), blackout (hiding specific regions), and RAM manipulation (temporarily modifying the game’s internal state to get modified observation).

In STORI, stochasticity types 1, 2.1, 2.2, and 3.2 are implemented as extensions of Type 3.1 environments. This is because screen-based observations serve as the default, well-studied ALE inputs for various reinforcement learning algorithms, providing a consistent foundation for comparing different types of stochasticity while also allowing for interpretable analysis of agent actions and behaviors.

A.1.3 ALGORITHMS ADDITIONAL DETAILS

- DreamerV3: The source implementation and default parameters for Atari100K config used from this code repository (MIT license): <https://github.com/NM512/dreamerv3-torch>
- STORM: The source implementation and default parameters (except eval num_episode was set to 100) used from this code repository: <https://github.com/weipu-zhang/STORM>

B ADDITIONAL BENCHMARK DETAILS

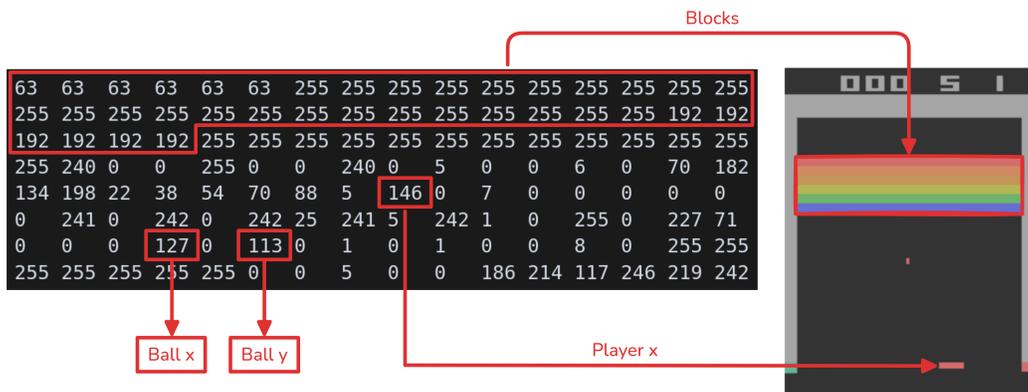


Figure 5: The figure shows the RAM state of Atari Breakout on the left and corresponding observation image from the emulator on the right, along with annotations for various state variables like ball position, blocks state etc.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

B.1 EXPERIMENT STOCHASTICITY MODES

B.1.1 STOCHASTICITY MODES USED IN BREAKOUT EXPERIMENTS

- Type 1: Random action executed from action space instead of predicted action with a probability of 0.3.
- Type 2.1: If a block is hit, there is probability of 0.25 that the hit is not considered and the block is not destroyed thereby returning 0 reward and the ball bounces back.
- Type 2.2: Episode starts with default setting (Type 3.1) and after 300 steps into the episode, the dynamics suddenly change to *Type 3.2A*.
- Type 3.1: Default Atari Breakout.
- Type 3.2A: The ball is only visible in a specific window between the blocks and the paddle and permanently hidden ($p = 1.0$) in rest of the space between them.
- Type 3.2B: Randomly hide left vertical half of the screen 75% ($p = 0.75$) of the episode.
- Type 3.2C: Only a random circular area of the screen is visible every frame ($p = 1.0$) similar to what someone will see when walking in a dark room with a torch.

B.1.2 STOCHASTICITY MODES USED IN BOXING EXPERIMENTS

- Type 1: Random action executed from action space instead of predicted action with a probability of 0.3.
- Type 2.1: Swaps the color of the enemy and player (character and score) with probability of 0.001 which results in 6-7 persistent swaps per episode (2 mins boxing round).
- Type 2.2: Episode starts with default setting (Type 3.1) and after 300 steps into the episode, the dynamics suddenly change to *Type 3.2C*.
- Type 3.1: Default Atari Boxing.
- Type 3.2A: Permanently hide ($p = 1.0$) scores and game clock.
- Type 3.2B: Randomly hide right vertical half of the screen 75% ($p = 0.75$) of the episode.
- Type 3.2C: Randomly hide enemy character 70% ($p = 0.7$) of the episode.

B.1.3 STOCHASTICITY MODES USED IN GOPHER EXPERIMENTS

- Type 1: Random action executed from action space instead of predicted action with a probability of 0.3.
- Type 2.1: Hole doesn't fill underground below the farmer and the reward is reverted to 0 whenever farmer digs, with probability of 0.3.
- Type 2.2: At the beginning of each episode, the environment is set to the default mode (Type 3.1). Every 600 steps, the dynamics transition *cyclically* between Type 3.2 and the default.
- Type 3.1: Default Atari Gopher.
- Type 3.2: Permanently hide ($p = 1.0$) underground gopher movement and holes and only hole openings are visible on surface (if any).

B.1.4 STOCHASTICITY MODES USED IN BANKHEIST EXPERIMENTS

- Type 1: With probability 0.3, a random action is executed from a restricted subset of the action space (0–9) instead of the predicted action. The restriction reduces the frequency of fire-based actions during random sampling, preventing the agent from instantly dying by triggering a bomb it drops on itself.
- Type 2.1: With probability 0.001, the robber is unexpectedly teleported to a different city.
- Type 2.2: At the beginning of each episode, the environment is set to the default mode (Type 3.1). Every 600 steps, the dynamics transition *cyclically* between Type 3.2 and the default.
- Type 3.1: Default Atari BankHeist.
- Type 3.2: Randomly hide city blocks 75% ($p = 0.75$) of the frames.

B.2 LEARNING CURVES FOR DIFFERENT STOCHASTICITY TYPES

Figures 6, 7, 8, and 9 illustrate the learning curves on Breakout, Boxing Gopher and BankHeist respectively, depicting the average evaluation return as a function of training steps up to 100K, for DreamerV3 and STORM.

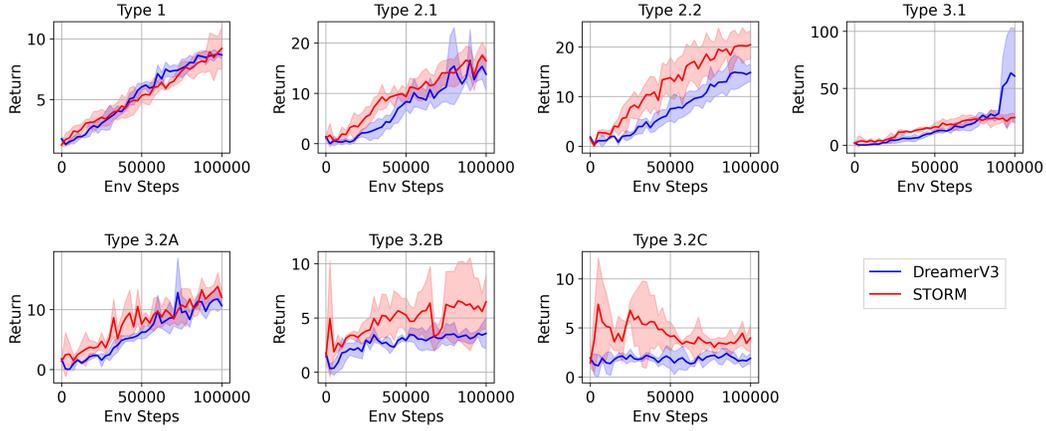


Figure 6: Breakout - learning curves

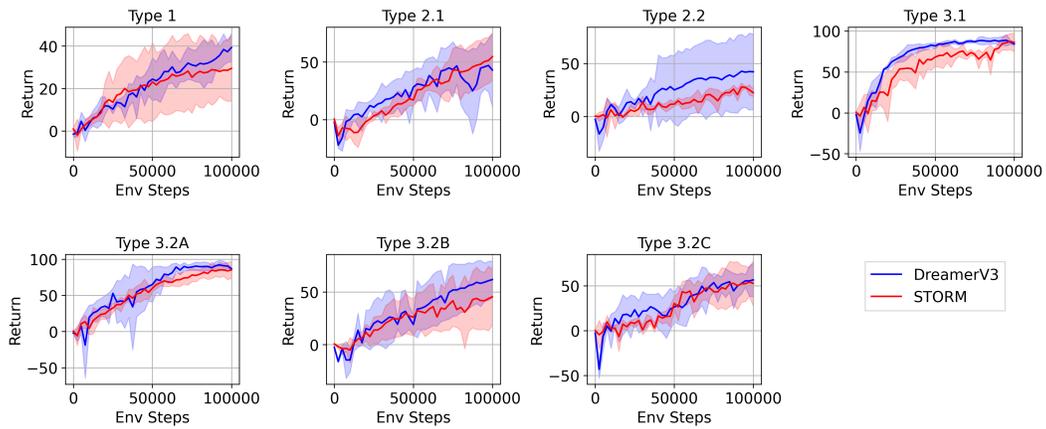


Figure 7: Boxing - learning curves

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

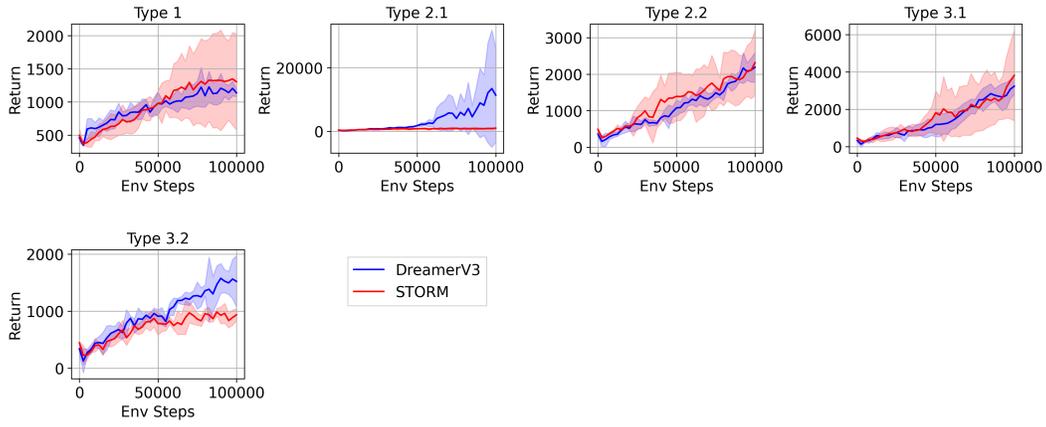


Figure 8: Gopher - learning curves

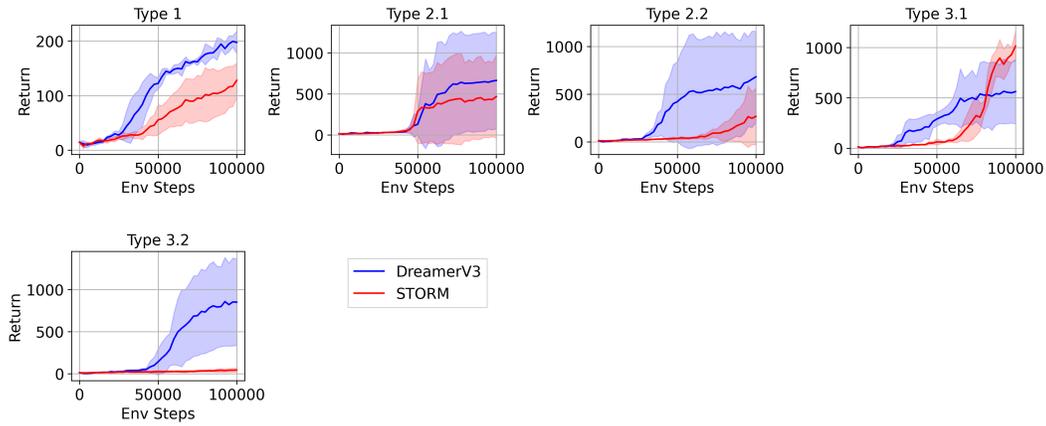


Figure 9: BankHeist - learning curves

B.3 OVERALL RESULTS FOR EVALUATION RETURN

Table 4: Overall Results for Evaluation Return

GAME NAME	STOCHASTICITY TYPE	DREAMERV3	STORM
Breakout	1	8.67 ± 0.30	9.20 ± 1.72
	2.1	13.80 ± 3.26	16.44 ± 2.03
	2.2	14.86 ± 1.74	20.45 ± 2.92
	3.1 (Default Baseline)	60.71 ± 41.89	24.17 ± 3.55
	3.2A	10.65 ± 0.41	12.05 ± 0.89
	3.2B	3.57 ± 1.43	6.48 ± 2.50
	3.2C	1.89 ± 0.49	3.96 ± 1.42
Boxing	1	39.32 ± 6.79	29.48 ± 15.41
	2.1	43.00 ± 31.98	54.69 ± 20.44
	2.2	42.21 ± 35.54	22.44 ± 3.54
	3.1 (Default Baseline)	84.22 ± 1.68	86.18 ± 11.29
	3.2A	86.90 ± 1.33	85.22 ± 11.77
	3.2B	61.74 ± 17.71	45.41 ± 26.56
	3.2C	56.52 ± 18.45	52.66 ± 25.68
Gopher	1	1137.00 ± 56.98	1303.53 ± 724.76
	2.1	11333.53 ± 14761.12	950.67 ± 188.37
	2.2	2190.87 ± 407.24	2315.40 ± 893.32
	3.1 (Default Baseline)	3235.27 ± 443.51	3811.67 ± 2431.85
	3.2	1521.13 ± 451.45	936.40 ± 106.83
BankHeist	1	197.63 ± 20.76	128.03 ± 31.41
	2.1	663.60 ± 587.85	467.80 ± 507.74
	2.2	682.67 ± 476.90	267.70 ± 295.86
	3.1 (Default Baseline)	562.30 ± 320.74	1015.73 ± 148.43
	3.2	849.80 ± 514.15	43.10 ± 22.85

1026 B.4 ADDITIONAL DETAILS: TYPE 2.1 ERROR ANALYSIS
10271028 We conducted a 1000-step probe in BANKHEIST under Type 2.1 stochasticity. During this probe,
1029 we executed a fixed action (action 3) for all 1000 steps and collected:

- 1030
1. the true environment frames, and
 - 1031 2. the corresponding world-model predictions.
-
- 1032

1033 For each sequence, we computed pixel-wise variances and compared:
1034

- 1035
- Type 3.1 (default, non-stochastic environment—baseline), and
 - 1036 • Type 2.1 (environment with stochastic teleportation events).
-
- 1037

1038 Across the 1000-step Type 2.1 run, four teleportation events occurred, producing large variance in
1039 the environment-state sequence.
10401041 DREAMERV3 RESULTS.
10421043 Env variance difference (2.1 – 3.1) \approx 299.10
1044 Model variance difference (2.1 – 3.1) \approx 0.0047
1045 Env–model variance gap (Type 2.1) \approx 300.34
10461047 STORM RESULTS.
10481049 Env variance difference (2.1 – 3.1) \approx 323.90
1050 Model variance difference (2.1 – 3.1) \approx 0.012
1051 Env–model variance gap (Type 2.1) \approx 325.21
10521053 **Interpretation.** The large environment variance difference between Type 2.1 and Type 3.1 high-
1054 lights that, under a fixed action sequence, the baseline environment (Type 3.1) exhibits extremely
1055 low variance, whereas the stochastic environment (Type 2.1) exhibits dramatically higher variance
1056 due to action-independent teleportation events. This demonstrates how sensitive the true environ-
1057 ment dynamics are to exogenous stochasticity.1058 In contrast, the world models fail to capture both the magnitude and the structure of this variance.
1059 Even when the true environment undergoes large, abrupt, action-independent deviations (e.g., tele-
1060 portation), the model predictions remain nearly unchanged from the deterministic baseline. Thus,
1061 while the environment distribution widens substantially under Type 2.1, the learned world mod-
1062 els continue to produce near-deterministic, low-variance prediction streams, indicating that they
1063 severely underestimate the true stochasticity.
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079

B.5 ADDITIONAL DETAILS: TYPE 3.2 ERROR ANALYSIS

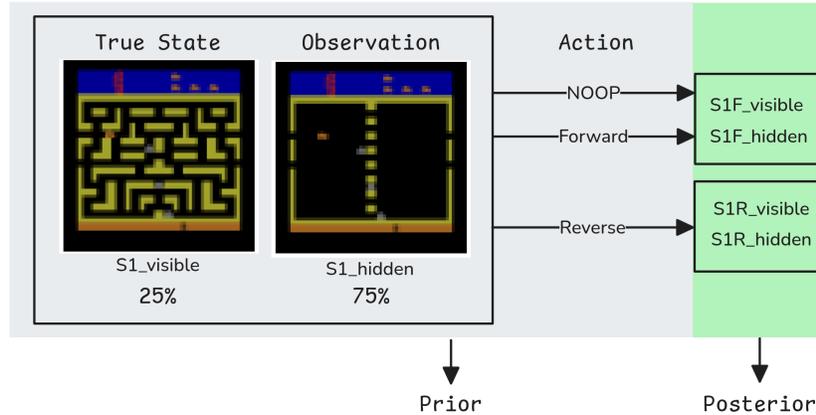


Figure 10: Detailed results from the BankHeist Type 3.2 partial observability probe. Each row corresponds to one of six carefully selected cases, showing the start state visibility, next state visibility, negative log-likelihood (NLL), and KL divergence for both DreamerV3 and STORM.

Table 5: Analysis on BankHeist (Type 3.2) - DreamerV3

Prior Obs	Metrics	Action — Posterior Observations (Visible/Hidden)					
		NOOP (0) S1F_visible	Forward (3) S1F_visible	Reverse (4) S1R_visible	NOOP (0) S1F_hidden	Forward (3) S1F_hidden	Reverse (4) S1R_hidden
S1_visible	- log d	41.72	38.43	44.73	19.13	21.56	29.81
	KL div	25.52	20.13	24.71	6.40	7.17	12.78
S1_hidden	- log d	41.36	38.74	45.14	25.15	22.63	29.25
	KL div	20.99	20.47	24.57	10.72	7.05	11.81

Table 6: Analysis on BankHeist (Type 3.2) - STORM

Prior Obs	Metrics	Action — Posterior Observations (Visible/Hidden)					
		NOOP (0) S1_visible	Forward (3) S1F_visible	Reverse (4) S1R_visible	NOOP (0) S1_hidden	Forward (3) S1F_hidden	Reverse (4) S1R_hidden
S1_visible	- log p	33.53	21.72	28.65	19.81	11.29	21.86
	KL div	115.63	116.00	114.13	114.94	115.47	113.79
S1_hidden	- log p	28.60	18.50	22.58	18.43	12.61	16.21
	KL div	115.55	115.49	115.05	115.03	115.19	114.73

1134 B.6 COMPUTE RESOURCES USED

1135
1136 The experiment runs were executed in several types of GPUs like A40, A100 and H100 depending
1137 on availability. Each node atleast had 32 vCPU and 50GB RAM. On GPUs with large memory,
1138 multiple runs were executed.

1139 DreamerV3 and STORM took around 24 hours and 12 hours respectively per run (training & evalu-
1140 ation) per seed when running on single GPU.

1141

1142 C INFORMATION-THEORETIC LEVERS

1143

1144 Information-theoretic measures provide quantitative levers to diagnose how stochasticity affects
1145 learning and planning:

1146

1147 **Action channel capacity.** For action-dependent noise, controllability is reduced. The effective
1148 capacity is measured by $I(A; \tilde{A} | S)$, which quantifies how much of the intended action A survives
1149 corruption into the executed action \tilde{A} .

1150

1151 **Predictive information of dynamics.** For action-independent randomness, the predictive struc-
1152 ture is measured by $I((S_t, A_t); S_{t+1})$, reflecting how much the next state depends on the current
1153 state-action pair. Under drift, temporal changes in this quantity indicate shifts in environment regu-
1154 larity.

1155

1156 **Representation sufficiency.** A latent Z_t should act as a sufficient statistic for planning. Ideally,
1157 $I(Z_t; S_t)$ is maximized, while $I(Z_t; O_t)$ remains bounded, ensuring that Z_t captures hidden states
1158 rather than surface-level noise, consistent with bisimulation invariance.

1159

1160 **Aliasing quantification.** In partially observable settings, the observation-state information gap
1161 can be written as $I(S_t; O_t) - I(S_t; O_t | A_t)$, capturing residual uncertainty after conditioning on
1162 actions. This disentangles sensor noise from genuine state ambiguity.

1163

1164 **Risk-sensitive planning.** Robust planning can be viewed through an information lens: risk-
1165 sensitive objectives such as Conditional Value at Risk (CVaR) optimize not the mean return but
1166 lower quantiles, effectively re-weighting information from rare but catastrophic outcomes.

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188 D ALL IMPLEMENTED STOCHASTICITY MODES

1189

1190 We define stochasticity modes along four Atari environments (Breakout, Boxing, Gopher,
1191 BankHeist), and the set of cropping modes are common to all games.

1192

1193 COMMON CROPPING MODES (ALL GAMES)

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

BREAKOUT

Action-independent random

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

BOXING

1237

1238

1239

1240

1241

Action-independent random

- 0: none
- 1: colorflip (swap player/enemy colors)
- 2: hit cancel (revert score; reward set to 0)

- 1242 • 3: displace to corners (swap player/enemy positions)
1243

1244 **Partial observation (blackout)**

- 1245 • 0: none
1246 • 1: all
1247 • 2: left boxing ring
1248 • 3: right boxing ring
1249 • 4: full boxing ring
1250 • 5: enemy score
1251 • 6: player score
1252 • 7: enemy+player score
1253 • 8: clock
1254 • 9: enemy+player score+clock
1255
1256
1257

1258 **Partial observation - RAM modification**

- 1259 • 0: none
1260 • 1: hide boxing ring
1261 • 2: hide enemy
1262 • 3: hide player
1263
1264
1265

1266 GOPHER

1267 **Action-independent random**

- 1268 • 0: none
1269 • 1: hole doesn't close (fill cancel; reward unchanged)
1270 • 2: hole doesn't close (fill cancel; reward set to 0)
1271 • 3: randomly remove one visible carrot (once per reset)
1272
1273
1274

1275 **Partial observation (blackout)**

- 1276 • 0: none
1277 • 1: all
1278 • 2: gopher attack (both sides)
1279 • 3: left gopher attack
1280 • 4: right gopher attack
1281 • 5: underground full (before-dug color)
1282 • 6: underground full offset (before-dug color)
1283 • 7: underground row 0 (before-dug)
1284 • 8: underground row 0 (dug color)
1285 • 9: underground row 1 (before-dug)
1286 • 10: underground row 1 (dug color)
1287 • 11: underground row 2 (before-dug)
1288 • 12: underground row 2 (dug color)
1289 • 13: underground row 3 (before-dug)
1290 • 14: underground row 3 (dug color)
1291 • 15: farmer (full)
1292
1293
1294
1295

- 1296 • 16: farmer below nose
- 1297 • 17: duck fly
- 1298 • 18: score
- 1299

1300 **Partial observation - RAM modification**

- 1301
- 1302 • 0: none
- 1303 • 1: hide left carrot
- 1304 • 2: hide middle carrot
- 1305 • 3: hide right carrot
- 1306 • 4: hide all carrots
- 1307 • 5: hide seed
- 1308
- 1309

1310 **BANKHEIST**

1311 **Action-independent random**

- 1312
- 1313 • 0: none
- 1314 • 1: dropped bomb is a dud
- 1315 • 2: fuel leaks (per city, once per episode)
- 1316 • 3: switch city mid-way (teleport)
- 1317 • 4: bank empty (reward suppressed when bank→police transition detected)
- 1318
- 1319

1320 **Partial observation (blackout)**

- 1321 • 0: none
- 1322 • 1: all
- 1323 • 2: city walls (all)
- 1324 • 3: top city wall
- 1325 • 4: left city wall
- 1326 • 5: bottom city wall
- 1327 • 6: right city wall
- 1328 • 7: left and right city walls together
- 1329 • 8: fuel region
- 1330 • 9: lives region
- 1331 • 10: score region
- 1332
- 1333
- 1334

1335 **Partial observation - RAM modification**

- 1336
- 1337 • 0: none
- 1338 • 1: hide robber's car
- 1339 • 2: hide change in fuel (always full)
- 1340 • 3: hide city blocks
- 1341 • 4: blend city blocks and wall (background color)
- 1342 • 5: hide banks (when currently a bank)
- 1343 • 6: hide police (when currently police)
- 1344
- 1345

1346 **CONCEPT DRIFT USAGE**

1347
1348 *All partial observation, action-independent, and action-dependent modes can also be used as a **second concept** in a concept drift setting, enabling controlled evaluation of robustness to non-stationary*
1349 *environments.*

1350 E THE USE OF LARGE LANGUAGE MODELS (LLMs)
1351

1352 We made use of large language models (LLMs) to assist with selected aspects of this work. Specifi-
1353 cally, LLMs were employed to improve the clarity and flow of writing, to summarize and condense
1354 long paragraphs during manuscript preparation, and to generate code snippets for repetitive compo-
1355 nents of the implementation. All outputs from the LLMs were carefully reviewed, validated, and
1356 edited by the authors to ensure accuracy and correctness.

1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403