

MORPHEUS: A Persistent Enterprise Benchmark for Continual Reinforcement Learning in the Big World

Anonymous authors

Paper under double-blind review

Abstract

We introduce MORPHEUS, a persistent enterprise simulation platform for continual reinforcement learning (CRL) research. Existing reinforcement learning (RL) benchmarks are episodic, low-dimensional, and stationary by design, properties antithetical to the challenges of real-world deployed systems. Grounded in the Big World Hypothesis, MORPHEUS provides environments in which the world never resets to an initial state, objectives shift over time, and past decisions have compounding consequences. MORPHEUS comprises four enterprise simulation environments; we evaluate two in this paper, drawn from outbound logistics and inbound warehouse operations, each exhibiting structured non-stationarity through a parameterisable failure injection engine and an asynchronous configuration shift controller. Policies are initialised via supervised fine-tuning on API-collected trajectories and subsequently trained with PPO-based reinforcement learning. Rewards are computed from operational verifiers embedded in the platform: structured failure event signals, financial ledger status, and resource throughput. We define a formal benchmark specification, propose a six-metric evaluation protocol covering per-configuration reward, adaptation speed, forgetting, recovery time, stability, and performance gap relative to a configuration-specific theoretical upper bound, and establish baseline results across four algorithm families: standard RL, replay-based, regularisation-based, and latent context modelling.

1 Introduction

Standard reinforcement learning (RL) benchmarks assume a Markov decision process (MDP) with a fixed transition function, a stationary reward signal, and an episode structure that resets the world to a known initial state distribution at the end of each episode. This paradigm has produced remarkable progress (Mnih et al., 2015; Schulman et al., 2017; Haarnoja et al., 2018) but studies a class of problems that excludes the most important challenges facing deployed intelligent systems.

Real decision-making systems operate continuously, without resets. Demand distributions shift, suppliers fail, regulatory requirements evolve, and the consequences of past decisions compound into current constraints. The Big World Hypothesis (Javed & Sutton, 2024) captures this precisely: intelligence learned in small, closed environments does not transfer to large, open, and evolving worlds. Here “large” and “evolving” mean that the world’s complexity exceeds what the agent can represent, making it appear non-stationary even when the underlying dynamics are fixed. Elelimy et al. (2025) formalise this argument and identify four foundations of traditional RL as antithetical to continual learning: the MDP formalism, a focus on optimal policies, expected cumulative reward as the evaluation metric, and episodic benchmarks.

In this work, our central question is: what would a benchmark look like that makes continual learning necessary? Drawing on these arguments, we identify three required properties. First, **persistence**: the world must not reset to an initial state; past decisions must compound into future dynamics. This follows directly from the Big World Hypothesis and from Elelimy et al. (2025)’s argument that continuing tasks are the appropriate setting for continual RL. Second, **non-stationarity**: the environment must change over time in ways that render any fixed policy eventually suboptimal. This follows from the Big World Hypothesis, where an environment larger than the agent will appear to change from the agent’s perspective. We additionally require that this non-stationarity be *parameterisable* as a design choice for benchmark reproducibility, so that different algorithms can be evaluated under identical conditions and results compared across papers. Third, **operational complexity**: the environment must be rich enough that no fixed optimal policy exists and adaptation requires genuine structural understanding. This follows from the Big World claim

40 that the environment exceeds the agent’s representational capacity. We ground this in enterprise operational processes
 41 where the action space, delayed consequences, and interacting domain systems collectively prevent any policy from
 42 remaining optimal across all conditions.

43 We present MORPHEUS, a persistent enterprise simulation platform satisfying all three requirements, with four envi-
 44 ronments (comprising of outbound order fulfilment and inbound warehouse receiving, among others). MORPHEUS
 45 provides four simulation environments in total and we evaluate two in this paper. Policies are initialised via supervised
 46 fine-tuning (SFT) on trajectories collected from an API-based LLM agent, then trained with PPO-based reinforcement
 47 learning. Configuration changes that implement regime shifts are triggered by an asynchronous controller running
 48 independently of the training loop.

49 We make the following contributions in this work.

- 50 1. We introduce MORPHEUS, a persistent enterprise simulation platform for continual RL research, based on the Big
 51 World Hypothesis and addressing the benchmark gaps identified by Elelimy et al. (2025) and Sutton et al. (2022).
- 52 2. We present a benchmark specification comprising two tasks over the `process-outbound` and
 53 `process-inbound` environments, stressing adaptation under structured drift and long-horizon credit assign-
 54 ment under delayed configuration shift effects.
- 55 3. We propose a six-metric evaluation protocol covering per-configuration reward, adaptation speed, forgetting, re-
 56 covery time, stability, and a performance gap relative to a configuration-specific theoretical upper bound, with
 57 relative adaptation advantage and plasticity tracking as supplementary diagnostics.
- 58 4. We establish baseline results across four algorithm families, providing an empirical study of algorithm behaviour
 59 under structured non-stationarity in persistent enterprise environments.

60 2 Background and Related Work

61 **Continual Reinforcement Learning.** Continual RL concerns agents that must learn and adapt over an indefinitely
 62 long horizon without treating each regime as an independent problem (Khetarpal et al., 2022). Core challenges include
 63 catastrophic forgetting (Kirkpatrick et al., 2017), loss of plasticity (Dohare et al., 2024), and regime detection without
 64 task boundary signals (Chandak et al., 2020). Abel et al. (2023) distinguish continual RL from multi-task RL and
 65 meta-learning by the absence of task labels and the requirement to handle non-stationarity arising from the world
 66 itself. Elelimy et al. (2025) argue that the MDP formalism fails in this setting and propose history processes and
 67 deviation regret as replacements. MORPHEUS is aligned with both: its environments are non-Markovian and we adopt
 68 a theoretically-grounded performance gap as a primary metric.

69 **The Big World Hypothesis and Continual RL.** The Big World Hypothesis (Javed & Sutton, 2024) asserts that
 70 the real world’s complexity exceeds any agent’s representational capacity ($C(\mathcal{E}) \gg C(\mathcal{A})$), causing the world to
 71 *appear* non-stationary even under fixed dynamics because the agent’s state representation omits causally relevant
 72 variables. We define a *regime* as a contiguous interval during which the dynamics, as observed through a given agent’s
 73 representation, are approximately stationary, and a *regime shift* as a transition inducing a meaningful change in the
 74 agent’s approximately optimal policy. Regime boundaries are therefore properties of the agent-environment pair, not
 75 the environment alone. In MORPHEUS, configuration changes serve as a reproducible experimental proxy; whether
 76 they constitute true regime shifts for a given agent is itself an empirical question. The Alberta Plan (Sutton et al., 2022)
 77 operationalises the Big World Hypothesis as a twelve-step programme, where Step 6 calls for a suite of continuing,
 78 non-episodic environments analogous to OpenAI Gym. MORPHEUS addresses this gap directly.

79 The Big World Hypothesis implies five CRL challenges MORPHEUS is designed to expose: **reward non-stationarity**
 80 (the reward mapping shifts across configurations); **catastrophic forgetting** (adapting to a new configuration risks
 81 overwriting prior policies, measurable via cyclic configuration sequences); **task-boundary-free detection** (no bound-
 82 ary signal is emitted and shifts must be inferred from experience); **non-stationary credit assignment** (value functions
 83 trained under one configuration produce biased targets after a shift, compounded by delayed outcomes); and **forward**
 84 **transfer** (shared domain structure across environments provides a falsifiable positive-transfer hypothesis).

85 **Existing Benchmarks.** ALE (Bellemare et al., 2013), MuJoCo (Todorov et al., 2012), and OpenAI Gym (Brock-
 86 man et al., 2016) are episodic by design. Procgen (Cobbe et al., 2020) introduces procedural variation but remains
 87 episodic. ContinualWorld (Wolczyk et al., 2021) uses explicit task boundaries and episodic resets. None tests the

88 core CRL property: detecting, adapting to, and retaining competence across unlabelled regime shifts in a world that
 89 never resets. AgarCL (Mohamed et al., 2025), based on Agar.io, provides emergent non-stationarity via multi-agent
 90 interactions well suited for open-ended robustness. MORPHEUS is complementary: it provides structured, repro-
 91 ducible configuration shifts enabling precise analysis of adaptation speed, forgetting, and recovery. Non-stationary
 92 MDP methods (Padakandla, 2021) largely preserve episodic structure, however, the MORPHEUS setting is harder be-
 93 cause regime shifts are unannounced and delayed reward produces biased temporal difference (TD) targets across
 94 configuration boundaries (Hung et al., 2019).

95 3 The MORPHEUS Platform

96 3.1 Architecture

97 Each environment in MORPHEUS is a self-contained *world plugin*: a TypeScript package exporting *Operational De-*
 98 *scriptors* (ODs), a simulation scheduler, seed data, and documentation. An OD defines the step-by-step execution
 99 plan for a capability (consisting of the sequence of available actions, tools called at each step, and data dependencies)
 100 and is the closest analogue to a transition structure in MORPHEUS. Agents interact via a capability API, where each
 101 call triggers an OD execution, with the failure injection engine (Section 3.3) introducing typed disruptions between
 102 steps. All results, failure events, state transitions, and ledger entries are logged to persistent, queryable collections.
 103 World instances are persistent (there are no episode resets). A world created at t_0 continues to evolve via the simula-
 104 tion scheduler independently of agent activity. An agent that fails to act faces real consequences: inventory depletes,
 105 orders go unfulfilled, and production schedules slip.

106 3.2 Simulation Environments

107 MORPHEUS provides four simulation environments, summarised in Table 1. The current study evaluates the
 108 two environments for which experiments have been completed. Evaluation on `inventory-optimizer` and
 109 `production-planning` is deferred to follow-on work.

Environment	World ID	Domains	Primary CRL Challenge
Inventory Optimisation	<code>inventory-optimizer</code>	WMS, ERP	Regime-dependent policy families
Process Inbound [†]	<code>process-inbound</code>	WMS, ERP, TMS	Sequential process accuracy
Process Outbound [†]	<code>process-outbound</code>	WMS, ERP, TMS	Delayed reward (OTIF)
Production Planning	<code>production-planning</code>	ERP, WMS, MES	Multi-step credit assignment

Table 1: Summary of MORPHEUS environments. Environments marked [†] are evaluated in this paper. WMS: Warehouse Management System. ERP: Enterprise Resource Planning. TMS: Transportation Management System. MES: Manufacturing Execution System.

110 Here, **Process Inbound** simulates the inbound logistics lifecycle from supplier arrival through quality control, put-
 111 away, and inventory reconciliation. Errors at receiving propagate through downstream steps with consequences only
 112 becoming observable later. To give a concrete example: a missing supplier advance shipping notice (ASN) (i.e., the
 113 electronic manifest accompanying an inbound shipment) forces the agent to proceed with a physical count alone, in-
 114 troducing discrepancies that surface as inventory record errors several steps later. The agent’s observation at any step
 115 does not reveal whether a prior discrepancy exists, making this a partially observable problem. **Process Outbound**
 116 simulates order fulfilment from creation through picking, dispatch, and payment settlement. The headline metric, On-
 117 Time In-Full (OTIF) delivery, is observable only at delivery, several simulated days after the dispatch decision. The
 118 agent must learn to use upstream process signals as proxies for this delayed outcome.

119 3.3 The Failure Injection Engine

120 Non-stationarity in MORPHEUS is implemented through a built-in failure injection engine that introduces typed
 121 disruptions into OD executions. Rather than injecting random noise, the engine draws from a structured taxon-
 122 omy of eleven failure types with specific operational semantics (e.g., `missing_data`, `dependency_failure`,
 123 `data_corruption`, `rate_limit`), making the non-stationarity interpretable and algorithm-designable. The en-
 124 gine operates at four preset failure rates: `light` (5%), `realistic` (8%), `moderate` (15%), and `aggressive`
 125 (30%). Each failure generates a structured *incident ticket*, where the reward penalises ticket generation with severity

126 weights (see Section 3.5 and Appendix A), directly linking operational failures to agent incentives. Crucially, different
127 presets constitute qualitatively different operating regimes, not merely different difficulty levels.

128 3.4 Configuration Shift Control

129 Non-stationarity between configurations is implemented by an **asynchronous configuration shift controller**, dis-
130 tinct from the simulation scheduler. The controller modifies the failure preset and demand parameters at pre-defined
131 simulation timestamps, running independently of the agent’s training loop so that configuration transitions are not
132 synchronised to gradient updates. This prevents the agent from using its own update periodicity as a proxy clock for
133 predicting shifts. Fixed configuration timestamps ensure all algorithms face identical shift points, prioritising repro-
134 ducibility over endogenous non-stationarity (where shifts would emerge from compounding operational decisions).
135 Experimental features for endogenous non-stationarity remains a compelling direction deferred to future work.

136 3.5 Reward Signal

137 Reward is computed from three categories of operational verifier logged natively by the platform, requiring no ex-
138 ternal annotation: (1) **failure event signals**, penalising incident ticket generation with severity weights that vary by
139 failure type; (2) **financial ledger status**, incorporating cost variance, payment timeliness, and freight efficiency; and
140 (3) **resource throughput**, measuring processed units and fulfilled orders relative to capacity. The composite reward
141 is a weighted sum $R(t) = w_f r_f(t) + w_l r_l(t) + w_p r_p(t)$ with default weights $w_f = 0.5$, $w_l = w_p = 0.25$. Full
142 formulation and theoretical upper bound derivation are in Appendix A.

143 3.6 Policy Initialisation and Training

144 The richness of the MORPHEUS action space makes pure RL from random initialisation impractical. We address
145 this with a two-stage pipeline. An API collection stage uses an LLM agent to produce observation-action-reasoning
146 trajectories across both environments and all four operating configurations. These trajectories then fine-tune a base
147 language model via imitation learning, producing a shared SFT checkpoint. All RL training runs are initialised from
148 this checkpoint, ensuring performance differences reflect continual learning behaviour rather than differences in basic
149 operational competence.

150 4 Benchmark Specification

151 We define two benchmark tasks, each stressing a distinct CRL failure mode. Both are continuing tasks with no episode
152 resets. Configuration shifts are applied by the asynchronous controller during training without any signal to the agent.

153 4.1 Task 1: Dynamic Resource Allocation Under Structured Drift

154 In `process-inbound`, an agent is tasked with allocating resources across incoming workloads w_t subject to ca-
155 pacity C_t , with configuration effects modulating job rate λ_t and capacity ceiling C_t . In `process-outbound`, an
156 agent distributes fulfilment capacity across K order categories subject to the capacity ceiling. Four configurations are
157 applied in sequence (Table 2). The **primary challenge** is the adaptation-stability trade-off, where fast adaptation to
158 bursty demand risks overwriting the low-demand policy, while excessive stability prevents timely response to change.

159 4.2 Task 2: Scheduling Under Drift with Delayed Effects

160 In `process-inbound` an agent is tasked with scheduling incoming shipments and putaway jobs across M receiving
161 lanes subject to regime-dependent arrival distributions and quality control constraints. In `process-outbound`, an
162 agent is tasked with scheduling outbound jobs across M processing lanes. Reward (financial component) is delayed:
163 OTIF is observable only at delivery, lagging dispatch decisions by several simulated days. Four configurations are
164 applied (Table 2). The **primary challenge** is long-horizon credit assignment under delayed configuration effects,
165 where the value functions calibrated for short-horizon reward produce systematically biased estimates under long-
166 horizon orders because the temporal structure of reward accumulation has changed. Together, Tasks 1 and 2 stress
167 complementary failure modes of continual RL within the same environment: adaptation stability and long-horizon
168 credit assignment respectively.

Config	Task 1	Task 2	Failure preset
1	Low demand	Short-horizon orders	light
2	Bursty demand	Long-horizon orders	realistic / moderate
3	Heavy SKU skew	Priority inversion	realistic
4	Constrained capacity	Long-horizon + disruption	aggressive / moderate

Table 2: Configuration sequences for Tasks 1 and 2.

169 **4.3 Extended Run: Recurrent Configurations**

170 Both tasks can be extended to include a return to the initial configuration ($A \rightarrow B \rightarrow C \rightarrow D \rightarrow A$). This enables
171 direct measurement of forgetting: performance on the return to configuration A is compared to the original encounter,
172 with any gap constituting the forgetting score.

173 **5 Evaluation Metrics**

174 We consider a continual RL setting (S, A, P_t, r_t) with latent configurations $\{z_t\}_{t=1}^T$. Cumulative reward $J(\pi) =$
175 $\mathbb{E}[\sum_t \gamma^{t-1} r_t(s_t, a_t)]$ alone is insufficient: a scalar sum over a non-stationary horizon conflates performance across
176 qualitatively different operating conditions. We propose a six-metric protocol, each targeting a distinct aspect of CRL.

177 Let t_k denote the onset of configuration k and \mathcal{T}_k its interval. R_k^{ub} is the *theoretical upper bound*: the maximum
178 expected reward achievable with full knowledge of the current configuration parameters, computed analytically (Ap-
179 pendix A). All metrics are computed post-hoc from experimental logs. An oracle agent receiving the true label z_t is
180 not included, with R_k^{ub} serving as a principled agent-independent reference.

181 **1. Per-Configuration Reward (all tasks).** Mean reward reported separately for each operating configuration:

$$J_k = \frac{1}{|\mathcal{T}_k|} \sum_{t \in \mathcal{T}_k} r_t(s_t, a_t) \quad (1)$$

182 Aggregate cumulative reward $R_{\text{total}} = \sum_{t=1}^T \gamma^{t-1} r_t(s_t, a_t)$ is reported only as a secondary summary for comparability
183 with prior work.

184 **2. Adaptation Speed (all tasks).** Our headline metric. Adaptation speed is the minimum number of steps from t_k
185 for the running average reward to reach αR_k^{ub} . We use $\alpha = 0.5$:

$$\tau_{\text{adapt}}^{(k)} = \min \left\{ \tau > 0 : \frac{1}{\tau} \sum_{j=0}^{\tau-1} r_{t_k+j} \geq \alpha R_k^{\text{ub}} \right\} \quad (2)$$

186 Using a running average rather than a single-step threshold makes the metric robust to reward noise. The threshold
187 αR_k^{ub} is now anchored to the theoretical maximum for each configuration, making it comparable across algorithms
188 and configurations. If the threshold is never reached within \mathcal{T}_k , $\tau_{\text{adapt}}^{(k)}$ is recorded as $|\mathcal{T}_k|$.

189 **3. Forgetting (extended recurrent runs only).** Applicable when a configuration recurs. Forgetting ($\bar{F}_k^{\text{policy}}$) is the
190 gap between interval mean rewards on first and second encounters, where negative values indicate positive backward
191 transfer.

192 **4. Recovery Time (all tasks).** The number of steps for the running average reward to stabilise within $\epsilon = 0.5 \cdot R_k^{\text{ub}}$
193 of the upper bound (Eq. 3). Recovery time is distinct from adaptation speed, where adaptation speed measures steps
194 to first cross the threshold, and recovery time measures steps to remain within ϵ of the upper bound.

$$\tau_{\text{rec}}^{(k)} = \min \left\{ \tau > 0 : \left| \frac{1}{\tau} \sum_{j=0}^{\tau-1} r_{t_k+j} - R_k^{\text{ub}} \right| \leq \epsilon \right\} \quad (3)$$

195 **5. Stability (all tasks).** Within-configuration reward variance σ_k^2 : high stability with low per-configuration reward
196 indicates a converged but suboptimal policy; high stability with high reward indicates genuine competence.

197 **6. Performance Gap (all tasks).** $\Delta_k^{(\pi)} = R_k^{\text{ub}} - \bar{r}_k^{(\pi)}$, where $\bar{r}_k^{(\pi)}$ is mean reward over the final 20% of \mathcal{T}_k . A large
 198 gap with low adaptation speed indicates failure to adapt; a large gap with high adaptation speed indicates convergence
 199 to a locally optimal but globally suboptimal policy.

200 **Relative Adaptation Advantage (supplementary).** $\text{RAA}_k^{(\pi)} = \tau_{\text{adapt,PPO}}^{(k)} - \tau_{\text{adapt}}^{(k),(\pi)}$ for $k > 1$. Positive values
 201 indicate the CL mechanism confers faster adaptation beyond incidental experience accumulation.

202 **Plasticity Tracking (diagnostic).** Effective (ER) of the policy network’s penultimate layer (Dohare et al., 2024);
 203 declining rank without performance recovery signals plasticity collapse.

204 6 Experimental Protocol and Algorithm Baselines

205 **Experimental Protocol:** Configuration sequences are fixed and identical across all algorithms and seeds. Five random
 206 seeds per algorithm per task is used, with results reported as mean \pm standard deviation. Upper bounds R_k^{ub} are
 207 computed analytically per configuration prior to experiments. Per-timestep logs include: reward r_t , component reward
 208 vector, configuration ID z_t , shift timestamps t_k , and effective rank of the penultimate layer.

209 **Baselines:** We evaluate four algorithm families, all initialised from the shared SFT checkpoint with PPO as the base
 210 optimiser. **(1) Standard RL (PPO)** (Schulman et al., 2017): no CL mechanism; establishes the failure baseline. **(2)**
 211 **Replay-Based (HER):** PPO with hindsight experience replay (Andrychowicz et al., 2017) retaining transitions across
 212 configuration boundaries; tests whether explicit retention prevents forgetting. **(3) Regularisation-Based (EWC):**
 213 PPO with Elastic Weight Consolidation (Kirkpatrick et al., 2017); tests the stability side of the plasticity-stability
 214 trade-off. **(4) Latent Context Modelling (LCM):** PPO with a latent context model (Rakelly et al., 2019) that infers
 215 the configuration from experience history; tests whether implicit detection enables fast adaptation.

216 7 Experimental Results

217 7.1 Process-Outbound Environment

218 **Task 1: Dynamic Resource Allocation Under Structured Drift.** Table 3 reports Task 1 results. EWC obtains the
 219 highest per-configuration reward ($\bar{J}_k = 0.00239$) and cumulative reward (0.5741), while LCM achieves the fastest
 220 adaptation ($\bar{\tau}^{\text{adapt}} = 21.8$ steps vs. 32.5 for EWC and 45.2 for PPO and HER), suggesting latent context inference is
 221 more effective at detecting structured drift. EWC and LCM both recover within 1.0 step on average; PPO and HER
 222 require 12.3 and 14.3 steps respectively. Task 1 thus exhibits a split: EWC is the best reward maximiser; LCM the
 223 strongest adaptor.

224 **Task 2: Scheduling Under Drift with Delayed Effects.** Table 4 reports Task 2 results for `process-outbound`.
 225 The winner pattern differs from Task 1. HER achieves the highest cumulative reward (0.5848) and per-configuration
 226 reward ($\bar{J}_k = 0.00127$), narrowly ahead of EWC (0.5753). Adaptation speed is slower across all methods relative to
 227 Task 1, with EWC achieving the fastest adaptation ($\bar{\tau}_{\text{adapt}} = 40.5$ steps) and a positive RAA of 9.3, while LCM is the
 228 slowest adaptor ($\bar{\tau}_{\text{adapt}} = 51.1$, $\text{RAA} = -1.4$), suggesting that latent context inference is less effective under delayed
 229 reward feedback. Recovery is fast across all methods (1.0–1.2 steps). Effective rank is broadly stable across families,
 230 with EWC maintaining the highest rank (36.80). Forgetting scores are small and similar across methods, ranging from
 231 0.00330 (PPO) to 0.00370 (HER). Task 2 in `process-outbound` therefore does not show the clean EWC/LCM
 232 split observed in Task 1: HER leads on reward, EWC leads on adaptation, and LCM offers no adaptation advantage
 233 under the delayed credit-assignment structure of this task.

234 7.2 Process-Inbound environment

235 **Task 1: Dynamic Resource Allocation Under Structured Drift.** Table 5 reports Task 1 results. HER obtains the
 236 highest cumulative reward (0.2474) and per-configuration reward ($\bar{J}_k = 0.00082$), while also achieving the highest
 237 effective rank (37.74). PPO is the second-best reward maximiser, with cumulative reward 0.2366 and $\bar{J}_k = 0.00079$.
 238 EWC and LCM obtain similar cumulative rewards (0.2278 and 0.2275 respectively). Adaptation speed saturates at
 239 60.0 steps for all methods, and recovery time is 1.0 step for all methods due to near-constant low per-configuration
 240 reward across regimes; consequently, RAA is 0.0 for HER, EWC, and LCM. Forgetting is small across all methods,
 241 with LCM obtaining the lowest forgetting score (0.00366). Task 1 therefore shows a clear reward and representation
 242 advantage for HER, while adaptation and recovery metrics do not distinguish the algorithm families.

Family	Cum. Reward	\bar{J}_k	$\bar{\tau}^{\text{adapt}}$	$\bar{\tau}^{\text{rec}}$	$\bar{\sigma}^2$	$\bar{\Delta}_k$	RAA	ER	$\bar{F}_k^{\text{policy}}$
PPO	0.5040±0.0213	0.00210±0.00007	45.2±9.5	12.3±0.0	0.000525 ±0.00002	0.9931±0.0005	—	35.77±1.84	0.00302±0.00038
HER	0.5347±0.0486	0.00223±0.00021	45.2±11.2	14.3±0.0	0.000597±0.00002	0.9931±0.0020	0.0±9.3	36.30 ±0.86	0.00328±0.00039
EWC	0.5741 ±0.0278	0.00239 ±0.00007	32.5±9.2	1.0 ±0.0	0.000623±0.00002	0.9927±0.0010	18.3±12.0	35.52±1.87	0.00334±0.00045
LCM	0.5210±0.0207	0.00217±0.00003	21.8 ±9.6	1.0 ±0.0	0.000573±0.00002	0.9929±0.0000	32.3 ±15.4	35.65±1.59	0.00359±0.00045

Table 3: Task 1 regime-aggregated metrics for the Process-Outbound Environment. Cumulative reward is reported as undiscounted return. Adaptation and recovery times use data-relative thresholds. **RAA** is averaged over configurations $k > 1$ relative to PPO. $\bar{F}_k^{\text{policy}}$ is policy-attributable forgetting (terminal bonus excluded). **ER** is mean effective rank of the penultimate-layer activations. Lower is better for $\bar{\tau}^{\text{adapt}}$, $\bar{\tau}^{\text{rec}}$, $\bar{\sigma}^2$, $\bar{\Delta}_k$, and $\bar{F}_k^{\text{policy}}$; higher is better for **cumulative reward**, \bar{J}_k , **RAA**, and **ER**.

Family	Cum. Reward	\bar{J}_k	$\bar{\tau}^{\text{adapt}}$	$\bar{\tau}^{\text{rec}}$	$\bar{\sigma}^2$	$\bar{\Delta}_k$	RAA	ER	$\bar{F}_k^{\text{policy}}$
PPO	0.5454±0.0093	0.00117±0.00002	49.7±14.5	1.0±0.0	0.000300±0.00002	0.9996±0.0009	—	36.04±1.82	0.00330±0.00038
HER	0.5848 ±0.0174	0.00127 ±0.00006	35.2±19.0	1.2±0.4	0.000300±0.00002	0.9990±0.0014	14.5 ±20.8	35.76±1.99	0.00370±0.00035
EWC	0.5753±0.0091	0.00123±0.00003	40.5 ±9.3	1.0 ±0.0	0.000290 ±0.00002	0.9991±0.0009	9.3±10.6	36.80 ±1.74	0.00365±0.00040
LCM	0.5438±0.0179	0.00114±0.00005	51.1±14.1	1.0±0.0	0.000300±0.00002	0.9999 ±0.0001	-1.4±10.1	35.83±1.07	0.00347±0.00023

Table 4: Task 2 regime-aggregated metrics for Process-Outbound Environment.

243 **Task 2: Scheduling Under Drift with Delayed Effects.** Table 6 reports Task 2 results. HER again obtains the highest
244 cumulative reward (0.2499) and per-configuration reward ($\bar{J}_k = 0.00083$), with LCM close behind on reward (0.2440,
245 $\bar{J}_k = 0.00081$). LCM is the clear adaptation leader, achieving the fastest adaptation speed ($\bar{\tau}_{\text{adapt}} = 51.0$) and the only
246 positive RAA (+1.56). HER has the highest effective rank (36.87), suggesting stronger representation preservation,
247 but adapts more slowly than PPO and LCM. EWC performs worst on cumulative reward (0.1861) and adaptation
248 speed (59.0), although it has the lowest forgetting score (0.00315). Task 2 thus separates reward maximisation from
249 adaptation: HER is the strongest reward maximiser, while LCM is the strongest adaptor.

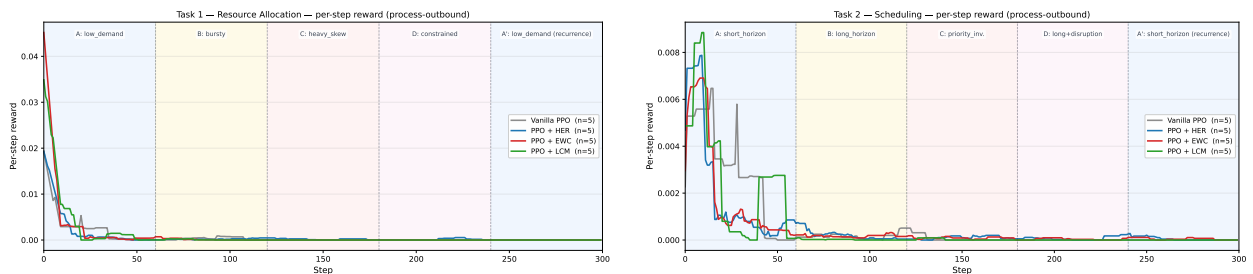
250 8 Discussion

251 We make the following observations from the performance of the algorithms.

252 8.1 Process-Outbound Environment

253 **All algorithms remain far below the theoretical upper bound.** Mean performance gaps $\bar{\Delta}_k$ sit near 1.0 for every
254 method on both tasks, indicating a large settled-state deficit rather than a minor tuning gap. This is not attributable to
255 isolated episodes or reset artefacts. MORPHEUS worlds are persistent with hidden shifts and delayed consequences, so
256 the gap reflects genuine difficulty maintaining competence as operational conditions drift.

257 Figure 1 shows the per-configuration reward of each task from one seed run in process-outbound. We find that the bulk
258 of variance in each algorithm family occurs on the first regime (regime A), when there is both a low amount of tickets
259 (and therefore, a higher chance of selecting the most impactful tasks to execute), with easier reward-wins via ticket
260 completion in the beginning of ticket completion compared to the tail end. During the more difficult regimes, this



(a) Task 1: dynamic resource allocation.

(b) Task 2: delayed-effect scheduling.

Figure 1: Per-configuration (or Per-step) reward over both tasks across circular regime shifts in process-outbound.

Family	Cumulative Reward	\bar{J}_k	$\bar{\tau}_{adapt}$	$\bar{\tau}_{rec}$	RAA	ER	\bar{F}_k^{policy}
PPO	0.2366 ± 0.0176	0.00079	60.0 ± 0.0	1.0 ± 0.0	–	35.06 ± 0.71	0.00397 ± 0.00027
HER	0.2474 ± 0.0302	0.00082	60.0 ± 0.0	1.0 ± 0.0	0.0	37.74 ± 1.02	0.00401 ± 0.00047
EWC	0.2278 ± 0.0125	0.00076	60.0 ± 0.0	1.0 ± 0.0	0.0	35.47 ± 0.31	0.00400 ± 0.00035
LCM	0.2275 ± 0.0215	0.00076	60.0 ± 0.0	1.0 ± 0.0	0.0	35.81 ± 0.39	0.00366 ± 0.00039

Table 5: Task 1 regime-aggregated metrics for the Process-Inbound Environment. Stability and performance gap are omitted as near-constant low per-configuration reward renders both metrics uninformative in `process-inbound`.

Family	Cumulative Reward	\bar{J}_k	$\bar{\tau}_{adapt}$	$\bar{\tau}_{rec}$	RAA	ER	\bar{F}_k^{policy}
PPO	0.2376 ± 0.0465	0.00079	52.2 ± 13.6	1.0 ± 0.0	–	34.88 ± 1.06	0.00399 ± 0.00085
HER	0.2499 ± 0.0191	0.00083	58.6 ± 2.5	1.0 ± 0.0	-8.56 ± 14.82	36.87 ± 1.07	0.00419 ± 0.00038
EWC	0.1861 ± 0.0807	0.00062	59.0 ± 1.7	1.0 ± 0.0	-9.11 ± 15.78	34.94 ± 0.50	0.00315 ± 0.00162
LCM	0.2440 ± 0.0348	0.00081	51.0 ± 7.8	1.0 ± 0.0	$+1.56 \pm 15.72$	35.58 ± 0.84	0.00396 ± 0.00054

Table 6: Task 2 regime-aggregated metrics for the Process-Inbound Environment.

261 ability falters, both along securing completion for previous tasks, as well as in selecting the most reward-maximizing
 262 tickets in a larger, more difficult set, given the environment changes up until that point. When returning to regime A at
 263 the end, the agent is not able to return to previous performance, hindered by the workload of the intermediate regimes.

264 **No single algorithm family dominates.** EWC leads on cumulative reward and per-configuration reward on Task 1,
 265 while HER leads on Task 2. LCM leads on adaptation speeds in Task 1 and lags behind in Task 2. This divergence
 266 validates the six-metric protocol: aggregate reward alone would crown EWC on Task 1 and HER on Task 2, obscuring
 267 the cross-task role reversal. The six metrics make visible that reward accumulation and adaptation speed identify
 268 different winners, and that the identity of each winner changes with the task’s CRL structure.

269 **Regime detection is hard even in controlled settings.** PPO and HER generally adapt only in the first configuration
 270 and fail to adapt in later regimes, even in a controlled four-regime sequence without label signals. Per-configuration
 271 reward decays toward zero despite continued experience, indicating agents cannot capture the regime changes.

272 8.2 Process-Inbound Environment

273 We make the following observations in addition to those provided for Process-Outbound.

274 **Forgetting remains small across all methods.** LCM has the lowest forgetting score on Task 1, while EWC has the
 275 lowest forgetting score on Task 2. However, the forgetting values are close in magnitude and should be read together
 276 with reward and effective rank. A low forgetting score without high reward does not imply strong continual-learning
 277 performance. It may instead be an artifact of uniformly low performance across the initial and recurrent configurations.

278 **The inbound results show no evidence of plasticity collapse.** Effective rank remains stable across all families, with
 279 HER maintaining the highest rank in both tasks. The main difficulty is therefore not representational collapse, but
 280 the ability to convert persistent representations into fast and high-reward adaptation under hidden configuration shifts.
 281 Overall, `process-inbound` favours HER for reward maximisation and representation preservation, while Task 2
 282 reveals a distinct adaptation advantage for LCM.

283 9 Conclusion

284 We presented MORPHEUS, a persistent enterprise simulation platform that makes continual RL necessary. Grounded
 285 in the Big World Hypothesis, MORPHEUS provides environments that are persistent, non-stationary by design, and
 286 operationally complex enough that no fixed policy remains optimal. We defined a two-task benchmark, a six-metric
 287 evaluation protocol anchored to configuration-specific theoretical upper bounds, and four algorithm baselines covering
 288 the major CRL paradigms. Baseline results reveal large settled-state deficits and task-dependent algorithm specialisa-
 289 tion. This demonstrates that no existing approach solves the benchmark and that multi-metric evaluation is necessary
 290 to characterise CRL behaviour. MORPHEUS will be released as an open-source infrastructure for researchers and
 291 practitioners. Limitations and future work are discussed in Appendix B.

292 **References**

- 293 David Abel, André Barreto, Benjamin Van Roy, Doina Precup, Hado P van Hasselt, and Satinder Singh. A Definition
294 of Continual Reinforcement Learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- 295 Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh
296 Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight Experience Replay. In *Advances in Neural Information
297 Processing Systems (NeurIPS)*, 2017.
- 298 Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The Arcade Learning Environment: An
299 Evaluation Platform for General Agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- 300 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech
301 Zaremba. OpenAI Gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 302 Yash Chandak, Georgios Theodorou, Shiv Shankar, Martha White, Sridhar Mahadevan, and Philip S. Thomas.
303 Optimizing for the Future in Non-Stationary MDPs. In *International Conference on Machine Learning (ICML)*,
304 2020.
- 305 Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging Procedural Generation to Benchmark
306 Reinforcement Learning. In *International Conference on Machine Learning (ICML)*, 2020.
- 307 Shibhansh Dohare, J. Fernando Hernandez-Garcia, Qingfeng Lan, Parash Rahman, A. Rupam Mahmood, and
308 Richard S. Sutton. Loss of Plasticity in Deep Continual Learning. *Nature*, 632(8026):768–774, 2024.
- 309 Esraa Elelimy, David Szepesvari, Martha White, and Michael Bowling. Rethinking the Foundations for Continual
310 Reinforcement Learning. *arXiv preprint arXiv:2504.08161*, 2025.
- 311 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft Actor-Critic: Off-Policy Maximum Entropy
312 Deep Reinforcement Learning with a Stochastic Actor. In *International Conference on Machine Learning (ICML)*,
313 2018.
- 314 Chia-Chun Hung, Timothy Lillicrap, Josh Abramson, Yan Wu, Mehdi Mirza, Federico Carnevale, Arun Ahuja, and
315 Greg Wayne. Optimizing Agent Behavior over Long Time Scales by Transporting Value. *Nature communications*,
316 10(1):5223, 2019.
- 317 Khurram Javed and Richard S. Sutton. The Big World Hypothesis and its Ramifications for Artificial Intelligence. In
318 *Finding the Frame: An RLC Workshop for Examining Conceptual Frameworks*, 2024.
- 319 Khimya Khetarpal, Matthew Riemer, Irina Rish, and Doina Precup. Towards Continual Reinforcement Learning: A
320 Review and Perspectives. *Journal of Artificial Intelligence Research*, 75:1401–1476, 2022.
- 321 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran
322 Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in
323 neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, 2017.
- 324 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves,
325 Martin Riedmiller, Andreas K. Fiedjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement
326 learning. *Nature*, 518(7540):529–533, 2015.
- 327 Mohamed A Mohamed, Kateryna Nekhomiaz, Vedant Vyas, Marcos M. Jose, Andrew Patterson, and Marlos C.
328 Machado. The Cell Must Go On: Agar.io for Continual Reinforcement Learning. *arXiv preprint arXiv:2505.18347*,
329 2025.
- 330 Sindhu Padakandla. A Survey of Reinforcement Learning Algorithms for Dynamically Varying Environments. *ACM
331 Computing Surveys (CSUR)*, 54(6):1–25, 2021.
- 332 Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient Off-Policy Meta-
333 Reinforcement Learning via Probabilistic Context Variables. In *International Conference on Machine Learning
334 (ICML)*, 2019.

- 335 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization
336 Algorithms. In *arXiv preprint arXiv:1707.06347*, 2017.
- 337 Richard S. Sutton, Michael Bowling, and Patrick M. Pilarski. The Alberta Plan for AI Research. *arXiv preprint*
338 *arXiv:2208.11173*, 2022.
- 339 Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A Physics Engine for Model-Based Control. In *2012*
340 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012.
- 341 Maciej Wołczyk, Michał Zając, Razvan Pascanu, Łukasz Kuciński, and Piotr Miłoś. Continual World: A Robotic
342 Benchmark For Continual Reinforcement Learning. In *Advances in Neural Information Processing Systems*
343 (*NeurIPS*), 2021.

Supplementary Materials

The following content was not necessarily subject to peer review.

344
345
346

347 A Reward Formulation and Theoretical Upper Bound Derivation

348 A.1 Composite Reward

349 The composite reward is:

$$R(t) = w_f \cdot r_f(t) + w_l \cdot r_l(t) + w_p \cdot r_p(t) \quad (4)$$

350 where r_f is the failure event signal, r_l is the financial ledger signal, and r_p is the resource throughput signal.

351 **Failure event signal.** Let \mathcal{T}_t denote the set of incident tickets generated at timestep t . Each ticket carries a severity
352 weight $s(\tau)$ depending on its failure type:

$$r_f(t) = - \sum_{\tau \in \mathcal{T}_t} s(\tau) \quad (5)$$

353 Severity weights: permission_denied, dependency_failure: $s = 1.0$; data_corruption,
354 missing_data, invalid_state: $s = 0.6$ – 0.8 ; rate_limit, format_change, partial_data: $s =$
355 0.3 – 0.4 ; stale_data, duplicate_data, timing_issue: $s = 0.1$ – 0.2 .

Financial ledger signal.

$$r_l(t) = 1 - \frac{\text{actual_cost}(t)}{\text{planned_cost}(t)} \quad (6)$$

356 clipped to $[-1, 1]$, incorporating cost variance, payment timeliness, and freight efficiency as logged by the platform
357 ledger.

Resource throughput signal.

$$r_p(t) = \frac{\text{units_processed}(t)}{\text{capacity}(t)} \quad (7)$$

358 clipped to $[0, 1]$.

359 Default weights: $w_f = 0.5$, $w_l = 0.25$, $w_p = 0.25$. These are research variables; researchers are encouraged to treat
360 them as such.

361 A.2 Theoretical Upper Bound Derivation

362 The theoretical upper bound R_k^{ub} for configuration k is computed as the maximum expected reward achievable by a
363 policy with perfect knowledge of the configuration parameters, under the assumption of zero failures and maximum
364 resource utilisation:

$$R_k^{\text{ub}} = w_f \cdot 0 + w_l \cdot r_l^{\text{max}} + w_p \cdot 1 \quad (8)$$

365 The failure signal contributes zero under the upper bound assumption (no tickets generated). The financial signal r_l^{max}
366 is derived from the minimum achievable cost ratio under the configuration’s demand parameters (planned cost equal
367 to actual cost, giving $r_l^{\text{max}} = 1$). The throughput signal is clipped at 1 by definition. This gives:

$$R_k^{\text{ub}} = w_l + w_p = 0.25 + 0.25 = 0.50 \quad (9)$$

368 under the default weight setting. The upper bound is therefore identical across configurations under the default weights,
369 with configuration difficulty reflected in how far agents fall below it. Researchers using alternative weight vectors
370 should recompute R_k^{ub} accordingly. When $w_l \cdot r_l^{\text{max}} < 1$ due to structural cost constraints in a given configuration
371 (e.g., freight costs that cannot be reduced below planned cost), the configuration-specific r_l^{max} should be used instead
372 of 1.

373 **B Limitations and Future Work**

374 There are four important limitations of this work. 1) **Environment scope.** Evaluation covers two of four environments,
375 with experiments on `inventory-optimizer` and `production-planning` deferred to future work. 2) **Upper**
376 **bound tightness.** R_k^{ub} assumes zero failures and perfect utilisation. In high-failure configurations it is an optimistic
377 ceiling rather than an attainable target. 3) **Endogenous non-stationarity.** The controller implements externally trig-
378 gered shifts. A study of endogenous regime changes driven by compounding decisions are deferred to future work.
379 4) **Reward grounding.** Component weights are research variables, not validated industry objectives. Future work
380 will evaluate the full suite, add an oracle baseline, establish a delayed-reward correction for zero-shot generalisation
381 in `process-outbound`, and explore endogenous non-stationarity.