

REINFORCEMENT LEARNING VIA LAZY-AGENT FOR ENVIRONMENTS WITH RANDOM DELAYS

Anonymous authors

Paper under double-blind review

ABSTRACT

Real-world reinforcement learning applications are often hampered by delayed feedback from environments, which violates the fundamental assumption of the Markovian property and introduces significant challenges. While numerous methods have been proposed for handling environments with constant delays, those with random delays remain largely unexplored owing to their inherent complexity and variability. **In this study, we explored environments with random delays and proposed a novel strategy to transform them into their equivalent constant-delay counterparts by introducing a simple agent called the *lazy-agent*. This approach naturally overcomes the challenges posed by the variability of random delays, enabling the application of state-of-the-art methods, originally designed for handling constant delays, to random-delay environments without any modification. Empirical results demonstrate that the lazy-agent-based algorithm significantly outperformed other baselines in terms of asymptotic performance and sample efficiency in random-delay environments.**

1 INTRODUCTION

Reinforcement learning (RL) has made remarkable progress in various domains, from gaming (Mnih et al., 2013; Silver et al., 2016) to robotic control systems (Haarnoja et al., 2018; Kalashnikov et al., 2018). However, real-world applications of RL often face challenges due to delays, which can take diverse forms such as latency in communication systems, delays in processing sensory data, or response delays from actuators. These delays can significantly degrade the performance of RL agents and may even cause instability in dynamic systems (Hwangbo et al., 2017; Mahmood et al., 2018).

While numerous methods have been proposed to address the challenges posed by delays within the RL framework, these efforts primarily focus on the unrealistic assumption of constant delays (Chen et al., 2021; Derman et al., 2021; Liotet et al., 2022; Kim et al., 2023; Wu et al., 2024), leaving random delays relatively unexplored owing to their inherent complexity and variability. However, in real-world, randomly varying delays present a more realistic challenge, exemplified by communication systems where diverse routing paths and the physical properties of the network can result in asynchronous data arrivals (Ge et al., 2013).

In this study, we explore environments with random delays and establish a connection to environments with constant delays by introducing a simple agent called the *lazy-agent*. **Specifically, we demonstrate that random-delay environments can be straightforwardly transformed into their equivalent constant-delay counterparts using lazy-agents, enabling state-of-the-art constant-delay approaches to be seamlessly applied to random-delay environments without any modification.** We train lazy-agents within the belief projection-based Q -learning (BPQL) framework (Kim et al., 2023), termed *lazy-BPQL*, to leverage its advantages in training agents in delayed environments. The efficacy of the proposed lazy-BPQL was evaluated on popular continuous control tasks in the MuJoCo benchmark (Todorov et al., 2012). Empirical results demonstrate that lazy-BPQL outperformed other baseline algorithms in terms of asymptotic performance and sample efficiency in random-delay environments, achieving performance comparable to agents trained in constant-delay environments.

2 BACKGROUNDS

2.1 STANDARD REINFORCEMENT LEARNING

A (*delay-free*) Markov decision process (MDP) (Bellman, 1957) can be defined with a five-tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$, where \mathcal{S} is the state space, and \mathcal{A} is the action space, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the transition kernel, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\gamma \in (0, 1)$ is a discount factor. Additionally, the policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ maps the state-to-action distribution.

Under this definition, at each discrete time t , an RL agent observes state s_t , makes a decision a_t based on a policy π , receives a reward r_t with respect to the action taken, and then observes the next state s_{t+1} from the environment. It repeats this process to find an optimal policy π^* that maximizes the expected discounted cumulative rewards, which is given as:

$$\pi^* := \arg \max_{\pi} \mathbb{E} \left[\sum_{k=0}^{H-1} \gamma^k r_{t+k+1} | \pi, \rho_0 \right] = \arg \max_{\pi} \mathbb{E} [G_0 | \pi, \rho_0], \quad (1)$$

where ρ_0 denotes the initial state distribution and G_0 is the discounted cumulative rewards starting from the initial state over a finite or infinite-horizon H under the policy π . Additionally, the values of states and actions at time t are defined as:

$$V^{\pi}(s) = \mathbb{E} \left[\sum_{k=0}^{H-1} \gamma^k r_{t+k+1} | \mathcal{S}_t = s, \pi \right], \quad Q^{\pi}(s, a) = \mathbb{E} \left[\sum_{k=0}^{H-1} \gamma^k r_{t+k+1} | \mathcal{S}_t = s, \mathcal{A}_t = a, \pi \right], \quad (2)$$

where $V^{\pi}(s)$ denotes the expected discounted cumulative rewards starting from state s under the policy π , and $Q^{\pi}(s, a)$ represents the expected discounted cumulative rewards starting from state s , taking action a , and then following the policy π .

Note that the dynamics governing MDPs assume the Markovian property, which indicates that the complete probability distribution in the dynamics can be fully determined by the present state and action, independent of their histories. However, this fundamental assumption can be violated by delayed feedback from the environment, leading to partially observable MDPs (Monahan, 1982), where the agent’s current state and action cannot capture sufficient information needed for timely decision-making. This can significantly degrade the performance of RL agents or even lead to complete failure in learning (Singh et al., 1994).

2.2 DELAYED REINFORCEMENT LEARNING

In MDP with delays, referred to as *delayed* MDP, the state of the environment may not be observed by the agent immediately (observation delay). The effect of the action applied to the environment may also be delayed (action delay). Additionally, the reward generated by the action taken may not reach the agent immediately (reward delay). These delays force the agent to make decisions based on outdated information, prevent timely and appropriate actions, or cause the agent to receive rewards that do not correspond to the actions taken, disrupting the learning process.

Delayed MDPs are typically categorized into *constant-delay* MDPs (CDMDPs), where feedback is delayed by a fixed number of time-steps; and *random-delay* MDPs (RDMDPs), where the number of delayed time-steps varies randomly. For example, in the case of random observation delay, the state s_t may be delayed by four time-steps and observed by the agent at time $t + 4$, whereas the next state s_{t+1} may be delayed by only one time-step and observed at time $t + 2$. Note that the subscript t explicitly indicates the times when states are generated by the actions applied to the environment.

In this study, we focus on randomly delayed observations under the assumption that the agent utilizes them for decision-making *in order*. This implies that any observed state can be used by the agent for decision-making only after all previously generated states have been both observed and utilized to ensure no state is omitted from the decision-making process. In the presence of randomness in observation delays, multiple states may become observable simultaneously, and their order may even be scrambled. To reduce confusion about the timing of state observations, we distinguish between a state being *observed* and *used*, under the aforementioned assumption of *ordering*, as follows:

Definition 2.1. A state is considered *observed* when the information about the state of the environment reaches the agent. A state is considered *used* when the agent utilizes the observed state information to make a decision by feeding it into the policy.

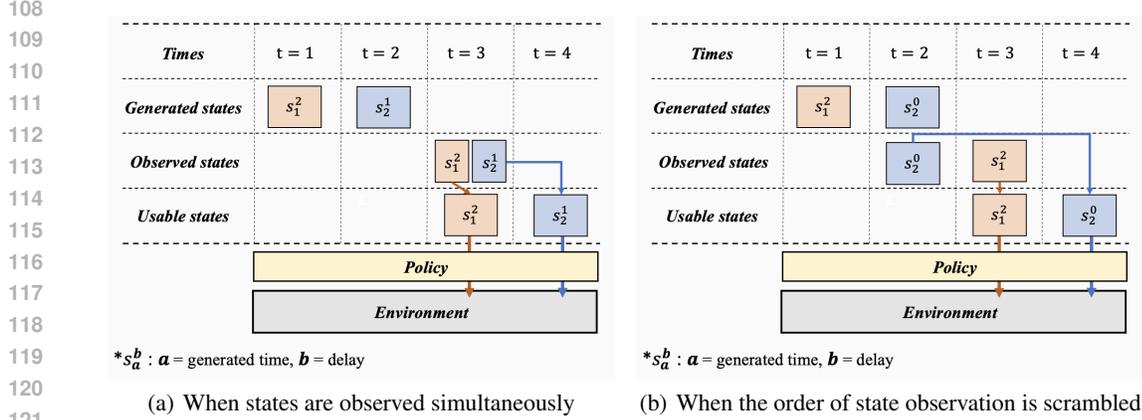


Figure 1: Visual examples illustrating cases where (a) states are observed simultaneously and (b) the order of state observation is scrambled. In both scenarios, an observed state is not available for use if a previously generated state has not yet been used in the decision-making process.

Suppose the state s_1 has an observation delay of 2, and the subsequent state s_2 has an observation delay of 1. These delayed states are *observed* by the agent at time 3 simultaneously. However, despite being observed, s_2 is not immediately *usable* at this time because the preceding state has not yet been used in the decision-making process. Thus, the agent uses the observed states in sequence to determine its actions: s_1 is used at time 3 and s_2 is used at time 4 (Fig. 1(a)). Throughout this study, the term ‘observe’ is utilized when a strict distinction between the two is unnecessary, specifically in cases of constant delays.

3 AUGMENTATION-BASED APPROACH

An augmentation-based approach is often preferred in delayed MDPs, as it retrieves the Markovian property and offers advantages for agents learning policies through conventional RL algorithms in such environments (Liotet et al., 2022; Kim et al., 2023; Wu et al., 2024). As demonstrated by Altman & Nain (1992); Katsikopoulos & Engelbrecht (2003), delayed MDPs can be reduced to equivalent MDPs without delays through this approach, known as *regular* MDPs, where the resulting optimal policies are optimal in the original delayed MDPs (Bander & White III, 1999; Katsikopoulos & Engelbrecht, 2003). The augmentation-based approach involves state augmentation, where the state is concatenated with additional **delay-related** information, similar to methods employed in conventional control theory (Kwon & Pearson, 1980; Park et al., 2008). In this section, we examine two types of delayed MDPs: CDMDPs and RDMDPs.

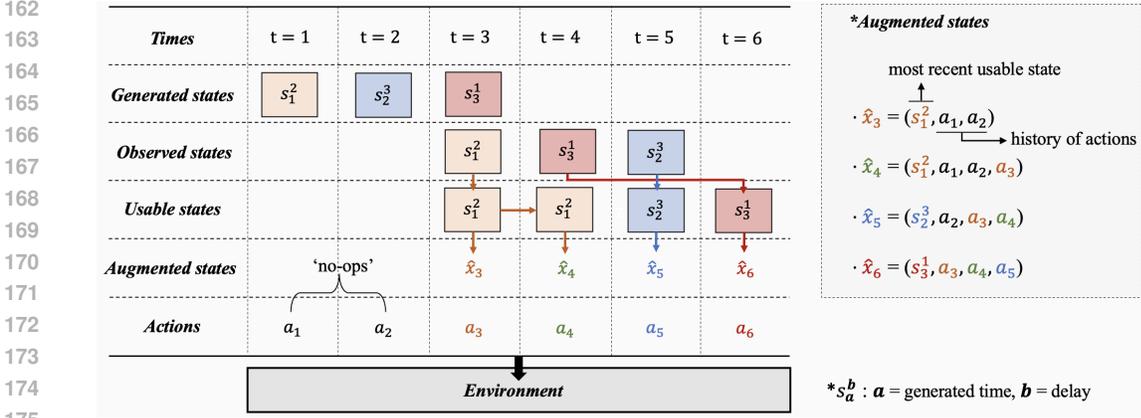
3.1 CONSTANT-DELAY MDPs

CDMDPs can be defined as a six-tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, o)$, where $o \in \mathbb{N}$ is a constant variable representing the observation delay. As demonstrated by Katsikopoulos & Engelbrecht (2003), it is reducible to regular MDPs $(\mathcal{X}_o, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$, where $\mathcal{X}_o = \mathcal{S} \times \mathcal{A}^o$ is the augmented state space with \mathcal{A}^o being the Cartesian product of \mathcal{A} with itself for o times, and $\mathcal{R} : \mathcal{X}_o \times \mathcal{A} \rightarrow \mathbb{R}$ is **the reward function with respect to the augmented state space, termed the augmented reward function**. Finally, the augmented state-based policy $\bar{\pi} : \mathcal{X}_o \times \mathcal{A} \rightarrow [0, 1]$ maps augmented state-to-action distribution.

To be specific, the augmented state at time t is defined as:

$$x_t = (s_{t-o}, a_{t-o}, a_{t-o+1}, \dots, a_{t-1}), \quad \forall t > o, \quad (3)$$

where s_{t-o} is the most recently *observed* state and $(a_{t-o}, \dots, a_{t-1})$ is the history of actions taken since s_{t-o} was generated. The agent implicitly estimates unobserved state s_t based on the augmented state x_t and selects action a_t accordingly. Note that since s_t is not explicitly known at time t , **the augmented reward** corresponding to the action a_t becomes a random variable that has to be determined based on the conditional expectation, which is given as $\bar{R}(x_t, a_t) := \mathbb{E}_{\mathbb{P}(s_t|x_t)} [R(s_t, a_t)]$.



177 Figure 2: A visual example of decision-making processes in environments with random delays. At
178 time 3, state s_1^1 is observed, and the agent immediately uses this observed state to make a decision.
179 At time 4, state s_3^1 is observed, whereas the preceding state s_2^3 remains unobserved and
180 consequently, s_3^1 cannot be used at this time because the previously generated state s_2^3 has not yet been observed and
181 used. Therefore, the agent continues to use s_1^1 until s_2^3 is observed. Eventually, s_2^3 is observed at time
182 5, and the observed states are then used in the correct order. Note that the states are reformulated as
183 augmented states before being fed into the policy, which subsequently determines the appropriate
184 actions.

185 3.2 RANDOM-DELAY MDPs

187 Using similar arguments as in Katsikopoulos & Engelbrecht (2003), RDMDPs can be defined as an
188 eight-tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, O, q_o, \tau)$, where $O \in \{0\} \cup \mathbb{N}$ is a random variable representing varying
189 observation delay, assumed to be sampled from an arbitrary discrete distribution q_o with support on
190 $\{0, 1, \dots, o_{\max}\}$, where o_{\max} denotes a maximum delay in time-steps, and $\tau : \mathcal{S} \rightarrow \mathbb{N}$ is a time-related
191 function that maps states to times at which they are used to make decisions for the first time. Under
192 this definition, we can define the states with observation delays and their corresponding augmented
193 states as follows:

194 **Definition 3.1.** Let $s_n^{o_n}$ be a state with an observation delay of o_n , where the subscript $n > 0$ denotes
195 the time at which the state was generated. Given that $s_n^{o_n}$ is the most recent *usable* state at time t ,
196 the augmented state at time t is defined as:

$$197 \hat{x}_t = (s_n^{o_n}, a_n, a_{n+1}, \dots, a_{t-1}), \quad \text{for } \tau(s_n^{o_n}) \leq t \leq n + o_{\max}, \quad (4)$$

199 where the sequence $(a_n, a_{n+1}, \dots, a_{t-1})$ represents the history of actions taken since $s_n^{o_n}$ was gen-
200 erated, and $\tau(s_n^{o_n}) \in \{n + o_n, \dots, t\}$ denotes the time when $s_n^{o_n}$ is used for decision-making for the
201 first time. Based on the augmented state \hat{x}_t , the agent implicitly estimates the unobserved state s_t
202 and selects action a_t accordingly. The notation $\hat{\cdot}$ indicates the augmented state defined in RDMDPs.

203 Given that $\tau(s_n^{o_n}) = t$ and $s_{n+1} \sim \mathcal{P}(\cdot | s_n, a_n)$, the usability of the next state $s_{n+1}^{o_{n+1}}$ at time $t + 1$
204 depends on its delay o_{n+1} . Consequently, the next augmented state \hat{x}_{t+1} is defined as:

$$205 \hat{x}_{t+1} = \begin{cases} (s_{n+1}^{o_{n+1}}, a_{n+1}, a_{n+2}, \dots, a_t) & \text{with prob. of } m(o_{n+1}), \\ (s_n^{o_n}, a_n, a_{n+1}, \dots, a_{t-1}, a_t) & \text{with prob. of } 1 - m(o_{n+1}), \end{cases} \quad (5)$$

206
207
208 where $m(o_{n+1}) = Pr(o_{n+1} \in \{0, 1, \dots, \tau(s_n^{o_n}) - n\})$.

209 Note that state $s_n^{o_n}$ continues to be used if the next state $s_{n+1}^{o_{n+1}}$ is not available for use at this time.
210 It remains in use until $s_{n+1}^{o_{n+1}}$ becomes usable, at which point the agent switches to use $s_{n+1}^{o_{n+1}}$ to
211 construct the augmented state (Fig. 2).
212

213 It is important to note that the dimension of the augmented state \hat{x}_{t+1} in equation 5 either remains
214 constant ($|\hat{x}_{t+1}| = |\hat{x}_t|$) or increases by 1 ($|\hat{x}_{t+1}| = |\hat{x}_t| + 1$). This implies that its dimension will
215 eventually reach infinity in infinite-horizon MDPs **without assuming a bounded maximum delay.**
Under the assumption of bounded maximum delays, Katsikopoulos & Engelbrecht (2003) proposed

a method called *freeze*. In this approach, the agent performs no actions ('no-ops') until the current state s is observed whenever the dimension of the augmented state reaches its maximum tolerable limit. Once the state s is observed, the augmented state is reset to $\hat{x} = (s)$, and the decision-making process resumes. However, this approach is known to be highly task-dependent, as the agent ignores environmental changes during the inactive periods, potentially leading to suboptimal policies.

4 BELIEF PROJECTION-BASED Q-LEARNING

While the augmentation-based approach provides a foundation for training agents with conventional RL algorithms in delayed MDPs, it has a known limitation: the state space grows exponentially as the number of delayed time-steps increases, resulting in sample inefficiency and slow convergence. This is called the state-space explosion issue, which makes the augmentation-based approach unfavorable for environments with long delays (Derman et al., 2021; Kim et al., 2023).

To *mitigate* this issue, belief projection-based Q-learning (BPQL), a model-free actor-critic framework designed to handle constant-delay environments, was proposed by Kim et al. (2023). It exhibited remarkable performance with simple modifications to the conventional augmentation-based approach, effectively *alleviating* the state-space explosion issue.

4.1 ALTERNATIVE REPRESENTATIONS FOR AUGMENTATION-BASED VALUES

First, a modified Bellman operator $\bar{\mathcal{T}}$, termed *delay Bellman operator*, was introduced to evaluate the values with respect to the augmented state space, which is given as:

$$\bar{\mathcal{T}}^{\bar{\pi}} \bar{V}^{\bar{\pi}}(x_t) \mapsto \mathbb{E}_{a_t \sim \bar{\pi}(\cdot|x_t)} \left[\mathbb{E}_{\mathbb{P}(s_t|x_t)} [R(s_t, a_t)] + \gamma \mathbb{E}_{x_{t+1} \sim \bar{\mathcal{P}}(\cdot|x_t, a_t)} [\bar{V}^{\bar{\pi}}(x_{t+1})] \right], \quad \forall t > o, \quad (6)$$

where $\bar{\pi}$ is the augmented state-based policy that receives augmented states as input, $\bar{V}^{\bar{\pi}}$ is the augmented state-based value representing the values of augmented states under the policy $\bar{\pi}$, o is a constant observation delay, and $\bar{\mathcal{P}} : \mathcal{X}_o \times \mathcal{A} \times \mathcal{X}_o \rightarrow [0, 1]$ is the transition kernel defined with respect to the augmented state space. By repeatedly applying the delay Bellman operator, $\bar{V}^{\bar{\pi}}$ converges, and then $\bar{\pi}$ is improved using conventional policy improvement methods. Similarly, the augmented state-based Q-value, $\bar{Q}^{\bar{\pi}}$, representing the values of augmented states for the given actions under the policy $\bar{\pi}$, can also be employed.

To *mitigate* the state-space explosion issue, the alternative representations for \bar{V} and \bar{Q} , referred to as *beta value* and *beta Q-value* (V_β and Q_β), are introduced. These values are evaluated with respect to the original state space rather than the augmented one, thereby naturally *alleviating* the state-space explosion issue. Beta-based values can be used as estimators for the augmentation-based values, which are given as:

$$\bar{V}^{\bar{\pi}}(x_t) = \mathbb{E}_{\mathbb{P}(s_t|x_t)} [V_\beta^{\bar{\pi}}(s_t)] + \Delta_{\text{residual}}^{\bar{\pi}}(x_t) \quad (7)$$

$$\bar{Q}^{\bar{\pi}}(x_t, a_t) = \mathbb{E}_{\mathbb{P}(s_t|x_t)} [Q_\beta^{\bar{\pi}}(s_t, a_t)] + \delta_{\text{residual}}^{\bar{\pi}}(x_t, a_t), \quad (8)$$

where $\Delta_{\text{residual}}^{\bar{\pi}}$ and $\delta_{\text{residual}}^{\bar{\pi}}$ represent the projection residuals.

In large or continuous spaces, a practical sampling-based reinforcement learning algorithm can be employed, where the beta Q-value and augmented state-based policy are parameterized by θ (beta critic) and ϕ (actor), respectively. The beta critic and actor are then trained by iteratively minimizing the following objective functions:

$$\mathcal{J}_{Q_\beta}(\theta) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}, x_{t+1}) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_{\theta, \beta}^{\bar{\pi}}(s_t, a_t) - R(s_t, a_t) - \gamma \mathbb{E}_{a_{t+1} \sim \bar{\pi}_\phi(\cdot|x_{t+1})} \left[Q_{\theta, \beta}^{\bar{\pi}}(s_{t+1}, a_{t+1}) - \alpha \log \bar{\pi}_\phi(a_{t+1}|x_{t+1}) \right] \right)^2 \right], \quad (9)$$

$$\mathcal{J}_{\bar{\pi}}(\phi) = \mathbb{E}_{(s_t, x_t) \sim \mathcal{D}} \left[\mathbb{E}_{a_t \sim \bar{\pi}_\phi(\cdot|x_t)} \left[\alpha \log \bar{\pi}_\phi(a_t|x_t) - Q_{\theta, \beta}^{\bar{\pi}}(s_t, a_t) \right] \right], \quad (10)$$

where $r_t = R(s_t, a_t)$, \mathcal{D} represents a replay buffer (Mnih et al., 2013), α is a temperature parameter (Haarnoja et al., 2018), and θ are the parameters of the target beta critic (Fujimoto et al., 2018).

Consequently, it demonstrated remarkable performance in constant-delay environments. However, since it was specifically designed to handle constant delays, it is not applicable to random-delay environments, which are the focus of our study. As discussed earlier, we demonstrate that conventional constant-delay approaches can be naturally applied to random-delay environments by establishing a connection between RDMDPs and CDMDPs, thereby facilitating the application of the BPQL framework in our study.

5 BRIDGING RDMDPs TO CDMDPs

In this section, we demonstrate that **RDMDPs can be transformed into their** equivalent CDMDPs by introducing a simple agent called lazy-agent, allowing **state-of-the-art** constant-delay approaches to be seamlessly extended to random-delay environments.

5.1 LAZY-AGENT

To address the challenges caused by variability in random delays, we let the agent assume that all states are delayed by the maximum number of time-steps. The agent then uses the observed states in decision-making processes at their maximum delayed times, that is, $\tau(s_n^{o_n}) = n + o_{\max}, \forall n > 0$. In this scheme, each state is consistently used in sequence exactly o_{\max} time-steps after being generated, regardless of its actual delay.

Suppose the state $s_1^{o_1}$ is observed at time $1 + o_1$. Since the agent assumes that all states are delayed by o_{\max} time-steps, it uses the observed state at time $1 + o_{\max}$ irrespective of its actual delay o_1 . Similarly, the subsequent state $s_2^{o_2}$ is observed at time $2 + o_2$; however, the agent uses it at time $2 + o_{\max}$. In short, the agent uses observed states at their maximum delayed times, regardless of their actual delays. This implies that the exact delays for each state may remain unknown to the agent, except for the maximum delay. **Consequently, this approach effectively circumvents the challenges associated with variability in random delays.** We refer to this agent as a *lazy-agent*. A visual representation is provided in Appendix E.

Formally, the augmented state \hat{x}_t in equation 4 is redefined for the lazy-agent as:

$$\hat{x}_t = (s_n^{o_n}, a_n, a_{n+1}, \dots, a_{t-1}), \quad \text{for } t = n + o_{\max}, \forall n > 0, \quad (11)$$

where $\tau(s_n^{o_n}) = n + o_{\max}$, resulting in the probability $m(o_{n+1})$ in equation 5 becoming 1. Consequently, given $s_{n+1} \sim \mathcal{P}(\cdot | s_n, a_n)$, the next augmented state \hat{x}_{t+1} in equation 5 is redefined as:

$$\hat{x}_{t+1} = (s_{n+1}^{o_{n+1}}, a_{n+1}, a_{n+2}, \dots, a_t), \quad \text{for } t = n + o_{\max}, \forall n > 0, \quad (12)$$

irrespective of its delay o_{n+1} .

Note that equation 11 and equation 12 align with the formulations for the augmented states defined in CDMDPs with a constant observation delay of o_{\max} . Furthermore, the dimension of the augmented state remains constant at $(o_{\max} + 1)$ at all times, addressing the issue of an exploding augmented state dimension in infinite-horizon MDPs.

Consequently, RDMDPs now become equivalent to CDMDPs with a constant delay of o_{\max} to the lazy-agents. The resulting CDMDPs can then be further reduced to regular MDPs, allowing the lazy-agents to be trained using conventional RL algorithms. To support our analysis, we present empirical results in Appendix B.2, demonstrating that the performance of lazy-agents trained in random-delay environments is comparable to that of agents trained in constant-delay environments. From these empirical results, we propose the following proposition:

Proposition 5.1. *Under the assumption of bounded observation delay and ordering, the RDMDPs $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, O, q_o, \tau)$ can be transformed into the equivalent CDMDPs $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, o_{\max})$ through the lazy-agent, where o_{\max} denotes the maximum delay in RDMDPs.*

Proof sketch. We begin by introducing the lazy-agent into RDMDPs, in which the agent assumes that all states are delayed by the maximum number of time-steps (o_{\max}) and uses the observed states for decision-making at their maximum delayed times, that is, $\tau(s_n^{o_n}) = n + o_{\max}, \forall n > 0$.

In this setting, the augmented state \hat{x}_t in equation 4 is redefined as:

$$\hat{x}_t = (s_n^{o_n}, a_n, a_{n+1}, \dots, a_{t-1}), \quad \text{for } t = n + o_{\max}, \forall n > 0, \quad (13)$$

324 indicating that $s_n^{o_n}$ is *used* for decision-making only at time $n + o_{\max}$, regardless of its actual delay.
 325 Under this assumption, the probability $m(o_{n+1})$ in equation 5 becomes:

$$326 \quad m(o_{n+1}) = Pr(o_{n+1} \in \{0, 1, \dots, \tau(s_n^{o_n}) - n\}) = Pr(o_{n+1} \in \{0, 1, \dots, o_{\max}\}) = 1. \quad (14)$$

328 Consequently, the next augmented state \hat{x}_{t+1} in equation 5 is redefined as:

$$329 \quad \hat{x}_{t+1} = (s_{n+1}^{o_{n+1}}, a_{n+1}, a_{n+2}, \dots, a_t), \quad \text{for } t = n + o_{\max}, \forall n > 0, \quad (15)$$

331 where $s_{n+1} \sim \mathcal{P}(\cdot | s_n, a_n)$. Evidently, these formulations for the augmented states are equivalent to
 332 those defined in CDMDPs with a constant delay of o_{\max} . To be more specific, by replacing n with
 333 $t - o_{\max}$ in equation 13, we obtain:

$$334 \quad \hat{x}_t = (s_{t-o_{\max}}^{o_{t-o_{\max}}}, a_{t-o_{\max}}, a_{t-o_{\max}+1}, \dots, a_{t-1}), \quad \forall t > o_{\max}, \quad (16)$$

336 which is exactly equivalent to equation 3 with $o = o_{\max}$, demonstrating that RDMDPs can be trans-
 337 formed into their equivalent CDMDPs through the lazy-agent. This completes the proof. \square

338 5.2 LAZY-BPQL

340 In the previous section, we demonstrated that RDMDPs can become equivalent to CDMDPs with
 341 lazy-agents. However, since the delays in the resulting CDMDPs are determined by the maximum
 342 delays in RDMDPs, the lazy-agents often encounter long-delay challenges, particularly the state-
 343 space explosion issue. Specifically, the augmented state space in derived CDMDPs would be defined
 344 as $\mathcal{X}_{o_{\max}} = \mathcal{S} \times \mathcal{A}^{o_{\max}}$, necessitating numerous samples for the augmented-based values to converge.

345 To address this issue, we employ lazy-agents within the BPQL framework to leverage its advan-
 346 tage in training agents in constant-delay environments while effectively **alleviating** the state-space
 347 explosion problem. We refer to this approach as *lazy-BPQL*, which is summarized in Algorithm 1.

349 6 EXPERIMENTS

350 6.1 BENCHMARKS AND BASELINE ALGORITHMS

351 We evaluated our algorithm on popular continuous control tasks in the MuJoCo benchmark by grad-
 352 ually increasing the maximum delay o_{\max} from 5 to 20, to assess its performance and robustness
 353 with respect to the degree of randomness in delays. In our experiments, we assumed that random
 354 delays are sampled from a discrete uniform distribution. Details of the benchmark environments and
 355 experiments are provided in Appendix D.

356 The following algorithms are included in experiments: normal SAC (Haarnoja et al., 2018), delayed-
 357 SAC (Derman et al., 2021; Kim et al., 2023), lazy-BPQL, and DC/AC (Bouteiller et al., 2020). The
 358 normal SAC adopts a naive approach that selects actions for currently usable states on a memoryless
 359 basis and performs ‘no-ops’ otherwise, without addressing the violation of the Markovian assump-
 360 tion in delayed MDPs. Delayed-SAC is a variant of delayed- Q (Derman et al., 2021) adapted by
 361 Kim et al. (2023) for application in continuous spaces. It employs an approximate forward model
 362 to explicitly predict unobserved states, which can be learned from transition samples collected in
 363 undelayed environments. With this model, the agent recursively predicts unobserved states through
 364 one-step predictions repeated over delayed time-steps and selects actions based on predicted states.
 365 Lastly, DC/AC is an improved version of SAC that implements an off-policy multi-step value estima-
 366 tion combined with a partial trajectory resampling method, significantly enhancing sample efficiency
 367 and demonstrating notable performance in environments with both constant and random delays.

370 6.2 RESULTS

371 6.2.1 PERFORMANCE COMPARISON

372 Table 1 and Fig. 3 show the performance of each algorithm on the MuJoCo tasks. The results indicate
 373 that the proposed lazy-BPQL demonstrates remarkable performance across all tasks, from relatively
 374 short delays ($o_{\max} = 5$) to long delays ($o_{\max} = 20$). In contrast, normal SAC performs poorly across
 375 all tasks, as it trains the agent directly in delayed environments without recovering the Markovian
 376 property, resulting in nearly random outcomes. Despite respectable performance in relatively simple
 377

tasks with small spaces, delayed-SAC exhibits unsatisfactory task-dependent performance, possibly due to the accumulation of nonnegligible prediction errors as the complexity of task or o_{\max} increases, underscoring the need for a more carefully designed dynamics model. Lastly, while DC/AC performs reasonably well in some tasks with short delays, its performance significantly deteriorates as the degree of randomness in delays increases. **To further highlight the performance achieved by lazy-BPQL compared to other baseline algorithms, we report the delay-free normalized scores (Wu et al., 2024) in Fig. 7 and Table 2 in Appendix B.1.**

Table 1: Results of the MuJoCo tasks with random delays of $o_{\max} \in \{5, 10, 20\}$. Each algorithm was evaluated for one million time-steps over five trials with different seeds. The standard deviations of average returns are denoted by \pm , and the best performance is in **bold**. Results for constant delays are in **blue**. **Additionally, delay-free SAC serves as the baseline performance in delay-free environments.**

Environment		Ant-v3	HalfCheetah-v3	Hopper-v3	Walker2d-v3	Humanoid-v3	InvertedPendulum-v2
o_{\max}	Algorithm						
x	Random policy	-58.7 \pm 4	-285.01 \pm 3	18.6 \pm 2	1.9 \pm 1	121.9 \pm 2	5.6 \pm 1
	Delay-free SAC	3279.2 \pm 180	8608.4 \pm 57	2435.2 \pm 23	3305.5 \pm 234	3228.1 \pm 410	964.3 \pm 29
5	Normal SAC	-76.6 \pm 4	-279.5 \pm 5	89.2 \pm 10	44.7 \pm 21	403.9 \pm 5	32.2 \pm 2
	DC/AC	907.5 \pm 90	2561.8 \pm 92	1931.6 \pm 192	2079.3 \pm 122	2798.4 \pm 452	854.7 \pm 30
	Delayed-SAC	986.4 \pm 128	4569.4 \pm 88	2200.4 \pm 190	1910.1 \pm 247	418.9 \pm 126	964.2 \pm 15
	Lazy-BPQL (proposed)	3679.8 \pm 167	5583.9 \pm 169	2174.1 \pm 155	2843.2 \pm 272	3157.7 \pm 292	958.8 \pm 14
	BPQL (constant-delay)	3761.9 \pm 112	5212.7 \pm 41	2136.3 \pm 158	2577.4 \pm 157	3194.9 \pm 374	955.9 \pm 28
10	Normal SAC	-84.6 \pm 9	-278.6 \pm 6	28.1 \pm 6	40.9 \pm 4	354.5 \pm 12	31.3 \pm 1
	DC/AC	342.9 \pm 34	1824.5 \pm 111	1262.4 \pm 261	1492.5 \pm 133	1023.8 \pm 359	4.9 \pm 0
	Delayed-SAC	966.9 \pm 180	2563.8 \pm 215	1878.5 \pm 176	1264.6 \pm 233	289.6 \pm 108	947.6 \pm 36
	Lazy-BPQL (proposed)	2744.5 \pm 112	4810.1 \pm 233	2300.9 \pm 164	2122.3 \pm 292	2820.5 \pm 348	936.9 \pm 38
	BPQL (constant-delay)	2831.9 \pm 103	4282.2 \pm 203	2129.2 \pm 184	2331.6 \pm 252	2891.5 \pm 357	934.7 \pm 20
20	Normal SAC	-83.1 \pm 9	-264.9 \pm 5	27.5 \pm 5	64.6 \pm 1	364.3 \pm 7	24.3 \pm 0
	DC/AC	258.3 \pm 42	860.9 \pm 288	12.8 \pm 6	-2.9 \pm 5	237.3 \pm 73	4.1 \pm 0
	Delayed-SAC	955.7 \pm 110	1377.8 \pm 140	1164.1 \pm 278	811.5 \pm 163	370.3 \pm 17	933.5 \pm 33
	Lazy-BPQL (proposed)	1976.5 \pm 248	3727.2 \pm 279	1346.7 \pm 245	1025.7 \pm 302	1143.8 \pm 371	566.9 \pm 88
	BPQL (constant-delay)	2078.9 \pm 157	3062.7 \pm 252	1526.7 \pm 227	846.7 \pm 443	1197.7 \pm 457	608.7 \pm 210

6.2.2 PERFORMANCE IN ENVIRONMENTS WITH CONSTANT DELAYS AND RANDOM DELAYS

To verify whether the proposed lazy-agents perform in random-delay environments as if they were in constant-delay environments, we evaluated the performance of lazy-lazy agents trained in environments with random delays of $o_{\max} \in \{5, 10, 20\}$ (lazy-BPQL), and compared it with normal agents trained in environments with constant delays of $o = o_{\max}$ (BPQL). The objective was to determine whether these two types of agents demonstrate similar performance. Table 1 presents the results for the MuJoCo tasks. The empirical findings confirm that both types of agents exhibited **almost identical** performance across all evaluated tasks, which strongly supports our argument that random-delay environments can be transformed into their equivalent constant-delay counterparts with the use of lazy-agents. Additional results are provided in Appendix B.2.

6.2.3 STATE-SPACE EXPLOSION ISSUE

To highlight the importance of mitigating the state-space explosion issue, we trained lazy-agents solely based on the augmentation-based approach without employing BPQL techniques, which we refer to as *lazy-augmented-SAC*. As presented in Table 5 in Appendix C.1, lazy-BPQL outperforms lazy-augmented-SAC for all evaluated tasks. Note that lazy-augmented-SAC completely failed to learn any useful policy even for tasks with $o_{\max} = 5$. These results clearly underscore the importance of alleviating the state-space explosion issue.

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

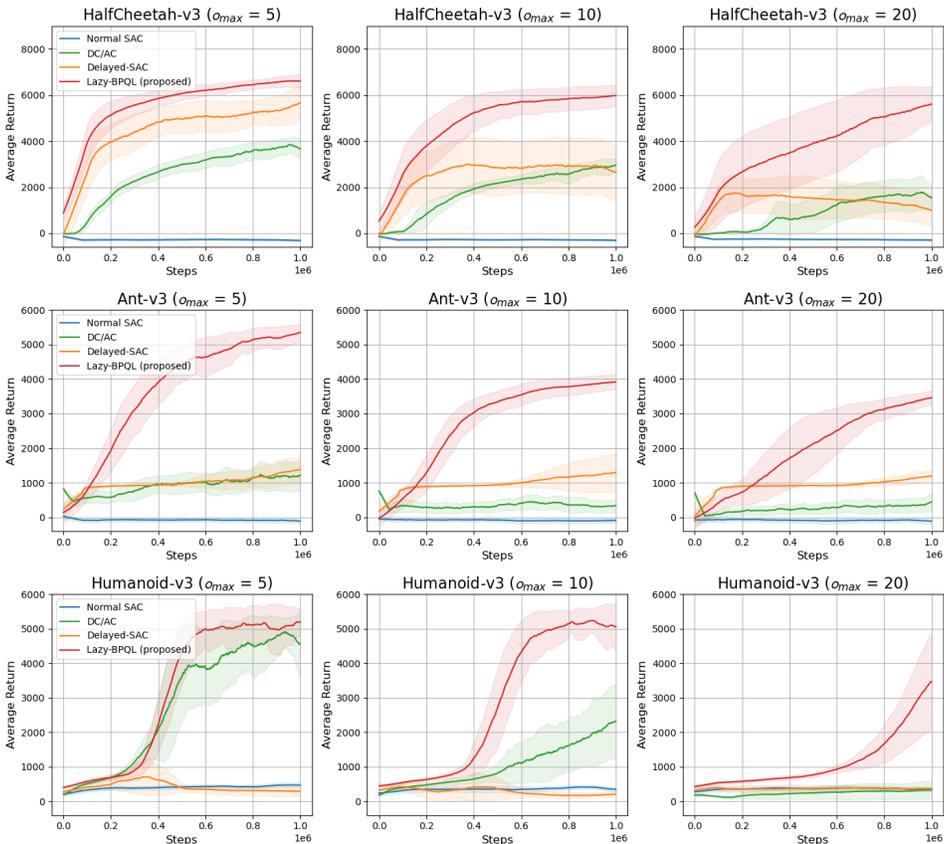


Figure 3: Performance curves of each algorithm on continuous control tasks in the MuJoCo benchmark with random delays of $o_{max} \in \{5, 10, 20\}$. All tasks were conducted with five different seeds for one million time-steps. The shaded regions represent the standard deviation of average returns across the trials. Across all tasks, the proposed lazy-BPQL exhibits remarkable performance, consistently outperforming other algorithms. **Additional results are provided in Appendix B.1.**

6.2.4 ENVIRONMENTS WITH HIGHER RANDOMNESS

We evaluated the performance of lazy-BPQL in random-delay environments with increased randomness ($o_{max} \in \{25, 30\}$) to empirically assess its robustness to greater randomness compared to other baseline algorithms. In the experiment, we included the second-best performing baseline, delayed-SAC, along with lazy-augmented-SAC to evaluate how effectively BPQL can mitigate the state-space explosion problem in such environments. The experiments were conducted on HalfCheetah-v3 and Ant-v3 tasks, and the result are listed in Table 6 in Appendix C.2. The results confirmed that lazy-BPQL exhibited performance degradation, but still maintained the best performance despite the increased randomness in delays, whereas other baselines were unable to learn any useful policies.

6.2.5 IMPACTS OF PROCESSING STATES IN ORDER

We investigated the impact of the assumption that states are used in order by comparing the performance of agents trained with and without this assumption. The results in Appendix C.3 reveal that the order in which observed states are used for decision-making can significantly affect the performance and learning stability of RL agents, with a notable drop in performance in the unordered case. Furthermore, the performance degradation becomes more pronounced as the randomness of delays increases. These findings seem to originate from the fact that both augmentation-based and model-based approaches heavily rely on preserving and understanding cause-and-effect relationships to restore the violated Markovian property caused by delays.

7 CONCLUSION

We investigated environments with random observation delays and proposed a novel approach to establish a connection to environments with constant delays by introducing a simple agent called the lazy-agent. **With the proposed lazy-agents, random-delay environments can be transformed into their equivalent constant-delay counterparts, facilitating the application of state-of-the-art constant-delay approaches to random-delay environments without any modifications.** We employed lazy-agents within the belief projection-based Q -learning (BPQL) framework, referred to as lazy-BPQL, to train our agents in equivalent constant-delay environments while effectively **mitigating** the state-space explosion issue of the augmentation-based approach. The empirical results demonstrated that the proposed lazy-BPQL significantly outperformed other baseline algorithms in terms of asymptotic performance and sample efficiency in random-delay environments, which strongly supports the efficacy of our approach.

It would be meaningful to employ our lazy-agents in real-world dynamic systems that suffer from random delays, where conventional constant-delay approaches are inadequate. In the future, we will extend the proposed algorithm to real-world applications, such as robotic locomotion and manipulation, by accounting for randomly varying sensor and actuator delays. We believe that the lazy-agents will play a pivotal role in extending conventional RL methods to real-world dynamic systems.

8 LIMITATIONS

Despite its notable advantages in constructing equivalent constant-delay environments from the original random-delay environments, employing the lazy-agent may encounter difficulties associated with long delays. This is because the constant delays in the equivalent environments are aligned with the maximum delay in the original random-delay environments. Consequently, the model-based approach may require more carefully designed dynamics models, as accumulated errors in recursive one-step predictions could result in significant performance degradation. On the other hand, the augmentation-based approach may confront the inherent state-space explosion issue. To circumvent this issue, we adopted a strategy of training lazy-agents within the BPQL framework, which can effectively mitigate the state-space explosion issue. Alternatively, the use of recurrent models, such as GRU (Cho, 2014), can also be considered, as explored in Firoiu et al. (2018). However, the necessity of knowing the maximum delay raises another concern, which may be unrealistic in some environments. This remains a challenge to be addressed in future work.

REFERENCES

- Eitan Altman and Philippe Nain. Closed-loop control with delayed information. *ACM Sigmetrics Performance Evaluation Review*, 20(1):193–204, 1992.
- James L Bander and Chelsea C White III. Markov decision processes with noise-corrupted and delayed state observations. *Journal of the Operational Research Society*, 50(6):660–668, 1999.
- Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, pp. 679–684, 1957.
- Yann Bouteiller, Simon Ramstedt, Giovanni Beltrame, Christopher Pal, and Jonathan Binas. Reinforcement learning with random delays. In *International Conference on Learning Representations*, 2020.
- Baiming Chen, Mengdi Xu, Liang Li, and Ding Zhao. Delay-aware model-based reinforcement learning for continuous control. *Neurocomputing*, 450:119–128, 2021.
- Kyunghyun Cho. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Esther Derman, Gal Dalal, and Shie Mannor. Acting in delayed environments with non-stationary Markov policies. *arXiv preprint arXiv:2101.11992*, 2021.
- Vlad Firoiu, Tina Ju, and Josh Tenenbaum. At human speed: deep reinforcement learning with action delay. *arXiv preprint arXiv:1810.07286*, 2018.

- 540 Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-
541 critic methods. In *International Conference on Machine Learning*, pp. 1587–1596. PMLR, 2018.
- 542
- 543 Yuan Ge, Qigong Chen, Ming Jiang, and Yiqing Huang. Modeling of random delays in networked
544 control systems. *Journal of Control Science and Engineering*, 2013(1):383415, 2013.
- 545
- 546 Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
547 Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-
548 cations. *arXiv preprint arXiv:1812.05905*, 2018.
- 549
- 550 Jemin Hwangbo, Inkyu Sa, Roland Siegwart, and Marco Hutter. Control of a quadrotor with rein-
551 forcement learning. *IEEE Robotics and Automation Letters*, 2(4):2096–2103, 2017.
- 552
- 553 Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre
554 Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforce-
555 ment learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pp. 651–
556 673. PMLR, 2018.
- 557
- 558 Konstantinos V Katsikopoulos and Sascha E Engelbrecht. Markov decision processes with de-
559 lays and asynchronous cost collection. *IEEE Transactions on Automatic Control*, 48(4):568–574,
560 2003.
- 561
- 562 Jangwon Kim, Hangyeol Kim, Jiwook Kang, Jongchan Baek, and Soohee Han. Belief projection-
563 based reinforcement learning for environments with delayed feedback. *Advances in Neural Infor-
564 mation Processing Systems*, 36:678–696, 2023.
- 565
- 566 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
567 2014.
- 568
- 569 Woosuk Kwon and A Pearson. Feedback stabilization of linear systems with delayed control. *IEEE
570 Transactions on Automatic control*, 25(2):266–269, 1980.
- 571
- 572 Pierre Liotet, Davide Maran, Lorenzo Bisi, and Marcello Restelli. Delayed reinforcement learning
573 by imitation. In *International Conference on Machine Learning*, pp. 13528–13556. PMLR, 2022.
- 574
- 575 A Rupam Mahmood, Dmytro Korenkevych, Brent J Komer, and James Bergstra. Setting up a rein-
576 forcement learning task with a real-world robot. In *2018 IEEE/RSJ International Conference on
577 Intelligent Robots and Systems (IROS)*, pp. 4635–4640. IEEE, 2018.
- 578
- 579 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wier-
580 stra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint
581 arXiv:1312.5602*, 2013.
- 582
- 583 George E Monahan. State of the art—a survey of partially observable Markov decision processes:
584 theory, models, and algorithms. *Management Science*, 28(1):1–16, 1982.
- 585
- 586 Jung Hun Park, Han Woong Yoo, Soohee Han, and Wook Hyun Kwon. Receding horizon controls
587 for input-delayed systems. *IEEE Transactions on Automatic Control*, 53(7):1746–1752, 2008.
- 588
- 589 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
590 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
591 the game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 592
- 593 Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Learning without state-estimation in
partially observable Markovian decision processes. In *Machine Learning Proceedings 1994*, pp.
284–292. Elsevier, 1994.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033.
IEEE, 2012.
- Thomas J Walsh, Ali Nouri, Lihong Li, and Michael L Littman. Learning and planning in environ-
ments with delayed feedback. *Autonomous Agents and Multi-Agent Systems*, 18:83–105, 2009.

594 Qingyuan Wu, Simon Sinong Zhan, Yixuan Wang, Yuhui Wang, Chung-Wei Lin, Chen Lv, Qi Zhu,
595 Jürgen Schmidhuber, and Chao Huang. Boosting reinforcement learning with strongly delayed
596 feedback through auxiliary short delays. In *Forty-First International Conference on Machine*
597 *Learning*, 2024. URL <https://openreview.net/forum?id=0IDaPnY5d5>.
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A RELATED WORK

Delays are prevalent in real-world reinforcement learning (RL) applications, arising from various factors such as computational time. Previous research on handling such delays within the RL framework can be categorized into two main approaches: augmentation-based and model-based.

Augmentation-based approach The augmentation-based approach is often preferred in delayed environments since it allows us to retrieve the violated Markovian property (Altman & Nain, 1992; Katsikopoulos & Engelbrecht, 2003) and offers advantages for training agents through conventional RL algorithms in such environments. Despite its notable advantages, the method of augmenting the state itself remains a challenge associated with sample complexity, commonly referred to as the state-space explosion issue or the curse of dimensionality, which results in learning inefficiency. To mitigate this issue, DC/AC (Bouteiller et al., 2020) introduces an off-policy multi-step value estimation combined with a partial trajectory resampling method, which greatly enhances sample efficiency and accelerates the learning process. BPQL (Kim et al., 2023) proposes another novel approach to overcome such difficulty by introducing alternative representations for augmentation-based values. These alternative values are evaluated with respect to the original state space rather than the augmented one, thereby inherently mitigating the state-space explosion issue and demonstrating remarkable performance. More recently, AD-RL (Wu et al., 2024) alleviates the performance degradation stemming from this issue by leveraging auxiliary tasks with shorter delays to learn tasks with relatively long delays, reducing sample complexity and achieving notable performance.

Model-based approach The model-based approach, also known as the state estimation method, aims to restore the Markovian property using learned dynamics models from underlying delay-free environments (Walsh et al., 2009; Firoiu et al., 2018; Chen et al., 2021; Derman et al., 2021). For example, delayed- Q (Derman et al., 2021) employs an approximate feed-forward model to explicitly predict unobserved states, which can be learned from transition samples collected in delay-free environments. Using this model, the agent recursively predicts unobserved states through one-step predictions repeated over delayed time-steps and selects actions based on the predicted states. Similarly, Firoiu et al. (2018) proposes an approach that utilizes recurrent neural networks (Cho, 2014) to model the dynamics. These approaches facilitate sample-efficient learning without being affected by the issues associated with the sample complexity posed by the augmented state. However, approximation errors in building dynamics models may induce accumulated prediction errors, leading to suboptimal performance. Furthermore, the presence of noise in observations can exacerbate inaccuracies in learning dynamics models.

While numerous methods have shown promise, most works focus on the unrealistic assumption of constant delays, exhibiting nonnegligible performance degradation when applied to environments with random delays. Building upon previously proposed state-of-the-art methods, this study makes its primary contribution by proposing a novel approach that enables the handling of constant delays and random delays in exactly the same manner. Specifically, we propose a method to construct equivalent constant-delay environments from the original random-delay environments by introducing a simple agent termed the lazy-agent. This approach offers valuable insight that there is no need to devise new methods for handling random delays, as the lazy-agent naturally facilitates the application of conventionally proposed state-of-the-art methods, originally designed for constant delays, to random-delay environments without any modifications.

B EXPERIMENTAL RESULTS

B.1 PERFORMANCE COMPARISON

Performance curves We present the performance curves of each algorithm on the MuJoCo tasks with random delays of $o_{\max} \in \{5, 10, 20\}$. All tasks were conducted with five different seeds for one million time-steps. The shaded regions represent the standard deviation of average returns. Empirical results demonstrate that lazy-BPQL exhibits remarkable performance across all evaluated tasks.

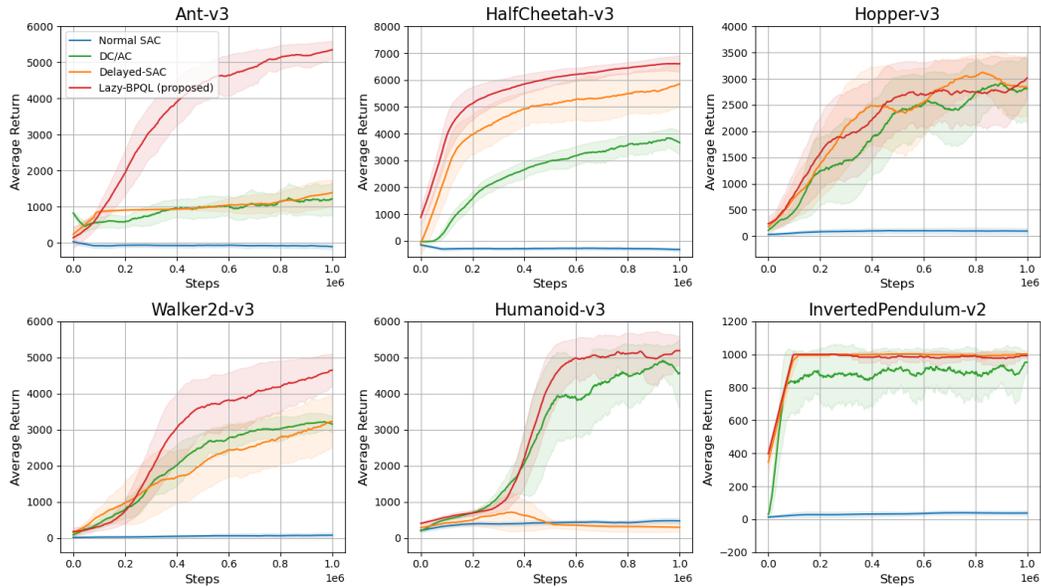


Figure 4: Performance curves of each algorithm on the MuJoCo tasks with $o_{\max} = 5$.

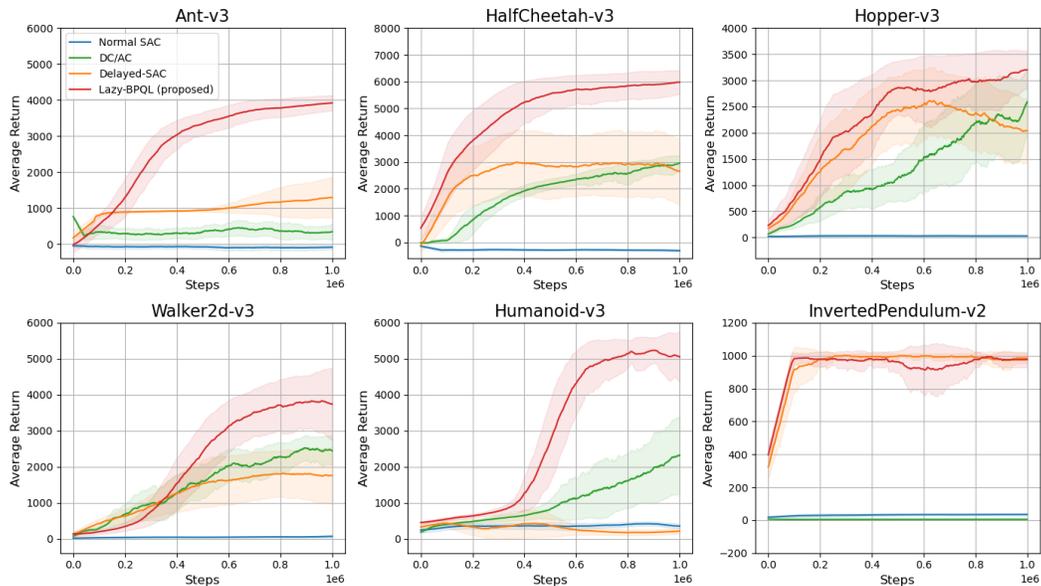


Figure 5: Performance curves of each algorithm on the MuJoCo tasks with $o_{\max} = 10$.

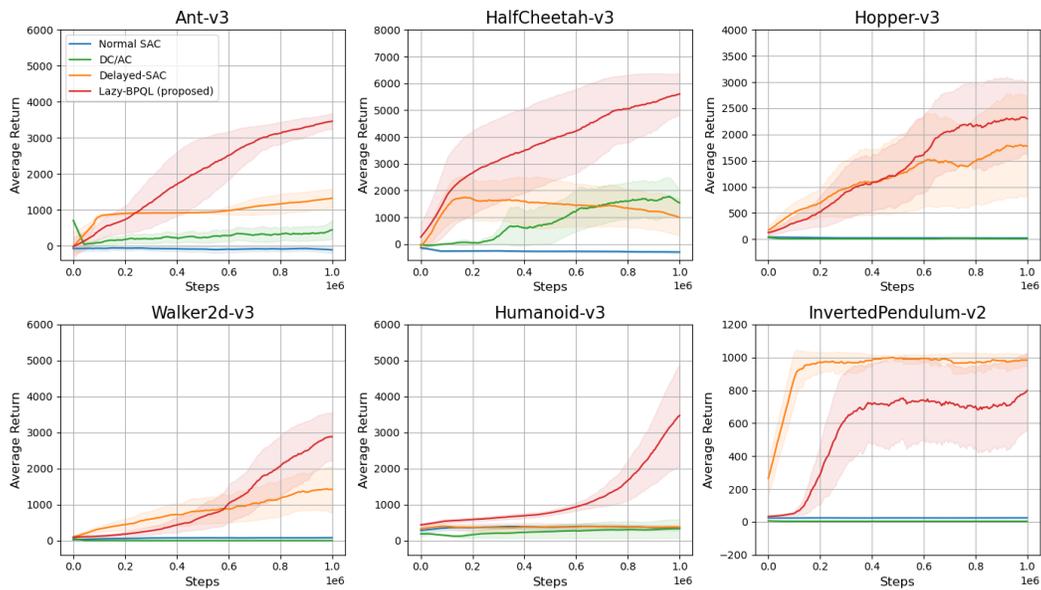


Figure 6: Performance curves of each algorithm on the MuJoCo tasks with $o_{\max} = 20$.

Delay-free normalized scores To clarify the performance achieved by lazy-BPQL compared to other baseline algorithms, we report the delay-free normalized scores for each algorithm on the MuJoCo tasks in Fig. 7 and Table 2, following Wu et al. (2024). The delay-free normalized score is defined as $R_{\text{normalized}} = (R_{\text{algorithm}} - R_{\text{random}}) / (R_{\text{delay-free}} - R_{\text{random}})$, where $R_{\text{algorithm}}$, $R_{\text{delay-free}}$, and R_{random} represent the average returns of the baselines, delay-free SAC, and random policy, respectively. Here, delay-free SAC serves as the baseline performance in delay-free environments.

From the results, we confirmed that for tasks with $o_{\max} = \{5, 10\}$, lazy-BPQL exhibits the best performance comparable to the delay-free performance, achieving average scores of 0.91 and 0.81, respectively. It outperforms the second-best performing baselines by wide average margins of 0.28 and 0.34 points, each. Even for tasks with the longest maximum delay of $o_{\max} = 20$, lazy-BPQL maintains the highest average score. These scores further highlight the effectiveness of lazy-BPQL, demonstrating its superiority over other baseline algorithms across all evaluated tasks in MuJoCo.

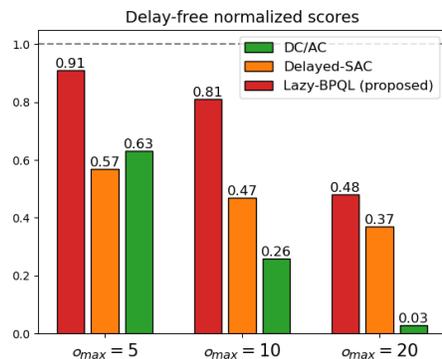


Figure 7: Delay-free normalized scores for each baseline algorithm, averaged across all the evaluated MuJoCo tasks. The dashed gray line represents the baseline score of delay-free SAC.

Table 2: Delay-free normalized scores of each algorithm with random delays of $o_{\max} \in \{5, 10, 20\}$. The best score is highlighted in bold.

Environment		Ant-v3	HalfCheetah-v3	Hopper-v3	Walker2d-v3	Humanoid-v3	InvertedPendulum-v2	Avg.
o_{\max}	Algorithm							
5	DC/AC	0.28	0.32	0.79	0.62	0.86	0.88	0.63
	Delayed-SAC	0.31	0.54	0.90	0.57	0.09	1.01	0.57
	Lazy-BPQL (proposed)	1.12	0.68	0.89	0.85	0.97	0.99	0.91
10	DC/AC	0.11	0.23	0.51	0.45	0.29	-0.01	0.26
	Delayed-SAC	0.30	0.32	0.77	0.38	0.05	0.99	0.47
	Lazy-BPQL (proposed)	0.84	0.57	0.94	0.64	0.87	0.98	0.81
20	DC/AC	0.09	0.12	-0.01	0.00	0.04	-0.01	0.03
	Delayed-SAC	0.30	0.18	0.48	0.24	0.08	0.97	0.37
	Lazy-BPQL (proposed)	0.61	0.45	0.55	0.30	0.32	0.58	0.48

B.2 EMPIRICAL RESULTS FOR PROPOSITION 5.1

To verify whether lazy-agents can perform in random-delay environments as if they were in constant-delay environments, we compared the performance of lazy-agents trained in **random-delay** environments (lazy-BPQL) with normal agents trained in **constant-delay** environments (BPQL) with constant delays set to $o = o_{\max}$. Each algorithm was evaluated for one million time-steps over five trials with different seeds on MuJoCo tasks, and the corresponding results are presented in Fig. 8, Table 3 and Table 4.

The results confirmed that both agents exhibited **almost identical** performance across all evaluated MuJoCo tasks with average margins of 0.1 in delay-free normalized scores (Wu et al., 2024). These empirical results strongly support our arguments that random-delay environments can be transformed into their equivalent constant-delay counterparts through the use of lazy-agents.

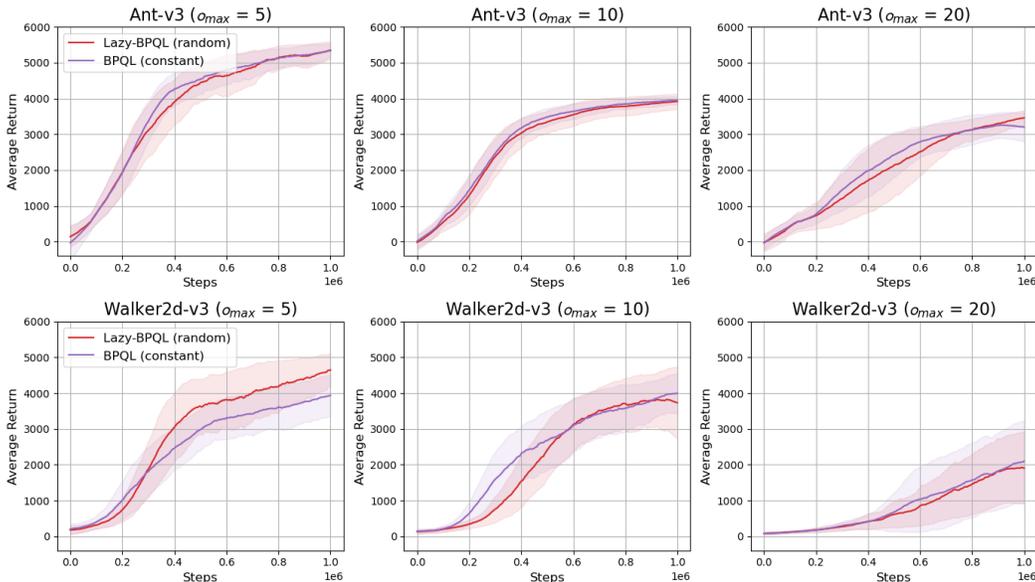


Figure 8: Performance curves of the proposed lazy-BPQL agents trained in random-delay environments with $o_{\max} \in \{5, 10, 20\}$ and the BPQL agent trained in constant-delay environments with $o = o_{\max}$ on continuous control tasks in the MuJoCo benchmark. All tasks were conducted with five different seeds for one million time-steps, and the shaded regions represent the standard deviation of average returns across the trials.

Table 3: Results of lazy-BPQL with random delays of $o_{\max} \in \{5, 10, 20\}$, and BPQL with constant delays of $o = o_{\max}$.

Environment		Ant-v3	HalfCheetah-v3	Hopper-v3	Walker2d-v3	Humanoid-v3	InvertedPendulum-v2
o_{\max}	Algorithm						
5	Lazy-BPQL (random-delay)	3679.8 \pm 167	5583.9 \pm 169	2174.1 \pm 155	2843.2 \pm 272	3157.7 \pm 292	958.8 \pm 14
	BPQL (constant-delay)	3761.9 \pm 112	5212.7 \pm 41	2136.3 \pm 158	2577.4 \pm 157	3194.9 \pm 374	955.9 \pm 28
10	Lazy-BPQL (random-delay)	2744.5 \pm 112	4810.1 \pm 233	2300.9 \pm 164	2122.3 \pm 292	2820.5 \pm 348	936.9 \pm 38
	BPQL (constant-delay)	2831.9 \pm 103	4282.2 \pm 203	2129.2 \pm 184	2331.6 \pm 252	2891.5 \pm 357	934.7 \pm 20
20	Lazy-BPQL (random-delay)	1976.5 \pm 248	3727.2 \pm 279	1346.7 \pm 245	1025.7 \pm 302	1143.8 \pm 371	566.9 \pm 88
	BPQL (constant-delay)	2078.9 \pm 157	3062.7 \pm 252	1526.7 \pm 227	846.7 \pm 443	1197.7 \pm 457	608.7 \pm 210

Table 4: Delay-free normalized scores of lazy-BPQL with random delays of $o_{\max} \in \{5, 10, 20\}$, and BPQL with constant delays of $o = o_{\max}$.

Environment		Ant-v3	HalfCheetah-v3	Hopper-v3	Walker2d-v3	Humanoid-v3	InvertedPendulum-v2	Avg.	Residue.
o_{\max}	Algorithm								
5	Lazy-BPQL (random-delay)	1.12	0.68	0.89	0.85	0.97	0.99	0.91	0.01
	BPQL (constant-delay)	1.14	0.62	0.88	0.74	0.99	1.00	0.90	
10	Lazy-BPQL (random-delay)	0.84	0.57	0.94	0.64	0.87	0.98	0.81	0.01
	BPQL (constant-delay)	0.86	0.51	0.88	0.70	0.89	0.97	0.80	
20	Lazy-BPQL (random-delay)	0.61	0.45	0.55	0.30	0.32	0.58	0.48	0.01
	BPQL (constant-delay)	0.64	0.35	0.63	0.25	0.34	0.58	0.47	

C ABLATION STUDY

C.1 STATE-SPACE EXPLOSION ISSUE

In this section, we present the performance of lazy-augmented-SAC and lazy-BPQL on the MuJoCo tasks with random delays of $o_{\max} \in \{5, 10, 20\}$. As listed in Table. 5, lazy-BPQL outperformed lazy-augmented-SAC across all evaluated tasks. Note that lazy-augmented-SAC completely failed to learn any useful policy even for tasks with $o_{\max} = 5$. These results clearly highlights the importance of mitigating the state-space explosion issue when employing augmentation-based approaches.

Table 5: Results of lazy-augmented-SAC and lazy-BPQL with random delays of $o_{\max} \in \{5, 10, 20\}$. Each algorithm was evaluated for one million time-steps over five trials with different seeds.

		o_{\max}	5	10	20
Environment	Algorithm				
Ant-v3	Lazy-BPQL (proposed)		3679.8\pm167	2744.5\pm112	1976.5\pm248
	Lazy-augmented-SAC		898.5 \pm 93	913.2 \pm 29	721.4 \pm 86
HalfCheetah-v3	Lazy-BPQL (proposed)		5583.9\pm169	4810.1\pm233	3727.2\pm279
	Lazy-augmented-SAC		2137.2 \pm 361	1068.3 \pm 122	500.9 \pm 137

C.2 ENVIRONMENTS WITH HIGHER RANDOMNESS

In this section, we provide the performance of lazy-BPQL in random-delay environments with increased randomness ($o_{\max} \in \{25, 30\}$) to empirically assess its robustness to greater randomness compared to other baseline algorithms. In experiment, we included the second-best performing baseline, delayed-SAC, along with lazy-augmented-SAC to verify how effectively BPQL can address the state-space explosion issue. The experiments were conducted in HalfCheetah-v3 and Ant-v3 tasks. Each algorithm was evaluated for one million time-steps over five trials with different seeds, and the results are listed in Table 6 and Table 7.

The results confirm that lazy-BPQL exhibited performance degradation, but still maintained the best performance despite the increased randomness in delays up to $o_{\max} = 30$, whereas other baselines were unable to learn any useful policies.

Table 6: Results of each baseline with random delays of $o_{\max} \in \{20, 25, 30\}$.

		o_{\max}		
Environment	Algorithm	20	25	30
Ant-v3	Lazy-BPQL (proposed)	1976.5 \pm 248	1944.3 \pm 176	1600.2 \pm 161
	Lazy-augmented-SAC	721.4 \pm 86	466.3 \pm 114	-34.3 \pm 81
	Delayed-SAC	955.7 \pm 110	949.9 \pm 141	961.2 \pm 154
HalfCheetah-v3	Lazy-BPQL (proposed)	3727.2 \pm 279	2492.1 \pm 379	1971.1 \pm 265
	Lazy-augmented-SAC	500.9 \pm 137	-5.8 \pm 131	-199.1 \pm 25
	Delayed-SAC	1377.8 \pm 140	1076.5 \pm 123	1194.8 \pm 73

Table 7: Delay-free normalized scores of each baseline with random delays of $o_{\max} \in \{20, 25, 30\}$.

		o_{\max}		
Environment	Algorithm	20	25	30
Ant-v3	Lazy-BPQL (proposed)	0.61	0.60	0.49
	Lazy-augmented-SAC	0.23	0.15	0.07
	Delayed-SAC	0.30	0.29	0.31
HalfCheetah-v3	Lazy-BPQL (proposed)	0.45	0.31	0.25
	Lazy-augmented-SAC	0.08	0.03	0.01
	Delayed-SAC	0.19	0.15	0.16

C.3 IMPACT OF PROCESSING STATES IN ORDER

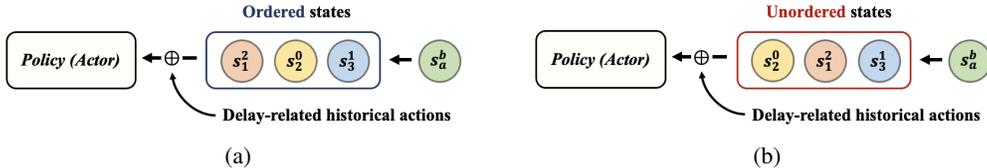


Figure 9: The visual examples illustrating cases where (a) the observed states are processed in order and (b) the observed states are processed out of order. \oplus denotes the concatenation operation.

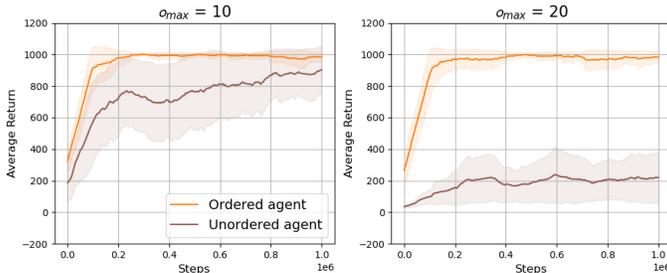


Figure 10: Results of ordered and unordered agent in InvertedPendulum-v2 MuJoCo task with random delays of $o_{max} \in \{10, 20\}$.

In the presence of randomness in observation delays, states may be observed simultaneously, and their order can even become scrambled. When utilizing these scrambled states for decision-making in random-delay environments, they can be used either in the observed order (unordered state processing) or in their original generated order (ordered state processing).

We investigated the impact of the assumption that states are used in order by comparing the performance of agents trained with and without this assumption (see Fig. 9). In the experiment, we utilized delayed-SAC for learning InvertedPendulum-v2 task in MuJoCo, as it demonstrated respectable and stable performance in relatively simple tasks. We aimed to verify how this assumption impacts such performance, even in such simple task. We refer to the delayed-SAC agent trained in an ordered manner as the *ordered agent*, and the agent trained in a disordered manner as the *unordered agent*. Each agent was evaluated for one million time-steps over five trials with random seeds, and the corresponding results are presented in Fig. 10 and Table 8.

Table 8: Results of ordered and unordered agent in InvertedPendulum-v2 MuJoCo task with random delays of $o_{max} \in \{10, 20\}$. The standard deviations of average returns are denoted by \pm .

Environment		InvertedPendulum-v2
o_{max}	Algorithm	
10	Unordered agent	739.5 \pm 36
	Ordered agent	947.6 \pm 36
20	Unordered agent	181.6 \pm 40
	Ordered agent	933.5 \pm 33

The results reveal that the order in which observed states are used can significantly affect the performance and learning stability of RL agents, with a notable drop in performance in the unordered case. Furthermore, the performance degradation becomes more pronounced as the randomness of delays increases. These findings seem to originate from the fact that both augmentation-based and model-based approaches heavily rely on preserving and understanding cause-and-effect relationships to restore the violated Markovian property caused by delays.

D EXPERIMENTAL DETAILS

D.1 ENVIRONMENTAL DETAILS

Table 9: Environmental details of the MuJoCo benchmark.

Task	State dimension	Action dimension	Time-step (s)
Ant-v3	27	8	0.05
HalfCheetah-v3	17	6	0.05
Walker2d-v3	17	6	0.008
Hopper-v3	11	3	0.008
Humanoid-v3	376	17	0.015
InvertedPendulum-v2	4	1	0.04

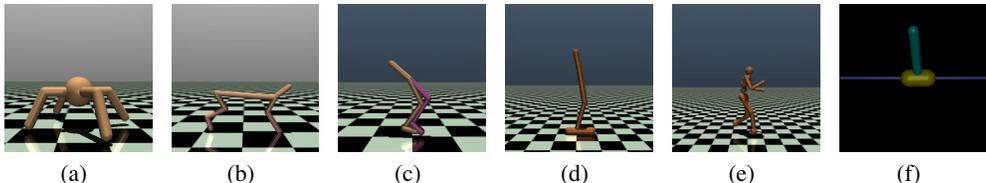


Figure 11: Experimental environments in the MuJoCo benchmark: (a) Ant-v3 (b) HalfCheetah-v3, (c) Walker2d-v3, (d) Hopper-v3, (e) Humanoid-v3, and (f) InvertedPendulum-v2

D.2 IMPLEMENTATION DETAILS

The implementation details of the proposed lazy-BPQL align with those presented in Kim et al. (2023), with the specific hyperparameters listed in Table 10. Since the baseline algorithms included in our experiments employ the SAC algorithm as their foundational learning algorithm, the hyperparameters are consistent across all approaches, except for the DC/AC algorithm.

Table 10: Hyperparameters for lazy-BPQL and the baselines.

Hyperparameters	Values
Actor network	256, 256
Critic network	256, 256
Learning rate (actor)	3e-4
Learning rate (critic)	3e-4
Temperature (α)	0.2
Discount factor (γ)	0.99
Replay buffer size	1e6
Mini-Batch size	256
Target entropy	$-\dim \mathcal{A} $
Target smoothing coefficient (ξ)	0.995
Optimizer	Adam (Kingma, 2014)
Total time-steps	1e6

D.3 PSEUDO CODE OF LAZY-BPQL

The proposed lazy-agent can be seamlessly integrated into the BPQL framework with minimal modifications by *using* the initial state for decision-making at its maximum delayed times. Subsequently, all states become naturally available for use at their respective maximum delayed times.

In the implementation, a temporary buffer \mathcal{B} has been employed, as utilized by Kim et al. (2023), to store *observed* states, corresponding rewards, and action histories, which enables the agent to access timely and relevant information for constructing augmented states. **Additionally, we have assumed that all feedback, including reward, is maximally delayed in equivalent constant-delay environments, similar to (Kim et al., 2023). Thus, the reward corresponding to the action a_t is assumed to be $r_{t-o_{\max}}$.**

Algorithm 1 Lazy Belief Projection-based Q -Learning (Lazy-BPQL)

```

1: Input: actor  $\bar{\pi}_\phi(a|\hat{x})$ , beta critic  $Q_{\theta,\beta}(s,a)$ , target beta critic  $Q_{\bar{\theta},\beta}(s,a)$ , replay buffer  $\mathcal{D}$ , temporary buffer  $\mathcal{B}$ , maximum delay  $o_{\max}$ , beta critic learning rate  $\lambda_Q$ , actor learning rate  $\lambda_\pi$ , soft update rate  $\xi$ , episodic length  $H$ , and total number of episodes  $E$ .
2: for episode  $e = 1$  to  $E$  do
3:   for time-step  $t = 1$  to  $H$  do
4:     if  $t < o_{\max}$  then
5:       select random or ‘no-ops’ action  $a_t$ 
6:       execute  $a_t$  on environment
7:       put  $a_t$ , observed states, rewards to  $\mathcal{B}$ 
8:     else if  $t = o_{\max}$  then ▷ wait for  $o_{\max}$  time-steps
9:       select random or ‘no-ops’ action  $a_t$ 
10:      execute  $a_t$  on environment
11:      put  $a_t$ , observed states, rewards to  $\mathcal{B}$ 
12:     else
13:       get  $s_{t-o_{\max}}, a_{t-o_{\max}}, \dots, a_{t-1}$  from  $\mathcal{B}$ 
14:       ▷ get most recent usable state and action histories
15:        $\hat{x}_t \leftarrow (s_{t-o_{\max}}, a_{t-o_{\max}}, \dots, a_{t-1})$  ▷ construct augmented state
16:        $a_t \leftarrow \bar{\pi}_\phi(\hat{x}_t)$ 
17:       execute  $a_t$  on environment
18:       put  $a_t$ , observed states, rewards to  $\mathcal{B}$ 
19:       if  $t > 2o_{\max}$  then
20:         get  $s_{t-2o_{\max}}, s_{t-2o_{\max}+1}, s_{t-o_{\max}}, r_{t-o_{\max}}, a_{t-2o_{\max}}, \dots, a_{t-o_{\max}}$  from  $\mathcal{B}$ 
21:          $\hat{x}_{t-o_{\max}} \leftarrow (s_{t-2o_{\max}}, a_{t-2o_{\max}}, \dots, a_{t-o_{\max}})$ 
22:          $\hat{x}_{t-o_{\max}+1} \leftarrow (s_{t-2o_{\max}+1}, a_{t-2o_{\max}+1}, \dots, a_{t-o_{\max}+1})$ 
23:         store  $(\hat{x}_{t-o_{\max}}, s_{t-o_{\max}}, a_{t-o_{\max}}, r_{t-o_{\max}}, \hat{x}_{t-o_{\max}+1}, s_{t-o_{\max}+1})$  in  $\mathcal{D}$ 
24:         pop  $s_{t-2o_{\max}}, a_{t-2o_{\max}}$  from  $\mathcal{B}$ 
25:       end if
26:     end if
27:   end for
28:   for each gradient step do
29:      $\theta \leftarrow \theta - \lambda_Q \nabla \mathcal{J}_{Q_\beta}(\theta)$  ▷ update beta critic
30:      $\phi \leftarrow \phi - \lambda_\pi \nabla \mathcal{J}_\pi(\phi)$  ▷ update actor
31:      $\bar{\theta} \leftarrow \xi\theta + (1-\xi)\bar{\theta}$  ▷ update target beta critic
32:   end for
33: end for
34: Output: actor  $\bar{\pi}_\phi$ 

```

As discussed in Section 3.1, the augmented reward for the action with respect to the augmented state is a random variable that has to be determined based on the conditional expectation as in equation 6. Fortunately, this expected value can be empirically obtained through the use of replay buffer \mathcal{D} :

$$\bar{r}_{t-o_{\max}} = \mathbb{E}_{(s,a) \sim \mathcal{D}}[r_{t-o_{\max}}] \quad (17)$$

where $\bar{r}_{t-o_{\max}}$ and $r_{t-o_{\max}}$ represent $\bar{R}(\hat{x}_{t-o_{\max}}, a_{t-o_{\max}})$ and $R(s = s_{t-o_{\max}}, a = a_{t-o_{\max}})$, each. Consequently, training the beta-critic and actor requires only the following set of experience tuples:

$$(\hat{x}_{t-o_{\max}}, s_{t-o_{\max}}, a_{t-o_{\max}}, r_{t-o_{\max}}, \hat{x}_{t-o_{\max}+1}, s_{t-o_{\max}+1}). \quad (18)$$

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

E VISUAL REPRESENTATION OF LAZY-AGENT

In this section, we provide a visual representation of the proposed lazy-agent employed in RDMDPs, where the maximum delay is set to $o_{\max} = 3$.

$*s_a^b$: a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states										
Observed states										
Usable states										
Augmented states										
Actions										

(a) Time $t = 0$

$*s_a^b$: a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2									
Observed states										
Usable states										
Augmented states	'no-ops'									
Actions	a_1									

(b) Time $t = 1$

$*s_a^b$: a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1								
Observed states										
Usable states										
Augmented states		'no-ops'								
Actions	a_1	a_2								

(c) Time $t = 2$

Figure 12: At times 1 and 2, the states s_1^2 and s_2^1 are generated but remain unobserved by the lazy-agent due to delays. In this scenario, the lazy-agent does nothing ('no-ops') until the initial state s_1^2 becomes usable.

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
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
1237
1238
1239
1240
1241

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2							
Observed states			s_1^2 s_2^1							
Usable states										
Augmented states		'no-ops'								
Actions	a_1	a_2	a_3							

(a) Time $t = 3$

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3						
Observed states			s_1^2 s_2^1							
Usable states				s_1^2						
Augmented states		'no-ops'		\hat{x}_4						
Actions	a_1	a_2	a_3	a_4						

(b) Time $t = 4$

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1					
Observed states			s_1^2 s_2^1		s_3^2					
Usable states				s_1^2	s_2^1					
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5					
Actions	a_1	a_2	a_3	a_4	a_5					

(c) Time $t = 5$

Figure 13: At time 3, states s_1^2 and s_2^1 are observed simultaneously. As the lazy-agent uses these observed states at their maximum delayed times, s_1^2 is used at time 4 and s_2^1 is used at time 5. These states are reformulated as augmented states before being fed into the policy, thereafter determining the appropriate actions. States s_3^2 , s_4^3 , and s_5^1 are generated at corresponding times, with s_3^2 being observed at time 5.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0				
Observed states			s_1^2 s_2^1		s_3^2	s_5^1 s_6^0				
Usable states				s_1^2	s_2^1	s_3^2				
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5	\hat{x}_6				
Actions	a_1	a_2	a_3	a_4	a_5	a_6				

(a) Time $t = 6$ * s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	s_7^3			
Observed states			s_1^2 s_2^1		s_3^2	s_5^1 s_6^0	s_4^3			
Usable states				s_1^2	s_2^1	s_3^2	s_4^3			
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5	\hat{x}_6	\hat{x}_7			
Actions	a_1	a_2	a_3	a_4	a_5	a_6	a_7			

(b) Time $t = 7$ * s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	s_7^3	...		
Observed states			s_1^2 s_2^1		s_3^2	s_5^1 s_6^0	s_4^3			
Usable states				s_1^2	s_2^1	s_3^2	s_4^3	s_5^1		
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5	\hat{x}_6	\hat{x}_7	\hat{x}_8		
Actions	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8		

(c) Time $t = 8$

Figure 14: States s_6^0 and s_7^3 are generated at respective times. At time 6, states s_5^1 and s_6^0 are observed simultaneously but are not immediately usable because the previously generated states, s_3^2 and s_4^3 , have not yet been used in decision-making processes. Instead, s_3^2 is used at this time. At time 7, state s_4^3 is observed and is available for use immediately. At time 8, state s_5^1 becomes usable, as all previously generated states have now been both observed and used.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	s_7^3	
Observed states			s_1^2 s_2^1		s_3^2	s_5^1 s_6^0	s_4^3			
Usable states				s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5	\hat{x}_6	\hat{x}_7	\hat{x}_8	\hat{x}_9	
Actions	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9	

(a) Time $t = 9$

* s_a^b : a = generated time, b = delay

Times	t = 1	t = 2	t = 3	t = 4	t = 5	t = 6	t = 7	t = 8	t = 9	t = 10
Generated states	s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	s_7^3
Observed states			s_1^2 s_2^1		s_3^2	s_5^1 s_6^0	s_4^3			s_7^3
Usable states				s_1^2	s_2^1	s_3^2	s_4^3	s_5^1	s_6^0	s_7^3
Augmented states		'no-ops'		\hat{x}_4	\hat{x}_5	\hat{x}_6	\hat{x}_7	\hat{x}_8	\hat{x}_9	\hat{x}_{10}
Actions	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9	a_{10}

(b) Time $t = 10$

Figure 15: At times 9 and 10, states s_6^0 and s_7^3 are used in sequence. Despite the state observations occurring simultaneously or being out of order, all the delayed states are consistently used in sequence at their maximum delayed times, i.e., $\tau(s_n^{o_n}) = n + o_{\max}, \forall n > 0$.