

---

# ADAPTIVE REINFORCEMENT LEARNING FOR UNOBSERVABLE RANDOM DELAYS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

In standard Reinforcement Learning (RL) settings, the interaction between the agent and the environment is typically modeled as a Markov Decision Process (MDP), which assumes that the agent observes the system state instantaneously, selects an action without delay, and executes it immediately. In real-world dynamic environments, such as cyber-physical systems, this assumption often breaks down due to delays in the interaction between the agent and the system. These delays can vary stochastically over time and are typically *unobservable* when deciding on an action. Existing methods deal with this uncertainty conservatively by assuming a known fixed upper bound on the delay, even if the delay is often much lower. In this work, we introduce the *interaction layer*, a general framework that enables agents to adaptively handle unobservable and time-varying delays. Specifically, the agent generates a matrix of possible future actions, anticipating a horizon of potential delays, to handle both unpredictable delays and lost action packets sent over networks. Building on this framework, we develop a model-based algorithm, *Actor-Critic with Delay Adaptation (ACDA)*, which dynamically adjusts to delay patterns. Our method significantly outperforms state-of-the-art approaches across a wide range of locomotion benchmark environments.

## 1 INTRODUCTION

State-of-the-art reinforcement learning (RL) algorithms, such as Proximal Policy Optimization (PPO) (Schulman et al., 2017) and Soft Actor-Critic (SAC) (Haarnoja et al., 2018), are typically built on the assumption that the environment can be modeled as a Markov Decision Process (MDP). This framework implicitly assumes that the agent observes the current state instantaneously, selects an action without delay, and executes it immediately.

This assumption often breaks down in real-world systems due to *interaction* delays that arise from various sources: the time taken to collect and transmit observations, the computation time needed for the agent to select an action, and the transmission and actuation delay when executing that action in the environment (as illustrated in Figure 1). Delays pose no issue if the state of the environment is not evolving between its observation and the execution of the selected action. But in continuously evolving systems—such as robots operating in the physical world—the environment’s state may have changed by the time the action is executed (Brooks & Leondes, 1972). Delays have been recognized as a key concern when applying RL to cyber-physical systems (Tan et al., 2018). Outside the scope of RL, delays have also been studied in classic control (Ray, 1988; Luck & Ray, 1990).

These interaction delays can be implicitly modeled by altering the transition dynamics of the MDP to form a partially observable Markov decision process (POMDP), in which the agent only receives outdated sensor observations. While this approach is practical and straightforward, it limits the agent’s access to information about the environment’s evolution during the delay period.

Another common approach to handling delays in RL is to enforce that actions are executed after a fixed delay (Katsikopoulos & Engelbrecht, 2003; Walsh et al., 2008). This is typically implemented by introducing an action buffer between the agent and the environment, ensuring that all actions are executed after a predefined delay. However, this method requires prior knowledge of the maximum possible delay and enforces that all actions incur this worst-case delay—even when most interactions in practice experience minimal or no delay. The advantage of this fixed-delay approach is that it provides the agent with perfect information about when its actions will take effect, simplifying

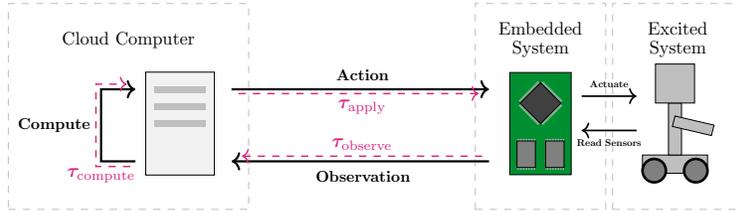


Figure 1: Illustration of a setup affected by interaction delays. Any delay between the embedded system and the excited system is considered negligible or otherwise accounted for (see Appendix F.1 for a detailed discussion). The factors contributing to interaction delay are  $\tau_{\text{observe}}$  ( $\tau_o$ ),  $\tau_{\text{compute}}$  ( $\tau_c$ ), and  $\tau_{\text{apply}}$  ( $\tau_a$ ). See Section 3.1 for more details about these factors.

decision-making. However, it is overly conservative and fails to adapt and account for variability in delay. Note that state-of-the-art algorithms for delayed MDPs, such as BPQL (Kim et al., 2023) and VDPO (Wu et al., 2024), rely on this fixed-delay paradigm.

Moving beyond this fixed-delay framework is challenging, especially because in real-world systems, delays are often unobservable. The agent does not know, at decision time, how long it will take for an action to be executed. One existing approach that attempts to address varying delays is DCAC (Bouteiller et al., 2021), but it does not offer any guarantees for when a generated action will be applied to the environment.

In this paper, we make the following contributions:

(i) We introduce a novel framework, the *interaction layer*, which allows agents to adapt to randomly varying delays—even when these delays are unobservable. In this setup, the agent generates a matrix of candidate actions ahead of time, each row in the matrix intended for a possible future arrival time (without knowing for certain which row will be selected). Specifically, the design handles both (a) that the future actions can have varying delays, and (b) that action packets sent over a network can be lost or arrive in incorrect order. The actual action is selected at the interaction layer once the delay is revealed. Similar to DCAC, we also report back the revealed delays in hindsight. This approach enables informed decision-making under uncertainty and robust behavior in the presence of stochastic, unobservable delays (Section 3).

(ii) We develop a new model-based reinforcement learning algorithm, *Actor-Critic with Delay Adaptation (ACDA)*, which leverages the interaction layer to adapt dynamically to varying delays. The algorithm provides two key concepts: (a) instead of using states as input to the policy, it uses a distribution of states as an embedding that enables the generation of more accurate time series of actions, and (b) an efficient heuristic to determine which of the previously generated actions are executed. These actions are needed to compute the state distributions. The approach is particularly efficient when delays are temporally correlated, something often seen in scenarios when communicating over transmission channels (Section 4).

(iii) We evaluate ACDA on a suite of MuJoCo locomotion tasks from the Gymnasium library (Towers et al., 2024), using randomly sampled delay processes designed to mimic real-world latency sources. Our results show that ACDA, equipped with the interaction layer, consistently outperforms state-of-the-art algorithms designed for fixed delays and for unobservable random delays. It achieves higher average returns across all benchmarks except one, where its performance remains within the standard deviation of the best constant-delay method (Section 5).

## 2 RELATED WORK

To our knowledge, there is no previous work that allows agents to make informed and controlled decisions under random unobservable delays in RL. Much of the existing work on how to handle delays in RL acts as if delays are constant equal to  $h$ , in which case, the problem can be modeled as an MDP with augmented state  $(s_t, a_t, a_{t+1}, \dots, a_{t+h-1})$  consisting of the last observed state and memorized actions to be applied in the future Katsikopoulos & Engelbrecht (2003). Even if the true delay is not constant, a construction used in previous work is to enforce constant interaction delay through *action buffering*, under the assumption that the maximum delay does not exceed  $h$  time-steps.

Through action buffering and state augmentation, one may, in principle, use existing RL techniques to deal with constant delays. However, it is hard to directly learn policies on augmented states in practice, which has prompted the development of algorithms that exploit the delayed dynamics. The real-time actor-critic by Ramstedt & Pal (2019) optimizes for a constant delay of one time step. Belief projection-based Q-learning (BPQL) by Kim et al. (2023) explicitly uses the delayed dynamics under constant delay to simplify the critic learning. BPQL achieves good performance during evaluation over longer delays, despite a simple structure of the learned functions. Our algorithm in Section 4.3 uses the same critic simplification, but applied to the randomly delayed setting.

Another approach explored for constant-delay RL is to have a delayed agent trying to imitate an undelayed expert, used in algorithms such as DIDA (Liotet et al., 2022) and VDPO (Wu et al., 2024). These assume access to the undelayed MDP, which in the real world can be applied in sim-to-real scenarios, but not when training directly on the real physical system.

DCAC by Bouteiller et al. (2021) is a framework that allows agents to make decisions under unobservable delays, but without any control over when an action is going to be applied to the environment. Like our approach, delays are available in hindsight, which DCAC uses for future decision making and value accreditation. Other approaches for random delays typically assume observability (Valensi et al., 2024; Wu et al., 2025), which is not applicable in our problem setting.

Model-based approaches have also been explored for delayed RL, as a way to plan into future horizons (Chen et al., 2021) or to estimate future states as policy inputs (Walsh et al., 2008; Firoiu et al., 2018). A commonly used dynamics model architecture is the recurrent state space model (RSSM) (Hafner et al., 2019) that combines a recurrent latent state with stochastically sampled states to transition in latent space. RSSM was designed for planning algorithms, but can also be used for state prediction. The model used by Firoiu et al. (2018) is similar to RSSM, but uses deterministic output of states from the latent representation. Another approach using RSSM is Dreamer (Hafner et al., 2020) that learns a latent state representation for the agent to make decisions in, originally in an undelayed setting but extended to the delayed setting by Karamzade et al. (2024). Wang et al. (2024) have explored further variations in model structures that can be used for delayed RL.

Our approach also learns a model to make decisions in latent space, but does not follow the RSSM structure. Instead, our model (introduced in Section 4.2) follows a simpler structure that learns a latent representation describing actual state distributions rather than uncertainty about an assumed existing true state. By the definitions of Moerland et al. (2023), our model is classified as a multi-step prediction model with state abstraction, even though we are only estimating distributions.

### 3 THE INTERACTION LAYER

In this section, we explain how random and unpredictable delays may affect the interaction between the agent and the system. To handle these delays, we introduce a new framework, called the *interaction layer* (Section 3.2), and model the way the agent and the system interact by a Partially Observable MDP (Section 3.3). The notation used for the interaction layer is explained as it appears in the text. See Appendix C for a more compact, formal description of the interaction layer.

#### 3.1 DELAYED MARKOV DECISION PROCESSES

We consider a controlled dynamical system modeled as an MDP  $\mathcal{M} = (S, A, p, r, \mu)$ , where  $S$  and  $A$  are the state and action spaces, respectively, where  $p(s'|s, a)_{(s', s, a) \in S \times S \times A}$  represents the system dynamics in the absence of delays,  $r$  is the reward function, and  $\mu$  is the distribution of the initial state.

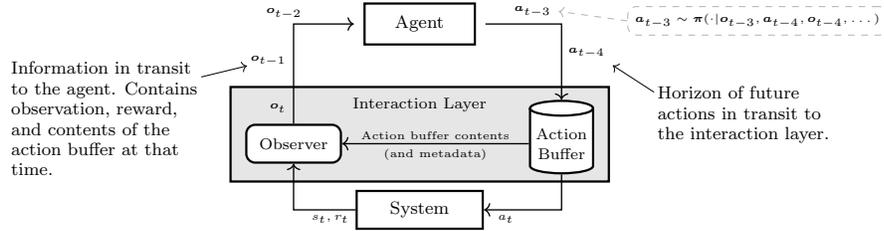
As in usual MDPs, we assume that at the beginning of each step  $t$ , the state of the system is sampled, but this information does not reach the agent immediately, but after an observation delay,  $\tau_o$ . After the agent receives the information, an additional computational delay,  $\tau_c$ , occurs due to processing the information and deciding on an appropriate action. The action created by the agent is then communicated to the system, with an additional final delay  $\tau_a$  before this action can be applied to the system. The delays<sup>1</sup>  $(\tau_o, \tau_c, \tau_a)$  are random variables that may differ across steps and can be

<sup>1</sup>We assume that The total delay  $\tau_o + \tau_c + \tau_a$  is measured in number of steps. In Appendix G.1, we present a similar model where delays can take any real positive value.

162 correlated. While it is possible to consider frameworks where  $\tau_o$  and  $\tau_c$  are observable, the action  
 163 delay  $\tau_a$  is inherently unobservable, as this delay may be caused by events taking place after the  
 164 action has been generated. Therefore, to simplify the problem in our framework, we consider the  
 165 sum of these three delays as a single delay  $d_t$ , which is unobservable. This single delay represents  
 166 the full round-trip delay of observing, computing an action, and applying the action to the system.  
 167 This is further explained in Section 3.3.

### 169 3.2 HANDLING DELAYS VIA THE INTERACTION LAYER

170  
 171 The unpredictable delays pose significant challenges from the agent’s perspective. First, the agent  
 172 cannot respond immediately to the newly observed system state at each step. Second, the agent  
 173 cannot determine when the selected action will be applied to the system. To address these issues,  
 174 we introduce the *interaction layer*, consisting of an *observer* and an *action buffer*, as illustrated in  
 175 Figure 2. The interaction layer is a direct part of the system that performs sensing and actuation,  
 176 whereas the agent can be far away, communicating over a network. Within the interaction layer,  
 177 the observer is responsible for sampling the system’s state and sending relevant information to the  
 178 agent. The agent generates a matrix of possible actions. These are sent back to the interaction layer  
 179 and stored in the action buffer. Depending on when the actions arrive in the action buffer, it selects  
 180 a row of actions, which are then executed in the following steps if no further decision is received.  
 181 The rest of this section gives technical details of the interaction layer, whereas Section 4 details the  
 182 policy for generating actions at the agent.



191 Figure 2: Illustration of the interaction layer and how the agent interacts with it from a global per-  
 192 spective. As the observation is received from the dynamical system, the next action is immediately  
 193 applied from the action buffer. Packets in transit with random delay imply partial observability.

194  
 195 **Action packet.** After that, the agent receives an observation packet  $o_t$  (generated at step  $t$  by the  
 196 interaction layer, described further below), the agent generates and sends an action packet  $a_t$ . The  
 197 packet includes a time stamp  $t$ , and a matrix of actions, as follows:

$$198 \mathbf{a}_t = \left( t, \begin{bmatrix} a_1^{t+1} & a_2^{t+1} & a_3^{t+1} & \dots & a_h^{t+1} \\ a_1^{t+2} & a_2^{t+2} & a_3^{t+2} & \dots & a_h^{t+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_1^{t+L} & a_2^{t+L} & a_3^{t+L} & \dots & a_h^{t+L} \end{bmatrix} \right). \quad (1)$$

204 The  $i$ -th row of the matrix of the action packet corresponds to the sequence of actions that would  
 205 constitute the action buffer if the packet reaches the interaction layer at time  $t + i$ . The reason for  
 206 using a matrix instead of a vector is that subsequent columns specify which actions to take if a new  
 207 action packet does not arrive at the interaction layer at a specific time step. For instance, if an action  
 208 packet arrives at time  $t + 2$ , then the interaction layer uses the first action in the buffer ( $a_1^{t+2}$   
 209 in this case). That is, the first column is always used when a new packet arrives at each time step. If  
 210 no packet arrives for a specific time step, the other columns are used instead (as explained more  
 211 below in the description of the action buffer). This approach enables adaptivity for the agent: it can  
 212 generate actions for specific delays without knowing what the delay is going to be ahead of time.  
 213 Figure 3 illustrates when an action packet arrives at the interaction layer and a row is inserted into  
 214 the action buffer (3rd row in this case because the packet arrived with a delay of 3).

215 The capacity of the action buffer is denoted by  $h$ , the horizon of actions to cover for gaps in the  
 interaction. The number of rows in the matrix, denoted by  $L$  (prediction length), is determined by

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

the agent. If the delay associated with the action packet exceeds the number of rows  $L$  in the matrix, that action packet is discarded. Additionally, if an action packet arrives out of order, where  $\mathbf{a}_t$  arrives after  $\mathbf{a}_{t'}$  and  $t < t'$ ,  $\mathbf{a}_t$  is discarded. This process is formally described in Appendix C.

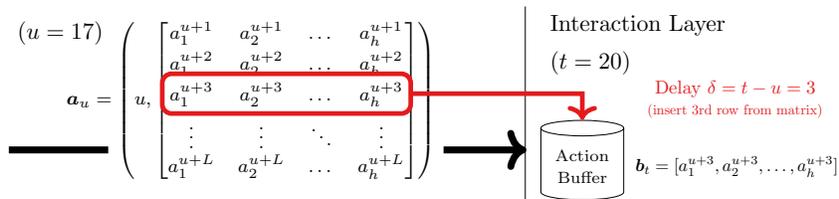


Figure 3: Example (see below for definitions of recovered delay  $\delta_t$  and shift counter  $c_t$ ): Suppose an action packet timestamped by the agent with time  $u = 17$ ,  $\mathbf{a}_u$ , arrives at the action layer at time 20. Then, at time  $t = 20$ ,  $\delta_{20} = 3$ , and  $c_{20} = 0$ . Now, suppose that 2 time units elapse without any new action packet arriving. Then, at time  $t = 22$ ,  $\delta_{22} = 3$  and counter  $c_{22} = 2$ . Hence, equation  $t = u + \delta_{22} + c_{22} = 17 + 3 + 2 = 22$  holds.

While this may appear as if action delays are observable, the action packet only allows us to specify what should happen if it arrives with a certain delay. If the action delay truly was observable, we could use information about delays for previous action packets to get perfect information about which actions will be applied to the underlying state prior to this action packet arriving.

**Action buffer.** The action buffer is responsible for executing an action each time step. If no new action arrives at a time step, the next item in the buffer is used. At the beginning of step  $t$ , the action buffer contains the following information:  $\mathbf{b}_t$ , a sequence of  $h$  actions to be executed next, and  $\delta_t$ , the delay of the action packet from which the actions  $\mathbf{b}_t$  were taken. For instance, if an action packet  $\mathbf{a}_u$  arrives at the action buffer at time  $t$ , then  $\delta_t = t - u$ , where  $u$  is the time stamp of the action packet  $\mathbf{a}_u$  that the agent created. If instead no new action packet arrived at time  $t$ , then  $\delta_t = \delta_{t-1}$ . To enable the use of an appropriate action even if no new packet arrives at a specific time step, the content of the buffer is shifted one step forward, as shown in Figure 4. Finally, the action buffer includes a counter  $c_t$  that records how many steps have passed since the action buffer was updated. The following invariant always holds:  $t = u + \delta_t + c_t$ . For a concrete example, see the caption of Figure 3.

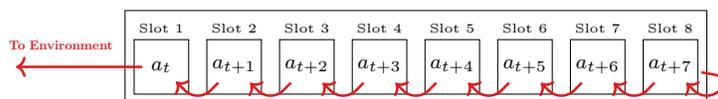


Figure 4: Action buffer shifting actions. Final slot is repeated. (Example: horizon  $h = 8$ )

**Observation packet.** The observer builds an observation packet  $\mathbf{o}_t$  at the beginning of step  $t$ . To this aim, it samples the system state  $s_t$ , collects information  $\mathbf{b}_t, \delta_t, c_t$  about the action buffer, forms the observation packet  $\mathbf{o}_t = (t, s_t, \mathbf{b}_t, \delta_t, c_t)$ , and sends it to the agent.

Enhancing the information contained in the observation and action packets (compared to the undelayed MDP scenario) allows the agent to make more informed decisions and ensures the system does not run out of actions when action packets experience delays. However, this is insufficient to model our delayed system as an MDP. This is because the agent does not have the knowledge of all the observation and action packets currently in transit. Therefore, we use the formalism of a POMDP to accurately describe the system dynamics.

### 3.3 THE POMDP MODEL

Next, we complete the description of our delayed MDP and model it as a POMDP. To this aim, we remark that the system essentially behaves as if in each step  $t$ , the agent immediately observes  $\mathbf{o}_t$  and selects an action packet  $\mathbf{a}_t$  that arrives at the interaction layer  $d_t$  steps after the observation

$\mathbf{o}_t$  was made, where  $d_t > 0^2$ . We assume that  $d_t$  is generated according to some distribution  $D^3$ . Furthermore, we assume that observation packets  $\mathbf{o}_t$  arrive in order at the agent. In this framework—where the agent selects an action packet  $\mathbf{a}_t$  as soon as the observation  $\mathbf{o}_t$  is generated—the single round-trip delay  $d_t$  represents the three delays ( $\tau_o, \tau_c, \tau_a$ ) as  $d_t = \tau_o + \tau_c + \tau_a$ .

The time step  $t$  is a local time tag from the perspective of the interaction layer. Our POMDP formulation does not assume a synchronized clock between the agent and the interaction layer. The agent acts asynchronously and generates an action packet upon receiving an observation packet.

We define  $\mathcal{I}_t$  as the set of action packets in transit at the beginning of step  $t$ , along with the times at which these packets will arrive at the interaction layer ( $\mathcal{I}_t$  is a set containing items on the form  $(u + d_u, \mathbf{a}_u)$ ). In reality, delays are observed only when action packets reach the interaction layer, and the agent does not necessarily know whether the action packets already generated have reached the interaction layer. Hence, we must assume that  $\mathcal{I}_t$  is not observable by the agent. The framework we just described corresponds to a POMDP, which we formalize in Appendix C in detail.

## 4 ACTOR-CRITIC WITH DELAY ADAPTATION

This section introduces *actor-critic with delay adaptation* (ACDA), a model-based RL algorithm using the interaction layer to adapt on-the-fly to varying unobservable delays, contrasting with state-of-the-art methods that enforce a fixed worst-case delay. A challenge with varying unobservable delays is that the agent lacks perfect information about the actions to be applied in the future. ACDA solves this with a heuristic (Section 4.1) that is effective when delays are temporally correlated.

The actions selected by ACDA will vary in length depending on the delay we are generating actions for. This lends itself poorly to commonly used policy function approximators in deep RL, such as multi-layer perceptrons (MLPs), that assume a fixed size of input. ACDA solves this with a model-based distribution agent (Section 4.2) that embeds the variable-length input into fixed-size embeddings of future state distributions, to which the generated action will be applied. The fixed-size embeddings are fed as input to an MLP to generate actions. ACDA learns a model of the environment dynamics online to compute these embeddings. Section 4.3 shows how we train ACDA.

### 4.1 HEURISTIC FOR ASSUMED PREVIOUS ACTIONS

A problem with unobservable delays is that we do not know when our previously sent action packets will arrive at the interaction layer. This means that we do not know which actions are going to be applied to the underlying system between generating the action packet and it arriving at the interaction layer. A naive assumption would be to assume the action buffer contents reported by the observation packet to be the actions that are going to be applied to the underlying system. However, this is unlikely to be true because the action buffer is going to be preempted by action packets already in transit.

ACDA employs a heuristic for estimating these previous actions to be applied to the system between  $\mathbf{o}_t$  being generated and  $\mathbf{a}_t$  arriving at the interaction layer. The heuristic assumes that, if  $\mathbf{a}_t$  arrives at time  $t + k$  (it having delay  $k$ ), then previous action packets will also have delay  $k$ . Such that  $\mathbf{a}_{t-1}$  will arrive at time  $t + k - 1$ ,  $\mathbf{a}_{t-2}$  at  $t + k - 2$ , etc.

Under this assumption, a new action packet will preempt the action buffer at every single time step. This means that, if we assume a delay of  $k$ , the action applied to the underlying system will be the action in the first column of the  $k$ -th row in the action packet last received by the interaction layer. By memorizing the action packets

---

#### Algorithm 1 Memorized Action Selection

---

**Input**  $k \in \mathbb{Z}^+$  (Delay assumption)  
 $\mathbf{a}_{t-1}, \mathbf{a}_{t-2}, \dots, \mathbf{a}_{t-k}$  (Memorized Packets)

- 1: **for**  $i \leftarrow 1$  to  $k$  **do**
- 2:      $(t - i, M^{t-i}) = \mathbf{a}_{t-i}$
- 3:      $\triangleright$  Unpacking action matrix  $M$  from packets
- 4: **return**  $(\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k}) = (M_{k,1}^{t-k}, \dots, M_{k,1}^{t-1})$

---

<sup>2</sup> $d_t$  corresponds to the value  $\delta_u$  if the action packet reaches the interaction layer at time  $u$ :  $\delta_{t+d_t} = d_t$

<sup>3</sup>For simplicity, we assume that the delay process is markovian and independent of the contents of the action packets and the state of the interaction layer. However, our POMDP formalism can be extended to delay distributions that depend on the previous state, which is more general and realistic.

previously sent, we can under this assumption select the actions that are going to be applied to the system as shown in Algorithm 1. When generating  $a_1^{t+k}$ , the first action on the  $k$ -th row in the action packet  $\mathbf{a}_t$ , we use Algorithm 1 to determine the actions  $(\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$  that will be applied to the observed state  $s_t$  before  $a_1^{t+k}$  is executed. For the action  $a_2^{t+k}$ , we know that this is only going to be executed if no new action packet arrived at  $t+k+1$ . We therefore extend the previous assumption and say that  $(\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k}, a_1^{t+k})$  are the actions applied to  $s_t$  before  $a_2^{t+k}$  is executed.

The main idea here is that the heuristic guesses the applied actions if the delay does not evolve too much over time. If the delay truly was constant, then all guesses would be accurate and ACDA would transform the POMDP problem to a constant-delay MDP. The heuristic’s accuracy is compromised during sudden changes in delay, such as network delay spikes. However, as we will see in the evaluation, occasional violations will not significantly impact overall performance.

## 4.2 MODEL-BASED DISTRIBUTION AGENT

The memorized actions used by ACDA are variable in length and therefore cannot be directly used as input to MLPs, which are often used in constant-delay approaches. Instead, ACDA constructs an embedding  $z_1^{t+k}$  of the distribution  $p(s_{t+k}|s_t, \hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$ , where  $\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k}$  are the memorized actions. We then provide  $z_1^{t+k}$  as input to an MLP to generate  $a_1^{t+k}$ . This allows the policy to reason about the possible states in which the generated action will be executed. Note that we are only concerned with the distribution itself and never explicitly sample from it. To compute these embeddings, we learn a model of the system dynamics using three components:  $\text{EMBED}_\omega$ ,  $\text{STEP}_\omega$ , and  $\text{EMIT}_\omega$ , where  $\omega$  represents learnable parameters.

- $\hat{z}_0 = \text{EMBED}_\omega(s_t)$  embeds a state  $s_t$  into a distribution embedding  $\hat{z}_0$ .
- $\hat{z}_{i+1} = \text{STEP}_\omega(\hat{z}_i, a_{t+i})$  updates the embedded distribution to consider what happens after also applying the action  $a_{t+i}$ . Such that if  $\hat{z}_i$  is an embedding of  $p(s_{t+i}|s_t, a_t, \dots, a_{t+i})$ , then  $\hat{z}_{i+1}$  is an embedding of  $p(s_{t+i+1}|s_t, a_t, \dots, a_{t+i}, a_{t+i+1})$ .
- The final component  $\text{EMIT}_\omega(s_{t+i}|\hat{z}_i)$  allows for a state to be sampled from the embedded distribution. This component is not used when generating actions, and is instead only used during training to ensure that  $\hat{z}_i$  is a good embedding of  $p(s_{t+i}|s_t, a_t, \dots, a_{t+i})$ .

The way these components are used to produce the embedding  $z_1^{t+k}$  is illustrated in Figure 5. We use the notation  $z_1^{t+k} = \hat{z}_k$  given that we are embedding the selected actions  $(\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$ . We use the notation  $\text{STEP}_\omega^k(\text{EMBED}_\omega(s_t), \hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$  to describe this multi-step embedding process. This notation is formalized in Appendix D.

The  $\text{EMBED}_\omega$  and  $\text{EMIT}_\omega$  components are implemented as MLPs, while  $\text{STEP}_\omega$  is implemented as a gated recurrent unit (GRU). We provide detailed descriptions of these components in Appendix D. We learn these components online by collecting information from trajectories about observed states  $s_t$  and  $s_{t+n}$  and their interleaved actions  $a_t, a_{t+1}, \dots, a_{t+n-1}$  in a replay buffer  $\mathcal{R}$ . The following loss function  $\mathcal{L}(\omega)$  is used to minimize the KL-divergence between the model and the underlying system dynamics:  $\mathcal{L}(\omega) = \mathbb{E}_{(s_t, a_t, a_{t+1}, \dots, a_{t+n-1}, s_{t+n}) \sim \mathcal{R}} [-\log \text{EMIT}_\omega(s_{t+n}|z_n)]$  where  $z_n = \text{STEP}_\omega^n(\text{EMBED}_\omega(s_t), a_t, a_{t+1}, \dots, a_{t+n-1})$ .

Given the embedding  $z_1^{t+k}$ , we produce  $a_1^{t+k}$  in the action packet  $\mathbf{a}_t$  using a policy  $\pi_\theta(a_1^{t+k}|z_1^{t+k})$ , i.e., generating actions given the (embedded) distribution over the state that the action will be ap-

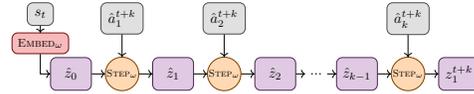


Figure 5: Illustration of the multi-step distribution model embedding  $p(s_{t+k}|s_t, \hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$  as  $\text{STEP}_\omega^k(\text{EMBED}_\omega(s_t), \hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$ .

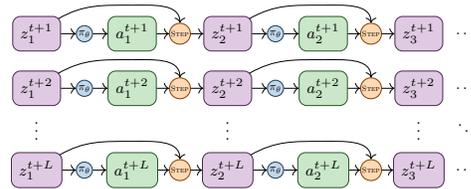


Figure 6: Generating the action packet from the embeddings. Each row in the figure corresponds to a row in the matrix of the action packet  $\mathbf{a}_t$ .

378 plied to. This policy structure allows the agent to reason about uncertainties in future states when  
 379 generating actions.

380  
 381 By extending the assumptions as shown in Section 4.1, we can also produce the embeddings  $z_2^{t+k}$ ,  
 382  $z_3^{t+k}$ , etc., as illustrated in Figure 6. This process of generating the matrix rows is similar to action  
 383 chunking Lai et al. (2022); Zhao et al. (2023); Li et al. (2025), though we use them to cover gaps in  
 384 the interaction, rather than with the expectation that they will all be executed. The complete process  
 385 of constructing the action packet is formalized in Appendix D. We also discuss the effect that this  
 386 has on the computational delay in Appendix F, why it is not a problem in our case, and how to  
 387 handle it if it should become a problem.

388 This model-based policy can also be applied in the constant-delay setting to achieve decent perfor-  
 389 mance. We evaluate how this compares against a direct MLP function approximator in Appendix  
 390 E.2, where the model-based policy is implemented in the BPQL algorithm.

### 391 4.3 TRAINING ALGORITHM

392 This section describes the training procedure in Algorithm 2, used to optimize the parameters of the  
 393 networks. It follows an actor-critic setup based on SAC. The training procedure of the critic  $Q_\phi$  is  
 394 similar to BPQL, where  $Q_\phi(s, a)$  evaluates the value of actions  $a$  on undelayed system states  $s$ .

395  
 396 Algorithm 2 is split into three parts: trajectory sampling (L3-L12), transition reconstruction (L13),  
 397 and training (L14-L15). We do this split to reduce the impact that the training procedure can have  
 398 on the computational delay  $\tau_c$  of the system. From the trajectory sampling, we collect POMDP  
 399 transition information  $(\mathbf{o}_t, \mathbf{a}_t, r_t, \mathbf{o}_{t+1})$  where  $\Gamma_t$  is used to discern if  $s_t$  is in a terminal state.

401 An important aspect of Algorithm 2 is how  
 402 trajectory information is reconstructed for  
 403 training. Specifically, we reconstruct the  
 404 trajectory  $(s_0, a_0, r_0, s_1, a_1, \dots)$  from the  
 405 perspective of the undelayed MDP, along  
 406 with the policy input used to generate each  
 407 action  $a_t$ . The policy input can be re-  
 408 trospectively recovered by examining the  
 409 current buffer action delay  $(\delta_t)$  and the  
 410 number of times the buffer has shifted  $(c_t)$ .  
 411 This trajectory reconstruction is necessary  
 412 since we follow the BPQL algorithm’s  
 413 actor-critic setup. The critic  $Q_\phi(s_t, a_t)$  es-  
 414 timates values in the undelayed MDP, and  
 415 we need to be able to regenerate actions  $a_t$   
 416 using the model-based policy to compute  
 417 the TD-error. Further details are provided  
 418 in Appendix D. The hyperparameters used  
 419 for ACDA, including the prediction length  $L$ , are located in Appendix B.4.

---

#### Algorithm 2 Actor-Critic with Delay Adaptation

---

- 1: Init. policy  $\pi_\theta$ , critic  $Q_\phi$ , model  $\omega$ , and replay  $\mathcal{R}$
  - 2: **for** each epoch **do**
  - 3:   Reset interaction layer state:  $s_0 \sim \mu, t = 0$
  - 4:   Collected trajectory:  $\mathcal{T} = \emptyset$
  - 5:   Observe  $\mathbf{o}_0$
  - 6:   **while** terminal state not reached **do**
  - 7:     **for**  $k \leftarrow 1$  to  $L$  **do**
  - 8:       Select  $\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k}$  by Alg. 1
  - 9:       Create the  $k$ -th row of  $\mathbf{a}_t$
  - 10:      Send  $\mathbf{a}_t$ , observe  $r_t, \mathbf{o}_{t+1}, \Gamma_{t+1}$
  - 11:      Add  $(\mathbf{o}_t, \mathbf{a}_t, r_t, \mathbf{o}_{t+1}, \Gamma_{t+1})$  to  $\mathcal{T}$
  - 12:       $t \leftarrow t + 1$
  - 13:    Reconstruct transition info from  $\mathcal{T}$ , add to  $\mathcal{R}$
  - 14:    **for**  $|\mathcal{T}|$  sampled batches from  $\mathcal{R}$  **do**
  - 15:      Update  $\pi_\theta, Q_\phi$  and  $\omega$  (by  $\mathcal{L}(\omega)$ )
- 

## 420 5 EVALUATION AND RESULTS

421  
 422 To assess the benefits of the interaction layer in a delayed setting, we simulate the POMDP described  
 423 in Section 3.3, wrapping existing environments from the Gymnasium library (Towers et al., 2024) as  
 424 the underlying system. Specifically, we aim to answer the question of whether our ACDA algorithm,  
 425 which uses information from the interaction layer, can outperform state-of-the-art algorithms under  
 426 random delay processes.

427 We evaluate on the three delay processes shown in Figure 7. The first two delay processes  $GE_{1,23}$  and  
 428  $GE_{4,32}$  follow Gilbert-Elliott models (Gilbert, 1960; Elliott, 1963) where the delay alternates between  
 429 good and bad states (e.g. a network or computational node being overloaded or having packets  
 430 dropped). The third delay process MM1 is modeled after an M/M/1 queue (Kleinrock, 1975), where  
 431 the sampled delay is the time spent in the queue by a network packet. The full definition of these  
 delay processes is located in Appendix B.3. We expect ACDA to perform well under the Gilbert-

Elliot processes that match the temporal assumptions of ACDA. In contrast, we anticipate a worse relative performance of ACDA compared to other baselines under M/M/1 queue delays that fluctuate more.

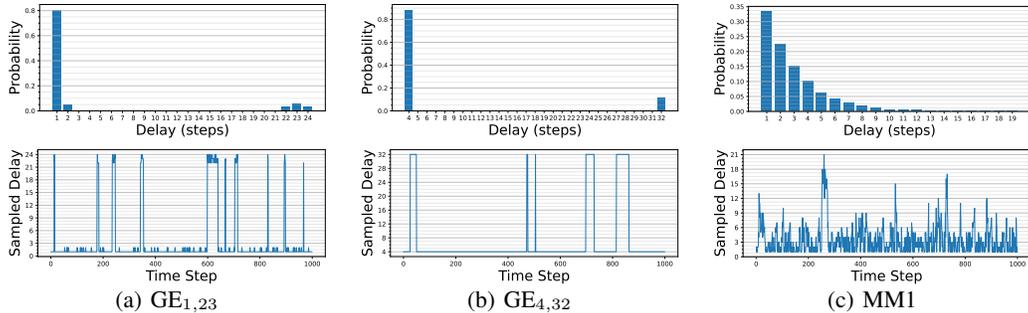


Figure 7: Evaluated delay processes, as a distribution histogram (above) and as a time series sampled delay (below) for each delay process. See Appendix B.3 for their definitions.

The state-of-the-art algorithms we compare against are DCAC (Bouteiller et al., 2021), BPQL (Kim et al., 2023), and VDPO (Wu et al., 2024). As BPQL and VDPO are designed to operate under constant delay, we apply a *constant-delay augmentation* (CDA) to allow them to operate with constant delay in random delay processes. CDA converts the interaction layer POMDP into a constant-delay MDP by making agents act under the worst-case delay of a delay process<sup>4</sup>, achieving the same effect as common constant-delay action buffers, as detailed in Appendix A. In addition to the state-of-the-art algorithms, we also evaluate the performance of SAC, both with CDA and when it acts directly on the state from the observation packet (implicitly modeling delays). In Appendix E.3, we evaluate when CDA uses a much lower assumed worst-case delay that holds most of the time, but is occasionally violated. Although this lower assumed delay yields increased performance for in some benchmarks, for other benchmarks the constant-delay algorithms performs worse or inhibit unexpected behavior, which is why we use the conservative worst-case for the evaluation here. We also evaluate the performance of Dreamer when implicitly modeling delays (Karamzade et al., 2024). Further details regarding the evaluation are presented in Appendix B.5.

We evaluate average return over a training period of 1 million steps on MuJoCo environments in Gymnasium, following the procedure from related work in delayed RL. However, an issue with the MuJoCo environments is that they have deterministic transition dynamics, rendering them theoretically unaffected by delay. To better evaluate the effect of delay, we make the transitions stochastic by imposing a 5% noise on the actions. We motivate and specify this in Appendix B.1.

The average return is computed every 10000 steps by freezing the training weights and sampling 10 trajectories under the current policy. We report the best achieved average return—where the return is the sum of rewards over a trajectory—for each algorithm, environment, and delay process in Table 1. All achieved average returns are also presented in Appendix E.1 as time series plots together with tables showing the standard deviation.

Table 1: Best evaluated average return for each algorithm.

Gymnasium env.	Ant-v4			Humanoid-v4			HalfCheetah-v4			Hopper-v4			Walker2d-v4		
	GE <sub>1,23</sub>	GE <sub>4,32</sub>	MM1												
SAC	14.22	-5.72	-0.58	862.18	494.43	921.04	2064.18	-158.78	20.69	306.91	279.74	333.06	708.33	60.86	604.80
SAC w/ CDA	69.28	18.93	102.00	414.05	230.45	613.03	128.47	591.32	550.84	426.92	315.47	627.59	428.44	257.18	2005.76
Dreamer	1111.73	1147.56	1121.11	1463.07	1091.48	981.38	1796.07	2493.19	584.40	334.30	515.36	975.72	1081.12	1233.79	1801.81
BPQL	2691.88	2509.52	<b>3074.17</b>	585.19	276.63	5435.29	4320.20	2136.36	4660.93	1328.71	433.29	3035.66	1215.91	875.09	3547.73
VDPO	2163.00	2266.99	2528.67	417.25	280.72	720.73	3144.23	3664.30	3831.96	709.20	330.44	1459.88	846.88	344.73	2144.25
DCAC	949.97	953.14	959.23	128.47	167.97	525.85	920.09	1123.47	35.60	16.99	57.98	1026.45	106.70	9.23	24.48
ACDA	<b>4112.78</b>	<b>2866.93</b>	2898.46	<b>4608.76</b>	<b>3725.59</b>	<b>5805.60</b>	<b>5984.25</b>	<b>4231.15</b>	<b>5898.36</b>	<b>2094.65</b>	<b>1727.79</b>	<b>3122.53</b>	<b>3863.59</b>	<b>1840.58</b>	<b>4562.33</b>

As shown in Table 1, ACDA outperforms state-of-the-art in all benchmarks except one, with a significant margin in most cases. The improvement is less substantial in Ant-v4, where performance often overlaps, as indicated by the standard deviation (Figure 8).

<sup>4</sup>The MM1 delay process does not have a maximum delay. We use a reasonable worst-case of 16 steps.

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

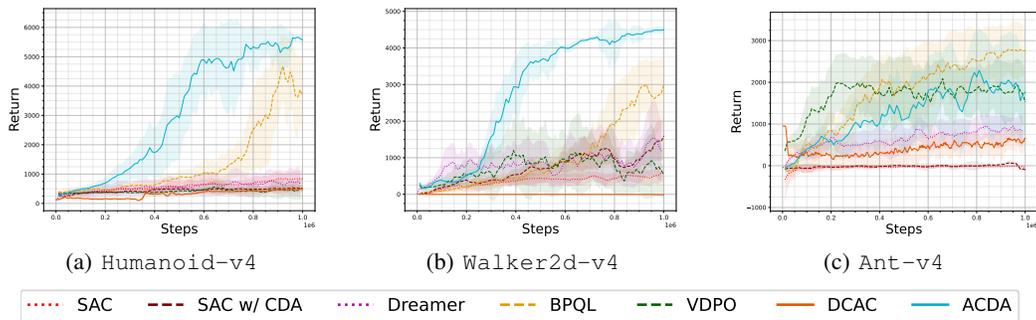


Figure 8: M/M/1 Queue results (all results in Appendix E.1.3). Shaded regions showing std. dev.

Lowering the assumed upper bound delay for constant-delay baselines can also yield better performance for some environments. We present results under lower (incorrect) assumed worst-case delay in Appendix E.3, where constant-delay baselines under  $GE_{4,32}$  perform better than ACDA in 3 of the environments. However, across all assumed delays for constant-delay baselines, whose results are presented in Appendix E.5, ACDA is still the best performing algorithm in 10 out of 15 benchmarks.

## 6 CONCLUSION

We introduced the interaction layer, a real-world viable POMDP framework for RL with random unobservable delays. Using the interaction layer, we described and implemented ACDA, a model-based algorithm that significantly outperforms state-of-the-art in delayed RL under random delay processes. Directions of future work include algorithms that can operate on wider areas of delay correlation and on alterations to the interaction layer to handle more complex interaction behavior. Example of alterations include more control of how action packets accepted or discarded, and optimizations to the action packet structure to reduce the amount of computation and transmitted information.

## REPRODUCIBILITY STATEMENT

We provide the source code for all experiments as supplementary material to the paper submission, including the raw measurements used to generate all plots and tables in the paper. The source code is accompanied by instructions for how to run the experiments and how to install the necessary dependencies. All experiments are performed using simulated environments, allowing anyone to reproduce the results shown here in the paper.

## REFERENCES

Yann Bouteiller, Simon Ramstedt, Giovanni Beltrame, Christopher Pal, and Jonathan Binas. Reinforcement learning with random delays. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=QFYnKlBJYR>.

D.M. Brooks and C.T. Leondes. Technical note - markov decision processes with state-information lag. *Operations Research*, 20(4):904–907, aug 1972. doi: 10.1287/opre.20.4.904.

Baiming Chen, Mengdi Xu, Liang Li, and Ding Zhao. Delay-aware model-based reinforcement learning for continuous control. *Neurocomputing*, 450:119–128, apr 2021. ISSN 0925-2312. doi: <https://doi.org/10.1016/j.neucom.2021.04.015>. URL <https://www.sciencedirect.com/science/article/pii/S0925231221005427>.

E. O. Elliott. Estimates of error rates for codes on burst-noise channels. *The Bell System Technical Journal*, 42(5):1977–1997, 1963. doi: 10.1002/j.1538-7305.1963.tb00955.x.

Vlad Firoiu, Tina Ju, and Josh Tenenbaum. At human speed: Deep reinforcement learning with action delay. *CoRR*, abs/1810.07286, 2018. URL <http://arxiv.org/abs/1810.07286>.

- 
- 540 E. N. Gilbert. Capacity of a burst-noise channel. *The Bell System Technical Journal*, 39(5):1253–  
541 1265, 1960. doi: 10.1002/j.1538-7305.1960.tb03959.x.
- 542
- 543 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
544 maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer Dy and An-  
545 dreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*,  
546 volume 80 of *Proceedings of Machine Learning Research*, pp. 1861–1870. PMLR, 10–15 Jul  
547 2018. URL <https://proceedings.mlr.press/v80/haarnoja18b.html>.
- 548 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James  
549 Davidson. Learning latent dynamics for planning from pixels. In Kamalika Chaudhuri and Ruslan  
550 Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*,  
551 volume 97 of *Proceedings of Machine Learning Research*, pp. 2555–2565. PMLR, 09–15 Jun  
552 2019. URL <https://proceedings.mlr.press/v97/hafner19a.html>.
- 553 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning  
554 behaviors by latent imagination. In *International Conference on Learning Representations*, 2020.  
555 URL <https://openreview.net/forum?id=S110TC4tDS>.
- 556
- 557 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. URL [https://](https://arxiv.org/abs/1606.08415)  
558 [arxiv.org/abs/1606.08415](https://arxiv.org/abs/1606.08415).
- 559
- 560 Armin Karamzade, Kyungmin Kim, Montek Kalsi, and Roy Fox. Reinforcement learning from  
561 delayed observations via world models, 2024. URL [https://arxiv.org/abs/2403.](https://arxiv.org/abs/2403.12309)  
562 12309.
- 563 K.V. Katsikopoulos and S.E. Engelbrecht. Markov decision processes with delays and asynchronous  
564 cost collection. *IEEE Transactions on Automatic Control*, 48(4):568–574, 2003. doi: 10.1109/  
565 TAC.2003.809799.
- 566
- 567 Jangwon Kim, Hangyeol Kim, Jiwook Kang, Jongchan Baek, and Soohye Han. Belief  
568 projection-based reinforcement learning for environments with delayed feedback. In A. Oh,  
569 T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in*  
570 *Neural Information Processing Systems*, volume 36, pp. 678–696. Curran Associates, Inc.,  
571 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/0252a434b18962c94910c07cd9a7fecc-Paper-Conference.pdf)  
572 [file/0252a434b18962c94910c07cd9a7fecc-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/0252a434b18962c94910c07cd9a7fecc-Paper-Conference.pdf).
- 573 Leonard Kleinrock. *Theory, Volume 1, Queueing Systems*. Wiley-Interscience, USA, 1975. ISBN  
574 0471491101.
- 575 Lucy Lai, Ann ZX Huang, and Samuel J Gershman. Action chunking as policy compression. In  
576 *PsyArXiv*, 2022.
- 577
- 578 Qiyang Li, Zhiyuan Zhou, and Sergey Levine. Reinforcement learning with action chunking. In  
579 *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL  
580 <https://openreview.net/forum?id=XUks1Y96NR>.
- 581
- 582 Pierre Liotet, Davide Maran, Lorenzo Bisi, and Marcello Restelli. Delayed reinforcement learning  
583 by imitation. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu,  
584 and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*,  
585 volume 162 of *Proceedings of Machine Learning Research*, pp. 13528–13556. PMLR, 17–23 Jul  
586 2022. URL <https://proceedings.mlr.press/v162/liotet22a.html>.
- 587
- 588 Rogelio Luck and Asok Ray. An observer-based compensator for distributed delays. *Auto-*  
589 *matica*, 26(5):903–908, 1990. ISSN 0005-1098. doi: [https://doi.org/10.1016/0005-1098\(90\)](https://doi.org/10.1016/0005-1098(90)90007-5)  
590 90007-5. URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/0005109890900075)  
591 0005109890900075.
- 592
- 593 Thomas M. Moerland, Joost Broekens, Aske Plaat, and Catholijn M. Jonker. Model-based rein-  
594forcement learning: A survey. *Foundations and Trends® in Machine Learning*, 16(1):1–118,  
595 2023. ISSN 1935-8237. doi: 10.1561/22000000086. URL [http://dx.doi.org/10.1561/](http://dx.doi.org/10.1561/22000000086)  
596 22000000086.

---

594 Simon Ramstedt and Chris Pal. Real-time reinforcement learning. In H. Wallach,  
595 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Ad-*  
596 *vances in Neural Information Processing Systems*, volume 32. Curran Associates,  
597 Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/54e36c5ff5f6a1802925ca009f3ebb68-Abstract.html)  
598 [54e36c5ff5f6a1802925ca009f3ebb68-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/54e36c5ff5f6a1802925ca009f3ebb68-Abstract.html).

599 Asok Ray. Distributed data communication networks for real-time process control. *Chemical En-*  
600 *gineering Communications*, 65(1):139–154, 1988. doi: 10.1080/00986448808940249. URL  
601 <https://doi.org/10.1080/00986448808940249>.

602 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy  
603 optimization algorithms. *CoRR*, abs/1707.06347, 2017. URL [http://arxiv.org/abs/](http://arxiv.org/abs/1707.06347)  
604 [1707.06347](http://arxiv.org/abs/1707.06347).

605 Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and  
606 Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots, June 2018.

607 Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U. Balis, Gianluca De Cola, Tristan Deleu,  
608 Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, Rodrigo Perez-Vicente, An-  
609 drea Pierré, Sander Schulhoff, Jun Jet Tai, Hannah Tan, and Omar G. Younis. Gymnasium: A  
610 standard interface for reinforcement learning environments, 2024. URL [https://arxiv.](https://arxiv.org/abs/2407.17032)  
611 [org/abs/2407.17032](https://arxiv.org/abs/2407.17032).

612 David Valensi, Esther Derman, Shie Mannor, and Gal Dalal. Tree search-based policy optimization  
613 under stochastic execution delay. In B. Kim, Y. Yue, S. Chaudhuri, K. Fragkiadaki, M. Khan, and  
614 Y. Sun (eds.), *International Conference on Representation Learning*, volume 2024, pp. 18475–  
615 18495, 2024. URL [https://proceedings.iclr.cc/paper\\_files/paper/2024/](https://proceedings.iclr.cc/paper_files/paper/2024/file/50e13537f46656a94a7acaf022921385-Paper-Conference.pdf)  
616 [file/50e13537f46656a94a7acaf022921385-Paper-Conference.pdf](https://proceedings.iclr.cc/paper_files/paper/2024/file/50e13537f46656a94a7acaf022921385-Paper-Conference.pdf).

617 Thomas J. Walsh, Ali Nouri, Lihong Li, and Michael L. Littman. Learning and planning in environ-  
618 ments with delayed feedback. *Autonomous Agents and Multi-Agent Systems*, 18(1):83, Jul 2008.  
619 ISSN 1573-7454. doi: 10.1007/s10458-008-9056-7. URL [https://doi.org/10.1007/](https://doi.org/10.1007/s10458-008-9056-7)  
620 [s10458-008-9056-7](https://doi.org/10.1007/s10458-008-9056-7).

621 Wei Wang, Dongqi Han, Xufang Luo, and Dongsheng Li. Addressing signal delay in deep rein-  
622 forcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.  
623 URL <https://openreview.net/forum?id=Z8UfDs4J46>.

624 Qingyuan Wu, Simon Sinong Zhan, Yixuan Wang, Yuhui Wang, Chung-Wei Lin, Chen Lv,  
625 Qi Zhu, and Chao Huang. Variational delayed policy optimization. In A. Globerson,  
626 L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in*  
627 *Neural Information Processing Systems*, volume 37, pp. 54330–54356. Curran Associates, Inc.,  
628 2024. URL [https://proceedings.neurips.cc/paper\\_files/paper/2024/](https://proceedings.neurips.cc/paper_files/paper/2024/file/61a18c8a7a1ea7445375dd7255905bc3-Paper-Conference.pdf)  
629 [file/61a18c8a7a1ea7445375dd7255905bc3-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/61a18c8a7a1ea7445375dd7255905bc3-Paper-Conference.pdf).

630 Qingyuan Wu, Yuhui Wang, Simon Sinong Zhan, Yixuan Wang, Chung-Wei Lin, Chen Lv, Qi Zhu,  
631 Jürgen Schmidhuber, and Chao Huang. Directly forecasting belief for reinforcement learning with  
632 delays. In *Forty-second International Conference on Machine Learning*, 2025. URL [https://](https://openreview.net/forum?id=S9unJQditt)  
633 [openreview.net/forum?id=S9unJQditt](https://openreview.net/forum?id=S9unJQditt).

634 Tony Z. Zhao, Vikash Kumar, Sergey Levine, and Chelsea Finn. Learning fine-grained bimanual  
635 manipulation with low-cost hardware. In *Robotics: Science and Systems*, 2023. URL [https://](https://roboticsconference.org/2023/program/papers/016/)  
636 [roboticsconference.org/2023/program/papers/016/](https://roboticsconference.org/2023/program/papers/016/).

637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

---

648	CONTENTS	
649		
650	<b>1 Introduction</b>	<b>1</b>
651		
652	<b>2 Related Work</b>	<b>2</b>
653		
654	<b>3 The Interaction Layer</b>	<b>3</b>
655		
656	3.1 Delayed Markov decision processes . . . . .	3
657		
658	3.2 Handling delays via the interaction layer . . . . .	4
659		
660	3.3 The POMDP model . . . . .	5
661		
662	<b>4 Actor-Critic with Delay Adaptation</b>	<b>6</b>
663		
664	4.1 Heuristic for Assumed Previous Actions . . . . .	6
665		
666	4.2 Model-Based Distribution Agent . . . . .	7
667		
668	4.3 Training Algorithm . . . . .	8
669		
670	<b>5 Evaluation and Results</b>	<b>8</b>
671		
672	<b>6 Conclusion</b>	<b>10</b>
673		
674	<b>A Constant-Delay Augmentation</b>	<b>16</b>
675		
676	<b>B Evaluation Details</b>	<b>17</b>
677		
678	B.1 Action Noise and Its Effect on Performance . . . . .	17
679		
680	B.2 Further Evaluation of Noise Effects . . . . .	18
681		
682	B.3 Evaluated Delay Distributions . . . . .	20
683		
684	B.4 Hyperparameters and Neural Network Structure . . . . .	22
685		
686	B.5 Practical Evaluation Details . . . . .	23
687		
688	<b>C Formal Description of the Interaction Layer POMDP</b>	<b>25</b>
689		
690	<b>D Detailed Model Description</b>	<b>27</b>
691		
692	D.1 Model Components and Objective . . . . .	27
693		
694	D.2 Training Algorithm . . . . .	28
695		
696	<b>E Additional Results</b>	<b>29</b>
697		
698	E.1 Time Series Results . . . . .	29
699		
700	E.1.1 Performance Evaluation under the $GE_{1,23}$ Delay Process . . . . .	30
701		
	E.1.2 Performance Evaluation under the $GE_{4,32}$ Delay Process . . . . .	31
	E.1.3 Performance Evaluation under the $M/M/1$ Queue Delay Process . . . . .	32
	E.2 Model-Based Distribution Agent vs. Adaptiveness . . . . .	33
	E.2.1 Performance of MDA under the $GE_{1,23}$ Delay Process . . . . .	34
	E.2.2 Performance of MDA under the $GE_{4,32}$ Delay Process . . . . .	35

---

702	E.2.3	Performance of MDA under the M/M/1 Queue Delay Process . . . . .	36
703			
704	E.3	Results when violating the upper bound assumptions . . . . .	37
705	E.3.1	GE <sub>1,23</sub> Delay Process with Low CDA . . . . .	38
706	E.3.2	GE <sub>4,32</sub> Delay Process with Low CDA . . . . .	39
707	E.3.3	M/M/1 Queue Delay Process with Low CDA . . . . .	40
708			
709	E.4	Results for average delays . . . . .	41
710	E.4.1	GE <sub>1,23</sub> Delay Process with Average CDA . . . . .	42
711	E.4.2	GE <sub>4,32</sub> Delay Process with Average CDA . . . . .	43
712	E.4.3	M/M/1 Queue Delay Process with Average CDA . . . . .	44
713			
714	E.5	Results from all evaluated baselines . . . . .	45
715	E.5.1	GE <sub>1,23</sub> Delay Process for All Baselines . . . . .	45
716	E.5.2	GE <sub>4,32</sub> Delay Process for All Baselines . . . . .	46
717	E.5.3	M/M/1 Queue Delay Process for All Baselines . . . . .	46
718			
719			
720			
721	<b>F</b>	<b>Practical Considerations of the Interaction Layer</b>	<b>47</b>
722			
723	F.1	Considerations for Non-interaction Delays . . . . .	47
724	F.2	Effect of Action Packet on Computational Delay . . . . .	47
725	F.3	Compensating for Excessive Computational Delays . . . . .	48
726	F.4	Effect of Action Packet on Communication . . . . .	49
727			
728			
729	<b>G</b>	<b>Interaction-Delayed Reinforcement Learning with Real-Valued Delay</b>	<b>50</b>
730			
731	G.1	Origin and Effect of Delay as Continuous Time . . . . .	50
732	G.2	Interaction Layer to Enforce Discrete Delay . . . . .	51
733			
734			
735			
736			
737			
738			
739			
740			
741			
742			
743			
744			
745			
746			
747			
748			
749			
750			
751			
752			
753			
754			
755			

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

---

## ICLR STATEMENT ON LLM USAGE

This paper has made use of Large Language Models (LLMs) to polish writing. More specifically, to check for grammatical errors and for phrasing suggestions.

## OUTLINE OF THE APPENDICES

Appendix A presents how we implement constant-delay augmentation (CDA) in our framework. This allows agents to act with constant delay using the interaction layer, even if the underlying delay process is stochastic. We primarily use this to provide a fair comparison against related work.

Appendix B presents the evaluation details. In Appendix B.1, we demonstrate that stochastic transitions are necessary to see the effects of delay, both theoretically and with an evaluated example. We also show in Appendix B.1 how we use action noise to convert deterministic transitions into stochastic ones. Further evaluation of the effect that stochasticity has on the best-performing algorithms is presented in Appendix B.2. Appendix B.3 describes the delay distributions used in the benchmarks. Lastly, in Appendix B.4, we present the hyperparameters used, and in Appendix B.5 we present the software and hardware used for running the benchmarks.

Appendix C formalizes the interaction layer as a POMDP. This POMDP is used to simulate the interaction layer in the benchmarks.

Appendix D formalizes the model and its objective, as well as providing an expanded version of the algorithm presented in Section 4.3.

Appendix E contains additional results. Appendix E.1 contains all results presented in the conclusion, with time series plots and standard deviation. In Appendix E.2, we evaluate the model-based distribution agent under CDA, showing that it is the adaptiveness that leads to gains in performance rather than the policy itself. In Appendix E.3, we evaluate the effect of using a lower bound for CDA that holds most of the time, but is occasionally violated. The latter two Appendices E.2 and E.3 show that the adaptiveness of the interaction layer offers gains and stability in performance that cannot be obtained by operating in a constant-delay manner.

Appendix F discusses practical considerations when deploying the interaction layer to real-world environments.

Appendix G shows an alternative, more realistic definition of delayed MDPs, which uses real-valued delays rather than discrete.

---

## A CONSTANT-DELAY AUGMENTATION

To be able to evaluate and compare fairly with state-of-the-art algorithms, we make it possible to have constant delay augmentation within our framework. Specifically, to allow algorithms such as BPQL and VDPO that expect a constant delay when the underlying delay process is stochastic, we apply a *constant-delay augmentation* (CDA) on top of the interaction layer. CDA converts the interaction layer POMDP into a constant-delay MDP, under the assumption that the maximum delay does not exceed  $h$  steps. This augmentation ensures that we evaluate state-of-the-art as intended when comparing their performance against ACDA.

CDA is implemented on top of the interaction layer by simply arranging the contents of the action packet matrix such that, no matter when action packet  $\mathbf{a}_t$  arrives (between  $t + 1$  and  $t + h$ ), each action will be executed  $h$  steps after it was generated. We illustrate this procedure of constructing the action packet in Algorithm 3.

---

### Algorithm 3 Constant-Delay Augmentation using the Interaction Layer

---

**Input**  $\mathbf{o}_t = (t, s_t, \mathbf{b}_t, \delta_t, c_t)$  (Observation packet)  
 $\pi$  (Constant-delay policy operating on the horizon  $h$ )

- 1:  $a_t, \dots, a_{t+h-1} = \mathbf{b}_t$
- 2:  $a_{t+h} \sim \pi(\cdot | s_t, a_t, \dots, a_{t+h-1})$

$$3: \mathbf{a}_t = \left( t, \begin{bmatrix} a_{t+1} & a_{t+2} & a_{t+3} & \cdots & a_{t+h-2} & a_{t+h-1} & a_{t+h} \\ a_{t+2} & a_{t+3} & a_{t+4} & \cdots & a_{t+h-1} & a_{t+h} & a_{t+h} \\ a_{t+3} & a_{t+4} & a_{t+5} & \cdots & a_{t+h} & a_{t+h} & a_{t+h} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ a_{t+h-1} & a_{t+h} & a_{t+h} & \cdots & a_{t+h} & a_{t+h} & a_{t+h} \\ a_{t+h} & a_{t+h} & a_{t+h} & \cdots & a_{t+h} & a_{t+h} & a_{t+h} \end{bmatrix} \right)$$

- 4: **return**  $\mathbf{a}_t$
- 

This states that if  $\mathbf{a}_t$  arrives at  $t + i$ , then the actions to be applied are  $a_{t+i}, a_{t+\min(h,i+1)}, a_{t+\min(h,i+2)}, \dots, a_{t+\min(h,i+(h-3))}, a_{t+\min(h,i+(h-2))}, a_{t+h}$ , which ensures that  $\mathbf{b}_t$  always is a correct guess of the actions to be applied next. For  $i > 1$ , we pad with  $a_{t+h}$  to the right on each row to represent the shifting behavior. Forming the action packets in this way ensures that  $a_{t+h}$  always gets executed at time  $t + h$ , given that the delay does not exceed  $h$ .

The policy  $\pi$  can be any constant-delay augmented policy. We can also apply the model-based distribution policy from the ACDA algorithm to the CDA setting, by letting  $\pi(a_{t+h} | s_t, a_t, \dots, a_{t+h-1}) = \pi_\theta(a_{t+h} | z_h)$ , where  $z_h = \text{STEP}_\omega^h(\text{EMBED}_\omega(s_t), a_t, \dots, a_{t+h-1})$ . We present the results of this policy in Appendix E.2.

This assumes the horizon  $h$  is a valid upper bound of the delay. We can still perform the augmentation if  $h$  is less than the upper bound, but then we are no longer guaranteed the MDP properties of constant delay. We present the results of this in Appendix E.3.

---

## 864 B EVALUATION DETAILS

865  
866 This section provides a more complete overview of the evaluation and the results. We provide  
867 justification for choosing the 5% noise on environments (Appendix B.1), the delay processes used  
868 (Appendix B.3), the hyperparameters used and neural network architectures used (Appendix B.4),  
869 as well as the software and hardware used during evaluation (Appendix B.5). The complete results  
870 for all benchmarks are presented separately in Appendix E.  
871

### 872 B.1 ACTION NOISE AND ITS EFFECT ON PERFORMANCE

873  
874 The benchmark environments used, as defined in Gymnasium, have fully deterministic transitions.  
875 As a result, they are theoretically unaffected by delay: the optimal value achievable in the delayed  
876 MDP is identical to that of the undelayed MDP. This follows trivially from the fact that, with a  
877 perfect deterministic model of the MDP dynamics, the agent can precisely predict the future state in  
878 which its action will be applied. Consequently, the agent can plan as if there were no delay at all.  
879

880 The same is not true for MDPs with stochastic transition dynamics. To show this, consider the MDP  
881 with  $S = \{H, T\}$ ,  $A = \{H, T\}$ ,  $r(H, H) = 1$ ,  $r(T, T) = 1$ ,  $r(H, T) = 0$ ,  $r(T, H) = 0$ , where  
882  $\forall s', s, a \quad p(s'|s, a) = 0.5$ . This MDP models flipping a fair coin, where the agent is given a reward  
883 of 1 if it can correctly identify the face of the current coin. Consider this MDP with a constant delay  
884 of 1 time step. Now, the agent instead has to guess the face of the next coin, on which it can do no  
885 better than a 50/50 guess. Therefore, in this example, the value of an optimal agent in the delayed  
886 MDP is half of the value of an optimal agent in the undelayed MDP.

887 To better highlight the practical issues with delay, we add uncertainty to transitions in the Gymna-  
888 sium environments by adding noise to the actions prior to being applied to the environment. Let  $\beta$   
889 be the noise factor indicating how much noise we add relative to the span of values that the action  
890 can take. Then we add noise to the actions  $a$  as follows:  
891

$$892 \text{ Assume } a = [a(1), a(2), \dots, a(n)] \quad (2)$$

$$893 a(i)_{\max} = \text{maximum value for } a(i) \quad (3)$$

$$894 a(i)_{\min} = \text{minimum value for } a(i) \quad (4)$$

$$895 \nu(i) = \beta \cdot (a(i)_{\max} - a(i)_{\min}) \cdot \xi \Big|_{\xi \sim \mathcal{N}(0,1)} \quad (5)$$

$$896 \tilde{a}(i) = \text{clip}(a(i) + \nu(i), a(i)_{\min}, a(i)_{\max}) \quad (6)$$

$$897 \tilde{a} = [\tilde{a}(1), \tilde{a}(2), \dots, \tilde{a}(n)] \quad (7)$$

902  
903 Here, we assume that the actions are continuous, which works since all environments in our evalua-  
904 tion are of this nature. Then the transitions become  $p(s'|s, \tilde{a})$ , with the noisy action applied instead  
905 of the original one. We use the noise factor  $\beta = 0.05$  in all our noisy environments evaluated here.

906 To see the effect that this noise has on delayed RL in practice, we evaluate the performance of BPQL  
907 when trained over different constant delays, with and without noise. The results are plotted in Figure  
908 9. The evaluation is done by training a BPQL policy on a specific constant delay and action noise,  
909 evaluating the policy’s average return every 10000 steps, and reporting the best achieved average  
910 return as the performance. This evaluation procedure, which is used by all evaluations in the paper,  
911 is further described in Appendix B.5.

912 The results without noise for constant delays of 3, 6, and 9, shown in Table 2, are representative  
913 of those reported by Kim et al. (2023) ( $8100 \pm 543.4$ ,  $6334.6 \pm 245.3$ , and  $5887.5 \pm 270.5$  for  
914 constant delays of 3, 6, and 9 respectively). This suggests that our implementation is faithful to their  
915 approach. Notably, we observe that in the deterministic setting, the impact of delay, while causing  
916 a significant initial drop in performance, does not lead to significant degradation over longer time  
917 horizons. This behavior contrasts with the noisy environments, where the performance declines  
more noticeably as the delay increases.

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

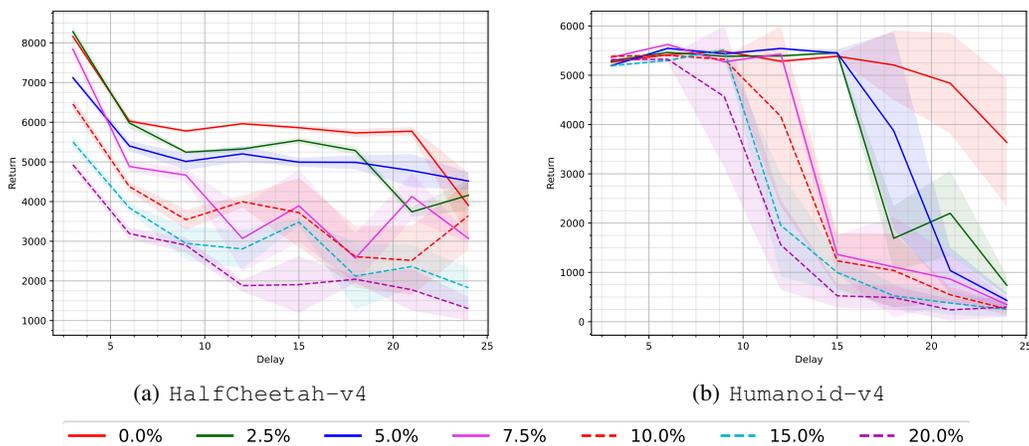


Figure 9: Best evaluated performance of BPQL after training over  $10^6$  timesteps when different noise is applied to the `HalfCheetah-v4` and `Humanoid-v4` environments. Each line represents when a specific action noise is induced on that environment, indicated by the % in the legend (e.g. 2.5% means  $\beta = 0.025$ ). Each plotted point represents the best evaluated average return when BPQL is trained on that noisy environment with that constant delay. The shaded regions represent the standard deviation.

Table 2: The noise evaluation measurements shown in Figure 9.

HalfCheetah-v4									
	Delay 3	Delay 6	Delay 9	Delay 12	Delay 15	Delay 18	Delay 21	Delay 24	
0% noise	8159.28 ± 50.12	<b>6027.13 ± 50.79</b>	<b>5778.11 ± 43.85</b>	5244.27 ± 43.00	5325.65 ± 62.94	5544.09 ± 78.04	5284.89 ± 103.65	3740.87 ± 127.23	3902.62 ± 822.27
2.5% noise	<b>8281.45 ± 143.90</b>	5981.20 ± 137.60	5244.27 ± 43.00	5325.65 ± 62.94	5544.09 ± 78.04	5284.89 ± 103.65	3740.87 ± 127.23	4158.25 ± 382.17	
5% noise	7121.49 ± 58.07	5399.96 ± 122.38	5010.51 ± 211.44	5201.90 ± 197.06	4992.80 ± 110.72	4988.70 ± 179.34	4778.82 ± 419.50	<b>4516.67 ± 222.24</b>	
7.5% noise	7839.51 ± 169.20	4883.81 ± 138.27	4665.25 ± 159.32	3069.77 ± 828.40	3893.18 ± 914.63	2575.36 ± 627.12	4128.39 ± 580.59	3071.16 ± 104.56	
10% noise	6459.33 ± 158.08	4376.52 ± 114.97	3543.70 ± 257.16	3996.73 ± 153.34	3721.40 ± 873.19	2610.57 ± 777.65	2517.84 ± 875.47	3641.14 ± 851.81	
15% noise	5500.11 ± 125.74	3841.51 ± 100.05	2946.32 ± 409.99	2805.37 ± 491.86	3483.07 ± 91.09	2122.89 ± 842.33	2363.06 ± 544.60	1830.24 ± 517.49	
20% noise	4929.20 ± 90.98	3192.02 ± 142.30	2903.07 ± 192.16	1881.11 ± 149.04	1907.26 ± 716.61	2042.48 ± 91.94	1776.34 ± 528.88	1300.96 ± 298.35	
Humanoid-v4									
	Delay 3	Delay 6	Delay 9	Delay 12	Delay 15	Delay 18	Delay 21	Delay 24	
0.0% noise	5265.75 ± 16.10	5416.97 ± 4.56	5483.19 ± 11.20	5283.45 ± 34.02	5383.76 ± 49.09	<b>5207.48 ± 704.16</b>	<b>4836.92 ± 1019.75</b>	<b>3639.32 ± 1295.49</b>	
2.5% noise	5292.31 ± 96.10	5462.32 ± 19.79	5385.21 ± 12.65	5392.97 ± 38.69	<b>5456.68 ± 18.84</b>	1692.16 ± 660.11	2199.76 ± 866.64	739.62 ± 198.24	
5.0% noise	5194.99 ± 8.40	5545.78 ± 46.43	5435.19 ± 40.38	<b>5543.46 ± 19.15</b>	5448.85 ± 77.24	3873.11 ± 2007.65	1038.41 ± 466.64	431.75 ± 154.26	
7.5% noise	5362.93 ± 10.58	<b>5625.55 ± 30.05</b>	5271.80 ± 24.04	5430.81 ± 33.14	1364.47 ± 391.72	1110.10 ± 1023.93	863.93 ± 490.55	346.29 ± 109.18	
10.0% noise	<b>5386.11 ± 18.15</b>	5403.12 ± 8.27	5322.27 ± 45.70	4167.17 ± 1831.92	1235.33 ± 563.62	1039.16 ± 747.36	545.63 ± 162.01	258.85 ± 116.50	
15.0% noise	5194.31 ± 44.01	5299.12 ± 29.14	<b>5516.88 ± 42.33</b>	1955.98 ± 1041.45	1001.01 ± 383.28	526.00 ± 235.22	378.92 ± 243.38	247.76 ± 135.52	
20.0% noise	5317.08 ± 29.38	5323.95 ± 58.41	4566.51 ± 1432.43	1561.88 ± 916.61	525.42 ± 236.13	490.00 ± 237.10	239.58 ± 214.93	297.82 ± 199.20	

We therefore conclude that a fair evaluation of delayed RL should be done in environments with stochastic dynamics. Further practical evaluation of the effect that noise has on performance across different training algorithms is shown in Appendix B.2.

## B.2 FURTHER EVALUATION OF NOISE EFFECTS

To see how noise affects the best performing delay-aware algorithms, BPQL, VDPO, and ACDA, we evaluate their performance as the noise varies from 0% to 25% on the same set of chosen environments across the three delay processes. As explained in Appendix B.5, each measurement represents the best average return for a policy trained on the combination of environment, action noise, and delay process. The results for delay processes  $GE_{1,23}$ ,  $GE_{4,32}$ , and MM1 are shown in Tables 3, 4, and 5, respectively.

Note that VDPO uses a fixed seed when resetting an environment. Therefore, VDPO has a standard deviation of 0 when evaluating without action noise, as all sampled trajectories are deterministic. This is due to us using the original VDPO implementation with as few modifications as possible, as further explained in Appendix B.5.

The results show a trend that ACDA performs even better as the noise increases. As ACDA adapts to the sampled delays, this performance increase is expected because noisy environments can still perform well for lower delays, as shown in Figure 9.

There are some outliers in the results, where the performance is slightly better for a higher action noise. We believe that these are due to randomness and that they would not be present if we sampled more trajectories during evaluation and averaged across multiple trained policies.

Table 3: Best returns from the  $GE_{1,23}$  delay process over different noise.

Ant-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	3736.70 ±	108.62	2691.88 ±	129.84	1421.78 ±	297.93	640.97 ±	316.38	69.69 ±	75.58	1.31 ±	18.92
VDPO	3492.27 ±	0.00	2163.00 ±	53.04	1162.88 ±	603.92	644.82 ±	373.70	296.89 ±	131.34	34.66 ±	39.05
ACDA	<b>4719.08 ±</b>	<b>658.29</b>	<b>4112.78 ±</b>	<b>818.44</b>	<b>2780.25 ±</b>	<b>761.75</b>	<b>1209.94 ±</b>	<b>832.61</b>	<b>536.10 ±</b>	<b>400.29</b>	<b>192.79 ±</b>	<b>105.91</b>
Humanoid-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	2462.64 ±	1341.26	585.19 ±	163.49	261.12 ±	125.97	365.19 ±	172.41	298.31 ±	200.46	237.62 ±	161.17
VDPO	464.39 ±	0.00	417.25 ±	210.09	312.26 ±	145.72	285.21 ±	186.51	276.49 ±	174.40	325.55 ±	130.45
ACDA	<b>4842.13 ±</b>	<b>861.55</b>	<b>4608.76 ±</b>	<b>1084.52</b>	<b>3751.03 ±</b>	<b>1552.10</b>	<b>3638.29 ±</b>	<b>1849.70</b>	<b>3852.50 ±</b>	<b>1237.91</b>	<b>1597.85 ±</b>	<b>1143.37</b>
HalfCheetah-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	4176.16 ±	897.01	4320.20 ±	1028.52	3908.33 ±	77.24	1810.59 ±	635.08	1899.01 ±	89.66	1206.78 ±	154.01
VDPO	4976.10 ±	0.00	3144.23 ±	1156.52	2240.86 ±	591.26	1664.11 ±	144.79	1049.28 ±	167.12	799.84 ±	212.99
ACDA	<b>6087.67 ±</b>	<b>1142.66</b>	<b>5984.25 ±</b>	<b>1885.78</b>	<b>5838.64 ±</b>	<b>724.34</b>	<b>4656.72 ±</b>	<b>693.20</b>	<b>4446.73 ±</b>	<b>627.70</b>	<b>2783.35 ±</b>	<b>194.57</b>
Hopper-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	3176.22 ±	48.33	1328.71 ±	937.67	549.60 ±	541.58	232.25 ±	176.51	135.24 ±	100.37	111.19 ±	92.11
VDPO	<b>3477.61 ±</b>	<b>0.00</b>	709.20 ±	522.01	181.60 ±	60.08	150.74 ±	94.33	96.75 ±	54.71	77.14 ±	73.00
ACDA	2381.98 ±	1226.41	<b>2094.65 ±</b>	<b>944.20</b>	<b>2344.23 ±</b>	<b>1167.03</b>	<b>1330.55 ±</b>	<b>895.65</b>	<b>1636.10 ±</b>	<b>1134.22</b>	<b>1057.21 ±</b>	<b>949.42</b>
Walker2d-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	1287.71 ±	754.84	1215.91 ±	776.93	652.90 ±	501.11	316.01 ±	216.53	595.48 ±	668.53	314.27 ±	417.49
VDPO	2005.31 ±	0.00	846.88 ±	808.67	810.89 ±	1173.36	283.58 ±	334.85	186.02 ±	307.16	199.08 ±	375.33
ACDA	<b>4030.01 ±</b>	<b>82.46</b>	<b>3863.59 ±</b>	<b>232.52</b>	<b>4295.73 ±</b>	<b>128.36</b>	<b>4045.24 ±</b>	<b>51.37</b>	<b>3199.49 ±</b>	<b>231.87</b>	<b>3234.59 ±</b>	<b>623.28</b>

Table 4: Best returns from the  $GE_{4,32}$  delay process over different noise.

Ant-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	3523.32 ±	146.78	2509.52 ±	117.37	<b>1456.09 ±</b>	<b>299.72</b>	547.34 ±	237.21	67.78 ±	122.29	0.28 ±	8.70
VDPO	<b>3574.87 ±</b>	<b>0.00</b>	2266.99 ±	90.89	1167.24 ±	473.41	<b>647.07 ±</b>	<b>346.39</b>	<b>244.78 ±</b>	<b>127.86</b>	26.36 ±	40.46
ACDA	2658.50 ±	285.82	<b>2866.93 ±</b>	<b>1172.46</b>	1406.31 ±	712.09	547.27 ±	329.31	138.69 ±	94.21	<b>29.26 ±</b>	<b>32.36</b>
Humanoid-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	876.89 ±	69.04	276.63 ±	131.70	301.68 ±	134.63	254.22 ±	131.77	201.02 ±	109.72	195.88 ±	123.20
VDPO	442.03 ±	0.00	280.72 ±	169.85	262.74 ±	131.86	233.62 ±	145.57	206.56 ±	159.37	194.34 ±	113.40
ACDA	<b>3877.77 ±</b>	<b>1776.04</b>	<b>3725.59 ±</b>	<b>1513.38</b>	<b>3454.60 ±</b>	<b>1567.23</b>	<b>3092.70 ±</b>	<b>1752.58</b>	<b>1043.06 ±</b>	<b>339.16</b>	<b>649.32 ±</b>	<b>270.74</b>
HalfCheetah-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	4894.08 ±	99.20	2136.36 ±	547.04	2019.46 ±	758.99	1785.40 ±	655.73	1920.19 ±	126.28	1286.61 ±	141.72
VDPO	<b>5059.93 ±</b>	<b>0.00</b>	3664.30 ±	929.25	1923.50 ±	379.20	1510.93 ±	435.71	1177.83 ±	283.21	790.87 ±	174.17
ACDA	4203.13 ±	279.18	<b>4231.15 ±</b>	<b>333.69</b>	<b>3239.61 ±</b>	<b>199.78</b>	<b>3149.08 ±</b>	<b>367.98</b>	<b>3218.57 ±</b>	<b>154.02</b>	<b>1826.06 ±</b>	<b>214.04</b>
Hopper-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	2668.14 ±	711.84	433.29 ±	381.79	190.79 ±	135.89	120.03 ±	137.44	76.54 ±	66.87	77.79 ±	72.13
VDPO	<b>3403.46 ±</b>	<b>0.00</b>	330.44 ±	263.74	138.13 ±	128.53	86.52 ±	81.11	77.02 ±	68.06	59.12 ±	60.31
ACDA	2947.10 ±	929.71	<b>1727.79 ±</b>	<b>959.50</b>	<b>1434.94 ±</b>	<b>805.19</b>	<b>1814.74 ±</b>	<b>1187.29</b>	<b>1128.22 ±</b>	<b>1015.20</b>	<b>532.97 ±</b>	<b>630.97</b>
Walker2d-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	1352.44 ±	328.51	875.09 ±	747.72	343.62 ±	314.86	275.07 ±	207.17	191.56 ±	383.58	134.20 ±	125.93
VDPO	1779.39 ±	0.00	344.73 ±	316.82	123.18 ±	160.70	147.64 ±	258.22	73.78 ±	142.23	56.70 ±	148.33
ACDA	<b>3945.33 ±</b>	<b>148.28</b>	<b>1840.58 ±</b>	<b>386.78</b>	<b>1409.42 ±</b>	<b>281.81</b>	<b>1322.32 ±</b>	<b>305.35</b>	<b>1149.82 ±</b>	<b>345.44</b>	<b>850.48 ±</b>	<b>172.82</b>

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

Table 5: Best returns from the MM1 delay process over different noise.

Ant-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	<b>3717.34 ± 125.79</b>		<b>3074.17 ± 106.78</b>		1680.23 ± 307.70		<b>764.00 ± 295.30</b>		123.27 ± 122.10		4.06 ± 13.70	
VDPO	3638.48 ± 9.30		2528.67 ± 144.63		1319.55 ± 537.44		709.82 ± 254.73		<b>334.61 ± 133.75</b>		32.48 ± 27.59	
ACDA	2593.23 ± 88.19		2898.46 ± 838.07		<b>1941.68 ± 567.39</b>		724.43 ± 487.25		306.75 ± 319.66		<b>76.13 ± 110.11</b>	
Humanoid-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	2926.45 ± 1042.65		5435.29 ± 68.34		733.25 ± 410.45		554.66 ± 205.97		423.77 ± 227.06		406.50 ± 195.39	
VDPO	762.75 ± 0.00		720.73 ± 634.35		544.09 ± 424.14		423.67 ± 252.84		526.95 ± 394.22		354.99 ± 129.91	
ACDA	<b>5238.97 ± 332.60</b>		<b>5805.60 ± 23.04</b>		<b>5548.29 ± 39.58</b>		<b>5343.60 ± 227.82</b>		<b>1752.02 ± 616.85</b>		<b>1585.00 ± 538.02</b>	
HalfCheetah-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	5627.41 ± 64.54		4660.93 ± 448.10		2291.63 ± 857.39		2171.03 ± 626.53		2026.23 ± 523.82		1574.84 ± 167.54	
VDPO	4684.17 ± 0.00		3831.96 ± 960.07		2454.03 ± 415.56		1896.91 ± 419.64		1277.00 ± 244.75		939.40 ± 251.26	
ACDA	<b>6309.62 ± 356.30</b>		<b>5898.36 ± 409.10</b>		<b>4998.40 ± 414.15</b>		<b>4173.81 ± 96.85</b>		<b>2547.28 ± 106.75</b>		<b>2243.67 ± 122.99</b>	
Hopper-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	3130.47 ± 29.44		3035.66 ± 103.80		1106.35 ± 490.33		397.27 ± 249.50		319.40 ± 198.93		196.96 ± 105.23	
VDPO	<b>3797.33 ± 0.00</b>		1459.88 ± 933.11		389.41 ± 247.52		201.10 ± 121.99		156.50 ± 87.59		152.03 ± 96.74	
ACDA	3029.85 ± 565.47		<b>3122.53 ± 417.37</b>		<b>2245.07 ± 1166.01</b>		<b>2250.64 ± 901.67</b>		<b>1079.67 ± 690.26</b>		<b>1189.84 ± 677.60</b>	
Walker2d-v4												
	0% noise		5% noise		10% noise		15% noise		20% noise		25% noise	
BPQL	3815.63 ± 39.18		3547.73 ± 133.51		2182.64 ± 763.72		883.42 ± 187.43		691.82 ± 431.84		758.42 ± 558.27	
VDPO	<b>5202.62 ± 0.00</b>		2144.25 ± 1650.85		1416.45 ± 1845.96		1538.01 ± 1867.52		1227.54 ± 1220.62		637.52 ± 1228.54	
ACDA	4485.74 ± 55.72		<b>4562.33 ± 87.98</b>		<b>3653.42 ± 35.72</b>		<b>3892.52 ± 52.42</b>		<b>3036.04 ± 631.82</b>		<b>1608.60 ± 877.13</b>	

### B.3 EVALUATED DELAY DISTRIBUTIONS

As mentioned in Section 5, we evaluate on delay processes following the Gilbert-Elliot and M/M/1 models. We formally define these processes in this section, as well as their conservative and optimistic worst-case delay assumptions (high and low CDA). Appendix E.2 and E.3 evaluate the performance under high and low CDA, respectively.

We consider  $GE_{1,23}$  and  $GE_{4,32}$ , two Gilbert-Elliot models. These are Markovian processes alternating between two states, a good state  $s_{\text{good}}$  and a bad state  $s_{\text{bad}}$ , as illustrated in Figure 10. We describe the models in Table 6 as a two-state Markov process, where they initially start in the good state. The notation  $D(d|s)$  is used to describe the probability of sampling the delay  $d$  in the Gilbert-Elliot state  $s$ . We set the opportunistic low CDA to be the maximum delay that can be sampled in the  $s_{\text{good}}$  state.

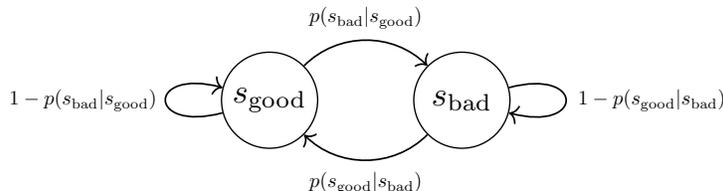


Figure 10: Illustration of transitions in a Gilbert-Elliot model.

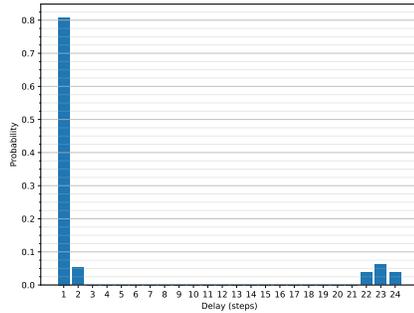
We plot the distribution histogram and time series over 1000 samples of the Gilbert-Elliot processes in Figure 11.

The M/M/1 queue process is described by simulating an M/M/1 queue according to the pseudocode in Algorithm 4. We set the arrival rate  $\lambda_{\text{arrive}} = 0.33$  and the service rate  $\lambda_{\text{service}} = 0.75$ . The arrivals and departures are dictated by independent Poisson processes parametrized by these values (e.g., the time between two arrivals is a r.v. with an exponential distribution of mean  $\lambda_{\text{arrive}}$ ). Note that there is no upper bound on this delay process, and it is therefore impossible to convert this to a constant-delay MDP through Constant Delay Augmentation (CDA). We can still apply the CDA

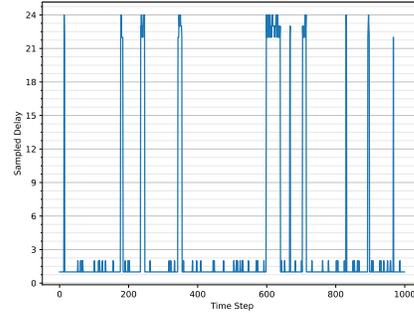
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

Table 6: Description of the Gilbert-Elliot delay processes used during evaluation.

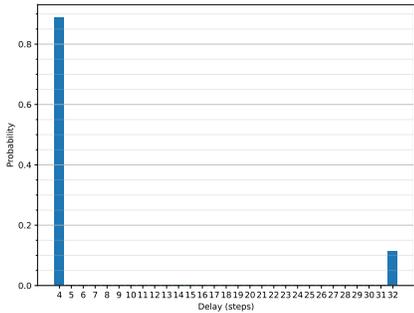
Property	GE <sub>1,23</sub>	GE <sub>4,32</sub>
$p(s_{\text{bad}} s_{\text{good}})$	$\frac{1}{125}$	$\frac{1}{250}$
$p(s_{\text{good}} s_{\text{bad}})$	$\frac{1}{20}$	$\frac{1}{32}$
$D(d s_{\text{good}})$	$\Pr[d = 1] = \frac{15}{16}$ $\Pr[d = 2] = \frac{1}{16}$	$\Pr[d = 4] = 1$
$D(d s_{\text{bad}})$	$\Pr[d = 22] = \frac{3}{11}$ $\Pr[d = 23] = \frac{5}{11}$ $\Pr[d = 24] = \frac{3}{11}$	$\Pr[d = 32] = 1$
Low CDA	2	4
High CDA	24	32



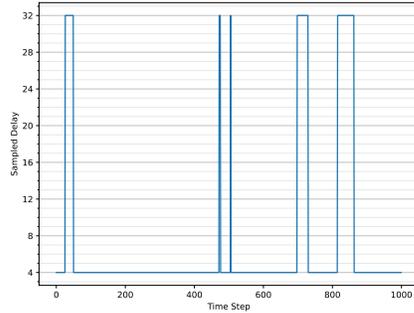
(a) GE<sub>1,23</sub> distribution histogram



(b) GE<sub>1,23</sub> distribution time series



(c) GE<sub>4,32</sub> distribution histogram



(d) GE<sub>4,32</sub> distribution time series

Figure 11: The Gilbert-Elliot delay processes.

conversion from Appendix A to apply constant-delay methodologies on this delay process, though they are no longer operating on an MDP, as the true delay may exceed the assumed upper bound.

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155

---

**Algorithm 4** M/M/1 Queue Delay Generator

---

```

Initial State:  $t_{\text{arrival}} \sim \text{Exp}(\cdot|\lambda_{\text{arrive}})$  (Time of arrival of the first packet)
                   $t_{\text{service}} \leftarrow \emptyset$  (Cannot serve anything yet)
                   $Q \leftarrow \text{FIFOqueue}()$  (Empty queue initially)

1: procedure SAMPLEDELAY
2:   if  $t_{\text{service}} = \emptyset$  then
3:      $t \leftarrow t_{\text{arrival}}$ 
4:      $Q.\text{insert}(t)$ 
5:      $t_{\text{arrival}} \sim \text{Exp}(\cdot|\lambda_{\text{arrive}}) + t$ 
6:      $t_{\text{service}} \sim \text{Exp}(\cdot|\lambda_{\text{service}}) + t$ 
7:   while  $t_{\text{arrival}} < t_{\text{service}}$  do
8:      $t \leftarrow t_{\text{arrival}}$ 
9:      $Q.\text{insert}(t)$ 
10:     $t_{\text{arrival}} \sim \text{Exp}(\cdot|\lambda_{\text{arrive}}) + t$ 
11:     $t \leftarrow t_{\text{service}}$ 
12:     $t_{\text{inserted}} \leftarrow Q.\text{pop}()$ 
13:     $d \leftarrow \lceil t - t_{\text{inserted}} \rceil$ 
14:    if  $Q.\text{isempty}()$  then
15:       $t_{\text{service}} \leftarrow \emptyset$ 
16:    else
17:       $t_{\text{service}} \sim \text{Exp}(\cdot|\lambda_{\text{service}}) + t$ 
18:    return  $d$ 

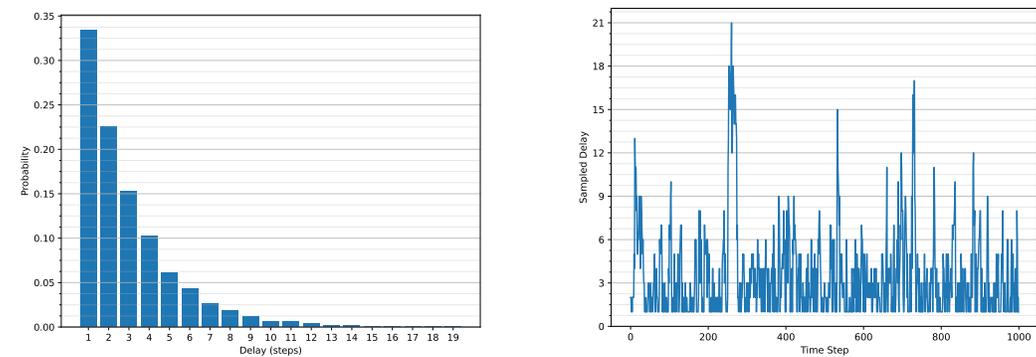
```

---

1156  
1157  
1158  
1159  
1160

We plot delays of the M/M/1 queue in Figure 12. We set the conservative delay (high CDA) to be 16, and the opportunistic delay (low CDA) to be 4 for the M/M/1 queue.

1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172



(a) MM1 distribution histogram over 5000 samples

(b) MM1 distribution time series

1173  
1174  
1175  
1176  
1177

Figure 12: Delays from an M/M/1 queue when  $\lambda_{\text{arrive}} = 0.33$  and  $\lambda_{\text{service}} = 0.75$ .

**B.4 HYPERPARAMETERS AND NEURAL NETWORK STRUCTURE**

1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

Hyperparameters used for training the policy  $\pi_{\theta}$  and the critic  $Q_{\phi}$  in ACDA follow the common learning rates used for SAC, and are therefore shared between SAC, BPQL, and ACDA. We show these together with the model hyperparameters used for ACDA in Tables 7(a) and 7(b). The hyperparameters for the model share the same replay size and batch size. We use slightly different parameters for the model when learning the dynamics on the 2D environments (HalfCheetah-v4, Hopper-v4, and Walker2d-v4) and the 3D environments (Ant-v4 and Humanoid-v4).

The  $\pi_{\theta}$ ,  $Q_{\phi}$ , and  $\text{EMBED}_{\omega}$  networks are all implemented as MLPs with 2 hidden layers of dimension 256 each. The policy outputs the mean and standard deviation of a Gaussian distribution, which is put through tanh and scaled to exactly cover the action space. The  $\text{EMIT}_{\omega}$  network consists of 2

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

Table 7(a): SAC Hyperparameters

Parameter	Value
Policy ( $\theta$ ) learning rate	$3 \cdot 10^{-4}$
Critic ( $\phi$ ) Learning rate	$3 \cdot 10^{-4}$
Temperature ( $\alpha$ ) learning rate	$3 \cdot 10^{-4}$
Starting temperature ( $\alpha$ )	0.2
Temperature threshold $\mathcal{H}$	$-\dim(A)$
Target smoothing coefficient	0.005
Replay buffer size	$10^6$
Discount $\gamma$	0.99
Minibatch size	256
Optimizer (policy, critic, temp.)	Adam
Activation (policy, critic)	ReLU

Table 7(b): Model Hyperparameters

Parameter	Value (2D)	Value (3D)
Model ( $\omega$ ) learning rate	$10^{-4}$	$5 \cdot 10^{-5}$
Model training window $n$	16	16
Latent GRU dimensionality	384	512
Angle clamping	Yes	No
Optimizer	Adam	Adam
Activation	ClipSiLU	ClipSiLU

common layers of dimension 256 each, with additional "head" layers of dimension 256 each for outputting the mean and standard deviation of a Gaussian distribution.

Both MLPs ( $\text{EMBED}_\omega$  and  $\text{EMIT}_\omega$ ) in the model make use of a clipped version of SiLU (Hendrycks & Gimpel, 2023) as their activation function, where  $\text{ClipSiLU}(x) = \text{SiLU}(\max(-20, x))$ . We found that the use of ClipSiLU significantly improved the model performance. Earlier experiments with models using ReLU activation did not manage to achieve good performance when used with ACDA.

The angle clamping mentioned in Table 7(b) constrains all components of the state space that represent an angle to reside in the range  $[-\pi, \pi)$ . We only apply this to 2D environments.

The model is trained using sub-trajectories  $T_n = (s_t, a_t, s_{t+1}, a_{t+1}, \dots, s_{t+n-1}, a_{t+n-1}, s_{t+n})$ , where  $n$  is the model training window in Table 7(b). We optimize the model parameters  $\omega$  with sub-trajectories using the following equation

$$\nabla_\omega \mathbb{E}_{T_n \sim \mathcal{R}} \left[ \frac{1}{n+1} \sum_{k=0}^n -\log \text{EMIT}_\omega \left( s_{t+k} | \text{STEP}_\omega^k (\text{EMBED}_\omega(s_t), a_t, \dots, a_{t+k-1}) \right) \right] \quad (8)$$

where for  $k = 0$ , we just evaluate the embedder as  $\text{EMIT}_\omega(s_t | \text{EMBED}_\omega(s_t))$ .

We set the prediction length  $L$  in ACDA to the assumed upper bound of the delay process in all benchmarks. For each baseline, the horizon  $h$  of the action buffer is also set to the assumed upper bound of the delay process, unless the baseline uses constant-delay augmentation (CDA), in which the assumed delay used for the CDA is also used for the action buffer horizon  $h$ .

## B.5 PRACTICAL EVALUATION DETAILS

The evaluation methodology used for all performance measurements is that of the maximum average return for a trained policy. A policy is trained on a total of  $10^6$  steps from the environment. Every 10000 training steps, all network weights are frozen, and the policy is evaluated by sampling 10 trajectories from the environment. These trajectories are discarded after evaluation and not used for training. Each policy is evaluated on the same underlying environment, delay process, and noise process as it was trained on. The average return (unweighted average sum of rewards  $\frac{1}{10} \sum_{i=1}^{10} \sum_{(s_t, a_t, r_t) \in \tau_i} r_t$ ) is reported as the performance of the policy. For time series plots with an x-axis *Steps* and a y-axis *Return*, we refer to the average return evaluated after that number of training steps. The maximum average return, as shown in the tables, is the maximum evaluated average return achieved at any point during the training process.

---

1242 The training algorithms for SAC, SAC /w CDA, BPQL, and ACDA are implemented in our own  
1243 framework. Dreamer, DCAC, and VDPO are evaluated using the authors’ own implementations<sup>567</sup>,  
1244 with small modifications to accommodate our delay processes, noise processes, and the interaction  
1245 layer.

1246 In our framework, we use PyTorch for deep learning functionality and Gymnasium for RL func-  
1247 tionality. The interaction layer is implemented as a Gymnasium environment wrapper, based on the  
1248 formalism described in Appendix C, that extends the Gymnasium API to support action and observa-  
1249 tion packets. The constant-delay augmentation (CDA) in Appendix A is implemented as a wrapper  
1250 on top of the interaction layer wrapper, reducing the API back to the original Gymnasium API. We  
1251 implement ACDA to explicitly make use of the extended interaction layer definitions, whereas we  
1252 implement BPQL and SAC w/ CDA to operate directly on the regular Gymnasium API using the  
1253 CDA wrapper. A pass-through wrapper for the interaction layer is used for evaluating SAC (without  
1254 CDA), where the action packet is filled with the provided action.

1255 All dependencies for our framework are provided as conda YAML files.

1256 For VDPO we reuse the code artifact from their original article. Modifications to their implementa-  
1257 tion include the addition of action noise, our interaction layer wrapper (with CDA), and additional  
1258 statistical reporting. These modifications are documented in the artifact. The reason for adding our  
1259 own interaction layer wrapper to VDPO is to capture effects from the M/M/1 delay process when  
1260 the maximum delay assumption is violated.

1261 The Dreamer baseline uses the Dreamer v3 implementation. We add the action noise and the in-  
1262 teraction layer wrappers on top of the environment, with the pass-through wrapper used to allow it  
1263 to operate using the regular Gymnasium API. This follows a similar procedure used by Karamzade  
1264 et al. (2024).

1265 As DCAC already has a framework for random delays, we do not add our interaction layer wrapper  
1266 to their implementation. Instead, we only add the action noise and the delay processes. To adapt our  
1267 single delay to their split observation and action delay, we set the observation delay to 0 and set the  
1268 sampled delay as the action delay. This matches the formalism in the interaction layer framework,  
1269 as we only consider the full round-trip delay.

1270 Each benchmark, meaning a single algorithm training on a single environment with a single delay  
1271 process, is run using a single Nvidia A40 GPU and 16 CPU cores. ACDA and VDPO benchmarks  
1272 take around 18-24 hours each to complete. BPQL and SAC benchmarks take around 6 hours each  
1273 to complete. Each Dreamer and DCAC benchmark takes roughly 3-5 days to complete.

1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293

---

<sup>5</sup><https://github.com/danijar/dreamerv3>

<sup>6</sup><https://github.com/rmst/rlrd>

<sup>7</sup><https://github.com/QingyuanWuNothing/VDPO>

---

## C FORMAL DESCRIPTION OF THE INTERACTION LAYER POMDP

This section describes the POMDP of the interaction layer introduced in Section 3.3. The POMDP formalizes the interaction between the agent and the interaction layer that wraps the underlying system. We assume that the underlying system can be described by an MDP  $\mathcal{M} = (S, A, r, p, \mu)$  where  $S$  is the state space,  $A$  the action space,  $r(s, a)$  the reward function,  $p(s'|s, a)$  the transition distribution, and  $\mu(s)$  the initial state distribution.

The interaction layer wraps the MDP  $\mathcal{M}$ , where the interaction delay is described by the delay process  $D$ . Samples  $d \sim D(\cdot)$  are not necessarily independent. To fully define the POMDP of the interaction layer, we also need the action buffer horizon  $h$  as well as the default action  $a_{\text{init}}$ . Given this information, the POMDP is described as the tuple  $\mathcal{P} = (S, A, p, r, \mu, \Omega, O)$ . We use the notation  $s \in S$ ,  $a \in A$ , and  $o \in \Omega$  to denote members of these sets. We also refer to items  $a \in A$  as *action packets* and items  $o \in \Omega$  as *observation packets*.

An observation is described as the tuple  $o = (t, s, \mathbf{b}, \delta, c)$ . This describes the state at the interaction layer at time  $t$ , where

- $t$  is the time at the interaction layer when the observation was generated,
- $s$  the underlying system state observed at the same time step,
- $\mathbf{b} = (b_1, b_2, \dots, b_h)$  are action buffer contents at time  $t$  ( $b_1$  is immediately applied to  $s$ ),
- $\delta$  is the delay of the action packet used to update the action buffer  $\mathbf{b}$ , and
- $c$  is the number of time steps without a new action packet replacing the action buffer contents.

When referring to the state of the action buffer at different time steps, e.g.,  $\mathbf{b}_t$  and  $\mathbf{b}_{t+1}$ , we use the notation  $b_{t,i}$  and  $b_{t+1,i}$  to refer to the  $i$ -th action in  $\mathbf{b}_t$  and  $\mathbf{b}_{t+1}$ , respectively.

Delays are referred to using different notations,  $d_t$  or  $\delta_t$ , depending on the context:

- $d_t$  is the unobserved delay sampled at time  $t$ . This is the delay of the action packet  $\mathbf{a}_t$ .
- $\delta_t$  is the delay recovered in hindsight. See description of the observation packet above for more information. This hindsight delay is related to  $d_t$  by the equation  $\delta_t = d_t - (\delta_t + c_t)$ .

The individual components of  $\mathcal{P}$  are defined in Equations 9-18:

$$\text{Action space } \mathbf{A} = \mathbb{N} \times \bigcup_{k=1}^{\infty} A^{k \times h} \quad (9)$$

$$\text{Observation space } \Omega = \mathbb{N} \times S \times A^h \times \mathbb{Z}^+ \times \mathbb{N} \quad (10)$$

$$\text{State space } \mathbf{S} = \Omega \times 2^{(\mathbb{Z}^+ \times \mathbf{A})} \quad (11)$$

$$\text{Initial state distribution } \boldsymbol{\mu}(s) = \begin{cases} \mu(s) & \text{if } s = ((0, s, (a_{\text{init}}, \dots), 1, 0), \emptyset) \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

$$\text{Observation distribution } O(o|s) = \begin{cases} 1 & \text{if } s = (o, \mathcal{I}) \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

$$\text{Reward function } r(s_t, \mathbf{a}_t) = r(s_t, b_1) \quad (14)$$

$$\text{where } s_t = ((t, s_t, (b_1, b_2, \dots, b_h), \delta_t, c_t), \mathcal{I}_t)$$

Equation 11 defines states as tuples  $s_t = (o_t, \mathcal{I}_t)$ , where  $o_t$  is the state of the interaction layer (observable) and  $\mathcal{I}_t$  is the set of action packets in transit (not observable). The transit set  $\mathcal{I}_t \in 2^{(\mathbb{Z}^+ \times \mathbf{A})}$  contains tuples of action packets and their arrival time.

Note that, by the definition of  $\boldsymbol{\mu}$  in Equation 12, we can always check if the action buffer contains the initial actions by  $t - (\delta_t + c_t) < 0$ . This holds until the first action packet is received.

The transition dynamics  $p(s_{t+1}|s_t, \mathbf{a}_t)$  are described in Equation 18 below. While this is simple to describe in text and with examples, it becomes complicated to define formally. We first define a couple of auxiliary functions below to help define the transition dynamics. We define the function

TRANSMIT( $\mathcal{I}_t, \mathbf{a}_t, d$ ) that adds the action  $\mathbf{a}_t$  with delay  $d$  to the transit set, together with the  $\min \mathcal{I}$  and  $\min_t \mathcal{I}$  operations to get the action packet with the nearest arrival:

$$\text{TRANSMIT}(\mathcal{I}_t, \mathbf{a}_t, d) = \{(t + d, \mathbf{a}_t)\} \cup \{(t', \mathbf{a}') \in \mathcal{I}_t : t' < t + d\} \quad (15)$$

$$\min_t \mathcal{I} = \min\{t' : (t', \mathbf{a}') \in \mathcal{I}\} \quad (16)$$

$$\min \mathcal{I} = \begin{cases} \emptyset & \text{if } \mathcal{I} = \emptyset \\ (t', \mathbf{a}') \in \mathcal{I} & \text{if } t' = \min_t \mathcal{I} \end{cases} \quad (17)$$

One aspect of the behavior of the interaction layer, modeled by the TRANSMIT( $\mathcal{I}_t, \mathbf{a}_t, d$ ) function in Equation 15, is that outdated action packets arriving at the interaction layer will be discarded. For example, if  $\mathbf{a}_t$  has delay  $d_t = 4$ , and  $\mathbf{a}_{t+1}$  has delay  $d_{t+1} = 2$ , then  $\mathbf{a}_{t+1}$  will arrive at time  $t + 3$ , whereas  $\mathbf{a}_t$  will arrive after at time  $t + 4$ . When  $\mathbf{a}_t$  arrives, the interaction layer will see that the contents of the action buffer are based on information from  $\mathbf{o}_{t+1}$ , whereas the action packet  $\mathbf{a}_t$  is based on information from  $\mathbf{o}_t$ . Therefore,  $\mathbf{a}_t$  is considered outdated and will be discarded. Also note that a consequence of this is that  $d_t$  will never be observed, not even in hindsight.

Using these functions, we define the transition probabilities below in Equation 18. The probabilities themselves are simple to describe as  $p(s_{t+1}|s_t, b_1) \times D(d)$ ; the complexity arises from checking that the new POMDP state is compatible with the possible sampled delays. The first case covers when no new action packet arrives at the interaction layer at time  $t + 1$ , the second case is when the received action packet  $\mathbf{a}_t$  has too few rows in the matrix to update the action buffer for the sampled delay  $d$ , and the third case is when a received action packet is used to update the action buffer.

$$p(s_{t+1}|s_t, \mathbf{a}_t) = \begin{cases} p(s_{t+1}|s_t, b_1) \cdot D(d) & \text{if } \mathcal{I}_{t+1} = \text{TRANSMIT}(\mathcal{I}_t, \mathbf{a}_t, d) \wedge \min_t \mathcal{I}_{t+1} > t + 1 \wedge \\ & c_{t+1} = c_t + 1 \wedge \delta_{t+1} = \delta_t \wedge \\ & \mathbf{b}_{t+1} = (b_2, b_3, \dots, b_{h-1}, b_h, b_h) \\ p(s_{t+1}|s_t, b_1) \cdot D(d) & \text{if } \mathcal{I}_{t+1} = \mathcal{I}_{\text{cand}} \setminus \{\min \mathcal{I}_{\text{cand}}\} \wedge \min_t \mathcal{I}_{\text{cand}} = t + 1 \wedge \\ & (t + 1 - u) > L \wedge c_{t+1} = c_t + 1 \wedge \delta_{t+1} = \delta_t \wedge \\ & \mathbf{b}_{t+1} = (b_2, b_3, \dots, b_{h-1}, b_h, b_h) \\ & \text{where } \mathcal{I}_{\text{cand}} = \text{TRANSMIT}(\mathcal{I}_t, \mathbf{a}_t, d) \\ & (t + 1, \mathbf{a}_u) = \min \mathcal{I}_{\text{cand}} \\ & (u, M^u) = \mathbf{a}_u \\ & M^u \in A^{L \times h} \\ p(s_{t+1}|s_t, b_1) \cdot D(d) & \text{if } \mathcal{I}_{t+1} = \mathcal{I}_{\text{cand}} \setminus \{\min \mathcal{I}_{\text{cand}}\} \wedge \min_t \mathcal{I}_{\text{cand}} = t + 1 \wedge \\ & (t + 1 - u) \leq L \wedge c_{t+1} = 0 \wedge \delta_{t+1} = (t + 1 - u) \wedge \\ & \mathbf{b}_{t+1} = M_{(t+1-u)}^u \\ & \text{where } \mathcal{I}_{\text{cand}} = \text{TRANSMIT}(\mathcal{I}_t, \mathbf{a}_t, d) \\ & (t + 1, \mathbf{a}_u) = \min \mathcal{I}_{\text{cand}} \\ & (u, M^u) = \mathbf{a}_u \\ & M^u \in A^{L \times h} \\ 0 & \text{otherwise} \end{cases} \quad (18)$$

where  $s_t = (\mathbf{o}_t, \mathcal{I}_t)$

$$s_{t+1} = (\mathbf{o}_{t+1}, \mathcal{I}_{t+1})$$

$$\mathbf{o}_t = (t, s_t, \mathbf{b}_t, \delta_t, c_t)$$

$$\mathbf{o}_{t+1} = (t + 1, s_{t+1}, \mathbf{b}_{t+1}, \delta_{t+1}, c_{t+1})$$

$$\mathbf{b}_t = (b_1, b_2, \dots, b_{h-2}, b_{h-1}, b_h)$$

## D DETAILED MODEL DESCRIPTION

This section provides formal definitions of the model introduced in Section 4.2, as well as a detailed definition of the training algorithm presented in Section 4.3. Appendix D.1 presents the formal model definition. Appendix D.2 presents the full training algorithm.

We use the variables  $\theta$ ,  $\phi$ , and  $\omega$  to denote the parameters of the policy, critic, and model, respectively. In practice, these are large vectors of real numbers where different parts of the vector contain the parameters for components in a deep neural network.

### D.1 MODEL COMPONENTS AND OBJECTIVE

The primary purpose of the model is to overcome the limitation on fixed-size inputs of MLPs. The idea is that, instead of generating actions directly with the augmented state input:

$$a_{t+k} \sim \pi_{\theta}(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1}), \quad (19)$$

we generate actions using the distribution over the state that the action will be applied to as policy input:

$$a_{t+k} \sim \pi_{\theta}(\cdot | p(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1})), \quad (20)$$

where  $p(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1})$  represents the distribution over states after applying the action sequence  $a_t, a_{t+1}, \dots, a_{t+k-1}$ , in order, to the state  $s_t$ . We represent  $p(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1})$  as a fixed-size latent representation, and thanks to this representation, we can generate actions with MLPs for variable-size inputs. The purpose of the model is to create these embeddings (defining the mapping between  $p(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1})$  and the corresponding latent representation). Since this kind of policy makes decisions using a distribution over the state, and that the distribution is embedded as a latent representation using a model, we refer to agents using this kind of policy as *model-based distribution agents* (MDA).

The model consists of three components:  $\text{EMBED}_{\omega}(s_t)$ ,  $\text{STEP}_{\omega}(z_i, a_{t+i})$ , and  $\text{EMIT}_{\omega}(\hat{s}_{t+i} | z_{t+i})$ .

$\text{EMBED}_{\omega}(s_t)$  embeds the state  $s_t$  into a latent representation  $z_0$ . In a perfect model,  $z_0$  would be an embedding of the Dirac delta distribution  $\delta(x - s_t)$ .

$\text{STEP}_{\omega}(z_i, a_{t+i})$  updates the latent representation  $z_i$  to include information about what happens if the action  $a_{t+i}$  is also applied. Such that, if  $z_i$  is a latent representation of  $p(\cdot | s_t, a_t, \dots, a_{t+i-1})$ , then  $z_{i+1} = \text{STEP}_{\omega}(z_i, a_{t+i})$  is a latent representation of  $p(\cdot | s_t, a_t, \dots, a_{t+i-1}, a_{t+i})$ .

$\text{EMIT}_{\omega}(\hat{s}_{t+i} | z_{t+i})$  converts the latent representation  $z_{t+i}$  back to a regular parameterized distribution. We use a normal distribution in our model, where  $\text{EMIT}_{\omega}$  outputs the mean and standard deviation for each component of the MDP state. We never sample from this distribution. This component is only used to ensure that we have a good latent representation.

To make the notation more compact, we use the multi-step notation  $\text{STEP}_{\omega}^k$  where

$$\text{STEP}_{\omega}^0(z) = z \quad (21)$$

$$\text{STEP}_{\omega}^k(z, a_0, a_1, \dots, a_{k-1}) = \text{STEP}_{\omega}^{k-1}(\text{STEP}_{\omega}(z, a_0), a_1, \dots, a_{k-1}) \quad (22)$$

With this notation, we say that  $\text{STEP}_{\omega}^k(\text{EMBED}_{\omega}(s_t), a_t, a_{t+1}, \dots, a_{t+k-1})$  embeds the distribution  $p(\cdot | s_t, a_t, a_{t+1}, \dots, a_{t+k-1})$ . We optimize the model to minimize the KL-divergence between the embedded distribution and the true distribution. That is

$$\min_{\omega} D_{\text{KL}}(p(\cdot | s_t, a_t, \dots, a_{t+n-1}) \| \text{EMIT}_{\omega}(\cdot | \text{STEP}_{\omega}^n(\text{EMBED}_{\omega}(s_t), a_t, \dots, a_{t+n-1}))) \quad (23)$$

for all possible states  $s_t$  and sequences of actions  $a_t, \dots, a_{t+n-1}$ . The loss function  $\mathcal{L}(\omega)$  from Section 4.2 is a Monte Carlo estimate of this objective:

$$\mathcal{L}(\omega) = \mathbb{E}_{(s_t, a_t, a_{t+1}, \dots, a_{t+n-1}, s_{t+n}) \sim \mathcal{R}} [-\log \text{EMIT}_{\omega}(s_{t+n} | z_n)] \quad (24)$$

where  $z_n = \text{STEP}_{\omega}^n(\text{EMBED}_{\omega}(s_t), a_t, a_{t+1}, \dots, a_{t+n-1})$

and  $\mathcal{R}$  is a replay buffer with experiences collected online.

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

## D.2 TRAINING ALGORITHM

This section presents the full version of the training algorithm from Section 4.3. As with SAC, we assume that  $\pi_\theta$  is represented as a reparameterizable policy that is a deterministic function with independent noise input.

---

### Algorithm 5 Actor-Critic with Delay Adaptation

---

- 1: Initialize policy  $\pi_\theta$ , critics  $Q_{\phi_1}, Q_{\phi_2}$ , model  $\omega$ , temperature  $\alpha$ , target networks  $\phi'_1, \phi'_2$ , and replay  $\mathcal{R}$
- 2: **for** each epoch **do**
- 3:   *// Stage 1: Sample trajectory*
- 4:   Collected trajectory:  $\mathcal{T} = \emptyset$
- 5:    $t \leftarrow 0$
- 6:   Reset interaction layer state:  $s_0 \sim \mu$ , Observe  $\mathbf{o}_0$
- 7:   **while** terminal state not reached **do**
- 8:      $(t, s_t, \mathbf{b}_t, \delta_t, c_t) = \mathbf{o}_t$
- 9:     **for**  $k \leftarrow 1$  to  $L$  **do**
- 10:       Select  $\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k}$  by Algorithm 1
- 11:        $y_0 \leftarrow (\hat{a}_1^{t+k}, \dots, \hat{a}_k^{t+k})$
- 12:       **for**  $i \leftarrow 1$  to  $h$  **do**
- 13:           $a_i^{t+k} \sim \pi_\theta(\cdot | \text{STEP}_\omega^{k+i-1}(\text{EMBED}_\omega(s_t), y_{i-1}))$
- 14:           $y_i \leftarrow (y_{i-1}, a_i^{t+k})$
- 15:       
$$\mathbf{a}_t \leftarrow \left( t, \begin{bmatrix} a_1^{t+1} & a_2^{t+1} & a_3^{t+1} & \dots & a_h^{t+1} \\ a_1^{t+2} & a_2^{t+2} & a_3^{t+2} & \dots & a_h^{t+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_1^{t+L} & a_2^{t+L} & a_3^{t+L} & \dots & a_h^{t+L} \end{bmatrix} \right)$$
- 16:       Send  $\mathbf{a}_t$  to interaction layer, observe  $r_t, \mathbf{o}_{t+1}, \Gamma_{t+1}$
- 17:       Add  $(\mathbf{o}_t, \mathbf{a}_t, r_t, \mathbf{o}_{t+1}, \Gamma_{t+1})$  to  $\mathcal{T}$
- 18:        $t \leftarrow t + 1$
- 19:     *// Stage 2: Reconstruct transition info*
- 20:     **for**  $(\mathbf{o}_i, \mathbf{a}_i, r_i, \mathbf{o}_{i+1}, \Gamma_{i+1}) \in \mathcal{T}$  **do**
- 21:        $(i, s_i, \mathbf{b}_i, \delta_i, c_i) = \mathbf{o}_i, \quad (i+1, s_{i+1}, \mathbf{b}_{i+1}, \delta_{i+1}, c_{i+1}) = \mathbf{o}_{i+1}$
- 22:        $a_i = b_{i,1}$
- 23:       **if**  $i - (\delta_i + c_i) \geq 0$  **then**
- 24:          *// (Recover the input used by  $\pi_\theta$  and  $\text{STEP}_\omega^k$  to generate  $b_{i,1}$  and  $b_{i+1,1}$ )*
- 25:           $j \leftarrow i - (\delta_i + c_i), \quad j' \leftarrow (i+1) - (\delta_{i+1} + c_{i+1})$
- 26:          Reconstruct  $\hat{a}_1^{j+\delta_i}, \dots, \hat{a}_{\delta_i}^{j+\delta_i}$ , choose  $a_1^{j+\delta_i}, \dots, a_{c_i}^{j+\delta_i}$  from  $\mathbf{a}_j$
- 27:          Reconstruct  $\hat{a}_1^{j'+\delta_{i+1}}, \dots, \hat{a}_{\delta_{i+1}}^{j'+\delta_{i+1}}$ , choose  $a_1^{j'+\delta_{i+1}}, \dots, a_{c_{i+1}}^{j'+\delta_{i+1}}$  from  $\mathbf{a}_{j'}$
- 28:           $y_i \leftarrow (\hat{a}_1^{j+\delta_i}, \dots, \hat{a}_{\delta_i}^{j+\delta_i}, a_{1,1}^{j+\delta_i}, \dots, a_{c_i}^{j+\delta_i})$
- 29:           $y_{i+1} \leftarrow (\hat{a}_1^{j'+\delta_{i+1}}, \dots, \hat{a}_{\delta_{i+1}}^{j'+\delta_{i+1}}, a_1^{j'+\delta_{i+1}}, \dots, a_{c_{i+1}}^{j'+\delta_{i+1}})$
- 30:          *// We denote their lengths as  $|y_i| = \delta_i + c_i$*
- 31:          Add  $(s_i, a_i, r_i, s_{i+1}, \Gamma_{i+1}, y_i, y_{i+1})$  to  $\mathcal{R}$
- 32:     *// Stage 3: Update network weights*
- 33:     **for**  $|\mathcal{T}|$  sampled batches of  $(s, a, r, s', \Gamma, y, y')$  from  $\mathcal{R}$  **do**
- 34:       *// These are computed in expectation of samples from  $\mathcal{R}$*
- 35:        $\hat{a}' \sim \pi_\theta(\cdot | \text{STEP}_\omega^{|y'|}(\text{EMBED}_\omega(s'), y'))$
- 36:        $x = r + \gamma(1 - \Gamma)(\min(Q_{\phi'_1}(s', \hat{a}'), Q_{\phi'_2}(s', \hat{a}')) - \alpha \log \pi_\theta(\hat{a}' | \dots))$
- 37:       Do gradient descent step on  $\nabla_{\phi_1}(Q_{\phi_1}(s, a) - x)^2$  and  $\nabla_{\phi_2}(Q_{\phi_2}(s, a) - x)^2$
- 38:        $\hat{a} \sim \pi_\theta(\cdot | \text{STEP}_\omega^{|y|}(\text{EMBED}_\omega(s), y))$
- 39:       Do gradient ascent step on  $\nabla_\theta(\min(Q_{\phi_1}(s, \hat{a}), Q_{\phi_2}(s, \hat{a})) - \alpha \log \pi_\theta(\hat{a} | \dots))$
- 40:       Update  $\alpha$  according to SAC
- 41:       Compute  $\nabla_\omega$  by Equation 8 and do gradient descent step
- 42:       Update target networks  $\phi'_1, \phi'_2$

---

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

---

## E ADDITIONAL RESULTS

This section presents additional results to complement those presented in Section 5. Appendix E.1 presents the results from Table 1 as time series plots, showing how the mean and standard deviation of the evaluated return change over the training process. In subsequent tables, the standard deviation is shown following the  $\pm$  symbol.

Appendix E.2 and E.3 answer two questions that are not part of the evaluation in Section 5. The questions that these appendices answer are:

- Appendix E.2 answers the question whether the gain in performance is due to the model-based policy introduced in Section 4.2 or due to the adaptiveness of ACDA. We evaluate this by modifying the BPQL algorithm such that it uses the MDA policy instead of a direct MLP policy, and compare how that performs against ACDA. The results clearly show that it is the adaptiveness of the interaction layer that is the reason for the strong performance of ACDA.
- Appendix E.3 answers the question whether acting with CDA under a different horizon  $h$  than the worst-case delay is better than trying to adapt to varying delays. We set up this demonstration by applying CDA with a horizon  $h$  that is greater than or equal to the sampled delay most of the time, but occasionally a sampled delay exceeds this horizon. This kind of CDA does not result in a constant-delay MDP that BPQL and VDPO expect, hence we did not include this in the main evaluation. The results in Appendix E.3 show that while VDPO and BPQL occasionally gain performance in this setting, ACDA is still the best performing algorithm in most cases. ACDA is always close to the highest performing algorithm in the few cases where ACDA does not achieve the best mean return. These results show that the adaptiveness of the interaction layer provides an increase in performance that cannot be achieved through constant-delay approaches.

In addition to the results from Appendix E.3, we also evaluate when the horizon  $h$  is based on the average sampled delay under each evaluated delay process. These average delay results are presented in Appendix E.4. Using the average delay as the  $h$  never results in the best average return for a benchmark. However, in a few instances, the average delay does result in the best performance amongst different horizon values within a specific constant-delay algorithm.

The results from all evaluated baselines are presented in Appendix E.5.

### E.1 TIME SERIES RESULTS

This section presents the results from the evaluation in Section 5 as time series plots, including standard deviation bands. Unless otherwise specified, the evaluation methodology follows that described in Section 5. To reduce noise and highlight trends, each time series is smoothed using a running average over five evaluation points.

The results in this appendix are split into the delay processes used. Appendix E.1.1 presents results for the  $GE_{1,23}$  delay process, Appendix E.1.2 for the  $GE_{4,32}$  delay process, and Appendix E.1.3 for the M/M/1 queue delay process.

1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

### E.1.1 PERFORMANCE EVALUATION UNDER THE $GE_{1,23}$ DELAY PROCESS

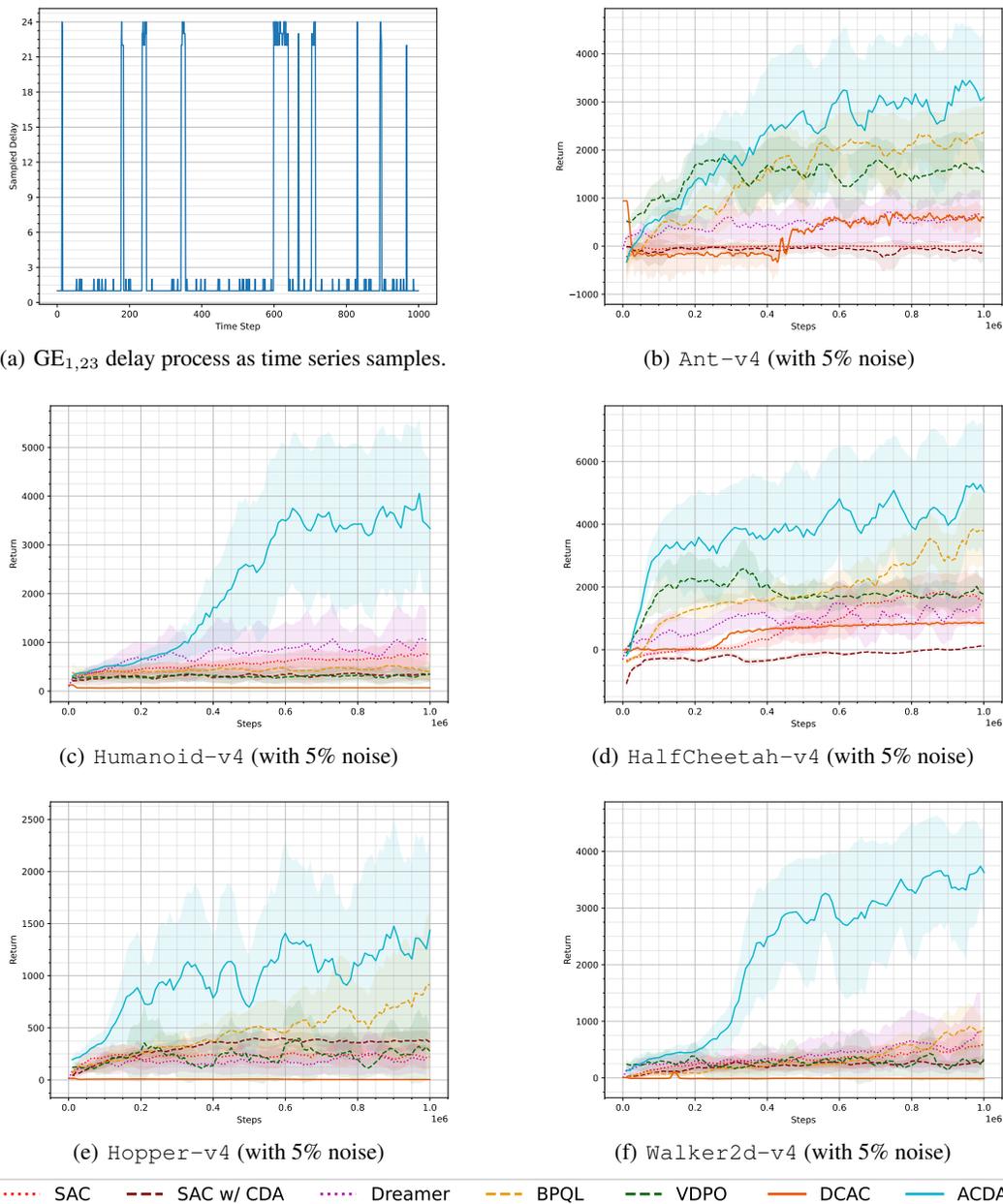
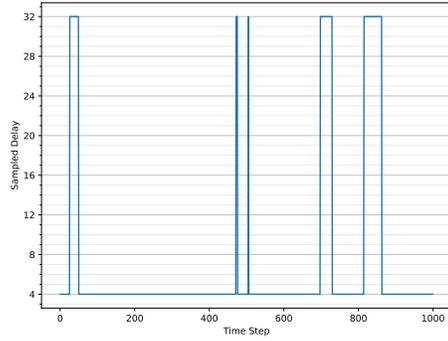


Figure 13: Time series evaluation during training on the  $GE_{1,23}$  delay process. All environments have added 5% noise to the actions.

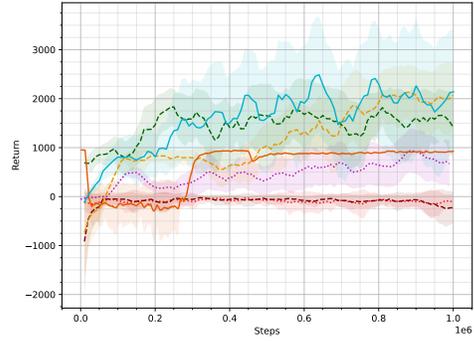
Table 8: Best returns from the  $GE_{1,23}$  delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	14.22 ± 14.89	862.18 ± 266.21	2064.18 ± 223.48	306.91 ± 51.26	708.33 ± 221.53
SAC w/ CDA	69.28 ± 114.47	414.05 ± 204.20	128.47 ± 9.55	426.92 ± 27.93	428.44 ± 509.44
Dreamer	1111.73 ± 412.35	1463.07 ± 649.56	1796.07 ± 381.47	334.30 ± 245.42	1081.12 ± 905.94
BPQL	2691.88 ± 129.84	585.19 ± 163.49	4320.20 ± 1028.52	1328.71 ± 937.67	1215.91 ± 776.93
VDPO	2163.00 ± 53.04	417.25 ± 210.09	3144.23 ± 1156.52	709.20 ± 522.01	846.88 ± 808.67
DCAC	949.97 ± 11.87	128.47 ± 36.09	920.09 ± 33.05	16.99 ± 15.94	106.70 ± 53.84
<b>ACDA</b>	<b>4112.78 ± 818.44</b>	<b>4608.76 ± 1084.52</b>	<b>5984.25 ± 1885.78</b>	<b>2094.65 ± 944.20</b>	<b>3863.59 ± 232.52</b>

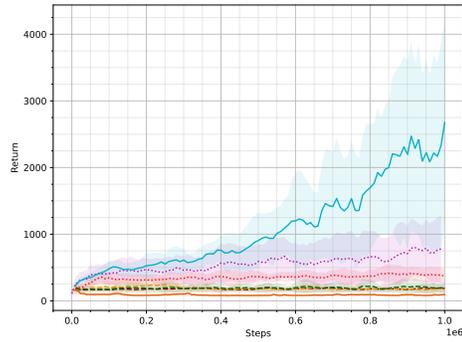
## E.1.2 PERFORMANCE EVALUATION UNDER THE $GE_{4,32}$ DELAY PROCESS



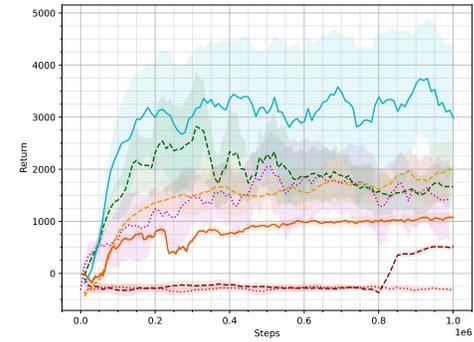
(a)  $GE_{4,32}$  delay process as time series samples.



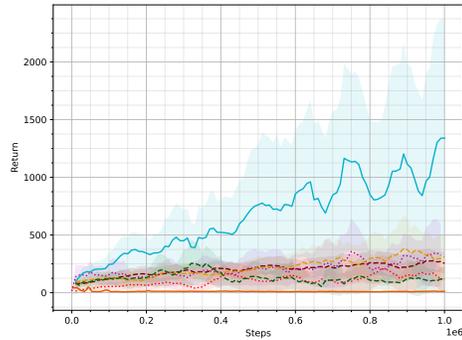
(b) Ant-v4 (with 5% noise)



(c) Humanoid-v4 (with 5% noise)



(d) HalfCheetah-v4 (with 5% noise)



(e) Hopper-v4 (with 5% noise)



(f) Walker2d-v4 (with 5% noise)



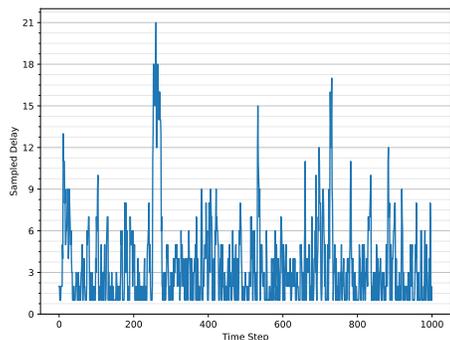
Figure 14: Time series evaluation during training on the  $GE_{4,32}$  delay process. All environments have added 5% noise to the actions.

Table 9: Best returns from the  $GE_{4,32}$  delay process.

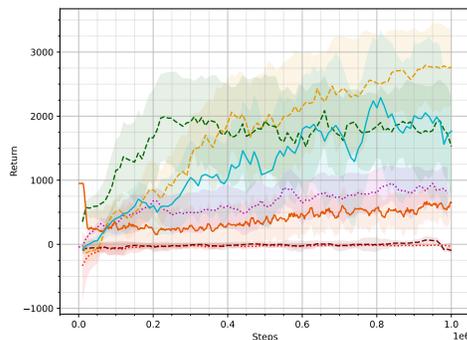
	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	-5.72 ± 19.62	494.43 ± 156.01	-158.78 ± 65.86	279.74 ± 109.83	60.86 ± 75.59
SAC w/ CDA	18.93 ± 23.64	230.45 ± 99.22	591.32 ± 36.39	315.47 ± 51.49	257.18 ± 73.42
Dreamer	1147.56 ± 371.15	1091.48 ± 577.52	2493.19 ± 231.22	515.36 ± 430.72	1233.79 ± 802.47
BPQL	2509.52 ± 117.37	276.63 ± 131.70	2136.36 ± 547.04	433.29 ± 381.79	875.09 ± 747.72
VDPO	2266.99 ± 90.89	280.72 ± 169.85	3664.30 ± 929.25	330.44 ± 263.74	344.73 ± 316.82
DCAC	953.14 ± 12.06	167.97 ± 82.14	1123.47 ± 100.34	57.98 ± 28.04	9.23 ± 20.35
<b>ACDA</b>	<b>2866.93 ± 1172.46</b>	<b>3725.59 ± 1513.38</b>	<b>4231.15 ± 333.69</b>	<b>1727.79 ± 959.50</b>	<b>1840.58 ± 386.78</b>

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

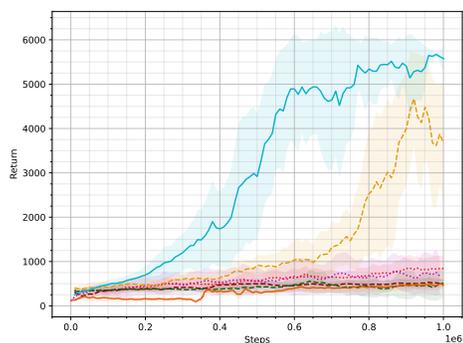
### E.1.3 PERFORMANCE EVALUATION UNDER THE M/M/1 QUEUE DELAY PROCESS



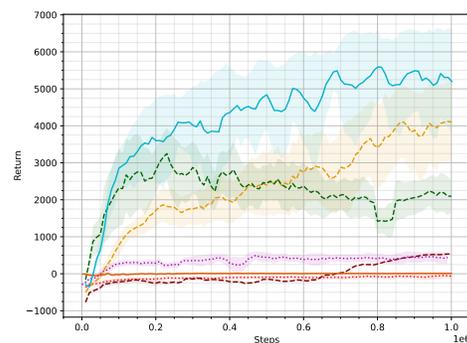
(a) MM1 delay process as time series samples.



(b) Ant-v4 (with 5% noise)



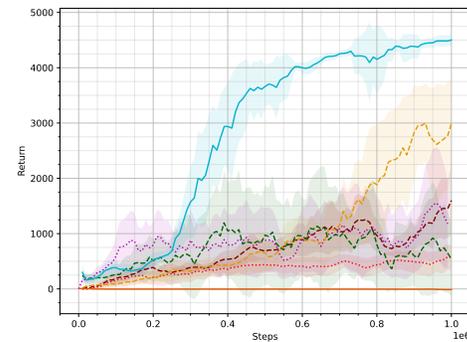
(c) Humanoid-v4 (with 5% noise)



(d) HalfCheetah-v4 (with 5% noise)



(e) Hopper-v4 (with 5% noise)



(f) Walker2d-v4 (with 5% noise)



Figure 15: Time series evaluation during training on the MM1 delay process. All environments have added 5% noise to the actions.

Table 10: Best returns from the MM1 delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	-0.58 ± 8.66	921.04 ± 299.47	20.69 ± 94.91	333.06 ± 96.04	604.80 ± 212.37
SAC w/ CDA	102.00 ± 33.77	613.03 ± 157.68	550.84 ± 16.28	627.59 ± 24.62	2005.76 ± 341.30
Dreamer	1121.11 ± 58.69	981.38 ± 597.01	584.40 ± 72.26	975.72 ± 650.05	1801.81 ± 1158.73
BPQL	<b>3074.17 ± 106.78</b>	5435.29 ± 68.34	4660.93 ± 448.10	3035.66 ± 103.80	3547.73 ± 133.51
VDPO	2528.67 ± 144.63	720.73 ± 634.35	3831.96 ± 960.07	1459.88 ± 933.11	2144.25 ± 1650.85
DCAC	959.23 ± 13.54	525.85 ± 135.36	35.60 ± 21.42	1026.45 ± 2.96	24.48 ± 45.46
ACDA	2898.46 ± 838.07	<b>5805.60 ± 23.04</b>	<b>5898.36 ± 409.10</b>	<b>3122.53 ± 417.37</b>	<b>4562.33 ± 87.98</b>

---

1728 E.2 MODEL-BASED DISTRIBUTION AGENT VS. ADAPTIVENESS  
1729

1730 The purpose of this Appendix is to answer the question of whether it is the model-based distribution  
1731 agent (MDA) or the adaptivity of ACDA that leads to its high performance. To answer this, we  
1732 modify the BPQL algorithm to use the MDA policy instead of the direct MLP policy that they used  
1733 in their original paper.

1734 It is necessary to modify the BPQL algorithm itself since the optimization of the MDA policy is  
1735 split into two steps, with different kinds of samples from the replay buffer. If the delay truly was  
1736 constant, then BPQL with MDA would be the same as the performance of ACDA, due to the perfect  
1737 conditions for the memorized action selection. However, it is necessary to split these into two  
1738 algorithms since this cannot capture the M/M/1 queue delay process, which cannot be represented  
1739 as a true constant-delay MDP.

1740 Like Appendix E.1, results are split based on the delay process used. Appendix E.2.1 presents results  
1741 for the  $GE_{1,23}$  delay process, Appendix E.2.2 for the  $GE_{4,32}$  delay process, and Appendix E.2.3 for  
1742 the M/M/1 queue delay process.

1743 These results show that, while BPQL sometimes performs better using the MDA policy, ACDA,  
1744 with its adaptivity, is still the best-performing algorithm.  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781

1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835

### E.2.1 PERFORMANCE OF MDA UNDER THE $GE_{1,23}$ DELAY PROCESS

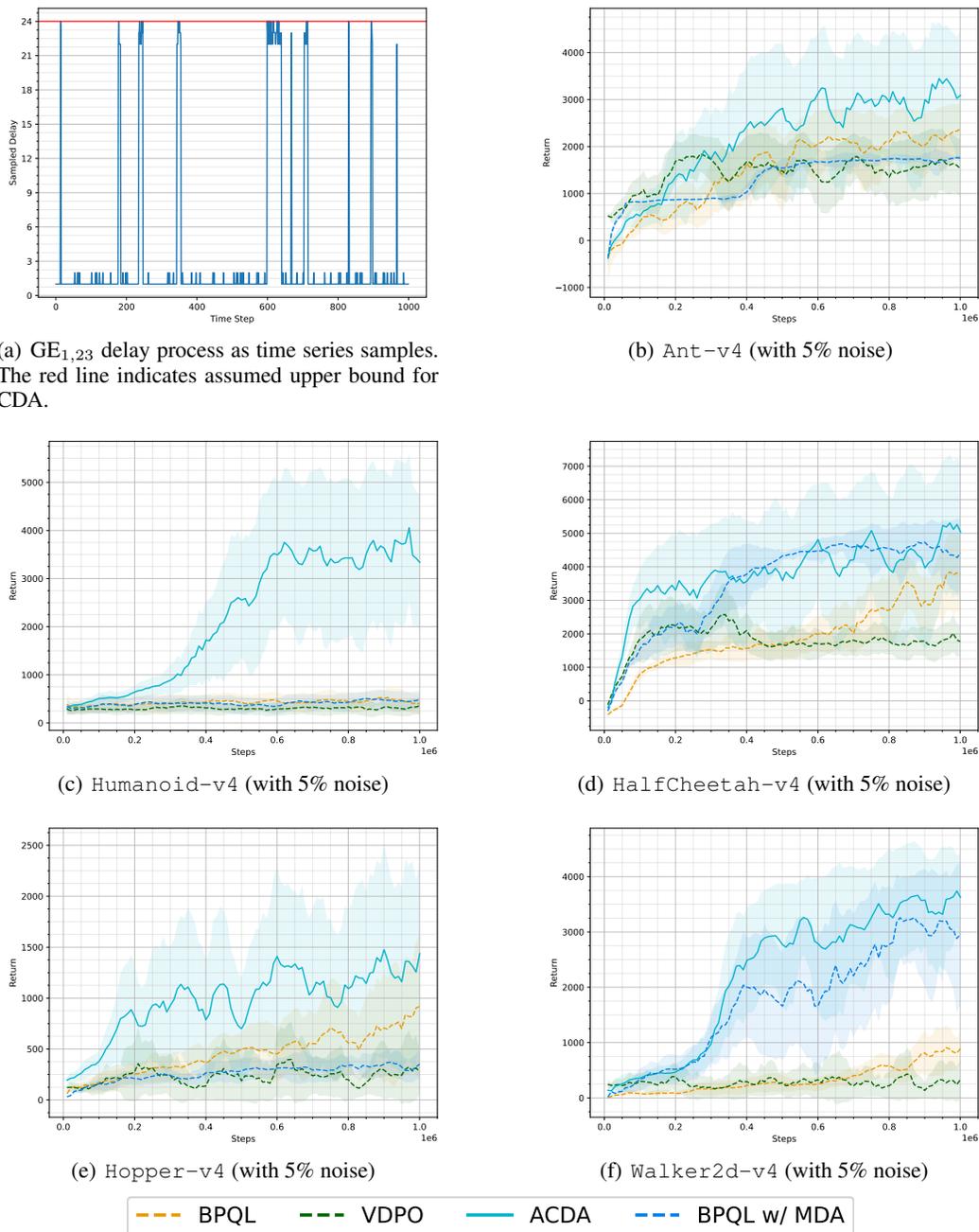


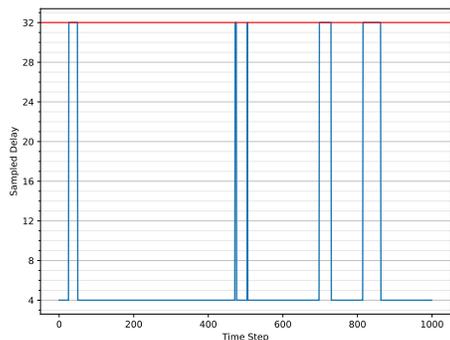
Figure 16: Time series evaluation during training on the  $GE_{1,23}$  delay process. All environments have added 5% noise to the actions.

Table 11: Best returns from the  $GE_{1,23}$  delay process.

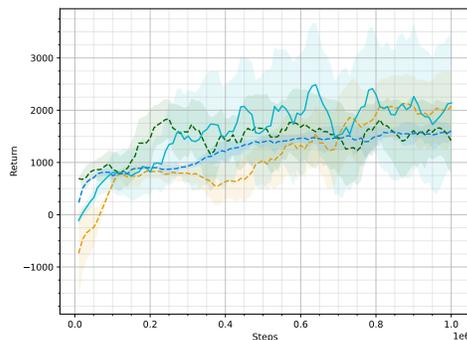
	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	2691.88 ± 129.84	585.19 ± 163.49	4320.20 ± 1028.52	1328.71 ± 937.67	1215.91 ± 776.93
VDPO	2163.00 ± 53.04	417.25 ± 210.09	3144.23 ± 1156.52	709.20 ± 522.01	846.88 ± 808.67
ACDA	<b>4112.78 ± 818.44</b>	<b>4608.76 ± 1084.52</b>	<b>5984.25 ± 1885.78</b>	<b>2094.65 ± 944.20</b>	<b>3863.59 ± 232.52</b>
BPQL w/ MDA	1795.29 ± 23.78	563.36 ± 96.11	4926.36 ± 60.08	465.14 ± 138.86	3681.39 ± 126.41

1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889

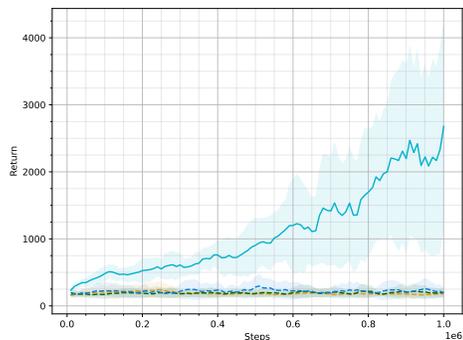
## E.2.2 PERFORMANCE OF MDA UNDER THE $GE_{4,32}$ DELAY PROCESS



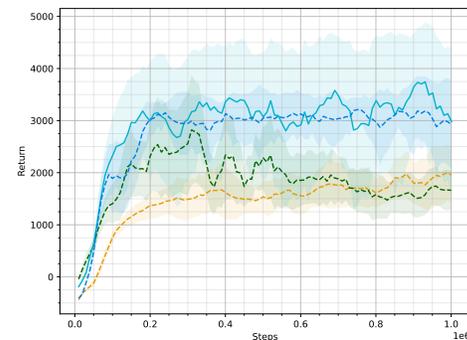
(a)  $GE_{4,32}$  delay process as time series samples. The red line indicates assumed upper bound for CDA.



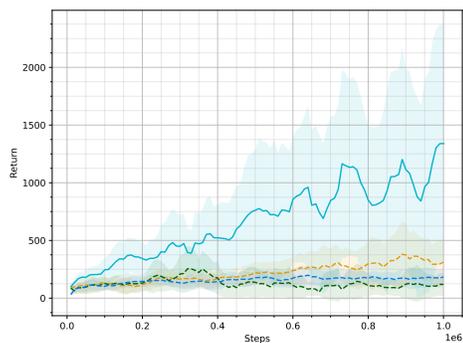
(b) Ant-v4 (with 5% noise)



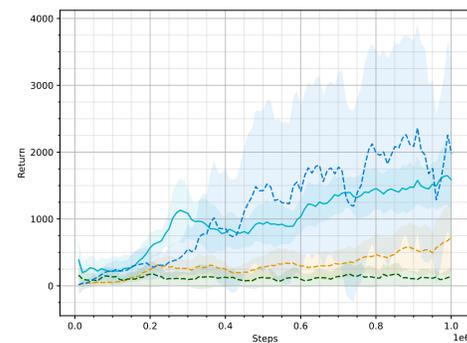
(c) Humanoid-v4 (with 5% noise)



(d) HalfCheetah-v4 (with 5% noise)



(e) Hopper-v4 (with 5% noise)



(f) Walker2d-v4 (with 5% noise)



Figure 17: Time series evaluation during training on the  $GE_{4,32}$  delay process. All environments have added 5% noise to the actions.

Table 12: Best returns from the  $GE_{4,32}$  delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	2509.52 ± 117.37	276.63 ± 131.70	2136.36 ± 547.04	433.29 ± 381.79	875.09 ± 747.72
VDPO	2266.99 ± 90.89	280.72 ± 169.85	3664.30 ± 929.25	330.44 ± 263.74	344.73 ± 316.82
ACDA	<b>2866.93 ± 1172.46</b>	<b>3725.59 ± 1513.38</b>	<b>4231.15 ± 333.69</b>	<b>1727.79 ± 959.50</b>	1840.58 ± 386.78
BPQL w/ MDA	1661.41 ± 43.42	359.46 ± 156.86	3609.45 ± 328.37	224.34 ± 90.12	<b>3015.45 ± 1428.05</b>

### E.2.3 PERFORMANCE OF MDA UNDER THE M/M/1 QUEUE DELAY PROCESS

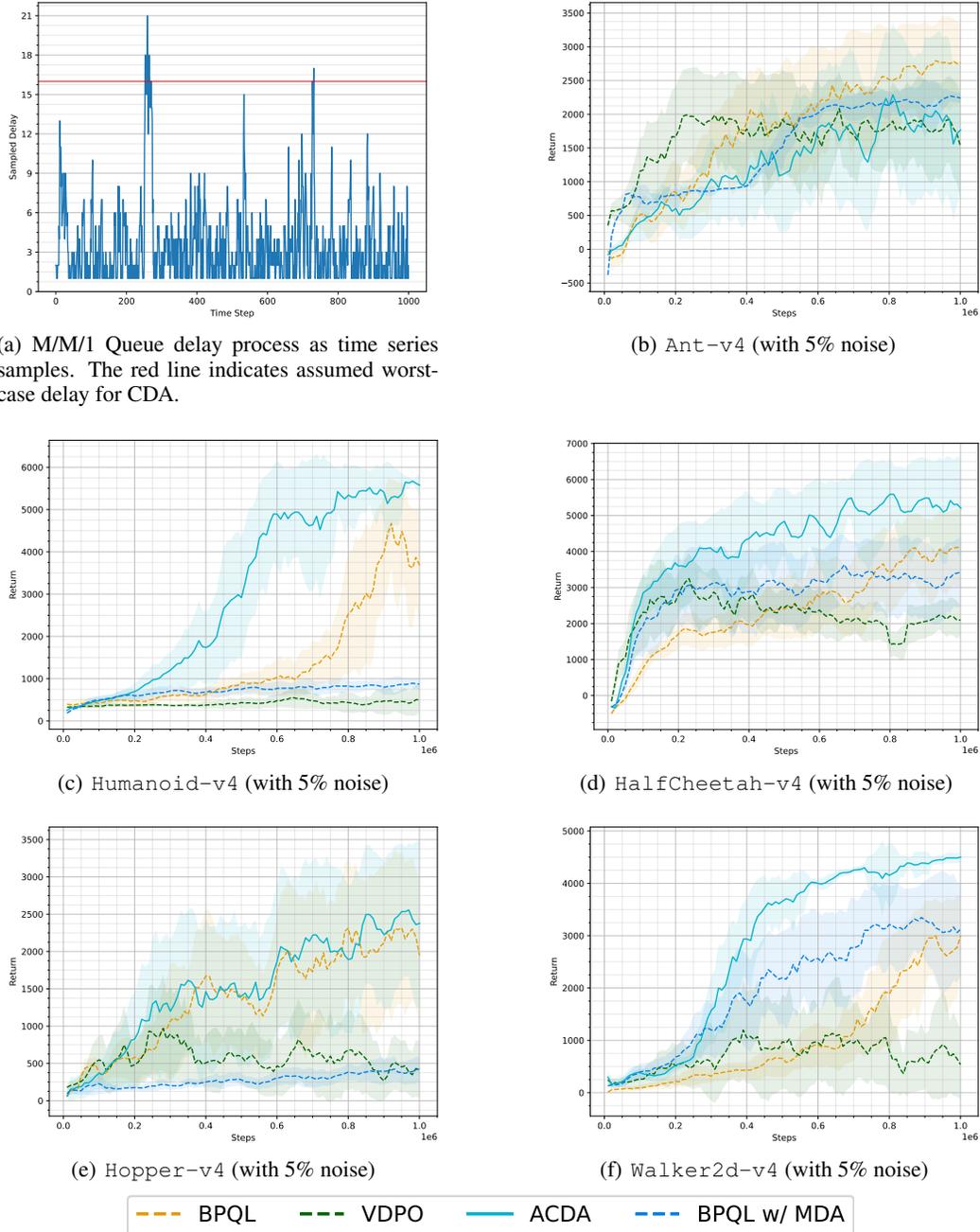


Figure 18: Time series evaluation during training on the MM1 delay process. All environments have added 5% noise to the actions.

Table 13: Best returns from the MM1 delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	<b>3074.17 ± 106.78</b>	5435.29 ± 68.34	4660.93 ± 448.10	3035.66 ± 103.80	3547.73 ± 133.51
VDPO	2528.67 ± 144.63	720.73 ± 634.35	3831.96 ± 960.07	1459.88 ± 933.11	2144.25 ± 1650.85
ACDA	2898.46 ± 838.07	<b>5805.60 ± 23.04</b>	<b>5898.36 ± 409.10</b>	<b>3122.53 ± 417.37</b>	<b>4562.33 ± 87.98</b>
BPQL w/ MDA	2308.18 ± 75.13	944.56 ± 92.79	3953.69 ± 351.34	512.43 ± 346.18	3667.61 ± 142.48

1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

---

### E.3 RESULTS WHEN VIOLATING THE UPPER BOUND ASSUMPTIONS

The  $GE_{1,23}$  and  $GE_{4,32}$  delay processes occasionally have a very high delay, but most of the time, the delays of these are very low. To convert these to a constant-delay MDP, we need to assume the worst-case possible delay of the process. We do this using the CDA method described in Appendix A. CDA can be applied regardless of whether the constant  $h$  that we wish to act under is a worst-case delay or not. Though CDA only guarantees the MDP property if  $h$  is a true upper bound of the delay process.

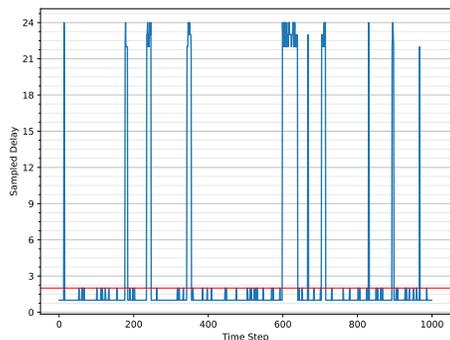
A natural question is how state-of-the-art approaches perform if we apply CDA to a more favorable constant  $h$ , which holds most of the time and is much lower than the worst-case delay. We answer this by evaluating BPQL and VDPO under more opportunistic constants  $h$ . These are compared against the performance of ACDA, which can still adapt to much larger delays. We also include an evaluation of BPQL with the MDA policy, as in Appendix E.2, but now when that acts under the opportunistic constant  $h$  instead. We present the results of this evaluation in Appendices E.3.1, E.3.2, and E.3.3, which are split based on the delay process used. The opportunistic constant  $h$  used is highlighted as a red line in a time series samples plot for each delay process.

The results show that ACDA still outperforms state-of-the-art in most benchmarks. While BPQL and VDPO can achieve better performance in some cases, ACDA is still performing close to the best algorithm. Also, operating under these opportunistic constants can have significant negative consequences. This is best highlighted by the results in Figure 19(d), where VDPO experiences a collapse in performance under the `HalfCheetah-v4` environment using the  $GE_{1,23}$  delay process.

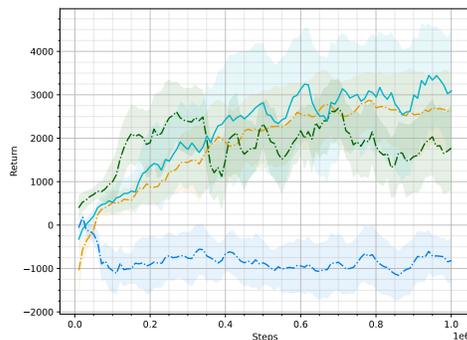
Based on these results, we conclude that the adaptiveness provided by the interaction layer is a necessity to be able to achieve high performance under random unobservable delays. While it is possible to sacrifice the MDP property to gain performance in the constant-delay setting, state-of-the-art still does not outperform the adaptive ACDA algorithm.

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

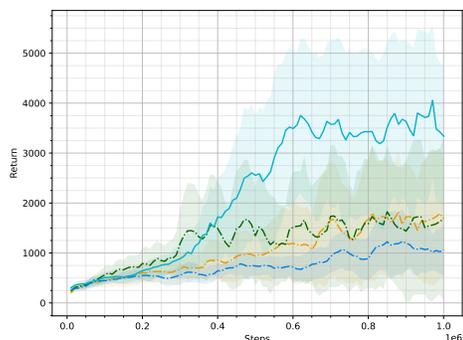
### E.3.1 $GE_{1,23}$ DELAY PROCESS WITH LOW CDA



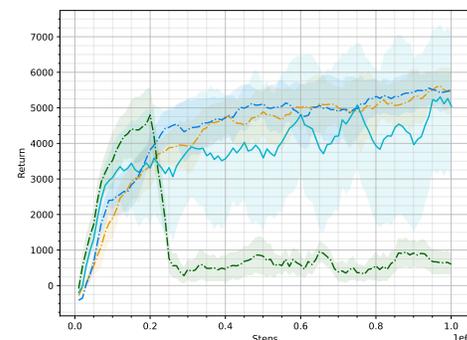
(a)  $GE_{1,23}$  delay process as time series samples. The red line indicates assumed upper bound for CDA.



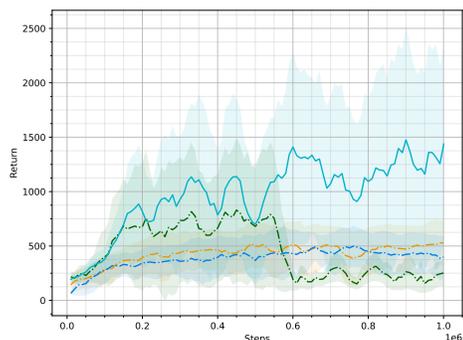
(b) Ant-v4 (with 5% noise)



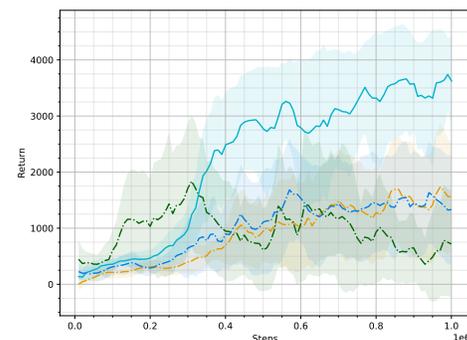
(c) Humanoid-v4 (with 5% noise)



(d) HalfCheetah-v4 (with 5% noise)



(e) Hopper-v4 (with 5% noise)



(f) Walker2d-v4 (with 5% noise)

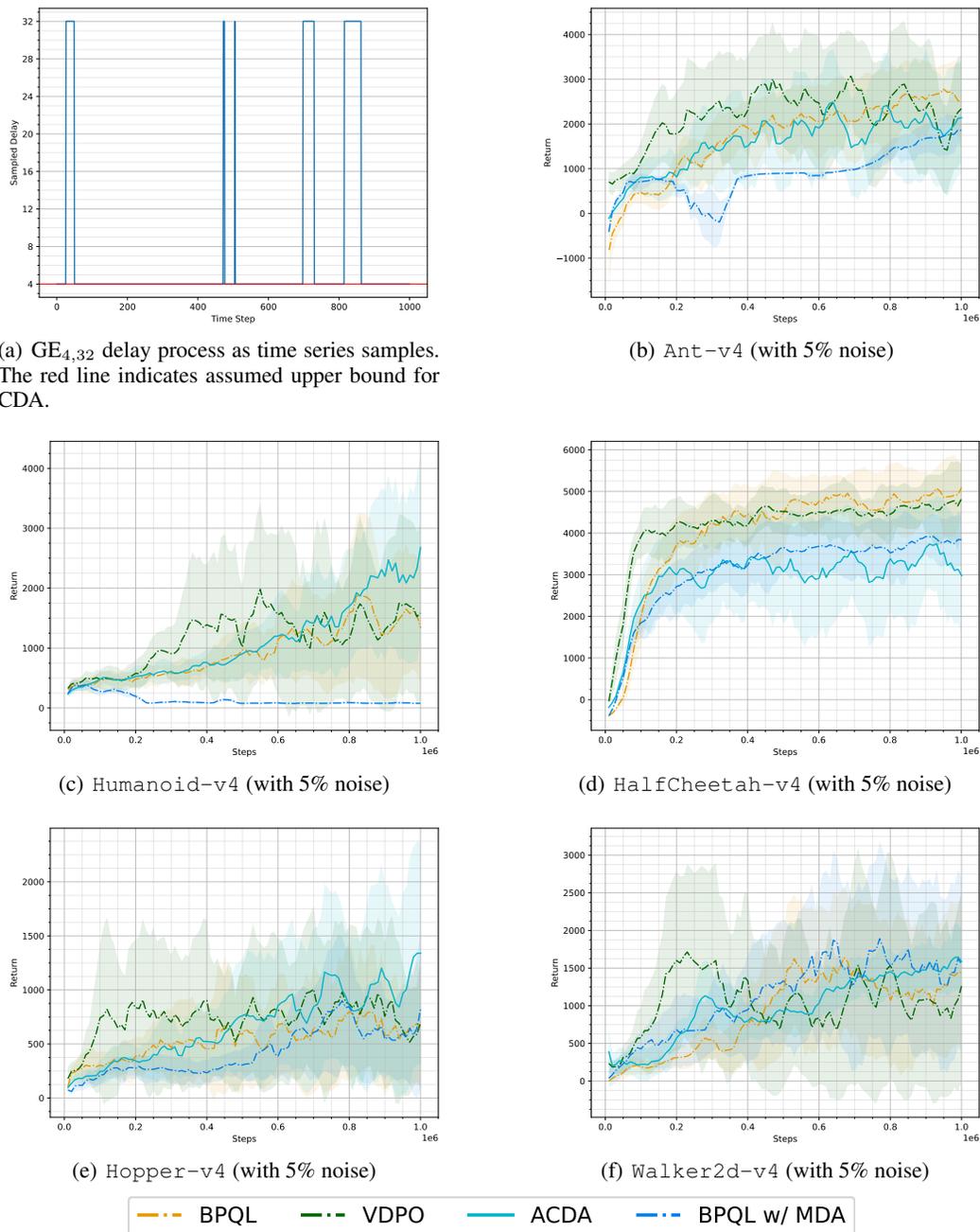


Figure 19: Time series evaluation during training on the  $GE_{1,23}$  delay process for opportunistic delay assumptions of  $h = 2$  (except for ACDA). All environments have added 5% noise to the actions.

Table 14: Best returns from the  $GE_{1,23}$  delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	3359.65 ± 288.24	2469.96 ± 1375.23	5944.67 ± 395.56	642.54 ± 420.04	2043.32 ± 923.00
VDPO	3103.10 ± 252.19	2658.48 ± 2044.40	5625.06 ± 524.23	1113.51 ± 836.72	2756.73 ± 1693.64
ACDA	<b>4112.78 ± 818.44</b>	<b>4608.76 ± 1084.52</b>	<b>5984.25 ± 1885.78</b>	<b>2094.65 ± 944.20</b>	<b>3863.59 ± 232.52</b>
BPQL w/ MDA	423.46 ± 73.63	1543.10 ± 536.14	5710.20 ± 566.48	545.01 ± 116.07	1923.28 ± 838.62

2052 E.3.2  $GE_{4,32}$  DELAY PROCESS WITH LOW CDA



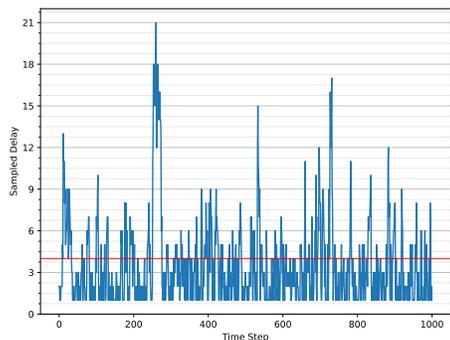
2096 Figure 20: Time series evaluation during training on the  $GE_{4,32}$  delay process for opportunistic delay  
 2097 assumptions of  $h = 4$  (except for ACDA). All environments have added 5% noise to the actions.  
 2098  
 2099

2100 Table 15: Best returns from the  $GE_{4,32}$  delay process.

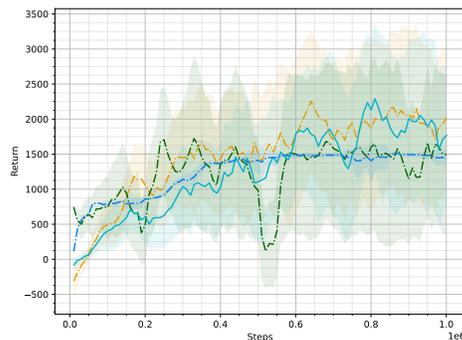
	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	3000.00 ± 754.09	2949.17 ± 1933.62	5315.27 ± 559.25	1243.58 ± 704.53	2107.39 ± 1219.55
VDPO	<b>3979.56 ± 331.76</b>	2956.52 ± 2322.74	<b>5424.96 ± 306.06</b>	1360.87 ± 627.11	2234.83 ± 1776.21
ACDA	2866.93 ± 1172.46	<b>3725.59 ± 1513.38</b>	4231.15 ± 333.69	<b>1727.79 ± 959.50</b>	1840.58 ± 386.78
BPQL w/ MDA	2155.47 ± 84.22	465.58 ± 82.59	4082.54 ± 411.47	1312.37 ± 868.12	<b>2355.23 ± 1145.73</b>

2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159

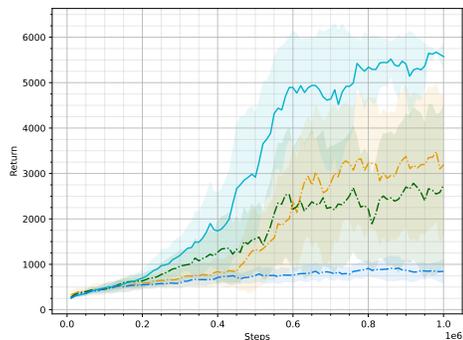
### E.3.3 M/M/1 QUEUE DELAY PROCESS WITH LOW CDA



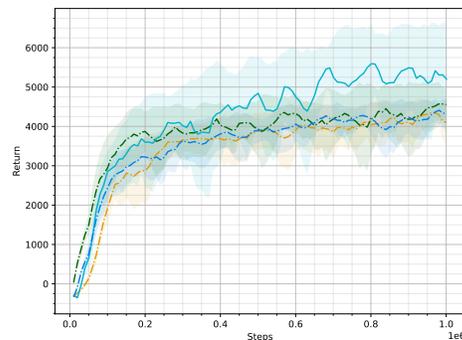
(a) M/M/1 Queue delay process as time series samples. The red line indicates assumed upper bound for CDA.



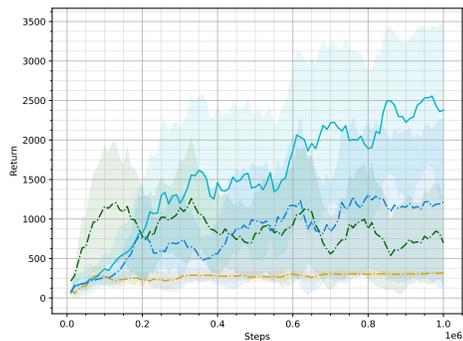
(b) Ant-v4 (with 5% noise)



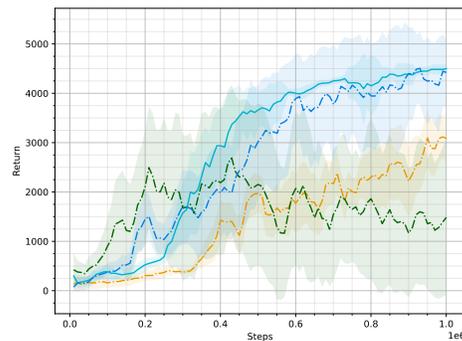
(c) Humanoid-v4 (with 5% noise)



(d) HalfCheetah-v4 (with 5% noise)



(e) Hopper-v4 (with 5% noise)



(f) Walker2d-v4 (with 5% noise)



Figure 21: Time series evaluation during training on the MM1 delay process for opportunistic delay assumptions of  $h = 4$  (except for ACDA). All environments have added 5% noise to the actions.

Table 16: Best returns from the MM1 delay process.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL	2577.34 ± 1217.11	4158.16 ± 981.22	4478.64 ± 193.82	407.91 ± 179.00	3475.80 ± 586.89
VDPO	2278.08 ± 726.85	3349.31 ± 1668.64	4857.65 ± 335.88	1628.83 ± 880.65	3554.77 ± 1278.22
ACDA	<b>2898.46 ± 838.07</b>	<b>5805.60 ± 23.04</b>	<b>5898.36 ± 409.10</b>	<b>3122.53 ± 417.37</b>	4562.33 ± 87.98
BPQL w/ MDA	1541.32 ± 16.71	1030.55 ± 185.86	4502.92 ± 214.25	1665.66 ± 1210.69	<b>4825.75 ± 126.28</b>

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

---

#### E.4 RESULTS FOR AVERAGE DELAYS

This sections performs an evaluation similar to that found in Appendix E.3, but using the average delay as the assumed delay. Although the results in Appendix E.3 are believed to be more favorable for the constant-delay processes, we also evaluate under the average delay to reduce bias in the results.

The average delay for each delay process are:

$$\mathbb{E}[\text{GE}_{1,23}] \approx 4.088 \tag{25}$$

$$\mathbb{E}[\text{GE}_{4,32}] \approx 7.177 \tag{26}$$

$$\mathbb{E}[\text{MM1}] \approx 2.916 \quad (\text{Monte-Carlo estimated from } 10^7 \text{ samples}) \tag{27}$$

As the average delays are not whole numbers, we evaluate both when the average delays are rounded up and when they are rounded down. Although the average delay for MM1 is slightly different to the theoretical average for M/M/1 queues ( $\frac{1}{0.75-0.33} \approx 2.381$ ) due to discretization, this does not affect the rounding.

We present the average CDA results for delay processes  $\text{GE}_{1,23}$ ,  $\text{GE}_{4,32}$ , and MM1 in Appendices E.4.1, E.4.2, and E.4.3 respectively. From these results, we see that ACDA is best in 13 out of 15 benchmarks, compared to 11 out of 15 benchmarks for low CDA (presented in Appendix E.3). Additionally, as shown later in Appendix E.5, the average CDA never yields the overall best-performing baseline for any benchmark.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

### E.4.1 $GE_{1,23}$ DELAY PROCESS WITH AVERAGE CDA

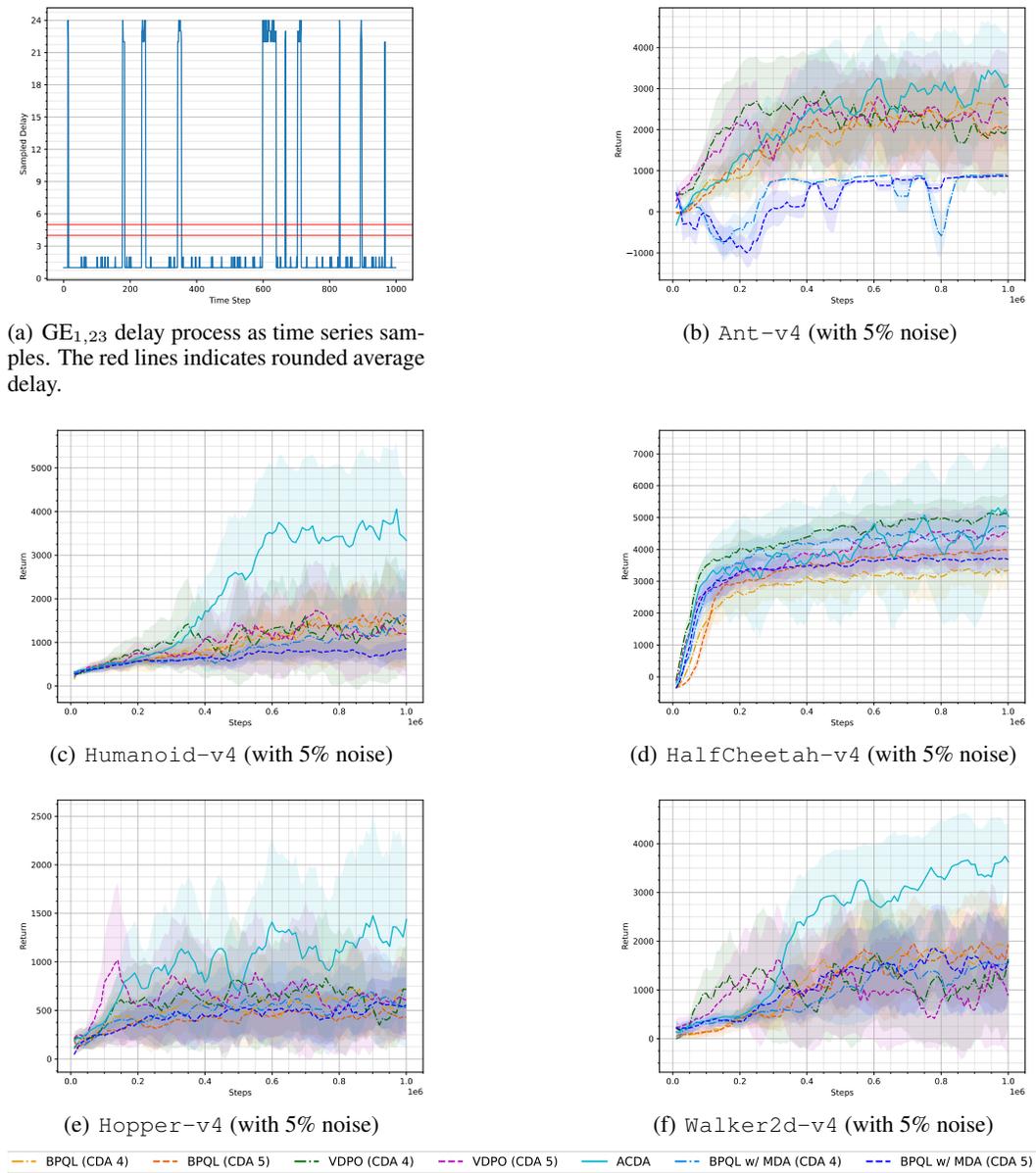


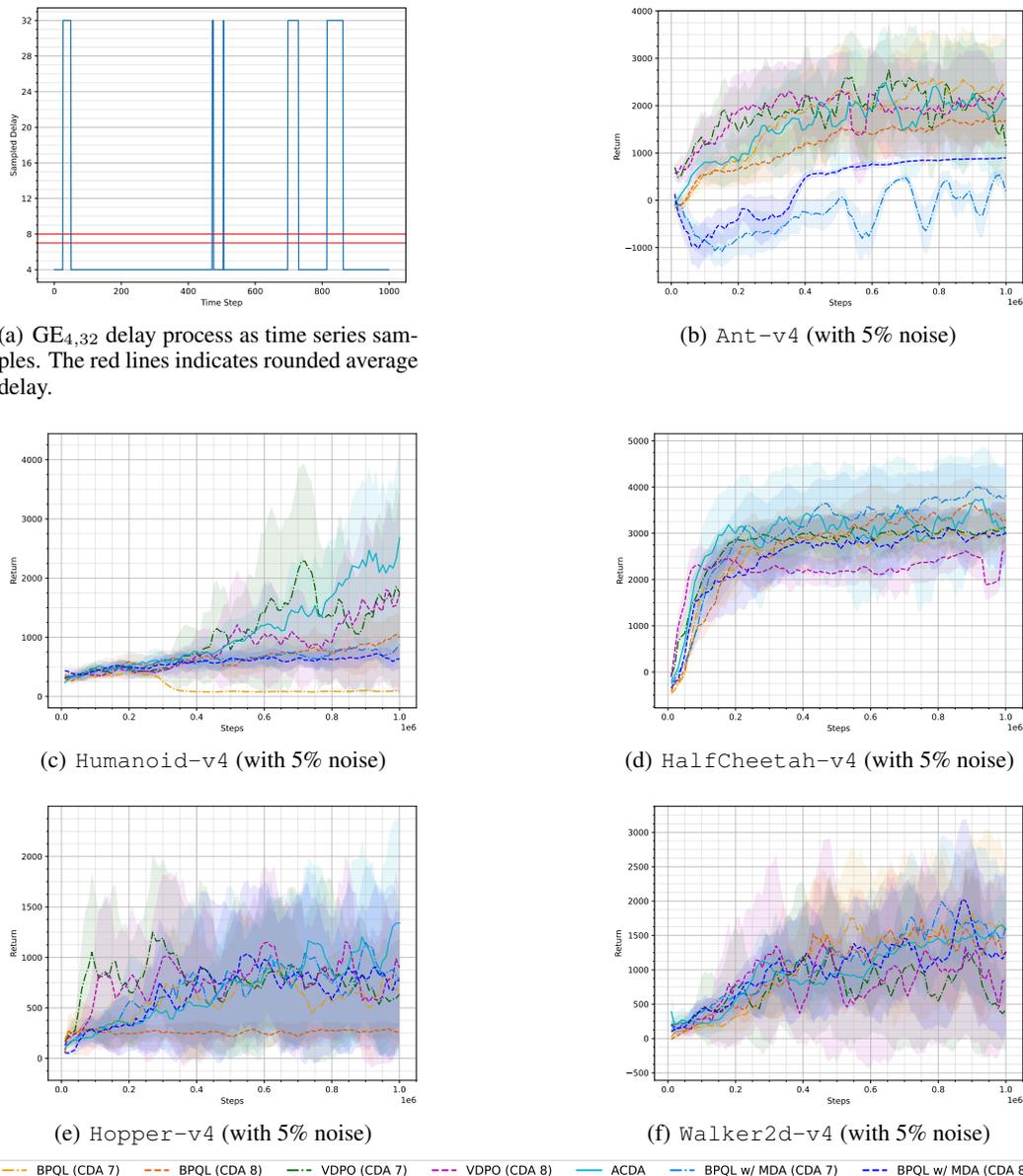
Figure 22: Time series evaluation during training on the  $GE_{1,23}$  delay process for average delay assumptions of  $h = 4$  and  $h = 5$ . All environments have added 5% noise to the actions.

Table 17: Best returns from the  $GE_{1,23}$  delay process. Average delay results below the solid line.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL (CDA 4)	3076.02 ± 443.41	2287.35 ± 1047.81	3740.65 ± 340.21	1097.92 ± 610.78	2317.55 ± 499.71
BPQL (CDA 5)	3127.63 ± 440.96	2204.08 ± 1004.49	4102.93 ± 555.25	603.92 ± 114.92	2704.21 ± 746.97
VDPO (CDA 4)	3240.02 ± 1026.94	2725.54 ± 1404.37	5327.23 ± 317.71	1254.58 ± 513.41	2608.30 ± 2130.91
VDPO (CDA 5)	3439.84 ± 1045.19	2072.25 ± 1421.87	4950.22 ± 602.94	1507.23 ± 955.98	2093.83 ± 1835.53
<b>ACDA</b>	<b>4112.78 ± 818.44</b>	<b>4608.76 ± 1084.52</b>	<b>5984.25 ± 1885.78</b>	<b>2094.65 ± 944.20</b>	<b>3863.59 ± 232.52</b>
BPQL w/ MDA (CDA 4)	918.30 ± 5.80	1971.09 ± 1698.64	4826.55 ± 339.02	951.58 ± 241.48	2211.10 ± 1042.52
BPQL w/ MDA (CDA 5)	881.15 ± 7.17	949.03 ± 282.09	3886.75 ± 289.39	726.35 ± 189.88	2219.83 ± 967.03

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

### E.4.2 $GE_{4,32}$ DELAY PROCESS WITH AVERAGE CDA



(a)  $GE_{4,32}$  delay process as time series samples. The red lines indicates rounded average delay.

(b) Ant-v4 (with 5% noise)

(c) Humanoid-v4 (with 5% noise)

(d) HalfCheetah-v4 (with 5% noise)

(e) Hopper-v4 (with 5% noise)

(f) Walker2d-v4 (with 5% noise)

Figure 23: Time series evaluation during training on the  $GE_{4,32}$  delay process for average delay assumptions of  $h = 7$  and  $h = 8$ . All environments have added 5% noise to the actions.

Table 18: Best returns from the  $GE_{4,32}$  delay process. Average delay results below the solid line.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL (CDA 7)	3295.31 ± 500.78	466.90 ± 92.62	3363.96 ± 440.71	1279.14 ± 751.20	<b>2573.74 ± 1335.71</b>
BPQL (CDA 8)	1936.41 ± 141.68	1168.42 ± 247.85	3815.76 ± 390.89	477.15 ± 319.04	2112.47 ± 938.36
VDPO (CDA 7)	<b>3449.78 ± 216.11</b>	3343.71 ± 1853.87	3465.86 ± 255.54	1577.83 ± 561.66	1845.23 ± 1890.44
VDPO (CDA 8)	2836.59 ± 391.05	2630.14 ± 1755.00	2743.14 ± 304.28	1673.32 ± 765.39	2443.61 ± 1789.56
ACDA	2866.93 ± 1172.46	<b>3725.59 ± 1513.38</b>	<b>4231.15 ± 333.69</b>	<b>1727.79 ± 959.50</b>	1840.58 ± 386.78
BPQL w/ MDA (CDA 7)	639.19 ± 50.33	972.62 ± 198.05	4149.78 ± 286.66	1430.83 ± 802.87	2279.09 ± 987.54
BPQL w/ MDA (CDA 8)	917.01 ± 30.49	841.34 ± 158.81	3324.03 ± 363.93	1664.18 ± 732.19	2456.04 ± 1239.15

2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375

### E.4.3 M/M/1 QUEUE DELAY PROCESS WITH AVERAGE CDA

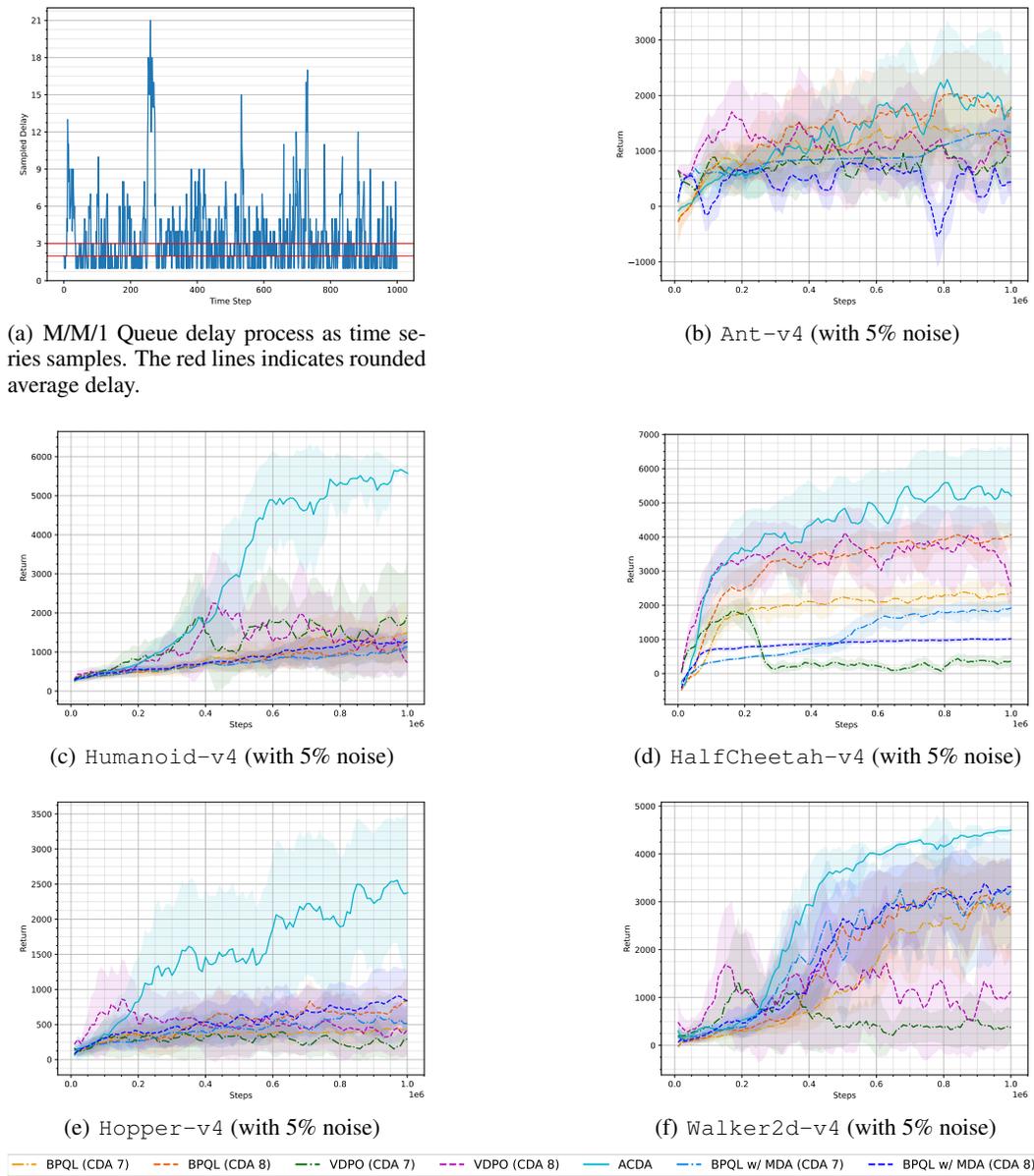


Figure 24: Time series evaluation during training on the MM1 delay process for average delay assumptions of  $h = 2$  and  $h = 3$ . All environments have added 5% noise to the actions.

Table 19: Best returns from the MM1 delay process. Average delay results below the solid line.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
BPQL (CDA 2)	1675.16 ± 613.61	1728.83 ± 946.59	2509.78 ± 159.44	518.81 ± 131.09	3377.82 ± 174.89
BPQL (CDA 3)	2504.20 ± 399.72	1507.54 ± 653.85	4204.65 ± 295.58	980.68 ± 331.18	3639.76 ± 106.58
VDPO (CDA 2)	1649.20 ± 522.32	2856.47 ± 1621.48	1937.06 ± 351.08	703.85 ± 558.47	1848.79 ± 1572.29
VDPO (CDA 3)	2187.20 ± 386.29	3217.84 ± 1761.53	4469.93 ± 489.60	1249.77 ± 716.85	2778.95 ± 1371.81
ACDA	<b>2898.46 ± 838.07</b>	<b>5805.60 ± 23.04</b>	<b>5898.36 ± 409.10</b>	<b>3122.53 ± 417.37</b>	<b>4562.33 ± 87.98</b>
BPQL w/ MDA (CDA 2)	1475.42 ± 31.45	1370.12 ± 474.63	2058.15 ± 133.25	864.57 ± 569.87	3892.97 ± 84.39
BPQL w/ MDA (CDA 3)	815.85 ± 5.94	1443.72 ± 299.19	1040.09 ± 46.75	994.34 ± 348.21	3541.48 ± 44.42

2376  
2377  
2378  
2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429

## E.5 RESULTS FROM ALL EVALUATED BASELINES

This section presents all results for every evaluated baseline, allowing for a wholistic comparison of performance. These results are previously presented in Appendices E.1, E.2, E.3, and E.4.

### E.5.1 $GE_{1,23}$ DELAY PROCESS FOR ALL BASELINES

Table 20: Best returns under the  $GE_{1,23}$  delay process for all baselines.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	14.22 ± 14.89	862.18 ± 266.21	2064.18 ± 223.48	306.91 ± 51.26	708.33 ± 221.53
SAC w/ CDA	69.28 ± 114.47	414.05 ± 204.20	128.47 ± 9.55	426.92 ± 27.93	428.44 ± 509.44
Dreamer	1111.73 ± 412.35	1463.07 ± 649.56	1796.07 ± 381.47	334.30 ± 245.42	1081.12 ± 905.94
BPQL (Low CDA)	3359.65 ± 288.24	2469.96 ± 1375.23	5944.67 ± 395.56	642.54 ± 420.04	2043.32 ± 923.00
BPQL (High CDA)	2691.88 ± 129.84	585.19 ± 163.49	4320.20 ± 1028.52	1328.71 ± 937.67	1215.91 ± 776.93
BPQL (CDA 4)	3076.02 ± 443.41	2287.35 ± 1047.81	3740.65 ± 340.21	1097.92 ± 610.78	2317.55 ± 499.71
BPQL (CDA 5)	3127.63 ± 440.96	2204.08 ± 1004.49	4102.93 ± 555.25	603.92 ± 114.92	2704.21 ± 746.97
VDPO (Low CDA)	3103.10 ± 252.19	2658.48 ± 2044.40	5625.06 ± 524.23	1113.51 ± 836.72	2756.73 ± 1693.64
VDPO (High CDA)	2163.00 ± 53.04	417.25 ± 210.09	3144.23 ± 1156.52	709.20 ± 522.01	846.88 ± 808.67
VDPO (CDA 4)	3240.02 ± 1026.94	2725.54 ± 1404.37	5327.23 ± 317.71	1254.58 ± 513.41	2608.30 ± 2130.91
VDPO (CDA 5)	3439.84 ± 1045.19	2072.25 ± 1421.87	4950.22 ± 602.94	1507.23 ± 955.98	2093.83 ± 1835.53
BPQL w/ MDA (Low CDA)	423.46 ± 73.63	1543.10 ± 536.14	5710.20 ± 566.48	545.01 ± 116.07	1923.28 ± 838.62
BPQL w/ MDA (High CDA)	1795.29 ± 23.78	563.36 ± 96.11	4926.36 ± 60.08	465.14 ± 138.86	3681.39 ± 126.41
BPQL w/ MDA (CDA 4)	918.30 ± 5.80	1971.09 ± 1698.64	4826.55 ± 339.02	951.58 ± 241.48	2211.10 ± 1042.52
BPQL w/ MDA (CDA 5)	881.15 ± 7.17	949.03 ± 282.09	3886.75 ± 289.39	726.35 ± 189.88	2219.83 ± 967.03
DCAC	949.97 ± 11.87	128.47 ± 36.09	920.09 ± 33.05	16.99 ± 15.94	106.70 ± 53.84
ACDA	<b>4112.78 ± 818.44</b>	<b>4608.76 ± 1084.52</b>	<b>5984.25 ± 1885.78</b>	<b>2094.65 ± 944.20</b>	<b>3863.59 ± 232.52</b>

2430  
2431  
2432  
2433  
2434  
2435  
2436  
2437  
2438  
2439  
2440  
2441  
2442  
2443  
2444  
2445  
2446  
2447  
2448  
2449  
2450  
2451  
2452  
2453  
2454  
2455  
2456  
2457  
2458  
2459  
2460  
2461  
2462  
2463  
2464  
2465  
2466  
2467  
2468  
2469  
2470  
2471  
2472  
2473  
2474  
2475  
2476  
2477  
2478  
2479  
2480  
2481  
2482  
2483

### E.5.2 GE<sub>4,32</sub> DELAY PROCESS FOR ALL BASELINES

Table 21: Best returns under the GE<sub>4,32</sub> delay process for all baselines.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	-5.72 ± 19.62	494.43 ± 156.01	-158.78 ± 65.86	279.74 ± 109.83	60.86 ± 75.59
SAC w/ CDA	18.93 ± 23.64	230.45 ± 99.22	591.32 ± 36.39	315.47 ± 51.49	257.18 ± 73.42
Dreamer	1147.56 ± 371.15	1091.48 ± 577.52	2493.19 ± 231.22	515.36 ± 430.72	1233.79 ± 802.47
BPQL (Low CDA)	3000.00 ± 754.09	2949.17 ± 1933.62	5315.27 ± 559.25	1243.58 ± 704.53	2107.39 ± 1219.55
BPQL (High CDA)	2509.52 ± 117.37	276.63 ± 131.70	2136.36 ± 547.04	433.29 ± 381.79	875.09 ± 747.72
BPQL (CDA 7)	3295.31 ± 500.78	466.90 ± 92.62	3363.96 ± 440.71	1279.14 ± 751.20	2573.74 ± 1335.71
BPQL (CDA 8)	1936.41 ± 141.68	1168.42 ± 247.85	3815.76 ± 390.89	477.15 ± 319.04	2112.47 ± 938.36
VDPO (Low CDA)	<b>3979.56 ± 331.76</b>	2956.52 ± 2322.74	<b>5424.96 ± 306.06</b>	1360.87 ± 627.11	2234.83 ± 1776.21
VDPO (High CDA)	2266.99 ± 90.89	280.72 ± 169.85	3664.30 ± 929.25	330.44 ± 263.74	344.73 ± 316.82
VDPO (CDA 7)	3449.78 ± 216.11	3343.71 ± 1853.87	3465.86 ± 255.54	1577.83 ± 561.66	1845.23 ± 1890.44
VDPO (CDA 8)	2836.59 ± 391.05	2630.14 ± 1755.00	2743.14 ± 304.28	1673.32 ± 765.39	2443.61 ± 1789.56
BPQL w/ MDA (Low CDA)	2155.47 ± 84.22	465.58 ± 82.59	4082.54 ± 411.47	1312.37 ± 868.12	2355.23 ± 1145.73
BPQL w/ MDA (High CDA)	1661.41 ± 43.42	359.46 ± 156.86	3609.45 ± 328.37	224.34 ± 90.12	<b>3015.45 ± 1428.05</b>
BPQL w/ MDA (CDA 7)	639.19 ± 50.33	972.62 ± 198.05	4149.78 ± 286.66	1430.83 ± 802.87	2279.09 ± 987.54
BPQL w/ MDA (CDA 8)	917.01 ± 30.49	841.34 ± 158.81	3324.03 ± 363.93	1664.18 ± 732.19	2456.04 ± 1239.15
DCAC	953.14 ± 12.06	167.97 ± 82.14	1123.47 ± 100.34	57.98 ± 28.04	9.23 ± 20.35
ACDA	2866.93 ± 1172.46	<b>3725.59 ± 1513.38</b>	4231.15 ± 333.69	<b>1727.79 ± 959.50</b>	1840.58 ± 386.78

### E.5.3 M/M/1 QUEUE DELAY PROCESS FOR ALL BASELINES

Table 22: Best returns under the MM1 delay process for all baselines.

	Ant-v4	Humanoid-v4	HalfCheetah-v4	Hopper-v4	Walker2d-v4
SAC	-0.58 ± 8.66	921.04 ± 299.47	20.69 ± 94.91	333.06 ± 96.04	604.80 ± 212.37
SAC w/ CDA	102.00 ± 33.77	613.03 ± 157.68	550.84 ± 16.28	627.59 ± 24.62	2005.76 ± 341.30
Dreamer	1121.11 ± 58.69	981.38 ± 597.01	584.40 ± 72.26	975.72 ± 650.05	1801.81 ± 1158.73
BPQL (Low CDA)	2577.34 ± 1217.11	4158.16 ± 981.22	4478.64 ± 193.82	407.91 ± 179.00	3475.80 ± 586.89
BPQL (High CDA)	<b>3074.17 ± 106.78</b>	5435.29 ± 68.34	4660.93 ± 448.10	3035.66 ± 103.80	3547.73 ± 133.51
BPQL (CDA 2)	1675.16 ± 613.61	1728.83 ± 946.59	2509.78 ± 159.44	518.81 ± 131.09	3377.82 ± 174.89
BPQL (CDA 3)	2504.20 ± 399.72	1507.54 ± 653.85	4204.65 ± 295.58	980.68 ± 331.18	3639.76 ± 106.58
VDPO (Low CDA)	2278.08 ± 726.85	3349.31 ± 1668.64	4857.65 ± 335.88	1628.83 ± 880.65	3554.77 ± 1278.22
VDPO (High CDA)	2528.67 ± 144.63	720.73 ± 634.35	3831.96 ± 960.07	1459.88 ± 933.11	2144.25 ± 1650.85
VDPO (CDA 2)	1649.20 ± 522.32	2856.47 ± 1621.48	1937.06 ± 351.08	703.85 ± 558.47	1848.79 ± 1572.29
VDPO (CDA 3)	2187.20 ± 386.29	3217.84 ± 1761.53	4469.93 ± 489.60	1249.77 ± 716.85	2778.95 ± 1371.81
BPQL w/ MDA (Low CDA)	1541.32 ± 16.71	1030.55 ± 185.86	4502.92 ± 214.25	1665.66 ± 1210.69	<b>4825.75 ± 126.28</b>
BPQL w/ MDA (High CDA)	2308.18 ± 75.13	944.56 ± 92.79	3953.69 ± 351.34	512.43 ± 346.18	3667.61 ± 142.48
BPQL w/ MDA (CDA 2)	1475.42 ± 31.45	1370.12 ± 474.63	2058.15 ± 133.25	864.57 ± 569.87	3892.97 ± 84.39
BPQL w/ MDA (CDA 3)	815.85 ± 5.94	1443.72 ± 299.19	1040.09 ± 46.75	994.34 ± 348.21	3541.48 ± 44.42
DCAC	959.23 ± 13.54	525.85 ± 135.36	35.60 ± 21.42	1026.45 ± 2.96	24.48 ± 45.46
ACDA	2898.46 ± 838.07	<b>5805.60 ± 23.04</b>	<b>5898.36 ± 409.10</b>	<b>3122.53 ± 417.37</b>	4562.33 ± 87.98

---

## 2484 F PRACTICAL CONSIDERATIONS OF THE INTERACTION LAYER

2485

2486 This section discusses practical considerations when deploying ACDA and the interaction layer to  
2487 real-world environments. Specifically, we discuss the delays not handled by the interaction layer in  
2488 Appendix F.1, the effect that ACDA has on computational delays in Appendix F.2, how to mitigate  
2489 computational delays in Appendix F.3, and the effect that ACDA has on transmission bandwidth in  
2490 Appendix F.4.

2491

### 2492 F.1 CONSIDERATIONS FOR NON-INTERACTION DELAYS

2493

2494 As illustrated in Figure 1 in the introduction, there are additional delays not handled by the interac-  
2495 tion layer. Namely, the delays between the interaction layer and the system itself. In this paper, we  
2496 presume that these delays are negligible or otherwise accounted for. If these delays are not negli-  
2497 gible and must be accounted for, then this can most likely be handled by prematurely sensing and  
2498 actuating the system. As the interaction layer by design is located close to the excited system, any  
2499 delays between them will likely take place over controlled channels (such as USB or SPI), meaning  
2500 that any delay over these channels is stable.

2501

### 2502 F.2 EFFECT OF ACTION PACKET ON COMPUTATIONAL DELAY

2503

2504 Although ACDA handles interaction delays, the computation of the action packet matrix itself does  
2505 add computational delay, compared to constant-delay approaches that only compute a single action.  
2506 This could cause concern for real-world applications if this additional delay is too significant. If  
2507 the computational delay is longer than the excitation period of the environment, the policy cannot  
2508 generate actions fast enough, and the interaction will stall. In this section, we discuss the effect that  
2509 the computation of the action packet can have on the delay, present measurements of computational  
2510 delay, and put that into the context of the evaluated environments.

2511

2512 There are a few remarks about the action packet itself:

2513

- 2514 • In ACDA, the computation of the rows is done in parallel. The effective computational time  
2515 is linear instead of quadratic. More specifically, the number of sequential computations is  
2516 proportional to the sum of the horizon and the prediction length,  $h + L$ .
- 2517 • The action packet contains horizons of actions, allowing for gaps in the interaction. For  
2518 a practical scenario, in the event of not having enough time to compute subsequent action  
2519 packets, it is possible to only generate action packets based on the latest observation packet  
2520 that has reached the agent.

2521

2522 With this in mind, we measure the average computational time of action packets for the network  
2523 structure used when evaluating the `Ant-v4` environment under the MM1 delay process. We mea-  
2524 sure the time taken in the training loop to generate a single action packet, as well as the execution  
2525 time of the individual network components. All measurements are done in our framework imple-  
2526 mented in PyTorch, collected on the same system used to run the benchmarks.

2527

2528

2529

2530

2531

2532

2533

2534

2535

2536

2537

Table 23: Execution time measurements.

Aspect	Measured Time
Generating a random action packet	24 ms
Generating an action packet using the MDA	68 ms
Randomly filling an action packet matrix	39 $\mu$ s
Single GRU forward pass	164 $\mu$ s
Single policy forward pass	61 $\mu$ s

2538

2539

2540

2541

2542

2543

2544

2545

2546

2547

A notable aspect of the measurements in Table 23 is the difference in generating a random action packet in the training loop, compared to directly filling an action packet matrix in PyTorch (24 ms vs 61  $\mu$ s). This hints at that our current implementation is not optimized, that there are further gains to

2538 be made, and that 68 ms is not a representative time for generating an action packet in an optimized  
 2539 implementation.

2540 The worst-case computation for a row in the action packet under the MM1 delay process is for  
 2541 the 16th row (delay 16). This consists of first embedding the observed state, then embedding the  
 2542 16 guessed preceding actions into a distribution, and then sequentially generating 16 actions and  
 2543 embedding the next distribution, excluding the distribution after the final action, as that is not needed.  
 2544 The effective computation time for the action packet is then as shown in Equation 28:

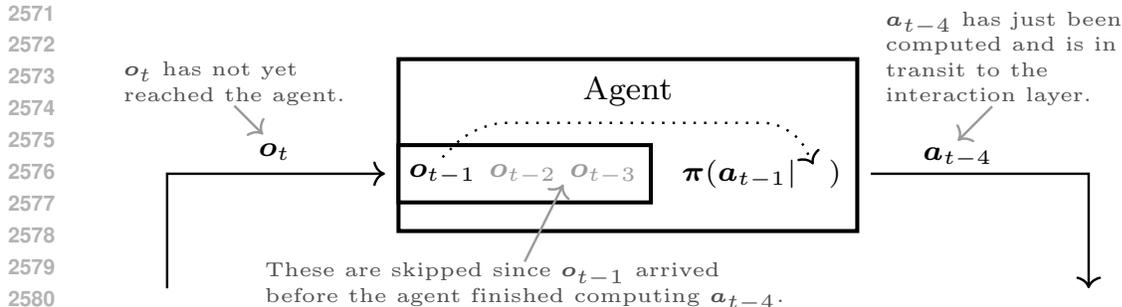
$$2546 t_{\text{action packet}} = t_{\text{embed}} + 16 \cdot t_{\text{GRU}} + 15 \cdot (t_{\text{policy}} + t_{\text{GRU}}) + t_{\text{policy}} \quad (28)$$

2548 As the embedding network is comparable in size to the policy network, we use the time for the policy  
 2549 as a proxy for the state embedding. Inserting the measurements from Table 23 into the formula from  
 2550 Equation 28, we get an estimated average execution time of 6.1 ms for an action packet. In the  
 2551 context of the Ant-v4 environment which uses a 50 ms actuation period, the computational delay  
 2552 is less than the time for a single step in the environment. The computational delay is also lower than  
 2553 the smallest actuation period for any environment used in the evaluation, the smallest period being  
 2554 8 ms for the Walker2d-v4 environment.

### 2556 F.3 COMPENSATING FOR EXCESSIVE COMPUTATIONAL DELAYS

2557 Even though the expected computational delay estimated in Appendix F.2 is smaller than the actua-  
 2558 tion periods of the environments evaluated in this paper, other environments may require an even  
 2559 shorter actuation period that is smaller than the computational delay. This shorter period will cause  
 2560 the ACDA agent to stall, as it cannot produce action packets faster than the rate at which observa-  
 2561 tion packets arrive.  
 2562

2563 A possible solution to this problem is for the ACDA agent to compute action packets based solely  
 2564 on the last received observation packet and to disregard all previous observation packets. These  
 2565 disregarded observation packets will not be used to compute any action packets, resulting in gaps in  
 2566 the interaction between the agent and the interaction layer. However, as each row in the matrix of an  
 2567 action packet describes a horizon of actions, ACDA already provides a solution for filling in these  
 2568 gaps in the interaction. The only aspect that has to change in the original algorithm is to make the  
 2569 memorized action prediction aware that a packet has been skipped. This process, which we refer to  
 2570 as *packet skipping*, is illustrated in Figure 25.



2582 Figure 25: Illustration of packet skipping when observations are reaching the agent faster than it can  
 2583 compute action packets. In the illustrated example, while  $a_{t-4}$  was being computed,  $o_{t-3}$ ,  $o_{t-2}$ ,  
 2584 and  $o_{t-1}$  had arrived at the agent. The action packets  $a_{t-3}$  and  $a_{t-2}$  are never computed, and the  
 2585 agent directly starts to compute  $a_{t-1}$  instead.

2586 To evaluate the plausibility of packet skipping, we train and evaluate ACDA when an action packet is  
 2587 computed based on every 2nd, 4th, and 8th incoming observation packet. These results are presented  
 2588 in Table 24, where they are put into the context of the results from the main body of the paper.

2590 Although there is a noticeable drop-off in performance for some benchmarks, such as  
 2591 Humanoid-v4 and Hopper-v4 under MM1 delays, ACDA still manages to achieve good performance  
 in most benchmarks, even when skipping the computation of action packets. It should also be

2592  
2593  
2594  
2595  
2596  
2597  
2598  
2599  
2600  
2601  
2602  
2603  
2604  
2605  
2606  
2607  
2608  
2609  
2610  
2611  
2612  
2613  
2614  
2615  
2616  
2617  
2618  
2619  
2620  
2621  
2622  
2623  
2624  
2625  
2626  
2627  
2628  
2629  
2630  
2631  
2632  
2633  
2634  
2635  
2636  
2637  
2638  
2639  
2640  
2641  
2642  
2643  
2644  
2645

Table 24: Best average return when only every 2nd, 4th, and 8th observation packet is used by ACDA to compute an action packet.

<i>Gymnasium env.</i>	Ant-v4			Humanoid-v4			HalfCheetah-v4			Hopper-v4			Walker2d-v4		
<i>Delay process</i>	GE <sub>1,23</sub>	GE <sub>4,32</sub>	MM1												
SAC	14.22	-5.72	-0.58	862.18	494.43	921.04	2064.18	-158.78	20.69	306.91	279.74	333.06	708.33	60.86	604.80
SAC w/ CDA	69.28	18.93	102.00	414.05	230.45	613.03	128.47	591.32	550.84	426.92	315.47	627.59	428.44	257.18	2005.76
Dreamer	1111.73	1147.56	1121.11	1463.07	1091.48	981.38	1796.07	2493.19	584.40	334.30	515.36	975.72	1081.12	1233.79	1801.81
BPQL	2691.88	2509.52	<b>3074.17</b>	585.19	276.63	5435.29	4320.20	2136.36	4660.93	1328.71	433.29	3035.66	1215.91	875.09	3547.73
VDPO	2163.00	2266.99	2528.67	417.25	280.72	720.73	3144.23	3664.30	3831.96	709.20	330.44	1459.88	846.88	344.73	2144.25
DCAC	949.97	953.14	959.23	128.47	167.97	525.85	920.09	1123.47	35.60	16.99	57.98	1026.45	106.70	9.23	24.48
ACDA	<b>4112.78</b>	2866.93	2898.46	<b>4608.76</b>	3725.59	<b>5805.60</b>	5984.25	<b>4231.15</b>	<b>5898.36</b>	2094.65	<b>1727.79</b>	<b>3122.53</b>	3863.59	1840.58	<b>4562.33</b>
ACDA (every 2nd)	3760.33	<b>3215.65</b>	2932.27	4108.88	<b>3986.12</b>	3270.50	<b>6205.55</b>	4026.24	3773.21	1983.83	1533.11	2022.95	3697.84	<b>3790.17</b>	3818.12
ACDA (every 4th)	3341.16	3180.04	2448.11	4044.15	2342.79	3688.68	4949.92	3838.81	3653.31	1940.75	1710.44	1383.08	<b>4634.92</b>	1682.69	3944.06
ACDA (every 8th)	3263.72	2246.90	2588.54	4311.94	2672.57	1315.06	4311.85	3418.11	3334.57	<b>2132.22</b>	1493.36	829.54	3379.25	3552.03	4519.21

noted that we do not believe the computational delay to be a practical problem for these benchmarks, and that scenarios such as only using every 8th incoming observation packet are excessive; however, these results demonstrate the viability of the approach. It should also be noted that although some benchmarks, such as `Walker2d-v4` under `GE4,32`, exhibit high variance, ACDA still achieves a high average return across all benchmarks.

#### F.4 EFFECT OF ACTION PACKET ON COMMUNICATION

While the action packet enables adaptability to interaction delays, it also introduces bandwidth overhead due to sending a matrix of actions rather than a single action. The worst-case bandwidth overhead in the evaluated environments is in the `Humanoid-v4` under the `GE4,32` delay, where each action in the  $32 \times 32$  action packet matrix consists of 17 32-bit floating point numbers, which are sent to the interaction layer every 15 ms. This benchmark requires a bandwidth of approximately 4.5 MiB/s, as opposed to sending a single action, which only requires a bandwidth of 4.5 kiB/s. While many communication channels, such as WiFi, support bandwidths well above 4.5 MiB/s, this becomes problematic for low-bandwidth communication channels such as Bluetooth.

$$32 \cdot 32 \cdot 17 \cdot 4 \cdot \frac{1000}{15} \approx 4.43 \cdot 2^{20} \tag{29}$$

We do not attempt to optimize the bandwidth in this paper; instead, we assume it is sufficient for all our benchmarks. It is possible to implement methods on the current interaction layer framework that mitigate this issue, such as the action skipping method described in Appendix F.3. It is also possible to reduce the size of the action packet by reducing the action buffer horizon  $h$ , which will not impact performance unless the gaps in the interaction exceed  $h$ .

There is also room for alterations to the interaction layer setup to reduce the size of the action packets, while still providing the adaptability to random delays. In the case of `GE4,32`, only rows 4 and 32 in the action packet matrix will ever be selected. Therefore, we could optimize the interaction by only sending the rows for the delays we believe will actually occur, discarding or artificially delaying the packet if it did not contain a row for the actual delay.

## G INTERACTION-DELAYED REINFORCEMENT LEARNING WITH REAL-VALUED DELAY

In Section 3.1, we introduced the delayed MDP using discrete delays measured in steps within the MDP. For real systems described by an MDP, each step corresponds to some amount of real-valued time, possibly controlled by a clock. Any interaction delay with the real system will also correspond to some amount of real-valued time that does not necessarily align with the time taken for a step in the MDP. Therefore, it makes sense to talk about delay directly as real-valued time when considering interaction delays for systems in the real world.

In Appendix G.1, we describe the effect that delays have when they are described as real-valued delays. We describe in Appendix G.2 how to implement the interaction layer to handle these real-valued delays.

### G.1 ORIGIN AND EFFECT OF DELAY AS CONTINUOUS TIME

MDPs usually assign a time  $t$  to states, actions, and rewards. This time  $t \in \mathbb{N}$  is merely a discrete ordering of events. We model the origin of delays in the real world as elapsed wall clock time in the domain of  $\mathbb{R}^+$ . We use the following notation to distinguish between them:

$$\begin{aligned} t \in \mathbb{N} & \quad (\text{Order of events in MDP.}) & (30) \\ \tau \in \mathbb{R}^+ & \quad (\text{Wall clock time elapsed in the real world.}) & (31) \end{aligned}$$

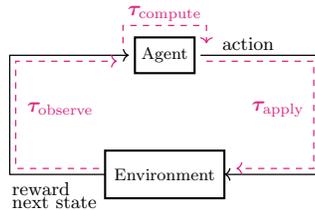
In the real world, there is an *interaction delay* in that it takes some time  $\tau_{\text{observe}} \in \mathbb{R}^+$  to observe a state, some time  $\tau_{\text{decide}} \in \mathbb{R}^+$  to generate the action, and some time  $\tau_{\text{apply}} \in \mathbb{R}^+$  to apply the action to the environment. In this time, the state  $s$  may evolve independently of an action being applied to the environment. Let this evolution process  $\Delta$  be defined as

$$\Delta : S \times \mathbb{R}^+ \times S \rightarrow \mathbb{R} \quad (32)$$

$$\text{such that } \tilde{s} \sim \Delta(\cdot | s, \tau) \quad (33)$$

$$\text{subject to } \forall s, \tilde{s}, \tau_1, \tau_2 : \Delta(\tilde{s} | \tilde{s}_i, \tau_2) \Big|_{\tilde{s}_i \sim \Delta(\cdot | s, \tau_1)} = \Delta(\tilde{s} | s, \tau_1 + \tau_2) \quad (34)$$

where  $s \in S$  is the state and  $\tau, \tau_1, \tau_2 \in \mathbb{R}^+$  are real wall-clock times in which the state has had time to evolve. The evolved state is unknown to the agent and is thus referred to as  $\tilde{s} \in S$ . The criterion in Equation 34 formally states that it should make no difference whether a state evolved for a single time period or if it is split into 2 time periods.



(a) Effects contributing to delay.

Standard (Assumed) RL Interaction:

$$\begin{aligned} a & \sim \pi(\cdot | s) \\ s' & \sim p(\cdot | s, a) \end{aligned}$$

With Interaction Delay:

$$\begin{aligned} a & \sim \pi(\cdot | s) \\ \tilde{s}_3 & \sim \Delta(\cdot | s, \tau_{\text{observe}} + \tau_{\text{decide}} + \tau_{\text{apply}}) \\ s' & \sim p(\cdot | \tilde{s}_3, a) \end{aligned}$$

(b) Violation of RL interaction assumption.

Figure 26: How delays violate the assumption used by state-of-the-art RL algorithms.

If the environment is sufficiently static, like a chess board, then this poses no issue because that  $\Delta(\cdot | s, \tau)$  will always evolve to the same state. If the environment is more dynamic, such as balancing

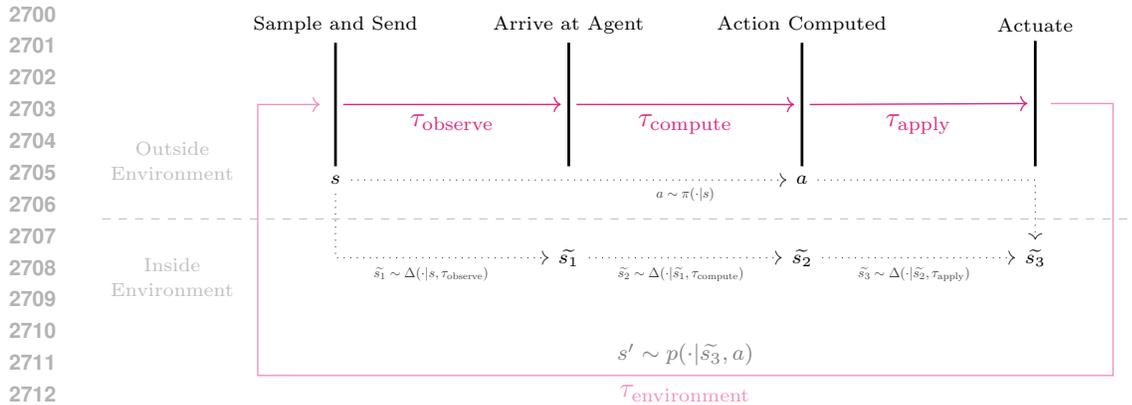


Figure 27: State evolution over the interaction process.

an inverted pendulum, then the interaction delay can result in the state we apply an action to has changed from the state that it was generated from. We illustrate the factors contributing to the interaction delay in Figure 26(a) and how they affect the evolving state in Figure 27. The violation of the equations is shown in Figure 26(b).

While there likely is some time passing within the environment in the real world (as illustrated by  $\tau_{\text{environment}}$  in Figure 27), we consider that time  $\tau_{\text{environment}}$  as part of the environment dynamics and not of the interaction delay.

These times may also be stochastic and unknown to the agent before generating the action. While they can be assumed identically and independently distributed (iid), effects such as clogging (over network or computation bandwidth) mean that a long delay is more likely to follow another long delay, resulting in a dependence in distributions. Delays can also be affected by how an agent interacts with the environment, for example, by controlling a system such that it moves to another access point on the network.

## G.2 INTERACTION LAYER TO ENFORCE DISCRETE DELAY

If we assume that the environment will be excited every  $\tau_{\text{environment}}$  seconds, then we can express the delay as a discrete number of steps, rounded up to the nearest multiple of  $\tau_{\text{environment}}$ .

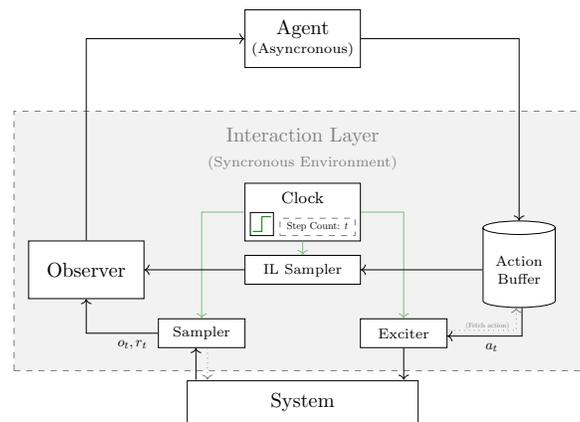


Figure 28: Illustration of the components that make up the interaction layer. A clock is used to ensure that the system is excited and sampled every  $\tau_{\text{environment}}$ .

Delay being expressed in real time makes it inconvenient to reason about with respect to the MDP describing the environment interaction. To resolve this, we introduce an *interaction layer* that sits

---

2754 between the agent and the system we want to control. The primary role of the interaction layer is  
2755 to discretize time and ensure that  $\tau_{\text{environment}}$  is constant. It operates under the assumptions that the  
2756 interaction layer can

- 2757
- 2758 1. observe the system at any time (read sensors),
  - 2759 2. excite the system at any time (apply actions), and
  - 2760 3. observe and excite with negligible real-world delay (assumed  $\tau = 0$ ).
- 2761

2762 Under these assumptions, the role of the interaction layer is primarily to

- 2763
- 2764 1. maintain an *action buffer* of upcoming actions to apply to the system,
  - 2765 2. accept incoming actions from an agent and insert them into the *action buffer*,
  - 2766 3. ensure that interaction with the system occurs periodically on a fixed interval, and
  - 2767 4. transmit state information back to the agent.
- 2768

2769

2770 The construction of an interaction layer is realistic for many real-world systems. Using the scenario  
2771 illustrated in Figure 1 as an example, the interaction layer could be implemented as a microcontroller  
2772 located on the vehicle itself. We illustrate the interaction layer in Figure 28.

2773 From this perspective, the agent acts reactively. The interaction layer manages the interaction with  
2774 the environment, and the agent generates new actions for the action buffer when triggered by emis-  
2775 sions from the interaction layer. We denote emitted data from the interaction layer as an *observation*  
2776 *packet*  $\mathbf{o}_t$ , and the data sent to the interaction layer as an *action packet*  $\mathbf{a}_t$ .

2777  
2778  
2779  
2780  
2781  
2782  
2783  
2784  
2785  
2786  
2787  
2788  
2789  
2790  
2791  
2792  
2793  
2794  
2795  
2796  
2797  
2798  
2799  
2800  
2801  
2802  
2803  
2804  
2805  
2806  
2807