

Figure 4: BPQL architecture: The agent selects an action using an augmented state, which is obtained from the temporary buffer, as an input of the policy. At every time $t > 2d$, a transition tuple $\bar{s}_{t-d}, s_{t-d}, a_{t-d}, r_{t-d}, \bar{s}_{t+1-d}, s_{t+1-d}$ is stored in the replay memory, and the policy and beta critic are trained by minimizing $J_{\bar{\pi}}$ and $J_{Q_{\beta}}$, respectively.

Algorithm 1 Belief-Projection-Based Q-Learning (BPQL)

```

1: initialize a policy  $\bar{\pi}_{\phi}(a|\bar{s})$ , beta critics  $Q_{\theta,\beta}(s, a)$ , target beta critics  $Q_{\theta,\beta}^{\text{tar}}(s, a)$ , an empty replay
   memory  $\mathcal{D}$ , a temporary buffer  $\mathcal{B}$ , delayed timesteps  $d$ , learning rate for the beta critic  $\lambda_{Q_{\beta}}$ ,
   learning rate for the policy  $\lambda_{\bar{\pi}}$ , and soft update ratio  $\tau$ 
2: for episode = 1 to  $E$  do
3:   for  $t = 1$  to  $H$  do
4:     if  $t < d$  then
5:       executes action  $a_t$  randomly or ‘No-Op’ action
6:       put  $a_t$  to  $\mathcal{B}$ 
7:     else if  $t = d$  then
8:       choose action  $a_d$  randomly or ‘No-Op’ action
9:        $s_1 \leftarrow \text{Env}(a_d)$ 
10:      put  $s_1$  and  $a_d$  to  $\mathcal{B}$ 
11:    else
12:      get  $s_{t-d}, a_{t-d}, \dots, a_{t-1}$  from  $\mathcal{B}$ 
13:       $\bar{s}_t \leftarrow (s_{t-d}, a_{t-d}, \dots, a_{t-1})$ 
14:       $a_t \leftarrow \pi_{\phi}(\bar{s}_t)$ 
15:       $s_{t+1-d}, r_{t-d} \leftarrow \text{Env}(a_t)$ 
16:      put  $s_{t+1-d}$  and  $a_t$  to  $\mathcal{B}$ 
17:      if  $t > 2d$  then
18:        get  $s_{t-2d}, s_{t-2d+1}, s_{t-d}, a_{t-2d}, \dots, a_{t-d}$  from  $\mathcal{B}$ 
19:         $\bar{s}_{t-d} \leftarrow (s_{t-2d}, a_{t-2d}, \dots, a_{t-d-1})$ 
20:         $\bar{s}_{t-d+1} \leftarrow (s_{t-2d+1}, a_{t-2d+1}, \dots, a_{t-d})$ 
21:        store  $(\bar{s}_{t-d}, s_{t-d}, a_{t-d}, r_{t-d}, \bar{s}_{t+1-d}, s_{t+1-d})$  in  $\mathcal{D}$ 
22:        pop  $s_{t-2d}, a_{t-2d}$  from  $\mathcal{B}$ 
23:      end if
24:    end if
25:  end for
26:  for each gradient step do
27:     $\omega \leftarrow \omega - \lambda_{Q_{\beta}} \nabla_{\theta} J_{Q_{\beta}}(\omega)$ 
28:     $\phi \leftarrow \phi - \lambda_{\bar{\pi}} \nabla_{\phi} J_{\bar{\pi}}(\phi)$ 
29:     $\bar{\omega} \leftarrow \tau \omega + (1 - \tau) \bar{\omega}$ 
30:  end for
31: end for

```

381 B Proof of Proposition 3.2

382 **Proposition 3.2** *Let the projection operator on the column-space of the belief matrix \mathbf{B} with respect*
 383 *to the weighted Euclidean norm $\|\cdot\|_{\mathbf{W}}$ be $\Pi_{\mathbf{W}}$, where $\mathbf{W} = [\rho^{\bar{\pi}}(\bar{s}_1), \rho^{\bar{\pi}}(\bar{s}_2), \dots, \rho^{\bar{\pi}}(\bar{s}_j)]$ is a steady-*
 384 *augmented state probability vector, and the Markov chain be irreducible i.e, $\rho^{\bar{\pi}}(\bar{s}_k) > 0$ for all*
 385 *$k \in \{1, 2, \dots, j\}$. Then the combined operator $\Pi_{\mathbf{W}}\bar{T}^{\bar{\pi}}$ is γ -contraction with respect to $\|\cdot\|_{\mathbf{W}}$.*

386 *Proof.* Let $\mathbf{V}_{\text{proj.}}^1$ and $\mathbf{V}_{\text{proj.}}^2$ be an arbitrary vector in \mathbb{R}^j , $\bar{P}^{\bar{\pi}}(\bar{s}'|\bar{s})$ be a transition probability of $\bar{s} \rightarrow \bar{s}'$
 387 when the agent follow the policy $\bar{\pi}$, and $V_k^{1,\text{proj.}}, V_k^{2,\text{proj.}}$ be the k-th element of the $\mathbf{V}_{\text{proj.}}^1$ and $\mathbf{V}_{\text{proj.}}^2$.
 388 respectively.

$$\|\Pi_{\mathbf{W}}\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^1 - \Pi_{\mathbf{W}}\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^2\|_{\mathbf{W}}^2 \quad (28)$$

$$= \|\Pi_{\mathbf{W}}(\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^1 - \bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^2)\|_{\mathbf{W}}^2 \quad (29)$$

$$\leq \|\Pi_{\mathbf{W}}(\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^1 - \bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^2)\|_{\mathbf{W}}^2 + \|(\mathbf{I} - \Pi_{\mathbf{W}})(\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^1 - \bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^2)\|_{\mathbf{W}}^2 \quad (30)$$

$$= \|\bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^1 - \bar{T}^{\bar{\pi}}\mathbf{V}_{\text{proj.}}^2\|_{\mathbf{W}}^2 \quad (31)$$

$$= \|\bar{\mathbf{R}} + \gamma\bar{\mathbf{P}}\mathbf{V}_{\text{proj.}}^1 - (\bar{\mathbf{R}} + \gamma\bar{\mathbf{P}}\mathbf{V}_{\text{proj.}}^2)\|_{\mathbf{W}}^2 \quad (32)$$

$$= \gamma \sum_{k=1}^j \rho^{\bar{\pi}}(\bar{s}_k) \left(\sum_{l=1}^j \bar{P}^{\bar{\pi}}(\bar{s}_l|\bar{s}_k) \left(V_l^{1,\text{proj.}} - V_l^{2,\text{proj.}} \right) \right)^2 \quad (33)$$

$$\leq \gamma \sum_{k=1}^j \rho^{\bar{\pi}}(\bar{s}_k) \sum_{l=1}^j \bar{P}^{\bar{\pi}}(\bar{s}_l|\bar{s}_k) \left(V_l^{1,\text{proj.}} - V_l^{2,\text{proj.}} \right)^2 \quad (34)$$

$$= \gamma \sum_{l=1}^j \sum_{k=1}^j \rho^{\bar{\pi}}(\bar{s}_k) \bar{P}^{\bar{\pi}}(\bar{s}_l|\bar{s}_k) \left(V_l^{1,\text{proj.}} - V_l^{2,\text{proj.}} \right)^2 \quad (35)$$

$$= \gamma \sum_{l=1}^j \rho^{\bar{\pi}}(\bar{s}_l) \left(V_l^{1,\text{proj.}} - V_l^{2,\text{proj.}} \right)^2 \quad (36)$$

$$\begin{aligned} & \because \sum_{k=1}^j \rho^{\bar{\pi}}(\bar{s}_k) \bar{P}^{\bar{\pi}}(\bar{s}_l|\bar{s}_k) = \rho^{\bar{\pi}}(\bar{s}_l) \text{ by the definition of } \rho^{\bar{\pi}} \\ & = \gamma \|\mathbf{V}_{\text{proj.}}^1 - \mathbf{V}_{\text{proj.}}^2\|_{\mathbf{W}}^2, \end{aligned} \quad (37)$$

389 where Equality (31) holds by the Pythagorean Theorem, and Inequality (34) follows from the Jensen's
 390 inequality. \square

391 C Environment Details

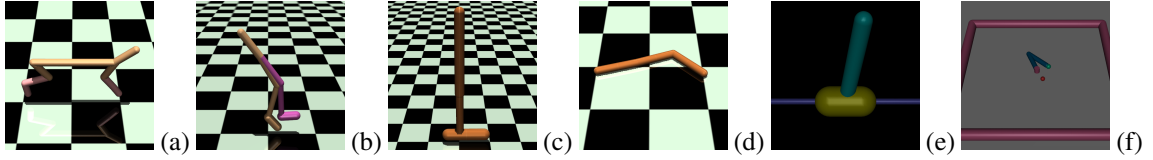


Figure 5: MuJoCo continuous control environments in the experiment: (a) HalfCheetah-v3, (b) Walker2d-v3, (c) Hopper-v3, (d) Swimmer-v3, (e) InvertedPendulum-v2, and (f) Reacher-v2.

Table 2: Details of the MuJoCo environment

	Action dimension	State dimension	Timestep (s)
HalfCheetah-v3	6	17	0.05
Walker2d-v3	6	17	0.008
Hopper-v3 (s)	3	11	0.008
Swimmer-v3	2	6	0.04
InvertedPendulum-v2	1	4	0.04
Reacher-v2	2	11	0.02

392 D Implementation Details

393 In this section, we provide the implementation details of the algorithms used in this study. All
 394 methods presented in the experiment used an SAC as their base learning algorithm with the following
 395 characteristics:

- 396 • Stochastic Gaussian policy approaches.
- 397 • Automated entropy adjustment [Haarnoja et al., 2018].
- 398 • Clipped-double Q-learning, which was introduced in the TD3 algorithm to prevent overesti-
 399 mating the Q-value [Fujimoto et al., 2018].
- 400 • Soft target update, which changes the target values slowly and improves the learning
 401 stability[Lillicrap et al., 2015].
- 402 • Adam optimizer, a variant of the stochastic gradient descent method [Kingma and Ba, 2014].

The details of the hyperparameters are presented in Table 3.

Table 3: Hyperparameters for BPQL and the baselines.

Hyperparameters	Values
Critic network	256, 256
Policy network	256, 256
Discount factor	0.99
Replay memory size	1 M
Minibatch size	256
Learning rate	0.0003
Target entropy	$-\dim \mathcal{A} $
Target smoothing coefficient	0.995
Optimizer	Adam

403

404 E Plots of Performance Comparison

405 In this section, we present several plots of the performance curves of BPQL and other baselines in
 406 the environments with delayed timesteps of 3, 6 and 9.



Figure 6: Performance curves for each algorithm in continuous tasks with three, six and nine delayed timesteps.

F Additional Experiments

F.1 Extremely Long Delayed Environment

We conducted additional experiments to determine how well BPQL could solve the control problem even in very long delayed environments. Figure 7 shows that in BPQL, the policy is improved through interaction with the environment in a very long delayed environment, but the conventional methods find it difficult improve their policy.

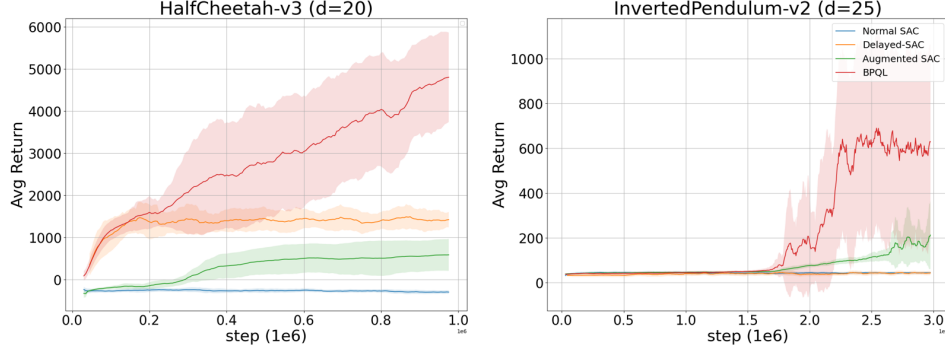


Figure 7: HalfCheetah-v3 environment with 20 delayed timesteps (*left*) and InvertedPendulum-v2 environment with 25 delayed timesteps (*right*). Each timestep is 0.05 s and 0.04 s for HalfCheetah-v3 and InvertedPendulum-v2 environments, respectively. In the InvertedPendulum-v2 environment, to prevent early failure, the agent uses a pretrained policy to determine the actions of the first 25 (=number of delayed timesteps) timesteps. We repeated the test for the tasks five times with different random seeds.

F.2 Action delay

We also evaluated BPQL and other baselines in *action delayed* environments. Figure 8 illustrates the interaction between an agent and action delayed environment.

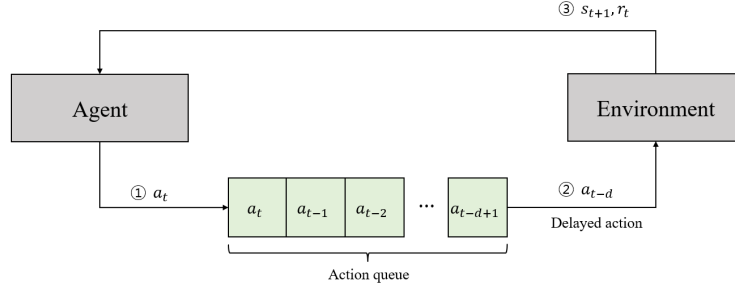


Figure 8: In an action delayed environment, the environment takes the delayed action a_{t-d} instead of the current action a_t , where d is the number of delayed timesteps.

F.2.1 Augmented State in an Action Delayed Environment

In an action delayed environment, the augmented state at time $t + d$ is defined as:

$$\bar{s}_{t+d} = (s_t, a_{t-d}, a_{t-d+1}, \dots, a_{t-1}), \quad (38)$$

where the d is the number of delayed timesteps for an action. The CDMDP for an action delayed environment is defined as a tuple $(\mathcal{X}, \mathcal{A}, R, P, \gamma, d_a)$, where \mathcal{X} is the original state space, \mathcal{A} is the action space, $R : \mathcal{X} \times \mathcal{A} \mapsto \mathbb{R}$ is the reward function, $P : \mathcal{X} \times \mathcal{A} \times \mathcal{X} \mapsto [0, 1]$ is the transition kernel, $\gamma \in (0, 1)$ is a discount factor, and d_a is the number of delayed timesteps for an action.

422 This CDMDP can be reduced to a MDP $(\mathcal{S}, \mathcal{A}, \bar{\mathcal{R}}, \bar{P}, \gamma)$, where $\mathcal{S} = \mathcal{X} \times \mathcal{A}^{d_a}$ is an augmented state
423 space, $\bar{\mathcal{R}} : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ is the reward mapping, and $\bar{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ is the transition kernel
424 as in the case of observation delayed environments [Katsikopoulos and Engelbrecht, 2003]. BPQL
425 for action delayed environments uses the augmented state defined in Equation (38) instead of the one
426 defined in Equation (6).

427 F.2.2 Results

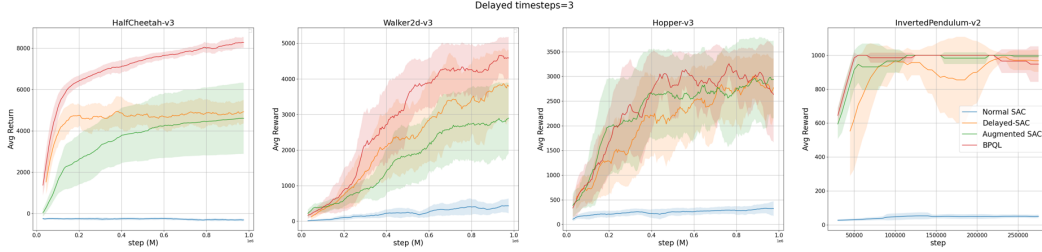


Figure 9: Performance curves for each algorithm in continuous tasks with 3 action delayed timesteps. We repeated the test for the tasks five times with different random seeds.

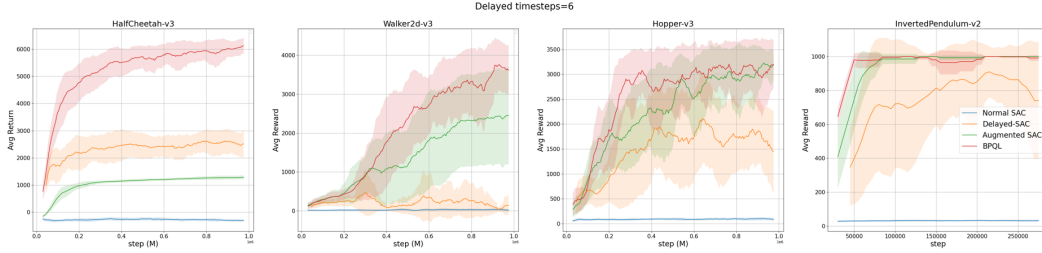


Figure 10: Performance curves for each algorithm in continuous tasks with 6 action delayed timesteps. We repeated the test for the tasks five times with different random seeds.

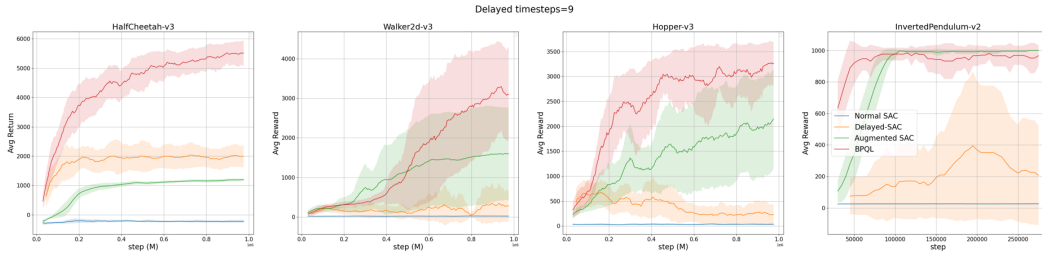


Figure 11: Performance curves for each algorithm in continuous tasks with 9 action delayed timesteps. We repeated the test for the tasks five times with different random seeds.