

BEYOND MARKOV ASSUMPTION: IMPROVING SAMPLE EFFICIENCY IN MDPs BY HISTORICAL AUGMENTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Under the Markov assumption of Markov Decision Processes (MDPs), an optimal stationary policy does not need to consider history and is no worse than any non-stationary or history-dependent policy. Therefore, existing Deep Reinforcement Learning (DRL) algorithms usually model sequential decision-making as an MDP and then try to optimize a stationary policy by single-step state transitions. However, such optimization is often faced with sample inefficiency when the causal relationships of state transitions are complex. To address the above problem, this paper investigates if augmenting the states with their historical information can simplify the complex causal relationships in MDPs and thus improve the sample efficiency for DRL. First, we demonstrate that a complex causal relationship of single-step state transitions may be inferred by a simple causal function of the historically augmented states. Then, we propose a convolutional neural network architecture to learn the representation of the current state and its historical trajectory. The main idea of this representation learning is to compress the high-dimensional historical trajectories into a low-dimensional space. In this way, we can extract the simple causal relationships from historical information and avoid the overfitting caused by high-dimensional data. Finally, we formulate Historical Augmentation Aided Actor-Critic (HA3C) algorithm by adding the learned representations to the actor-critic method. The experiment on standard MDP tasks demonstrates that HA3C outperforms current state-of-the-art methods in terms of both sample efficiency and performance.

1 INTRODUCTION

Sequential decision-making widely exists in real-world control tasks, such as robot control and autonomous driving (Dorf & Bishop, 2011; Ibarz et al., 2021; Sallab et al., 2017). Generally speaking, it can be modelled as a Markov Decision Process (MDP), where an agent iteratively takes action in an environment for transiting from one state to another (Puterman, 1990). Each transition is evaluated by a reward signal passing from the environment to the agent so that Reinforcement Learning (RL) can learn the optimal policy by maximizing the cumulative reward (Sutton & Barto, 2018). The Markov Assumption of MDPs asserts that the probability distributions of the reward and next state depend only on the current state and action. Under the Markov assumption of MDPs, there exists an optimal stationary policy which does not need to consider history and is no worse than any non-stationary or history-dependent policy (Puterman, 2014). Therefore, existing RL algorithms usually try to optimize a stationary policy by single-step state transitions.

With advances in deep learning, many effective Deep RL (DRL) methods were proposed (Fujimoto et al., 2018; Haarnoja et al., 2018; Lillicrap et al., 2016; Mnih et al., 2016; 2015). Under the Markov assumption of MDPs, they are usually based on the actor-critic method where the critic estimates the Q -function, i.e., the expected cumulative reward after taking action at each state, while the actor updates the policy to choose the action which can maximize the estimated Q -function (Schulman et al., 2015; Silver et al., 2014). However, such optimization may miss the useful causal relationships of state transitions, leading to sample inefficiency (Allen et al., 2021; Buckman et al., 2018; Du et al., 2020; Guo et al., 2020). An existing partial solution to this issue is representation learning in which a neural network is trained to infer the causal relationships of state transitions by predicting the

reward or future state of each state-action pair (Munk et al., 2016; Ni et al., 2023; Ravindran, 2004; Rezaei-Shoshtari et al., 2022). Then, the sample efficiency of DRL can be improved by adding the learned representations to the actor-critic method. Unfortunately, it is hard to train the neural networks which can infer complex causal relationships, e.g., polynomial causal relationships and the basic laws of physics (Andoni et al., 2014; Cranmer et al., 2020). Standard complexity-theoretic results strongly suggest that there is no algorithm efficient enough for learning arbitrary target functions, even for target functions representable by very low-depth networks (Applebaum et al., 2006). Therefore, the sample efficiency for DRL is still limited in complex MDP tasks.

This paper addresses the above problem by augmenting the states with their historical information. Based on the analysis in Section 3, we believe that historical augmentation can simplify the causal relationships of state transitions by its inherent contextual information and increasing the search space of the causal functions (Hallak et al., 2015; Sprunger & Jacobs, 2019). Therefore, we focus on optimizing a history-dependent stationary policy in an MDP. Our DRL approach comprises two key components: 1) Learning the state representations to capture the causal relationships in an MDP and 2) finding the optimal stationary policy by the learned representations. Given an action and the historically augmented current state, our representation learning utilizes a Convolutional Neural Network (CNN) architecture to compress the high-dimensional historical trajectory of the given state into a low-dimensional space while predicting the future state. The compressed historical trajectories can be seen as the abstracted features which can represent the simple causal relationships and avoid the overfitting caused by high-dimensional data (Andre & Russell, 2002). To keep the Markov assumption of MDPs, our representation learning does not compress the current state. We add the learned state representations to the actor-critic method. In this way, the causal relationships captured by our representation learning can be utilized to estimate the Q -function and update policy. Therefore, our new DRL approach can optimize the policy in a complex MDP with high sample efficiency. We combine historical augmentation, state representations, and TD3 in our approach to formulate a new DRL algorithm, Historical Augmentation Aided Actor-Critic (HA3C). The experiment on standard MDP tasks, i.e. Mujoco control tasks and Deep Mind Control (DMC) suite, empirically demonstrates that HA3C outperforms current state-of-the-art methods in terms of both sample efficiency and performance (Brockman et al., 2016; Todorov et al., 2012; Tassa et al., 2018).

Our contributions are as follows: 1) Existing RL methods usually utilize historical information to recover Markov assumption in dynamics. It is the first time in the literature that historical augmentation can be used to improve sample efficiency when Markov assumption is satisfied. 2) We propose a new DRL approach to address the problem of how to effectively utilize the historical information in MDPs. 3) Based on this approach, we formulate a new RL algorithm, HA3C, which outperforms existing state-of-the-art DRL algorithms, e.g. TD7 (Fujimoto et al., 2023). 4) Our examples, experiment, and discussion illustrate that in fact, DRL needs to consider historical information in complex MDP tasks.

2 BACKGROUND

An MDP can be written as a 5-tuple $\mathbb{M} = \langle \mathcal{S}, \mathcal{A}, R, \mathbf{P}, \gamma \rangle$ with state space \mathcal{S} , action space \mathcal{A} , reward function R , transition model \mathbf{P} , and discount factor γ . In an MDP, RL can maximize the discounted cumulative reward by learning how to map the states to the actions (Baird, 1995; Duan et al., 2016; Williams, 1992). For a given state $s_t \in \mathcal{S}$ at time step t , the agent executes an action $\mathbf{a}_t \in \mathcal{A}$ w.r.t. a policy $\pi : \mathcal{S} \mapsto \mathcal{A}$, to obtain a reward $r_t = R(s_t, \mathbf{a}_t)$ and transfer to a new state s_{t+1} . The return of the agent is defined as the discounted cumulative reward $G_t = \sum_{i=t}^{+\infty} \gamma^{i-t} r_i$. Based on the Markov assumption of MDPs, RL can find the optimal policy to maximize the following value function which is the expected return when $s_t = s$ and following π thereafter.

$$V^\pi(s) = \mathbb{E}^\pi [G_t | s_t = s] = \mathbb{E}^\pi \left[\sum_{i=0}^{+\infty} \gamma^i r_{t+i} | s_t = s \right],$$

where $\mathbb{E}^\pi [*]$ denotes the expected value of a random variable given that the agent follows policy π .

With advances in deep learning, combining neural networks into RL has drawn significant attention in the literature. Many DRL algorithms learn the optimal policy by the actor-critic method (Kaelbling et al., 1996), where the critic network estimates the Q -function which is the expected return when

108 $\mathbf{s}_t = \mathbf{s}$, $\mathbf{a}_t = \mathbf{a}$, and following policy π thereafter.

$$109$$

$$110 Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}^\pi [G_t | \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}] = \mathbb{E}^\pi \left[\sum_{i=0}^{+\infty} \gamma^i r_{t+i} | \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} \right],$$

111

112 while the actor network updates the policy to maximize the estimated Q -function.

113

114 To improve sample efficiency, some DRL methods learn the state representations of the collected

115 state transitions and then add the learned representations to the actor-critic method (Anand et al.,

116 2019; Dayan, 1993; Gelada et al., 2019; Li et al., 2006). This representation learning aims to

117 capture the causal relationships in MDPs, and thus improves sample efficiency (Liu et al., 2020;

118 Van Hoof et al., 2016; Ye et al., 2023; Zhang et al., 2021). For example, ML-DDPG learns the state

119 representations by predicting the next state representation and the reward in MDPs (Munk et al.,

120 2016). As an improvement of ML-DDPG, OFENet learns the high-dimensional state representations

121 by predicting the next state in DenseNet architecture (Ota et al., 2020). TD7 improves the learning

122 of state representations by AvgL1Norm and then combines the learned representations with TD3,

123 checkpoints, and prioritized replay buffer (Fujimoto et al., 2023).

124 DRL algorithms need to consider historical information when the Markov assumption of MDPs

125 is violated (Eysenbach et al., 2020; Majeed & Hutter, 2018; Hafner et al., 2019b). For Partially

126 Observable MDPs (POMDPs), in which the states are partially observable, deep recurrent Q -network

127 uses LSTMs to encode state transition trajectories in the Q -function estimation (Hausknecht &

128 Stone, 2015). As an improvement of deep recurrent Q -network, deep transformer Q -network uses

129 transformers to encode the state transition trajectories (Esslinger et al., 2022). As a famous DRL

130 algorithm, Dreamer encodes the historical information into the state at every time step in POMDPs (Ha

131 & Schmidhuber, 2018; Hafner et al., 2019a). In delayed MDPs, in which the current state and reward

132 may arrive at the agent with a delay (Katsikopoulos & Engelbrecht, 2003), researchers usually recover

133 the Markov assumption of MDPs by considering the historical actions (Bouteiller et al., 2020; Derman

134 et al., 2021). When the Markov assumption of MDPs is violated by the state abstraction, it is possible

135 to find a history-based policy which works in the abstracted space and is of the same quality as

136 optimal policy (Patil et al., 2024). However, the history-based DRL for the dynamics which are under

137 Markov assumption is largely absent from the literature.

138 3 MOTIVATION

139

140 Let $\mathbf{h}_t = \{\mathbf{s}_0, \mathbf{a}_0, \dots, \mathbf{s}_t\}$ as the history up to time step t in a sequential decision-making task. The

141 optimal policy may change the decision rule in different time steps and select actions based on

142 historical information. In this case, we should optimize a history-dependent policy $\pi = \{d_t | t =$

143 $0, 1, \dots\}$ which selects action at time step t by decision-rule $d_t(\mathbf{a}_t | \mathbf{h}_t)$. Fortunately, based on the

144 Markov assumption of MDPs, there is an optimal stationary policy $\pi(\mathbf{a}_t | \mathbf{s}_t)$ which is unrelated to

145 time and selects action \mathbf{a}_t by only the state \mathbf{s}_t . This Markov assumption asserts that the probability

146 distributions of state \mathbf{s}_{t+1} and reward r_t depend only on the \mathbf{s}_t and \mathbf{a}_t as

$$147 P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | \mathbf{s}_0, \mathbf{a}_0, r_0, \dots, \mathbf{s}_t, \mathbf{a}_t\} = P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | \mathbf{s}_t, \mathbf{a}_t\},$$

148 where P is the probability distribution in \mathcal{P} . Let HR and SR denote the class of history-dependent

149 and stationary policies, respectively. Lemma 3.1 is the key of most existing RL algorithms (Puterman,

150 2014)[Thm. 6.2.10]. The different types of policies are detailed in Appendix A.

151 **Lemma 3.1.** *Under the Markov assumption of MDPs, there exists a policy $\pi^* \in SR$ such that, for*

152 *all t , $V_{\pi^*}(\mathbf{s}_t) = \sup_{\pi \in HR} V_\pi(\mathbf{h}_t)$.*

153

154 Based on Lemma 3.1, existing DRL algorithms for MDPs usually optimize a stationary policy by

155 single-step transitions. If the causal relationships in the modelled MDP are simple, e.g., there are

156 only linear causal relationships in this MDP, such optimization effectively finds the optimal policy.

157 A classical result is that a neural network with a single hidden layer can successfully learn a linear

158 function (Andoni et al., 2014). However, it is still hard to capture complex causal relationships by

159 neural networks. Standard complexity-theoretic results strongly suggest that there is no algorithm

160 efficient enough for learning arbitrary functions, even for target functions representable by very

161 low-depth networks (Applebaum et al., 2006). In fact, a more complex causal function requires

neural networks to approximate with more parameters, samples, and time consumption (Bianchini & Scarselli, 2014).

Historical augmentation has the potential to address the above problem by simplifying the causal relationships in MDPs as it can increase the search space of the causal functions and provide much contextual information on state transitions (Hallak et al., 2015; Sodhani et al., 2022).

Example 3.1. For example, if we model the state transitions with Fibonacci sequence as $s_0 = 1, s_1 = 1, s_2 = 2, s_3 = 3, s_4 = 5, \dots$, when $t > 2$, the state transitions in this model will satisfy the Markov assumption of Markov Processes as (Dynkin, 1965)

$$P\{s_{t+1} = s' | s_0, \dots, s_t\} = P\{s_{t+1} = s' | s_t\}.$$

Without considering history, at s_t , s_{t+1} will be predicted by a complex time-related formula

$$s_{t+1} = \frac{1}{\sqrt{5}} \left[\left(\frac{1 + \sqrt{5}}{2} \right)^{t+1} - \left(\frac{1 - \sqrt{5}}{2} \right)^{t+1} \right].$$

Fortunately, when considering history, we can predict s_{t+1} by a simple linear function

$$s_{t+1} = s_{t-1} + s_t.$$

In Appendix B, we give another example to illustrate that by historical augmentation, a non-linear causal relationship in single-step transitions may be simplified as a linear causal relationship. Fig. 1(a) presents the original MDP causal relationships and Fig. 1(b) demonstrates the MDP causal relationships with state augmentation. When inferring the causal relationships in a trajectory, the causal

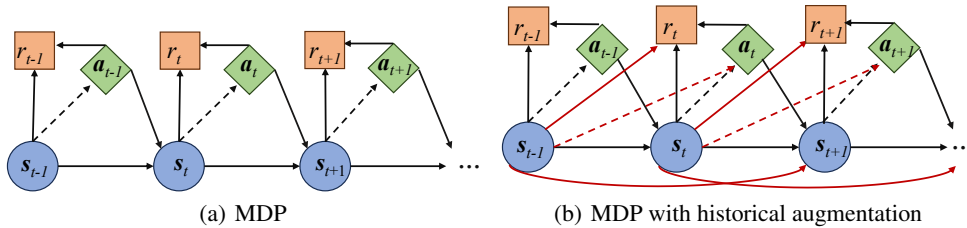


Figure 1: Causal diagrams of an MDP with or without historical augmentation. The black lines index the original MDP causal relationships and the red lines index the added causal relationships, e.g., the causal relationships from historical augmentation. The dashed lines indicate the information needed in the optimization.

function in Fig 1(b) can be simpler than the causal function in Fig 1(a).

From the analysis above, the motivation of our work is that historical information can simplify the complex causal relationships in MDPs and thus has the potential to improve the sample efficiency of DRL. However, the challenges are 1) how to ensure that the causal relationships learned from historical augmentation are simple and 2) avoiding overfitting caused by the high-dimensional historical data.

4 METHOD

In this section, we propose a new DRL approach by the representation learning of historically augmented states. Then, we formulate a new DRL algorithm, HA3C, and illustrate the advantage of this algorithm with a visual example.

4.1 REPRESENTATION LEARNING ON HISTORICALLY AUGMENTED STATES

To address the problem of how to effectively utilize the historical information in MDPs, we propose a new DRL approach by the representation learning of historically augmented states. The main idea of this representation learning is to compress the high-dimensional historical trajectories into a low-dimensional representation space (Andre & Russell, 2002; Li et al., 2006). On the one hand, the compressed historical trajectories can be seen as the abstracted features of the historical information to extract the simple causal relationships. On the other hand, this compression will avoid the overfitting caused by the high-dimensional historical data (Ying, 2019).

216 To keep the Markov assumption of MDPs, our
 217 representation learning does not compress the
 218 current state. Let $\mathbf{s}_{k,t} = \{\mathbf{s}_{t-k+1}, \dots, \mathbf{s}_t\}$. If
 219 $t < k$, one can set each $\mathbf{s}_i \in \mathbf{s}_{k-t,-1}$ by the
 220 zero vector $\mathbf{0}$. The causal diagram of MDP with
 221 our state abstraction is in Fig. 2. As we can
 222 see, when predicting \mathbf{s}_{t+1} by $\mathbf{s}_{k,t}$ and \mathbf{a}_t , the
 223 dimensionality reduction is only performed on
 224 $\mathbf{s}_{k-1,t-1}$.

225 Let $S_k D$ denote the class of the stationary de-
 226 terministic policies based on k -order state tra-
 227 jectories. Theorem 4.1 forms the basis of our
 228 DRL approach. This theorem can be implied by
 229 Lemma 3.1. For completeness, we provide a proof in Appendix D.

230 **Theorem 4.1.** *Under the Markov assumption of MDPs, there exists a stationary deterministic policy*
 231 $\mu^* \in S_k D$ *such that, for all t , $V^{\mu^*}(\mathbf{s}_{k,t}) = \sup_{\pi \in HR} V^\pi(\mathbf{h}_t)$.*

232 To capture the simplified causal relationships in MDPs by historical augmentation, we define a pair
 233 of encoders $\mathbf{z}^{\mathbf{s}_{k,t}} = f(\mathbf{s}_{k,t})$ and $\mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t} = g(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)$. Based on the Markov assumption in MDPs,
 234 we can predict $\mathbf{z}^{\mathbf{s}_{k,t+1}}$, i.e., the representation of $\mathbf{s}_{k,t+1}$, by $\mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t}$. Thus, the two encoders
 235 are trained by minimizing the following predicting loss:

$$236 L(f, g) = \|g(f(\mathbf{s}_{k,t}), \mathbf{a}_t) - |f(\mathbf{s}_{k,t+1})|_{\times}\|_2^2 = \|\mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t} - |\mathbf{z}^{\mathbf{s}_{k,t+1}}|_{\times}\|_2^2, \quad (1)$$

237 where $|\cdot|_{\times}$ denotes the stop-gradient operation. As presented in Fig. 3, a simple yet effective CNN
 238 network architecture is utilized in our representation learning. In the network of $f(\mathbf{s}_{k,t})$, we first use
 239 a CNN layer to mine the historical information in $\mathbf{s}_{k-1,t-1}$. This layer produces the feature maps of
 240 $\mathbf{s}_{k-1,t-1}$ by the multiple filters, which are as wide as the state dimensionality. Second, we utilize a
 241 max pooling layer to capture the most important features and an average pooling layer to capture
 242 the tendency features. Third, we compress the captured features into a low-dimensional space and
 243 learn the features of \mathbf{s}_t . Finally, we concatenate the compressed features of $\mathbf{s}_{k-1,t-1}$ and the learned
 244 features of \mathbf{s}_t . The concatenated features are the input of the next fully connected layer to obtain the
 245 representation $\mathbf{z}^{\mathbf{s}_{k,t}}$. In the network of $g(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)$, we put the concatenation of $\mathbf{z}^{\mathbf{s}_{k,t}}$ and
 246 \mathbf{a}_t into the two fully connected layers to obtain the representation $\mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t}$.

247 We combine our learned representations with the actor-critic method and thus the Q -function can
 248 be defined as $\hat{Q}(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)$ and the policy can be defined as $\mu(\mathbf{z}^{\mathbf{s}_{k,t}}) \in S_k D$. Define the probability
 249 distribution of $\mathbf{z}^{\mathbf{s}_{k,t+1}}$ under μ as

$$250 P^\mu\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,t}} = \mathbf{z}^{\mathbf{s}_{k,:}}\} = \int_{\mathcal{Z}} \mathbb{E}_{\mathbf{a} \sim \mu(\mathbf{z}^{\mathbf{s}_{k,:}})} [p(\mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})] d\mathbf{z}^{\mathbf{s}_{k,:}},$$

251 where $\mathbf{s}_{k,:}$ is a k -order state trajectory $\{\mathbf{s}_0, \dots, \mathbf{s}_{k-1}\}$ ending with \mathbf{s} , i.e., $\mathbf{s}_{k-1} = \mathbf{s}$, \mathcal{Z} is the set of all
 252 possible $\mathbf{z}^{\mathbf{s}_{k,:}}$, and $p(\mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})$ is the probability of transferring to $\mathbf{z}^{\mathbf{s}'_{k,:}}$ with taking \mathbf{a} at $\mathbf{z}^{\mathbf{s}_{k,:}}$.
 253 Our optimization is based on a Bellman optimality operator B for μ as

$$254 B_\mu \hat{Q}(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a}) = \max_{\mu} \mathbb{E}_{\mathbf{a}_{t+1} \sim \mu, \mathbf{z}^{\mathbf{s}_{k,t+1}} \sim P^\mu} [r_t + \gamma \hat{Q}(\mathbf{z}^{\mathbf{s}_{k,t+1}}, \mathbf{a}_{t+1})]. \quad (2)$$

255 The following theorem gives the conditions to find the optimal stationary policy in our approach. The
 256 proof of this theorem is given in Appendix D.

257 **Theorem 4.2.** *Given a finite MDP, if 1) $f(\cdot)$ and $g(\cdot)$ are fixed, 2) $\forall \mathbf{s}_{k,:}, \mathbf{s}'_{k,:} \in S_{k,:}, \mathbf{s} \neq \mathbf{s}' \Leftrightarrow$
 258 $\mathbf{z}^{\mathbf{s}_{k,:}} \neq \mathbf{z}^{\mathbf{s}'_{k,:}}$, and 3) $L(f, g) \rightarrow 0$, then $\hat{Q}(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)$ converges to the optimal $Q^*(\mathbf{s}_t, \mathbf{a}_t)$ by the
 259 Bellman optimality operator in equation 2.*

260 This theorem illustrates that no matter whether different historical trajectories lead to different
 261 representations on the state \mathbf{s} , we can still find the optimal stationary policy in the representation
 262 space. To make condition 2) hold, we can increase the dimensionality of \mathbf{s} in representation learning.
 263 This operation also can improve sample efficiency (Ota et al., 2020). To see condition 3) hold, there
 264 should exist a $\mathbf{s}'_{k,:}$ that satisfies

$$265 p\{\mathbf{s}_{k,t+1} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t}, \mathbf{a}_t\} \rightarrow 1.$$

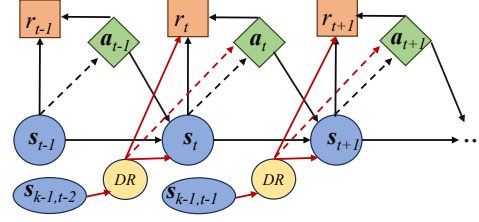


Figure 2: Causal diagram of a historically augmented MDP with state abstraction. DR represents the operation of dimensionality reduction.

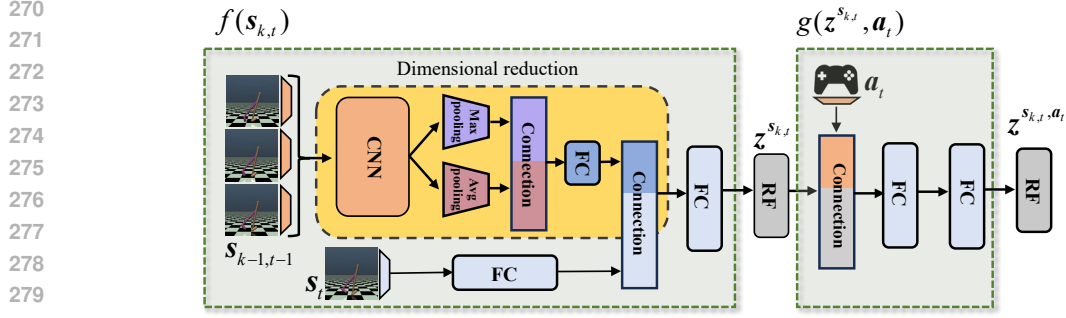


Figure 3: Network architecture of our representation learning. FC represents a fully connected layer and RF represents the state representation features.

There is an analysis of the function approximation error in Appendix D. We add $z^{s_{k,t}, a_t}$ to \hat{Q} to consider the learned relationship between a_t and $z^{s_{k,t}}$ in the representation space. We also add s_t to \hat{Q} and μ to consolidate the relationships in single-step transitions. Thus Q and μ can be written as $\hat{Q}(z^{s_{k,t}, a_t}, z^{s_{k,t}}, s_t, a_t)$ and $\mu(z^{s_{k,t}}, s_t)$, respectively. The operations in \hat{Q} and μ are shown in Fig. 4.

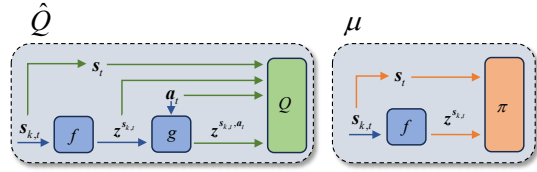


Figure 4: The operations in Q and μ .

Our approach can be connected with POMDPs, High-order MDPs (HMDPs), and state abstraction. A detailed analysis of the connections between our approach and the related work is shown in Appendix C.

4.2 HA3C ALGORITHM

In this subsection, we propose HA3C algorithm which is a combination of TD3, representation learning, historical augmentation. HA3C has several networks as follows. Two critic networks ($\hat{Q}_{\phi_1}, \hat{Q}_{\phi_2}$), two target critic networks ($\hat{Q}_{\phi_1^T}, \hat{Q}_{\phi_2^T}$), an actor network μ_{θ} , a target actor network μ_{θ^T} , two encoders (f_{σ}, g_{σ}), two fixed encoders ($f_{\sigma^F}, g_{\sigma^F}$), two target encoders ($f_{\sigma^T}, g_{\sigma^T}$), a checkpoint actor network π_{θ^C} , and a checkpoint encoder f_{σ^C} .

To learn the representations with historical augmentation, f_{σ} and g_{σ} are trained by the transitions in buffer $\mathcal{B} = \{s_{k,i}, a_i, r_i, s_{k,i+1}\}$ to minimize the predicting loss in equation 1. For any parameter set α , we define

$$z_{\alpha}^{s_{k,t}} = f_{\alpha}(s_{k,t}), \quad z_{\alpha}^{s_{k,t}, a_t} = g_{\alpha}(z^{s_{k,t}}, a_t).$$

Based on the assumption that f_{σ^F} and g_{σ^F} satisfy the conditions in Theorem 4.2 on the most transitions in \mathcal{B} , the Q -function is estimated by the following Huber loss function (Huber, 1992).

$$L(\phi_i, \mathcal{B}) = \text{Huber}(s_{k,t}, a_t, r_t, s_{k,t+1}) \sim_{\mathcal{B}} [x_t - (\hat{Q}_{\phi_i}(z_{\sigma^F}^{s_{k,t}, a_t}, z_{\sigma^F}^{s_t}, s_t, a_t))], \quad (3)$$

$$x_t = r_t + \gamma \text{clip}(\min(\hat{Q}_{\phi_i^T}(z_{\sigma^T}^{s_{k,t+1}, a'}, z_{\sigma^T}^{s_{t+1}}, s_{t+1}, a')), \hat{Q}^{\min}, \hat{Q}^{\max}),$$

$$a' = \mu_{\theta^T}(z_{\sigma^T}^{s_{k,t+1}}, s_{t+1}) + \epsilon_T, \epsilon_T \sim \mathcal{N},$$

where ϵ_T is target policy noise (Fujimoto et al., 2018), \mathcal{N} is a Gaussian distribution $\mathcal{N}(0, \sigma)$, and \hat{Q}^{\min} and \hat{Q}^{\max} are updated at each time step as

$$\hat{Q}^{\max} \leftarrow \max(x_t, \hat{Q}^{\max}), \quad \hat{Q}^{\min} \leftarrow \min(x_t, \hat{Q}^{\min}).$$

Based on the learned Q -function, the policy network π_{θ} is updated by

$$\max_{\theta} \mathbb{E}_{s_{k,t} \sim \mathcal{B}} \left[\sum_{i=1,2} \hat{Q}_{\phi_i}(z^{s_{k,t}, a}, z^{s_t}, s_t, a) \right], \quad (4)$$

$$a = \mu_{\theta}(z_{\sigma^F}^{s_{k,t}}, s_t).$$

To explore the new actions and thus generate new transitions in \mathcal{B} , exploration noise ϵ is added as

$$\mathbf{a}_t \leftarrow \mathbf{a}_t + \epsilon_e, \epsilon_e \sim \mathcal{N}.$$

In our TD learning, σ^F , σ^T , ϕ^T , and θ^T are updated by

$$\sigma^F \leftarrow \sigma^T, \quad \sigma^T \leftarrow \sigma, \quad \phi^T \leftarrow \phi, \quad \theta^T \leftarrow \theta. \quad (5)$$

Because DRL algorithms are unstable (Henderson et al., 2018; Teh et al., 2017), we use the checkpoint policy to obtain the cumulative reward in our evaluation (Vaswani et al., 2017). In the training of HA3C, if the current policy outperforms the checkpoint policy, we will update the checkpoint policy with the current policy, then $\sigma^C \leftarrow \sigma$ and $\theta^C \leftarrow \theta$. The checkpoint policy can give a more accurate evaluation by maintaining the high-performance policy unchanged. Furthermore, the LAP replay buffer is utilized to store and replay the transitions (Fujimoto et al., 2023; 2020). The algorithm of online HA3C is presented in Algorithm 1.

Algorithm 1 Online HA3C

```

Initialize the hyper-parameters and networks
Initialize replay buffer  $\mathcal{B}$ 
for  $episode = 0$  to  $episode_{max}$  do
    Collect  $k$ -order transitions by  $\mu_\theta$  and store them in LAP buffer  $\mathcal{B}$ 
    if Checkpoint condition then
        if  $\mu_\theta$  outperforms  $\mu_{\theta^c}$  then then
            Update checkpoint networks  $\mu_{\theta^c} \leftarrow \mu_\theta$  and  $f_{\sigma^c} \leftarrow f_\sigma$ 
        end if
    end if
    Sample  $k$ -order transitions from LAP buffer  $\mathcal{B}$ 
    Train the encoder  $f_\sigma$  and  $g_\sigma$  by equation 1
    Train  $\hat{Q}_{\phi_1}$  and  $\hat{Q}_{\phi_2}$  by equation 3
    Train  $\pi_\theta$  by equation 4
    if Target update frequency steps have passed then
        Update target networks by equation 5
    end if
end for

```

Fig. 5 is an example to illustrate the advantage of learning the policy in HA3C. We first collect the obtained states of Walker2d MuJoCo control task by learning the policy with and without historical augmentation, respectively. The max learning step is 4×10^5 . Then we map the collected states in 2D space together by UMAP. Finally, we show the reached states without the learning of historical augmentation in the left subfigure of Fig. 5 and the reached states with the learning of historical augmentation in the right subfigure of Fig. 5. Each state is coloured by the reward of reaching it. As we can see, although the actions to obtain the states in high-reward regions (indexed by the red circles) can be explored, without historical augmentation, it is hard to learn the policy which can regenerate these explored actions. Therefore, in the left subfigure, there are only a few states in the high-reward regions. Fortunately, as shown in the right subfigure, there are a lot of states in the high-reward regions when learning the policy with historical augmentation. The visual results of other environments are shown in Appendix F.

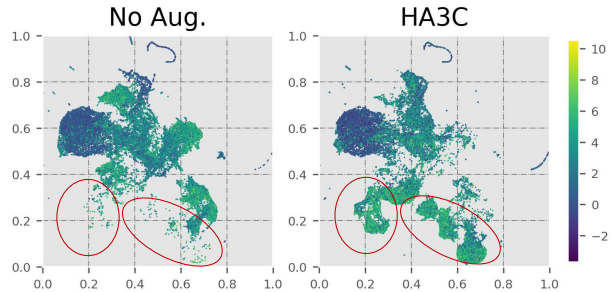


Figure 5: Visual results of the obtained states in Walker2d environment. Each state is coloured by the reward of reaching this state.

5 EXPERIMENTAL RESULT

In this section, first, we compare HA3C to five existing RL algorithms on five Mujoco control tasks (Todorov et al., 2012). Second, we give the ablation study of HA3C to illustrate that historical augmentation is the real source of the improvement in sample efficiency. Third, we analyze the parameter sensitivity on the length of the historical state trajectory and the number of dimensions of compressed historical trajectories. Finally, we give the running times of the different RL algorithms. The experimental setting is in Appendix E. Appendix F has some supplementary experiments including the state visualization and DMC experiment (Tassa et al., 2018).

5.1 COMPARATIVE EVALUATION

In this subsection, we evaluate our HA3C on five MuJoCo control tasks including Walker2d, HalfCheetah, Ant, Humanoid, and Hopper. The compared algorithms are TD3 (Fujimoto et al., 2018), SAC (Haarnoja et al., 2018), TQC (Kuznetsov et al., 2020), TD3+OFE (Ota et al., 2020), and TD7 (Fujimoto et al., 2023). For all algorithms, each task runs 10 instances with different random seeds. In each instance, the evaluation is performed every 5000 time steps. The learning curves are shown in Fig. 6 and the numerical results at 400K time step and 1M time step are shown in Table 1.

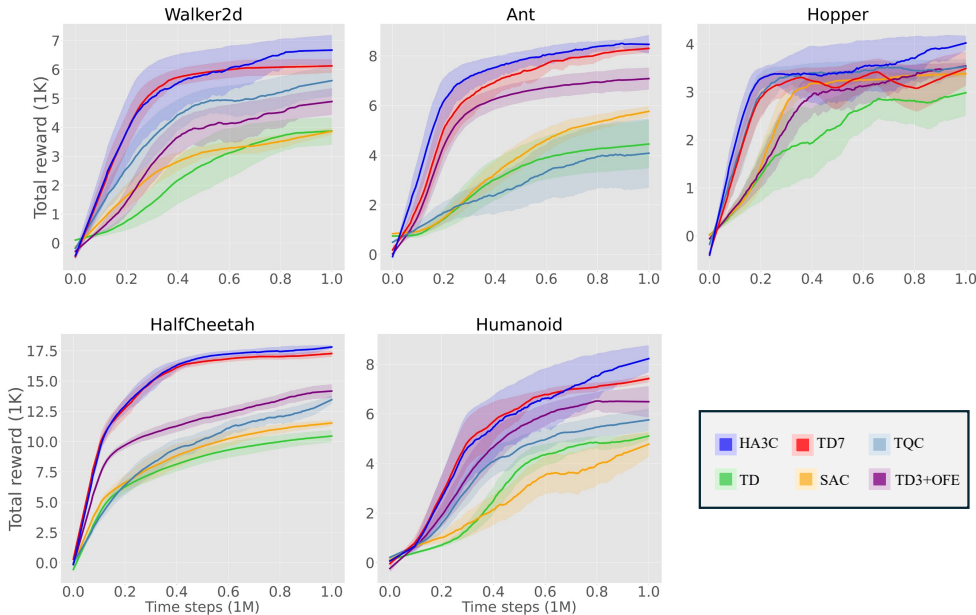


Figure 6: Learning curves of different RL algorithms on the MuJoCo control tasks. The shaded area captures a 90% confidence interval around the average performance.

From Fig. 6 and Table 1, we can see that 1) With the help of historical augmentation, HA3C significantly outperforms the compared algorithms in terms of the early average highest returns (400K time step) and final average highest returns (1M time step); 2) as shown in Fig. 6, because of the instability in rapidly learning complex causal relationships, the early average returns of HA3C on Walker2d and Humanoid are a little lower than the early average returns of TD7, however, HA3C can get the highest final average returns on all of the control tasks.

5.2 ABLATION STUDY

Our ablation study aims to prove that our historical augmentation is the real source of the improvement in sample efficiency. Therefore, we compare HA3C to the following two ablations: 1) Copy Aug. copies the current state k times instead of augmenting with k steps of history in our CNN; 2) No Aug. is TD3 with single-step representation learning and LAP. Our ablation study is performed on Ant, Hopper, and Walker2d. All of the comparison methods have the same parameter setting.

Table 1: The average highest returns over 10 instances on the MuJoCo control tasks at 400K and 1M time steps. \pm captures the standard deviation over trials. The best score is highlighted by cyan and the second best score is highlighted by orange.

Algorithm	Time step	Walker2d	HalfCheetah	Ant	Humanoid	Hopper
TD3	400K	2636 \pm 933	8229 \pm 757	3297 \pm 1084	1384 \pm 282	2876 \pm 859
	1M	4198 \pm 516	10560 \pm 675	4617 \pm 1287	5308 \pm 105	3387 \pm 137
SAC	400K	3122 \pm 156	8945 \pm 1368	3893 \pm 569	2268 \pm 905	3276 \pm 86
	1M	3921 \pm 163	11729 \pm 258	5956 \pm 2209	5498 \pm 131	3422 \pm 87,
TQC	400K	4994 \pm 397	9644 \pm 1006	3307 \pm 939	4061 \pm 703	3534 \pm 91
	1M	5895 \pm 552	13431 \pm 561	5258 \pm 1165	6140 \pm 426	3602 \pm 117
TD3+OFE	400K	4329 \pm 550	11508 \pm 635	6406 \pm 549	5193 \pm 797	3471 \pm 45
	1M	4574 \pm 551	14759 \pm 696	7246 \pm 497	7262 \pm 209	3616 \pm 28
TD7	400K	5787 \pm 444	15625 \pm 559	7305 \pm 197	5823 \pm 231	3440 \pm 92
	1M	6354 \pm 209	17343 \pm 359	8346 \pm 291	7405 \pm 236	3757 \pm 214
HA3C	400K	6441 \pm 366	16652 \pm 323	7838 \pm 138	6099 \pm 305	3783 \pm 153
	1M	7143 \pm 456	18108 \pm 294	8687 \pm 128	8584 \pm 273	4143 \pm 170

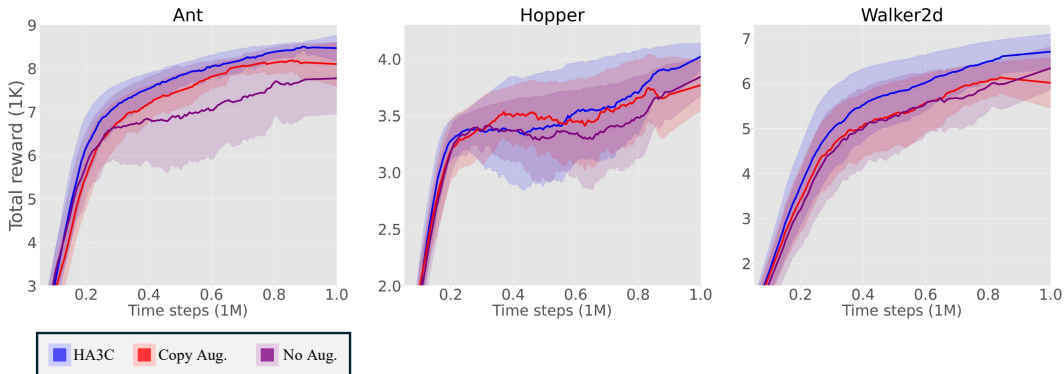


Figure 7: Learning curves of the ablation study on the MuJoCo benchmark. The shaded area captures a 90% confidence interval around the average performance.

As we can see from Fig. 7, HA3C significantly outperforms the compared algorithms in terms of both sample efficiency and performance on Ant and Walker2d. HA3C also significantly outperforms the compared algorithms in final performance on Hopper. This phenomenon illustrates that historical augmentation is the real source for improving sample efficiency.

5.3 PARAMETER SENSITIVITY ANALYSIS

In Fig. 8, we analyze the sensitivities of two important parameters, k and N , on Ant. k is the length of the historical state trajectory and N is the number of dimensions of compressed historical trajectories. Both of the above parameters are not used in the previous representation-based RL algorithms. k is set from $\{6, 12, 18, 24\}$ and N is set from $\{8, 16, 64, 256\}$.

As we can see, HA3C is a little sensitive to k and N . When $k \leq 12$ and $N \leq 16$, our historical augmentation will significantly improve the sample efficiency. When $N = 256$, the historical information cannot improve neither sample efficiency nor final performance. This phenomenon illustrates that compressing the historical trajectories into a low-dimensional space is the key to effectively utilize the historical information in MDP tasks.

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

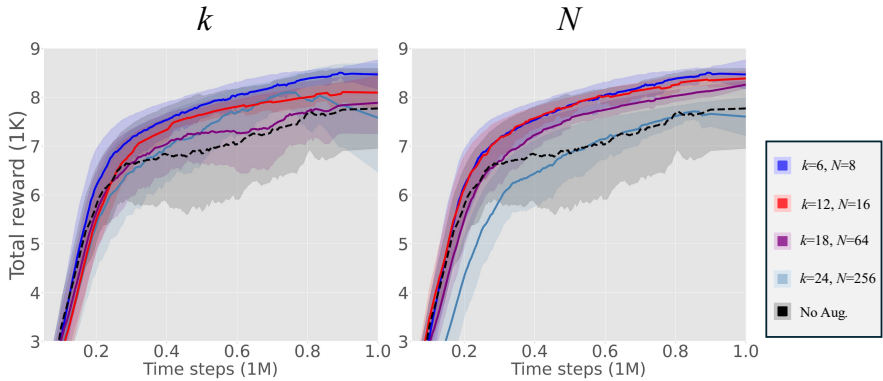


Figure 8: Learning curves of the parameter sensitivity analysis on the MuJoCo benchmark. The shaded area captures a 90% confidence interval around the average performance.

5.4 RUNNING TIME

To understand the computational cost of HA3C, we compare the running times of different algorithms with identical computational resources in HalfCheetah control task. The result is shown in Fig. 9. As we can see, the computational cost of HA3C is less than the computational costs of TD3+OFE and TQC but is more than the computational costs of TD3, SAC, and TD7.

6 CONCLUSION

Under the Markov assumption of MDPs, the probability distributions of the next state and reward depend only on the current state and action. Therefore, given a finite Q -table, we can find the optimal policy in an MDP by a heuristic algorithm which only considers single-step transitions. Different from the heuristic algorithm, DRL algorithms need to approximate the causal functions by learning the causal relationships in MDPs. In this case, DRL is often faced with sample inefficiency from complex causal relationships, as a more complex causal function requires neural networks to approximate with more parameters, samples, and time consumption.

This paper addresses the above problem by augmenting the current state with historical information. We believe that historical augmentation can simplify the causal relationships of state transitions by its inherent contextual information and increasing the search space of the causal functions. Therefore, we focus on optimizing a history-dependent stationary policy in MDPs and propose a new RL algorithm, HA3C. The main idea of HA3C is to learn the state representations by compressing the high-dimensional historical trajectories into a low-dimensional space. In this way, we can extract the simple causal relationships from historical trajectories and avoid the overfitting caused by high-dimensional historical data. Our experiment demonstrates the superior performance of HA3C over five state-of-the-art RL algorithms on MuJoCo control tasks.

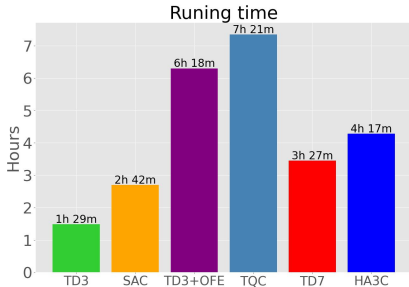


Figure 9: Running times of different algorithms for 1M time steps.

REFERENCES

Cameron Allen, Neev Parikh, Omer Gottesman, and George Konidaris. Learning markov state abstractions for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 34:8229–8241, 2021.

- 540 Ankesh Anand, Evan Racah, Sherjil Ozair, Yoshua Bengio, Marc-Alexandre Côté, and R Devon
541 Hjelm. Unsupervised state representation learning in atari. *Advances in Neural Information*
542 *Processing Systems*, 32, 2019.
- 543 Alexandr Andoni, Rina Panigrahy, Gregory Valiant, and Li Zhang. Learning polynomials with neural
544 networks. In *International Conference on Machine Learning*, pp. 1908–1916. PMLR, 2014.
- 545 David Andre and Stuart J Russell. State abstraction for programmable reinforcement learning agents.
546 In *AAAI*, pp. 119–125, 2002.
- 547 Benny Applebaum, Yuval Ishai, and Eyal Kushilevitz. Cryptography in nc^0 . *SIAM Journal on*
548 *Computing*, 36(4):845–888, 2006.
- 549 Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In
550 *Machine Learning Proceedings 1995*, pp. 30–37. Elsevier, 1995.
- 551 Maurice Stevenson Bartlett. *An introduction to stochastic processes*. University Press Cambridge,
552 1966.
- 553 Monica Bianchini and Franco Scarselli. On the complexity of neural network classifiers: A compari-
554 son between shallow and deep architectures. *IEEE Transactions on Neural Networks and Learning*
555 *Systems*, 25(8):1553–1565, 2014.
- 556 Yann Bouteiller, Simon Ramstedt, Giovanni Beltrame, Christopher Pal, and Jonathan Binas. Rein-
557 forcement learning with random delays. In *International Conference on Learning Representations*,
558 2020.
- 559 Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and
560 Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 561 Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-efficient
562 reinforcement learning with stochastic ensemble value expansion. *Advances in Neural Information*
563 *Processing Systems*, 31, 2018.
- 564 Miles Cranmer, Sam Greycanus, Stephan Hoyer, Peter Battaglia, David Spergel, and Shirley Ho.
565 Lagrangian neural networks. *arXiv preprint arXiv:2003.04630*, 2020.
- 566 Peter Dayan. Improving generalization for temporal difference learning: The successor representation.
567 *Neural Computation*, 5(4):613–624, 1993.
- 568 Esther Derman, Gal Dalal, and Shie Mannor. Acting in delayed environments with non-stationary
569 markov policies. In *International Conference on Learning Representations*, 2020.
- 570 Esther Derman, Gal Dalal, and Shie Mannor. Acting in delayed environments with non-stationary
571 markov policies. *arXiv preprint arXiv:2101.11992*, 2021.
- 572 Richard C Dorf and Robert H Bishop. *Modern control systems*. Pearson, 2011.
- 573 Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for
574 sample efficient reinforcement learning? In *International Conference on Learning Representations*,
575 2020.
- 576 Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep
577 reinforcement learning for continuous control. In *International Conference on Machine Learning*,
578 volume 28, pp. 1329–1338, New York, USA, 2016.
- 579 EB Dynkin. *Markov processes*. Springer, 1965.
- 580 Kevin Esslinger, Robert Platt, and Christopher Amato. Deep transformer q-networks for partially
581 observable reinforcement learning. In *NeurIPS 2022 Foundation Models for Decision Making*
582 *Workshop*, 2022.
- 583 Ben Eysenbach, Xinyang Geng, Sergey Levine, and Russ R Salakhutdinov. Rewriting history with
584 inverse rl: Hindsight inference for policy improvement. *Advances in Neural Information Processing*
585 *Systems*, 33:14783–14795, 2020.

- 594 Scott Fujimoto, Herke Van Hoof, and David Meger. Addressing function approximation error in actor-
595 critic methods. In *International Conference on Machine Learning*, volume 80, pp. 1587–1596,
596 Stockholm, Sweden, 2018.
- 597 Scott Fujimoto, David Meger, and Doina Precup. An equivalence between loss functions and non-
598 uniform sampling in experience replay. *Advances in Neural Information Processing Systems*, 33:
599 14219–14230, 2020.
- 600 Scott Fujimoto, Wei-Di Chang, Edward J Smith, Shixiang Shane Gu, Doina Precup, and David Meger.
601 For SALE: State-action representation learning for deep reinforcement learning. In *Thirty-seventh*
602 *Conference on Neural Information Processing Systems*, 2023.
- 603 Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G Bellemare. Deepmdp:
604 Learning continuous latent space models for representation learning. In *International Conference*
605 *on Machine Learning*, pp. 2170–2179. PMLR, 2019.
- 606 Ruo Cheng Guo, Lu Cheng, Jundong Li, P Richard Hahn, and Huan Liu. A survey of learning causality
607 with data: Problems and methods. *ACM Computing Surveys (CSUR)*, 53(4):1–37, 2020.
- 608 David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- 609 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
610 maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference*
611 *on Machine Learning*, volume 80, pp. 1861–1870, Stockholm, Sweden, 2018.
- 612 Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning
613 behaviors by latent imagination. In *International Conference on Learning Representations*, 2019a.
- 614 Danijar Hafner, Timothy Lillicrap, Ian Fischer, Ruben Villegas, David Ha, Honglak Lee, and James
615 Davidson. Learning latent dynamics for planning from pixels. In *International Conference on*
616 *Machine Learning*, pp. 2555–2565. PMLR, 2019b.
- 617 Assaf Hallak, Dotan Di Castro, and Shie Mannor. Contextual markov decision processes. *arXiv*
618 *preprint arXiv:1502.02259*, 2015.
- 619 Matthew Hausknecht and Peter Stone. Deep recurrent q-learning for partially observable mdps. In
620 *AAAI Fall Symposium Series*, 2015.
- 621 Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger.
622 Deep reinforcement learning that matters. In *Proceedings of the AAAI Conference on Artificial*
623 *Intelligence*, volume 32, 2018.
- 624 Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in Statistics: Methodology*
625 *and Distribution*, pp. 492–518. Springer, 1992.
- 626 Julian Ibarz, Jie Tan, Chelsea Finn, Mrinal Kalakrishnan, Peter Pastor, and Sergey Levine. How to
627 train your robot with deep reinforcement learning: lessons we have learned. *The International*
628 *Journal of Robotics Research*, 40(4-5):698–721, 2021.
- 629 Tommi Jaakkola, Michael Jordan, and Satinder Singh. Convergence of stochastic iterative dynamic
630 programming algorithms. *Advances in Neural Information Processing Systems*, 6, 1993.
- 631 Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. Reinforcement learning: A survey.
632 *Journal of Artificial Intelligence Research*, 4:237–285, 1996.
- 633 Markus Kalisch and Peter Bühlman. Estimating high-dimensional directed acyclic graphs with the
634 pc-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.
- 635 Konstantinos V Katsikopoulos and Sascha E Engelbrecht. Markov decision processes with delays
636 and asynchronous cost collection. *IEEE Transactions on Automatic Control*, 48(4):568–574, 2003.
- 637 DP Kingma. Adam: a method for stochastic optimization. In *International Conference on Learning*
638 *Representations*, 2015.

- 648 Arsenii Kuznetsov, Pavel Shvechikov, Alexander Grishin, and Dmitry Vetrov. Controlling overesti-
649 mation bias with truncated mixture of continuous distributional quantile critics. In *International*
650 *Conference on Machine Learning*, pp. 5556–5566. PMLR, 2020.
- 651 Lihong Li, Thomas J Walsh, and Michael L Littman. Towards a unified theory of state abstraction for
652 mdps. In *AI&M*, 2006.
- 654 Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
655 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In
656 *International Conference on Learning Representations*, San Juan, Puerto Rico, 2016.
- 657 Guoqing Liu, Chuheng Zhang, Li Zhao, Tao Qin, Jinhua Zhu, Li Jian, Nenghai Yu, and Tie-Yan
658 Liu. Return-based contrastive representation learning for reinforcement learning. In *International*
659 *Conference on Learning Representations*, 2020.
- 660 Sultan Javed Majeed and Marcus Hutter. On q-learning convergence for non-markov decision
661 processes. In *IJCAI*, volume 18, pp. 2546–2552, 2018.
- 663 Francisco S Melo. Convergence of q-learning: A simple proof. *Institute Of Systems and Robotics,*
664 *Tech. Rep.*, pp. 1–4, 2001.
- 665 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
666 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
667 through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- 668 Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim
669 Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement
670 learning. In *International Conference on Machine Learning*, pp. 1928–1937. PMLR, 2016.
- 671 Alfred Müller. Integral probability metrics and their generating classes of functions. *Advances in*
672 *applied probability*, 29(2):429–443, 1997.
- 673 Jelle Munk, Jens Kober, and Robert Babuška. Learning state representation for deep actor-critic
674 control. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pp. 4667–4673. IEEE,
675 2016.
- 676 Chengzhuo Ni, Yaqi Duan, Munther Dahleh, Mengdi Wang, and Anru R Zhang. Learning good
677 state and action representations for markov decision process via tensor decomposition. *Journal of*
678 *Machine Learning Research*, 24(115):1–53, 2023.
- 681 Kei Ota, Tomoaki Oiki, Devesh Jha, Toshisada Mariyama, and Daniel Nikovski. Can increasing input
682 dimensionality improve deep reinforcement learning? In *International Conference on Machine*
683 *Learning*, pp. 7424–7433. PMLR, 2020.
- 684 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
685 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style,
686 high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32,
687 2019.
- 688 Gandharv Patil, Aditya Mahajan, and Doina Precup. On learning history-based policies for controlling
689 markov decision processes. In *International Conference on Artificial Intelligence and Statistics*,
690 pp. 3511–3519. PMLR, 2024.
- 691 Martin L Puterman. Markov decision processes. *Handbooks in operations research and management*
692 *science*, 2:331–434, 1990.
- 693 Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John
694 Wiley & Sons, 2014.
- 695 Balaraman Ravindran. *An algebraic approach to abstraction in reinforcement learning*. University
696 of Massachusetts Amherst, 2004.
- 697 Sahand Rezaei-Shoshtari, Rosie Zhao, Prakash Panangaden, David Meger, and Doina Precup. Con-
698 tinuous mdp homomorphisms and homomorphic policy gradient. *Advances in Neural Information*
699 *Processing Systems*, 35:20189–20204, 2022.
- 700
701

- 702 Ahmad EL Sallab, Mohammed Abdou, Etienne Perot, and Senthil Yogamani. Deep reinforcement
703 learning framework for autonomous driving. *Electronic Imaging*, 2017(19):70–76, 2017.
704
- 705 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
706 policy optimization. In *International Conference on Machine Learning*, volume 37, pp. 1889–1897,
707 Lille, France, 2015.
- 708 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
709 Deterministic policy gradient algorithms. In *International Conference on Machine Learning*,
710 volume 32, pp. 1387–1395, Beijing, China, 2014.
711
- 712 Shagun Sodhani, Franziska Meier, Joelle Pineau, and Amy Zhang. Block contextual mdps for
713 continual learning. In *Learning for Dynamics and Control Conference*, pp. 608–623. PMLR, 2022.
- 714 David Sprunger and Bart Jacobs. The differential calculus of causal functions. *arXiv preprint*
715 *arXiv:1904.10611*, 2019.
716
- 717 Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- 718 Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden,
719 Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint*
720 *arXiv:1801.00690*, 2018.
721
- 722 Yee Teh, Victor Bapst, Wojciech M Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas
723 Heess, and Razvan Pascanu. Distral: Robust multitask reinforcement learning. *Advances in Neural*
724 *Information Processing Systems*, 30, 2017.
- 725 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
726 In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033,
727 Algarve, Portugal, 2012. IEEE.
- 728 Herke Van Hoof, Nutan Chen, Maximilian Karl, Patrick van der Smagt, and Jan Peters. Stable rein-
729 forcement learning with autoencoders for tactile and visual data. In *2016 IEEE/RSJ International*
730 *Conference on Intelligent Robots and Systems (IROS)*, pp. 3928–3934. IEEE, 2016.
731
- 732 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
733 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing*
734 *Systems*, 30, 2017.
- 735 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
736 learning. *Machine Learning*, 8(3-4):229–256, 1992.
737
- 738 Mingxuan Ye, Yufei Kuang, Jie Wang, Rui Yang, Wengang Zhou, Houqiang Li, and Feng Wu.
739 State sequences prediction via fourier transform for representation learning. In *Thirty-seventh*
740 *Conference on Neural Information Processing Systems*, 2023.
- 741 Xue Ying. An overview of overfitting and its solutions. In *Journal of Physics: Conference Series*,
742 volume 1168, pp. 022022. IOP Publishing, 2019.
743
- 744 Amy Zhang, Rowan McAllister, Roberto Calandra, Yarín Gal, and Sergey Levine. Learning invariant
745 representations for reinforcement learning without reconstruction. In *International Conference on*
746 *Learning Representations*, 2021.
747
748
749
750
751
752
753
754
755

A DIFFERENT POLICIES

Time-related policies can be History-dependent (H) or k -order Markov (M_k) (Derman et al., 2020; Puterman, 2014). Denote \mathcal{H}_t as the set of possible histories up to time step t . A history-dependent policy $\pi = \{d_t | t = 0, 1, \dots\}$ at t maps histories to actions as $d_t : \mathcal{H}_t \mapsto \mathcal{A}$. A k -order Markov policy $\pi = \{d_t | t = 0, 1, \dots\}$ at t maps k -order state transition trajectories to actions as $d_t : \mathcal{S}_{k,t} \mapsto \mathcal{A}$. A k -order stationary (S_k) policy is unrelated to time as $\pi : \mathcal{S}_{k,:} \mapsto \mathcal{A}$. In general, a randomized (R) policy selects the actions by a probability distribution as $\pi(\mathbf{a} | *)$. π is a deterministic (D) policy if and only if $\pi(\mathbf{a} | *) \in \{0, 1\}$. Based on the above analysis, we can obtain History-dependent Random (HR) policies, History-dependent Deterministic (HD) policies, k -order Markov Random (M_kR) policies, k -order Markov Deterministic (M_kD) policies, k -order Stationary Random (S_kR) policies, and k -order Stationary Deterministic (S_kD) policies.

The above policies are summarized in Table 2. The relationships among them are demonstrated in Fig. 10. It is noteworthy that sometimes historical actions will be considered in decision-making. In this case, without loss of generality, a historical state $\mathbf{s}_{i|i \leq t-1}$ can be updated by $\mathbf{s}_i \leftarrow \mathbf{s}_i \cup \mathbf{a}_i$.

Table 2: Different types of policies.

Policy	Abbreviation	Action
History-dependent Random	HR	$\mathbf{a}_t \sim d_t(\mathbf{s}_{0,t}), d_t \in \pi$
History-dependent Deterministic	HD	$\mathbf{a}_t = d_t(\mathbf{s}_{0,t}), d_t \in \pi$
k -order Markov Random	M_kR	$\mathbf{a}_t \sim d_t(\mathbf{s}_{k-t+1,t}), d_t \in \pi$
k -order Markov Deterministic	M_kD	$\mathbf{a}_t = d_t(\mathbf{s}_{k-t+1,t}), d_t \in \pi$
k -order Stationary Random	S_kR	$\mathbf{a}_t \sim \pi(\mathbf{s}_{k-t+1,t})$
k -order Stationary Deterministic	S_kD	$\mathbf{a}_t = \pi(\mathbf{s}_{k-t+1,t})$

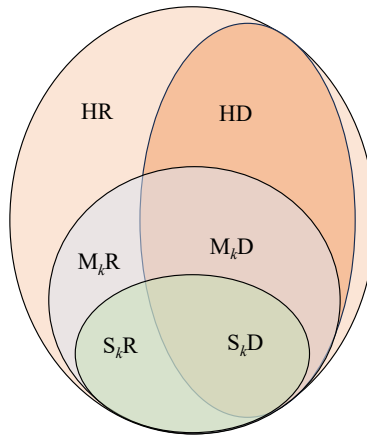


Figure 10: The relations among different policies.

B AN EXAMPLE OF IMPROVING SAMPLE EFFICIENCY IN MDPs BY HISTORICAL AUGMENTATION

Define a sequence as follows: 1) $|\beta_0| \neq 1$; 2) If $i > 1$, then $\beta_{i+1} = \beta_i^2$.

Based on the sequence above, we can define an MDP $\mathbb{M} = \langle \mathcal{S}, \mathcal{A}, R, \mathbf{P}, \gamma \rangle$. At time step t , state $\mathbf{s}_t = [\beta_t, \beta_{t+2}]^\top$ and action \mathbf{a}_t is computed by a linear function $f(\cdot)$ on state \mathbf{s}_t or augmented state $\mathbf{s}_{k,t}$. Without considering historical information, reward r_t is defined as

$$r_t = -|f(\mathbf{s}_t) - (\beta_t + \sqrt{\beta_{t+2} + \beta_{t+2}})| = -|\mathbf{w}\mathbf{s}_t + b - (\beta_t + \sqrt{\beta_{t+2} + \beta_{t+2}})|, \quad (6)$$

where \mathbf{w} is a two-dimensional vector and b is a constant. In transition model \mathbf{P} , \mathbf{s}_0 can be defined as $[\beta_0, \beta_2]^\top$ and \mathbf{s}_{t+1} can be computed by \mathbf{s}_t as

$$\mathbf{s}_{t+1} = [\beta_t^2, \beta_{t+2}^2]^\top = \mathbf{s}_t \odot \mathbf{s}_t, \quad (7)$$

where \odot is Hadamard product. $\gamma = 0.99$.

From equation 6 and equation 7, it is easy to see that \mathbb{M} satisfies the Markov assumption of MDPs. To maximize the discounted cumulative reward in \mathbb{M} , we should minimize

$$\arg \min_{\mathbf{w}, b} \|f(\mathbf{s}_t) - (\beta_t + \sqrt{\beta_{t+2} + \beta_{t+2}})\|_2 = \arg \min_{\mathbf{w}, b} \|\mathbf{w}\mathbf{s}_t + b - (\beta_t^2 + \sqrt{\beta_{t+2} + \beta_{t+2}})\|_2 \quad (8)$$

at each time step t . However, it is hard to minimize equation 8 by $f(\mathbf{s}_t)$, which is a linear model on \mathbf{s}_t .

The above problem can be solved by the historical augmentation of \mathbf{s}_t . When considering the historical augmentation of \mathbf{s}_t , $f(\cdot)$ on $\mathbf{s}_{2,t}$ can be defined as

$$f(\mathbf{s}_{2,t}) = \mathbf{w}_0\mathbf{s}_t + \mathbf{w}_1\mathbf{s}_{t-1} + b.$$

Instead of minimizing equation 8, we can minimize

$$\begin{aligned} & \arg \min_{\mathbf{w}_0, \mathbf{w}_1, b} \|f(\mathbf{s}_{2,t}) - (\beta_t + \sqrt{\beta_{t+2} + \beta_{t+2}})\|_2 \\ &= \arg \min_{\mathbf{w}_0, \mathbf{w}_1, b} \|\mathbf{w}_0\mathbf{s}_t + \mathbf{w}_1\mathbf{s}_{t-1} + b - (\beta_t^2 + \sqrt{\beta_{t+2} + \beta_{t+2}})\|_2. \end{aligned}$$

Let $\mathbf{w}_0 = [1, 1]$, $\mathbf{w}_1 = [0, 1]$, and $b = 0$. From $\beta_{t+1} = \sqrt{\beta_{t+2}}$, we have

$$\begin{aligned} & \|\mathbf{w}_0\mathbf{s}_t + \mathbf{w}_1\mathbf{s}_{t-1} + b - (\beta_t + \sqrt{\beta_{t+2} + \beta_{t+2}})\|_2 \\ &= \|\mathbf{w}_0\mathbf{s}_t + \mathbf{w}_1\mathbf{s}_{t-1} + b - (\beta_t + \beta_{t+1} + \beta_{t+2})\|_2 \\ &= \|([1, 1][\beta_t, \beta_{t+2}]^\top + [0, 1][\beta_{t-1}, \beta_{t+1}]^\top - (\beta_t + \beta_{t+1} + \beta_{t+2}))\|_2 \\ &= 0. \end{aligned}$$

In this case, the cumulative reward in \mathbb{M} can be maximized.

C CONNECTED TO RELATED WORK

C.1 CONNECTED TO HMDPs

In HMDPs, the probability distributions of the reward and next state depend not only on the current state and action but also on the historical states and actions. For a k -order HMDPs, we have

$$P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | \mathbf{s}_0, \mathbf{a}_0, r_0, \dots, \mathbf{s}_t, \mathbf{a}_t\} = P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | \mathbf{s}_{t-k+1}, \mathbf{a}_{t-k+1}, \dots, \mathbf{s}_t, \mathbf{a}_t\}.$$

The causal diagram of HMDP is presented in Fig. 11(a). Our approach optimizes the policy by a simplified HMDP model in which the probability distributions of the reward and next state depend on the current state-action pair and compressed historical trajectory as

$$P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | \mathbf{s}_0, \mathbf{a}_0, r_0, \dots, \mathbf{s}_t, \mathbf{a}_t\} = P\{\mathbf{s}_{t+1} = \mathbf{s}', r_t = r | DR(\mathbf{s}_{t-1, k-1}), \dots, \mathbf{s}_t, \mathbf{a}_t\}.$$

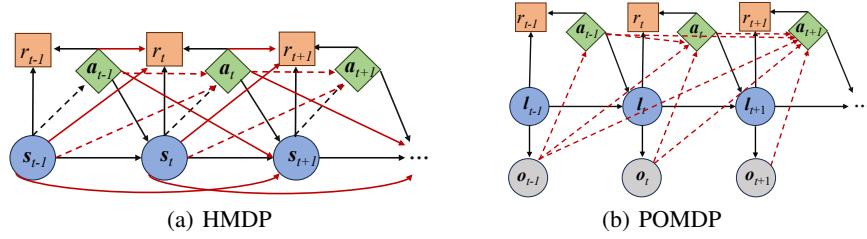


Figure 11: Causal diagram of HMDPs and POMDPs.

C.2 CONNECTED TO POMDPs

In POMDPs, the states are partially observable. Define the partially observable state at time step t as l_t and the observable part of l_t as o_t . The causal diagram of POMDPs is shown in Fig. 11(b). Under the faithfulness assumption, o_t and o_{t+k} are mutually dependent conditional on $\forall k > 1$, $\{o_i, a_i\}_{t < i < t+k}$ (Kalisch & Bühlman, 2007). Therefore, in a POMDP, the optimal policy π at time step t should consider not only o_t but also the historical information $\{o_i, a_i\}_{0 \leq i < t}$. When k is large, long-length rollout estimation is needed in POMDPs.

RL algorithms of world models, such as Dream, model the sequential decision-making as a POMDP (Ha & Schmidhuber, 2018; Hafner et al., 2019a). They usually encode the historical information at t by an encoder f^t to construct s_{t+1} as

$$s_{t+1} = f^t(o_t, a_t, \dots, f^1(o_1, a_1, f^0(o_0, a_0))).$$

When t is large, some partially observable states will be encoded many times, leading to the loss of some important discriminative information.

Compared with POMDP-based RL algorithms, our HA3C can better adjust the considered steps in history according to the actual task and thus effectively find the optimal policy in history-based sequential decision-making.

C.3 CONNECTED TO STATE ABSTRACTION

State abstraction aims to reduce ground MDPs with large state spaces to abstract MDPs with smaller state spaces by aggregating states according to some notion of equality or similarity (Bartlett, 1966). Through abstraction, intelligent agents may need to consider only the salient distinguishing information of their environments. Given an abstraction function as $F: \mathcal{S} \rightarrow \bar{\mathcal{S}}$, we can define the abstract version of MDP \mathbb{M} as $\bar{\mathbb{M}} = \langle \bar{\mathcal{S}}, \mathcal{A}, \bar{R}, \bar{P}, \gamma \rangle$. A Q -irrelevance abstraction function F^Q is that for any action a , $F^Q(s) = F^Q(s')$ implies $Q(s, a) = Q(s', a)$. Then we have the following theorem.

Theorem C.1. Define an MDP as $\mathbb{M}_k = \langle \mathcal{S}_{k,:}, \mathcal{A}, R, \mathbf{P}_k, \gamma \rangle$. Under the conditions 1), 2), and 3) in Theorem 4.2, encoder f is a Q -irrelevance abstraction on $s_{k,:}$.

Theorem C.1 illustrates that our representation learning can be seen as the Q -irrelevance abstraction of the historically augmented states. The proof of this theorem is given in Appendix D.

D THEORETICAL ANALYSIS

D.1 PROOF OF THEOREM 4.1

Now we give the proof to Theorem 4.1. The different types of policies in this proof are summarized in Table 2. The relationships between these policies are shown in Fig. 10.

Based on the Markov assumption of MDPs, we have

$$\begin{aligned} & P\{s_{t+1} = s', r_t = r | s_0, a_0, r_0, \dots, s_t, a_t\} \\ &= P\{s_{t+1} = s', r_t = r | s_{t-k+1,t}, a_t\} \\ &= P\{s_{t+1} = s', r_t = r | s_t, a_t\}. \end{aligned} \quad (9)$$

For any $\pi \in HR$, we can define $V_\pi(\mathbf{h}_t)$ by

$$V^\pi(\mathbf{h}_t) = \mathbb{E}^\pi \left[\sum_{i=t}^{+\infty} \gamma^i R(\mathbf{h}_{t+i}, \mathbf{a}_{t+i}) \right].$$

From Fig. 10, we have $S_k D \in M_k D \in M_k R \in HR$. In view of equation 9, we see for all t that

$$\sup_{\pi \in HR} V^\pi(\mathbf{h}_t) = \sup_{\pi \in S_k D} V^\pi(\mathbf{s}_{k,t}).$$

First, for all t , we demonstrate that

$$\sup_{\pi \in HR} V^\pi(\mathbf{h}_t) = \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}). \quad (10)$$

This is a direct result of Theorem D.1. The proof of this theorem is presented in D.1.1.

Theorem D.1. *Let $\pi = \{d_t | t = 0, 1, \dots\} \in HR$. Then $\forall \mathbf{s}_{k,:} \in S_{k,:}$, based on equation 9, there exists a policy $\pi' = \{d'_t | t = 0, 1, \dots\} \in M_k R$ satisfying*

$$p^\pi(\mathbf{a}_{t+i} = \mathbf{a}', \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) = p^{\pi'}(\mathbf{a}_{t+i} = \mathbf{a}', \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}),$$

where $p^\pi(*)$ denotes the probability of $*$ provided that the agent follows policy π .

Then Theorem D.2 illustrates that the value functions of $\pi \in M_k D$ and $\pi \in M_k R$ have the same upper bound. The proof of this theorem is demonstrated in D.1.2.

Theorem D.2. *If a bounded function V on $S_{k,:}$ satisfies the optimal Bellman equation that*

$$V(\mathbf{s}_{k,t}) = \sup_{\mathbf{a} \in \mathcal{A}} \left\{ R(\mathbf{s}_{k,t}, \mathbf{a}) + \gamma \int_{S_{k,:}} V(\mathbf{s}_{k,t+1} | \mathbf{s}_{t+1} = \mathbf{s}') p(\mathbf{s}' | \mathbf{s}_{k,t}, \mathbf{a}) d\mathbf{s}'_{k,:} \right\},$$

then

$$\sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}) = \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}).$$

Finally, based on equation 9, for all $\mathbf{s}_{k,:} \in S_{k,:}$, if $\mathbf{s}_{k,t} = \mathbf{s}_{k,:}$, then

$$\sup_{\mathbf{a} \in \mathcal{A}} V(\mathbf{s}_{k,t}) = \sup_{\mathbf{a} \in \mathcal{A}} V(\mathbf{s}_{k,:}). \quad (11)$$

Let $\mathbf{a} = \pi(\mathbf{s}_{k,:})$, where $\pi \in S_k D$. It follows that

$$\sup_{\pi \in S_k D} V^\pi(\mathbf{s}_{k,:}) = \sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}). \quad (12)$$

Under equation 10, equation 11 and equation 12, $\forall t$, if $\mathbf{s}_{k,t} = \mathbf{s}_{k,:}$, then

$$\sup_{\pi \in HR} V^\pi(\mathbf{h}_t) = \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}) = \sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}) = \sup_{\pi \in S_k D} V^\pi(\mathbf{s}_{k,:}).$$

D.1.1 PROOF OF THEOREM D.1

We assume that Theorem D.1 holds for $i = 1, 2, 3, \dots, n-1$. Given a policy $\pi \in HR$, based on equation 9, we see that there exists a policy $\pi' \in M_k R$ satisfying

$$\begin{aligned} & p^\pi(\mathbf{s}_{k,t+i} = \mathbf{s}''_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) \\ &= \int_{S_{k,:}} \int_{\mathcal{A}} p^\pi(\mathbf{s}_{k,t+i-1} = \mathbf{s}'_{k,:}, \mathbf{a}_{t+i-1} = \mathbf{a}' | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) p(\mathbf{s}'' | \mathbf{s}'_{k,:}, \mathbf{a}') d\mathbf{a}' d\mathbf{s}'_{k,:} \\ &= \int_{S_{k,:}} \int_{\mathcal{A}} p^{\pi'}(\mathbf{s}_{k,t+i-1} = \mathbf{s}'_{k,:}, \mathbf{a}_{t+i-1} = \mathbf{a}' | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) p(\mathbf{s}'' | \mathbf{s}'_{k,:}, \mathbf{a}') d\mathbf{a}' d\mathbf{s}'_{k,:} \\ &= p^{\pi'}(\mathbf{s}_{k,t+i} = \mathbf{s}''_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}). \end{aligned}$$

The above equality follows from the induction hypothesis. The π' also can satisfy

$$p^{\pi'}(\mathbf{a}_{t+i} = \mathbf{a}' | \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:}) = p^\pi(\mathbf{a}_{t+i} = \mathbf{a}' | \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:}).$$

Therefore,

$$\begin{aligned} & p^{\pi'}(\mathbf{a}_{t+i} = \mathbf{a}', \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) \\ &= p^{\pi'}(\mathbf{a}_{t+i} = \mathbf{a}' | \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:}) p^{\pi'}(\mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) \\ &= p^\pi(\mathbf{a}_{t+i} = \mathbf{a}' | \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:}) p^\pi(\mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}) \\ &= p^\pi(\mathbf{a}_{t+i} = \mathbf{a}', \mathbf{s}_{k,t+i} = \mathbf{s}'_{k,:} | \mathbf{s}_{k,t} = \mathbf{s}_{k,:}). \end{aligned}$$

972 D.1.2 PROOF OF THEOREM D.2

973 In view of $M_k D \in M_k R$, we have

$$974 \sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}) \leq \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}). \quad (13)$$

975 It follows that

$$976 \sup_{\mathbf{a} \in \mathcal{A}} \left\{ R(\mathbf{s}_{k,t}, \mathbf{a}) + \gamma \int_{\mathcal{S}_{k,:}} V(\mathbf{s}_{k,t+1} | \mathbf{s}_{k,t+1} = \mathbf{s}') p(\mathbf{s}' | \mathbf{s}_{k,t}, \mathbf{a}) d\mathbf{s}' \right\}$$

$$977 \geq \int_{\mathcal{A}} p(d_t(\mathbf{s}_{k,t}) = \mathbf{a}) \left[R(\mathbf{s}_{k,t}, \mathbf{a}) + \gamma \int_{\mathcal{S}_{k,:}} V(\mathbf{s}_{k,t+1} | \mathbf{s}_{k,t+1} = \mathbf{s}') p(\mathbf{s}' | \mathbf{s}_{k,t}, \mathbf{a}) d\mathbf{s}'_{k,:} \right] d\mathbf{a},$$

978 where $d_t \in M_k R$. Thus

$$979 \sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}) \geq \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}). \quad (14)$$

980 Combining equation 13 and equation 14, we have

$$981 \sup_{\pi \in M_k D} V^\pi(\mathbf{s}_{k,t}) = \sup_{\pi \in M_k R} V^\pi(\mathbf{s}_{k,t}).$$

982 D.2 PROOF OF THEOREM 4.2

983 To prove Theorem 4.2, we give the proof of Theorem C.1 first. Under the condition 1) of Theorem 4.2,
984 one sees that there are only two independent variables $\mathbf{s}_{k,:}$ and \mathbf{a} . Under the Markov assumption and
985 the condition 2) of Theorem 4.2, we have

$$986 P\{\mathbf{s}_{k,t+1} = \mathbf{s}'_{k,:} | \mathbf{s}_0, \mathbf{a}_0, r_0, \dots, \mathbf{s}_t, \mathbf{a}_t\} = P\{\mathbf{s}_{k,t+1} = \mathbf{s}'_{k,:} | \mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t}\}. \quad (15)$$

987 Then, under the condition 3) of Theorem 4.2, we have

$$988 P\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{s}_{k,t}, \mathbf{a}_t\} \doteq P\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t}\} \quad (16)$$

$$989 = P\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | g(f(\mathbf{s}_{k,t}), \mathbf{a}_t)\}$$

$$990 = P\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | g(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)\}$$

$$991 = P\{\mathbf{z}^{\mathbf{s}_{k,t+1}} = \mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,t}, \mathbf{a}_t}\}.$$

1000 Define an MDP as $\mathbb{M}_k = \langle \mathcal{S}_{k,:}, \mathcal{A}, R, \mathbf{P}_k, \gamma \rangle$. From equation 15 and equation 16, we obtain

$$1001 \mathbf{z}^{\mathbf{s}_{k,:}} = \mathbf{z}^{\mathbf{s}'_{k,:}} \rightarrow Q(\mathbf{s}_{k,:}, \mathbf{a}) = Q(\mathbf{s}'_{k,:}, \mathbf{a})$$

1002 Because $\mathbf{z}^{\mathbf{s}_{k,:}} = f(\mathbf{s}_{k,:})$, we see that encoder f is a Q -irrelevance abstraction on $\mathbf{s}_{k,:}$.

1003 Define an abstracted MDP of \mathbb{M}_k as $\overline{\mathbb{M}}_k = \langle \mathcal{Z}, \mathcal{A}, R, \mathbf{P}_k, \gamma \rangle$, where \mathcal{Z} is the encoded space of $\mathcal{S}_{k,:}$.
1004 Operator B_μ can be written as

$$1005 B_\mu \hat{Q}(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a}) = R(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a}) + \max_{\mu} \gamma \int_{\mathcal{Z}} \hat{Q}(\mathbf{z}^{\mathbf{s}_{k,:}}, \mu(\mathbf{z}^{\mathbf{s}_{k,:}})) p(\mathbf{z}^{\mathbf{s}'_{k,:}} | \mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a}) d\mathbf{z}^{\mathbf{s}'_{k,:}}.$$

1006 Now we provide a proof (sketch) to Theorem 4.2. Since the optimality of μ follows from the optimal
1007 actions as well as their Q -values are preserved after abstraction, we see that B is a contraction in the
1008 sup-norm and the optimal Q -function \hat{Q}^* is the unique fixed point of B . Thus we can finally find the
1009 optimal policy μ^* by B_μ (Melo, 2001). When the agent estimates the optimal Q -function based on
1010 experience, we have the following update rule in each time step T by Lemma D.3 (Jaakkola et al.,
1011 1993; Melo, 2001).

$$1012 \hat{Q}_{t+1}(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t) = \hat{Q}_t(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t) + \alpha_t (r_t + \gamma \max_{\mu} \hat{Q}_t(\mathbf{z}^{\mathbf{s}_{k,t+1}}, \mu(\mathbf{z}^{\mathbf{s}_{k,t+1}})) - \hat{Q}_t(\mathbf{z}^{\mathbf{s}_{k,t}}, \mathbf{a}_t)).$$

1013 \hat{Q}_t converges to Q^* as long as

$$1014 \sum_{t=0}^{\infty} \alpha_t = \infty, \quad \sum_{t=0}^{\infty} \alpha_t^2 < \infty.$$

Lemma D.3. *The random process $\{\Delta_t\}$ taking values in \mathbb{R}^n and defined as*

$$\Delta_{t+1}(\mathbf{y}) = (1 - \alpha_t)\Delta_t(\mathbf{y}) + \alpha_t F_t(\mathbf{y})$$

converges to zero under the following assumptions:

1) $\sum_{t=0}^{\infty} \alpha_t = \infty$ and $\sum_{t=0}^{\infty} \alpha_t^2 < \infty$,

2) $\mathbb{E}[\|F_t(\mathbf{y})\|_{\mathcal{F}_t}|_w] \leq \gamma\|\Delta_t\|_w$ with $\gamma < 1$, and

3) $\text{Var}[F_t(\mathbf{y})|\mathcal{F}_t] \leq C(1 + \|\Delta_t\|_w^2)$ for $C > 0$,

where $\mathcal{F} = \{\Delta_t, \Delta_{t-1}, \dots, F_{t-1}, \dots, \alpha_{t-1}, \dots\}$ strands for the past at step t and $\|\cdot\|_w$ refers to some weighted maximum norm.

D.3 APPROXIMATION ERROR ANALYSIS

Define the value function in \mathcal{Z} as \hat{V} . The bound of the approximation error between the transition probabilities in space $\mathcal{S}_{k,:}$, and \mathcal{Z} based on the optimal value function \hat{V}^* can be defined as (Müller, 1997)

$$\max_{\mathbf{s}_{k,:}, \mathbf{a}} \left| \int_{\mathcal{S}_{k,:}} \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{s}'_{k,:}|\mathbf{s}_{k,:}, \mathbf{a})d\mathbf{s}'_{k,:} - \int_{\mathcal{Z}} \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{z}'_{k,:}|\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})d\mathbf{z}'_{k,:} \right| = \delta.$$

Based on δ , we analyze the approximation error in Theorem D.4.

Theorem D.4. *The worst-case difference between $V^\mu(\mathbf{z}^{\mathbf{s}_{k,:}})$ and optimal value function $V^*(\mathbf{s})$ is bounded as:*

$$\|V^*(\mathbf{s}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty \leq \frac{\gamma\delta}{1-\gamma}.$$

We provide the proof to the above theorem as follows. Based on the Markov assumption of MDPs, we have

$$\|V^*(\mathbf{s}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty = \|V^*(\mathbf{s}_{k,:}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty.$$

Now we prove that

$$\|V^*(\mathbf{s}_{k,:}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty \leq \frac{\gamma\delta}{1-\gamma}. \quad (17)$$

In view of $R(\mathbf{s}, \mathbf{a}) = R(\mathbf{s}_{k,:}, \mathbf{a}) = R(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})$ in the value function approximation, we have

$$\begin{aligned} & \|V^*(\mathbf{s}_{k,:}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty \\ & \leq \max_{\mathbf{s}_{k,:}, \mathbf{a}} \|Q^*(\mathbf{s}_{k,:}, \mathbf{a}) - \hat{Q}^*(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})\| \\ & = \max_{\mathbf{s}_{k,:}, \mathbf{a}} \left| R(\mathbf{s}_{k,:}, \mathbf{a}) + \gamma \int_{\mathcal{S}_{k,:}} V^*(\mathbf{s}'_{k,:})p(\mathbf{s}'_{k,:}|\mathbf{s}_{k,:}, \mathbf{a})d\mathbf{s}'_{k,:} \right. \\ & \quad \left. - R(\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a}) - \gamma \int_{\mathcal{Z}} \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{z}'_{k,:}|\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})d\mathbf{z}'_{k,:} \right| \\ & \leq \gamma \max_{\mathbf{s}_{k,:}, \mathbf{a}} \left| \int_{\mathcal{S}_{k,:}} V^*(\mathbf{s}'_{k,:})p(\mathbf{s}'_{k,:}|\mathbf{s}_{k,:}, \mathbf{a})d\mathbf{s}'_{k,:} - \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{s}'_{k,:}|\mathbf{s}_{k,:}, \mathbf{a})d\mathbf{s}'_{k,:} \right| \\ & \quad + \gamma \max_{\mathbf{s}_{k,:}, \mathbf{a}} \left| \int_{\mathcal{S}_{k,:}} \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{s}'_{k,:}|\mathbf{s}_{k,:}, \mathbf{a})d\mathbf{s}'_{k,:} - \int_{\mathcal{Z}} \hat{V}^*(\mathbf{z}'_{k,:})p(\mathbf{z}'_{k,:}|\mathbf{z}^{\mathbf{s}_{k,:}}, \mathbf{a})d\mathbf{z}'_{k,:} \right| \\ & \leq \gamma \left(\|V^*(\mathbf{s}_{k,:}) - \hat{V}^*(\mathbf{z}^{\mathbf{s}_{k,:}})\|_\infty + \delta \right). \end{aligned}$$

This proves equation 17. Thus Theorem D.4 holds.

D.4 ANALYZING SAMPLE EFFICIENCY IN EXPLORATION AND EXPLOITATION

In this subsection, we illustrate the benefit of sample efficiency from history augmentation based on two facts:

- 1) Historical augmentation can improve exploration in DRL. The policy can generate different actions for different transition trajectories that end with the same state;
- 2) Historical augmentation can also improve exploitation in DRL. History augmentation may simplify the causal relationships between the state and the explored high-reward action, thus the policy network can effectively learn and then regenerate this action.

The detailed analysis of these two facts is as follows. In the previous DRL methods for MDPs, when the policy μ and $s_t = s$ are fixed, we can get only one action by

$$a_t = \mu(s_t), \quad \mu \in S_1 D.$$

However, based on our history-based policy

$$a_t = \mu(s_{k,t}), \quad \mu \in S_k D|_{k \geq 2}.$$

a_t can be changed by the change of the $s_{k-1,t-1}$. We define the set of possible actions from policy $\mu \in S_k D$ at state s as A_μ^s and the set of possible k -order trajectories end with state s as S_k^s . As we can see, $|A_\mu^s| \leq |S_k^s|$.

Fig. 12 is the causal diagram of regenerating a high-reward action with or without historical augmentation. For a policy network $\mu_\theta \in S_1 D$ and $a = \mu_\theta(s)$, we may get $a^* = a + \epsilon$ with $R(s, a^*) > R(s, a)$. However, it may be hard to regenerate a^* by the policy network $\mu_\theta(s)$ because the noise ϵ is independent of parameter θ . Fortunately, the causal relationship between $s_{k,t}|_{k \geq 2}$ and a^* may be simpler than the causal relationship between s_t and a^* (See the example in Appendix B). In this case, we can effectively learn the policy $\mu_\theta \in S_k D$ to regenerate the a^* at state s by $a^* = \mu_\theta(s_{k,t})$ (See the example in Fig. 5).

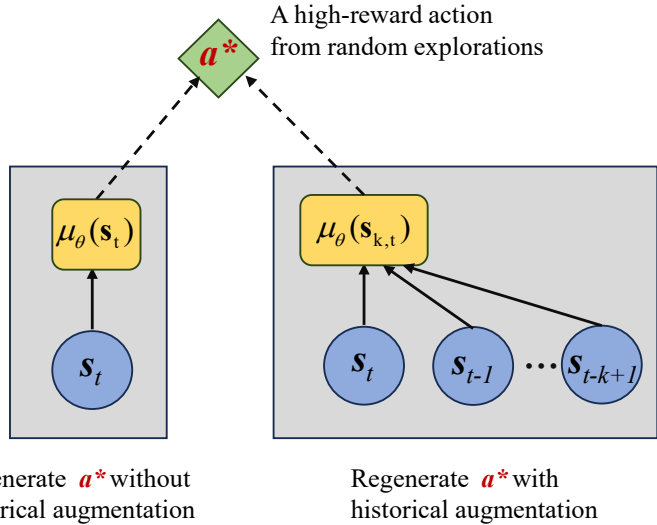


Figure 12: The causal diagram of regenerating a high-reward action with or without historical augmentation. The dashed lines indicate the information needed in the optimization.

E EXPERIMENTAL SETTING

All experiments are run on a single Nvidia 3090 GPU and AMD 5900X CPU. We use the following software versions:

- Python 3.9.12
- Pytorch 2.0.0 (Paszke et al., 2019)
- CUDA 12.2
- Gymnasium 0.29.1 (Brockman et al., 2016)
- MuJoCo 3.2.3 (Todorov et al., 2012)

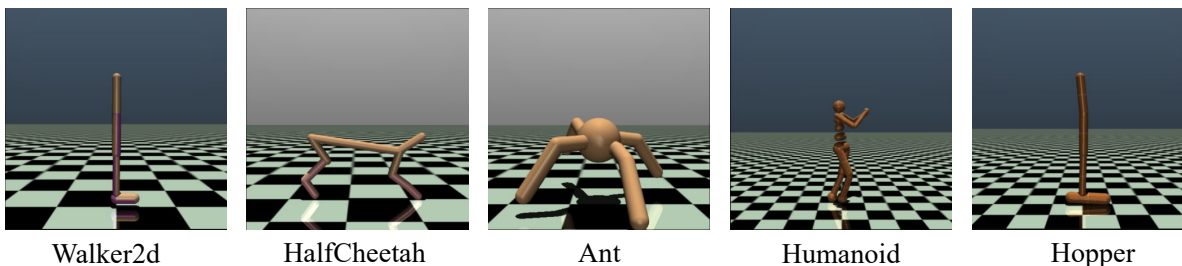


Figure 13: The environments in our experiments.

The environments in our experiment are shown in Fig. 13 and detailed as follows:

- 1) Walker2d aims to walk in the forward direction as fast as possible.
- 2) HalfCheetah aims to run forward as fast as possible.
- 3) Ant aims to coordinate the four legs to move in the forward direction as fast as possible.
- 4) Humanoid aims to walk forward as fast as possible without falling over.
- 5) Hopper aims to make hops that move in the forward direction as fast as possible.

The compared RL algorithms in our experiment are detailed as follows.

- Online:
 - 1) TD3 takes the minimum value between a pair of critic networks to address the overestimation of Q -value and reduces per-update error by delaying policy updates (Fujimoto et al., 2018).
 - 2) SAC is an actor-critic algorithm based on the maximum entropy approach. The objective encourages policy stochasticity by augmenting the reward with the entropy at each step (Haarnoja et al., 2018).
 - 3) OFE-TD3 increases the input dimensionality of the networks by representation learning to improve the sample efficiency of TD3 (Ota et al., 2020).
 - 4) TQC addresses the overestimation of Q -value by the combination of the distributional representation of a critic, truncation of critic prediction, and ensembling of multiple critics (Kuznetsov et al., 2020).
 - 5) TD7 is an effective DRL algorithm which combines TD3, state representation learning, checkpoints, prioritized experience replay, and a behaviour cloning term (only used for offline RL) (Fujimoto et al., 2023).

The hyper-parameters of HA3C are shown in Table 3. For Hopper, γ is set as 0.992. Network architecture details are described in Pseudocode 1-3. The optimizer of our networks is Adam Kingma (2015).

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 3: Hyper-parameters.

Parameter	Value	Brief explanation
Start-timesteps	25000	Time steps of the initial random policy is used
Batch-size	512	Batch size for both actor and critic
t_{pol}	2	Policy update frequency
t_{tar}	250	Target update rate
t_{ear}	1	Early assessment episodes for checkpoint operation
t_{lat}	3	Late assessment episodes for checkpoint operation
T_{ear}	750K	Early time steps for checkpoint operation
σ_e	0.1	Std of exploration noise
σ_T	0.005	Std of target policy noise
c	(-0.11,0.11)	Target policy noise clipping
k	6	The length of the considering state sequences
γ	0.99	Discount factor
l_e	0.0006	The learning rate of the encoder network
l_p	0.0003	The learning rate of the actor-network
l_Q	0.0003	The learning rate of the network of the Q -functions
α	0.25	Controlling the amount of prioritization in LAP
P_m	1.1	Minimum priority in LAP

Pseudocode 1: Critic network Details**Critic network:**

L1 = Linear(state-dim + action-dim, 256)

L2 = Linear(z^s -dim * 2 + 256, 256)

L3 = Linear(256, 256)

L4 = Linear(256, 1)

Critic forward pass: $\mathbf{x} = \text{Concatenate}([s_t, \mathbf{a}_t])$ $\mathbf{x} = \text{AvgL1Norm}(\text{L1}(\mathbf{x}))$ $\mathbf{x} = \text{Concatenate}([z^{s_{k,t}, \mathbf{a}_t}, z^{s_{k,t}}, \mathbf{x}])$ $\mathbf{x} = \text{Elu}(\text{L2}(\mathbf{x}))$ $\mathbf{x} = \text{Elu}(\text{L3}(\mathbf{x}))$ $\tau(s_{k,t}, \mathbf{a}_t) = \text{L4}(\mathbf{x})$ **Pseudocode 2: Actor network Details****Actor network:**

L1 = Linear(state-dim, 256)

L2 = Linear(z^s -dim + 256, 256)

L3 = Linear(256, 256)

L4 = Linear(256, action-dim)

Actor forward pass: $\mathbf{x} = \text{AvgL1Norm}(\text{L1}(s_t))$ $\mathbf{x} = \text{Concatenate}([z^{s_{k,t}}, \mathbf{x}])$ $\mathbf{x} = \text{ReLU}(\text{L1}(\mathbf{x}))$ $\mathbf{x} = \text{ReLU}(\text{L2}(\mathbf{x}))$ $\mathbf{a}_t = \text{Tanh}(\text{L3}(\mathbf{x}))$

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

Pseudocode 3: Encoder Details

State Encoder f Network:

Conv = Conv2d(kernel-num=64, kernel-size=(3, state-dim), stride=1)
Pool = MaxPool2d((1, 1))
L1 = Linear(64, 16)
L2 = Linear(state-dim, 256)
L3 = Linear(256+16, 256)
L4 = Linear(256, zs-dim)

State Encoder f Forward Pass:

$x = \text{Conv}(s_{k-1, t-1})$
 $x = \text{Pool}(x)$
 $x = \text{Elu}(L1(x))$
 $x = \text{AvgL1Norm}(x)$
 $y = \text{Elu}(L2(s_t))$
 $x = \text{Concatenate}([x, y])$
 $x = \text{Elu}(L3(x))$
 $z^{s_{k,t}} = \text{AvgL1Norm}(L4(x))$

State-Action Encoder g Network:

L1 = Linear(action-dim + z^s -dim, 256)
L2 = Linear(256, 256)
L3 = Linear(256, z^s -dim)

State-Action Encoder g Forward Pass:

$x = \text{Concatenate}([a_t, z^{s_{k,t}}])$
 $x = \text{Elu}(L1(x))$
 $x = \text{Elu}(L2(x))$
 $z^{s_{k,t}, a_t} = L3(x)$

F SUPPLEMENTARY EXPERIMENT

F.1 BIPEDAWALKER EXPERIMENT

To illustrate the benefit of history augmentation for complex MDP tasks, we test HA3C and No Aug. (HA3C without historical augmentation) on BipedalWalker and BipedalWalker-hardcore tasks. In BipedalWalker a robot is trained to move forward with slightly uneven terrain. Compared with BipedalWalker, BipedalWalker-hardcore is a more complex task, where the above robot is trained to move forward with ladders, stumps, and pitfalls. Therefore, the causal relationships in the transitions of BipedalWalker-hardcore are more complex than those in the transitions of BipedalWalker. The environments and learning curves are shown in Fig. 14 and the numerical results are shown in Table 4.

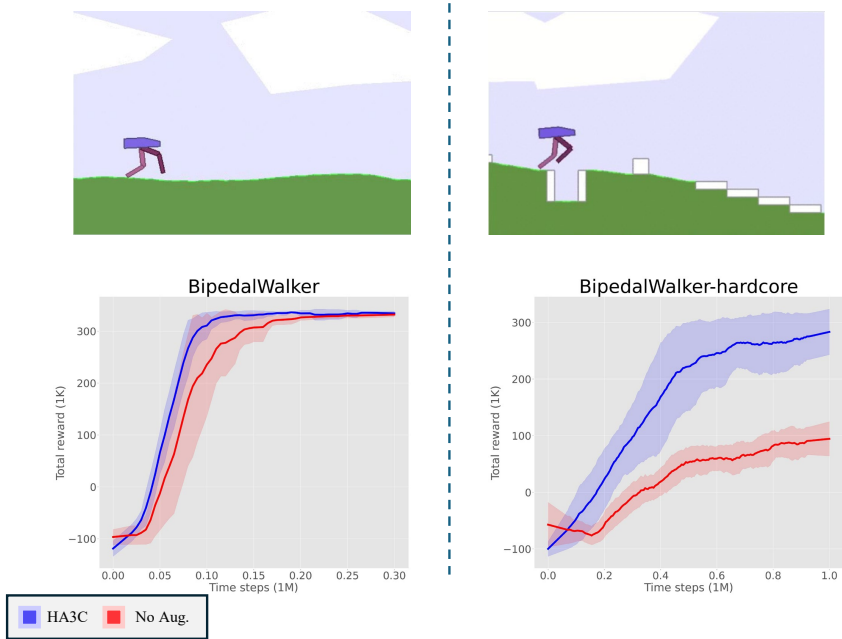


Figure 14: The environments and learning curves on BipedalWalker and BipedalWalker-hardcore tasks.

Table 4: The average highest returns of HA3C and No Aug. on BipedalWalker and BipedalWalker-hardcore tasks.

Algorithm	BipedalWalke	BipedalWalker-hardcore
HA3C	332 ±27	316 ±19
No Aug.	325 ±31	171 ±21

As we can see, although, both HA3C and No Aug. can get the high cumulative rewards in BipedalWalker, only HA3C can get the high cumulative rewards in BipedalWalker-hardcore. This is because by historical augmentation our HA3C can simplify the causal relationships in the transitions of BipedalWalker-hardcore.

F.2 VISUALIZED RESULTS OF HA3C

Fig. 15 presents the visual results of the transitions in HA3C and No Aug. The collected states of each control task are mapped together by UMAP. The max learning step is 4×10^5 and each state is coloured by the reward of reaching it.

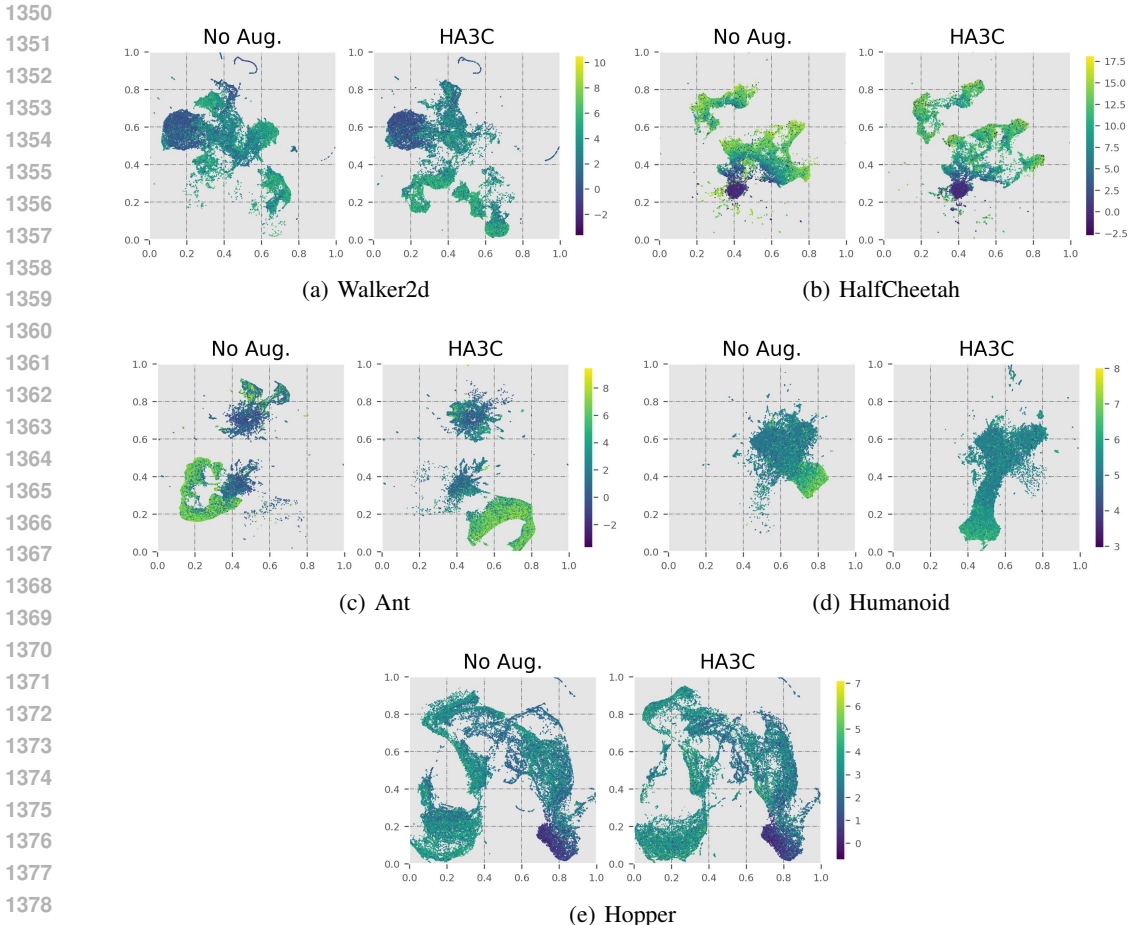


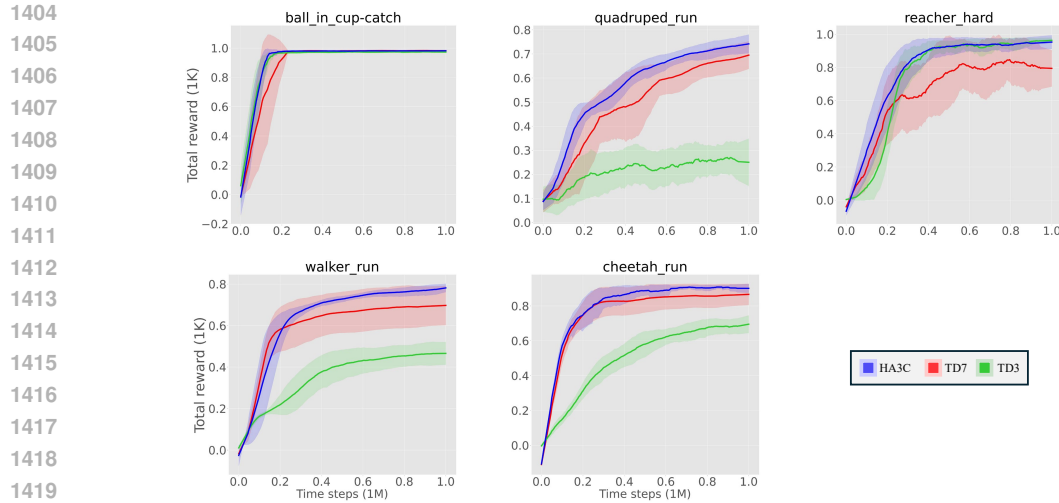
Figure 15: Visualized results of the explored states in No Aug. and HA3C.

As we can see, in Walker2d, Ant, and Humanoid, the high-reward states from HA3C are more than those from No Aug. This result illustrates that the sample efficiency of DRL can be effectively improved by learning the state representations with historical augmentation.

F.3 DEEP MIND CONTROL SUITE EXPERIMENT

In this subsection, we evaluate our HA3C on five DMC tasks including ball_in_cup-catch, walker-run, quadruped-run, cheetah-run, and reacher-hard (Tassa et al., 2018). The compared algorithms are TD3 (Fujimoto et al., 2018) and TD7 (Fujimoto et al., 2023). For all algorithms, each task runs 10 instances with 10^6 time steps with different random seeds. In each instance, the evaluation is performed every 5000 time steps. Some parameters are changed as follows. For quadruped-run, l_e is set as 0.0006, σ_T is set as 0.06, and c is set as $(-0.12, 0.12)$. For other tasks, l_e is set as 0.0005 and c is set as $(-0.1, 0.1)$. The learning curves are shown in Fig. 18 and the numerical results at 300K time step and 1M time step are shown in Table 5.

As we can see, in most cases, HA3C has higher cumulative rewards than the compared algorithms. For walker-run, quadruped-run, and reacher-hard, HA3C outperforms the compared algorithms in terms of both the early performance and the final performance. For ball_in_cup-catch and cheetah-run, HA3C outperforms all of the compared algorithms in the final performance but the average return of HA3C is lower than the average return of TD7 in the early performance.



1420
1421 Figure 16: Learning curves of different RL algorithms on the deep mind control suite tasks.

1422
1423 Table 5: The average highest returns over 10 instances on the deep mind control suite tasks at 400K
1424 and 1M time steps.

1425

Algorithm	Time step	ball_in_cup-catch	walker-run	quadruped-run	cheetah-run	reacher-hard
TD3	400K	981±2	387±71	331±65	550±76	971±3
	1M	985±1	481±54	444±22	729±39	979±1
TD7	400K	990±2	654±96	531±69	836±75	879±91
	1M	991±1	706±95	703±54	868±56	979±5
HA3C	400K	989±2	713±41	598±36	834±108	976±5
	1M	992±1	789±19	758±24	916±5	985±5

1426
1427
1428
1429
1430
1431
1432
1433
1434

1435
1436 F.4 LONGER TRAINING RUNS

1437
1438 In this section, we compare our HA3C with TD7 on five Mujoco control tasks with 3M training steps.
1439 The learning curves are shown in Fig. 17 and the numerical results are shown in Table 6.

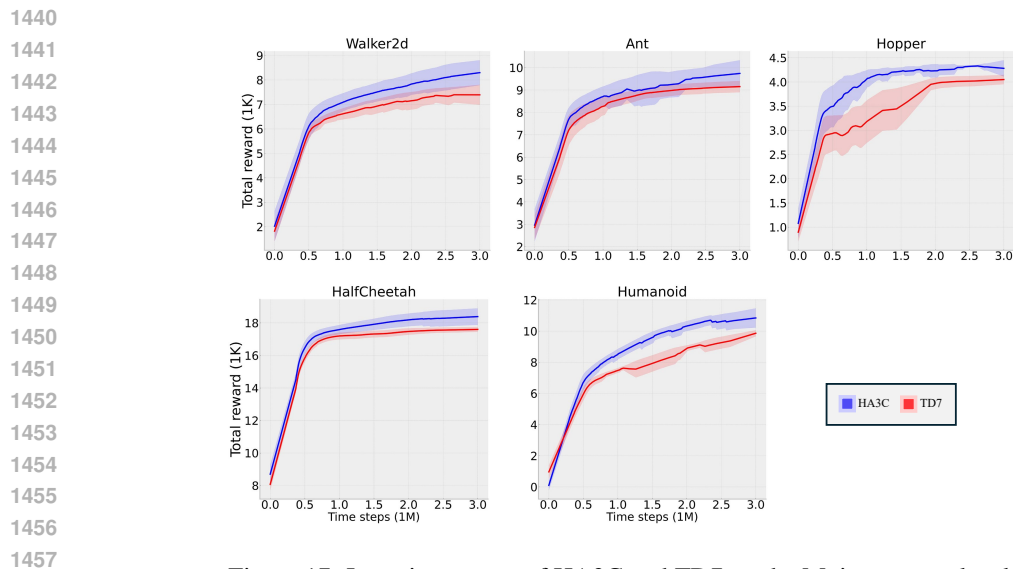


Figure 17: Learning curves of HA3C and TD7 on the Mujoco control tasks.

Table 6: The average highest returns over 10 instances of HA3C and TD7 at 3M time steps. \pm captures the standard deviation over trials.

Algorithm	Walker2d	HalfCheetah	Ant	Humanoid	Hopper
TD7	7570 \pm 321	17787 \pm 286	9225 \pm 450	9850 \pm 226	4049 \pm 156
HA3C	8463 \pm 829	18687 \pm 683	9794 \pm 891	11381 \pm 344	4413 \pm 59

As we can see, HA3C outperforms TD7 on the five Mujoco control tasks. It is noteworthy that the cumulative rewards of HA3C are significantly higher than the cumulative rewards of TD7 on Walker2d, Humanoid, and Hopper.

F.5 COMBINING HISTORICAL REPRESENTATION LEARNING WITH SAC

In this section, we combine our historical representation learning with SAC to construct HA3C-SAC method (Haarnoja et al., 2018). Then we evaluate HA3C-SAC on three MuJoCo control tasks including Walker2d, Humanoid, and Hopper. The compared methods includes the original SAC and SALE-SAC, which combines the representation learning with SAC without historical augmentation (Fujimoto et al., 2023). The learning curves are shown in Fig. 17 and the numerical results are shown in Table 7.

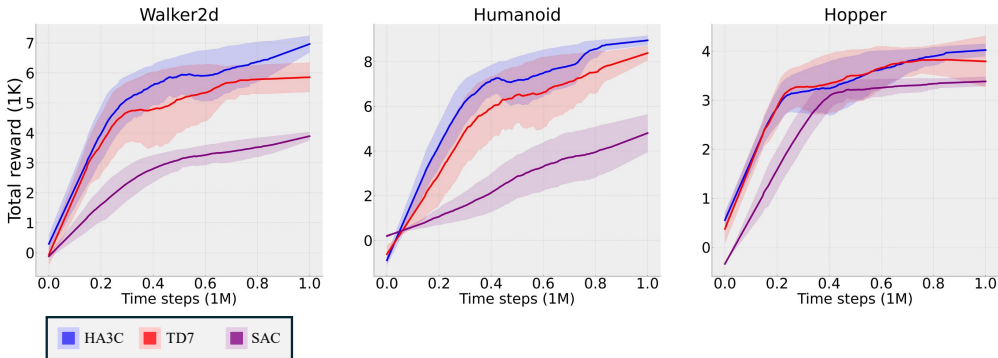


Figure 18: Learning curves of different RL algorithms on the deep mind control suite tasks.

Table 7: The average highest returns on Mujoco control tasks at 400K and 1M time steps.

Algorithm	Time step	Walker2d	Humanoid	Hopper
SAC	400K	2843 \pm 148	2268 \pm 905	3195 \pm 33
	1M	3921 \pm 163	5498 \pm 131	3422 \pm 86
SALE-SAC	400K	5414 \pm 377	6430 \pm 191	3515 \pm 125
	1M	6021 \pm 492	8368 \pm 330	4038 \pm 126
HA3C-SAC	400K	5796 \pm 395	7112 \pm 339	3566 \pm 39
	1M	6950 \pm 623	9047 \pm 238	4131 \pm 48

As we can see, HA3C-SAC outperforms SAC and SALE-SAC on the three Mujoco control tasks. The above results and the results Section 5.1 illustrate that our historical representation learning is robust to different algorithms and tasks.