
Stationary Deep Reinforcement Learning with Quantum K-spin Hamiltonian Equation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A foundational issue in deep reinforcement learning (DRL) is that *Bellman’s optimality equation has multiple fixed points*—failing to return a consistent one. A
2 direct evidence is the instability of existing DRL algorithms, namely, the high
3 variance of cumulative rewards over multiple runs. As a fix of this problem, we
4 propose a quantum K-spin Hamiltonian regularization term (H-term) to help a
5 policy network stably find a *stationary* policy, which represents the lowest energy
6 configuration of a system. First, we make a novel analogy between a Markov
7 Decision Process (MDP) and a *quantum K-spin Ising model* and reformulate the
8 objective function into a quantum K-spin Hamiltonian equation, a functional of
9 policy that measures its energy. Then, we propose a generic actor-critic algorithm
10 that utilizes the H-term to regularize the policy/actor network and provide Hamilto-
11 nian policy gradient calculations. Finally, on six challenging MuJoCo tasks over
12 20 runs, the proposed algorithm reduces the variance of cumulative rewards by
13 65.2% \sim 85.6% compared with those of existing algorithms.
14

15 1 Introduction

16 Deep reinforcement learning (DRL) [36] algorithms are quite unstable, namely, agents trained with
17 different random seeds may have dramatically different performance. Existing works [1, 13, 6, 24]
18 reported a high variance over multiple runs, thus it requires to train tens of agents and then pick the
19 best one. Such a high variance puts a challenge on reliability and reproducibility [15, 14], limiting
20 the wider adoption in real-world tasks.

21 The difficulty of guaranteeing stability lies in a foundational issue that **Bellman’s optimality equation**
22 **has multiple fixed points** [4, 31, 22, 18]—failing to return a consistent one. Consider MDP examples
23 with an terminal state 0, as shown in Fig. 1 (we adapt dynamic programming examples [4, 31] into
24 reinforcement learning settings. Observational experiments [are given](#) in Section 2.2),

- 25 • **Shortest path problem (deterministic)** in Fig. 1(a): At state 1, an agent transits to either state
26 1 or 0 with reward 0 or b , respectively. Assume the value function for state 0 is $V(0) = 0$. The
27 Bellman’s optimality equation for state 1 is $V(1) = \max\{V(1), b\}$, where any $V(1) \geq b$ is a
28 feasible solution. If initialize $V_0(1) \geq b$, a resulting policy is that an agent at state 1 always transits
29 back to state 1; otherwise, drives to terminal state 0 (always returns back to itself with reward 0).
- 30 • **Blackmailer’s problem (stochastic)** in Fig. 1(b): At state 1, a profit maximizing blackmailer
31 demands a cash amount $a \in (0, 1]$; a victim transits to state 1 with probability a or state 0 with
32 probability $1 - a$, respectively. At state 0, a victim always refuses to yield, i.e., $V(0) = 0$. The
33 Bellman’s optimality equation for state 1 is $V(1) = \max_a \{a + (1 - a)V(1)\}$, where any $V(1) \geq 1$
34 is a feasible solution. If initialize $V_0(1) > 1$, the blackmailer’s policy is demanding $a \rightarrow 0$ to keep
35 the victim at state 1; otherwise, demanding $a = 1$ that drives the victim to terminal state 0.

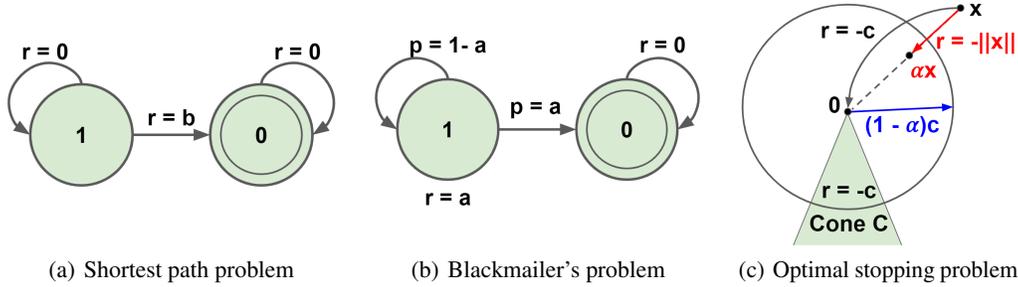


Figure 1: MDP examples where $\gamma = 1$. Examples with $\gamma < 1$ are given in Fig. 5 in Appx. A.

36 • **Optimal stopping problem (terminating policies)** in Fig. 1(c): In a space \mathbb{R}^2 with terminal
 37 state of point 0, an agent at point $x \neq 0$ moves to either point 0 with negative reward $-c$ or
 38 point αx with reward $-||x||$, respectively, where $\alpha \in (0, 1)$. The Bellman’s optimality equation
 39 is $V(x) = \max\{-c, -||x|| + V(\alpha x)\}$ and the optimal policy is to continue inside the sphere of
 40 radius $(1 - \alpha)c$ and to stop outside. If add a cone region C within which an agent always receives
 41 a reward $-c$, a second policy is jumping to point 0 at any point in region C .

42 The instability problem has been partially addressed, such as ensemble methods [2, 8], regularization
 43 approaches [38, 9], and baseline-correction approaches [33, 42]. In particular, Generalized Advantage
 44 Estimation (GAE) [33] is a widely used one that significantly reduces the variance of the advantage
 45 function. However, they did NOT fix the foundational issue of Bellman’s optimality equation in Fig.
 46 1, thus the objective function inherently fails to search for a consistent policy. For practical usage, we
 47 often expect a DRL algorithm stably converges to a certain policy independent of initialization and
 48 noises.

49 As a fix of the problem, we make a novel analogy between an MDP and a *quantum K-spin Ising model*
 50 [26, 12], and reformulate a reward-maximization RL task into an energy-minimization problem,
 51 namely, finding the lowest-energy configuration of a quantum spin system. We hypothesize that
 52 a *physically stationary policy would have the lowest energy*. Different from our quantum K-spin
 53 perspective, several recent papers utilized the (classical) Hamiltonian equation to endow RL agents
 54 the capability of inductive biases. For example, [21, 40] used Hamiltonian mechanics to train an
 55 agent that learns and respects *conservation laws*; [43] applied a Hamiltonian Monte Carlo (HMC)
 56 simulator to approximate the posterior action probability; and [28] proposed an unbiased estimator
 57 for the stochastic Hamiltonian gradient methods for min-max optimization problems.

58 In this paper, we propose a quantum-inspired *K-spin Hamiltonian regularization term (H-term)* that
 59 helps policy gradient algorithms stably converge to a physically stationary policy. Our contributions
 60 are as follows: 1) we derive a novel regularizer, *H-term*, by a novel reformulation of the RL’s objective
 61 as a K-spin Hamiltonian equation that measures the energy of a policy; 2) we propose a generic
 62 actor-critic algorithm that utilizes the H-term to regularize the policy/actor network and provide
 63 Hamiltonian policy gradient calculations on both deterministic and stochastic cases; and 3) on six
 64 challenging MuJoCo tasks over 20 runs, we show that the proposed algorithm reduces the variance of
 65 cumulative rewards by 65.2% \sim 85.6% compared with those of existing algorithms.

66 2 Bellman’s Optimality Equation and Multiple Fixed Points

67 We describe Bellman’s optimality equation and provide observational experiments to verify the
 68 existence of multiple policies.

69 2.1 Bellman’s Optimality Equation

70 A reinforcement learning (RL) [36] agent interacts with an unknown environment and learns an
 71 optimal policy that maximizes the cumulative reward. Mathematically, the environment can be
 72 formulated as a Markov Decision Process (MDP) with the five-tuple $\langle \mathcal{S}, \mathcal{A}, \mathbb{P}, R, \gamma \rangle$. Here \mathcal{S} and
 73 \mathcal{A} denote the state and action spaces; $\mathbb{P} : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ denotes a transition probability function,
 74 where Δ is a probability simplex; $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ denotes a reward function; and $\gamma \in (0, 1]$

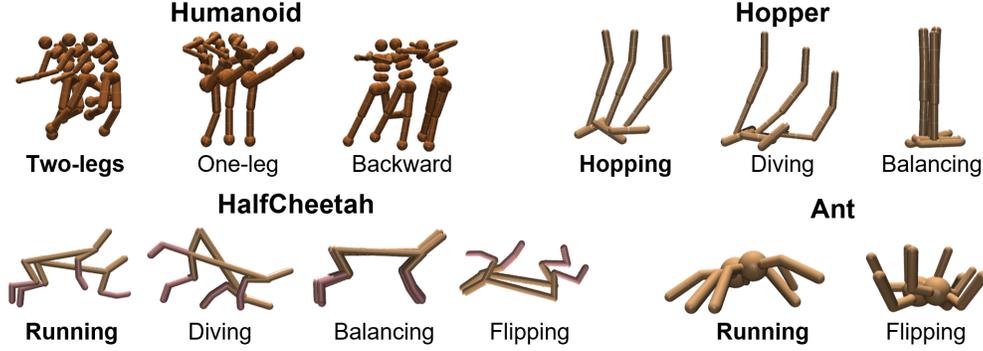


Figure 2: Different policy types for four MuJoCo [39] tasks. **Two-legs running**, one-leg running, and running backward for Humanoid; **running**, diving, and balancing for Hopper; **running**, diving, balancing and flipping for HalfCheetah; **running** and flipping for Ant. The bold ones indicate physically stationary policies.

75 denotes a discount factor. The objective is to find an optimal policy $\pi^* : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ that maximizes
 76 (discounted) expected reward.

77 Consider a discrete, finite, discounted MDP with infinite horizon, one can define the Q-value function
 78 of a state-action pair (s, a) under policy π as follows

$$Q^\pi(s, a) = \mathbb{E}_{S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k), A_{k+1} \sim \pi(S_{k+1}, \cdot)} \left[\sum_{k=0}^{\infty} \gamma^k \cdot R(S_k, A_k, S_{k+1}) | S_0 = s, A_0 = a \right], \quad (1)$$

79 where $R(S_k, A_k, S_{k+1})$ denotes the immediate reward when taking action A_k at state S_k and arriving
 80 at state S_{k+1} , capital letters denote random variables and lowercase letters denote values. The
 81 conventional objective function $J(\theta)$ of reinforcement learning [36] takes the following form

$$J(\theta) \triangleq \mathbb{E}_{S_0, A_0} [Q^{\pi_\theta}(S_0, A_0)] = \mathbb{E}_{\tau \sim \pi} [R(\tau) \cdot P(\tau | \pi_\theta)], \quad (2)$$

82 where τ is a trajectory, i.e., $\tau = (S_0, A_0, \dots)$, and

$$83 \quad P(\tau | \pi_\theta) = d_0(s_0) \cdot \prod_{k=0}^T \mathbb{P}(s_{k+1} | s_k, a_k) \pi_\theta(a_k | s_k).$$

84 The Bellman equation [36] converts (1) into a recursive form as follows

$$\begin{aligned} Q^\pi(s, a) &= \sum_{s' \in \mathcal{S}} \mathbb{P}(s' | s, a) \left[R(s, a, s') + \gamma \sum_{a' \in \mathcal{A}} \pi(s', a') Q^\pi(s', a') \right] \\ &= R(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathbb{P}(s' | s, a) \sum_{a' \in \mathcal{A}} \pi(s', a') Q^\pi(s', a'), \end{aligned} \quad (3)$$

85 which expresses the expected reward as a summation of immediate reward $R(s, a)$ and discounted fu-
 86 ture rewards, and the immediate reward $R(s, a)$ is defined as $R(s, a) = \sum_{s' \in \mathcal{S}} \mathbb{P}(s' | s, a) R(s, a, s')$.

87 The Bellman's optimality equation [36] is

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} \mathbb{P}(s' | s, a) \left[R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]. \quad (4)$$

88 The optimal policy π^* is given by a greedy strategy such that $\pi^* = \arg \max_{\pi} Q^\pi(s, a)$.

89 2.2 Observational Experiments for Multiple Policies and Physically Stationary Policy

90 As we pointed out in Fig. 1, theoretically there exists multiple policies. Here, we perform observa-
 91 tional experiments on four challenging MuJoCo tasks [39], namely, Humanoid, Hopper, HalfCheetah,
 92 and Ant (details given in Appx. B.1), which are typical examples of the locomotion control of a robot.

93 We render the obtained policies over multiple runs and then identify physically stationary ones. We
 94 observe various types of moving strategies, as shown in Fig. 2, which verifies that multiple policies

Table 1: Analogy between MDP and quantum K-spin Ising model.

MDP (Our formulation in (7))		Quantum K-spin Ising Model [26, 12] in (5)	
State-action pairs	μ_0, \dots, μ_{K-1}	Spins	j_0, \dots, j_{K-1}
Optimal policy	$\pi_{\mu_0}^* \times \pi_{\mu_1}^* \times \dots \times \pi_{\mu_{K-1}}^*$	Optimal configuration	$\sigma_{j_0} \times \sigma_{j_1} \times \dots \times \sigma_{j_{K-1}}$
Policy	$\pi_{\mu_0} \times \pi_{\mu_1} \times \dots \times \pi_{\mu_{K-1}}$	Configuration	$\sigma_{j_0} \times \sigma_{j_1} \times \dots \times \sigma_{j_{K-1}}$
Discounted reward	$L_{\mu_0 \dots \mu_{K-1}}$	Density function	$L_{j_0 \dots j_{K-1}}$
Functional of policy	$H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}})$	Functional of spins	$H(\sigma_{j_0}, \dots, \sigma_{j_{K-1}})$
Stationary condition	$\frac{\delta H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}})}{\delta \pi_{\mu}} = 0$	Stationary condition	$\frac{\delta H(\sigma_{j_0}, \dots, \sigma_{j_{K-1}})}{\delta \sigma_j} = 0$

95 are very common. For example, the Humanoid agent learns either jumping with a single leg or
 96 running with two legs, as shown in Fig. 2 (top-left); another interesting example is HalfCheetah, in
 97 which an agent can run normally or in a flipped manner, as shown in Fig. 2 (bottom-left). Among
 98 the obtained policies, one can easily identify the physically stationary policies that control the robot
 99 moving forward with a *stable gait* (defined as gait that does not lead to fall).

100 3 Reinforcement Learning as Quantum K-spin Hamiltonian Equation

101 First, we make a novel analogy between the MDP and a quantum K-spin Ising model, and reformulate
 102 the objective function (2) into a quantum K-spin Hamiltonian equation. Then, we show that this
 103 Hamiltonian equation helps fix the issue of multiple fixed points in Fig. 1.

104 3.1 Motivation through Analogy with Quantum K-spin Ising Model

105 Given the multiple policies in Section 2.2, it is critical to employ a proper criteria to search for a
 106 physically stationary policy. In Table 1, we make a novel analogy between an MDP and a quantum
 107 K-spin Ising model [26, 12].

108 The Hamiltonian equation for a quantum K-spin Ising model [26, 12] measures the energy of a
 109 particular configuration of a quantum K-spin system and takes the following form

$$H = - \sum_{k=0}^{K-1} \sum_{j_0=1}^N \dots \sum_{j_k=1}^N L_{j_0 \dots j_k} \sigma_{j_0} \dots \sigma_{j_k}, \quad (5)$$

110 where N is the number of spins in the k -th configuration, $\sigma_{j_k} = \pm 1$ are spin variables, and $L_{j_0 \dots j_k}$ is
 111 an energy density function for k nearest spins' configuration $(\sigma_{j_0}, \dots, \sigma_{j_k})$.

112 **Analogy in Table 1.** Starting from an analogy between a state-action pair $\mu_k = (S_k, A_k)$ and
 113 a spin j_k , we can map an optimal policy $\pi^*(\mu_k) \in \{0, 1\}$ in (4) to a single-qubit spin operator
 114 $\sigma_{j_k} \in \{-1, 1\}$ via $\pi^*(\mu_k) \longleftrightarrow (\mathbb{1}_{\mu_k} - \sigma_{\mu_k})/2$, where $\pi_{\theta}(\mu_k)$ denotes the probability of taking
 115 action A_k at state S_k , following policy π_{θ} . The energy density function $L_{j_0 \dots j_k}$ can be defined as the
 116 discounted reward on a path $(\mu_0, \dots, \mu_{k-1})$ of length k , i.e.,

$$L_{\mu_0, \dots, \mu_k} = \gamma^k \cdot R(\mu_k) \cdot d_0(s_0) \cdot \prod_{\ell=0}^{k-1} \mathbb{P}(s_{\ell+1} | \mu_{\ell}), \quad (\text{obtained via Monte Carlo simulation}) \quad (6)$$

117 where $d_0(s_0)$ denotes the distribution of initial state s_0 . Analogy to the quantum K-spin Ising model,
 118 we can derive a functional of policy $H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}})$ in the context of MDP, which measures the
 119 energy of an RL policy. In this case, a physically stationary policy satisfies the stationary condition
 120 $\delta H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}}) / \delta \pi_{\mu} = 0$.

121 **3.2 Reformulation into Quantum K-spin Hamiltonian Equation**

122 Inspired by [17], we formally reformulate (2) into a K -spin Hamiltonian equation

$$\begin{aligned}
 H(\theta) &\triangleq -\mathbb{E}_{S_0, A_0} [Q^{\pi_\theta}(S_0, A_0)] \\
 &= -\lim_{K \rightarrow \infty} \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} \pi_\theta(\mu_0) \cdots \pi_\theta(\mu_k), \\
 &= -\lim_{K \rightarrow \infty} \mathbb{E}_{\mu_0, \mu_1, \dots, \mu_K} \left[\sum_{k=0}^{K-1} L_{\mu_0, \dots, \mu_k} \right],
 \end{aligned} \tag{7}$$

123 where $K \rightarrow \infty$, the expectation is taken over $S_0 \sim d_0(\cdot)$, $A_0 \sim \pi_\theta(S_0, \cdot)$, and the density function
 124 L_{μ_0, \dots, μ_k} is given in (6).

125 The following blue part will be put into the appendix.

126 **Monte Carlo Estimator** [30]: Consider a general probabilistic objective \mathcal{F} of the form:

$$\mathcal{F} \triangleq \mathbb{E}_{p(\mathbf{x}; \theta)} [f(\mathbf{x}; \phi)], \tag{8}$$

127 in which a function f of an input variable \mathbf{x} with *structural parameters* ϕ is evaluated on average
 128 with respect to an input distribution $p(\mathbf{x}; \theta)$ with *distributional parameters* θ .

129 A Monte Carlo method evaluates the function by first drawing independent samples $\hat{\mathbf{x}}^{(1)}, \dots, \hat{\mathbf{x}}^{(N)}$
 130 from the distribution $p(\mathbf{x}; \theta)$, and then computing the average:

$$\hat{\mathcal{F}}_N = \frac{1}{N} \sum_{i=1}^N f(\hat{\mathbf{x}}^{(i)}), \text{ where } \hat{\mathbf{x}}^{(i)} \sim p(\mathbf{x}; \theta) \text{ for } i = 1, \dots, N. \tag{9}$$

131

132 The Monte Carlo estimator for (2) is

$$\hat{J}(\theta) = \frac{1}{N} \sum_{i=1}^N R(\tau^{(i)}), \text{ where } \tau^{(i)} \sim P(\tau^{(i)} | \pi_\theta) \text{ for } i = 1, \dots, N, \tag{10}$$

133 and

134

$$P(\tau^{(i)} | \pi_\theta) = d_0(s_0^{(i)}) \cdot \prod_{k=0}^{T-1} \mathbb{P}(s_{k+1}^{(i)} | s_k^{(i)}, a_k^{(i)}) \pi_\theta(a_k^{(i)} | s_k^{(i)}). \tag{11}$$

135

136 The Monte Carlo estimator for (7) is

$$\hat{H}(\theta) = \frac{1}{N'} \sum_{i=1}^{N'} \sum_{k=0}^{K-1} L_{\mu_0^{(i)}, \dots, \mu_k^{(i)}}, \text{ for } i = 1, \dots, N', \tag{12}$$

137 and

$$L_{\mu_0^{(i)}, \dots, \mu_k^{(i)}} = \gamma^k \cdot R(\mu_k^{(i)}) \cdot d_0(s_0^{(i)}) \cdot \prod_{\ell=0}^{k-1} \mathbb{P}(s_{\ell+1}^{(i)} | \mu_\ell^{(i)}). \tag{13}$$

138

139 **Remark:** The above two Monte Carlo estimators are quite different in the simulation process. (11)
 140 samples a random trajectory by following an environment's stochastic transition and a policy. In
 141 contrast, (13) measures a random path's discounted reward (the "energy") without following any
 142 policy, and the Hamiltonian equation (7) combinatorially enumerates all possible paths of length K
 143 over the state-action space. In other words, the simulation process of the Hamiltonian term does not
 144 rely on any policy. Therefore, the Hamiltonian term is a suitable regularizer for both on-policy and
 145 off-policy algorithms.

146 This fundamental difference is due to the Ising model in (5), which combinatorially enumerates all
 147 paths and separates the environment and the policy.

148 **Physical interpretation:** Analogy to a quantum K-spin system, $H(\theta)$ in (7) measures the energy of
 149 policy π . We hypothesize that the energy of a policy is a favorable criteria, since a stationary policy
 150 with minimum energy: 1). *achieves a relative high reward independent of the initialization;* and 2). *is*
 151 *robust to interference/noise in the inference stage.*

152 **K-step truncation in practice.** Minimizing (7) is NP-hard [11]. Since $\gamma \in (0, 1)$, γ^K monotonically
 153 decreased with look-ahead steps K , therefore, we truncate (7) to finite K terms. One can show
 154 that these K terms in (7) is a geometric sequence with a truncation error ratio $1 - \gamma^K$. Assuming
 155 $1 - \gamma^K \leq 1 - \epsilon$, where $\epsilon > 0$ is small, thus we have the look-ahead steps satisfies $K \geq \log_{\gamma} \epsilon$.

156 3.3 Revisiting Examples in Fig. 1

157 We elaborate how adding the energy measured by (7) onto each state can help drive to the terminal
 158 state (a stationary policy), which fixes the foundational issue of multiple fixed points in Fig. 1 where
 159 $\gamma = 1$. We have $H(0) = 0$ for the terminal state 0.

- 160 • (a) **Shortest path problem (deterministic):** $H(1) = -\sum_{k=1}^{\infty} b = -\infty$. At state 1, the Bellman’s
 161 optimality equation becomes $V(1) = \max\{V(1) + \lambda H(1), b\}$. Independent of the initial value
 162 $V_0(1)$, an agent obtains a policy that always transits back to terminal state 0.
- 163 • (b) **Blackmailer’s problem (stochastic):** $H(1) = -\infty$. The Bellman’s optimality equation
 164 becomes $V(1) = \max_a\{a + (1 - a)(V(1) + \lambda H(1))\}$ for state 1. For any $V_0(1) < \infty$, the optimal
 165 policy becomes $a = 1$ that drives to the terminal state 0.
- 166 • (c) **Optimal stopping problem (terminating policies):** any policy that takes infinite steps will
 167 have $H(x) = -\infty$, since at each step number k , there are always trajectories that jump to point 0
 168 with reward $-c$; and a direct jumping policy will have $H(x) = -c$. Therefore, adding $H(x)$ to
 169 each point $x \neq 0$ will lead to a policy of *jumping back to point 0*.

170 4 Actor-Critic Algorithm with Quantum K-spin Hamiltonian Regularization

171 We propose a generic actor-critic algorithm with a H-term, derive two Hamiltonian’s policy gradient
 172 theorems for both deterministic and stochastic cases, and present the Monte Carlo gradient estimation.

173 4.1 Stationary Actor-Critic Algorithm with H-term

174 Actor-critic algorithms in reinforcement learning perform a bilevel optimization, namely alternating
 175 between approximating a value function and optimizing a policy. In practice, a critic network with
 176 parameter ϕ approximates the Q -value function, and an actor network with parameter θ approximates
 177 the policy π , details given in Appx. F. However, since the critic’s update is governed by the Bellman’s
 178 optimality equation, actor-critic algorithms suffer the multiple fixed points problem.

179 Motivated by Section 3.3, we propose a novel H-term for both deterministic and stochastic actor-critic
 180 algorithms. Similar to the entropy term in [23], the proposed H-term is an add-on term to regularize
 181 the actor network and help it converge to a stationary policy. Specifically, the objective functions of
 182 actor and critic networks become:

$$\begin{cases} \text{Actor : } \max_{\theta} J_{\pi}(\theta, \phi) \triangleq (1 - \gamma) \mathbb{E}_{S_0 \sim d_0, A_0 \sim \pi_{\theta}(S_0, \cdot)} [Q_{\phi}(S_0, A_0)] - \lambda H(\theta), \\ \text{Critic : } \min_{\phi} J_Q(\theta, \phi) \triangleq \frac{1}{2} \mathbb{E}_{S \sim d_{\theta}(\cdot), A \sim \pi_{\theta}(S, \cdot)} [(Q_{\phi}(S, A) - y(S, A))^2], \end{cases} \quad (14)$$

183 where a target Q-value is $y(S_k, A_k) = R(S_k, A_k) + \gamma Q_{\phi}(S_{k+1}, A_{k+1})$, and $\lambda > 0$ is a temperature
 184 parameter. As an interpretation, the second term $-\lambda H(\theta)$ in the maximization objective function of
 185 actor network aims to find a minimum energy configuration for the MDP problem, namely, a policy
 186 π that will add a minimum amount of energy to each state’s value function (as in Section 3.3).

187 **New algorithm.** In Alg. 1, an agent interacts with an environment and alternatively updates its actor
 188 network and critic network. The algorithm has M episodes and each episode consists of a (Monte
 189 Carlo) simulation process and a learning process (gradient estimation) as follows:

- 190 • During the (Monte Carlo) simulation process (lines 5-10 of Alg. 1), an agent takes action a_t ac-
 191 cording to a policy $\pi_{\theta}(\cdot | s_t)$, $t = 0, \dots, T - 1$, generating a trajectory of T steps/transitions.

Algorithm 1 Stationary Actor-Critic Algorithm with H-term

```
1: Input: learning rate  $\alpha$ , temperature  $\lambda$ , look-ahead step  $K$ , and parameters  $M, T, G, B, B'$ 
2: Initialize actor network  $\pi$  and critic network  $Q$  with parameters  $\theta, \phi$ , and replay buffers  $\mathcal{D}_1, \mathcal{D}_2$ 
3: for episode = 1,  $\dots$ ,  $M$  do
4:   Initialize state  $s_0$ 
5:   for  $t = 0, \dots, T - 1$  do
6:     Select action  $a_t \sim \pi_\theta(\cdot|s_t)$ 
7:     Execute action  $a_t$ , receive reward  $r_t$ , and observe new state  $s_{t+1}$ 
8:     Store a transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}_1$ 
9:   end
10:  Store a trajectory  $\tau$  of length  $T$  in  $\mathcal{D}_2$ 
11:  for  $g = 1, \dots, G$  do
12:    Randomly sample a mini-batch of  $B$  transitions  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^B$  from  $\mathcal{D}_1$ 
13:    Randomly sample a mini-batch of  $B'$  trajectories (of length  $K$ )  $\{\tau_j\}_{j=1}^{B'}$  from  $\mathcal{D}_2$ 
14:    Update critic network using a conventional method
15:    Update actor network as  $\theta \leftarrow \theta + \alpha \left( \nabla_\theta \hat{J}(\theta) - \lambda \nabla_\theta \hat{H}(\theta) \right)$ .
16:  end
17: end
```

192 Then, these T transitions are stored into a replay buffer \mathcal{D}_1 , while the full trajectory $\tau =$
193 $(s_0, a_0, r_0, s_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$ is stored in replay buffer \mathcal{D}_2 .

194 • During the learning process ($G \geq 1$ updates in one episode) (lines 11-16 of Alg. 1), a mini-
195 batch of B transitions $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^B$ and a mini-batch of B' trajectories (of length K)
196 $\{\tau_j = (s_0^j, a_0^j, r_0^j, s_1^j, \dots, s_{K-1}^j, a_{K-1}^j, r_{K-1}^j, s_K^j)\}_{j=1}^{B'}$ are sampled from \mathcal{D}_1 and \mathcal{D}_2 , respectively.
197 The critic network is updated by a conventional method, e.g., minimizing the mean squared error
198 (MSE) between an estimated Q-value and a target value. The actor is updated by a Monte Carlo
199 gradient estimator over B transitions and B' trajectories.

200 **Two new hyperparameters.** We introduce two hyperparameters: a temperature $\lambda > 0$ that is a
201 relative weight of the H-term, and a look-ahead step $K \leq T$ that defines the horizon of the H-term.

202 **Implementation of replay buffer \mathcal{D}_2 .** After a full trajectory τ of length T is generated, it is
203 partitioned into $T - K + 1$ trajectories of length K . We rank them according to the cumulative reward
204 and store the top portion, say 80%, into a new replay buffer \mathcal{D}_2 (line 10 of Alg. 1). We randomly
205 sample a mini-batch of B' trajectories from \mathcal{D}_2 (line 13 of Alg. 1) to compute the H-term.

206 4.2 Hamiltonian Policy Gradient and Monte Carlo-based Gradient Estimator

207 We provide the policy gradient of the quantum K-spin Hamiltonian equation in (7), which are variants
208 of the well-known policy gradient theorem [36]. We provide detailed derivations in Appx. D.

209 **Stochastic version.** The Hamiltonian stochastic gradient of (7) w.r.t. parameter θ is

$$\nabla_\theta H(\theta) = -\mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_\theta \log (\pi_\theta(\mu_0) \cdot \pi_\theta(\mu_1) \cdots \pi_\theta(\mu_k)) \right]. \quad (15)$$

210 **Deterministic version.** Let $\eta_\theta(\cdot) : \mathcal{S} \rightarrow \mathcal{A}$ denote a deterministic policy, while we use $\tilde{\pi}_{\theta, \delta}(\mu)$ to
211 represent that a Gaussian noise (a.k.a, an exploration noise) with standard deviation $\delta > 0$ is added in
212 the exploration process. The Hamiltonian deterministic gradient of (7) w.r.t. parameter θ is

$$\nabla_\theta H'(\theta) = -\mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_\theta \log (\tilde{\pi}_{\theta, \delta}(\mu_0) \cdot \tilde{\pi}_{\theta, \delta}(\mu_1) \cdots \tilde{\pi}_{\theta, \delta}(\mu_k)) \right]. \quad (16)$$

213 The quantum K-spin Hamiltonian equation in (7) is a reformulation of (2). We verify the gradient
214 calculation by showing that: when $K \rightarrow \infty$, the Hamiltonian stochastic and deterministic policy
215 gradient $\nabla_\theta H(\theta)$ and $\nabla_\theta H'(\theta)$ are equal to the stochastic policy gradient $\nabla_\theta J(\theta)$ in [37] and
216 deterministic policy gradient $\nabla_\theta J'(\theta)$ in [35], respectively.

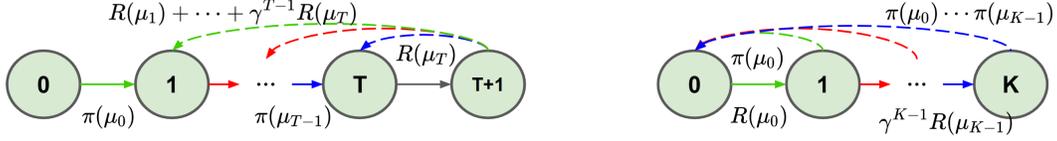


Figure 3: Comparison of REINFORCE’s (left) and Hamiltonian’s policy gradient (right).

217 Note that the gradient $\nabla_{\theta}H(\theta)$ in (15) and $\nabla_{\theta}H'(\theta)$ in (16) w.r.t. a distributional parameter θ takes
 218 an expectation form. Thus, a Monte Carlo gradient estimator is practically useful. We obtain the
 219 Monte Carlo gradient estimator of $\nabla_{\theta}H(\theta)$, illustrated in Fig. 3 (right), as follows

$$\nabla_{\theta}\hat{H}(\theta) = -\frac{1}{N'} \sum_{i=1}^{N'} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k^i) \cdot \nabla_{\theta} \log [\pi_{\theta}(\mu_0^i) \cdots \pi_{\theta}(\mu_k^i)] \right]. \quad (17)$$

220 As a contrast, we provide the Monte Carlo gradient estimator of REINFORCE’s [37] policy gradient,
 221 as illustrated in Fig. 3 (left), as follows

$$\nabla_{\theta}\hat{J}(\theta) = \frac{1}{NT} \sum_{i=1}^N \left[\sum_{t=0}^{T-1} G_t^i \cdot \nabla_{\theta} \log \pi_{\theta}(\mu_t^i) \right], \quad \text{where } G_t^i = \sum_{t'=t+1}^T \gamma^{t'-t-1} R(\mu_{t'}^i). \quad (18)$$

222 An interesting observation is that both gradient calculations follow a similar pattern as shown in Fig.
 223 3. REINFORCE’s policy gradient [37] in Fig. 3 (left) employs an estimate of future rewards, while
 224 Hamiltonian policy gradient in Fig. 3 (right) uses trajectories in replay buffer \mathcal{D}_2 .

225 **Computational complexity:** we measure the computation complexity by the times of computing one
 226 $\nabla_{\theta} \log \pi_{\theta}(\mu)$. Assume $N = B$ and $N' = B'$, since most DRL algorithms use a mini-batch stochastic
 227 gradient decent methods. REINFORCE’s [37] policy gradient in (18) takes $O(BT)$ computations,
 228 while Alg. 1 adds $O(B'K(K+1)/2)$ computations in each gradient update step, thus a total
 229 complexity of $O(BT + B'K(K+1)/2)$.

230 5 Performance Evaluation

231 We evaluate the proposed H-term from four aspects: 1) increasing cumulative reward, 2) reducing
 232 variance, 3) driving to physically stationary policy, and 4) the impact of trajectory length K . All
 233 experiments were executed on an NVIDIA DGX-2 server [10]. The server contains 8 A100 GPUs,
 234 320 GB GPU memory, and 128 CPU cores running at 2.25 GHz.

235 5.1 Experimental Settings

236 **Environments (tasks).** We consider six challenging MuJoCo tasks [39] as in Section 2.2. The agent
 237 learns to control the locomotion of a robot and aims to move forward as quickly as possible. These
 238 tasks have high-dimensional continuous state and action spaces, in which there exists multiple locally
 239 optimal polices as revealed in Section 2.2.

240 **Compared algorithms.** To evaluate deterministic and stochastic algorithms, we choose Deep
 241 Deterministic Policy Gradient (DDPG) [27] and Proximal Policy Optimization (PPO) [34]. Since the
 242 H-term is compatible with existing variance reduction techniques, we implement the PPO algorithm
 243 with GAE [33]. For a fair comparison, we keep the hyperparameters (listed in Appx. H) the same
 244 and make sure that the obtained results reproduce existing benchmark tests [13].

245 **Performance metrics.** We employ two performance metrics, the cumulative rewards and variance,
 246 while in Section 5.4, we further consider different policies and report the number of convergence. We
 247 run each experiment with 20 random seeds and in each run we test 100 episodes.

248 5.2 H-term Increases Cumulative Reward

249 Experience replay is crucial in improving performance in terms of cumulative reward. The proposed
 250 H-term in (17) can be viewed as a novel experience replay technique for an actor network. Here, we

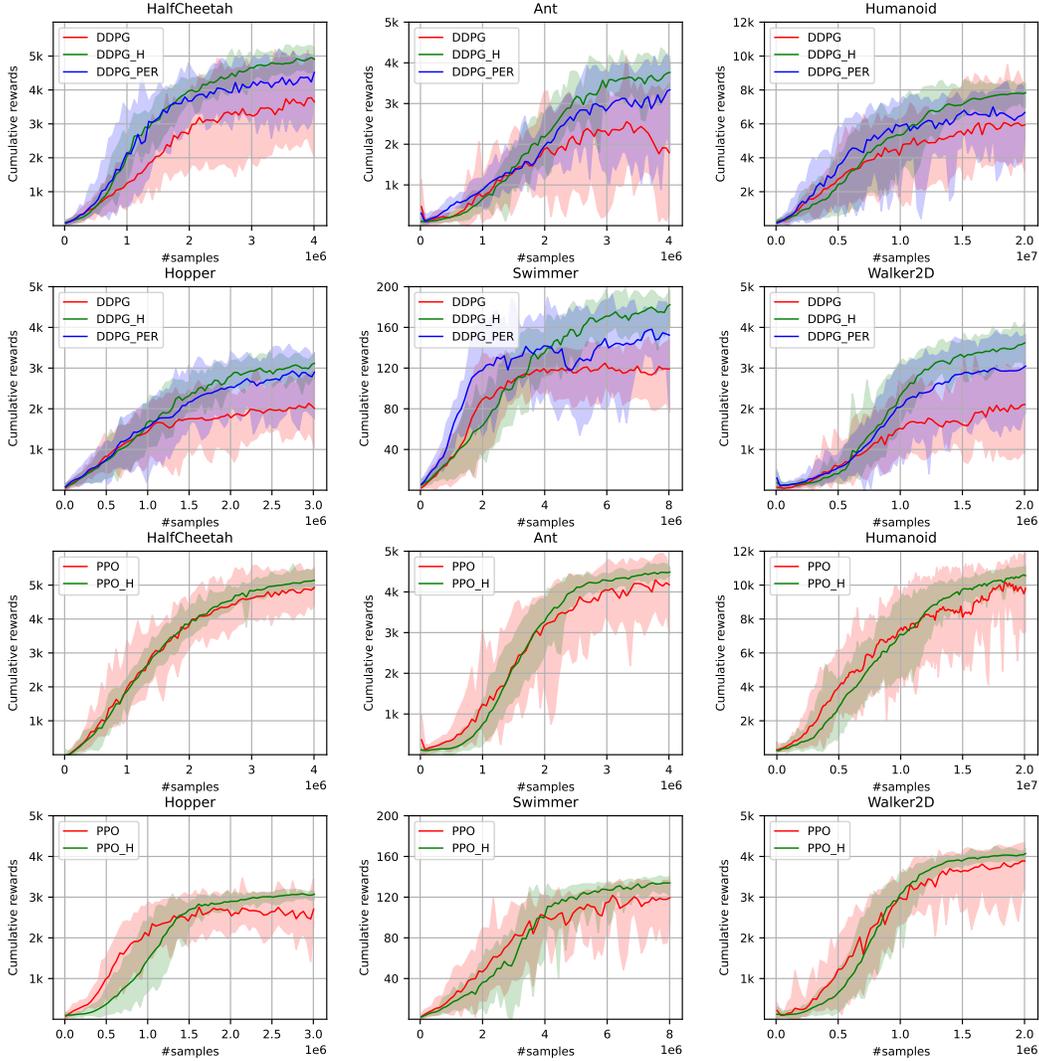


Figure 4: Cumulative rewards vs. #samples for compared DRL algorithms on six MuJoCo tasks.

251 add a compared algorithm, DDPG with Prioritized Experience Replay [32] (DDPG+PER), where
 252 PER prioritizes experience by the TD error to update a critic network.

253 In Fig. 4, both DDPG+PER and DDPG+H achieve a substantial improvement of cumulative reward. In
 254 particular, DDPG+H achieves the highest cumulative rewards in all six tasks, which are comparable to
 255 PPO’s performance in Fig. 4. It is worthwhile to discuss the advantage of DDPG+H over DDPG+PER.
 256 DDPG+PER utilizes a prioritized replay strategy to obtain a more accurate critic network, however,
 257 it is updated via the Bellman equation with the trouble of multiple fixed points. In contrast, the
 258 H-term in DDPG+H is performed on the actor network. Our results indicate that an experience replay
 259 technique on actor network may be much more powerful.

260 5.3 H-term Reduces Variance

261 The PPO algorithm with GAE is regarded as the state-of-the-art algorithm in MuJoCo environments.
 262 However, it still has a very high variance (the shaded area) after the policies have converged, as
 263 shown in Fig. 4. We observe that, at the end of training, the PPO algorithm has a variance of 969.2,
 264 1563.4, 2513.5, 905.3, 60.7, 1290.1 in the six tasks, respectively. Such a high variance is mainly due
 265 to the fact that the agent may converge to a random one of multiple policies.

Table 2: Experimental results on six challenging MuJoCo tasks.

Tasks	Policies	PPO		PPO+H ($K = 8$)		PPO+H ($K = 16$)		PPO+H ($K = 24$)	
HalfCheetah	running	13		19		20		20	
	flipping	5	4720.8	0	5028.4	0	5104.3	0	4995.1
	diving	1	± 969.2	1	± 211.3	0	± 228.4	0	± 383.3
	balancing	1		0		0		0	
Ant	running	17		20		20		20	
	jumping	0	4164	0	4505.3	0	4645.6	0	4662.5
	flipping	3	± 1563.4	0	± 253.6	0	± 225.4	0	± 277.5
Humanoid	two-legs	7		17		16		17	
	one-leg	12	9433.4	3	9670.3	4	10189.1	3	9942.2
	backward	1	± 2513.5	0	± 497.2	0	± 683.7	0	± 538.4
Hopper	hopping	10		18		20		20	
	diving	8	2659.3	2	3116.5	0	3300.1	0	3340.7
	balancing	2	± 905.3	0	± 289.4	0	± 184.2	0	± 191.5
Swimmer	moving	14		20		20		19	
	balancing	6	110.7	0	130.6	0	132.5	1	132.2
Walker	walking	5		16		16		15	
	diving	8	5461.7	2	5819.9	4	5927.2	5	6089.3
	balancing	7	± 1290.1	2	± 315.6	0	± 296.8	0	± 314.7

266 In Fig. 4, the shaded areas of PPO+H ($K = 16$) are dramatically smaller, i.e., a variance of 228.4,
267 225.4, 683.7, 184.2, 31.6, and 296.8, respectively. The variance has been reduced by 65.2% \sim 85.6%,
268 which verifies the effectiveness of the proposed H-term. In Fig. 4, we also observe that the H-term
269 can help the DDPG algorithm reduce variances, namely, the variances of DDPG+H are much smaller
270 than those of vanilla DDPG and DDPG+PER. Therefore, we may conclude that the H-term guides
271 the agent to search for a stationary policy among multiple feasible ones.

272 More experimental results are given in Appx. H due to the space limit, including the cases of $K = 8$
273 and $K = 24$, and the curves of the H-term value during the training process. One may verify that the
274 stationary policies have relative lower H-values.

275 5.4 H-term Drives to Physically Stationary Policy

276 A key question needs to be answered: *is H-term really guiding the agent converge to a physically*
277 *stationary policy?* Similar to Section 2.2, we perform observational experiments on the six MuJoCo
278 tasks and measure the number of convergence to the different policies over 20 runs. As shown in
279 Table 2, the vanilla PPO algorithm converges to the physically stationary policy (**bold**) with 13, 17, 7,
280 10, 14, and 5 times for the six tasks, while the PPO+H ($K = 16$) converges to the stationary policy
281 with 20, 20, 16, 20, 20, and 16 times, respectively. From the empirical observation, we find that the
282 PPO gets stuck in locally optimal policies, failing to find a consistent one. As expected, PPO+H can
283 converge to the stationary policy with a substantially higher ratio, which verifies the effectiveness of
284 the proposed H-term in finding a physically stationary policy.

285 5.5 Impact of Trajectory Length K

286 We investigate the impact of trajectory length K that defines the horizon of the H-term. Based on (17),
287 we foresee that a large K means a more accurate estimation of $\nabla_{\theta} H(\theta)$ but at a price of computations.
288 Here, we evaluate the PPO+H with $K = 8, 16$ and set the size of replay buffer \mathcal{D}_2 to 1, 000. In Table
289 2, we observe that the cumulative reward increases and the variance decreases as K increases from 8
290 to 16. However, for the case $K = 24$, both metrics get worse due to the out-of-memory issue. For
291 $K = 24$, we reduce the replay buffer size to 800 to meet the memory limit. The smaller replay buffer
292 size hurts the diversity of the trajectories and may lead to a small performance drop. [This hypothesis](#)
293 [is verified in Appx. H.2, where all trials use a consistent replay buffer size 800.](#)

294 **6 Conclusions**

295 In this paper, we have addressed a foundational issue of existing deep reinforcement learning
296 algorithms that **Bellman’s optimality equation has multiple fixed points**. This issue leads to the
297 instability of DRL algorithms, puts a challenge on their reliability and reproducibility, and thus
298 limits the wider adoption in real-world tasks. As a fix of the problem, we propose a novel H-term, a
299 physically inspired regularizer, by making a novel analogy between a MDP and a quantum K-spin
300 Ising model. Experimentally, we show that the H-term helps DRL algorithms find a stationary
301 policy that has the lowest energy in a system and reduces the variance of cumulative rewards by
302 65.2% ~ 85.6% compared with those of existing algorithms.

303 For future works, we will explore the potential of directly training a policy network using (7) as in
304 Appx. I, quantum simulator [25] and quantum reinforcement learning [5][16]. It is interesting to
305 apply Monte Carlo estimator for unbiased policy gradient calculations. We would like to show that
306 the proposed H-term can help distributional RL algorithms [3] find a stationary policy, since the
307 distributional Bellman optimality operator is not a contraction and thus there is also no unique policy.

308 **Broader Impact Statement**

309 In the field of DRL research, this paper points out the issue of multiple fixed points and may attract
310 more attention from various perspectives. In terms of applications, this paper moves toward the
311 reliable and reproducible research by improving the stability of existing DRL algorithms. On the other
312 hand, the obtained stationary policy may have broad practical impact in many real-world application
313 areas, including but not limited to robotics, transportation, and finance.

314 From an interdisciplinary perspective, our approach lies at the intersection of (quantum) Hamiltonian
315 mechanics and DRL, especially the explicit analogy between an MDP and a quantum K-spin Ising
316 model. We hope that our study will attract more attention to bring together the strengths of both
317 approaches and yield new insights in both fields.

318 In terms of broader societal impact of this work, we do not see any foreseeable strongly negative
319 impacts. However, this work essentially makes a tradeoff between the computational resource and
320 stability, which may lead to future works with higher computational cost.

321 **References**

- 322 [1] Rishabh Agarwal, Max Schwarzer, Pablo Samuel Castro, Aaron C Courville, and Marc Belle-
323 mare. Deep reinforcement learning at the edge of the statistical precipice. *Advances in Neural*
324 *Information Processing Systems*, 2021.
- 325 [2] Oron Anschel, Nir Baram, and Nahum Shimkin. Averaged-DQN: Variance reduction and
326 stabilization for deep reinforcement learning. In *International Conference on Machine Learning*,
327 pages 176–185. PMLR, 2017.
- 328 [3] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforce-
329 ment learning. In *International Conference on Machine Learning*, pages 449–458. PMLR,
330 2017.
- 331 [4] Dimitri Bertsekas. *Reinforcement Learning and Optimal Control*. Athena Scientific, 2019.
- 332 [5] Jacob Biamonte, Peter Wittek, Nicola Pancotti, Patrick Rebentrost, Nathan Wiebe, and Seth
333 Lloyd. Quantum machine learning. *Nature*, 549(7671):195–202, 2017.
- 334 [6] Stephanie CY Chan, Samuel Fishman, Anoop Korattikara, John Canny, and Sergio Guadarrama.
335 Measuring the reliability of reinforcement learning algorithms. 2019.
- 336 [7] Tianyi Chen, Yuejiao Sun, and Wotao Yin. Closing the gap: Tighter analysis of alternating
337 stochastic gradient methods for bilevel problems. *Advances in Neural Information Processing*
338 *Systems*, 34, 2021.
- 339 [8] Xinyue Chen, Che Wang, Zijian Zhou, and Keith Ross. Randomized ensembled double q-
340 learning: Learning fast without a model. *ICLR*, 2021.
- 341 [9] Richard Cheng, Abhinav Verma, Gabor Orosz, Swarat Chaudhuri, Yisong Yue, and Joel Burdick.
342 Control regularization for reduced variance reinforcement learning. In *International Conference*
343 *on Machine Learning*, pages 1141–1150. PMLR, 2019.
- 344 [10] Jack Choquette, Wishwesh Gandhi, Olivier Giroux, Nick Stam, and Ronny Krashinsky. NVIDIA
345 A100 Tensor Core GPU: Performance and innovation. *IEEE Micro*, 41(2):29–35, 2021.
- 346 [11] Barry A Cipra. The ising model is np-complete. *SIAM News*, 33(6):1–3, 2000.
- 347 [12] Vasil S Denchev, Sergio Boixo, Sergei V Isakov, Nan Ding, Ryan Babbush, Vadim Smelyanskiy,
348 John Martinis, and Hartmut Neven. What is the computational value of finite-range tunneling?
349 *Physical Review X*, 6(3):031015, 2016.
- 350 [13] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking
351 deep reinforcement learning for continuous control. In *International Conference on Machine*
352 *Learning*, pages 1329–1338. PMLR, 2016.

- 353 [14] Gabriel Dulac-Arnold, Nir Levine, Daniel J Mankowitz, Jerry Li, Cosmin Paduraru, Sven Gowal,
354 and Todd Hester. An empirical investigation of the challenges of real-world reinforcement
355 learning. *arXiv preprint arXiv:2003.11881*, 2020.
- 356 [15] Gabriel Dulac-Arnold, Daniel Mankowitz, and Todd Hester. Challenges of real-world reinforce-
357 ment learning. 2019.
- 358 [16] Vedran Dunjko, Jacob M Taylor, and Hans J Briegel. Quantum-enhanced machine learning.
359 *Physical review letters*, 117(13):130501, 2016.
- 360 [17] Eliot Kapit Wesley Jones Eric B. Jones, Peter Graf. K-spin Hamiltonian for quantum-resolvable
361 markov decision processes. *ArXiv*, abs/2003.11881, 2020.
- 362 [18] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you
363 need: Learning skills without a reward function. *ICLR*, 2019.
- 364 [19] Evan Greensmith, Peter Bartlett, and Jonathan Baxter. Variance reduction techniques for
365 gradient estimates in reinforcement learning. *Advances in Neural Information Processing*
366 *Systems*, 14, 2001.
- 367 [20] Evan Greensmith, Peter L Bartlett, and Jonathan Baxter. Variance reduction techniques for
368 gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5(9),
369 2004.
- 370 [21] Samuel Greydanus, Misko Dzamba, and Jason Yosinski. Hamiltonian neural networks. *Advances*
371 *in Neural Information Processing Systems*, 32, 2019.
- 372 [22] Agrim Gupta, Silvio Savarese, Surya Ganguli, and Li Fei-Fei. Embodied intelligence via
373 learning and evolution. *Nature Communications*, 12(1):1–12, 2021.
- 374 [23] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
375 policy maximum entropy deep reinforcement learning with a stochastic actor. In *International*
376 *Conference on Machine Learning*, pages 1861–1870. PMLR, 2018.
- 377 [24] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David
378 Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on*
379 *artificial intelligence*, volume 32, 2018.
- 380 [25] Eric B Jones, Peter Graf, Eliot Kapit, and Wesley Jones. K-spin Hamiltonian for quantum-
381 resolvable markov decision processes. *Quantum Machine Intelligence*, 2(2):1–11, 2020.
- 382 [26] Scott Kirkpatrick, C Daniel Gelatt Jr, and Mario P Vecchi. Optimization by simulated annealing.
383 *Science*, 220(4598):671–680, 1983.
- 384 [27] Timothy P Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval
385 Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning.
386 *ICLR*, 2016.
- 387 [28] Nicolas Loizou, Hugo Berard, Alexia Jolicoeur-Martineau, Pascal Vincent, Simon Lacoste-
388 Julien, and Ioannis Mitliagkas. Stochastic Hamiltonian gradient methods for smooth games. In
389 *International Conference on Machine Learning*, pages 6370–6381. PMLR, 2020.
- 390 [29] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
391 Wierstra, and Martin Riedmiller. Playing Atari with deep reinforcement learning. *ICLR*, 2013.
- 392 [30] Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte carlo gradient
393 estimation in machine learning. *J. Mach. Learn. Res.*, 21(132):1–62, 2020.
- 394 [31] Alexey B Piunovskiy. *Examples in Markov Decision Processes*, volume 2. World Scientific,
395 2013.
- 396 [32] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay.
397 *ICLR*, 2016.

- 398 [33] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-
399 dimensional continuous control using generalized advantage estimation. *ICLR*, 2016.
- 400 [34] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
401 policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- 402 [35] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
403 Deterministic policy gradient algorithms. In *International Conference on Machine Learning*,
404 pages 387–395. PMLR, 2014.
- 405 [36] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT press,
406 2018.
- 407 [37] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient meth-
408 ods for reinforcement learning with function approximation. *Advances in Neural Information*
409 *Processing Systems*, 12, 1999.
- 410 [38] Pierre Thodoroff, Audrey Durand, Joelle Pineau, and Doina Precup. Temporal regularization
411 for markov decision process. *Advances in Neural Information Processing Systems*, 31, 2018.
- 412 [39] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based
413 control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
414 5026–5033. IEEE, 2012.
- 415 [40] Peter Toth, Danilo Jimenez Rezende, Andrew Jaegle, Sébastien Racanière, Aleksandar Botev,
416 and Irina Higgins. Hamiltonian generative networks. *ICLR*, 2019.
- 417 [41] Junfeng Wen, Saurabh Kumar, Ramki Gummadi, and Dale Schuurmans. Characterizing the gap
418 between actor-critic and policy gradient. In *International Conference on Machine Learning*,
419 pages 11101–11111. PMLR, 2021.
- 420 [42] Cathy Wu, Aravind Rajeswaran, Yan Duan, Vikash Kumar, Alexandre M Bayen, Sham Kakade,
421 Igor Mordatch, and Pieter Abbeel. Variance reduction for policy gradient with action-dependent
422 factorized baselines. *ICLR*, 2018.
- 423 [43] Duo Xu and Faramarz Fekri. Improving actor-critic reinforcement learning via Hamiltonian
424 Monte Carlo method. *Deep Reinforcement Learning Workshop at NeurIPS*, 2021.

425 Checklist

- 426 1. For all authors...
- 427 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
428 contributions and scope? [\[Yes\]](#)
- 429 (b) Did you describe the limitations of your work? [\[Yes\]](#) More computational cost.
- 430 (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) May lead
431 to future works with higher computational cost.
- 432 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
433 them?
- 434 2. If you are including theoretical results...
- 435 (a) Did you state the full set of assumptions of all theoretical results? [\[Yes\]](#)
- 436 (b) Did you include complete proofs of all theoretical results? [\[Yes\]](#) In the appendix.
- 437 3. If you ran experiments...
- 438 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
439 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
- 440 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
441 were chosen)? [\[Yes\]](#)
- 442 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
443 ments multiple times)? [\[Yes\]](#) We reported the variance over 20 random seeds as shaded
444 area in the figures, and also list it in a table.

- 445 (d) Did you include the total amount of compute and the type of resources used (e.g., type
446 of GPUs, internal cluster, or cloud provider)? [N/A]
- 447 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 448 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 449 (b) Did you mention the license of the assets? [N/A]
- 450 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 451 (d) Did you discuss whether and how consent was obtained from people whose data you're
452 using/curating? [N/A]
- 453 (e) Did you discuss whether the data you are using/curating contains personally identifiable
454 information or offensive content? [N/A]
- 455 5. If you used crowdsourcing or conducted research with human subjects...
- 456 (a) Did you include the full text of instructions given to participants and screenshots, if
457 applicable? [N/A]
- 458 (b) Did you describe any potential participant risks, with links to Institutional Review
459 Board (IRB) approvals, if applicable? [N/A]
- 460 (c) Did you include the estimated hourly wage paid to participants and the total amount
461 spent on participant compensation? [N/A]

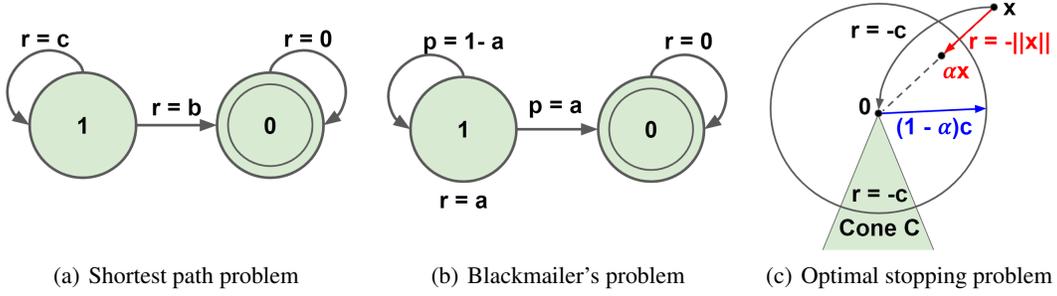


Figure 5: Revisiting Fig. 1 for the discounted cases where $\gamma \in (0, 1)$.

462 A More Examples with Multiple Fixed Points

463 First, we consider the discounted formulations of the three examples (shown in Fig. 1), as shown in
 464 Fig. 5 where $\gamma \in (0, 1)$. The differences are marked in red.

- 465 • (a) **Shortest path problem (deterministic, discounted case)**: Given two states 1 and 0, an agent
 466 at state 1 transits to either state 1 or 0 with rewards $r = c$ or $r = b$, respectively. $c > (1 - \gamma) \cdot b$.
 467 At state 0, the value function is $V(0) = 0$. At state 1, the Bellman's optimality equation is
 468 $V(1) = \max_a \{c + \gamma \cdot V(1), b\}$, where any $V(1) \geq (b - c)/\gamma$ is a solution. If initialize $V_0(1) \geq b$,
 469 an agent obtains a policy that always transits back to state 1; otherwise, a result policy drives to
 470 terminal state 0.
- 471 • (b) **Blackmailer's problem (stochastic, discounted case)**: Different from (a), a profit maximizing
 472 blackmailer/agent at 1 demands a cash amount $a \in (0, 1]$ (an action), while a victim transits to state
 473 1 with probability a or to state 0 with probability $1 - a$, respectively. At state 0, a victim always
 474 refuses to yield to the blackmailer's demand, i.e., $V(0) = 0$. The Bellman's optimality equation
 475 is $V(1) = \max_a \{a + \gamma \cdot (1 - a)V(1)\}$ for state 1, where any $V(1) \geq 1$ is a feasible solution. If
 476 initialize $V_0(1) = c > 1$, the blackmailer's policy is demanding $a \rightarrow 0$ at the k -th step to keep
 477 the victim stay at state 1, for any $k \leq K_0 = -\lfloor \log_\gamma c \rfloor$; and taking $a = 1$ to transit to terminal
 478 state 0 at the k -th step, for any $k \geq K_0 + 1$; otherwise initialize $V_0(1) = c \leq 1$, the result policy is
 479 demanding the maximum $a = 1$ that drives the victim to a refusal state 0 (a terminal state).
- 480 • (c) **Optimal stopping problem (terminating policies, discounted case)**: In a space \mathbb{R}^2 with
 481 terminating state at point 0, at point $x \neq 0$ an agent moves to either point 0 with negative reward $-c$
 482 or point αx with reward $-||x||$, respectively, where $\alpha \in (0, 1)$. The Bellman's optimality equation
 483 is $V(x) = \max\{-c, -||x|| + \gamma \cdot V(\alpha x)\}$ and the optimal policy is to continue inside the sphere of
 484 radius $(1 - \alpha)c$ and to stop outside. If add a cone region C within which an agent always receives
 485 a reward $-c$, a second policy is jumping to point 0 at any point in region C .

486 Then, we elaborate how the proposed H-term fixes the problems in Fig. 5.

487 (a) Shortest path problem (deterministic, discounted case)

488 Assume $V_0(1) \geq b$ and $c > (1 - \gamma)b$, we have

$$\begin{aligned}
 V_1(1) &= c + \gamma \cdot V_0(1) \geq c + \gamma \cdot b > b \\
 V_2(1) &= c + \gamma \cdot c + \gamma^2 \cdot V_0(1) \geq (1 + \gamma)c + \gamma^2 b > b \\
 V_3(1) &= c + \gamma \cdot c + \gamma^2 c + \gamma^3 \cdot V_0(1) \geq (1 + \gamma + \gamma^2)c + \gamma^3 b > b \\
 &\dots \\
 V_k(1) &= \sum_{i=0}^{k-1} \gamma^i \cdot c + \gamma^k \cdot V_0(1) \geq \sum_{i=0}^{k-1} \gamma^i \cdot c + \gamma^k b > b \\
 &\dots \\
 V^*(1) &= \sum_{i=0}^{\infty} \gamma^i \cdot c = \frac{1}{1 - \gamma} c > b
 \end{aligned} \tag{19}$$

489 The values of $H(0)$ and $H(1)$ are as follows:

$$H(0) = 0, \quad H(1) = -b - \sum_{k=2}^{\infty} \left(\sum_{i=1}^{k-1} \gamma^{i-1} \cdot c + \gamma^k b \right) = -\infty. \quad (20)$$

490 Adding the above H-values to state 1 and 0, respectively, we have

$$\begin{aligned} V^*(1) + H(1) &= \sum_{i=0}^{\infty} \gamma^i \cdot c - \infty = -\infty \\ V^*(0) + H(0) &= b. \end{aligned} \quad (21)$$

491 Therefore, $V^*(1) + H(1) < V^*(0) + H(0)$, independent of the initial value $V_0(1)$. That is, an agent
492 always obtains a policy that drives to terminal state 0 at step 1.

493 **(b) Blackmailer's problem (stochastic, discounted case)**

494 If initialize $V_0(1) = c > 1$, the blackmailer's policy is demanding $a \rightarrow 0$ at the k -th step to keep the
495 victim stay at state 1, for any $k \leq K_0 = -\lfloor \log_{\gamma} c \rfloor$; and taking $a = 1$ to transit to terminal state 0 at
496 the k -th step, for any $k \geq K_0 + 1$; otherwise initialize $V_0(1) = c \leq 1$, the result policy is demanding
497 the maximum $a = 1$ that drives the victim to a refusal state 0 (a terminal state).

498 The values of $H(0)$ and $H(1)$ are as follows:

$$H(0) = 0, \quad H(1) = - \sum_{k=1}^{\infty} \sum_{i=1}^{k-1} \gamma^{i-1} \cdot a = -\infty. \quad (22)$$

499 For arbitrary initial value of $V_0(1)$, $V_1(1) = a + (1 - a) \cdot \gamma(V_0(1) + H(1))$ take maximum $V_1(1) = 1$
500 when $a = 1$. Therefore, the policy always drives to terminal state 0 at step 1.

501 **(c) Optimal stopping problem (terminating policies, discounted case)**

502 Any policy that takes infinite steps will have

$$H(x) = -c - \sum_{k=2}^{\infty} \left[\sum_{i=1}^{k-1} \gamma^i \cdot \alpha^i \cdot \|x\| + \gamma^k \cdot (-c) \right] = -\infty \quad (23)$$

503 and a direct jumping policy will have $H(x) = -c$. Therefore, the H-term drives to a terminating
504 policy.

505 **B MuJoCo Tasks with Multiple Policies**

506 **B.1 Description of MuJoCo Taks**

507 We selected six challenging robotic locomotion tasks from MuJoCo, namely, Swimmer-v3, Hopper-
 508 v3, Walker2D-v3, HalfCheetah-v3, Ant-v3, Humanoid-v3. Table 3 lists the action space and state
 509 space of each task.

Table 3: The state and action spaces of six challenging MuJoCo tasks.

Tasks	Agent	Action Space	State Space
Swimmer-v3	Three-link swimming robot	2	8
Hopper-v3	Two-dimensional one-legged robot	3	11
Walker2d-v3	Two-dimensional bipedal robot	6	17
HalfCheetah-v3	Two-dimensional robot	6	17
Ant-v3	Four-legged creature	8	111
Humanoid-v3	Three-dimensional bipedal robot	17	376

510 **B.2 Multiple policies in MuJoCo tasks**

511 In the supplementary files, we includes rendered videos of different policies, as given in Table 4.

- 512 • Different policies are obtained over 20 runs of the PPO algorithm. We rendered theses polices and
 513 classified them by physical gaits.
 514 • The policies in bold texts are physically stationary.

Table 4: List of video files for different policies.

Task	Different Policies	Video Name
Hopper	hopping	hopper_hopping.mp4
	diving	hopper_diving.mp4
	standing	hopper_standing.mp4
Ant	running	ant_running.mp4
	standing	ant_standing.mp4
	flipping	ant_flipping.mp4
Walker	walking	walker_walking.mp4
	diving	walker_diving.mp4
	standing	walker_standing.mp4
Humanoid	two-legs	humanoid_two_legs.mp4
	one-leg	humanoid_one_leg.mp4
	backward	humanoid_backward.mp4
HalfCheetah	running	halfcheetah_running.mp4
	diving	halfcheetah_diving.mp4
	flipping	halfcheetah_flipping.mp4
	standing	halfcheetah_standing.mp4
Swimmer	moving	swimmer_moving.mp4
	standing	swimmer_standing.mp4

515 **C Quantum K-Spin Hamiltonian Formulation of Reinforcement Learning**

516 We provide the detailed steps of reformulating (1) into a K -spin Hamiltonian equation

$$\begin{aligned}
H(\theta) &\triangleq -\mathbb{E}_{S_0, A_0} [Q^{\pi_\theta}(S_0, A_0)] \\
&= -\mathbb{E}_{S_0, A_k \sim \pi_\theta(S_k, \cdot), S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k)} \left[\sum_{k=0}^{\infty} \gamma^k \cdot R(S_k, A_k) \right] \\
&= -\sum_{k=0}^{K-1} \mathbb{E}_{S_0, A_0, \dots, S_k \sim \mathbb{P}(\cdot | S_{k-1}, A_{k-1}), A_k \sim \pi_\theta(S_k, \cdot)} [\gamma^k \cdot R(S_k, A_k)] \\
&= -\sum_{k=0}^{K-1} \sum_{\mu_0}^{S \times A} \dots \sum_{\mu_k}^{S \times A} \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \pi_\theta(\mu_0) \prod_{i=0}^{k-1} [\mathbb{P}(S_{i+1} | \mu_i) \cdot \pi_\theta(\mu_{i+1})] \quad (24) \\
&= -\sum_{k=0}^{K-1} \sum_{\mu_0}^{S \times A} \dots \sum_{\mu_k}^{S \times A} \left[\gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i) \right] \cdot \pi_\theta(\mu_0) \dots \pi_\theta(\mu_k) \\
&= -\sum_{k=0}^{K-1} \sum_{\mu_0}^{S \times A} \dots \sum_{\mu_k}^{S \times A} L_{\mu_0, \dots, \mu_k} \cdot \pi_\theta(\mu_0) \dots \pi_\theta(\mu_k),
\end{aligned}$$

517 where $K \rightarrow \infty$, and the density function is

$$L_{\mu_0, \dots, \mu_k} = \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i). \quad (25)$$

Table 5: Revisiting the analogy between MDP and quantum K-spin Ising model.

MDP (Our formulation in (7))		Quantum K-spin Ising Model [26, 12] in (5)	
State-action pairs	μ_0, \dots, μ_{K-1}	Spins	j_0, \dots, j_{K-1}
Optimal policy	$\pi_{\mu_0}^* \times \pi_{\mu_1}^* \times \dots \times \pi_{\mu_{K-1}}^*$	Quantum field	$\sigma_{j_0} \times \sigma_{j_1} \times \dots \times \sigma_{j_{K-1}}$
Cumulative rewards	$L_{\mu_0 \dots \mu_{K-1}}$	Density function	$L_{j_0 \dots j_{K-1}}$
Functional of policy	$H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}})$	Functional of spins	$H(\sigma_{j_0}, \dots, \sigma_{j_{K-1}})$
Stationary condition	$\frac{\delta H(\pi_{\mu_0}, \dots, \pi_{\mu_{K-1}})}{\delta \pi_\mu} = 0$	Stationary condition	$\frac{\delta H(\sigma_{j_0}, \dots, \sigma_{j_{K-1}})}{\delta \sigma_j} = 0$

518 D Derivation Steps for Section 4.2: Hamiltonian's Policy Gradients

519 We provide the policy gradient of the quantum K-spin Hamiltonian equation in (7) for both stochastic
520 and deterministic cases, which are variants of the well-known policy gradient theorem [36].

521 **Theorem 1. (Hamiltonian's stochastic policy gradient)** *The stochastic gradient of the K-spin*
522 *Hamiltonian equation (7) w.r.t. parameter θ is*

$$\nabla_{\theta} H(\theta) = -\mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_{\theta} \log (\pi_{\theta}(\mu_0) \cdot \pi_{\theta}(\mu_1) \cdots \pi_{\theta}(\mu_k)) \right]. \quad (26)$$

523 **Corollary 1.** *When $K \rightarrow \infty$, the Hamiltonian's stochastic policy gradient $\nabla_{\theta} H(\theta)$ in (26) is equal*
524 *to the stochastic policy gradient $\nabla_{\theta} J(\theta)$ in [37],*

$$\lim_{K \rightarrow \infty} \nabla_{\theta} H(\theta) = -\nabla_{\theta} J(\theta) = -\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [Q^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)]. \quad (27)$$

525 Let $\eta_{\theta}(\cdot) : \mathcal{S} \rightarrow \mathcal{A}$ denote a deterministic policy, while we use $\tilde{\pi}_{\theta, \delta}(\mu)$ to represent that a Gaussian
526 noise (a.k.a, an exploration noise) with standard deviation $\delta > 0$ is added in the exploration process.

527 **Theorem 2. (Hamiltonian's deterministic policy gradient)** *The deterministic gradient of the K-spin*
528 *Hamiltonian equation (7) w.r.t. parameter θ is*

$$\nabla_{\theta} H'(\theta) = -\mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_{\theta} \log (\tilde{\pi}_{\theta, \delta}(\mu_0) \cdot \tilde{\pi}_{\theta, \delta}(\mu_1) \cdots \tilde{\pi}_{\theta, \delta}(\mu_k)) \right]. \quad (28)$$

529 **Corollary 2.** *When $K \rightarrow \infty$, the Hamiltonian's deterministic policy gradient $\nabla_{\theta} H'(\theta)$ in (28) is*
530 *equal to the deterministic policy gradient $\nabla_{\theta} J'(\theta)$ in [35],*

$$\lim_{K \rightarrow \infty} \nabla_{\theta} H'(\theta) = -\nabla_{\theta} J'(\theta) = -\mathbb{E}_{s \sim d_{\theta}} \left[\nabla_a Q^{\tilde{\pi}_{\theta, \delta}}(s, a) |_{a=\eta_{\theta}} \nabla_{\theta} \eta_{\theta}(s) \right]. \quad (29)$$

531 **Corollary 3.** *When the variance of the exploration noise approaches zero, i.e., $\delta \rightarrow 0$, the determin-*
532 *istic policy gradient $\nabla_{\theta} H'(\theta)$ is the limiting case of the stochastic policy gradient $\nabla_{\theta} H(\theta)$,*

$$\nabla_{\theta} H'(\theta) = \lim_{\delta \rightarrow 0} \nabla_{\theta} H(\theta). \quad (30)$$

533 D.1 Proof of Theorem 1: Hamiltonian's Stochastic Policy Gradient

Proof.

$$\begin{aligned} \nabla_{\theta} H(\theta) &= - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} \nabla_{\theta} [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \\ &= - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \nabla_{\theta} \log [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \\ &= - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \pi_{\theta}(\mu_0) \prod_{i=0}^{k-1} [\mathbb{P}(S_{i+1} | \mu_i) \cdot \pi_{\theta}(\mu_{i+1})] \cdot \nabla_{\theta} \log [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \\ &= -\mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_{\theta} \log [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \right], \end{aligned} \quad (31)$$

534 where $\mu_k = (S_k, A_k)$, $S_0 \sim d_0(\cdot)$, $A_k \sim \pi_{\theta}(S_k, \cdot)$, $S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k)$ for $k = 0 \cdots K$. \square

535 **D.2 Proof of Corollary 1**

Proof.

$$\begin{aligned}
\nabla_{\theta} H(\theta) &\stackrel{(a)}{=} - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} \nabla_{\theta} [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)] \\
&\stackrel{(b)}{=} - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} \sum_{i=0}^k \pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_{i-1}) \pi_{\theta}(\mu_{i+1}) \cdots \pi_{\theta}(\mu_k) \nabla_{\theta} \pi_{\theta}(\mu_i) \\
&\stackrel{(c)}{=} - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i) \sum_{i=0}^k \left[\prod_{j=0}^{i-1} \pi_{\theta}(\mu_j) \cdot \nabla_{\theta} \pi_{\theta}(\mu_i) \cdot \prod_{j=i+1}^k \pi_{\theta}(\mu_j) \right] \\
&\stackrel{(d)}{=} - \sum_{k=0}^{K-1} \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} d_0(S_0) \sum_{i=0}^k \left[\gamma^i \prod_{j=0}^{i-1} \pi_{\theta}(\mu_j) \mathbb{P}(S_{j+1} | \mu_j) \right] \nabla_{\theta} \pi_{\theta}(\mu_i) \left[\prod_{j=i+1}^{k-1} \pi_{\theta}(\mu_j) \mathbb{P}(S_{j+1} | \mu_j) \pi_{\theta}(\mu_k) \gamma^{k-i} R(\mu_k) \right] \\
&\stackrel{(e)}{=} - \sum_{k=0}^{K-1} \sum_{i=0}^k \sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_{S_i}^{\mathcal{S}} \rho(S_0, S_i, i) \sum_{A_i}^{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(S_i, A_i) \cdot \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} \rho(S_i, S_k, k-i) \cdot \pi_{\theta}(\mu_k) \cdot R(\mu_k) \\
&\stackrel{(f)}{=} - \sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_S^{\mathcal{S}} \sum_{i=0}^{K-1} \rho(S_0, S, i) \sum_A^{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(S, A) \cdot \left[\sum_{S'}^{\mathcal{S}} \sum_{k=i}^{K-1} \rho(S, S', k-i) \cdot \sum_{A'}^{\mathcal{A}} \pi_{\theta}(S', A') \cdot R(S', A') \right] \\
&\stackrel{(g)}{=} - \sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_S^{\mathcal{S}} \sum_{i=0}^{\infty} \rho(S_0, S, i) \sum_A^{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(S, A) \cdot Q^{\pi_{\theta}}(S, A) \\
&\stackrel{(h)}{=} - \left[\sum_S^{\mathcal{S}} \sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_{i=0}^{\infty} \rho(S_0, S, i) \right] \cdot \sum_S^{\mathcal{S}} \frac{\sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_{i=0}^{\infty} \rho(S_0, S, i)}{\sum_S^{\mathcal{S}} \sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_{i=0}^{\infty} \rho(S_0, S, i)} \sum_A^{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(S, A) \cdot Q_{\theta}(S, A) \\
&\stackrel{(i)}{\propto} - \sum_S^{\mathcal{S}} d_{\pi_{\theta}}(S) \sum_A^{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(S, A) \cdot Q^{\pi_{\theta}}(S, A) \\
&\stackrel{(j)}{=} - \mathbb{E}_{S \sim d_{\theta}, A \sim \pi_{\theta}(S, \cdot)} [Q^{\pi_{\theta}}(S, A) \nabla_{\theta} \log \pi_{\theta}(S, A)],
\end{aligned} \tag{32}$$

536 where $\rho(S, S', i)$ denotes the probability of state S transfer to S' in i steps.

537 We provide detailed explanations step-by-step:

538 • Equality (a) holds by definition.

539 • In equality (b), using the chain rule, we take derivative of $\nabla_{\theta} [\pi_{\theta}(\mu_0) \cdots \pi_{\theta}(\mu_k)]$ with respect to
540 $\pi_{\theta}(\mu_i)$, $i = 1, \dots, k$.

541 • In equality (c), we plug in L_{μ_0, \dots, μ_k} in (6).

542 • In equality (d), we insert $\mathbb{P}(S_{i+1} | \mu_i) \mathbb{P}(S_{i+1} | \mu_i)$ between $\pi_{\theta}(\mu_i)$ and $\pi_{\theta}(\mu_{i+1})$, $i = 1, \dots, k$.

543 • In equality (e), we split trajectory $\mu_0, \dots, \mu_i, \dots, \mu_k$ into two trajectories μ_0, \dots, μ_i and
544 μ_i, \dots, μ_k . Therefore, we can classify all trajectories μ_0, \dots, μ_k by μ_0, μ_i, μ_k , and i .

545 • In equality (f), we reorganize $\sum_{k=0}^{K-1} \sum_{i=0}^k$ into $\sum_{i=0}^{K-1} \sum_{k=i}^{K-1}$. The former one first traverses the
546 length k of a trajectory, and then traverses the i -th step on it. The latter one first traverses the i -th
547 step of a trajectory, and then traverses the length k of it.

548 • In equality (g), we calculate the limit of (f) when K approaches ∞ .

549 • In equality (h), we normalize $\sum_{S_0}^{\mathcal{S}} d_0(S_0) \sum_{i=0}^{\infty} \rho(S_0, S, i)$ to be a probability distribution.

550 • In equality (i), we remove the constant $\sum_S \sum_{S_0}^S d_0(S_0) \sum_{i=0}^{\infty} \rho(S_0, S, i)$ and replace the fraction
 551 with $d_{\pi_\theta}(S)$, the stationary distribution of state S under policy π_θ .

552 • In equality (j), we reformulate (i) as expectation.

553

□

554 D.3 Proof of Theorem 2: Hamiltonian's Deterministic Policy Gradient

555 *Proof.* Let $\eta_\theta(\cdot) : \mathcal{S} \rightarrow \mathcal{A}$ denote a deterministic policy, while we use $\tilde{\pi}_{\theta,\delta}(\mu)$ to represent that a
 556 Gaussian noise (a.k.a, an exploration noise) with standard deviation $\delta > 0$ is added in the exploration
 557 process. In the inference stage, there is no exploration noise, the policy is deterministic, i.e., $\delta = 0$
 558 and $A_k = \eta_\theta(S_k)$.

$$\begin{aligned}
 H'(\theta) &\triangleq -\mathbb{E}_{S_0 \sim d_0, A_0 \sim \tilde{\pi}_{\theta,\delta}} \left[Q^{\tilde{\pi}_{\theta,\delta}}(S_0, A_0) \right] \\
 &= -\mathbb{E}_{S_0, A_k \sim \tilde{\pi}_{\theta,\delta}(S_k, \cdot), S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k)} \left[\sum_{k=0}^{\infty} \gamma^k \cdot R(S_k, A_k) \right] \\
 &= -\sum_{k=0}^K \mathbb{E}_{S_0, A_k \sim \tilde{\pi}_{\theta,\delta}(S_k, \cdot), S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k)} \left[\gamma^k \cdot R(S_k, A_k) \right] \\
 &= -\sum_{k=0}^K \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \tilde{\pi}_{\theta,\delta}(\mu_0) \prod_{i=0}^{k-1} [\mathbb{P}(S_{i+1} | \mu_i) \cdot \tilde{\pi}_{\theta,\delta}(\mu_{i+1})] \quad (33) \\
 &= -\sum_{k=0}^K \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} \left[\gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i) \right] \cdot \tilde{\pi}_{\theta,\delta}(\mu_0) \cdots \tilde{\pi}_{\theta,\delta}(\mu_k) \\
 &= -\sum_{k=0}^K \sum_{\mu_0}^{\mathcal{S} \times \mathcal{A}} \cdots \sum_{\mu_k}^{\mathcal{S} \times \mathcal{A}} L_{\mu_0, \dots, \mu_k} \cdot \tilde{\pi}_{\theta,\delta}(\mu_0) \cdots \tilde{\pi}_{\theta,\delta}(\mu_k),
 \end{aligned}$$

559 where $K \rightarrow \infty$, and

$$L_{\mu_0, \dots, \mu_k} = \gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i). \quad (34)$$

560

□

561 **D.4 Proof of Corollary 2**

Proof.

$$\begin{aligned}
\nabla_{\theta} H'(\pi_{\theta}) &= - \sum_{k=0}^K \sum_{\mu_0}^{S \times A} \cdots \sum_{\mu_k}^{S \times A} (L_{\mu_0, \dots, \mu_k} \cdot \nabla_{\theta} [\tilde{\pi}_{\theta}(\mu_0) \cdots \tilde{\pi}_{\theta}(\mu_k)] + \nabla_{\theta} L_{\mu_0, \dots, \mu_k} \cdot \tilde{\pi}_{\theta}(\mu_0) \cdots \tilde{\pi}_{\theta}(\mu_k)) \\
&= - \sum_{k=0}^K \sum_{\mu_0}^{S \times A} \cdots \sum_{\mu_k}^{S \times A} [\tilde{\pi}_{\theta}(\mu_0) \cdots \tilde{\pi}_{\theta}(\mu_k)] \cdot \nabla_{\theta} L_{\mu_0, \dots, \mu_k} \\
&= - \sum_{k=0}^K \sum_{\mu_0}^{S \times A} \cdots \sum_{\mu_k}^{S \times A} \nabla_{\theta} \left[\gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i) \right] \\
&= - \sum_{k=0}^K \sum_{\mu_0}^{S \times A} \cdots \sum_{\mu_k}^{S \times A} \nabla_A \left[\gamma^k \cdot R(\mu_k) \cdot d_0(S_0) \cdot \prod_{i=0}^{k-1} \mathbb{P}(S_{i+1} | \mu_i) \right] \nabla_{\theta} \eta_{\theta}(S) \\
&= - \sum_{S_0}^S d_0(S_0) \nabla_A \mathbb{E}_{S_{t+1} \sim \mathbb{P}(\cdot | S_t, A_t)} \left[\sum_{t=0}^{\infty} \gamma^k R(S_t, A_t) \right] \cdot \nabla_{\theta} \eta_{\theta}(S) \\
&= - \sum_{S_0}^S d_0(S_0) \nabla_A Q(S_0, A_0) \cdot \nabla_{\theta} \eta_{\theta}(S) \\
&= - \mathbb{E}_{S_0 \sim d_0(\cdot)} [\nabla_A Q(S_0, A_0) \cdot \nabla_{\theta} \eta_{\theta}(S)]
\end{aligned} \tag{35}$$

562 where $\mu_k = (S_k, A_k)$, $S_0 \sim d_0(\cdot)$, $A_k \sim \pi_{\theta}(S_k, \cdot)$, $S_{k+1} \sim \mathbb{P}(\cdot | S_k, A_k)$, for $k = 0 \cdots K$. \square

563 **D.5 Proof of Corollary 3**

564 *Proof.* In Corollary 2 and Corollary 1, we have

$$\begin{aligned}
\nabla_{\theta} H'(\theta) &= -\nabla_{\theta} J'(\theta), \\
\nabla_{\theta} H(\theta) &= -\nabla_{\theta} J(\theta),
\end{aligned} \tag{36}$$

565 when $K \rightarrow \infty$.

566 [35] proved that

$$\nabla_{\theta} J'(\theta) = \lim_{\delta \rightarrow 0} \nabla_{\theta} J(\theta), \tag{37}$$

567 where δ is the standard deviation of the Gaussian noise of stochastic policy π_{θ} .

568 Therefore,

$$\nabla_{\theta} H'(\theta) = \lim_{\delta \rightarrow 0} \nabla_{\theta} H(\theta) \tag{38}$$

569 \square

570 **E Variance Reduction (Newly Added)**

571 For the general function in (8), one simple but effective variance reduction technique is to subtract a
572 baseline term as follows:

$$\mathbb{E}_{p(\mathbf{x};\theta)} [(f(\mathbf{x}) - \beta)\nabla_{\theta} \log p(\mathbf{x};\theta)], \quad (39)$$

573 where β is the baseline term.

574 **Our reasoning logic:**

575 1). We first briefly describe a high-level idea [19] that adding a baseline term, like the proposed
576 H-term, will help reduce the gradient variance.

577 2). We sketch the steps to show how the proposed H-term will mathematically reduce the gradient
578 variance, following the framework in Section 5.2 of [20].

579 **High-level IDEA.** One generic approach to reduce the variance of Monte Carlo estimates is to use an
580 additive control variate. Suppose we wish to estimate the integral of the function $f : \mathcal{X} \rightarrow \mathbb{R}$, and we
581 know the value of the integral of another function on the same space $\phi : \mathcal{X} \rightarrow \mathbb{R}$.

582 We have

$$\int_{\mathcal{X}} f(x) = \int_{\mathcal{X}} (f(x) - \phi(x)) + \int_{\mathcal{X}} \phi(x), \quad (40)$$

583

584 and the integral of $f(x) - \phi(x)$ can be estimated. If $\phi(x) = f(x)$, meaning that , then we have
585 managed to reduce our variance to zero [19]. More generally,

$$\text{Var}(f - \phi) = \text{Var}(f) - 2\text{Cov}(f, \phi) + \text{Var}(\phi). \quad (41)$$

586 If ϕ and f are strongly correlated, so that the covariance term on the right hand side is greater than
587 the variance of ϕ , i.e., $-2\text{Cov}(f, \phi) + \text{Var}(\phi) \leq 0$. then a variance reduction has been made over the
588 original estimation problem [19], i.e., $\text{Var}(f - \phi) \leq \text{Var}(f)$.

589 **Our reasoning.** Then, we present our reasoning.

590 Note that the gradient of the new objective function of the actor network in (14) consists of two
591 components, namely $\nabla_{\theta} J(\theta)$ and $\nabla_{\theta} H(\theta)$. Here, we consider

$$\nabla_{\theta} J(\theta) - \lambda \nabla_{\theta} H(\theta), \quad \text{where } \lambda > 0 \text{ is a temperature parameter,} \quad (42)$$

592 where $\nabla_{\theta} J(\theta)$ in (27) is the above function $f(\cdot)$ and $\lambda \nabla_{\theta} H(\theta)$ in (15) is the above function $\phi(\cdot)$.

593 The Hamiltonian stochastic gradient in (15) has the optimal value

$$\nabla_{\theta} H^*(\theta) = - \lim_{K \rightarrow \infty} \mathbb{E}_{\mu_0, \dots, \mu_{K-1}} \left[\sum_{k=0}^{K-1} \gamma^k \cdot R(\mu_k) \cdot \nabla_{\theta} \log (\pi_{\theta}(\mu_0) \cdot \pi_{\theta}(\mu_1) \cdots \pi_{\theta}(\mu_k)) \right]. \quad (43)$$

594 According to Theorem 8 of [19] that is proved via (41), we have

$$\begin{aligned} \text{Var} [\nabla_{\theta} J(\theta) - \lambda \nabla_{\theta} H^*(\theta)] &= \text{Var} [\nabla_{\theta} J(\theta)] - \frac{1}{\lambda} \mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} \left[\frac{(\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2 \nabla_{\theta} J(\theta)])^2}{\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2]} \right] \\ &\leq \text{Var} [\nabla_{\theta} J(\theta)], \end{aligned} \quad (44)$$

595 where the second term is positive, and

$$\nabla_{\theta} H^*(\theta) = \frac{\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2 \nabla_{\theta} J(\theta)]}{\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2]}. \quad (45)$$

596 In Alg. 1 and Alg. 2, we used a general H-term $\nabla_{\theta} H(\theta)$, not the optimal one in (43). Next, we
597 provide a general characterization for this case..

598 According to Theorem 10 of [19], we have

$$\begin{aligned} &\text{Var} [\nabla_{\theta} J(\theta) - \lambda \nabla_{\theta} H(\theta)] - \text{Var} [\nabla_{\theta} J(\theta) - \lambda \nabla_{\theta} H^*(\theta)] \\ &= \lambda^2 \mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2 (\nabla_{\theta} H(\theta) - \nabla_{\theta} H^*(\theta))^2] \end{aligned} \quad (46)$$

599 Assume Lipschitz continuity of the gradient $\nabla_{\theta}H(\theta)$ such that

$$\|\nabla_{\theta}H(\theta) - \nabla_{\theta}H^*(\theta)\|_2 \leq 2L(H(\theta) - H^*(\theta)) \leq 2L\epsilon, \quad (47)$$

600 given $K \geq \log_{\gamma} \epsilon$ with $L > 0, \epsilon > 0$, as pointed out in the end of Section 3.2.

601 Therefore, combining (46), (48) with (48), we obtain that

$$\begin{aligned} \text{Var}[\nabla_{\theta}J(\theta) - \lambda \nabla_{\theta}H(\theta)] &= \text{Var}[\nabla_{\theta}J(\theta)] - \frac{1}{\lambda} \mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} \left[\frac{(\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2 \nabla_{\theta}J(\theta)])^2}{\mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2]} \right] \\ &\quad + \lambda^2 (2L\epsilon)^2 \mathbb{E}_{s \sim d_{\theta}, a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(s, a))^2] \\ &\leq \text{Var}[\nabla_{\theta}J(\theta)], \end{aligned} \quad (48)$$

602 when both $|\nabla_{\theta} \log \pi_{\theta}(s, a)|$ and $|\nabla_{\theta}J(\theta)|$ are upper bounded, e.g., $|\nabla_{\theta} \log \pi_{\theta}(s, a)| < C_1$ and
603 $|\nabla_{\theta}J(\theta)| < C_2$; and we set ϵ, λ properly such that

$$\begin{aligned} -\frac{1}{\lambda} C_2^2 + 4\lambda^2 L^2 \epsilon^2 C_1^2 &< 0 \\ \lambda^3 \epsilon^2 &< \frac{C_2^2}{4L^2 C_1^2}, \end{aligned} \quad (49)$$

604 which can be easily satisfied by properly selecting λ and $K \geq \log_{\gamma} \epsilon$.

605 **Conclusion:**

606 To sum up, we show that it is easy to achieve $\text{Var}[\nabla_{\theta}J(\theta) - \lambda \nabla_{\theta}H(\theta)] \leq \text{Var}[\nabla_{\theta}J(\theta)]$, which
607 means adding the H-term can lead to smaller variance than that of the conventional gradient.

608 **F Conventional Actor-Critic Algorithms for Deep Reinforcement Learning**

609 The gradient of (2) is [36]

$$\nabla_{\theta} J(\theta) \triangleq \sum_S d_{S,\theta}(S) \sum_A Q_{\theta}(S, A) \nabla_{\theta} \pi_{\theta}(S, A). \quad (50)$$

610 Since Q_{θ} in (50) is unknown [41] (the stationary distribution d_{θ} is unknown), one can plug in a critic
611 network with parameter ϕ as an estimator of Q_{θ} and obtain

$$\nabla_{\theta}^{\phi} J(\theta, \phi) = \sum_S d_{S,\theta}(S) \sum_A Q_{\phi}(S, A) \nabla_{\theta} \pi_{\theta}(S, A), \quad (51)$$

612 where $d_{S,\theta} \in \mathbb{R}_+^{|\mathcal{S}||\mathcal{A}| \times 1}$ denotes the stationary distribution over the states instead of state-action
613 pairs.

614 (51) is a bi-level optimization problem [7], and a natural solution is an iterative algorithm that
615 alternates between estimating Q_{ϕ} with parameter ϕ and improving policy π_{θ} with parameter θ .
616 Therefore, a family of actor-critic algorithms are proposed with following objective functions:

$$\begin{cases} \text{Actor : } \max_{\theta} J_{\pi}(\theta, \phi) = (1 - \gamma) \mathbb{E}_{S_0 \sim d_0, A_0 \sim \pi_{\theta}(S_0, \cdot)} [Q_{\phi}(S_0, A_0)] \\ \text{Critic : } \max_{\phi} J_Q(\theta, \phi) = \frac{1}{2} \mathbb{E}_{S \sim d_{\theta}(\cdot), A \sim \pi_{\theta}(S, \cdot)} [(Q_{\phi}(S, A) - y(S, A))^2]. \end{cases} \quad (52)$$

617 The gradient of (52) can be estimated as follows

$$\begin{aligned} \nabla_{\theta} \widehat{J}_{\pi}(\theta, \phi) &= \frac{1}{N} \sum_{i=1}^N Q_{\phi}(\mu) \cdot \nabla_{\theta} \log \pi_{\theta}(\mu) \\ \nabla_{\phi} \widehat{J}_Q(\theta, \phi) &= \frac{1}{N} \sum_{i=1}^N [Q_{\phi}(S, A) - y(S, A)] \cdot \nabla_{\phi} Q_{\phi}(S, A) \end{aligned} \quad (53)$$

618 The parameters ϕ and θ are updated as follows:

$$\begin{cases} \text{Actor : } \theta \leftarrow \theta + \alpha \nabla_{\theta}^{\phi} \widehat{J}_{\pi}, \\ \text{Critic : } \phi \leftarrow \phi - \alpha \nabla_{\phi} \widehat{J}_Q. \end{cases} \quad (54)$$

619 G Stationary Deterministic Policy Gradient Algorithm with H-term

620 For completeness, we present the details of the deterministic actor-critic algorithm with H-term.

Algorithm 2 Stationary Actor-Critic Algorithm with H-term

```

1: Input: learning rate  $\alpha$ , temperature  $\lambda$ , look-ahead step  $K$ , and parameters  $\delta, M, T, G, B, B'$ 
2: Initialize actor network  $\eta$  and critic network  $Q$  with parameters  $\theta, \phi$ , and replay buffers  $\mathcal{D}_1, \mathcal{D}_2$ 
3: for episode = 1,  $\dots$ ,  $M$  do
4:   Initialize state  $s_0$ 
5:   for  $t = 0, \dots, T - 1$  do
6:     Take action  $a_t = \eta_\theta(s_t) + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \delta^2)$ 
7:     Execute action  $a_t$ , receive reward  $r_t$ , and observe new state  $s_{t+1}$ 
8:     Store a transition  $(s_t, a_t, r_t, s_{t+1})$  in  $\mathcal{D}_1$ 
9:   end
10:  Store a trajectory  $\tau$  of length  $T$  in  $\mathcal{D}_2$ 
11:  for  $g = 1, \dots, G$  do
12:    Randomly sample a mini-batch of  $B$  transitions  $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1}^B$  from  $\mathcal{D}_1$ 
13:    Randomly sample a mini-batch of  $B'$  trajectories (of length  $K$ )  $\{\tau_j\}_{j=1}^{B'}$  from  $\mathcal{D}_2$ 
14:    Update critic network using a conventional method
15:    Update actor network as  $\theta \leftarrow \theta + \alpha \left( \nabla_\theta \widehat{J}'(\theta) - \lambda \nabla_\theta \widehat{H}'(\theta) \right)$ .
16:  end
17: end

```

621 We apply the proposed Hamiltonian equation (7) to regularize the actor network. Specifically, $H'(\theta)$
622 in (7) is added to the actor's objective with weight $\lambda > 0$. The objective functions of actor and critic
623 networks become:

$$\begin{cases} \text{Actor : } \max_{\theta} J'_\pi(\theta, \phi) = (1 - \gamma) \mathbb{E}_{S_0 \sim d_0, A_0 = \eta_\theta(S_0)} [Q_\phi(S_0, A_0)] - \lambda H'(\theta), \\ \text{Critic : } \min_{\phi} J_Q(\theta, \phi) = \frac{1}{2} \mathbb{E}_{S \sim d_\theta(\cdot), A = \eta_\theta(S)} [(Q_\phi(S, A) - y(S, A))^2]. \end{cases} \quad (55)$$

624 The gradient of (55) is

$$\nabla_\theta J'_\pi(\theta, \phi) = (1 - \gamma) \sum_S d_{S, \theta}(S) \nabla_A Q_\phi(S, A) \cdot \nabla_\theta \eta_\theta(S) - \lambda \nabla_\theta H'(\theta), \quad (56)$$

625

$$\nabla_\phi J_Q(\theta, \phi) = \sum_S d_{S, \theta}(S) \cdot [Q_\phi(S, A) - y(S, A)] \cdot \nabla_\phi Q_\phi(S, A)|_{A = \eta_\theta(S)}. \quad (57)$$

626 To estimate $\nabla_\theta H'(\theta)$, the Monte Carlo gradient estimator in (17) is used. Therefore, (56) and (57)
627 can be estimated as follows:

$$\nabla_\theta \widehat{J}'_\pi(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N [\nabla_A Q_\phi(S, A)|_{A = \eta_\theta(S)} \nabla_\theta \eta_\theta(S)] - \frac{1}{N'} \sum_{i=1}^{N'} \left[\lambda \sum_{k=0}^K \gamma^k R(\mu_k) \nabla_\theta \log [\tilde{\pi}_\theta(\mu_0) \cdots \tilde{\pi}_\theta(\mu_k)] \right], \quad (58)$$

628

$$\nabla_\phi \widehat{J}_Q(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N [Q_\phi(S, A) - y(S, A)] \cdot \nabla_\phi Q_\phi(S, A)|_{A = \eta_\theta(S)}. \quad (59)$$

629 H Experiments: Hyperparameters and More Results

630 H.1 Hyperparameters in Experiments

Table 6: Hyperparameters used for the PPO and PPO + H in MuJoCo tasks

Parameters	Values
Optimizer	Adam
Learning rate	$3 \cdot 10^{-4}$
Discount (γ)	0.99
GAE parameter	0.95
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Mini-batch size	32
Importance rate of H-term (λ)	2^{-3}
Truncation step of H-term (K)	16

Table 7: Hyperparameters used for the DDPG and DDPG + H in MuJoCo tasks

Parameters	Values
Optimizer	Adam
Learning rate	$5 \cdot 10^{-4}$
Target Update Rate (τ)	10^{-3}
Discount (γ)	0.995
Replay buffer size	10^6
Number of hidden layers for all networks	3
Number of hidden units per layer	256
Batch size	64
Importance rate of H-term (λ)	2^{-3}
Truncation step of H-term (K)	16

631 H.2 More Results

632 Fig. 6 shows the H-value (average over 20 runs) during the training process, which verified that the
 633 trained agents have converged to policies with small H-values.

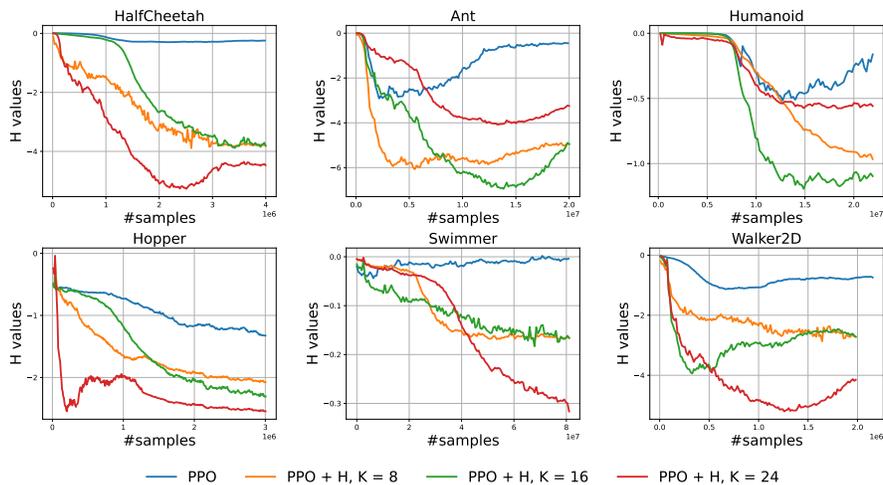


Figure 6: H values during the training process.

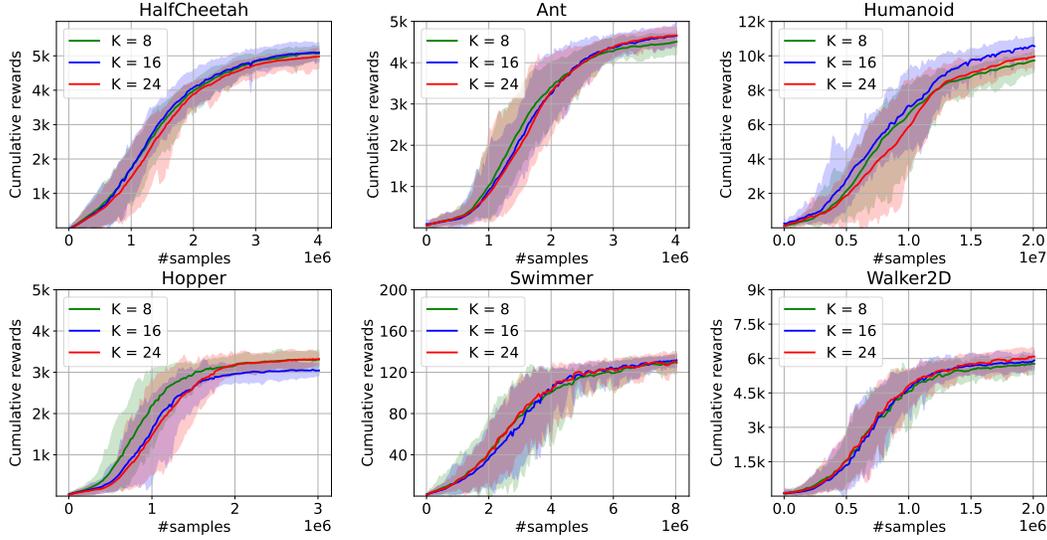


Figure 7: For the proposed PPO+H algorithm, the performance with different K values.

Tasks	Policies		PPO	PPO+H ($K = 8$)	PPO+H ($K = 16$)	PPO+H ($K = 24$)			
HalfCheetah	running	13		17	20	20			
	flipping	5	4720.8	0	4839.5	0	5001.8	0	4995.1
	diving	1	± 969.2	3	± 392.5	0	± 321.5	0	± 383.3
	balancing	1		0		0		0	
Ant	running	17		20	20	20			
	jumping	0	4164	0	4351.6	0	4553.1	0	4662.5
	flipping	3	± 1563.4	0	± 294.5	0	± 276.4	0	± 277.5
Humanoid	two-legs	7		15	15	17			
	one-leg	12	9433.4	5	9583.4	5	9873.2	3	9942.2
	backward	1	± 2513.5	0	± 753.4	0	± 653.7	0	± 538.4
Hopper	hopping	10		18	20	20			
	diving	8	2659.3	2	3014.9	0	3254.2	0	3340.7
	balancing	2	± 905.3	0	± 304.7	0	± 246.1	0	± 191.5
Swimmer	moving	14		19	20	19			
	balancing	6	110.7	1	121.5	0	142.3	1	132.2
Walker	walking	5		16	16	15			
	diving	8	5461.7	2	5794.4	3	5921.8	5	6089.3
	balancing	7	± 1290.1	2	± 341.8	1	± 304.5	0	± 314.7

Table 8: Experimental results on six challenging MuJoCo tasks.

634 Fig. 7 shows more performance of the PPO+H algorithm, for $K = 8, 16, 24$. We run each experiment
 635 with 20 random seeds and each run we test 100 episodes.

636 To verify the hypothesize that smaller replay buffer hurts the performance, we rerun the trials of
 637 $K = 8, 16$ with a replay buffer size 800.

638 I Hamiltonian Policy Network

639 I.1 Hamiltonian Policy Network

640 Since Hamiltonian equation in (7) is a functional of policy π_θ , a natural question would be: can
641 we use the Hamiltonian equation replace existing Bellman’s equation (3) or the policy gradient’s
642 objective function (2)?

643 As a verification, we test the capability of Hamiltonian equation in (7) as a loss function to train a
644 policy network. The algorithm is first given as follows.

Algorithm 3 Hamiltonian Policy Network

```
1: Input: learning rate  $\alpha$ , look-ahead step  $K$ , and parameters  $M, T, G, B$ 
2: Initialize policy network with parameters  $\theta$ , and replay buffer  $\mathcal{D}$ 
3: for episode = 1,  $\dots$ ,  $M$  do
4:   Initialize state  $s_0$ 
5:   for  $t = 0, \dots, T - 1$  do
6:     Select action  $a_t \sim \pi_\theta(\cdot|s_t)$ 
7:     Execute action  $a_t$ , receive reward  $r_t$ , and observe new state  $s_{t+1}$ 
8:   end
9:   Store a trajectory  $\tau$  of length  $T$  in  $\mathcal{D}$ 
10:  for  $g = 1, \dots, G$  do
11:    Randomly sample a mini-batch of  $B$  trajectories (of length  $K$ )  $\{\tau_j\}_{j=1}^B$  from  $\mathcal{D}$ 
12:    Update policy network as  $\theta \leftarrow \theta - \alpha \nabla_\theta \hat{H}(\theta)$ .
13:  end
14: end
```

645 In Alg. 3, an agent interacts with an environment and updates its policy network. The algorithm has
646 M episodes and each episode consists of a (Monte Carlo) simulation process and a learning process
647 (gradient estimation) as follows:

- 648 • During the (Monte Carlo) simulation process (lines 5-9 of Alg. 3), an agent takes action a_t
649 according to a policy $\pi_\theta(\cdot|s_t)$, $t = 0, \dots, T - 1$, generating a trajectory of T steps/transitions.
650 Then, the full trajectory $\tau = (s_0, a_0, r_0, s_1, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$ is stored in replay buffer
651 \mathcal{D} .
- 652 • During the learning process ($G \geq 1$ updates in one episode) (lines 10-12 of Alg. 1), a mini-batch of
653 B trajectories (of length K) $\{\tau_j = (s_0^j, a_0^j, r_0^j, s_1^j, \dots, s_{K-1}^j, a_{K-1}^j, r_{K-1}^j, s_K^j)\}_{j=1}^B$ are sampled
654 from \mathcal{D} , respectively. The policy network is updated by a Monte Carlo gradient estimator over B
655 trajectories.

656 **Implementation of replay buffer \mathcal{D} .** After a full trajectory τ of length T is generated, it is partitioned
657 into $T - K + 1$ trajectories of length K . We rank them according to the cumulative reward and
658 store the top portion, say 80%, into a new replay buffer \mathcal{D} (line 9 of Alg. 3). We randomly sample a
659 mini-batch of B trajectories from \mathcal{D} (line 11 of Alg. 3) to compute the H-term.

660 I.2 Frozenlake Task

661 **Environment:** Frozenlake 8×8 , a game in OpenAI Gym.

662 **Rules:** As shown in Fig. 8 (left), the Frozenlake task has 8×8 states with 4 optional actions to move
663 around. The agent needs to go from the start point and find the way to the destination in limited steps.
664 There are 8 holes which can cause the agent to fail the game.

665 **Experiment settings:** We take Deep Q-learning (DQN) [29] as our baseline and use the implementa-
666 tion from the ElegantRL library. We use a 4-layer fully connected neural network as the deep policy
667 network both in DQN and DHN. We use the Adam optimizer with a learning rate 1×10^{-3} and a
668 batch size 100.

669 **Evaluation:** We evaluate the performance of policy by computing the success rate, in which we use
670 50 agents to walk 100 steps and compute the rates of agents who successfully arrive the destination.

671 **Results for the Frozenlake task:** Fig. 9 (left) shows the success rate of agents with increasing the
 672 number of transitions learned by the network. compared with DQN, DHN has a more stable training
 673 process. It is easy for DQN to quickly obtain a good policy to win the game. But with increasing the
 674 number of transitions fed to the network, the performance of DQN shows a large and frequent shock
 675 while the performance of DHN shows the strong stability.

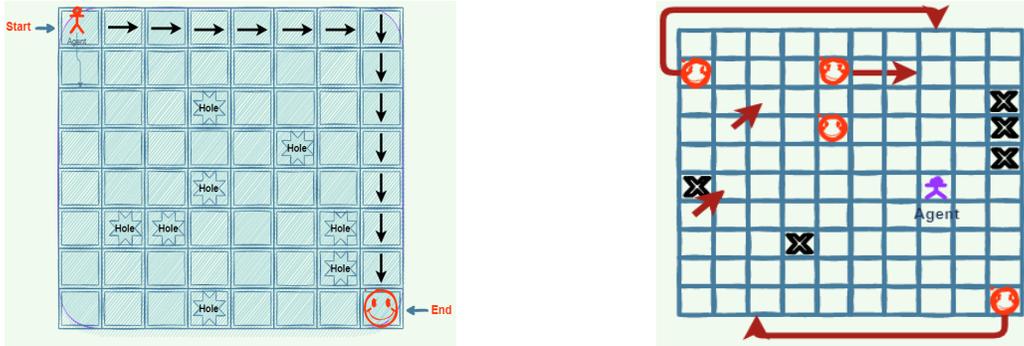


Figure 8: The Frozenlake task (left) and Gridworld task (right).

676 I.3 Gridworld Task

677 **Environment:** a Gridworld of size 10×10 , a game available in our code.

678 **Rules:** As shown in the Fig. 8 (right), the Gridworld has 10×10 states with 4 optional actions to
 679 move around. The agent will initialize at a random locations and it needs to find the smiley as many
 680 as possible which has 10 reward in turn. It should be noted that there are some endpoints which may
 681 cause the agent game over and some transfer-points which transfer the agent to certain location.

682 **Experiment settings and evaluation:** Both the experiment settings and evaluation method are the
 683 same with that on Frozenlake 8×8 game.

684 **Results for the Gridworld task:** Fig. 9 (right) shows the mean reward obtained by the agents with
 685 increasing the training time. Compared with DQN, DHN has a faster training process. It only needs
 686 massive random parallel samples of trajectories and do not need any policy for guided sampling while
 687 DQN needs guided exploration in the training process which costs a large time consumption.

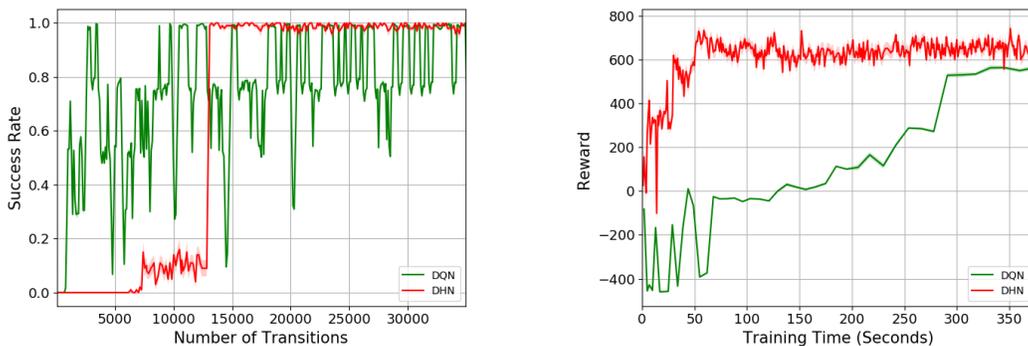


Figure 9: Comparison between the DQN and DHN algorithms. The Frozenlake task (left) and Gridworld task (right).