

## A Broader Impact

Designing appropriate rewards is a challenging task in real-world applications of reinforcement learning (RL). Often, researchers or algorithm engineers possess prior knowledge, such as rules and constraints, about the problem they are trying to solve, but are unable to represent this knowledge precisely with numeric values. Improper reward settings can lead to unexpected behaviors learned by an RL algorithm, resulting in exploitation of the rewards for higher returns. In this paper, we propose a potential-based exploration bonus as shaping reward that helps avoid repetitive and tedious tuning of rewards in RL applications, such as autonomous driving, robot learning control, video games and online recommendation/ad optimization. Our approach liberates researchers and algorithm engineers from the laborious task of reward tuning, while also providing insight into the quality of the designed rewards through the shaping weights learned by our methods.

## B Bisimulation Relation and Wasserstein Distances

To delve deeply into the formulation of bisimulation metric, we present a detailed definition of both bisimulation relation and Wasserstein distances.

### B.1 Definition of Bisimulation Relation

**Definition 4.** (*Bisimulation Relations* [Givan et al. \[2003\]](#)) Given an MDP  $\mathcal{M}$ , an equivalence relation  $B$  between states is a bisimulation relation if, for all states  $s_i, s_j \in \mathcal{S}$  that are equivalent under  $B$  (denoted  $s_i \equiv_B s_j$ ) the following conditions hold:

$$\mathcal{R}(s_i, a) = \mathcal{R}(s_j, a) \quad \forall a \in \mathcal{A}, \quad (11)$$

$$\mathcal{P}(G \mid s_i, a) = \mathcal{P}(G \mid s_j, a) \quad \forall a \in \mathcal{A}, \quad \forall G \in \mathcal{S}_B \quad (12)$$

where  $\mathcal{S}_B$  is the partition of  $\mathcal{S}$  under the relation  $B$  (the set of all groups  $G$  of equivalent states), and  $\mathcal{P}(G \mid s, a) = \sum_{s' \in G} \mathcal{P}(s' \mid s, a)$

### B.2 Wasserstein Distances

**Definition 5.** (*Wasserstein metric* [Villani et al. \[2009\]](#)) Let  $d : X \times X \rightarrow [0, \infty)$  be a distance function and  $\Omega$  the set of all joint distributions with marginals  $\mu$  and  $\lambda$  over the space  $X$ ;

$$W_p(d)(\mu, \lambda) = \left( \inf_{\omega \in \Omega} \mathbb{E}_{(x_1, x_2) \sim \omega} [d(x_1, x_2)^p] \right)^{\frac{1}{p}} \quad (13)$$

**Definition 6.** (*Dual formulation of the Wasserstein metric* [Villani et al. \[2009\]](#)) Let  $d : X \times X \rightarrow [0, \infty)$  be a distance function and  $\mu$  and  $\lambda$  marginals over the space  $X$ ;

$$W_p(d)(\mu, \lambda) = \left( \sup_{\phi \oplus \psi \leq d^p} \mathbb{E}_{x_1 \sim \mu} [\phi(x_1)] + \mathbb{E}_{x_2 \sim \lambda} [\psi(x_2)] \right)^{\frac{1}{p}} \quad (14)$$

where  $\phi \oplus \psi \leq d^p \iff \phi(x) + \psi(y) \leq d(x, y)^p, \forall (x, y) \in X \times X$

This dual formulation takes a simple form for  $p = 1$  :

$$W_1(d)(\mu, \lambda) = \sup_{f \in \text{Lip}_{1,d}(X)} \mathbb{E}_{x_1 \sim \mu} [f(x_1)] - \mathbb{E}_{x_2 \sim \lambda} [f(x_2)] \quad (15)$$

where  $\text{Lip}_{1,d}(X)$  denotes 1-Lipschitz functions  $f : X \rightarrow \mathbb{R}$  such that  $|f(x_1) - f(x_2)| \leq d(x_1, x_2)$ .

It is notable that the 2-Wasserstein metric  $W_2(\|\cdot\|_2)$  (or simply  $W_2$ ) has a closed-form for Gaussians:

$$W_2(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mu_j, \Sigma_j))^2 = \|\mu_i - \mu_j\|_2^2 + \|\Sigma_i - \Sigma_j\|_{\mathcal{F}}^2 \quad (16)$$

where  $\|\cdot\|_{\mathcal{F}}$  denotes the Frobenius norm. We can observe that for point masses (i.e.,  $\Sigma_i, \Sigma_j \rightarrow 0$ ), the 2-Wasserstein metric is equivalent to the Euclidean distance between the two points.

**Lemma 1.** (*p-Wasserstein Inequality* [Villani et al. \[2009\]](#)) For any two distributions  $\mu, \lambda$ , if  $p \leq q$  :

$$W_p(\mu, \lambda) \leq W_q(\mu, \lambda) \quad (17)$$

Lemma 1 shows the inequality between  $p$ -Wasserstein.

## C Theorems and Proofs

**Theorem 1.** Let  $\text{met}$  be the space of bounded pseudometrics on  $\mathcal{S}$ ,  $\gamma \in [0, 1)$  and  $\pi$  a policy that is continuously improving. Define  $\mathcal{H} : \text{met} \mapsto \text{met}$  by:

$$\begin{aligned} \mathcal{H}(d, \pi)(s_i, s_j) &= |r_{s_i}^\pi - r_{s_j}^\pi| + \gamma W(d)(\mathcal{P}_{s_i}^\pi, \mathcal{P}_{s_j}^\pi) \\ &\quad + \|I^\pi(\cdot | s_i, s_{i+1}) - I^\pi(\cdot | s_j, s_{j+1})\|_1 \end{aligned} \quad (18)$$

Then  $\mathcal{H}$  has a least fixed point  $\tilde{d}$  which is an inverse dynamic bisimulation metric.

*Proof.* Let  $(X, \preceq)$  be a partial order. An  $\omega$ -chain of this partial order is an increasing sequence  $\{x_n\}$ . The partial order is said to be an  $\omega$ -complete partial order ( $\omega$ -cpo) if it contains least upper bounds of all  $\omega$ -chains. It is called an  $\omega$ -cpo with bottom if it additionally contains a least element,  $\perp$ , called bottom. A function  $f : X \rightarrow Y$  between  $\omega$ -cpo's is said to be monotonic if  $x \preceq x' \implies f(x) \preceq f(x')$ . It is continuous if for every  $\omega$ -chain  $\{x_n\}$ ,  $f(\sqcup_{n \in \mathbb{N}} \{x_n\}) = \sqcup_{n \in \mathbb{N}} \{f(x_n)\}$ . A point  $x \in X$  is said to be a prefixed-point of  $f$  if  $f(x) \preceq x$ . It is a fixed-point if  $x = f(x)$ .

**Lemma 2.** Given a fixed policy  $\hat{\pi}$ . Define  $\mathcal{H}^{\hat{\pi}} : \text{met} \mapsto \text{met}$  by  $\mathcal{H}(d, \hat{\pi})(s_i, s_j) = |r_{s_i}^{\hat{\pi}} - r_{s_j}^{\hat{\pi}}| + \gamma W(d)(\mathcal{P}_{s_i}^{\hat{\pi}}, \mathcal{P}_{s_j}^{\hat{\pi}}) + \|I^{\hat{\pi}}(\cdot | s_i, s_{i+1}) - I^{\hat{\pi}}(\cdot | s_j, s_{j+1})\|_1$ , then  $\mathcal{H}^{\hat{\pi}}$  has a least fixed point  $d^{\hat{\pi}}$ .

*Proof.* This proof mimics the proof of Theorem 4.5 from Ferns et al. [2004]. We make use of the same pointwise ordering on  $\text{met}$ :  $d \leq d'$  iff  $d(s, t) \leq d'(s, t)$  for all  $s, t \in \mathcal{S}$ , which gives us an  $\omega$ -cpo with bottom  $\perp$ , which is the everywhere-zero metric. Since Lemma 4.4 from Ferns et al. [2004] (wasserstein metric  $W$  is continuous) also applies in our definition, it only remains to show that  $\mathcal{H}(d, \hat{\pi})$  is continuous:

$$\begin{aligned} \mathcal{H}^{\hat{\pi}}\left(\sqcup_{n \in \mathbb{N}} \{x_n\}\right)(s_i, s_j) &= |r_{s_i}^{\hat{\pi}} - r_{s_j}^{\hat{\pi}}| + \gamma W\left(\sqcup_{n \in \mathbb{N}} \{x_n\}\right)(\mathcal{P}_{s_i}^{\hat{\pi}}, \mathcal{P}_{s_j}^{\hat{\pi}}) \\ &\quad + \|I^{\hat{\pi}}(\cdot | s_i, s_{i+1}) - I^{\hat{\pi}}(\cdot | s_j, s_{j+1})\|_1 \\ &= |r_{s_i}^{\hat{\pi}} - r_{s_j}^{\hat{\pi}}| + \gamma \sup_{n \in \mathbb{N}} W(x_n)(\mathcal{P}_{s_i}^{\hat{\pi}}, \mathcal{P}_{s_j}^{\hat{\pi}}) \\ &\quad + \|I^{\hat{\pi}}(\cdot | s_i, s_{i+1}) - I^{\hat{\pi}}(\cdot | s_j, s_{j+1})\|_1 \\ &\quad \text{by continuity of } W \\ &= \sup_{n \in \mathbb{N}} (|r_{s_i}^{\hat{\pi}} - r_{s_j}^{\hat{\pi}}| + \gamma W(x_n)(\mathcal{P}_{s_i}^{\hat{\pi}}, \mathcal{P}_{s_j}^{\hat{\pi}}) \\ &\quad + \|I^{\hat{\pi}}(\cdot | s_i, s_{i+1}) - I^{\hat{\pi}}(\cdot | s_j, s_{j+1})\|_1) \\ &= \sup_{n \in \mathbb{N}} \{\mathcal{H}^{\hat{\pi}}(x_n)(s_i, s_j)\} \\ &= \left(\sqcup_{n \in \mathbb{N}} \{\mathcal{H}^{\hat{\pi}}(x_n)\}\right)(s_i, s_j) \end{aligned} \quad (19)$$

The rest of the proof follows in the same way as in Ferns et al. [2004].  $\square$

Ideally, to prove this theorem we show that  $\mathcal{H}$  is monotonically increasing and continuous, and apply Fixed Point Theorem to show the existence of a fixed point that  $\mathcal{H}$  converges to. Unfortunately, we can show that  $\mathcal{H}$  under  $\pi$  as  $\pi$  monotonically converges to  $\pi^*$  is not also monotonic, unlike the original bisimulation metric setting Ferns et al. [2004] and the policy evaluation setting Castro [2020]. We start the iterates  $\mathcal{H}^n$  from bottom  $\perp$ , denoted as  $\mathcal{H}^n(\perp)$ . In Ferns et al. [2004] the  $\max_{a \in \mathcal{A}}$  can be thought of as learning a policy between every two pairs of states to maximize their distance, and therefore this distance can only stay the same or grow over iterations of  $\mathcal{H}$ . In Castro [2020],  $\pi$  is fixed, and under a deterministic MDP it can also be shown that distance between states  $d_n(s_i, s_j)$  will only expand, not contract as  $n$  increases. In the policy iteration setting, however, with  $\pi$  starting from initialization  $\pi_0$  and getting updated:

$$\pi_k(s) = \arg \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} [r_{ss'}^a + \gamma V^{\pi_{k-1}}(s')] \quad (20)$$

there is no guarantee that the distance between two states  $d_{n-1}^{\pi_{k-1}}(s_i, s_j) < d_n^{\pi_k}(s_i, s_j)$  under policy iterations  $\pi_{k-1}, \pi_k$  and distance metric iterations  $d_{n-1}, d_n$  for  $k, n \in \mathbb{N}$ , which is required for monotonicity.

Instead, we show that using the policy improvement theorem which gives us:

$$V^{\pi_k}(s) \geq V^{\pi_{k-1}}(s), \forall s \in \mathcal{S} \quad (21)$$

$\pi$  will converge to a fixed point using the Fixed Point Theorem, and taking the result of lemma 2 that  $\mathcal{H}^\pi$  has a fixed point for every  $\pi \in \Pi$ , we can show that a fixed point inverse dynamic bisimulation metric will be found with policy iteration.  $\square$

**Theorem 2.** (Value difference bound) Given any two states  $s_i, s_j \in \mathcal{S}$  in an MDP  $\mathcal{M}$ , let  $V^\pi(s)$  be the value function of policy  $\pi$ , we can get:

$$|V^\pi(s_i) - V^\pi(s_j)| \leq d_{inv}(s_i, s_j) \quad (22)$$

where  $d_{inv}$  is a inverse dynamic bisimulation metric.

*Proof.* We will use the standard value function update:  $V_n^\pi(s) = \mathcal{R}_s^\pi + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_s^\pi(s') V_{n-1}^\pi(s')$  with  $V_0^\pi \equiv 0$  and our update operator from Theorem 1:  $d_n^\pi(s, t) = \mathcal{H}^\pi(d_{n-1}^\pi)(s, t)$  with  $d_0^\pi \equiv 0$ , and prove this by induction, showing that for all  $n \in \mathbb{N}$  and  $s, t \in \mathcal{S}$ ,  $|V_n^\pi(s) - V_n^\pi(t)| \leq d_n^\pi(s, t)$ .

The base case holds vacuously:  $0 = V_0^\pi(s, t) \leq d_0^\pi(s, t) = 0$ , so assume true up to  $n$ .

$$\begin{aligned} |V_{n+1}^\pi(s) - V_{n+1}^\pi(t)| &= \left| \mathcal{R}_s^\pi + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_s^\pi(s') V_n^\pi(s') - \left( \mathcal{R}_t^\pi + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_t^\pi(s') V_n^\pi(s') \right) \right| \\ &\leq |\mathcal{R}_s^\pi - \mathcal{R}_t^\pi| + \left| \gamma \sum_{s' \in \mathcal{S}} V_n^\pi(s') (\mathcal{P}_s^\pi(s') - \mathcal{P}_t^\pi(s')) \right| \\ &\leq |\mathcal{R}_s^\pi - \mathcal{R}_t^\pi| + |\gamma W(d_n^\pi)(\mathcal{P}_s^\pi, \mathcal{P}_t^\pi)| \\ &= |\mathcal{R}_s^\pi - \mathcal{R}_t^\pi| + \gamma W(d_n^\pi)(\mathcal{P}_s^\pi, \mathcal{P}_t^\pi) \text{ since } W(d_n^\pi) \text{ is a metric} \\ &\leq |\mathcal{R}_s^\pi - \mathcal{R}_t^\pi| + \gamma W(d_n^\pi)(\mathcal{P}_s^\pi, \mathcal{P}_t^\pi) + \|I^\pi(\cdot | s, s') - I^\pi(\cdot | t, t')\|_1 \\ &= \mathcal{H}^\pi(d_n^\pi)(s, t) \\ &= d_{n+1}^\pi(s, t) \end{aligned} \quad (23)$$

where the second inequality follows from noticing that, by induction, for all  $s, t \in \mathcal{S}$ ,  $V_n^\pi(s) - V_n^\pi(t) \leq d_n(s, t)$ , and the third inequality shows that the property as value difference bound holds for both bisimulation metric and inverse dynamic bisimulation metric.  $\square$

**Theorem 3.** The potential function  $d_{inv}(s, s_0)$  is an approximation of absolute value of optimal value function  $V^*(s)$ .

*Proof.* Based on theorem 2, we can get:

$$\begin{aligned} d_{inv}(s, s_0) &\geq |V^\pi(s) - V^\pi(s_0)| \\ &\geq |V^\pi(s)| - |V^\pi(s_0)| \\ &= |V^\pi(s)| - C_1 \end{aligned} \quad (24)$$

where  $C_1$  is a constant. Assume that the potential function  $d_{inv}(s, s_0)$  is defined under policy  $\pi$ , we use the closed form induction for the 2-Wasserstein metric  $W_2(\|\cdot\|_2)$  as Equation 16:

$$\begin{aligned} \lim_{t \rightarrow \infty} d_{inv}(s_t, s_0) - |V^*(s_t)| &= \lim_{t \rightarrow \infty} [|r_t^\pi - r_0^\pi| + \gamma \|\mu_t^\pi - \mu_0\|_2^2 \\ &\quad + \gamma \|\Sigma_t^{(1/2)} - \Sigma_0^{(1/2)}\|_F^2 \\ &\quad + \gamma \|a_t^\pi - a_0\|_1 - V^*(s_t)] \end{aligned} \quad (25)$$

The action  $a$  belongs to the set of bounded action space  $\mathcal{A}$ , so the fourth term on the right of Equation (25) is bounded, as the timestep increases, the distribution  $\mathcal{P}^\pi(\cdot | s)$  converges to the transition function  $\mathcal{N}(\mu_c, \Sigma_c)$ , thus the second term and third term are bounded. So we can get:

$$\begin{aligned} \lim_{t \rightarrow \infty} d_{inv}(s_t, s_0) - |V^*(s_t)| &= \lim_{t \rightarrow \infty} [|r_t^\pi| - |V^*(s_t)|] + C_2 \\ &= \lim_{t \rightarrow \infty} [|r_t^\pi| - |\max_a r_t(s, a)|] \\ &\quad - \gamma \left| \sum_{s_{t+1} \in \mathcal{S}} \mathcal{P}_{s_t s_{t+1}}^a V^*(s_{t+1}) \right| + C_2 \\ &\leq C_2 \end{aligned} \quad (26)$$

where  $C_2$  is a constant. Based on equation 24 and equation 26, we can get:

$$\lim_{t \rightarrow \infty} |V^*(s_t)| - C_1 \leq \lim_{t \rightarrow \infty} d_{inv}(s_t, s_0) \leq \lim_{t \rightarrow \infty} |V^*(s_t)| + C_2 \quad (27)$$

so the potential function  $d_{inv}(s, s_0)$  is an approximation of absolute value of optimal value function  $V^*(s)$   $\square$

**Theorem 4.** Suppose that the shaping reward function  $\mathcal{F}$  takes the form of Equation (7), the optimal value function of the modified MDP  $\mathcal{M}'$ , the potential function  $\Phi(s)$  and the optimal value function of original MDP  $\mathcal{M}$  holds the condition that:

$$V_{\mathcal{M}'}^*(s, a) = V_{\mathcal{M}}^*(s, a) - \Phi(s) \quad (28)$$

*Proof.* According to the Bellman equation, the optimal value function of the original MDP  $\mathcal{M}$  has the following form:

$$V_{\mathcal{M}}^*(s) = \max_a \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [\mathcal{R}(s, a, s') + \gamma V_{\mathcal{M}}^*(s')] \quad (29)$$

We can obtain the following by subtracting the potential function from both sides:

$$\begin{aligned} V_{\mathcal{M}}^*(s) - \Phi(s) &= \max_a \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [\mathcal{R}(s, a, s') + \gamma \Phi(s') - \Phi(s) \\ &\quad + \gamma (V_{\mathcal{M}}^*(s') - \Phi(s'))] \end{aligned} \quad (30)$$

Based on Equation (1), we can get:

$$\begin{aligned} V_{\mathcal{M}'}^*(s) &= \max_a \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [\mathcal{R}(s, a, s') + \mathcal{F}(s, a, s') + \gamma V_{\mathcal{M}'}^*(s')] \\ &= \max_a \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a)} [\mathcal{R}'(s, a, s') + \gamma \hat{V}_{\mathcal{M}'}(s')] \end{aligned} \quad (31)$$

So the value function of the modified MDP  $\mathcal{M}'$  is:

$$V_{\mathcal{M}'}^*(s) \triangleq V_{\mathcal{M}}^*(s) - \Phi(s) \quad (32)$$

$\square$

## D Additional Experiments

### D.1 Delayed Reward Setting

We provide the full result for the delayed reward setting, where the accumulated reward is given every 10, 20, 30 or 40 steps. Table 2 demonstrates that LIBERTY outperforms other methods in 19 out of 24 delayed reward tasks. This suggests that even with sparse rewards, LIBERTY can achieve effective exploration while maintaining high performance. Since DPBA only receives delayed rewards, its performance decreases as the delay period increases. In contrast, RIDE and RND can generate more exploration for the agent at each step, resulting in better performance than DPBA. Moreover, curiosity-driven methods perform similarly to delayed reward-based methods, despite the lack of explicit rewards. The result provides further evidence that LIBERTY's potential-based exploration bonus encourages efficient exploration in the environment, even in scenarios with only delayed rewards.

Table 2: Full Table of quantitative results comparison between LIBERTY and other baseline methods in different environments of Mujoco with the delayed reward setting. The best and the runner-up results are (**bold**) and (underline)

Delay = 10						
Methods	HalfCheetah	Hopper	Walker2d	Ant	Humanoid	Swimmer
ICM	1374 $\pm$ 368	1258 $\pm$ 325	1127 $\pm$ 225	-105 $\pm$ 43	462 $\pm$ 54	27 $\pm$ 11
RND	1694 $\pm$ 495	1976 $\pm$ 458	1405 $\pm$ 262	143 $\pm$ 17	532 $\pm$ 29	32 $\pm$ 15
NGU	1180 $\pm$ 513	989 $\pm$ 262	1275 $\pm$ 480	-164 $\pm$ 35	413 $\pm$ 78	24 $\pm$ 12
RIDE	2467 $\pm$ 456	1876 $\pm$ 431	1651 $\pm$ 325	92 $\pm$ 31	570 $\pm$ 45	65 $\pm$ 16
DPBA	1514 $\pm$ 365	2103 $\pm$ 129	1997 $\pm$ 115	<b>592 <math>\pm</math> 67</b>	518 $\pm$ 23	43 $\pm$ 17
LIBERTY	<b>2973 <math>\pm</math> 437</b>	<b>2479 <math>\pm</math> 315</b>	<b>2766 <math>\pm</math> 487</b>	292 $\pm$ 68	<b>681 <math>\pm</math> 73</b>	<b>73 <math>\pm</math> 21</b>
LIBERTY w/o I.D.	1783 $\pm$ 412	1676 $\pm$ 275	1732 $\pm$ 392	131 $\pm$ 22	505 $\pm$ 37	46 $\pm$ 11
Delay = 20						
Methods	HalfCheetah	Hopper	Walker2d	Ant	Humanoid	Swimmer
ICM	1185 $\pm$ 287	1097 $\pm$ 275	995 $\pm$ 201	-175 $\pm$ 23	434 $\pm$ 48	23 $\pm$ 10
RND	1595 $\pm$ 415	<u>1925 <math>\pm</math> 401</u>	1379 $\pm$ 193	127 $\pm$ 12	519 $\pm$ 25	29 $\pm$ 12
NGU	1198 $\pm$ 492	978 $\pm$ 177	1158 $\pm$ 375	-177 $\pm$ 31	416 $\pm$ 57	21 $\pm$ 11
RIDE	2285 $\pm$ 402	1621 $\pm$ 382	1724 $\pm$ 307	105 $\pm$ 27	509 $\pm$ 21	59 $\pm$ 13
DPBA	1337 $\pm$ 315	1783 $\pm$ 89	1543 $\pm$ 137	<b>381 <math>\pm</math> 45</b>	501 $\pm$ 18	32 $\pm$ 15
LIBERTY	<b>2619 <math>\pm</math> 354</b>	<b>2112 <math>\pm</math> 208</b>	<b>2345 <math>\pm</math> 414</b>	263 $\pm$ 55	<b>617 <math>\pm</math> 53</b>	<b>67 <math>\pm</math> 18</b>
LIBERTY w/o I.D.	1554 $\pm$ 256	1527 $\pm$ 312	1447 $\pm$ 363	103 $\pm$ 19	497 $\pm$ 36	37 $\pm$ 15
Delay = 30						
Methods	HalfCheetah	Hopper	Walker2d	Ant	Humanoid	Swimmer
ICM	1017 $\pm$ 276	965 $\pm$ 213	798 $\pm$ 199	-198 $\pm$ 25	417 $\pm$ 45	19 $\pm$ 11
RND	1483 $\pm$ 393	<u>1773 <math>\pm</math> 391</u>	1038 $\pm$ 191	99 $\pm$ 13	501 $\pm$ 27	24 $\pm$ 11
NGU	1113 $\pm$ 435	901 $\pm$ 172	1013 $\pm$ 366	-181 $\pm$ 29	399 $\pm$ 45	20 $\pm$ 13
RIDE	1973 $\pm$ 369	1405 $\pm$ 315	1345 $\pm$ 305	87 $\pm$ 21	487 $\pm$ 25	41 $\pm$ 15
DPBA	1021 $\pm$ 297	1461 $\pm$ 91	1179 $\pm$ 135	<b>241 <math>\pm</math> 37</b>	453 $\pm$ 18	21 $\pm$ 11
LIBERTY	<b>2273 <math>\pm</math> 317</b>	<b>1873 <math>\pm</math> 228</b>	<b>2077 <math>\pm</math> 398</b>	215 $\pm$ 48	<b>587 <math>\pm</math> 63</b>	<b>52 <math>\pm</math> 15</b>
LIBERTY w/o I.D.	1305 $\pm$ 243	1297 $\pm$ 238	1197 $\pm$ 367	64 $\pm$ 21	468 $\pm$ 31	25 $\pm$ 11
Delay = 40						
Methods	HalfCheetah	Hopper	Walker2d	Ant	Humanoid	Swimmer
ICM	919 $\pm$ 199	857 $\pm$ 175	697 $\pm$ 172	-213 $\pm$ 27	403 $\pm$ 34	13 $\pm$ 7
RND	1276 $\pm$ 387	<b>1683 <math>\pm</math> 338</b>	968 $\pm$ 168	71 $\pm$ 15	483 $\pm$ 25	17 $\pm$ 11
NGU	1028 $\pm$ 405	879 $\pm$ 155	997 $\pm$ 280	-198 $\pm$ 27	387 $\pm$ 27	11 $\pm$ 6
RIDE	<u>1798 <math>\pm</math> 355</u>	1235 $\pm$ 269	<u>1025 <math>\pm</math> 282</u>	63 $\pm$ 18	468 $\pm$ 23	<b>32 <math>\pm</math> 11</b>
DPBA	883 $\pm$ 275	1382 $\pm$ 85	1016 $\pm$ 129	105 $\pm$ 31	405 $\pm$ 15	9 $\pm$ 3
LIBERTY	<b>2039 <math>\pm</math> 315</b>	<u>1612 <math>\pm</math> 215</u>	<b>1921 <math>\pm</math> 372</b>	<b>142 <math>\pm</math> 45</b>	<b>566 <math>\pm</math> 35</b>	31 $\pm$ 13
LIBERTY w/o I.D.	1231 $\pm$ 253	1213 $\pm$ 207	1012 $\pm$ 358	58 $\pm$ 13	455 $\pm$ 27	17 $\pm$ 8

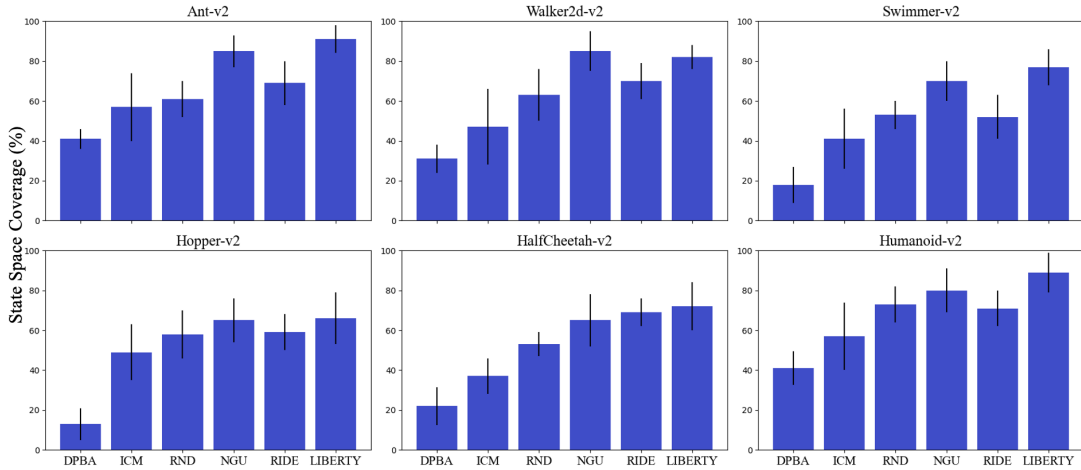


Figure 6: Reward-free exploration results on Ant, Walker2d, Swimmer, Hopper, HalfCheetah and Humanoid. Error bars represent std, deviations over 10 seeds.

## D.2 Reward-free Exploration

We discretize the state-space of all six environments in MuJoCo into bins and compare the number of bins explored, in terms of coverage percentage. As Figure 6 shows, LIBERTY demonstrates its robustness across different tasks by achieving the highest number of bins in 5 out of 6 environments. RIDE and NGU also perform well, closely following LIBERTY. On the other hand, RND outperforms ICM in all six cases, while DPBA performs the worst among the evaluated methods.

## D.3 Episodic Counts for Continuous Control

For the continuous control experiments, we present an additional comparison against the NGU [Badia et al. [2019]] and RIDE [Raileanu and Rocktäschel [2020]], as shown in Table 3, which also employ episodic counts for their intrinsic reward signals. Exact state counts are poorly informative in continuous state spaces, as identical states are rare. Therefore, we adopt episodic counts as pseudo-counts following the procedure presented in NGU for both methods. For NGU, we use inverse-dynamics features to compute pseudo-counts, while for RIDE, we directly use the states. Figure 7 presents the results, including ablations with respect to episodic count modulation. The analysis shows a significant drop in the performance of RIDE in all six environments, while NGU without count consistently underperforms NGU itself. Overall, the episodic count component improves performance for both NGU and RIDE across all environments. It is worth noting that our method eliminates the need for episodic count, which greatly improves scalability.

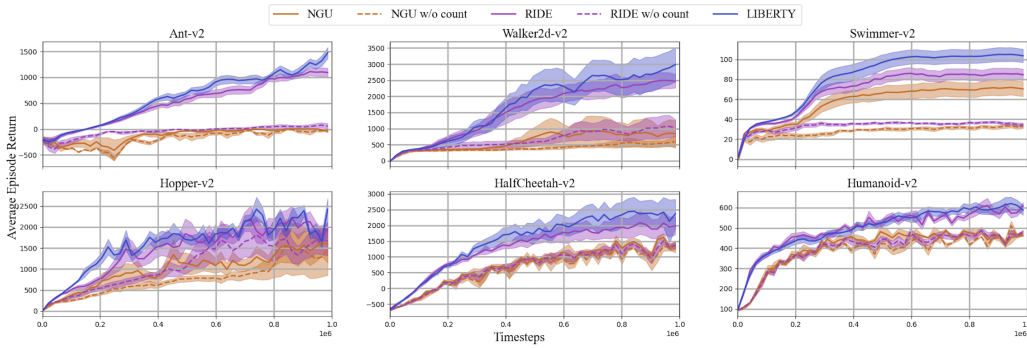


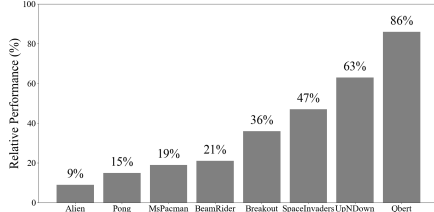
Figure 7: Episodic Counts for Continuous Control. Comparison between LIBERTY, NGU and RIDE, with the latter two being tested both with (original methods) and without (ablations) the episodic-count modulation component. The x-axis represents the number of steps ( $1e6$ ) in training. The y-axis represents the average episode return over the last 100 training episodes (standard deviations in shade). All of the experiments were run using 10 different seeds.

## D.4 Improvements Over Benchmarks

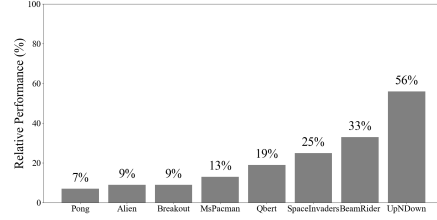
To compare the improvements of LIBERTY over different methods, we select three types of agents: PPO as the baseline agent, DPBA as the PBRs agent, and RIDE as the exploration (via state difference) agent. We ran each agent for 5 separate runs each for 50 million frames on each game for both the baseline agents and LIBERTY-augmented agents. We plotted the best results in Figure 8. The results show that the LIBERTY exploration bonus is effective across different game settings and significantly improves the performance of different types of agents.

## D.5 Scale of Shaping Reward

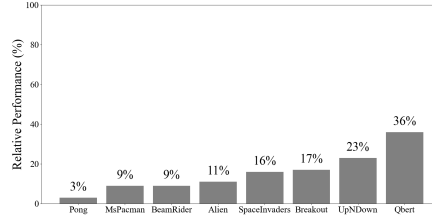
We consider the additive form of reward shaping, which can be defined as  $\mathcal{R}' = \mathcal{R} + \mathcal{F}$ , where  $\mathcal{R}$  is the original reward function,  $\mathcal{F}$  is the shaping reward function, and  $\mathcal{R}'$  is the modified reward function. The scale of the shaping reward can significantly impact the performance of the trained policy. Therefore, we conducted a study on the scale parameter  $\lambda$ , where the modified reward function can be expressed as  $\mathcal{R}' = \mathcal{R} + \lambda\mathcal{F}$ . Note that the default setting in our experiments is  $\lambda = 1$ . To investigate the impact of different scales of shaping rewards on exploration efficiency, we selected  $\lambda = 0.01$ ,  $\lambda = 0.1$ ,  $\lambda = 10$  and  $\lambda = 100$  as different hyper-parameters to be tested on Atari games.



(a) Improvements of LIBERTY augmented agents over baseline (PPO) agents



(b) Improvements of LIBERTY augmented agents over PBRS (DPBA) agents



(c) Improvements of LIBERTY augmented agents over exploration (RIDE) agents

Figure 8: In all figures, the columns correspond to different games labeled on the x-axes and the y-axes show human score normalized improvements.

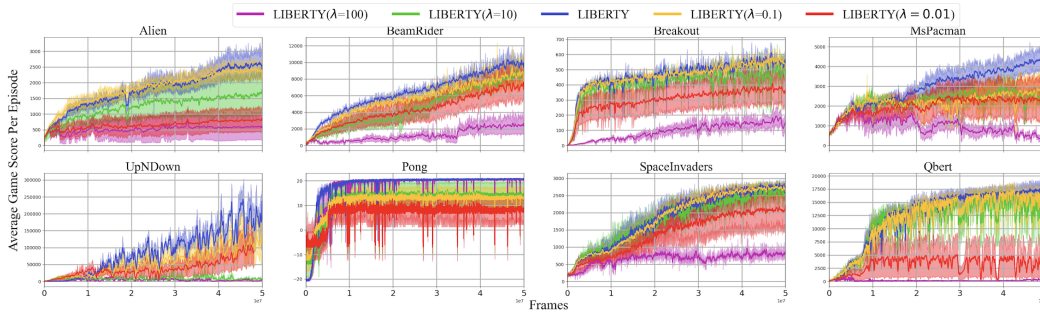


Figure 9: Comparison between LIBERTY with different hyper-parameter  $\lambda$  in the atari games. The x-axis represents the number of frames ( $1e7$ ) in training. The y-axis represents the average game score per episode over the last 100 training episodes (standard deviations in shade). All of the experiments were run using 10 different seeds.

Figure 9 demonstrates that the performance of LIBERTY with  $\lambda = 100$  and  $\lambda = 0.01$  significantly drops. On the other hand, the performance of LIBERTY with  $\lambda = 0.1$  and  $\lambda = 10$  are comparable and almost at the same level as LIBERTY. These results indicate that too much exploration bonus ( $\lambda = 100$ ) or too little exploration bonus ( $\lambda = 0.01$ ) can destroy the algorithm’s performance.

## D.6 Goal-conditioned Tasks

Goal-conditioned reinforcement learning setting is represented by a tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, G)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $\mathcal{P}$  is the probabilistic transition model,  $G \subset \mathcal{S}$  is the goal space which is a set of assignment of values to states, and  $R(s, g) = \mathbb{I}(s = g) \in \mathcal{R}$  is the sparse deterministic reward function that returns 1 only if the state  $s$  match the goal  $g$ . We evaluate LIBERTY in the robot-arm control environment, including three tasks, `FetchPush`, `FetchPickAndPlace`, `FetchSlide`. As Figure 10 shows, we can see that our method LIBERTY still achieves the best performance in the challenging tasks of goal-conditioned environments compared with other baselines, indicating its scalability and superiority.



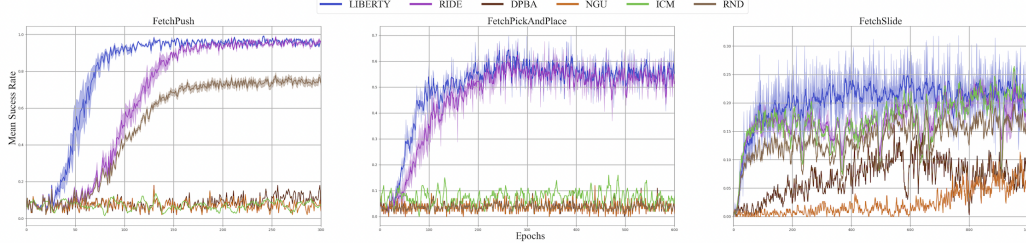


Figure 10: Comparison between LIBERTY with baseline methods in three goal-conditioned tasks, *FetchPush*, *FetchPickAndPlace*, *FetchSlide*. The x-axis represents the number of epochs in training. The y-axis represents the mean success rate (standard deviations in shade). All of the experiments were run using 10 different seeds.

## D.7 Discussion and Ablation Study on $\ell_1$ and $\ell_2$ Norms

We choose  $\ell_1$  norm to measure the discrepancy between actions output by the inverse dynamic module. The reason behind is that actions are scalar or low dimensional vectors in most RL environments. In discrete-action environments, such as *Breakout*, the action space is scalar so there is no difference between  $\ell_1$  and  $\ell_2$  norm. For the case of continuous-action, such as the 4-dimensional vector action space in the *FetchReach* task where the actions' value range is  $[-1, 1]$ , since the action's value range is small, the difference between  $\ell_1$  and  $\ell_2$  is almost negligible. And we had carried ablation study on  $\ell_1$  and  $\ell_2$  norms on six mujoco continuous control tasks, as Figure 11 shows, the performance between  $\ell_1$  and  $\ell_2$  norm is very closed, so we choose the low computational-cost  $\ell_1$  norm.

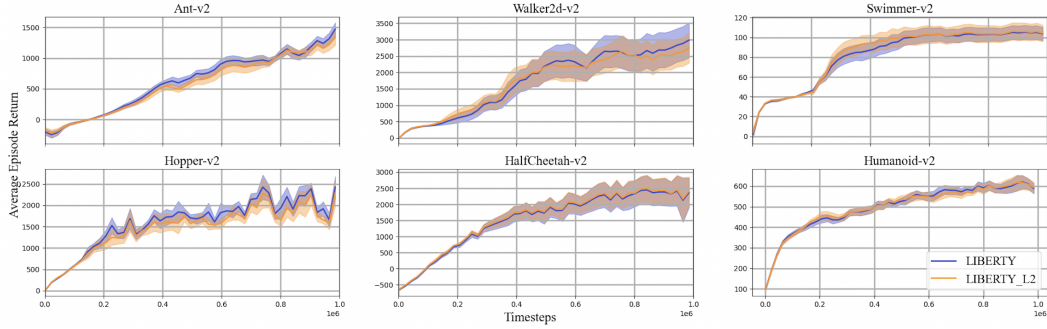


Figure 11: Comparison between LIBERTY and LIBERTY using  $\ell_2$  norm for inverse dynamic module. The x-axis represents the number of timesteps in training. The y-axis represents the average episode reward (standard deviations in shade). The experiments were run using 20 different seeds.

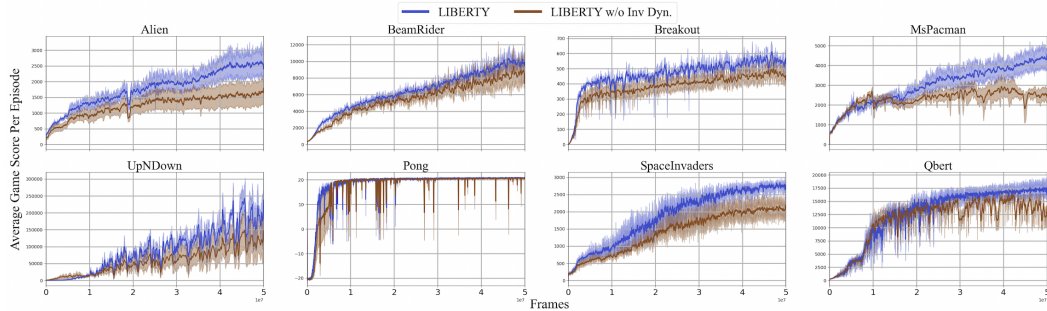


Figure 12: Comparison between LIBERTY and LIBERTY without inverse dynamic module in the atari games. The x-axis represents the number of frames ( $1e7$ ) in training. The y-axis represents the average game score per episode over the last 100 training episodes (standard deviations in shade). All of the experiments were run using 10 different seeds.



## D.8 Ablation study of Inverse Dynamic on Atari Games

To investigate the necessity of inverse dynamic, we denote the variant of LIBERTY as "LIBERTY w/o I.D." which is trained without inverse dynamic. The variant uses the setting of bisimulation metric [Castro \[2020\]](#) only as the potential function. As [Figure 12](#) shows, the performance of LIBERTY on Atari games without the inverse dynamic module had a significant drop in all eight games, which further supports the necessity of inverse dynamic module in promoting efficient exploration.

## D.9 Results on SuperMarioBros

We present the results of level 1 in SuperMarioBros in [Figure 13](#). The x-axis represents the number of frames (1e6) in training. The y-axis represents the average episode reward (standard deviations in shade). The experiments were run using 10 different seeds. We use environment wrappers to preprocess each frame to resize the input image, and stack observations to accelerate the training speed. We can see that our method LIBERTY achieves the best performance among all competitive baselines, which demonstrates the superiority and efficiency of our method. The explanation is that the shaping reward based on state discrepancy and value difference can be helpful in the exploration stage of training. Thus, with the assistance of the shaping reward, LIBERTY helps the agent explore novel states at the start of training, and achieves better or comparable performance at the end of training. NGU achieves the second-best result, because it focus on a long-term exploration. The results of RIDE and RND are comparable, followed by ICM. which suggests that the curiosity-driven methods may fail to explore efficiently compared with our method. DPBA perform the worst which means that the shaping reward is not always useful.

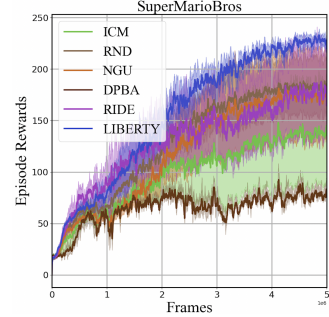


Figure 13: Results of level 1 on SuperMarioBros.

## E Implementation Details

### E.1 Baseline Method Comparison

We make a comparison between the benchmarked methods across three key aspects. Firstly, we consider whether the method includes an episodic count term in the exploration bonus. Secondly, we evaluate whether it ensures policy invariance from the original MDP. Finally, we examine whether the method relies on prior knowledge. The results are demonstrated in [Table 3](#).

Algorithm	Objective	Episodic Count	Policy Invariance	Prior Knowledge
ICM	$-\log p(\phi(s_{t+1})   \phi(s_t), a_t)$	✗	✗	✗
RND	$-\log p(\phi(s_{t+1})   s_{t+1})$	✗	✗	✗
RIDE	$\ \phi(s_{t+1}) - \phi(s_t)\ _2 / \sqrt{N_{ep}(s_{t+1})}$	✓	✗	✗
NGU	$\approx \alpha \cdot b_{\text{RND}}(s) / \sqrt{N_{ep}(s_{t+1})}$	✓	✗	✗
PBRS	$\gamma \Phi(s_{t+1}) - \Phi(s_t)$	✗	✓	✓
LIBERTY (ours)	$\gamma d_{inv}(s_{t+1}, s_0) - d_{inv}(s_t, s_0)$	✗	✓	✗

Table 3: Summary and comparison of several benchmark methods.  $\phi = f(s)$ : features encoder;  $\theta$ : model parameters;  $b_{\text{RND}}(s)$ : the RND bonus;  $N_{ep}(s)$ : episodic (pseudo) count of visits to state  $s$ ;  $\alpha$ : normalized coefficient;  $\Phi(s)$ : potential function defined by prior knowledge;  $d_{inv}(s, s_0)$ : the inverse dynamic bisimulation metric-based potential function.

### E.2 Evaluation Setup

We use the encoder architecture in [Zhang et al. \[2020\]](#) with two more convolutional layers to the convnet trunk as the potential function network. The encoder has kernels of size  $3 \times 3$  with 32 channels for all the convolutional layers and set stride to 1 everywhere, except of the first convolutional layer, which has stride 2, and interpolate with ReLU activations. Finally, we add tanh nonlinearity to the 50 dimensional output of the fully-connected layer. Each deconvolutional layer has kernels of size  $3 \times 3$

with 32 channels and stride 1 , except of the last layer, where stride is 2 . We use ReLU activations after each layer. The forward dynamic model are MLPs with two hidden layers with 200 neurons each and ReLU activations. The inverse dynamic module has 3 hidden layers with 64 neurons with ReLU activation functions. The output is a 64 dimensional embedding feature. Adam is chosen as the optimizer with learning rate  $\alpha=0.0001$ .

### E.3 Algorithm Details

---

#### Algorithm 1 LIBERTY

---

```

1: while not converged do
2:   for  $t = 1$  to MAX_STEP_PER_EPISODE do
3:     Sample action  $a_t \sim \pi_\theta(a_t | s_t)$ 
4:     Step environment  $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$ 
5:     Record transition in the buffer:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, s_{t+1}, r_{t+1}\}$ 
6:     Sample transitions batch  $B \sim \mathcal{D}$  from the buffer
7:     Permute batch:  $\hat{B} = \text{permute}(B)$ .
8:     Set shaping intrinsic reward  $F(s_t, a, s_{t+1}) = \gamma d_{inv}(s_{t+1}, s_0) - d_{inv}(s_t, s_0)$ 
9:     Train policy  $\pi(\theta)$  via policy gradient:  $\mathbb{E}_B[J(\theta)]$ 
10:    Train the potential function  $d_{inv}$ :  $\mathbb{E}_{B, \hat{B}}[J(\phi)]$  ▷Equation (5)
11:    Train inverse dynamics:  $J(\theta_I)$  ▷Equation (3)
12:    Train forward dynamics:  $J(\eta) = (\mathcal{P}(\cdot | s_t, a_t; \eta) - s_{t+1})^2$ 
13:  end for
14: end while

```

---

### E.4 Codebase Used

Our codebase was built atop the following codebases:

- The official NGU codebase (including the implementation of ICM and RND): <https://github.com/opendilab/DI-engine>
- The official RIDE codebase: <https://github.com/facebookresearch/impact-driven-exploration>

#### E.4.1 Hyperparameters of Benchmarks

Algorithm	Hyperparameters
Base Learner (PPO)	policy module: three 64 -unit FC hidden layers, relu activation threshold of probability ratio clipping ( $\epsilon$ ) : 0.2 update timesteps: 1e6 number of epoches per update: 50 number of minibatches: 4 batch size: 1024 GAE parameter ( $\lambda$ ) : 0.95 optimizer: Adam learning rate: $10^{-4}$ policy gradient clip norm: 1.0 discount rate ( $\gamma$ ) : 0.999
DPBA	potential network: two 64-unit FC layers, tanh activation optimizer: Adam learning rate of potential network: $5 \times 10^{-4}$ potential network gradient clip norm: 10.0
ICM	forward model loss coefficient: 1.0 inverse model loss coefficient: 0.1 entropy cost: 0.005 intrinsic reward coefficient $\beta$ : 0.01
RND	proportion of experience used for training predictor: 0.05 predictor Model updates per PPO epoch: 3 entropy cost: 0.005 intrinsic reward coefficient $\beta$ : 0.01
NGU	episodic memory capacity: 2000 action prediction network filter sizes: (3,3) action prediction network filter strides: (1, 1) RND clipping factor $L$ : 3 intrinsic reward coefficient $\beta$ : 0.1
RIDE	forward model loss coefficient: 1.0 inverse model loss coefficient: 0.1 entropy cost: 0.0005 intrinsic reward coefficient $\beta$ : 0.1

Table 4: The hyperparameters of the tested benchmark algorithms in the MuJoCo experiment

Algorithm	Hyperparameters
Base Learner (PPO)	<p>policy network: two CNN with 3 convolutional layers and 1 FC layer with 512 hidden units</p> <p>the 1st convolutional layer has thirty-two <math>8 \times 8</math> filters with stride 4</p> <p>the 2nd convolutional layer has sixty-four <math>4 \times 4</math> filters with stride 2</p> <p>The 3rd convolutional layer has sixty-four <math>3 \times 3</math> filters with stride 1</p> <p>threshold of probability ratio clipping (<math>\epsilon</math>) : 0.5</p> <p>update timesteps: <math>5e7</math></p> <p>number of epoches per update: 50</p> <p>number of minibatches: 16</p> <p>batch size: 2048</p> <p>GAE parameter (<math>\lambda</math>) : 0.95</p> <p>optimizer: Adam</p> <p>learning rate: <math>7 \times 10^{-4}</math></p> <p>policy gradient clip norm: 0.9</p> <p>discount rate (<math>\gamma</math>) : 0.999</p>
DPBA	<p>potential network: 32-unit FC layer + 16-unit FC layer</p> <p>tanh activation</p> <p>optimizer: Adam</p> <p>learning rate of potential network: <math>5 \times 10^{-4}</math></p> <p>potential network gradient clip norm: 1.0</p>
ICM	<p>forward model loss coefficient: 0.2</p> <p>inverse model loss coefficient: 0.08</p> <p>entropy cost: <math>10^{-4}</math></p> <p>intrinsic reward coefficient <math>\beta</math>: 0.005</p>
RND	<p>proportion of experience used for training predictor: 0.25</p> <p>predictor Model updates per PPO epoch: 6</p> <p>entropy cost: 0.005</p> <p>intrinsic reward coefficient <math>\beta</math>: 0.1</p>
NGU	<p>episodic memory capacity: 5000</p> <p>action prediction network filter sizes: (3,3)</p> <p>action prediction network filter strides: (1, 1)</p> <p>RND clipping factor <math>L</math>: 5</p> <p>intrinsic reward coefficient <math>\beta</math>: 0.01</p>
RIDE	<p>forward model loss coefficient: 0.5</p> <p>inverse model loss coefficient: 0.8</p> <p>entropy cost: 0.0005</p> <p>intrinsic reward coefficient <math>\beta</math>: 0.1</p>

Table 5: The hyperparameters of the tested benchmark algorithms in the Atari experiment