# Efficient Potential-based Exploration in Reinforcement Learning using Inverse Dynamic Bisimulation Metric

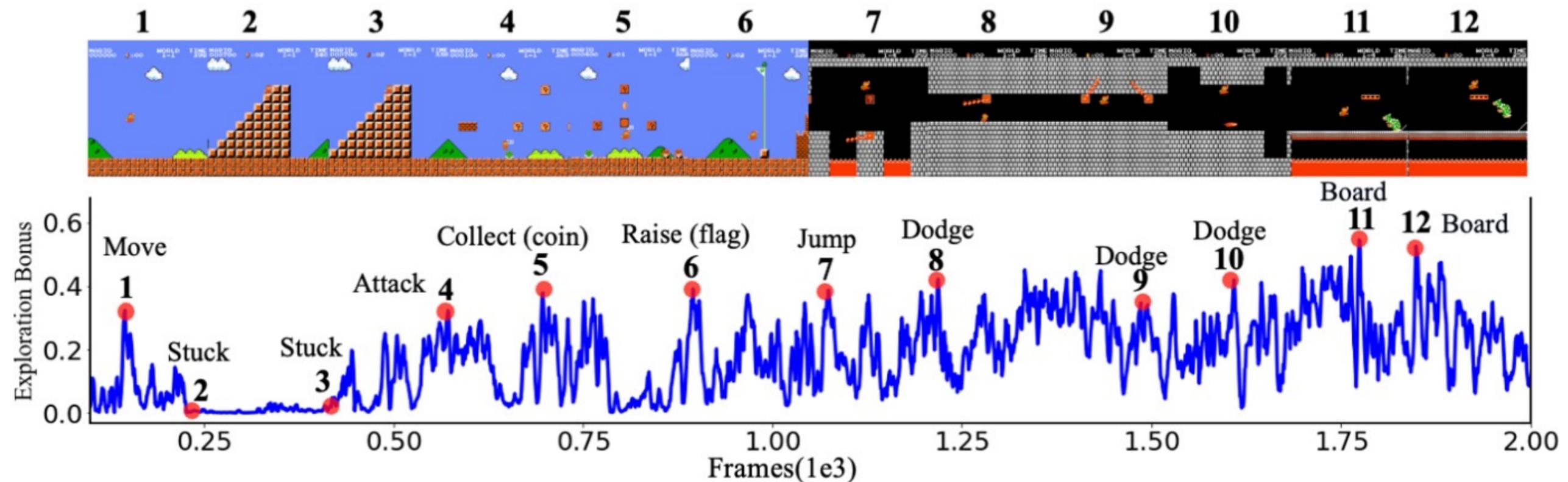Yiming Wang[1], Ming Yang[1], Renzhi Dong[1], Binbin Sun[1], Furui Liu[2], Leong Hou U[1*]

[1]State Key Laboratory of Internet of Things for Smart City, University of Macau, Macao SAR, China

[2]Zhejiang Lab, Hangzhou, China

# Introduction to exploration in RL

- Reward shaping methods

  - Rely on human prior knowledge
  - Introduce human cognitive biases

- Curiosity-driven exploration methods

  - Lack of scalability
  - Rely on count-based episodic term
  - Cause policy variance of original MDP

- LIBERTY: exp**L**oration v**I**a **B**isimulation m**Et**R**ic-based s**T**ate discrepanc**Y**

  - Our method (LIBERTY) uses the bisimulation metric to measure state discrepancy and propose a potential function based on the inverse dynamic bisimulation metric, which promotes effective exploration while preserving the optimal policy of the original MDP

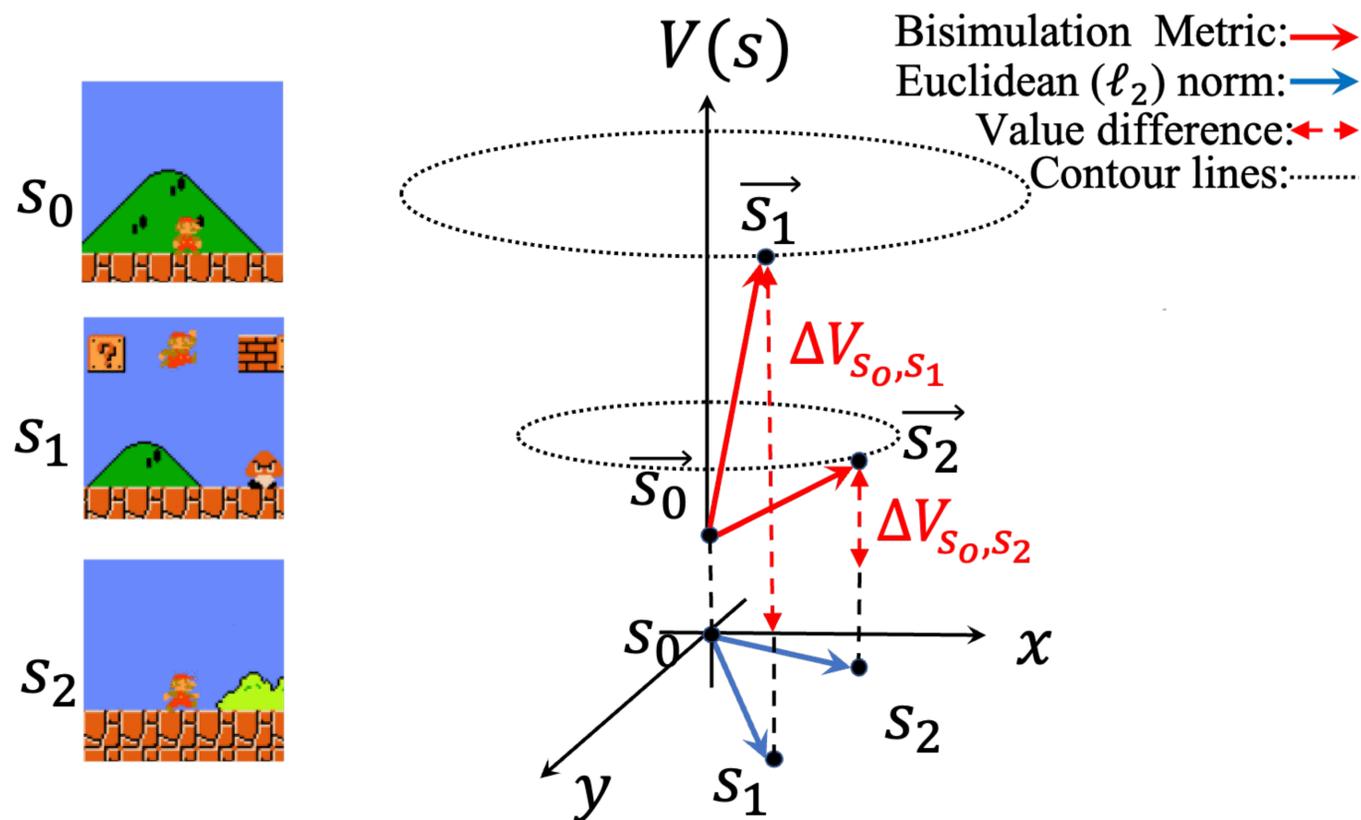# Motivation: state discrepancy as exploration bonus



- Using the difference between states $d(s_i, s_j)$ as exploration bonus
- Many spikes are related to significant occurrences, e.g., moving forward (1), attacking enemies (4), collecting coins (5) etc. The reward is close to 0 when the agent is stuck (2,3)

# Method: bisimulation metric measuring state discrepancy

- Bisimulation metric:

$$d_\pi(s_i, s_j) = \left| r_i^\pi - r_j^\pi \right| + \gamma W_1(d_\pi)\left( \mathcal{P}^\pi(\cdot \mid s_i), \mathcal{P}^\pi(\cdot \mid s_j) \right)$$



> Project the state into 3D latent space and Z axis denotes value
> Bisimulation metric identifies the value differences between states, enabling the agent to reach state $s_1$ with a higher value compared to $s_2$, starting from initial state $s_0$

# Method: inverse dynamic bisimulation metric

- **Meaningless exploration**: state difference is caused by background changing without taking actions



Get reward

No actions

- Add inverse dynamic module ($I: S \times S \rightarrow A$) to avoid meaningless exploration

$$d_{inv}(s_i, s_j) = \left| r_i^\pi - r_j^\pi \right| + \gamma W_2(d_{inv})\left( \mathcal{P}^\pi(\cdot \mid s_i), \mathcal{P}^\pi(\cdot \mid s_j) \right)$$
$$+ \gamma \left\| I(\cdot \mid s_i, s_{i+1}) - I(\cdot \mid s_j, s_{j+1}) \right\|_1$$

# Method: inverse dynamic bisimulation metric (cont.)
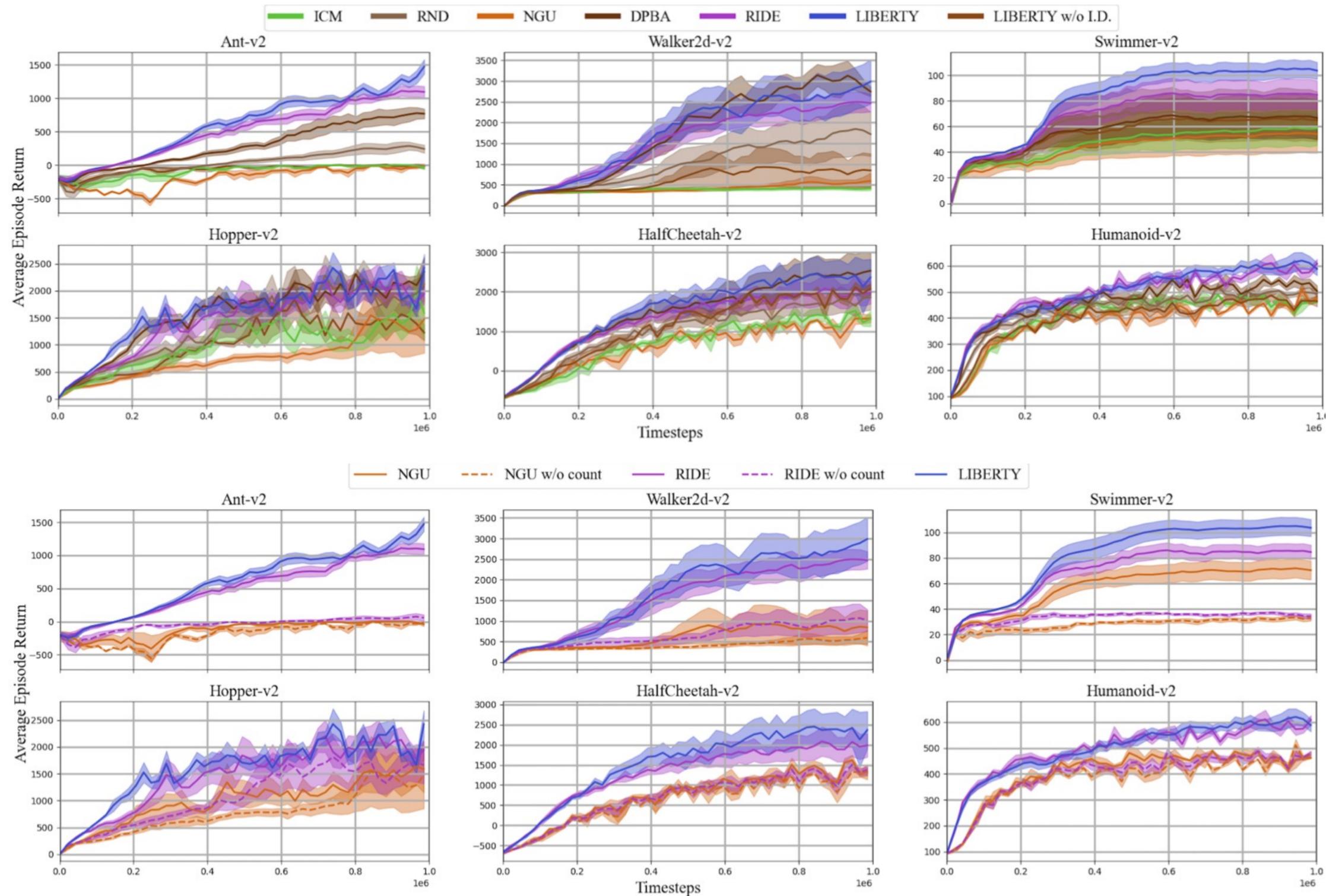
- Potential function based on $d_{inv}$

$$\Phi(s) = d_{inv}(s, s_0)$$

- Potential-based shaping reward function:

$$\mathcal{F}(s_t, a, s_{t+1}) = \gamma d_{inv}(s_{t+1}, s_0) - d_{inv}(s_t, s_0)$$
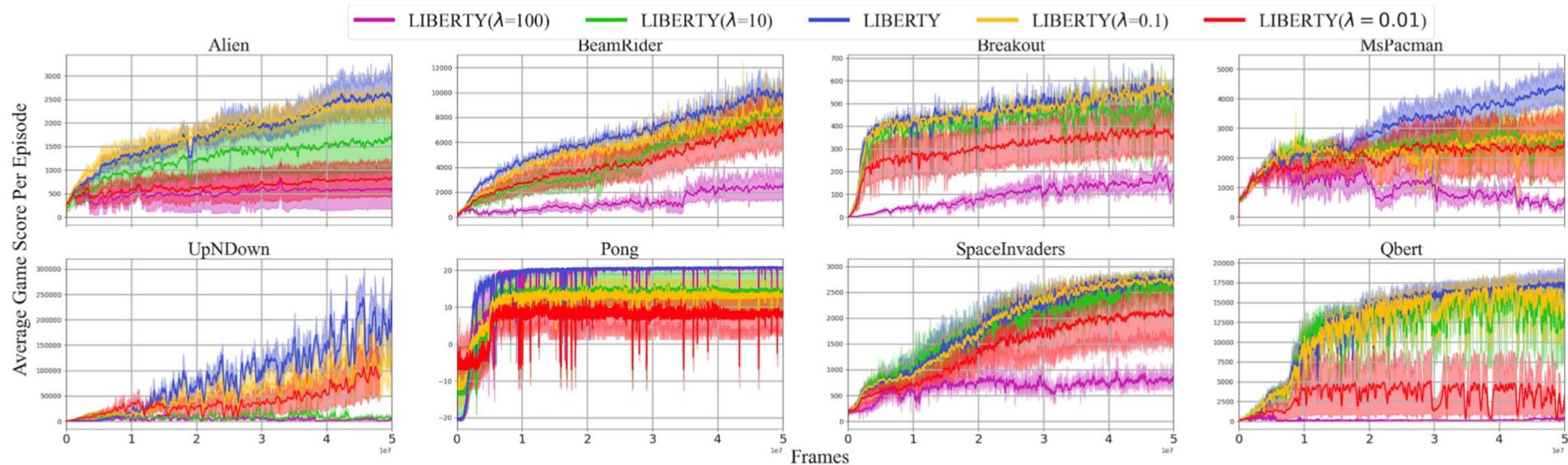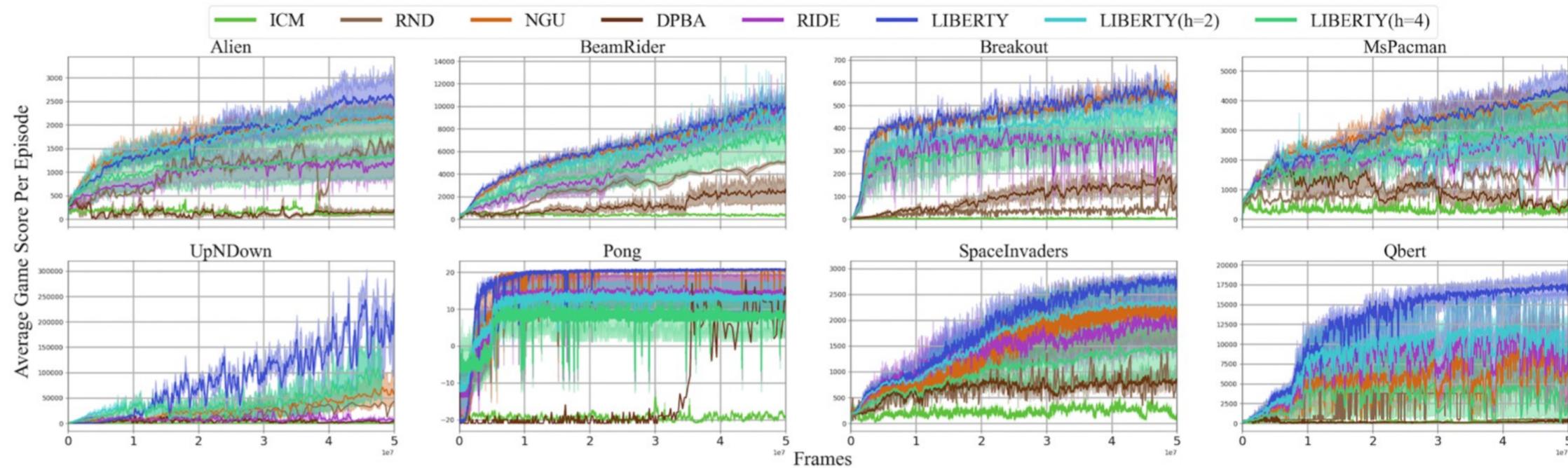
# Experiment

- Results on MuJoCo continuous control

# Experiment

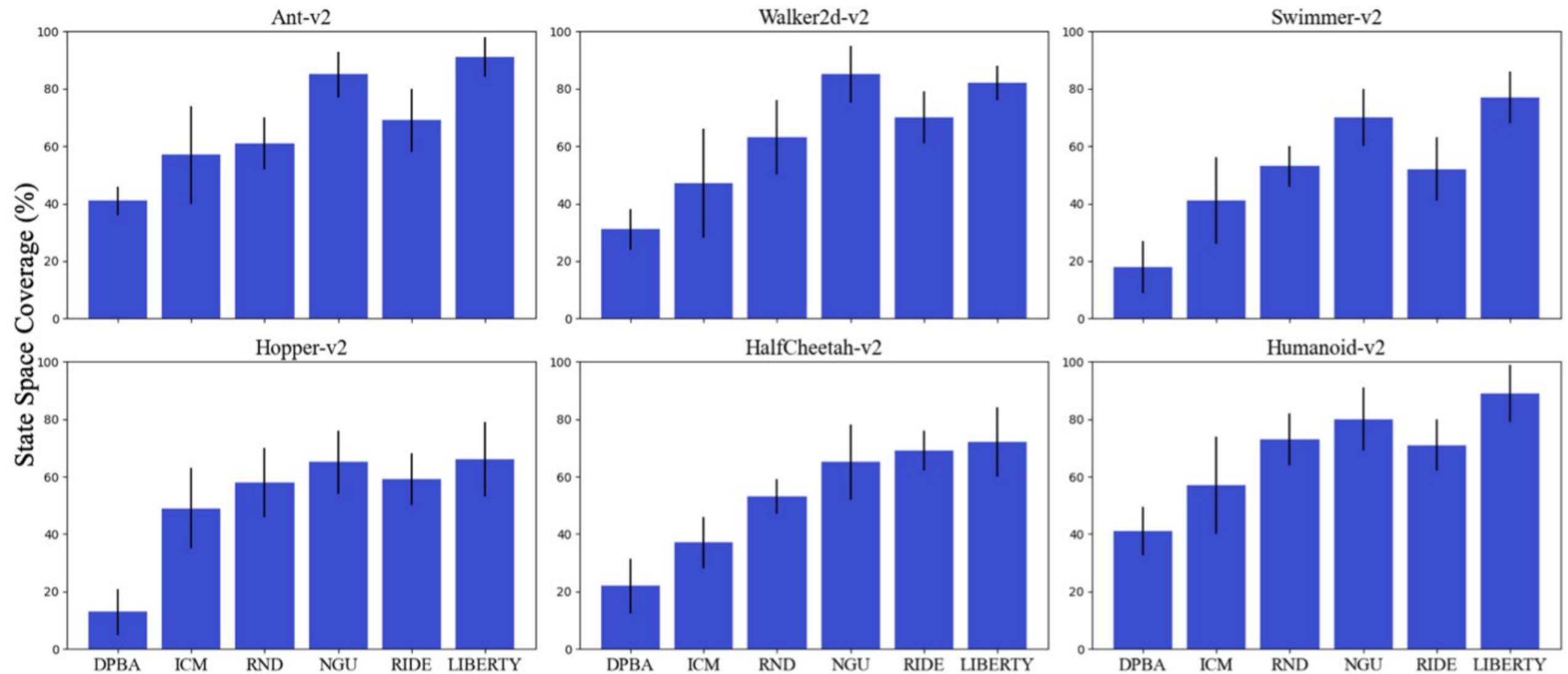- Results on Atari games with discrete actions

# Experiment

- Results on the delayed reward setting

Table 1: 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)

| Methods | Delay = 10 | | | | | |
|---|---|---|---|---|---|---|
| | 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 | $\underline{2467 \pm 456}$ | $1876 \pm 431$ | $1651 \pm 325$ | $92 \pm 31$ | $\underline{570 \pm 45}$ | $\underline{65 \pm 16}$ |
| DPBA | $1514 \pm 365$ | $\underline{2103 \pm 129}$ | $\underline{1997 \pm 115}$ | $\mathbf{592 \pm 67}$ | $518 \pm 23$ | $43 \pm 17$ |
| LIBERTY | $\mathbf{2973 \pm 437}$ | $\mathbf{2479 \pm 315}$ | $\mathbf{2766 \pm 487}$ | $\underline{292 \pm 68}$ | $\mathbf{681 \pm 73}$ | $\mathbf{73 \pm 21}$ |
| LIBERTY w/o I.D. | $1783 \pm 412$ | $1676 \pm 275$ | $1732 \pm 392$ | $131 \pm 22$ | $505 \pm 37$ | $46 \pm 11$ |

| Methods | Delay = 40 | | | | | |
|---|---|---|---|---|---|---|
| | 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 | $1276 \pm 387$ | $\mathbf{1683 \pm 338}$ | $968 \pm 168$ | $71 + 15$ | $\underline{483 \pm 25}$ | $17 \pm 11$ |
| NGU | $1028 \pm 405$ | $879 \pm 155$ | $997 \pm 280$ | $-198 \pm 27$ | $387 \pm 27$ | $11 \pm 6$ |
| RIDE | $\underline{1798 \pm 355}$ | $1235 \pm 269$ | $\underline{1025 \pm 282}$ | $63 \pm 18$ | $468 \pm 23$ | $\mathbf{32 \pm 11}$ |
| DPBA | $883 \pm 275$ | $1382 \pm 85$ | $1016 \pm 129$ | $\underline{105 \pm 31}$ | $405 \pm 15$ | $9 \pm 3$ |
| LIBERTY | $\mathbf{2039 \pm 315}$ | $\underline{1612 \pm 215}$ | $\mathbf{1921 \pm 372}$ | $\mathbf{142 \pm 45}$ | $\mathbf{566 \pm 35}$ | $\underline{31 \pm 13}$ |
| LIBERTY w/o I.D. | $1231 \pm 253$ | $1213 \pm 207$ | $1012 \pm 358$ | $58 \pm 13$ | $455 \pm 27$ | $17 \pm 8$ |

# Experiment

- Results on the reward-free setting

# Contribution

- We develop a new potential function to ensure policy invariance without the need for <span style="color:red">prior human knowledge</span>

- Our approach achieves <span style="color:red">more efficient</span> exploration by encouraging agents to explore states with higher TD-error

  ➢ Theorem 1(value difference bound):
  $$\left|V^{\pi}(s_i) - V^{\pi}(s_j)\right| \leq d_{inv}(s_i, s_j)$$
  ➢ Theorem 2(approximation of optimal value function):
  $$d_{inv}(s, s_0) \approx V^*(s)$$

- Our method achieves best performance across various settings in extensive environments, which demonstrate its <span style="color:red">scalability</span> and <span style="color:red">superiority</span> compared with other methods

# Thank you for your attention!

Code is available