



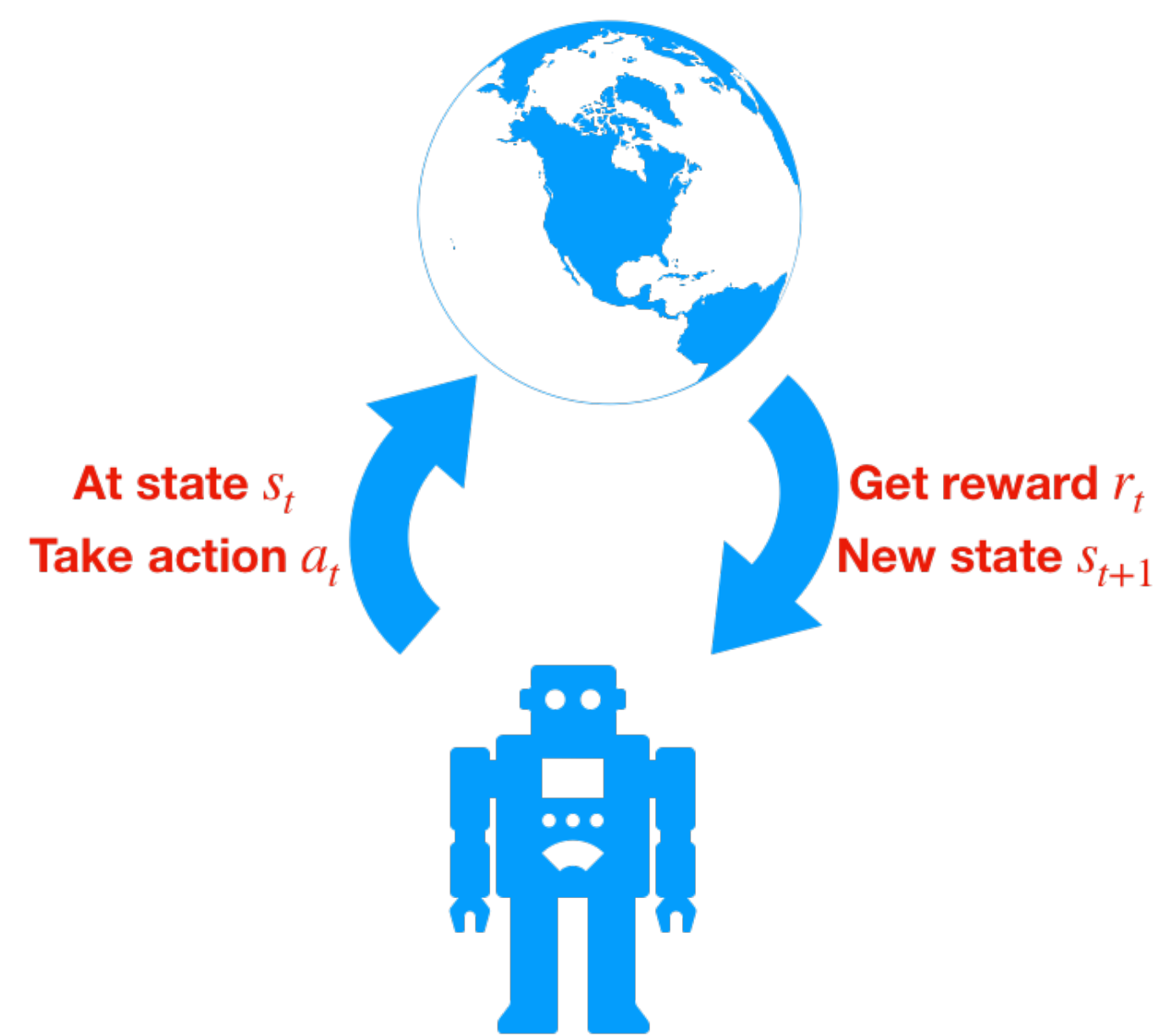
Sample Efficient Policy Gradient Methods with Recursive Variance Reduction

Pan Xu and Felicia Gao and Quanquan Gu
University of California, Los Angeles

Policy Gradient

Markov decision process (MDP)

- State $s \in \mathcal{S}$
- Action $a \in \mathcal{A}$
- Transition transition $P(s'|s, a)$
- Reward $r(s, a)$
- Policy $\pi_{\theta}(a|s)$



Goal: maximize the performance function

$$J(\theta) = \mathbb{E}_{\tau \sim p(\cdot|\theta)} \left[\sum_{h=0}^{H-1} \gamma^h r(s_h, a_h) \right]$$

- Trajectory: $\tau = (s_0, a_0, s_1, r_1, \dots, s_{H-1}, a_{H-1}, s_H)$
- Trajectory distribution: $p(\tau|\theta) = \rho(s_0) \prod_{h=0}^{H-1} \pi_{\theta}(a_h|s_h) P(s_{h+1}|s_h, a_h)$

GPOMDP policy gradient estimator:

$$g(\tau_i|\theta) = \sum_{h=0}^{H-1} \gamma^h r(s_h^i, a_h^i) \sum_{t=0}^h \nabla_{\theta} \log \pi_{\theta}(a_t^i|s_t^i) \approx \nabla J(\theta)$$

SRVR-PG

- Input: number of epochs S , epoch length m , step size η , batch size N , mini-batch size B , gradient estimator $g(\tau|\theta)$, initial parameter $\tilde{\theta}^0 = \theta_0$

```

01 for  $s = 0, \dots, S - 1$  do
02    $\theta_0^{s+1} = \theta^s$ 
03   sample  $N$  trajectories  $\{\tau_i\}$ 
04    $\mathbf{v}_0^{s+1} = \frac{1}{N} \sum_{i=1}^N g(\tau_i|\theta_0^{s+1})$ 
05    $\theta_1^{s+1} = \theta_0^{s+1} + \eta \mathbf{v}_0^{s+1}$ 
06   for  $t = 1, \dots, m - 1$  do
07     sample  $B$  trajectories  $\{\tau_j\}$ 
08      $\mathbf{v}_t^{s+1} = \mathbf{v}_{t-1}^{s+1} + \frac{1}{B} \sum_{j=1}^B (g(\tau_j|\theta_t^{s+1}) - g_w(\tau_j|\theta_{t-1}^{s+1}))$ 
09      $\theta_{t+1}^{s+1} = \theta_t^{s+1} + \eta \mathbf{v}_t^{s+1}$ 
10   end for
11 end for
12 output  $\theta_{\text{out}}$  uniformly from  $\{\theta_t^s\}$ 

```

Step-wise Importance Sampling

Policy gradient estimator

$$g_w(\tau|\theta_{t-1}^{s+1}) = \sum_{h=0}^{H-1} \omega_{0:h}(\tau) \left[\sum_{t=0}^h \nabla_{\theta} \log \pi_{\theta_{t-1}^{s+1}}(a_t|s_t) \right] \gamma^h r(s_h, a_h)$$

Importance weight

$$\omega_{0:h}(\tau) := \omega_{0:h}(\tau|\theta_{t-1}^{s+1}, \theta_t^{s+1}) = \prod_{k=0}^h \frac{\pi_{\theta_{t-1}^{s+1}}(a_k|s_k)}{\pi_{\theta_t^{s+1}}(a_k|s_k)}$$

Convergence Analysis

Assumptions

- Smoothness:** for all θ, s, a , the policy satisfies $\|\nabla_{\theta} \log \pi_{\theta}(a|s)\| \leq G$, $\|\nabla_{\theta}^2 \log \pi_{\theta}(a|s)\|_2 \leq M$
- Bounded variance:** there exists a constant $\xi > 0$ such that $\text{Var}(g(\tau|\theta)) \leq \xi^2$, for all policy π_{θ}
- Importance weight:** there is a constant $W < \infty$ such that $\omega(\cdot|\theta_1, \theta_2) = p(\cdot|\theta_1)/p(\cdot|\theta_2)$ satisfies $\text{Var}(\omega(\tau|\theta_1, \theta_2)) \leq W$, $\forall \theta_1, \theta_2 \in \mathbb{R}^d, \tau \sim p(\cdot|\theta_2)$

Convergence of SRVR-PG

Let step size $\eta \leq 1/(4L)$ and batch size $B \geq Cm\eta G^4 MW \gamma(1-\gamma)^{-3}$

$$\mathbb{E}[\|\nabla J(\theta_{\text{out}})\|_2^2] \leq \frac{8(J(\theta^*) - J(\theta_0))}{\eta S m} + \frac{6\xi^2}{N}$$

Remark: independent of H due to the step-wise importance weight

Comparison on Sample Complexity

Sample complexity

- total number of trajectories needed to achieve $\|\nabla J(\theta)\|_2^2 \leq \epsilon$
- Setting $N = O(1/\epsilon)$, $B = m = S = O(1/\epsilon^{1/2})$, we obtain the sample complexity of SRVR-PG: $O(1/\epsilon^{3/2})$

Algorithms	Complexity
REINFORCE (Williams, 1992)	$O(1/\epsilon^2)$
PGT (Sutton et al., 2000)	$O(1/\epsilon^2)$
GPOMDP (Baxter and Bartlett, 2001)	$O(1/\epsilon^2)$
SVRPG (Papini et al., 2018)	$O(1/\epsilon^2)$
SVRPG (Xu et al., 2019)	$O(1/\epsilon^{5/3})$
SRVR-PG	$O(1/\epsilon^{3/2})$

Experiments

Setup

- Environments: Cartpole, Mountain Car, Pendulum
- Deep Gaussian policy

$$\pi_{\theta}(a|s) = 1/\sqrt{2\pi} \exp(- (f(\theta; \phi(s)) - a)^2 / (2\sigma^2))$$

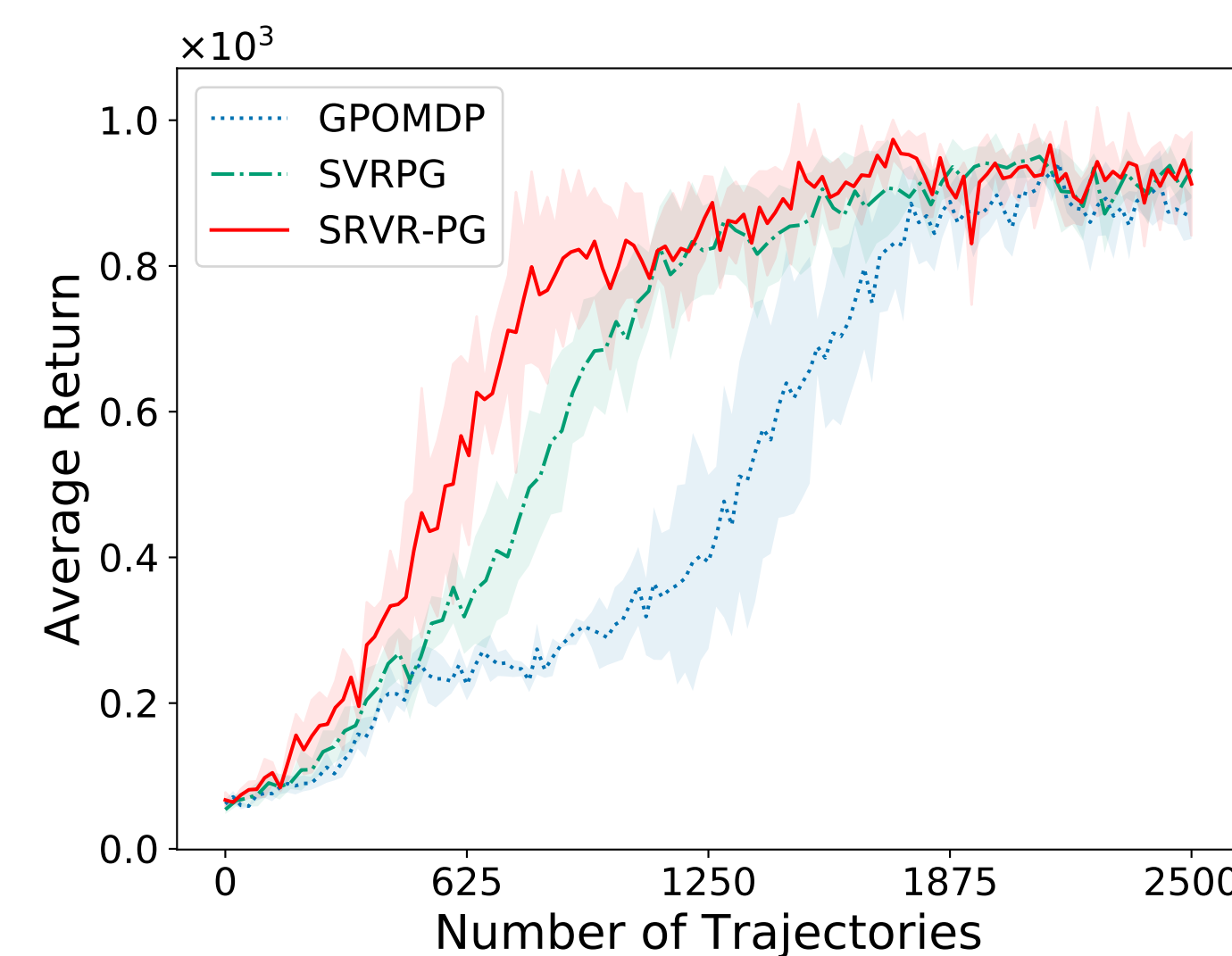
- σ^2 is a fixed standard deviation parameter
- $\phi: \mathcal{S} \mapsto \mathbb{R}^d$ is a neural network function.

Problem dependent parameters

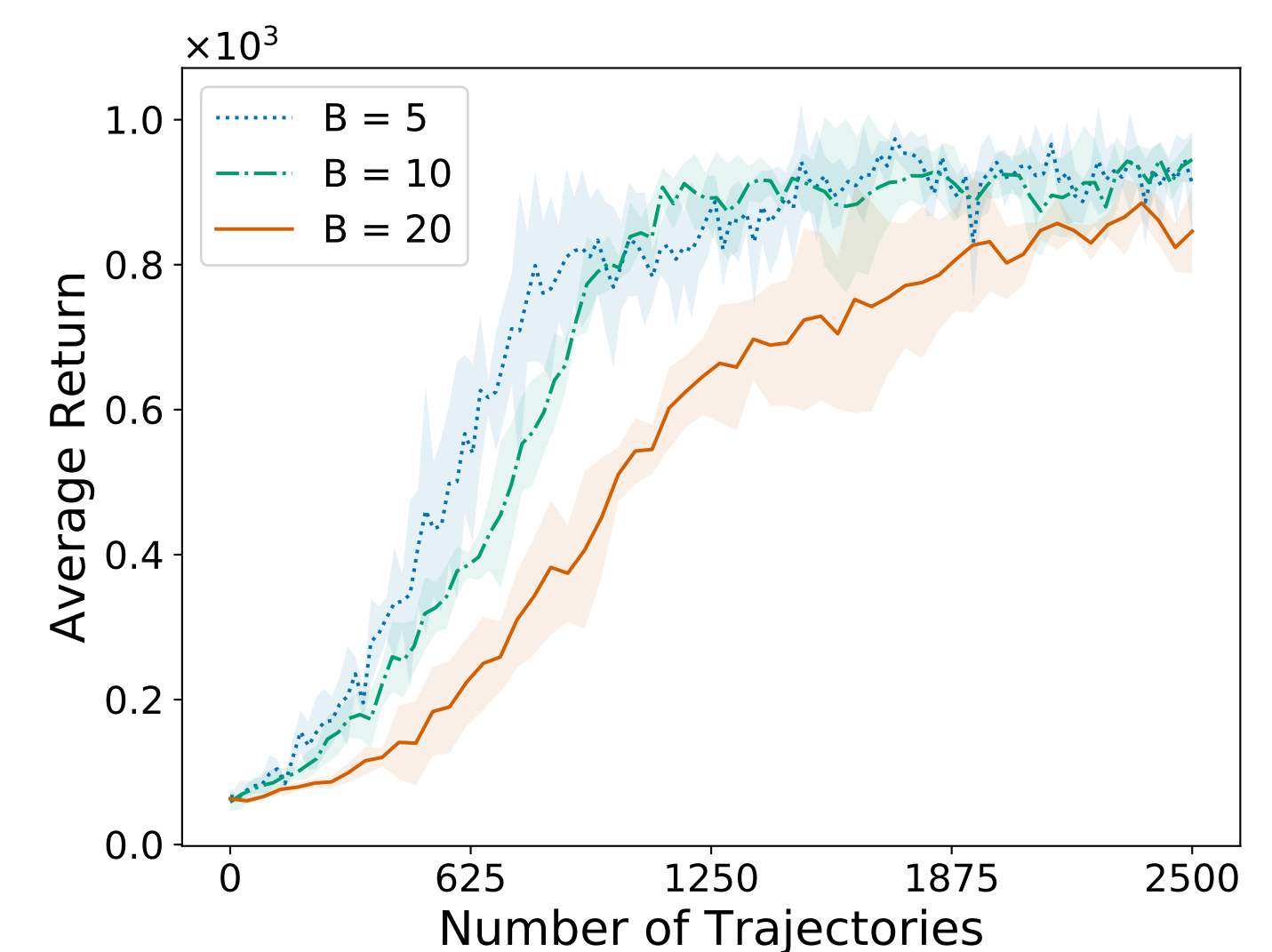
Parameters	Cartpole	Mountain Car	Pendulum
NN size	64	64	8×8
NN activation function	Tanh	Tanh	Tanh
Task horizon	100	1000	200
Total trajectories	2500	3000	2 × 10 ⁵

Comparison of different algorithms

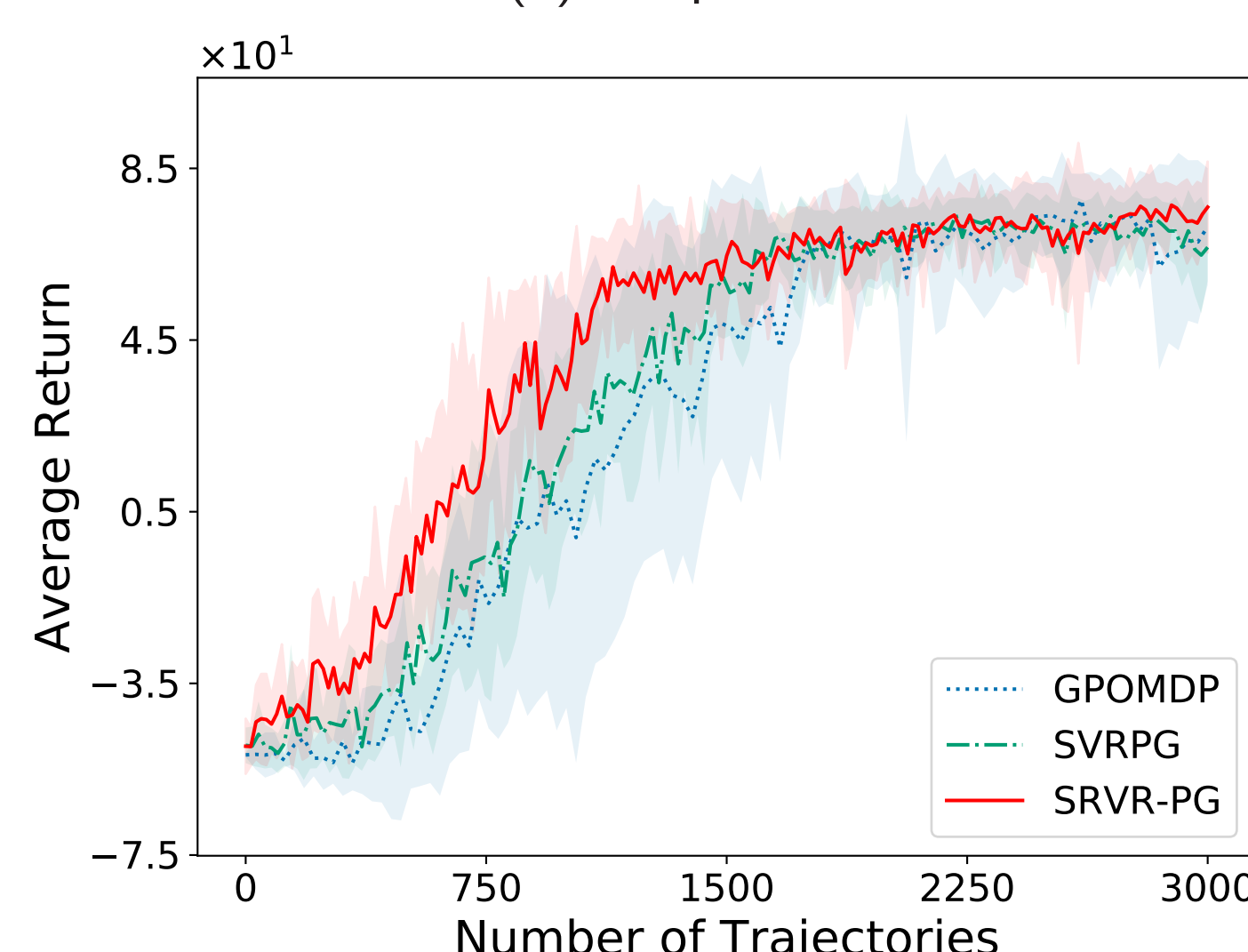
- Experimental results are averaged over 10 repetitions
- Figures (a), (c) and (e): comparison with baselines GPOMDP and SVRPG
- Figures (b), (d) and (f): comparison of different batch size B on the performance of SRVR-PG



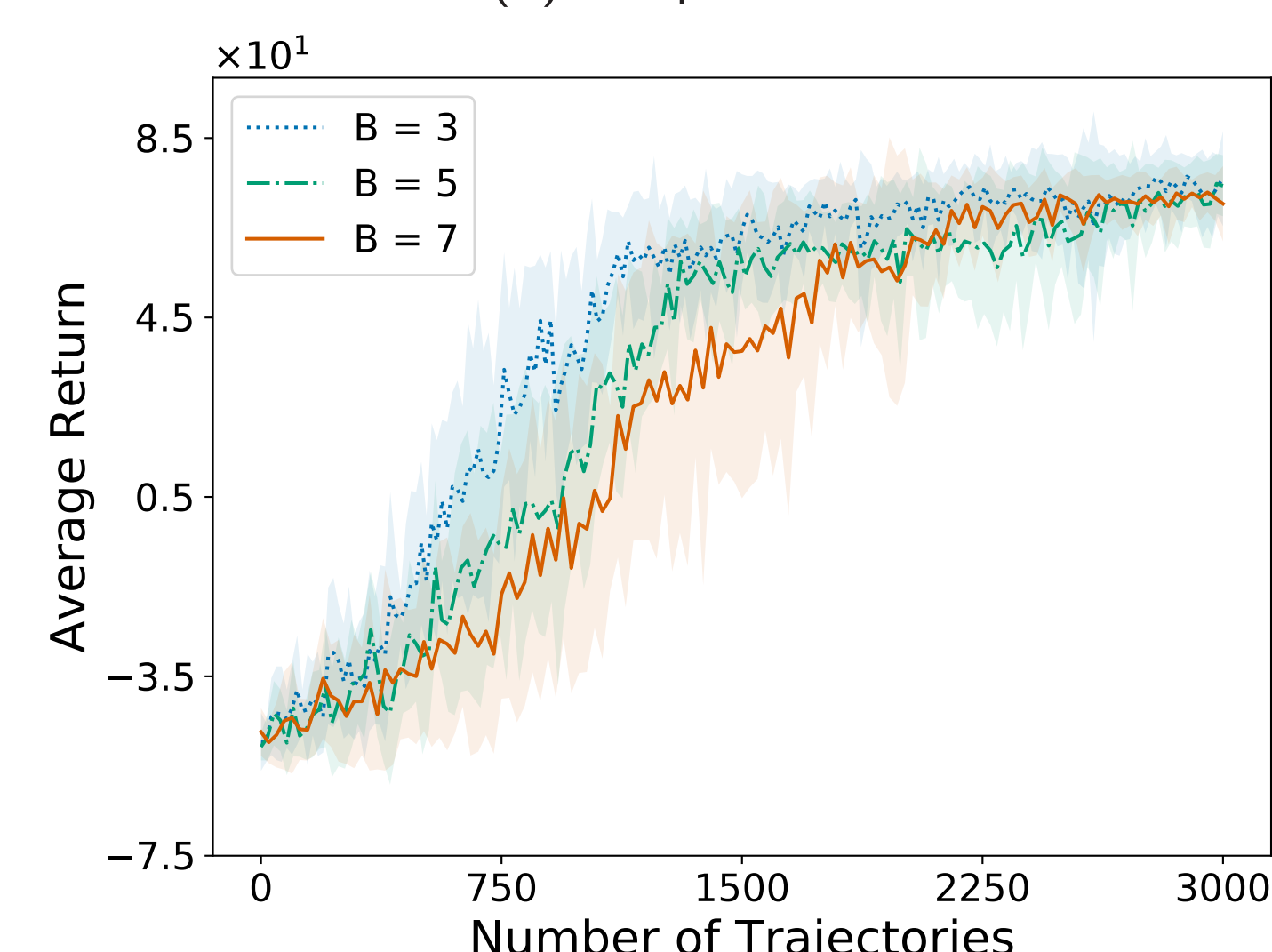
(a) Cartpole



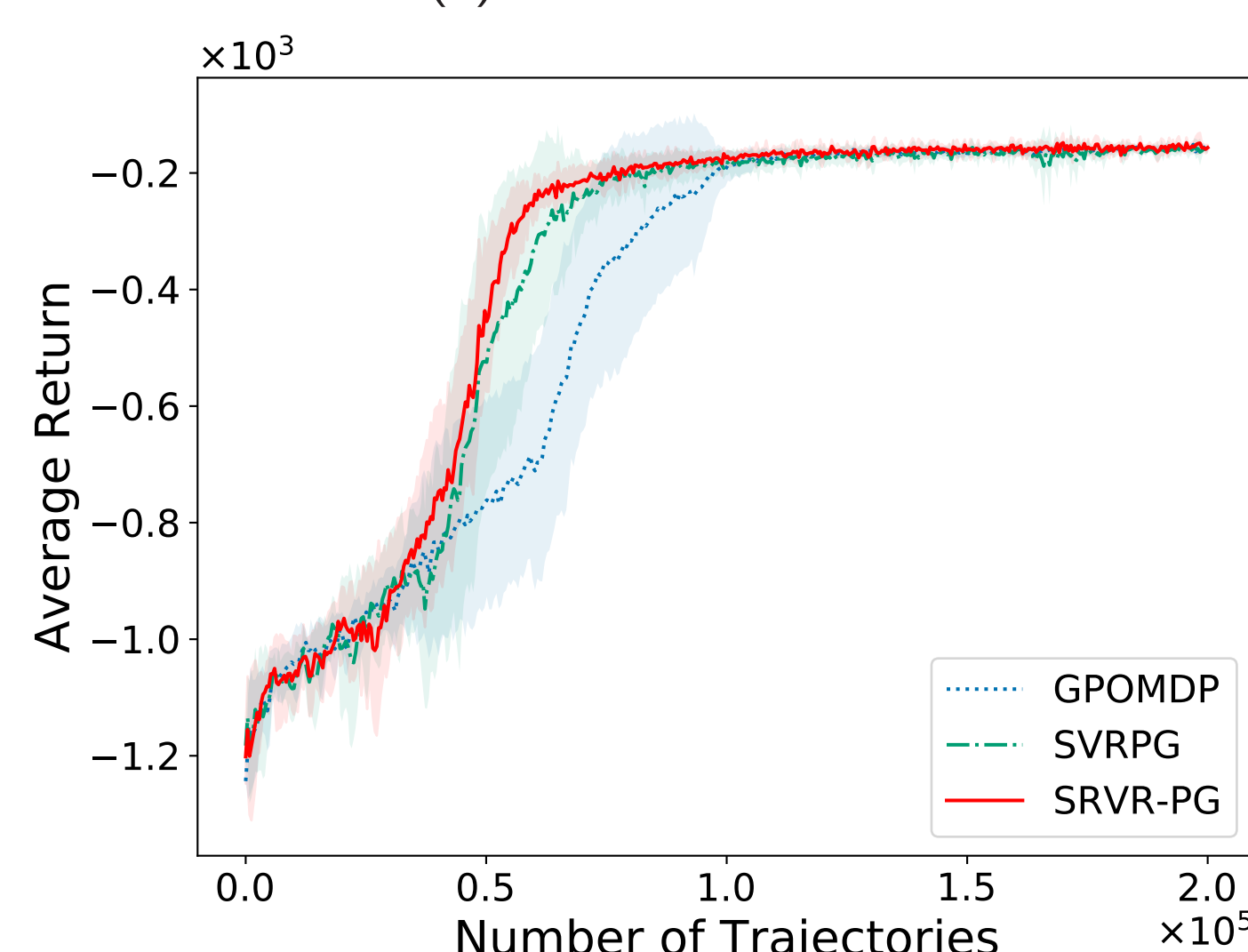
(b) Cartpole



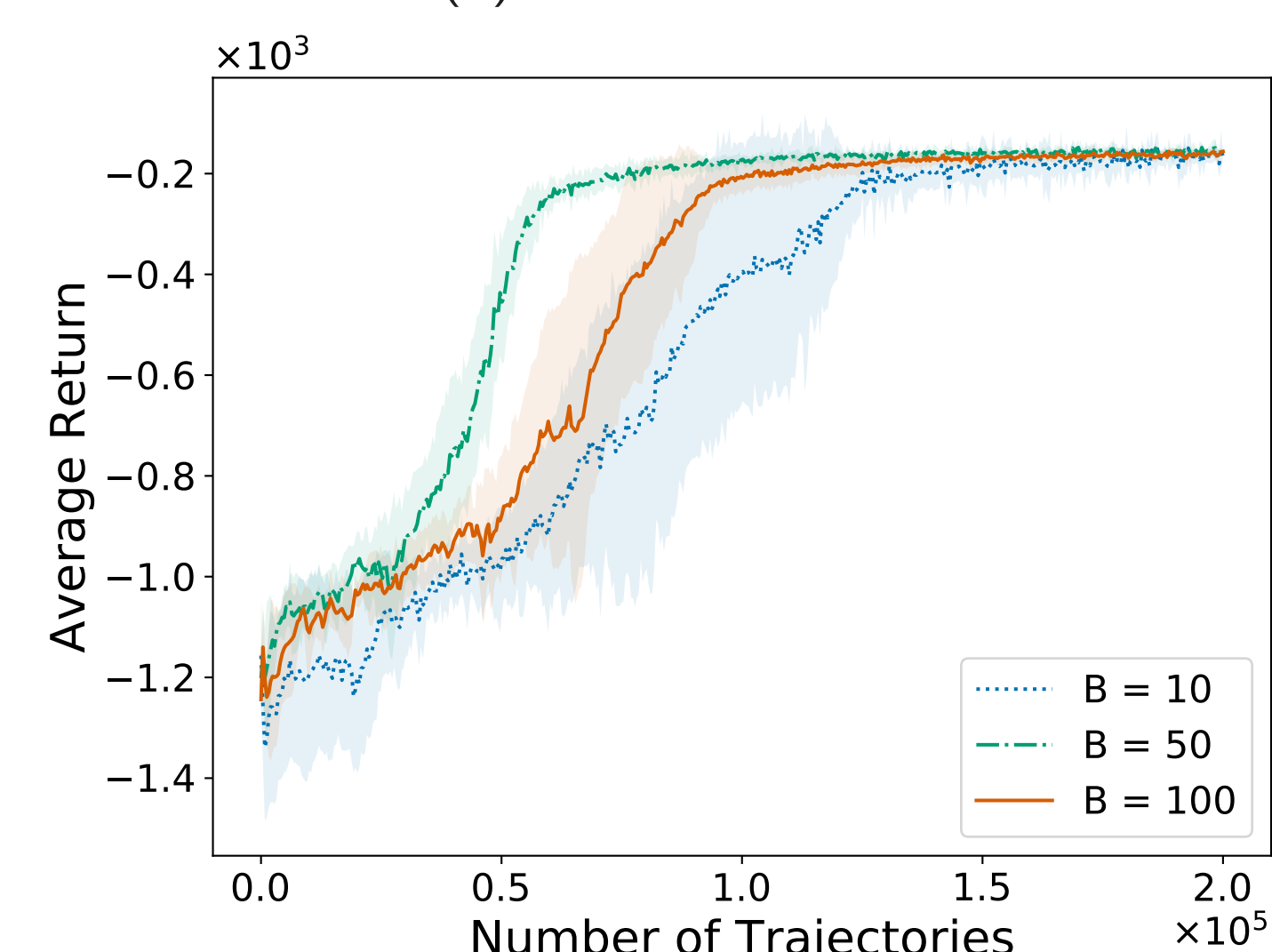
(c) Mountain Car



(d) Mountain Car



(e) Pendulum



(f) Pendulum