

# **SUNRISE**

## **A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning**

Lee, Kimin, Michael Laskin, Aravind Srinivas, Pieter Abbeel

UC Berkeley (ICML 2021)

발표자: 이동진  
2023년 03월 06일

## About authors 주저자



Kimin Lee

Research Scientist, [Google](#)

Verified email at google.com - [Homepage](#)

[Machine learning](#) [Reinforcement learning](#) [Deep learning](#)

 FOLLOW

TITLE	CITED BY	YEAR
<a href="#">A simple unified framework for detecting out-of-distribution samples and adversarial attacks</a> K Lee, K Lee, H Lee, J Shin Advances in neural information processing systems 31	1160	2018
<a href="#">Training confidence-calibrated classifiers for detecting out-of-distribution samples</a> K Lee, H Lee, K Lee, J Shin International Conference on Learning Representations	669	2017
<a href="#">Using pre-training can improve model robustness and uncertainty</a> D Hendrycks, K Lee, M Mazeika International Conference on Machine Learning, 2712-2721	506	2019
<a href="#">Reinforcement learning with augmented data</a> M Laskin, K Lee, A Stooke, L Pinto, P Abbeel, A Srinivas Advances in neural information processing systems	382	2020
<a href="#">Decision transformer: Reinforcement learning via sequence modeling</a> L Chen, K Lu, A Rajeswaran, K Lee, A Grover, M Laskin, P Abbeel, ... Advances in neural information processing systems	370	2021

## About authors 주저자



Kimin Lee

Research Scientist, [Google](#)

Verified email at google.com - [Homepage](#)

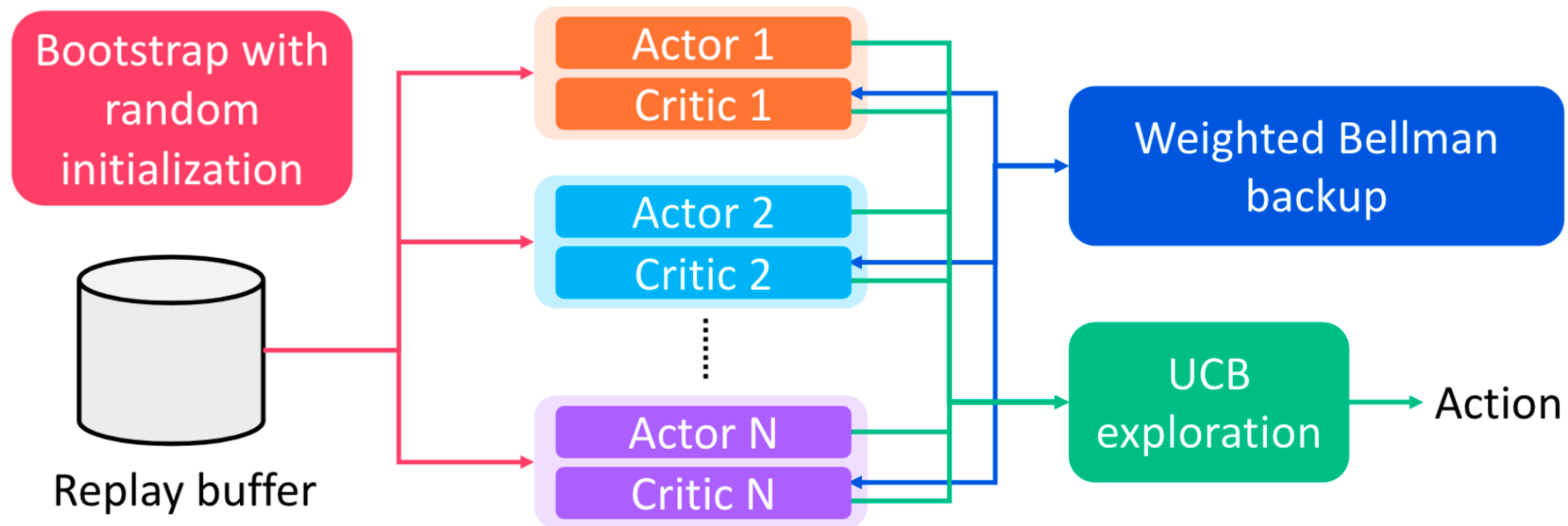
[Machine learning](#) [Reinforcement learning](#) [Deep learning](#)

 FOLLOW

TITLE	CITED BY	YEAR
<a href="#">Overcoming catastrophic forgetting with unlabeled data in the wild</a> K Lee, K Lee, J Shin, H Lee Proceedings of the IEEE/CVF International Conference on Computer Vision, 312-321	129	2019
<a href="#">Network randomization: A simple technique for generalization in deep reinforcement learning</a> K Lee, K Lee, J Shin, H Lee International Conference on Learning Representations	128	2019
<a href="#">Sunrise: A simple unified framework for ensemble learning in deep reinforcement learning</a> K Lee, M Laskin, A Srinivas, P Abbeel International Conference on Machine Learning, 6131-6141	111	2021
<a href="#">Robust inference via generative classifiers for handling noisy labels</a> K Lee, S Yun, K Lee, H Lee, B Li, J Shin International conference on machine learning, 3763-3772	86	2019
<a href="#">Context-aware dynamics model for generalization in model-based reinforcement learning</a> K Lee, Y Seo, S Lee, H Lee, J Shin International Conference on Machine Learning, 5757-5766	63	2020

# 논문을 선택한 동기

- 지도학습에서 앙상블은 성능 향상의 치트키
- 강화학습에서 앙상블이란 무엇일까?
  - Replay buffer를 공유하는  $N$ 개의 actor와  $N$ 개의 critic



(a) SUNRISE: actor-critic version

Quoted from Kimin Lee et al., SUNRISE: A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning. (2021)

# SUNRISE Overview

## ➤ 알고리즘 분류

- Model-free, off-policy learning, continuous/discrete action space

# SUNRISE Overview

## ➤ 알고리즘 분류

- Model-free, off-policy learning, continuous/discrete action space

## ➤ 기반 알고리즘

- Soft Actor-Critic (연속 행동 공간), Rainbow (이산 행동 공간)
- $N$ 개의 actor와  $N$ 개의 critic

# SUNRISE Overview

## ➤ 알고리즘 분류

- Model-free, off-policy learning, continuous/discrete action space

## ➤ 기반 알고리즘

- Soft Actor-Critic (연속 행동 공간), Rainbow (이산 행동 공간)
- $N$ 개의 actor와  $N$ 개의 critic

## ➤ 사용한 환경

- Mujoco, DeepMind Control Suite, Atari

# N개의 actor와 critic이 있으면 좋은 점

- 현재 주어진 상태  $s_t$ 에서  $N$ 개의 행동 후보  $\{a_{t,i}\}_{i=1}^N$  생성
- 각 행동마다  $N$ 개의 행동 가치 함수가 있어서 분포를 생각할 수 있음
  - Weighted Bellman backups
  - Bootstrap with random initialization
  - UCB exploration



# Single model-free/off-policy 알고리즘의 문제점

- Sample inefficient
  - 경험 데이터를 충분히 많이 활용하지 못함
- Error propagation in Q-learning
  - Function approximation error 및 overestimation을 내재한 타겟을 사용하여 행동 가치 함수 학습
- Exploration Exploitation trade-off
  - 주로 랜덤에 의한 exploration을 함

# 앙상블을 사용한 문제 완화

## ➤ Sample inefficient

- 경험 데이터를 충분히 많이 활용하지 못함 ➡ 경험 데이터를 많은 네트워크 학습에 사용

## ➤ Error propagation in Q-learning

- Function approximation error 및 overestimation을 내재한 타겟을 사용하여 행동 가치 함수 학습 ➡ Double Q-learning으로 완화 (예) TD3, Double DQN

## ➤ Exploration Exploitation trade-off

- 주로 랜덤에 의한 exploration을 함 ➡ 여러 네트워크로부터 예측 불확실성을 고려

# Weighted Bellman backups (1)

- 일반적인 soft Q-learning

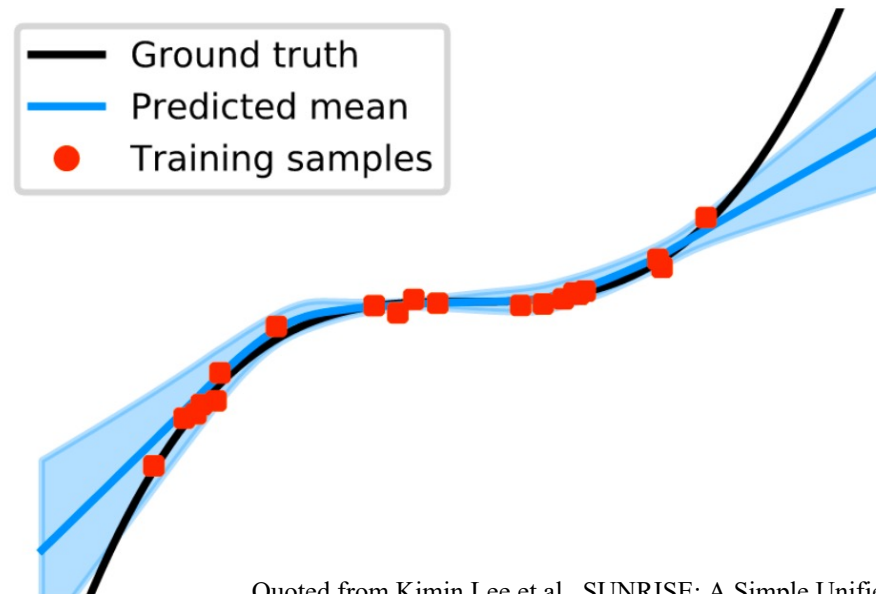
$$\mathcal{L}_Q(\theta) = (Q_\theta(s_t, a_t) - r_t - \gamma \bar{V}(s_{t+1}))^2$$

$$\bar{V}(s_t) = \mathbb{E}_{a_t \sim \pi_\phi} [Q_{\bar{\theta}}(s_t, a_t) - \alpha \log \pi_\phi(a_t | s_t)]$$

- Function approximation error가 내재된 target을 사용
  - ➡ inconsistency and unstable convergence

# Weighted Bellman backups (2)

- $N$ 개의 actor와 critic  $\{Q_{\theta_i}, \pi_{\phi_i}\}_{i=1}^N$  ※ 실험에서는  $N = 5$
- 행동 가치 함수 추정치  $\{Q_{\theta_i}(s, a)\}_{i=1}^N$ 의 분포를 고려
- Q-네트워크간 분산이 클수록 불확실하고 에러를 갖고 있는 추정치로 생각



Quoted from Kimin Lee et al., SUNRISE: A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning. (2021)

# Weighted Bellman backups (3)

➤ Weighted Bellman backups

$$\mathcal{L}_{WQ}(\theta_i) = w(s_{t+1}, a_{t+1}) \left( Q_{\theta_i}(s_t, a_t) - r_t - \gamma \bar{V}(s_{t+1}) \right)^2,$$
$$w(s, a) = \sigma(-\bar{Q}_{\text{std}}(s, a) * T) + 0.5,$$

- $a_{t+1} \sim \pi_{\phi_i}(a|s_{t+1})$
- $\sigma$ : 시그모이드 함수
- $\bar{Q}_{\text{std}}(s, a)$ 는  $Q_{\bar{\theta}_i}(s_t, a_t)$ 의 표본 표준편차
- $T$ : temperature

➤  $\bar{Q}_{\text{std}}(s, a)$ 가 클수록 Q-네트워크 학습에 적은 가중치를 부여

# Bootstrap with random initialization

- 에이전트간 다양성 (diversity) 확보 방법 2가지
- 1. 파라미터 랜덤 초기화  $\{\theta_i, \phi_i\}_{i=1}^N$
- 2. Transition masking을 통해 서로 다른 데이터를 사용
  - Transition  $\tau_t = (s_t, a_t, r_t, s_{t+1})$  저장할 때, 에이전트마다  $m_{t,i} \sim \text{Bernoulli}(\beta)$ 도 함께 저장
  - Replay buffer는 공유하지만, 에이전트마다 각 transition을 사용할 수도 있고 아닐 수도 있음
  - 하지만, 실험에서는  $\beta = 1$ 일 때, 즉 데이터를 다 사용할 때 성능이 제일 좋았음.

# UCB exploration (1)

## ➤ Upper confidence bound 방법

- Multi-armed bandit (MAB) 문제에서 밴딧이 줬던 보상의 평균 뿐만 아니라 밴딧을 선택한 횟수 까지 고려하여 밴딧을 선택하는 알고리즘.

- $$A_t = \operatorname{argmax}_a Q(a) + c \sqrt{\log t / N_t(a)}$$

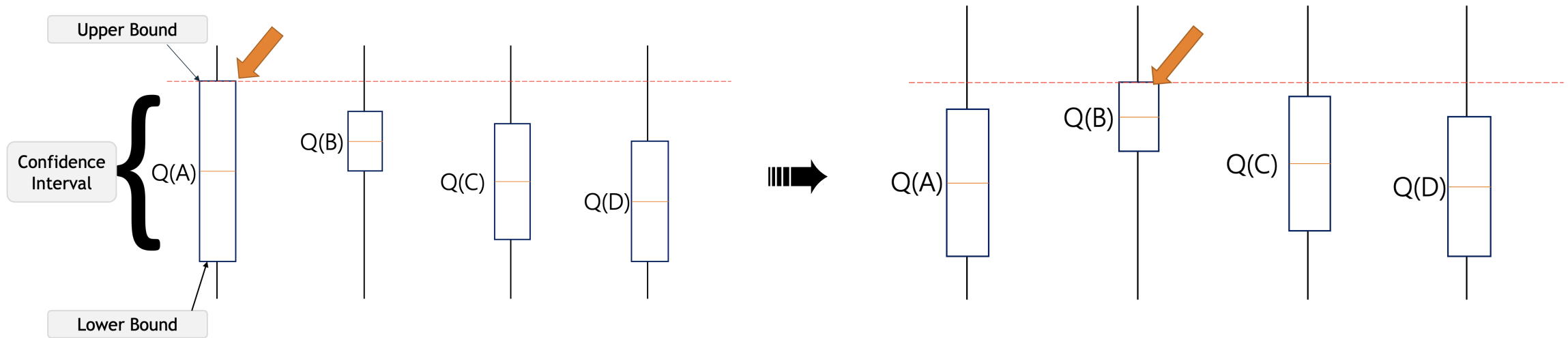


Image from <https://www.geeksforgeeks.org/upper-confidence-bound-algorithm-in-reinforcement-learning/>

## UCB exploration (2)

- 현재 주어진 상태  $s_t$ 에서  $N$ 개의 행동 후보  $\{a_{t,i}\}_{i=1}^N$  생성
- 각 행동마다  $N$ 개의 행동 가치 함수가 있어서 분포를 생각할 수 있음
- UCB가 가장 큰 행동을 선택

$$a_t = \max_a Q_{\text{mean}}(s_t, a) + \lambda Q_{\text{std}}(s_t, a)$$

- 구현에서는  $a \in \{a_{t,i}\}_{i=1}^N$



# SUNRISE Algorithm

$$\mathcal{L}_{\text{actor}}^{\text{SAC}}(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}} [\mathcal{L}_{\pi}(s_t, \phi)], \quad (3)$$

$$\mathcal{L}_{\pi}(s_t, \phi) = \mathbb{E}_{a_t \sim \pi_{\phi}} [\alpha \log \pi_{\phi}(a_t | s_t) - Q_{\theta}(s_t, a_t)]. \quad (4)$$

---

## Algorithm 1 SUNRISE: SAC version

---

```

1: for each iteration do
2:   for each timestep  $t$  do
3:     // UCB EXPLORATION
4:     Collect  $N$  action samples:  $\mathcal{A}_t = \{a_{t,i} \sim \pi_{\phi_i}(a|s_t) | i \in \{1, \dots, N\}\}$ 
5:     Choose the action that maximizes UCB:  $a_t = \arg \max_{a_{t,i} \in \mathcal{A}_t} Q_{\text{mean}}(s_t, a_{t,i}) + \lambda Q_{\text{std}}(s_t, a_{t,i})$ 
6:     Collect state  $s_{t+1}$  and reward  $r_t$  from the environment by taking action  $a_t$ 
7:     Sample bootstrap masks  $M_t = \{m_{t,i} \sim \text{Bernoulli}(\beta) \mid i \in \{1, \dots, N\}\}$ 
8:     Store transitions  $\tau_t = (s_t, a_t, s_{t+1}, r_t)$  and masks in replay buffer  $\mathcal{B} \leftarrow \mathcal{B} \cup \{(\tau_t, M_t)\}$ 
9:   end for
10:  // UPDATE AGENTS VIA BOOTSTRAP AND WEIGHTED BELLMAN BACKUP
11:  for each gradient step do
12:    Sample random minibatch  $\{(\tau_j, M_j)\}_{j=1}^B \sim \mathcal{B}$ 
13:    for each agent  $i$  do
14:      Update the Q-function by minimizing  $\frac{1}{B} \sum_{j=1}^B m_{j,i} \mathcal{L}_{WQ}(\tau_j, \theta_i)$  in (5)
15:      Update the policy by minimizing  $\frac{1}{B} \sum_{j=1}^B m_{j,i} \mathcal{L}_{\pi}(s_j, \phi_i)$  in (4)
16:    end for
17:  end for
18: end for

```

---

# OpenAI Gym Mujoco

➤ 로봇의 각 관절의 각도 및 각속도를 나타내는 벡터 입력

	Cheetah	Walker	Hopper	Ant	SlimHumanoid-ET
PETS	2288.4 $\pm$ 1019.0	282.5 $\pm$ 501.6	114.9 $\pm$ 621.0	1165.5 $\pm$ 226.9	2055.1 $\pm$ 771.5
POPLIN-A	1562.8 $\pm$ 1136.7	-105.0 $\pm$ 249.8	202.5 $\pm$ 962.5	1148.4 $\pm$ 438.3	-
POPLIN-P	4235.0 $\pm$ 1133.0	597.0 $\pm$ 478.8	2055.2 $\pm$ 613.8	2330.1 $\pm$ 320.9	-
METRPO	2283.7 $\pm$ 900.4	-1609.3 $\pm$ 657.5	1272.5 $\pm$ 500.9	282.2 $\pm$ 18.0	76.1 $\pm$ 8.8
TD3	3015.7 $\pm$ 969.8	-516.4 $\pm$ 812.2	1816.6 $\pm$ 994.8	870.1 $\pm$ 283.8	1070.0 $\pm$ 168.3
SAC	4474.4 $\pm$ 700.9	299.5 $\pm$ 921.9	1781.3 $\pm$ 737.2	979.5 $\pm$ 253.2	1371.8 $\pm$ 473.4
SUNRISE	4501.8 $\pm$ 443.8	1236.5 $\pm$ 1123.9	2643.2 $\pm$ 472.3	1502.4 $\pm$ 483.5	1926.6 $\pm$ 375.0

Table 1. Performance on OpenAI Gym at 200K timesteps. The results show the mean and standard deviation averaged over ten runs. For baseline methods, we report the best number in prior works (Wang & Ba, 2020; Wang et al., 2019).

# DeepMind Control Suite

## ➤ 이미지 입력

500K step	PlaNet	Dreamer	SLAC	CURL	DrQ	RAD	SUNRISE
Finger-spin	561 $\pm$ 284	796 $\pm$ 183	673 $\pm$ 92	926 $\pm$ 45	938 $\pm$ 103	975 $\pm$ 16	983 $\pm$ 1
Cartpole-swing	475 $\pm$ 71	762 $\pm$ 27	-	845 $\pm$ 45	868 $\pm$ 10	873 $\pm$ 3	876 $\pm$ 4
Reacher-easy	210 $\pm$ 44	793 $\pm$ 164	-	929 $\pm$ 44	942 $\pm$ 71	916 $\pm$ 49	982 $\pm$ 3
Cheetah-run	305 $\pm$ 131	570 $\pm$ 253	640 $\pm$ 19	518 $\pm$ 28	660 $\pm$ 96	624 $\pm$ 10	678 $\pm$ 46
Walker-walk	351 $\pm$ 58	897 $\pm$ 49	842 $\pm$ 51	902 $\pm$ 43	921 $\pm$ 45	938 $\pm$ 9	953 $\pm$ 13
Cup-catch	460 $\pm$ 380	879 $\pm$ 87	852 $\pm$ 71	959 $\pm$ 27	963 $\pm$ 9	966 $\pm$ 9	969 $\pm$ 5
100K step							
Finger-spin	136 $\pm$ 216	341 $\pm$ 70	693 $\pm$ 141	767 $\pm$ 56	901 $\pm$ 104	811 $\pm$ 146	905 $\pm$ 57
Cartpole-swing	297 $\pm$ 39	326 $\pm$ 27	-	582 $\pm$ 146	759 $\pm$ 92	373 $\pm$ 90	591 $\pm$ 55
Reacher-easy	20 $\pm$ 50	314 $\pm$ 155	-	538 $\pm$ 233	601 $\pm$ 213	567 $\pm$ 54	722 $\pm$ 50
Cheetah-run	138 $\pm$ 88	235 $\pm$ 137	319 $\pm$ 56	299 $\pm$ 48	344 $\pm$ 67	381 $\pm$ 79	413 $\pm$ 35
Walker-walk	224 $\pm$ 48	277 $\pm$ 12	361 $\pm$ 73	403 $\pm$ 24	612 $\pm$ 164	641 $\pm$ 89	667 $\pm$ 147
Cup-catch	0 $\pm$ 0	246 $\pm$ 174	512 $\pm$ 110	769 $\pm$ 43	913 $\pm$ 53	666 $\pm$ 181	633 $\pm$ 241

Table 2. Performance on DeepMind Control Suite at 100K and 500K environment steps. The results show the mean and standard deviation averaged five runs. For baseline methods, we report the best numbers reported in prior works (Kostrikov et al., 2021).

# Weighted Bellman backups 효과 검증 (1)

➤ Mujoco 환경의 보상에 노이즈가 추가하여 에이전트 학습. 평가는 원래 보상으로.

■  $r'(s, a) = r(s, a) + \mathcal{N}(0,1)$

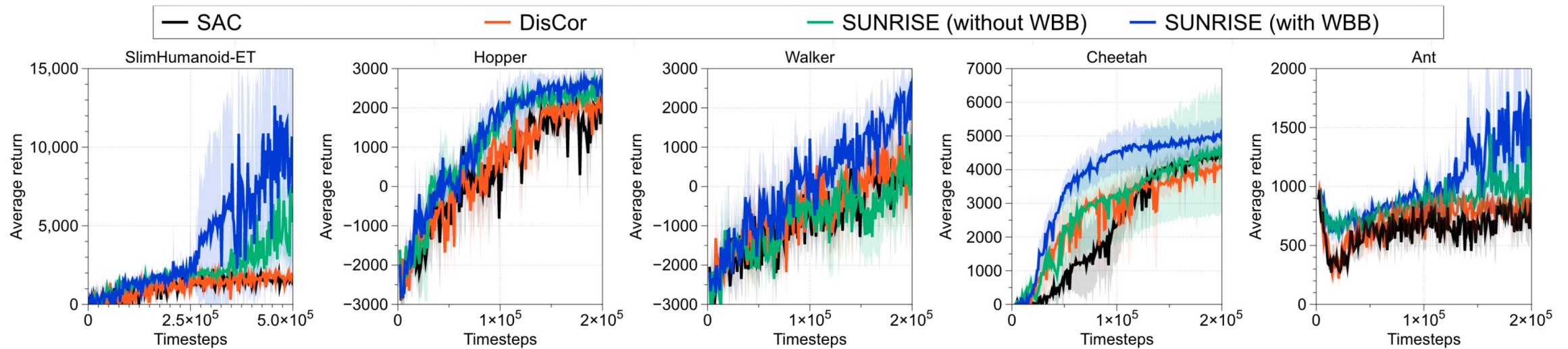
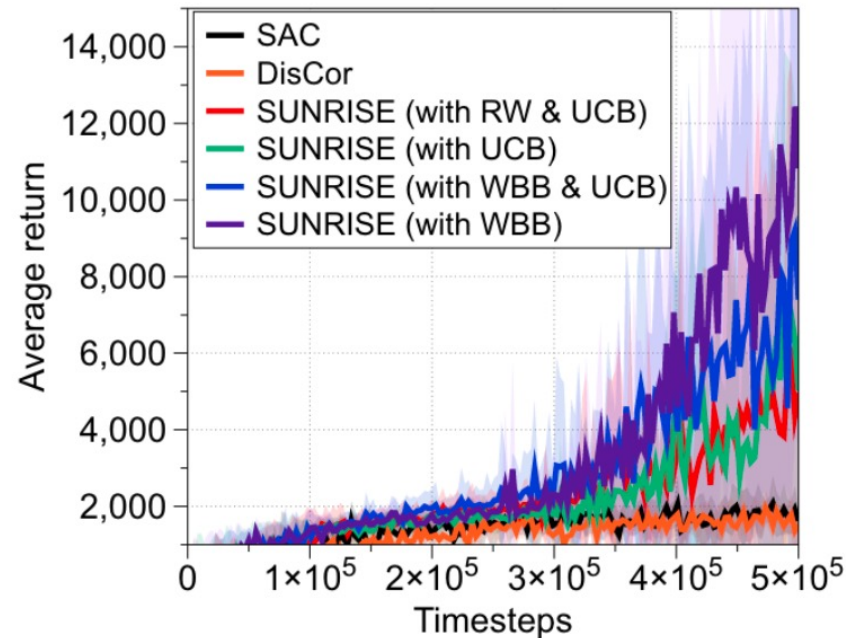


Figure 2. Learning curves on OpenAI Gym with noisy rewards. To verify the effects of the weighted Bellman backups (WBB), we consider SUNRISE with WBB and without WBB. The solid line and shaded regions represent the mean and standard deviation, respectively, across four runs.

Quoted from Kimin Lee et al., SUNRISE: A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning. (2021)

## Weighted Bellman backups 효과 검증 (2)

- 제일 복잡한 Humanoid 환경에 더 큰 노이즈를 추가한 보상으로 실험
  - $r'(s, a) = r(s, a) + \mathcal{N}(0, 5)$

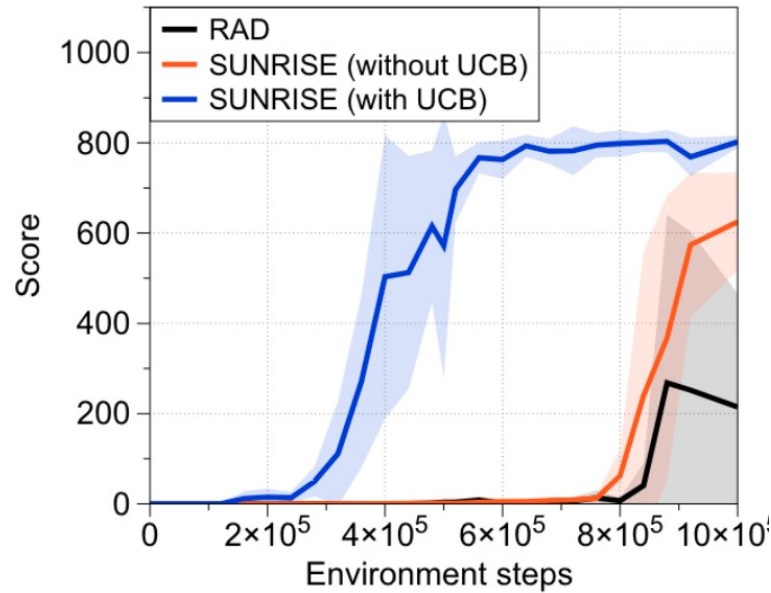


Mujoco의 Humanoid 환경  
RW: Random weight  
WBB: Weighted Bellman backups

(a) Large noise

# UCB exploration의 효과

- CartPole with sparse reward에 실험
- Single agent가 더 많이 학습되면 SUNRISE를 뛰어 넘을 수 있는지 확인

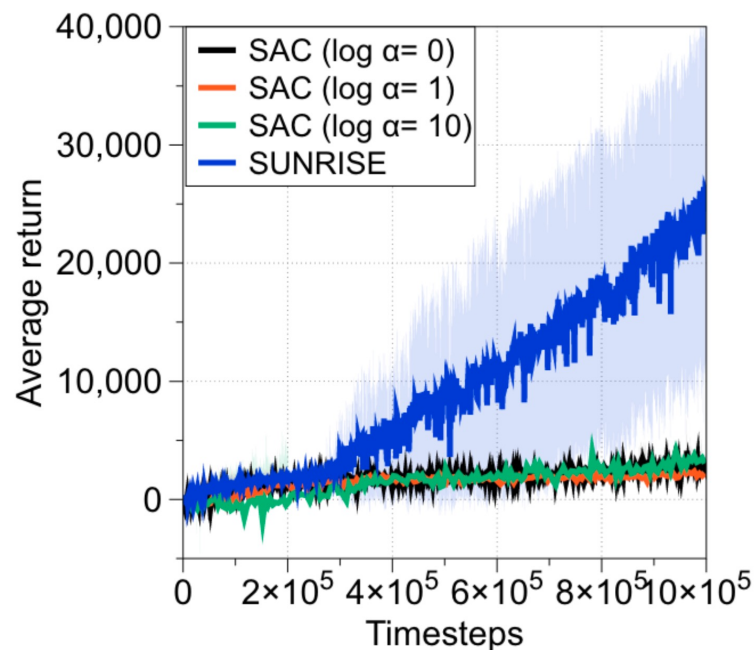


DMC의 CartPole with sparse reward

(b) Sparse reward

## 더 많이 학습한 single agent와 비교

- 네트워크가  $N$ 배 더 많은 만큼 더 많은 파라미터 업데이트를 한다.
- Single agent가 더 많이 학습되면 SUNRISE를 뛰어 넘을 수 있는지 확인



Mujoco의 Humanoid 환경

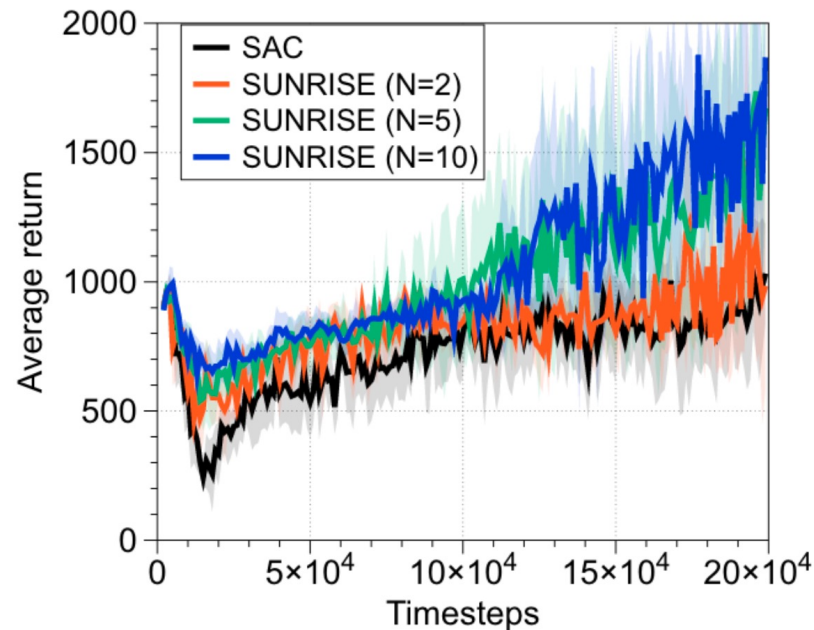
(c) Gradient update

Quoted from Kimin Lee et al., SUNRISE: A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning. (2021)



## 네트워크 개수의 영향

- 네트워크 개수가 증가할 수록 성능이 증가하지만,
- $N = 5$ 에서 saturation 된다.



Mujoco의 Ant 환경

(d) Ensemble size

Quoted from Kimin Lee et al., SUNRISE: A Simple Unified Framework for Ensemble Learning in Deep Reinforcement Learning. (2021)



# Computation overhead

- 각 네트워크 forward pass가 병렬 처리할 수 있기 때문에 충분히 efficient 하다고 주장
- But, 시간에 따른 실험 결과는 없음