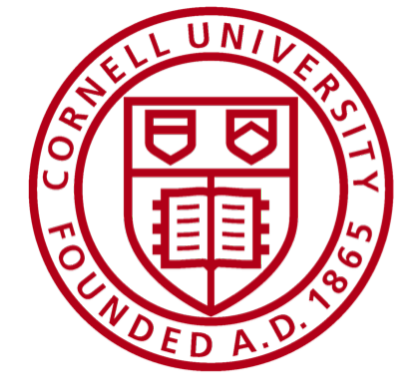


Improving Stability in Deep Reinforcement Learning with Weight Averaging

Evgenii Nikishin¹, Pavel Izmailov², Ben Athiwaratkun², Dmitrii Podoprikin^{1,3}
Timur Garipov⁴, Pavel Shvechikov¹, Dmitry Vetrov^{1,3}, Andrew Gordon Wilson²
¹National Research University Higher School of Economics, ²Cornell University
³Samsung-HSE Laboratory, ⁴Samsung AI Center in Moscow



Outline

- Deep reinforcement learning (RL) methods are notoriously unstable during training.
- Stochastic weight averaging (SWA) is a technique based on averaging the weights collected during training with an SGD-like method.
- We propose to apply SWA, in order to reduce the effect of noise on training.
- We show that SWA stabilizes the model solutions and improves the average rewards.

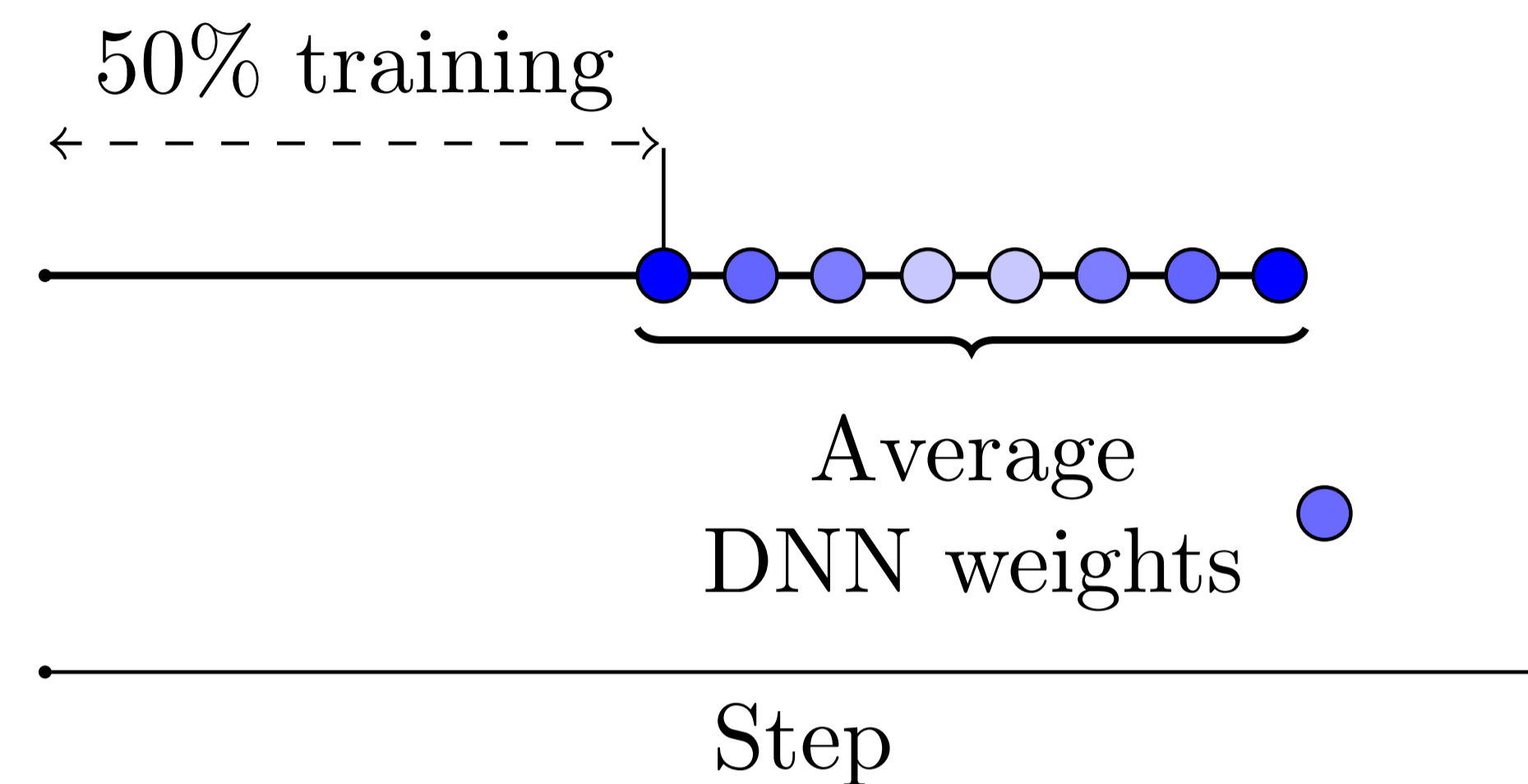
Methods

Reinforcement learning

- Advantage Actor-Critic (A2C) is a standard RL algorithm, often applied to problems with discrete action spaces.
- Deep Deterministic Policy Gradient (DDPG) is another standard RL algorithm, but suitable for continuous action spaces.

Stochastic weight averaging

SWA was shown to find solutions with better generalization in both supervised and semi-supervised learning.



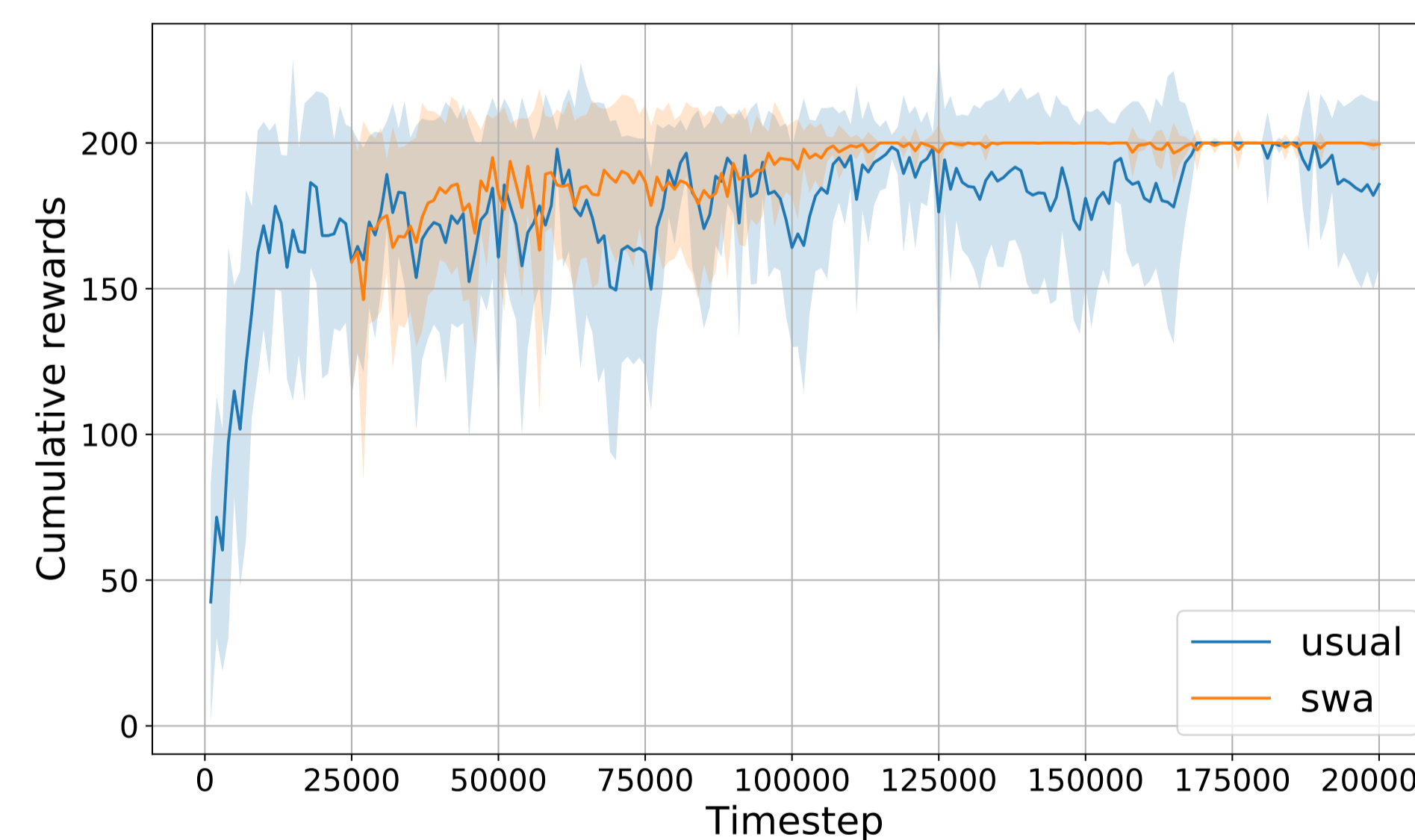
After initial pretraining phase (for example, 50% of computational budget), SWA collects weights every c timesteps.

Note that there is no need to store all the weights in memory:

$$w_{SWA} \leftarrow \frac{n_{SWA} \cdot w_{SWA} + w}{n_{SWA} + 1}, \quad n_{SWA} \leftarrow n_{SWA} + 1$$

Results

Average cumulative rewards of A2C for CartPole environment with and without SWA



A2C on Atari environments

ENV NAME	A2C	A2C + SWA
Breakout	522 ± 34	703 ± 60
Qbert	18777 ± 778	21272 ± 655
SpaceInvaders	7727 ± 1121	21676 ± 8897
Seaquest	1779 ± 4	1795 ± 4
CrazyClimber	147030 ± 10239	139752 ± 11618
BeamRider	9999 ± 402	11321 ± 1065

DDPG on MuJoCo environments

ENV NAME	DDPG	DDPG + SWA
Hopper	613 ± 683	1615 ± 1143
Walker2d	1803 ± 96	2457 ± 241
Half-Cheetah	3825 ± 1187	4228 ± 1117
Ant	865 ± 899	1051 ± 696

Discussion

- Currently, the averaging does not affect the training procedure. Modification of the training procedure based on weight averaging can potentially help in stabilization and training acceleration.
- Theoretical justification of weight averaging in RL context.
- Analysis of the RL loss surface can reveal new approaches to weight averaging for further stabilization.