# Soft Actor-Critic

## Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor
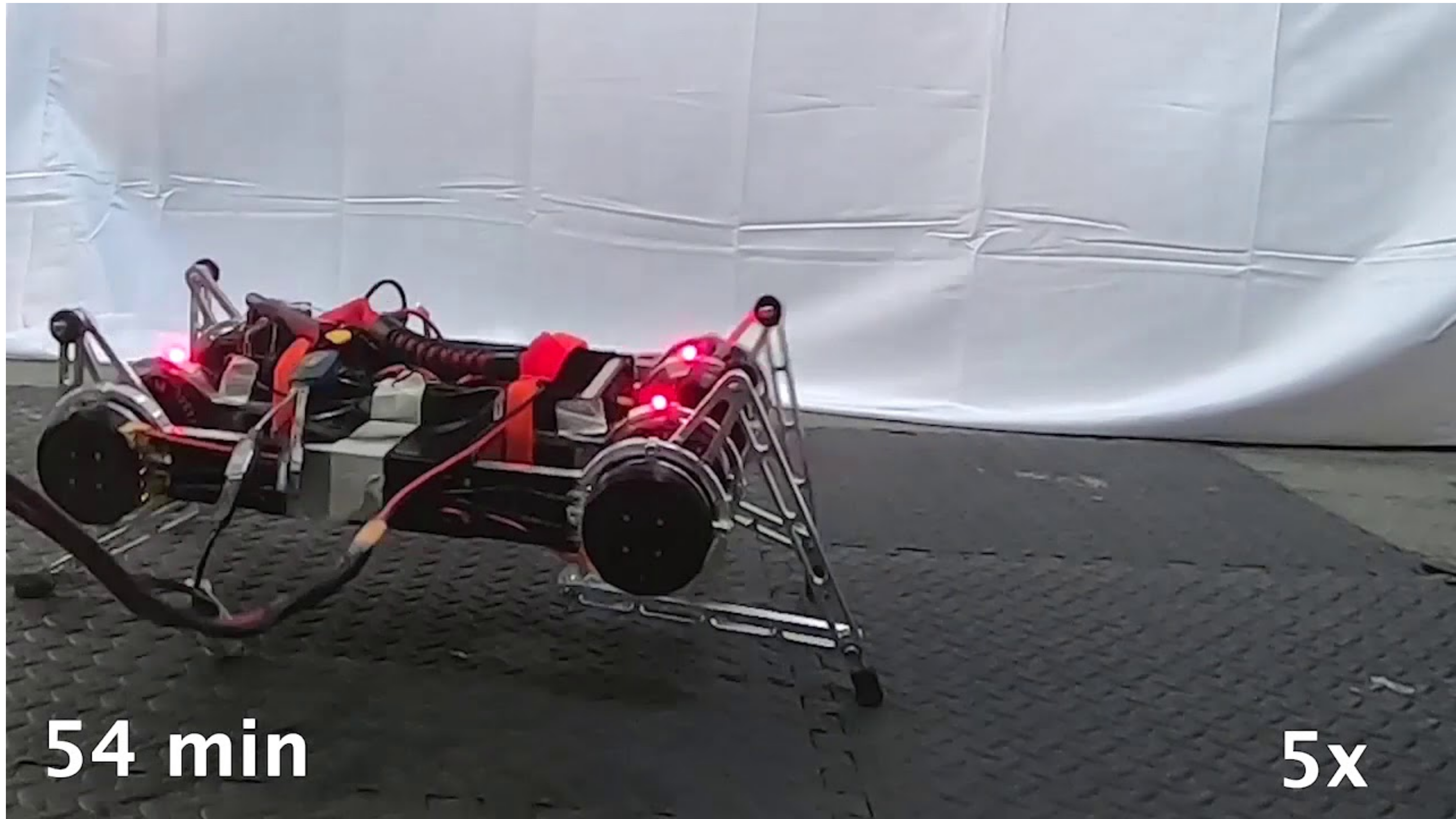
Tuomas Haarnoja      Aurick Zhou      Pieter Abbeel      Sergey Levine

- Why Soft Actor-Critic ?

—> sample efficient

—> very stable

—> exploration more efficient

54 min      5x

- Some general points:

- actor-critic

- off-policy algorithm

- continuous state and action spaces

# Entropy-regularized RL setting

$$H(X) = -\sum_{i=1}^{n} P(x_i) \log P(x_i) = -\mathbb{E}\left[\log P(X)\right]$$

$$\pi_{\text{old}}^* = \arg\max_{\pi} \mathop{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \left( R\left(s_t, a_t, s_{t+1}\right) \right) \right]$$

$$\pi^* = \arg\max_{\pi} \mathop{E}_{\tau \sim \pi} \left[ \underbrace{\sum_{t=0}^{\infty} \gamma^t \left( R\left(s_t, a_t, s_{t+1}\right) + \alpha H\left(\pi\left(\,\cdot\mid s_t\right)\right) \right)}_{V^\pi(s_0)} \right]$$

# Entropy-regularized RL setting

$$V^\pi(s) = \mathop{\mathrm{E}}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \left( R\left(s_t, a_t, s_{t+1}\right) + \alpha H\left(\pi\left(\,\cdot\mid s_t\right)\right)\right) \mid s_0 = s \right]$$

$$Q^\pi(s, a) = \mathop{\mathrm{E}}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R\left(s_t, a_t, s_{t+1}\right) + \alpha \sum_{t=1}^{\infty} \gamma^t H\left(\pi\left(\,\cdot\mid s_t\right)\right) \mid s_0 = s, a_0 = a \right]$$

Equation de Bellman pour $Q^\pi(s, a)$ :

$$Q^\pi(s, a) = \mathop{\mathrm{E}}_{\substack{s' \sim P \\ a' \sim \pi}} \left[ R\left(s, a, s'\right) + \gamma V^\pi(s') \right]$$

$$Q^\pi(s, a) = \mathop{\mathrm{E}}_{\substack{s' \sim P \\ a' \sim \pi}} \left[ R\left(s, a, s'\right) + \gamma \left( Q^\pi\left(s', a'\right) + \alpha H\left(\pi\left(\,\cdot\mid s'\right)\right)\right)\right]$$

# Critic

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)
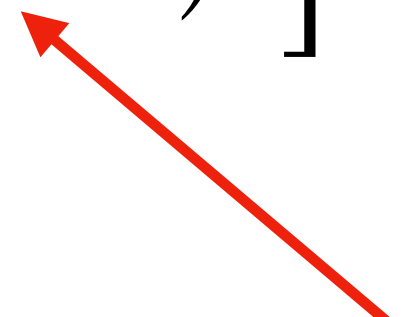
$$L(\phi, \mathscr{D}) = \mathop{\mathrm{E}}_{(s,a,r,s',d)\sim\mathscr{D}} \left[ \left( Q_\phi(s,a) - Q^\pi(s,a) \right)^2 \right]$$

approximator

target

# Critic

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

$$L(\phi, \mathscr{D}) = \underset{(s,a,r,s',d) \sim \mathscr{D}}{\mathrm{E}} \left[ \left( Q_\phi(s, a) - Q^\pi(s, a) \right)^2 \right]$$

approximator

target

$$Q^\pi(s, a) \approx R\left(s, a, s'\right) + \gamma \left( Q^\pi\left(s', \tilde{a}'\right) - \alpha \log \pi \left( \tilde{a}' \mid s' \right) \right) \quad \text{avec} \quad \tilde{a}' \sim \pi_\theta \left( \cdot \mid s' \right)$$

# Critic

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

$$L(\phi, \mathscr{D}) = \mathop{\mathrm{E}}_{(s,a,r,s',d)\sim\mathscr{D}} \left[ \left( Q_\phi(s,a) - y(r,s',d) \right)^2 \right]$$
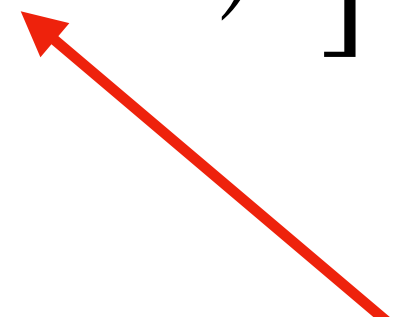
approximator

target

$$y(r,s',d) = R(s,a,s') + \gamma(1-d)\left( Q_\phi(s',\tilde{a}') - \alpha \log \pi\left(\tilde{a}' \mid s'\right) \right) \quad \text{avec} \quad \tilde{a}' \sim \pi_\theta\left( \cdot \mid s' \right)$$

# Critic

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

$$L(\phi, \mathscr{D}) = \underset{(s,a,r,s',d) \sim \mathscr{D}}{\mathrm{E}} \left[ \left( Q_\phi(s, a) - y(r, s', d) \right)^2 \right]$$

approximator

target

$$y(r, s', d) = R(s, a, s') + \gamma(1 - d)\left( Q_{\phi_{targ}}(s', \tilde{a}') - \alpha \log \pi\left( \tilde{a}' \mid s' \right) \right) \quad \text{avec} \quad \tilde{a}' \sim \pi_\theta\left( \cdot \mid s' \right)$$

$$\phi_{\text{targ}} \leftarrow \rho \phi_{\text{targ}} + (1 - \rho)\phi$$

Silver, David, et al. "Deterministic policy gradient algorithms." *International conference on machine learning.* PMLR, 2014.

# Critic

1<sup>st</sup> trick: Target Networks
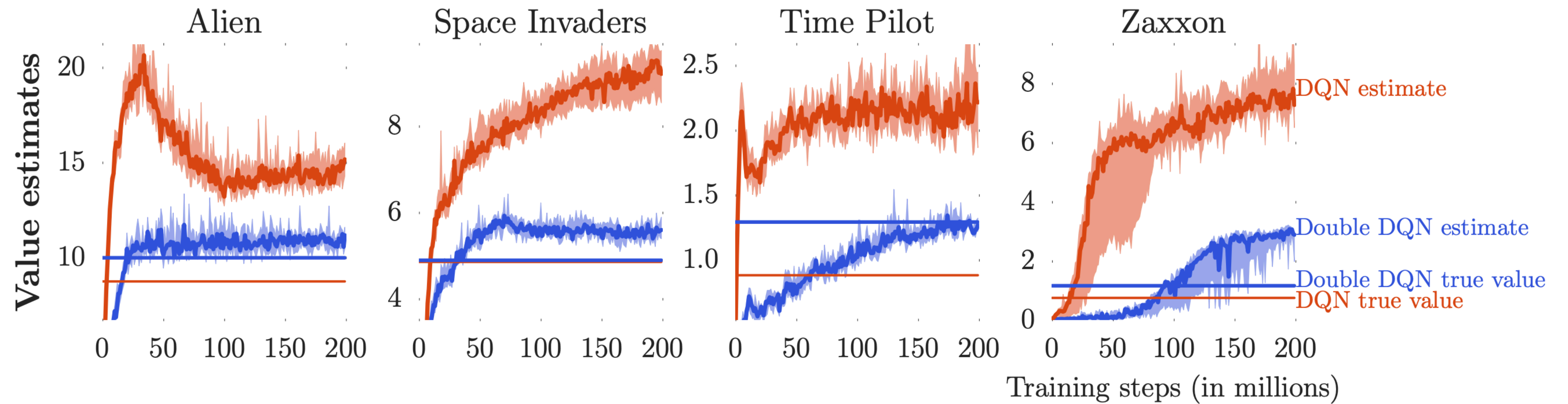
Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

$$L(\phi, \mathscr{D}) = \underset{(s,a,r,s',d)\sim\mathscr{D}}{\mathrm{E}} \left[ \left( Q_\phi(s,a) - y(r,s',d) \right)^2 \right]$$

approximator                                    target

$$y(r,s',d) = R(s,a,s') + \gamma(1-d)\left( Q_{\phi_{targ}}(s',\tilde{a}') - \alpha \log \pi\left(\tilde{a}' \mid s'\right) \right) \quad \text{avec} \quad \tilde{a}' \sim \pi_\theta\left( \cdot \mid s' \right)$$
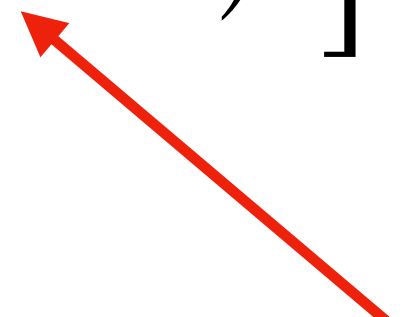
# Critic



Source: Deep Reinforcement Learning with Double Q-learning (Hasselt et al., 2015)

# Critic

1<sup>st</sup> trick: Target Networks

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

$$L(\phi, \mathscr{D}) = \underset{(s,a,r,s',d)\sim\mathscr{D}}{\mathrm{E}} \left[\left(Q_\phi(s,a) - y(r,s',d)\right)^2\right]$$

approximator

target

$$y(r,s',d) = R(s,a,s') + \gamma(1-d)\left(Q_{\phi_{targ}}(s',\tilde{a}') - \alpha \log \pi\left(\tilde{a}' \mid s'\right)\right) \quad \text{avec} \quad \tilde{a}' \sim \pi_\theta\left(\cdot \mid s'\right)$$

# Critic

2<sup>nd</sup> trick: Clipped Double Q-network $\longrightarrow$ $\phi_1$ and $\phi_2$

Loss function of the Critic: MSBE (*Mean Squared Bellman Error*)

For $i = 1,2$ :

$$L(\phi_i, \mathscr{D}) = \underset{(s,a,r,s',d)\sim\mathscr{D}}{\mathrm{E}} \left[ \left( Q_{\phi_i}(s,a) - y(r,s',d) \right)^2 \right]$$

approximator

target

$$y(r,s',d) = R(s,a,s') + \gamma(1-d)\left( \min_{j=1,2} Q_{\phi_{\mathrm{targ},j}}(s',\tilde{a}') - \alpha \log \pi\left(\tilde{a}' \mid s'\right) \right) \quad \text{avec} \ \ \tilde{a}' \sim \pi_\theta\left( \cdot \mid s' \right)$$

# Actor

$$\pi^* = \arg\max_{\pi} V^\pi(s)$$

Objective function of the Actor for each state s:

$$V^\pi(s) = \mathop{\mathrm{E}}_{a \sim \pi} \left[ Q^\pi(s, a) - \alpha \log \pi(a \mid s) \right]$$

# Actor

$$\pi^* = \arg\max_{\pi} V^{\pi}(s)$$

Objective function of the Actor for each state s:

$$V^{\pi}(s) = \mathop{\mathbb{E}}_{\tilde{a}\sim\pi_\theta(\cdot\,|\,s)} \left[ \min_{j=1,2} Q_{\phi_{\text{targ},j}}(s,\tilde{a}) - \alpha \log \pi_\theta(\tilde{a}\,|\,s) \right]$$

# Actor

$$\pi^* = \arg\max_{\pi} V^{\pi}(s)$$

Objective function of the Actor for each state s:

$$V^{\pi}(s) = \underset{\tilde{a}\sim\pi_{\theta}(\,\cdot\,|\,s)}{\mathrm{E}} \left[ \min_{j=1,2} Q_{\phi_{\mathrm{targ},j}}(s,\tilde{a}) - \alpha \log \pi_{\theta}(\tilde{a}\mid s) \right]$$

Pain point: the distribution depends on policy params

# Actor

Objective function of the Actor for each state s:

$$V^{\pi}(s) = \mathop{\mathrm{E}}_{\xi \sim \mathcal{N}} \left[ Q^{\pi_\theta} \left( s, \tilde{a}_\theta(s, \xi) \right) - \alpha \log \pi_\theta \left( \tilde{a}_\theta(s, \xi) \mid s \right) \right]$$

$$\tilde{a}_\theta(s, \xi) = \tanh \left( \mu_\theta(s) + \sigma_\theta(s) \odot \xi \right), \quad \xi \sim \mathcal{N}(0, I)$$
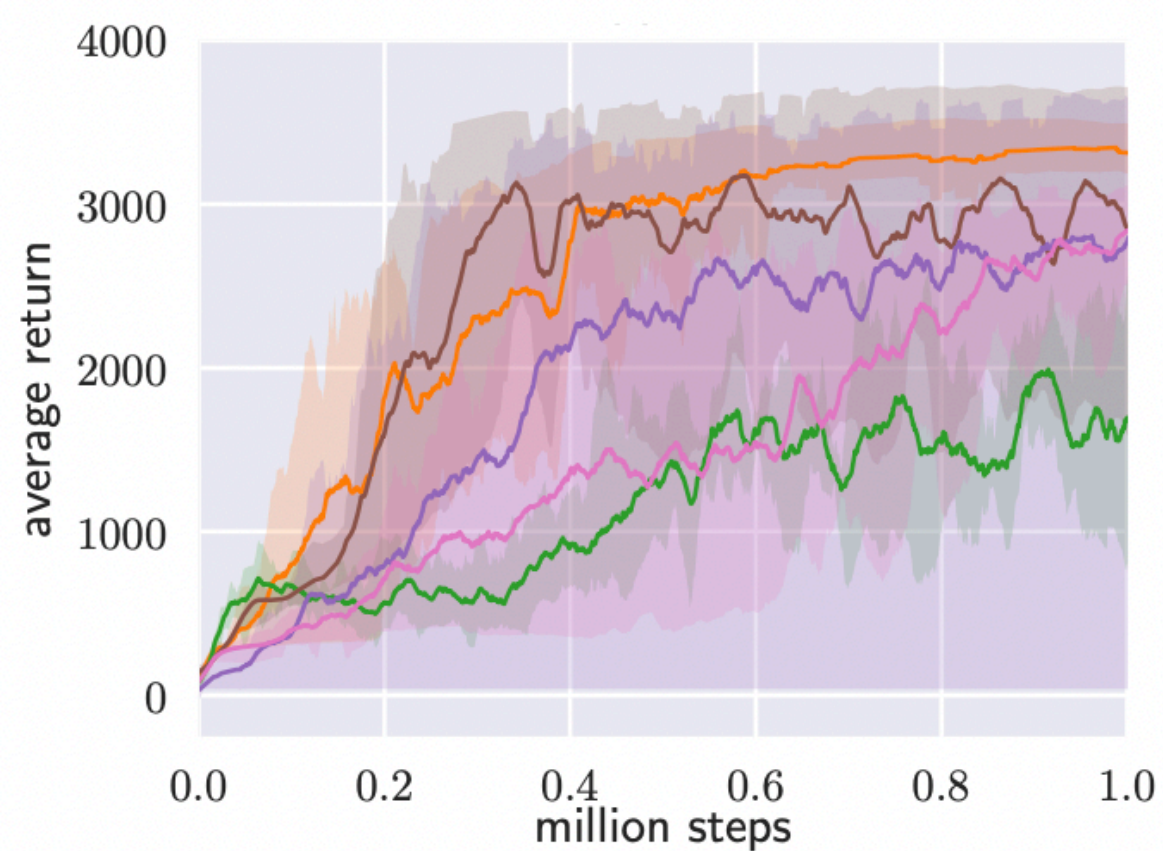
# Actor

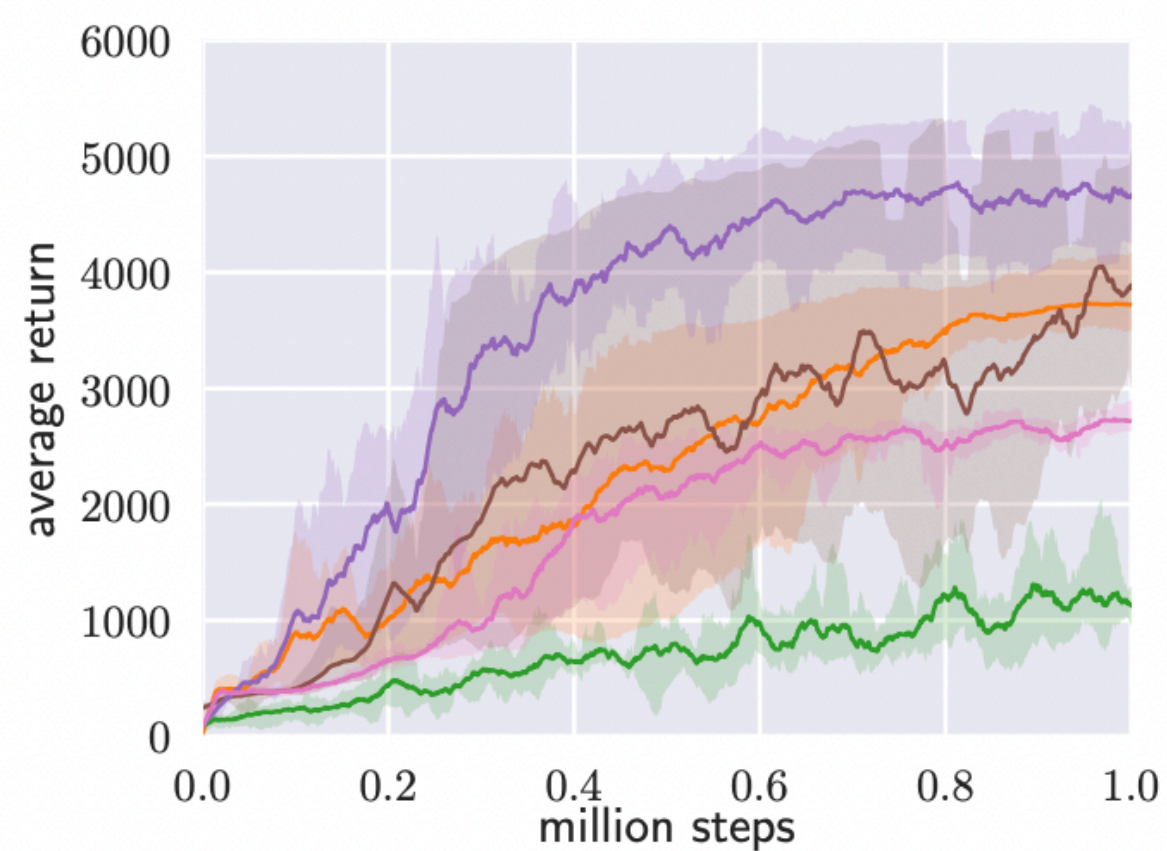$3^{rd}$ Trick : reparameterization

Objective function of the Actor:

$$\max_{\theta} \operatorname*{E}_{\substack{s \sim \mathcal{D} \\ \xi \sim \mathcal{N}}} \left[ \min_{j=1,2} Q_{\phi_j}\left(s, \tilde{a}_\theta(s, \xi)\right) - \alpha \log \pi_\theta\left(\tilde{a}_\theta(s, \xi) \mid s\right) \right]$$

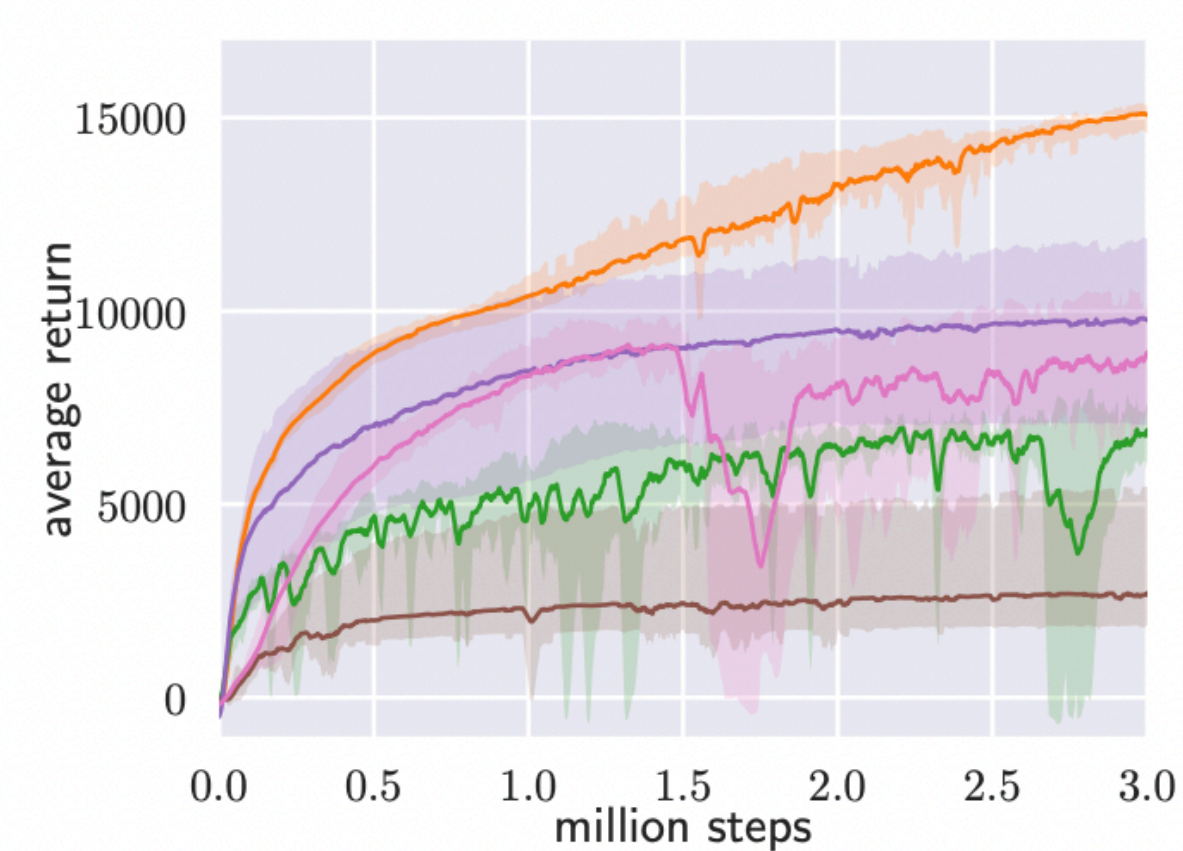$$\tilde{a}_\theta(s, \xi) = \tanh\left(\mu_\theta(s) + \sigma_\theta(s) \odot \xi\right), \quad \xi \sim \mathcal{N}(0, I)$$
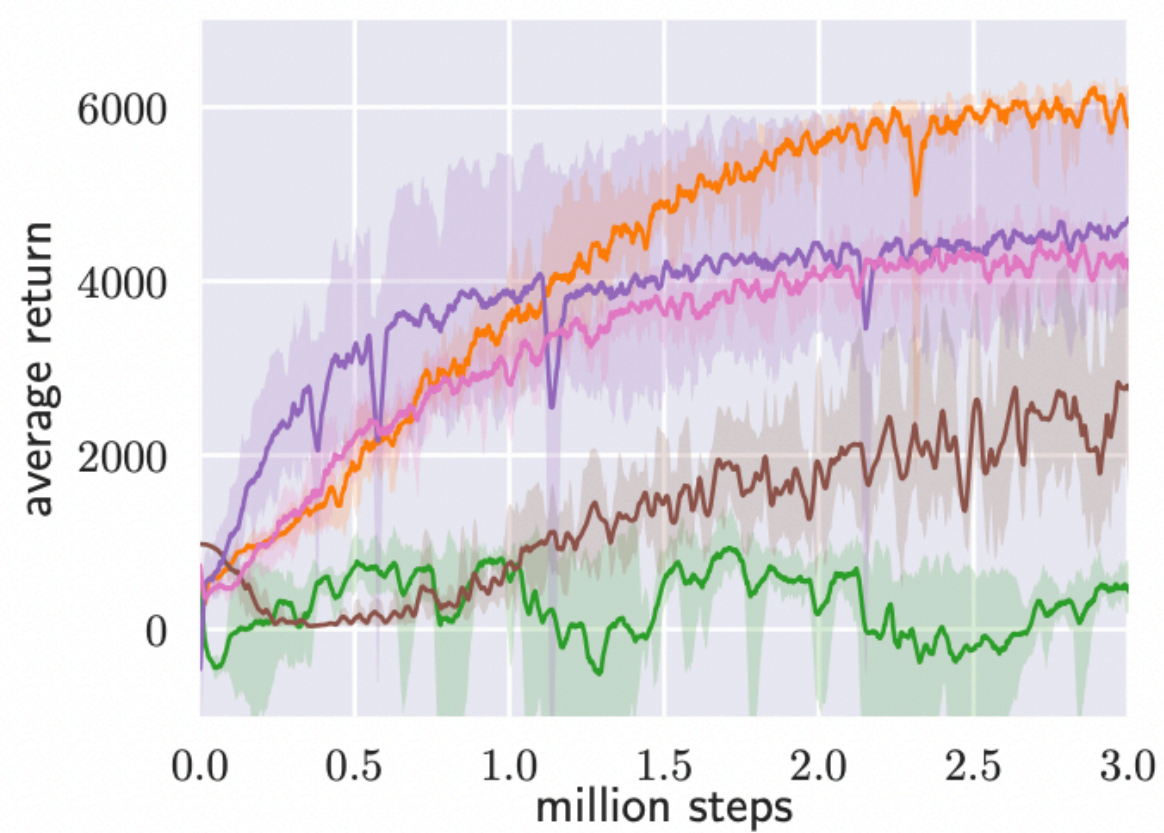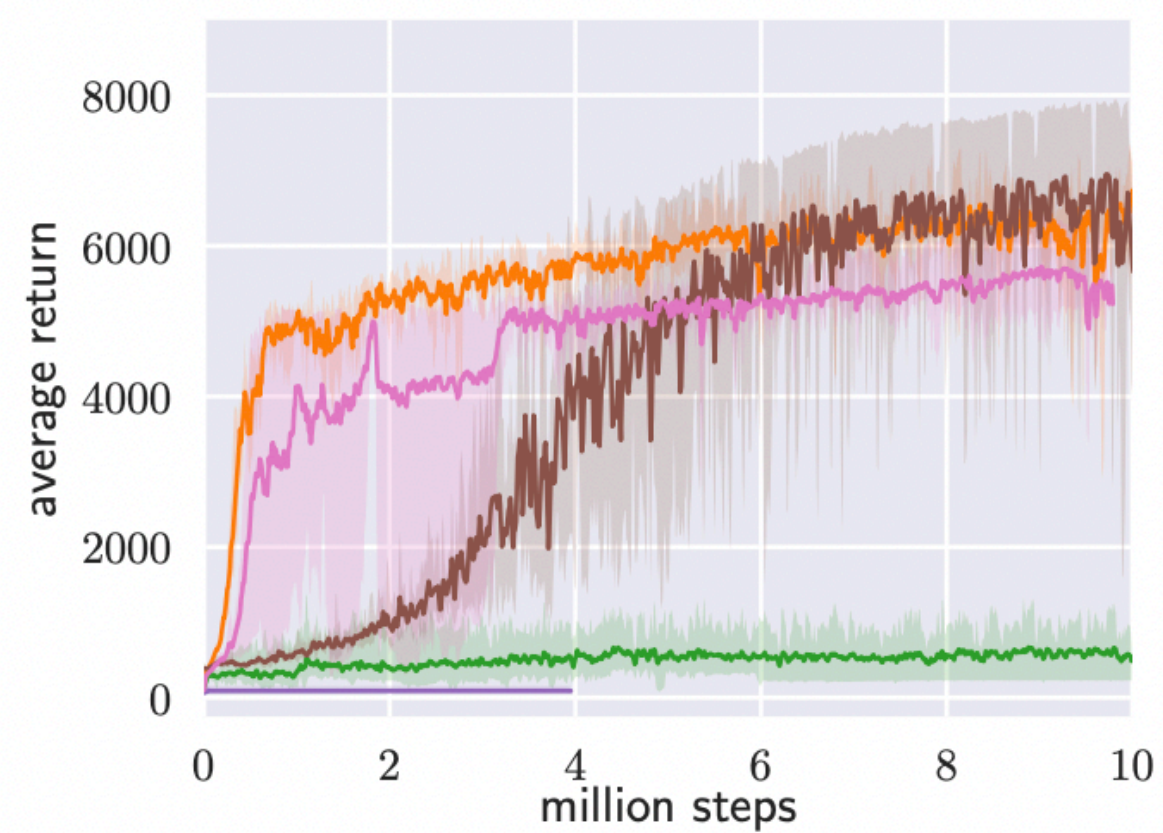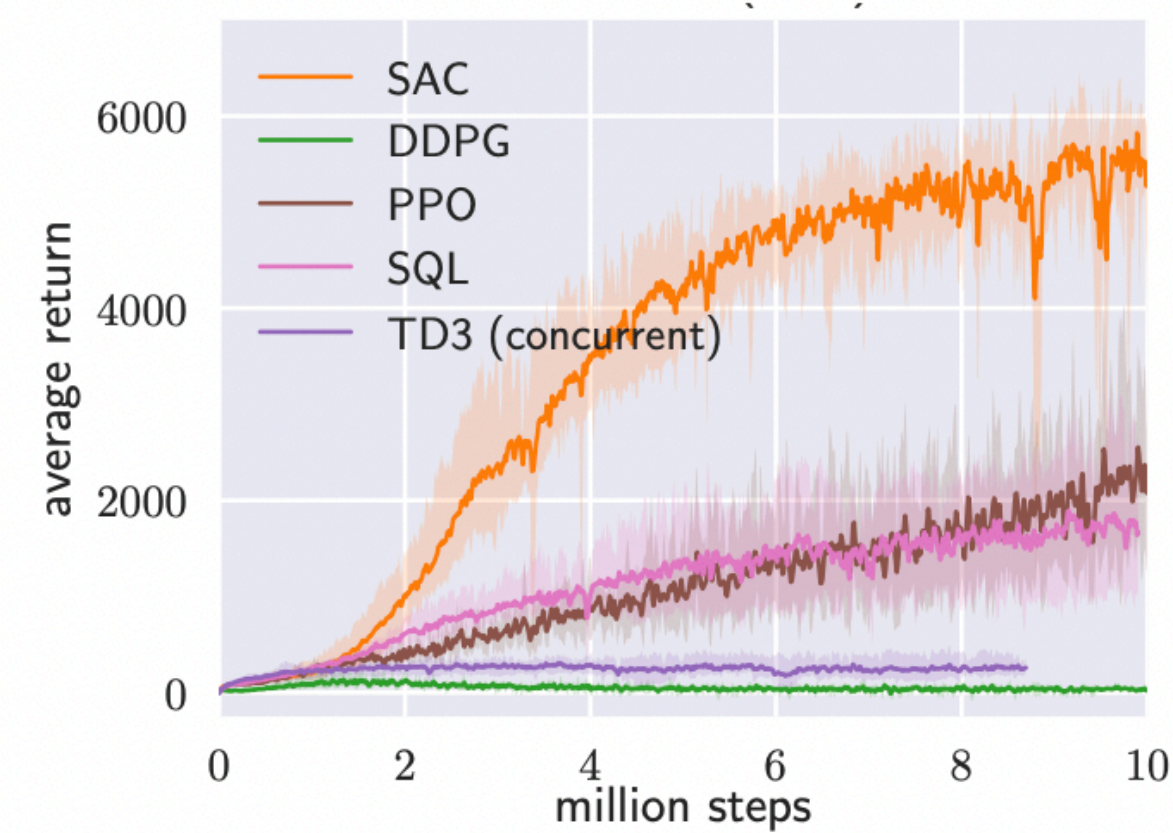
# Results



(a) Hopper-v1

(b) Walker2d-v1

(c) HalfCheetah-v1

(d) Ant-v1

(e) Humanoid-v1

(f) Humanoid (rllab)

# References

- Haarnoja, Tuomas, et al. "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor." *International conference on machine learning.* PMLR, 2018.

- Silver, David, et al. "Deterministic policy gradient algorithms." *International conference on machine learning.* PMLR, 2014.

- Fujimoto, Scott, Herke Hoof, and David Meger. "Addressing function approximation error in actor-critic methods." *International conference on machine learning.* PMLR, 2018.

- https://spinningup.openai.com/en/latest/algorithms/sac.html