

# Continuous Control (Part 1)

**Continuous control with deep reinforcement learning**

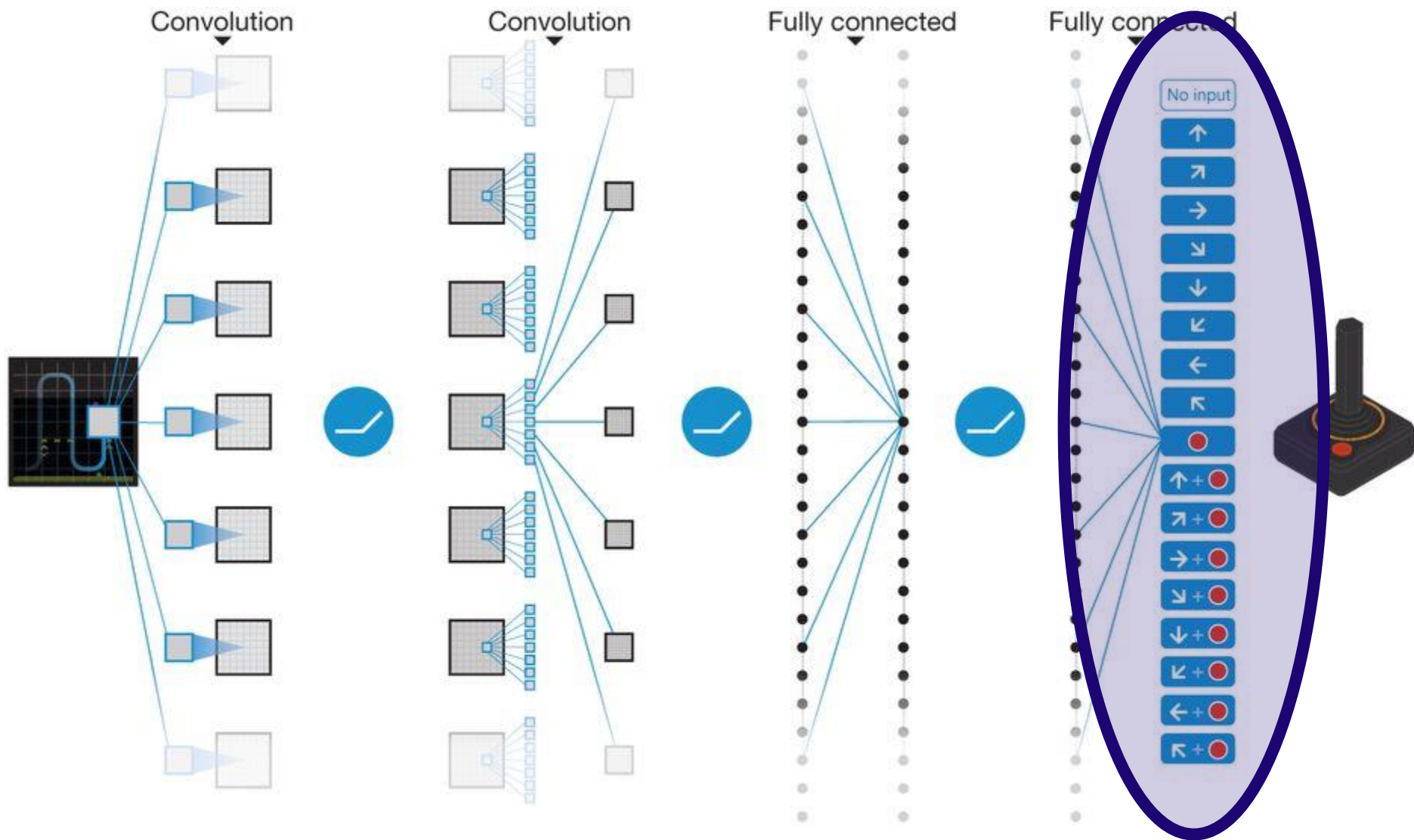
Timothy P. Lillicrap et. al. ICLR 2016.

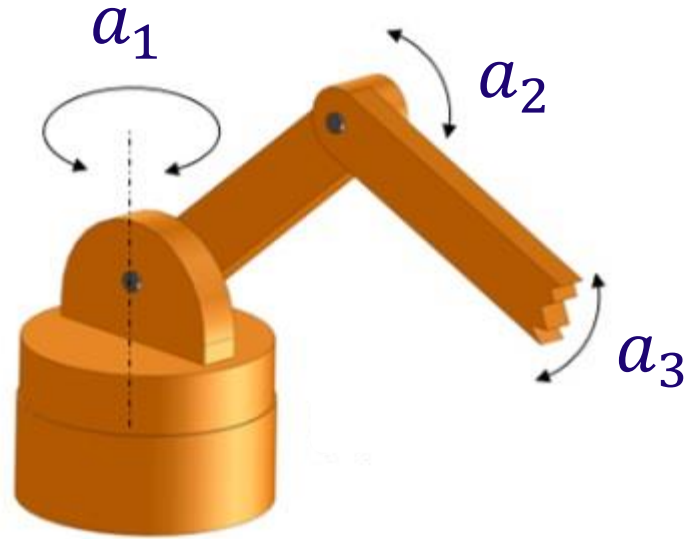
**Trust Region Policy Optimization**

John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, Pieter Abbeel. ICML 2015.

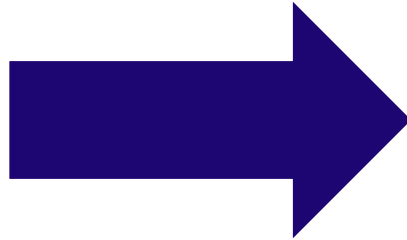
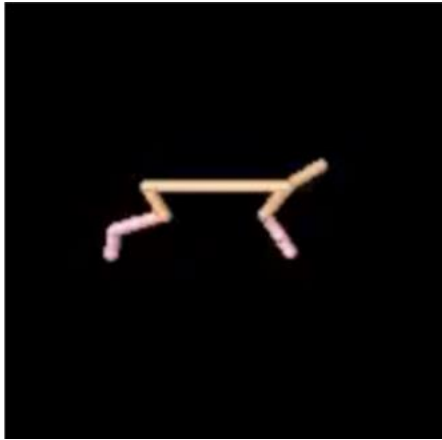
Presented by

Mark Arnold

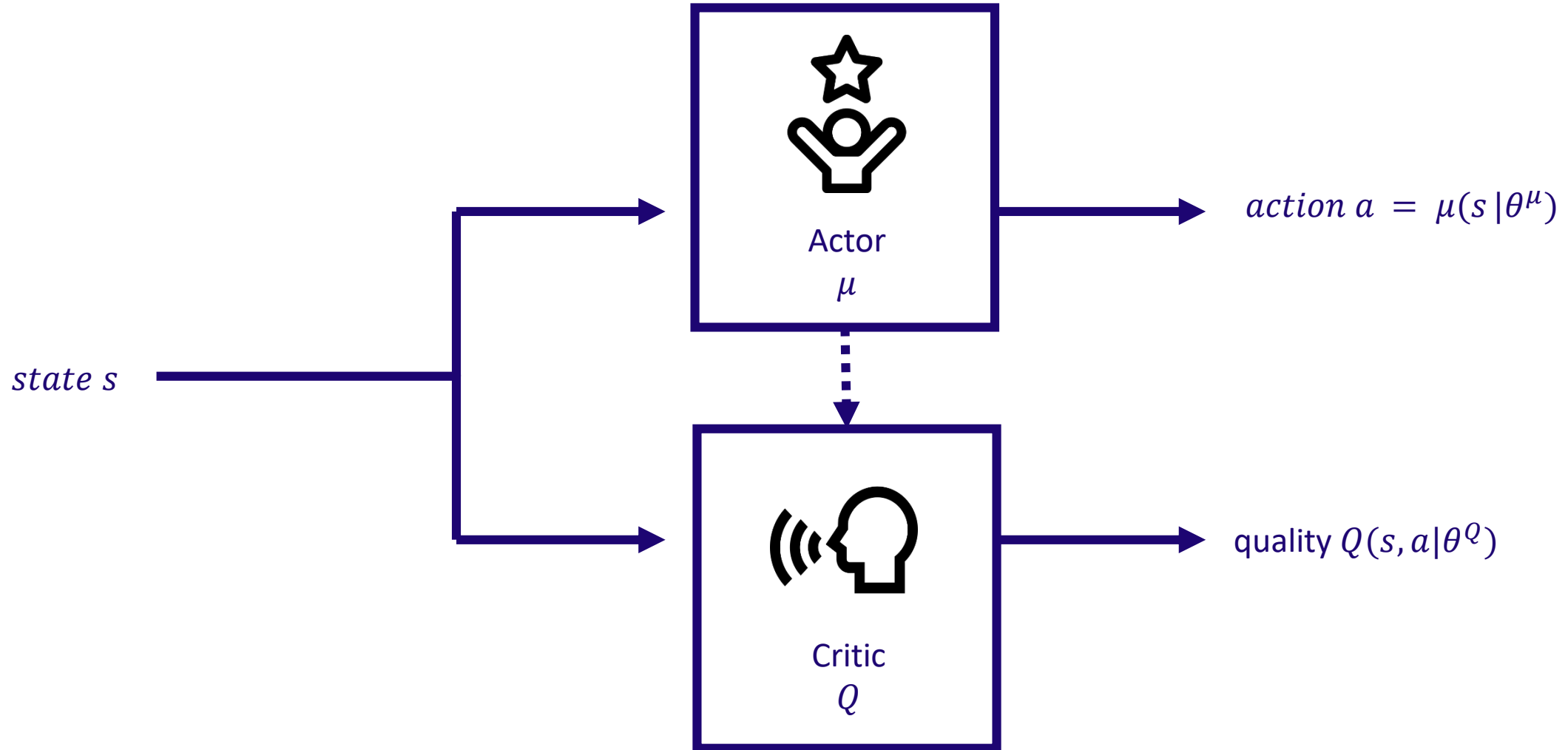


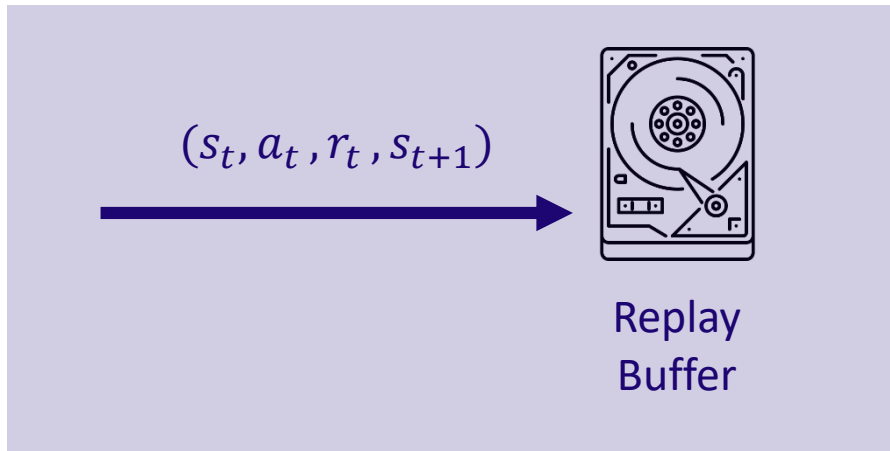


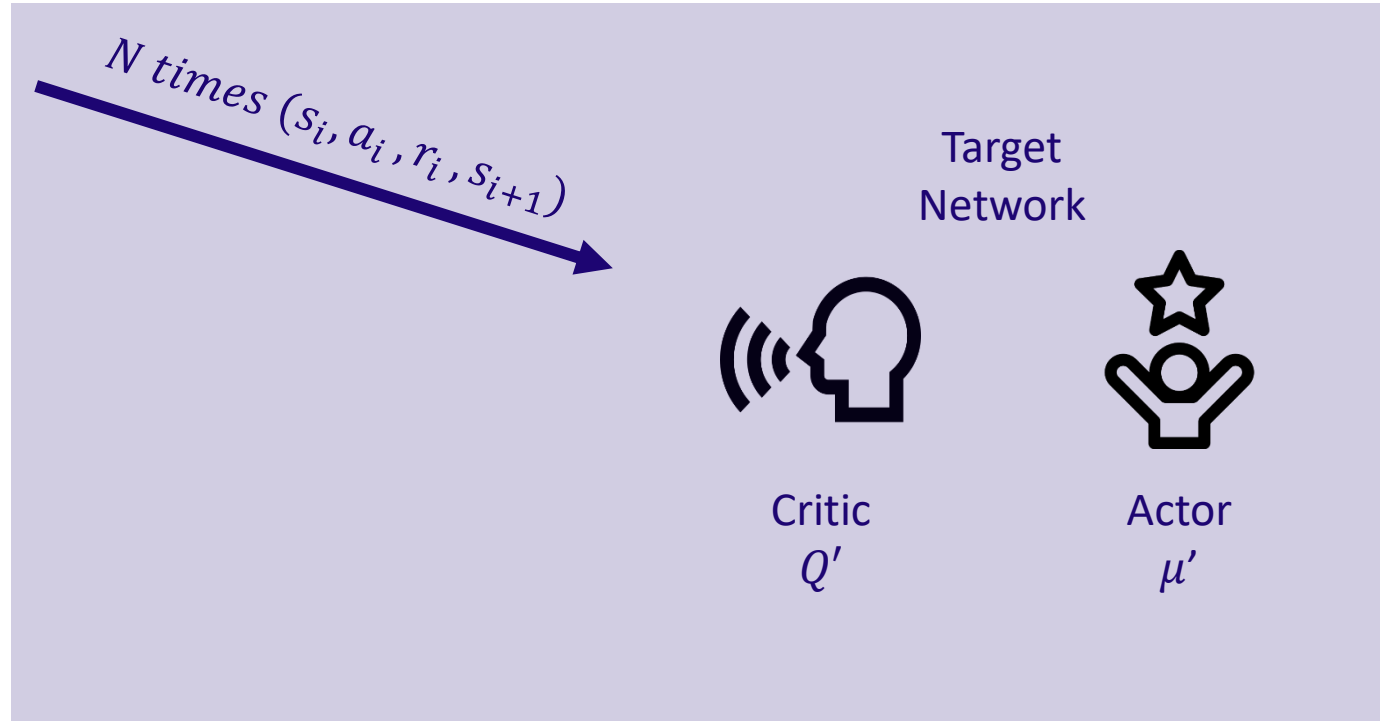
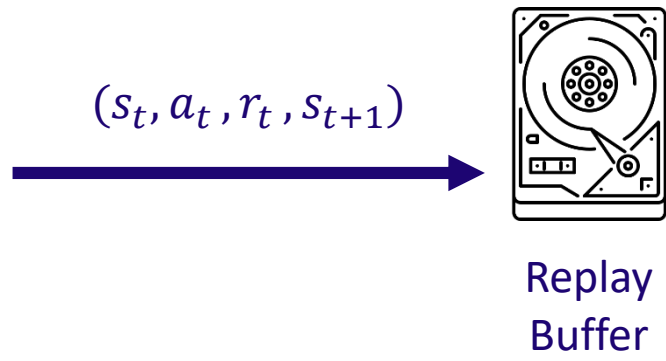
$$a_i \in \{-k_2, -k_1, 0, k_1, k_2\} \quad \longrightarrow \quad 5^3 = 125$$

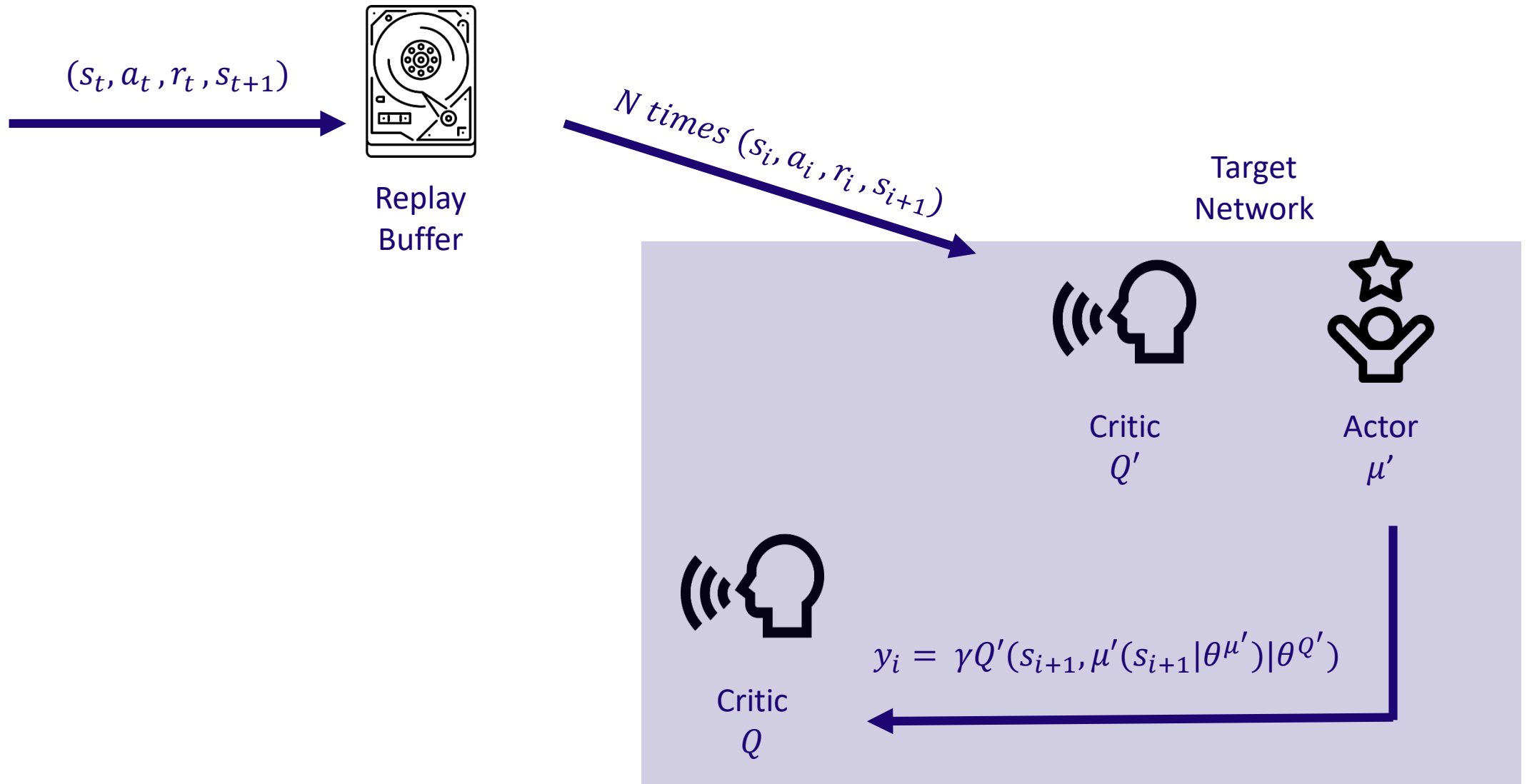


# Deep Deterministic Policy Gradient (DDPG)

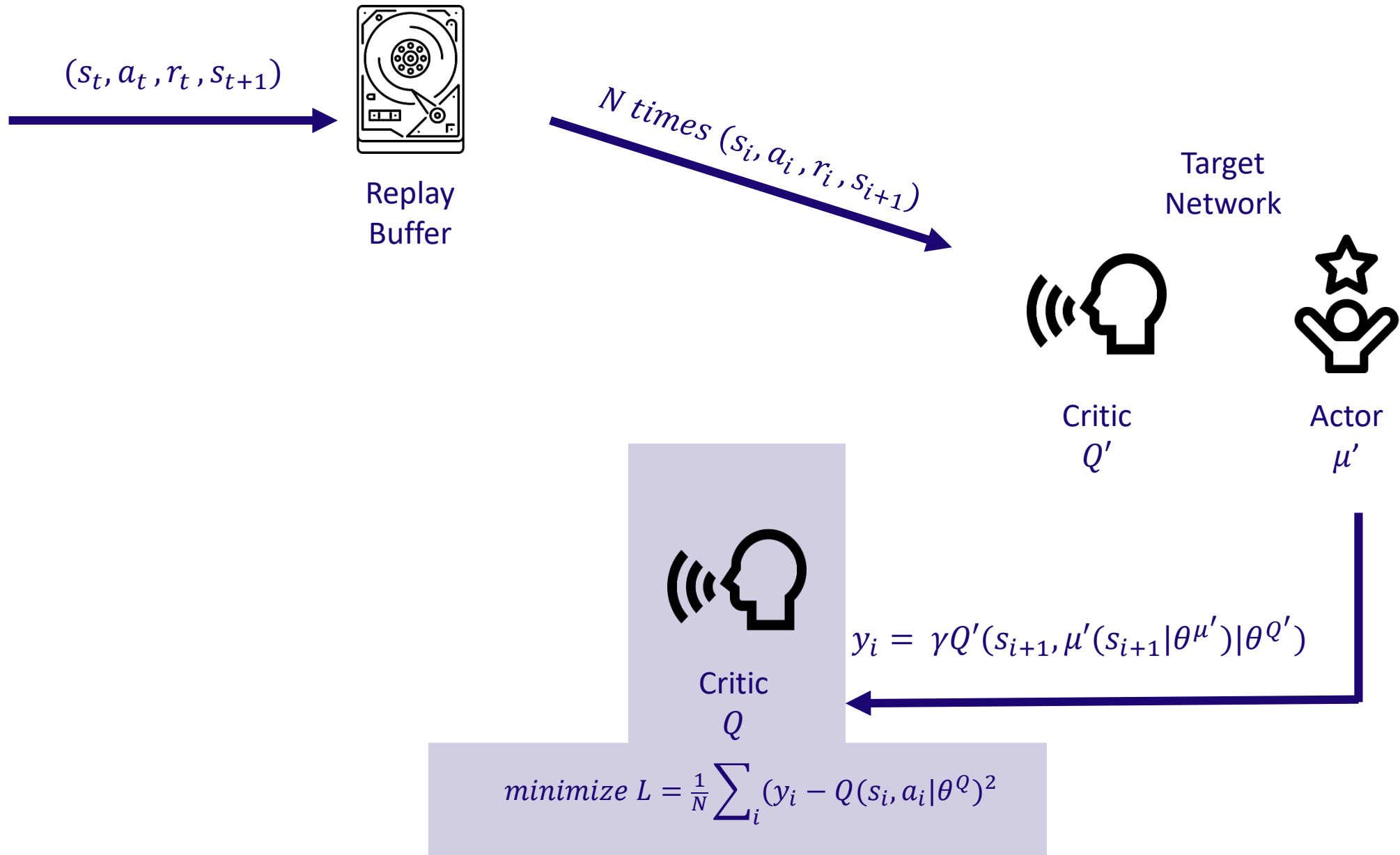


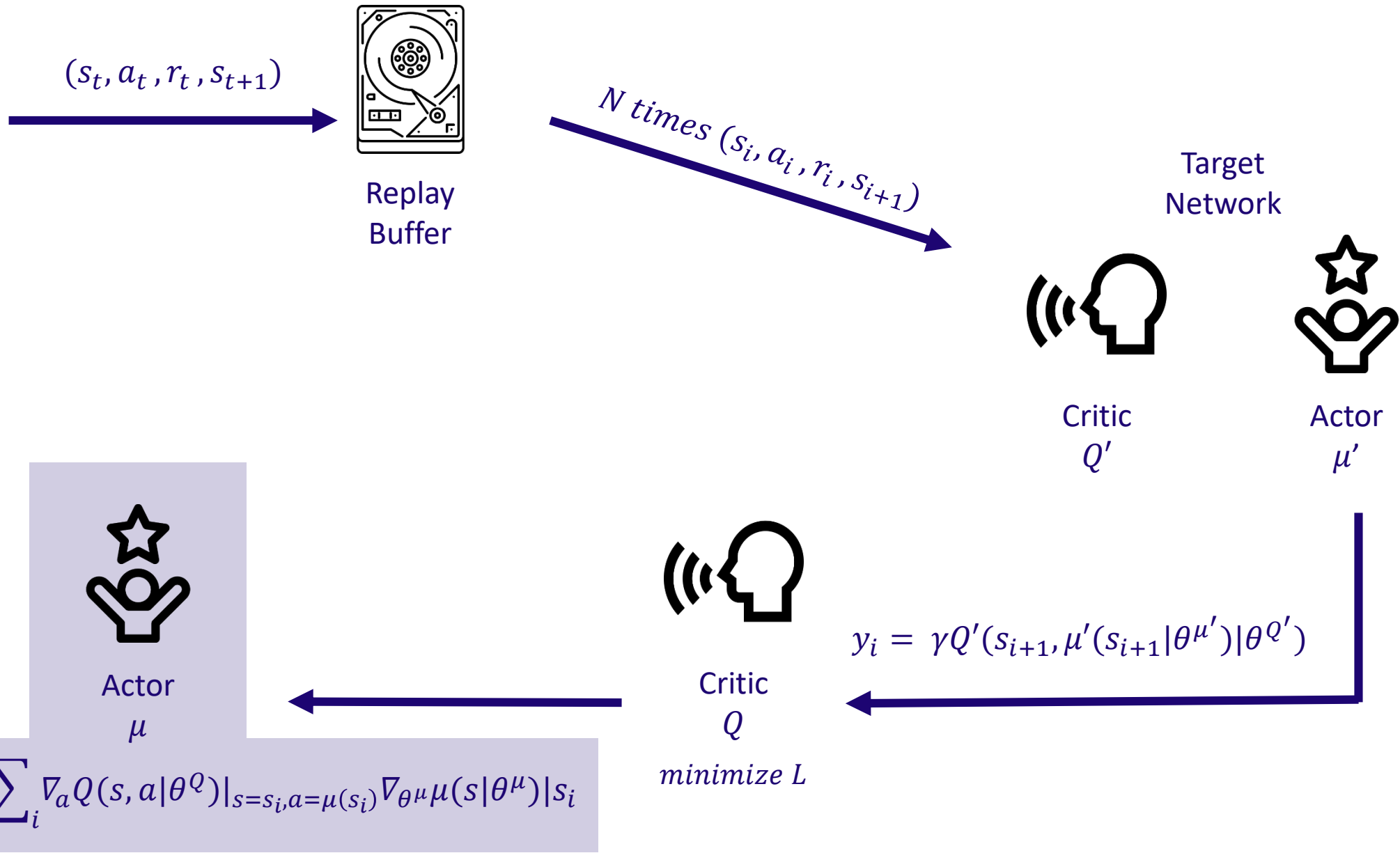


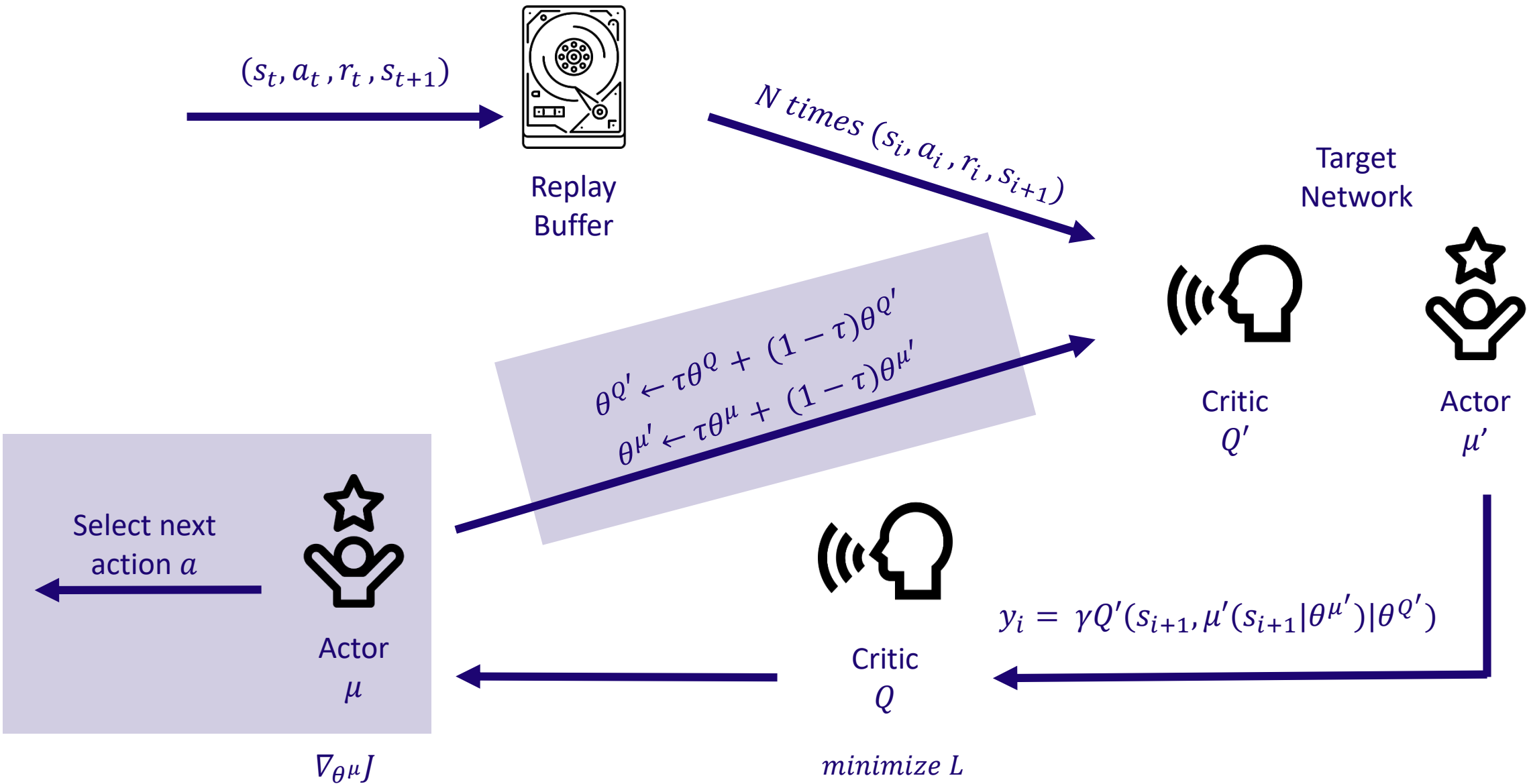


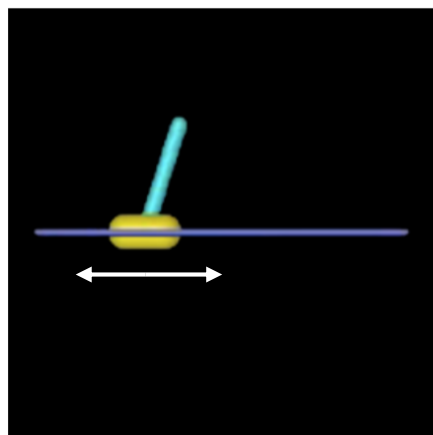
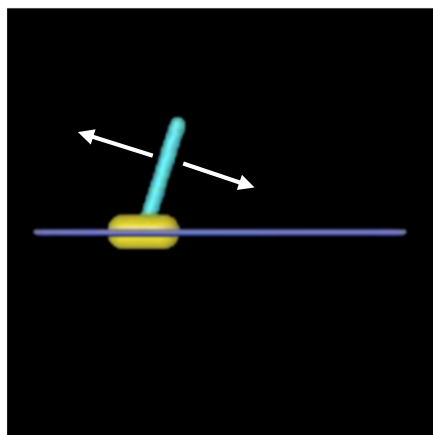
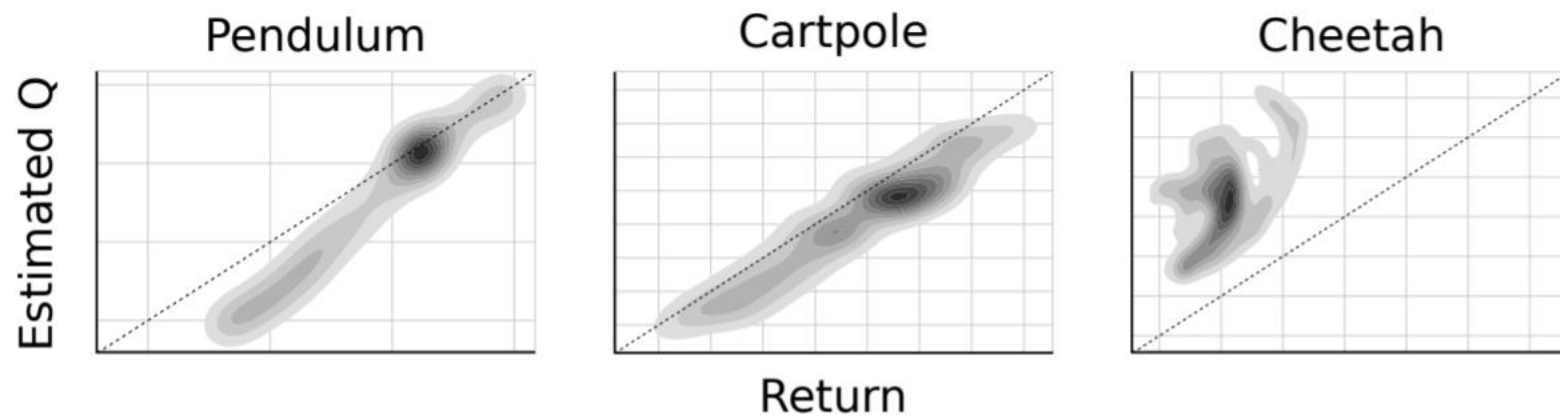


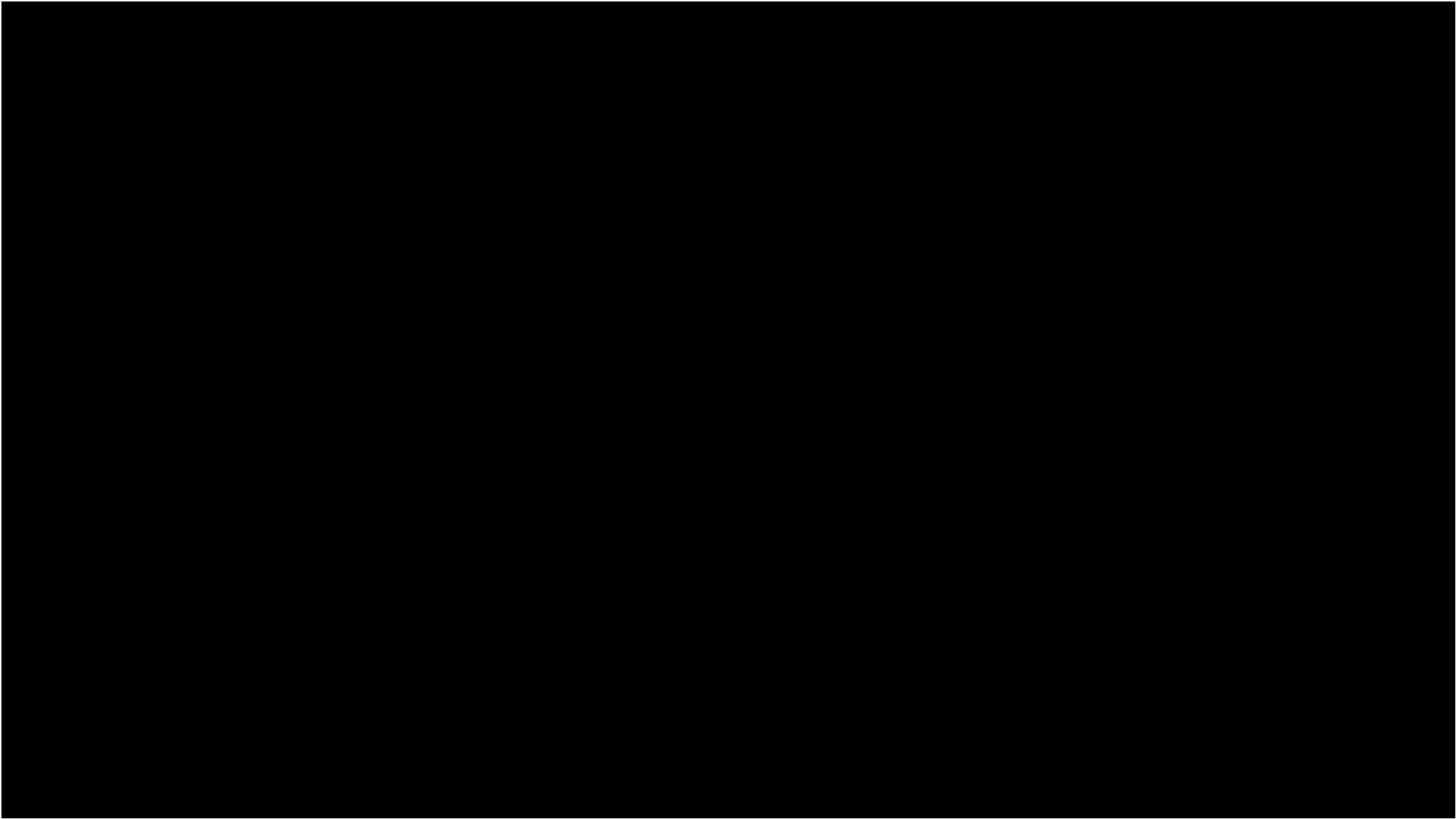












# Trust Region Policy Optimization (TRPO)





$$A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(a)$$



$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(a)$$

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where  $s_0 \sim \rho_0(s_0)$ ,  $a_t \sim \pi(a_t | s_t)$ ,  
 $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(a)$$

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where  $s_0 \sim \rho_0(s_0)$ ,  $a_t \sim \pi(a_t | s_t)$ ,  
 $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$$

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(a)$$

$$\rho_\pi(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots$$

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where  $s_0 \sim \rho_0(s_0)$ ,  $a_t \sim \pi(a_t | s_t)$ ,  
 $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$$

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(a)$$

$$\rho_\pi(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_0 = s) + \dots$$

$$\eta(\pi) = E_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where  $s_0 \sim \rho_0(s_0)$ ,  $a_t \sim \pi(a_t | s_t)$ ,  
 $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$

$$\eta(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a | s) A_\pi(s, a)$$

$$\eta(\tilde{\pi}) = \eta(\pi) + E_{s_0, a_0, \dots \sim \tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$$

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(a)$$

$$\rho_\pi(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_0 = s) + \dots$$

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

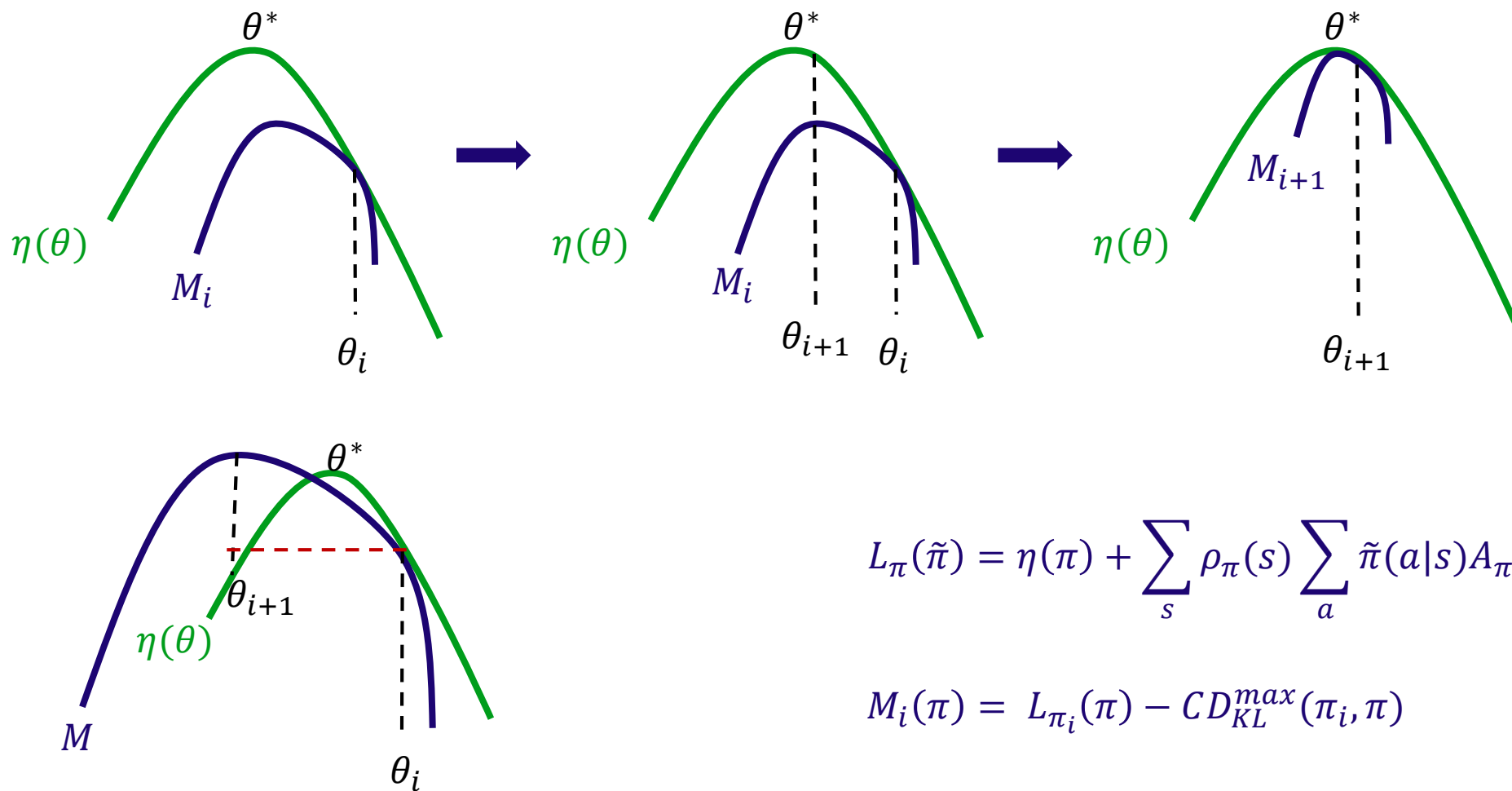
where  $s_0 \sim \rho_0(s_0)$ ,  $a_t \sim \pi(a_t | s_t)$ ,  
 $s_{t+1} \sim P(s_{t+1} | s_t, a_t)$

$$\eta(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a | s) A_\pi(s, a)$$

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$$

$$L_\pi(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_\pi(s) \sum_a \tilde{\pi}(a | s) A_\pi(s, a)$$

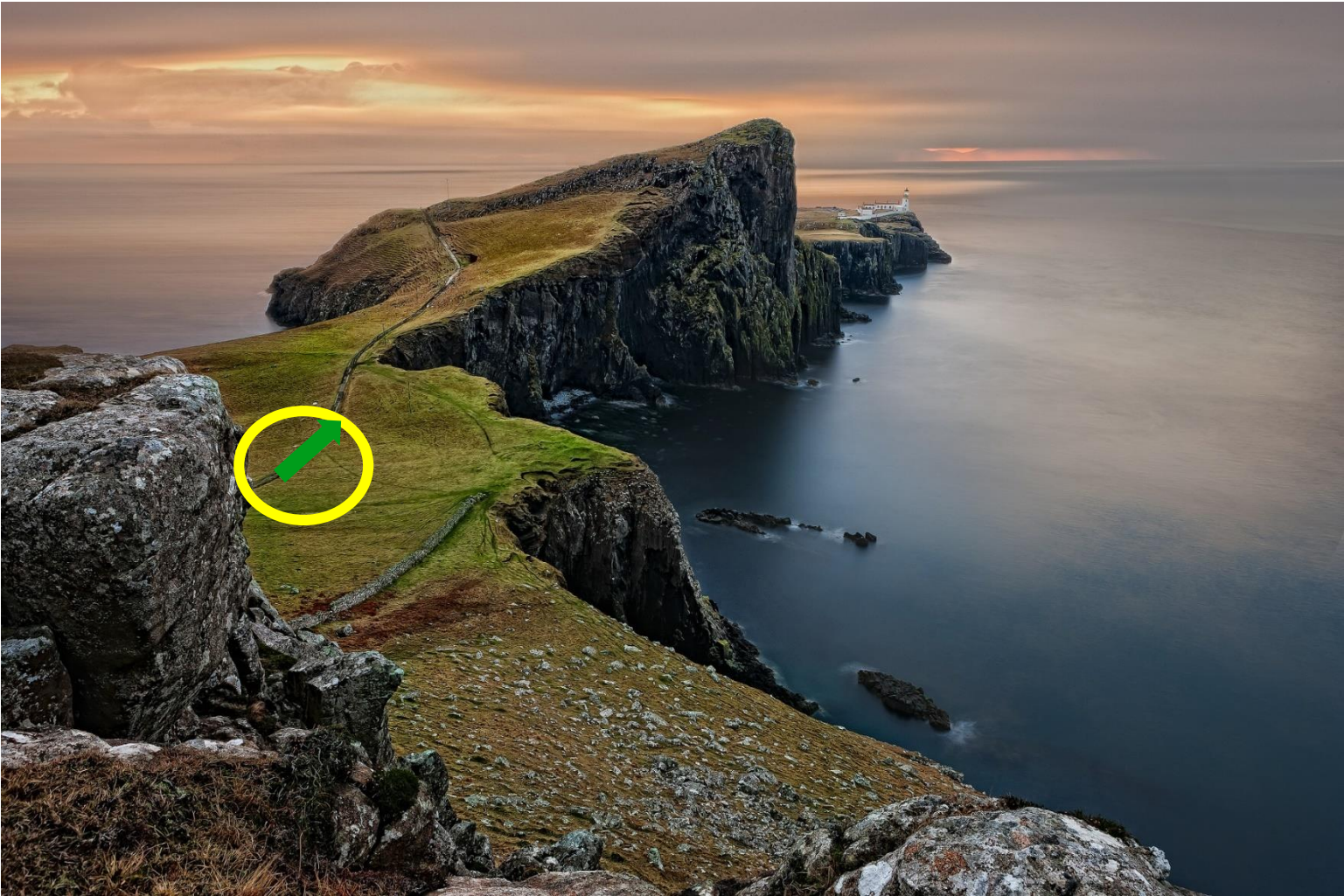
# Minorize-Maximization



$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

$$M_i(\pi) = L_{\pi_i}(\pi) - CD_{KL}^{max}(\pi_i, \pi)$$

# Trust Region



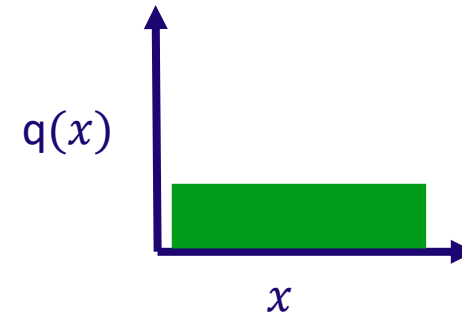
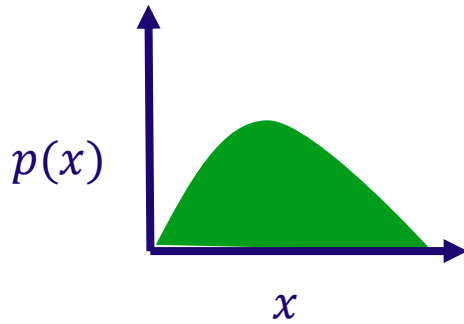




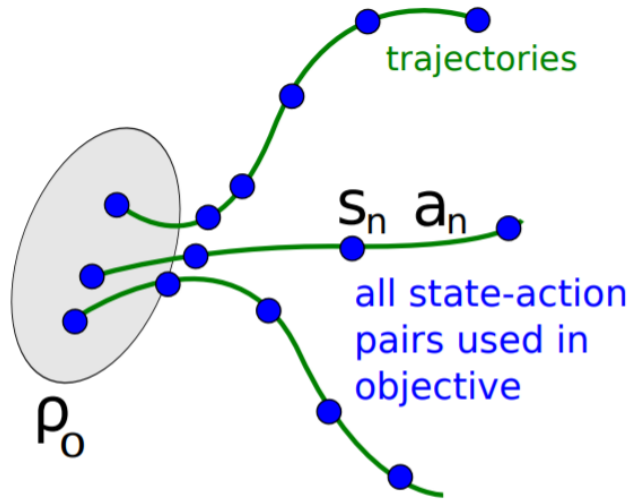


$$\begin{aligned} & \text{maximize}_{\theta} L_{\theta_{old}}(\theta) \\ & \text{subject to } \overline{D}_{KL}^{\rho_{\theta_{old}}}(\theta_{old}, \theta) \leq \delta \end{aligned}$$

# Importance Sampling



$$E_q[f^*(x)] = E_q \left[ \frac{p(x)}{q(x)} f(x) \right] = \sum_x q(x) \frac{p(x)}{q(x)} f(x) = \sum_x p(x) f(x) = E_p[f(x)]$$



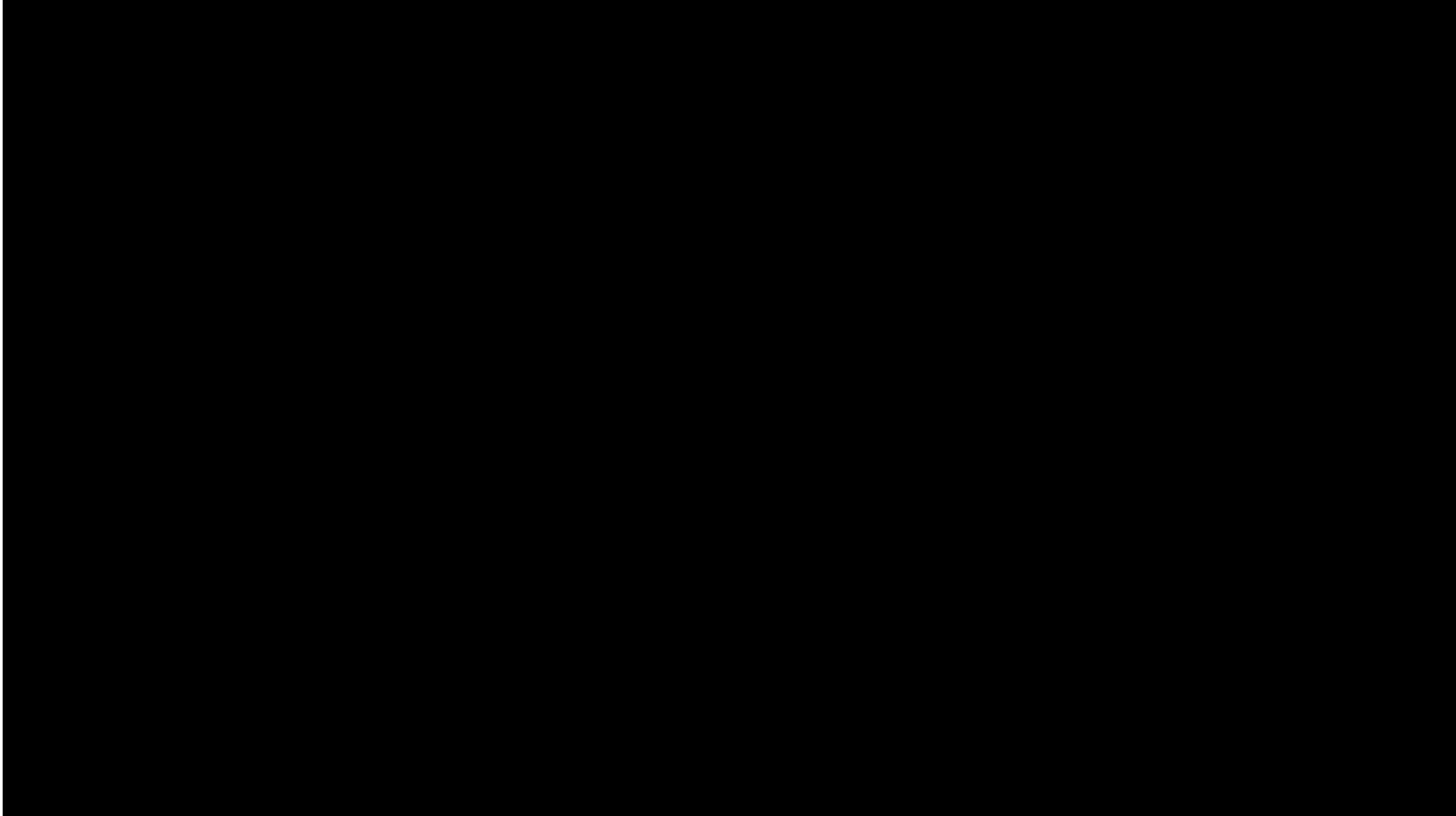
single path procedure

$$\begin{aligned} & \text{maximize}_{\theta} \mathbb{E}_{s \sim \rho_{\theta_{old}}, a \sim q} \left[ \frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{old}}(s, a) \right] \\ & \text{subject to } \mathbb{E}_{s \sim \rho_{\theta_{old}}} \left[ D_{KL}(\pi_{\theta_{old}}(\cdot|s) \parallel \pi_{\theta}(\cdot|s)) \right] \leq \delta \end{aligned}$$

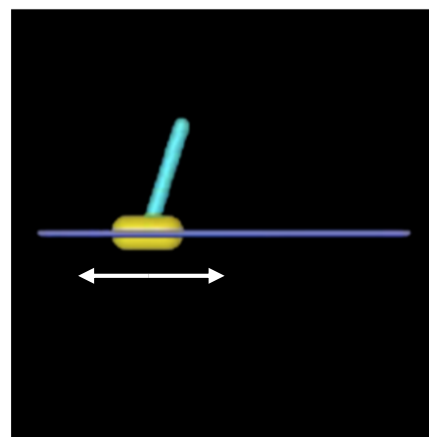
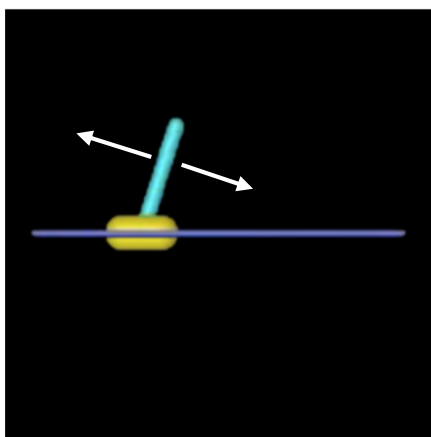
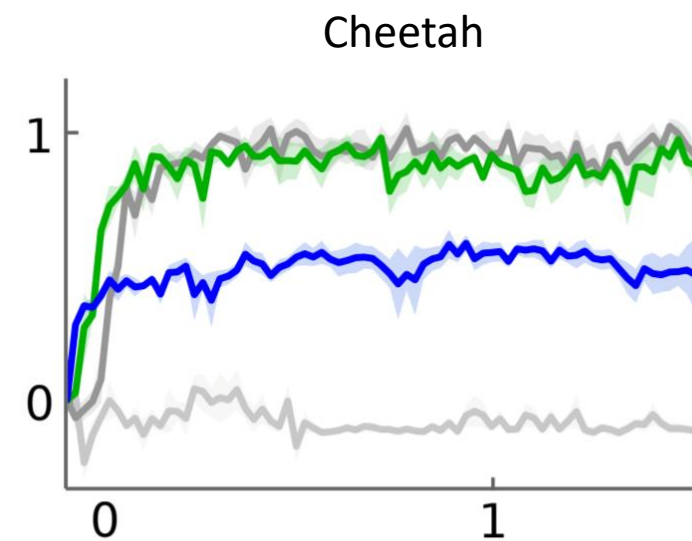
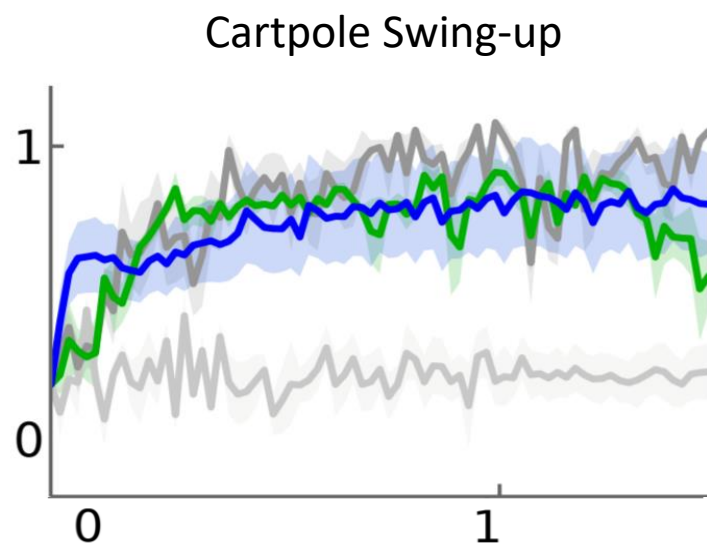
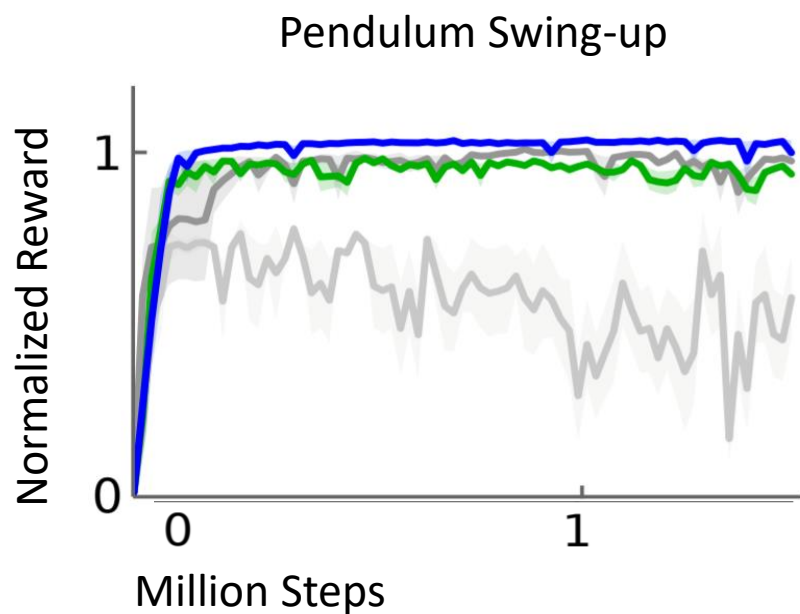


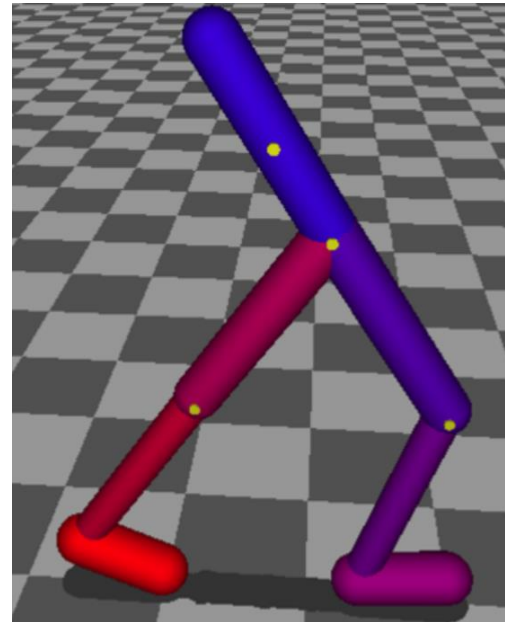
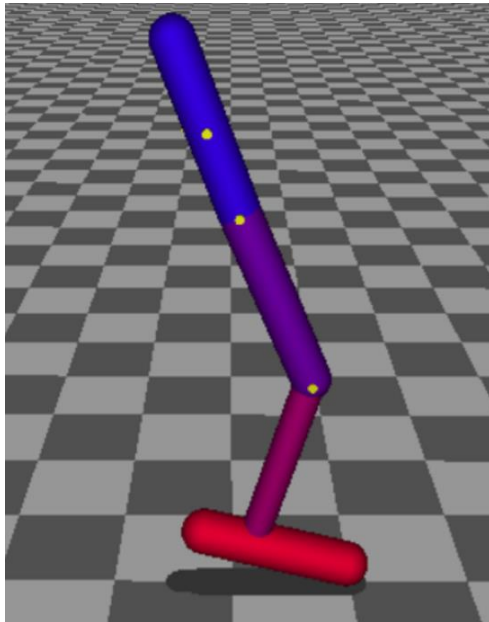
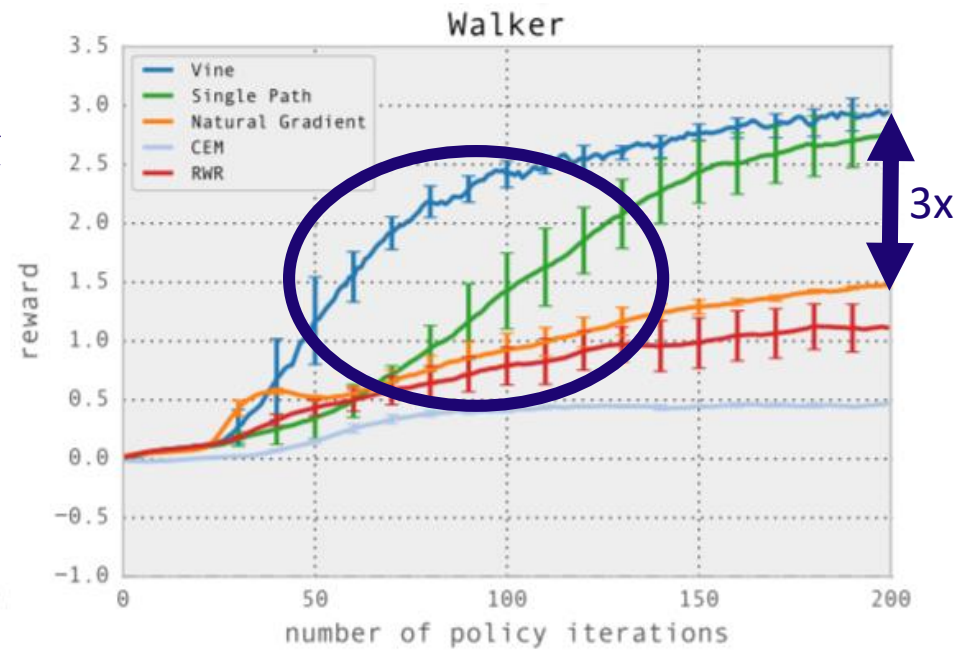
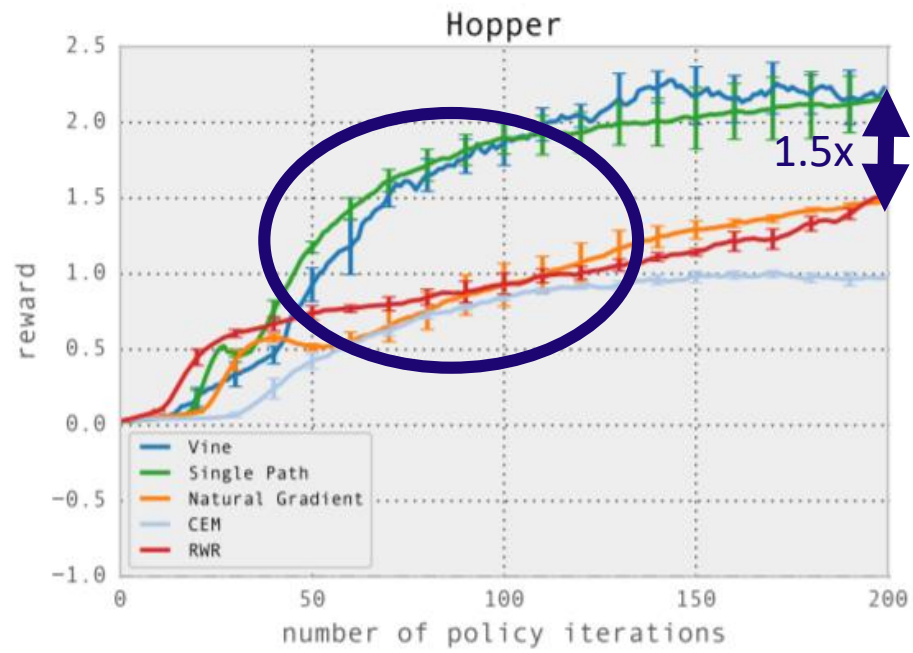
Solve by Approximation  
(Fisher Information Matrix)

# Results

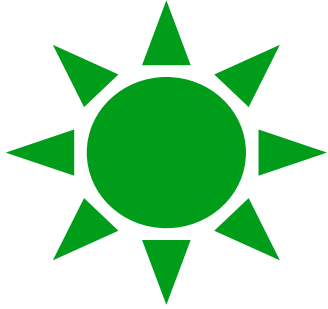


(light grey) – minibatch NFQCA , (grey) – DDPG, (green) – DDPG minibatch, (blue) – DDPG pixel only inputs

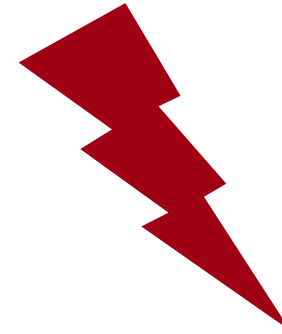




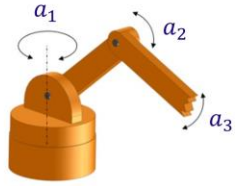
# Conclusion



- General approaches
- Basis for further research activities

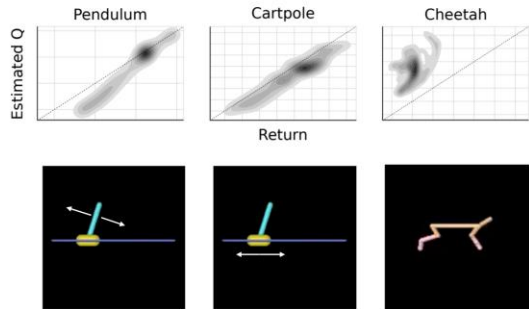
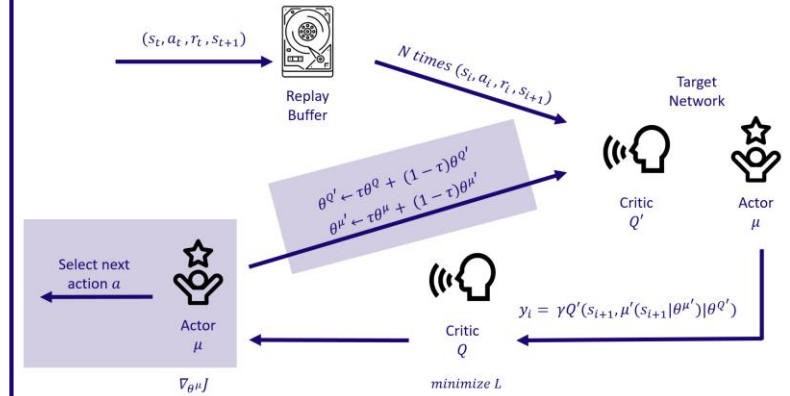
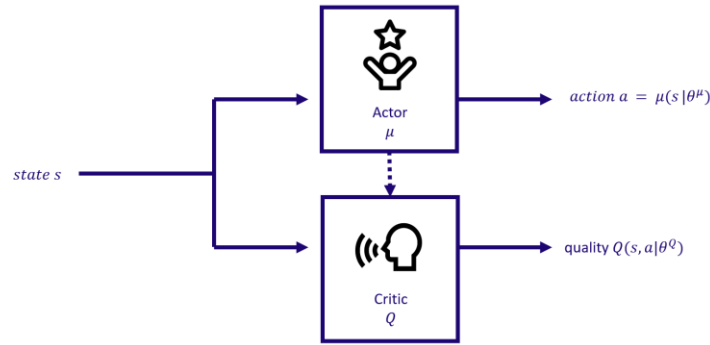


- DDPG → reliability on step size
- TRPO → issues with scalability (computation)



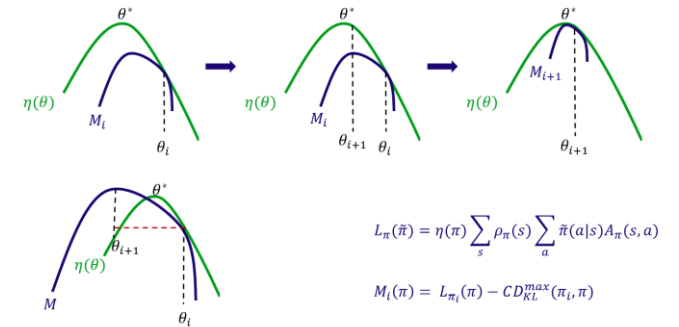
$$a_i \in \{-k_2, -k_1, 0, k_1, k_2\} \rightarrow 5^3 = 125$$

## Deep Deterministic Policy Gradient (DDPG)



# ? Questions ?

## Minorize-Maximization



## Importance Sampling



$$E_q[f^*(x)] = E_q\left[\frac{p(x)}{q(x)} f(x)\right] = \sum_x q(x) \frac{p(x)}{q(x)} f(x) = \sum_x p(x) f(x) = E_p[f(x)]$$

