# Continuous Control with Deep Reinforcement Learning

Timothy P. Lillicrap, et al. · 14 p · 2015

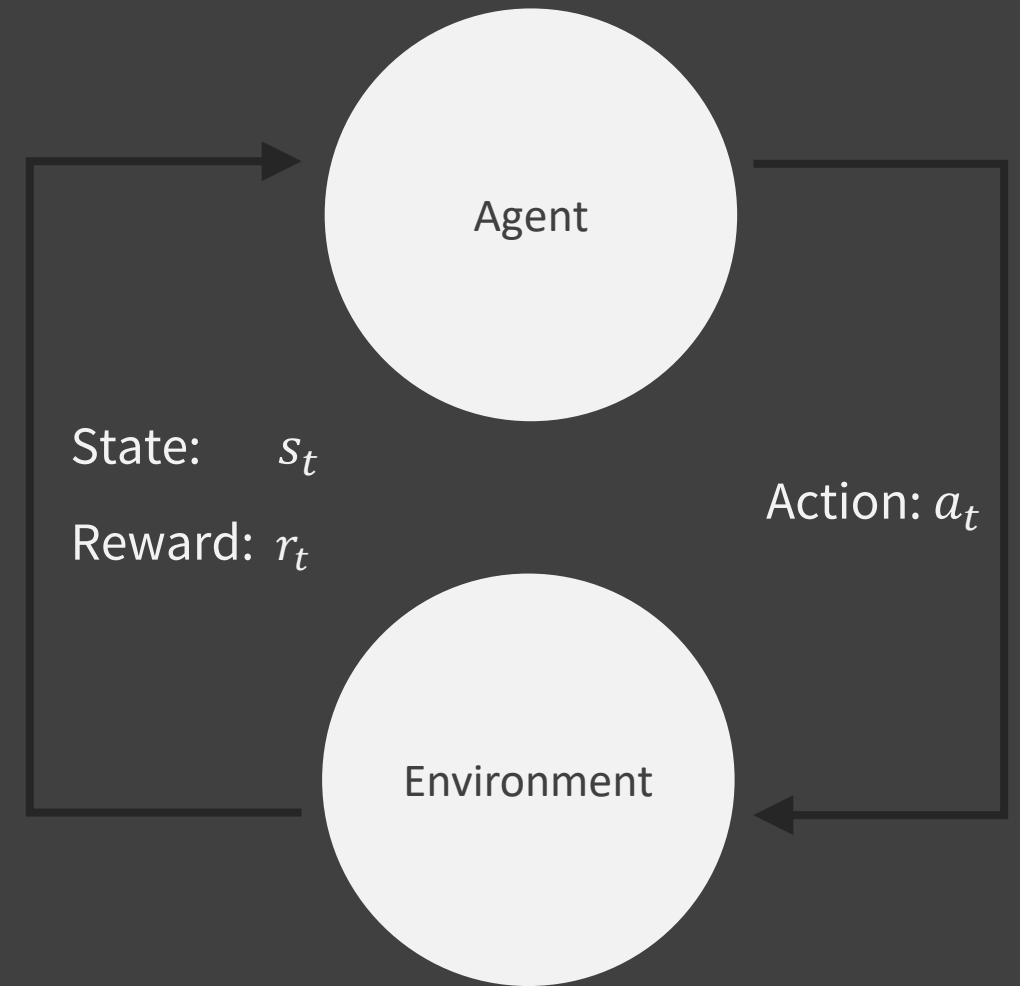By Flavio Schneider · November 2019 · **ETH** *zürich*

# Index

- What is reinforcement learning (RL)?

- How can we describe RL formally?

- How does a RL *algorithm* look like?

# What is Reinforcement Learning?

- One of the three basic paradigms of ML

- Learns intelligent behavior from reward

- Many applications

- Exploitation and exploration

Agent-Environment loop



Agent

State: $s_t$
Reward: $r_t$

Action: $a_t$

Environment

# What is Reinforcement Learning?

State · Action · Reward · Policy · Goal


Pendulum Example

# What is Reinforcement Learning?

State · Action · Reward · Policy · Goal

State: $s_t \in S = R^{n+1}$ describes to the agent the

environment completely at time $t$.

State Vector

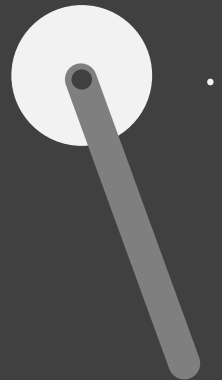$$s_t = \begin{bmatrix} Pendulum\ angle \\ Pendulum\ speed \end{bmatrix}$$

$t = 0$      $t = 5$      $t = 10$

...  ...  ...

$$s_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \qquad s_5 = \begin{bmatrix} 20 \\ -5 \end{bmatrix} \qquad s_{10} = \begin{bmatrix} -20 \\ 4 \end{bmatrix}$$

# What is Reinforcement Learning?

State · **Action** · Reward · Policy · Goal

Action: $a_t \in A$ represents what the agent is going to do at

time $t$.

Discrete Action Space

$$a_t \in A = (Left, Right) = (a^{(0)}, a^{(1)})$$

Continuous Action Space

$$a_t \in A = [-1, 1]$$

# What is Reinforcement Learning?

State · Action · **Reward** · Policy · Goal

Reward: $r_t = R(s_t)$ or $r_t = R(s_t, a_t)$ is a score that tells the agent how good was the last move $a_t$ in state $s_t$.

Return: $R(\tau) = \sum_{t=0}^{T} \gamma^t r_t$ where $\tau = (s_0, a_0, s_1, a_1, \dots)$ is a sequence of state/action pairs and $\gamma \in (0,1)$.
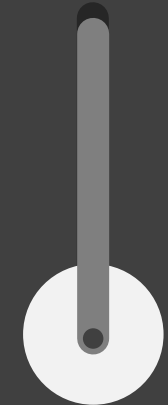
Reward Function

$$r_t = R(s_t) = \begin{cases} 1 & if\ s_t[0] == 180 \\ 0 & otherwise \end{cases}$$

$t = 0$ \qquad\qquad $t = 100$



$\dots$

$r_0 = 0$ \qquad\qquad $r_{100} = 1$

# What is Reinforcement Learning?

State · Action · Reward · **Policy** · Goal

Policy: $\mu(s_t) = a_t$ is the brain of the agent, a function that

maps states to actions.



$$s_5 = \begin{bmatrix} 20 \\ -5 \end{bmatrix}$$

$$\mu\left(\begin{bmatrix} 20 \\ -5 \end{bmatrix}\right) = a_5^{(0)} = Left$$

$$s_5 = \begin{bmatrix} 20 \\ -5 \end{bmatrix}$$

$$\mu\left(\begin{bmatrix} 20 \\ -5 \end{bmatrix}\right) = a_5^{(i)} = 0.83$$

# What is Reinforcement Learning?

State · Action · Reward · Policy · **Goal**

Trajectory Probability:

$$P(\tau \mid \mu) := \prod_{t=0}^{T-1} P(s_{t+1} \mid s_t, a_t) P(\mu(s_t) = a_t)$$
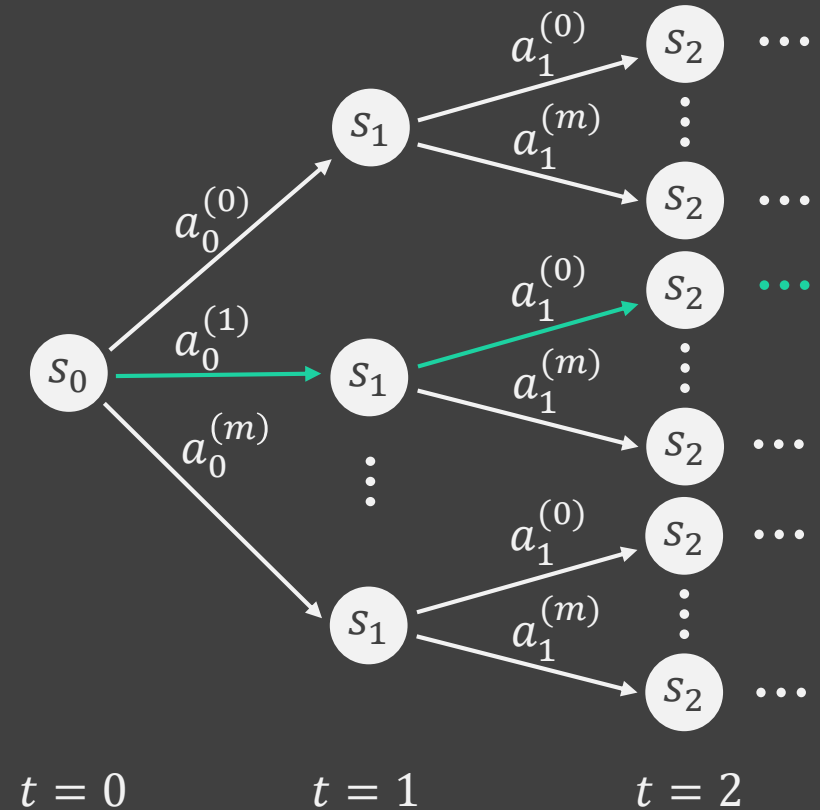
Expected Return:

$$\mathbb{E}_{a_t = \mu(s_t)}[R(\tau)] = \int_{\tau} P(\tau \mid \mu) R(\tau)$$

Goal · Find Optimal Policy:

$$\mu^* := \operatorname*{argmax}_{\mu} \mathbb{E}_{a_t = \mu(s_t)}[R(\tau)]$$

Trajectory

$$\tau = \left( s_0, a_0^{(1)}, s_1, a_1^{(0)}, \dots \right)$$



$t = 0 \qquad t = 1 \qquad t = 2$

# Important Functions

Quality · Bellmann

# Important Functions
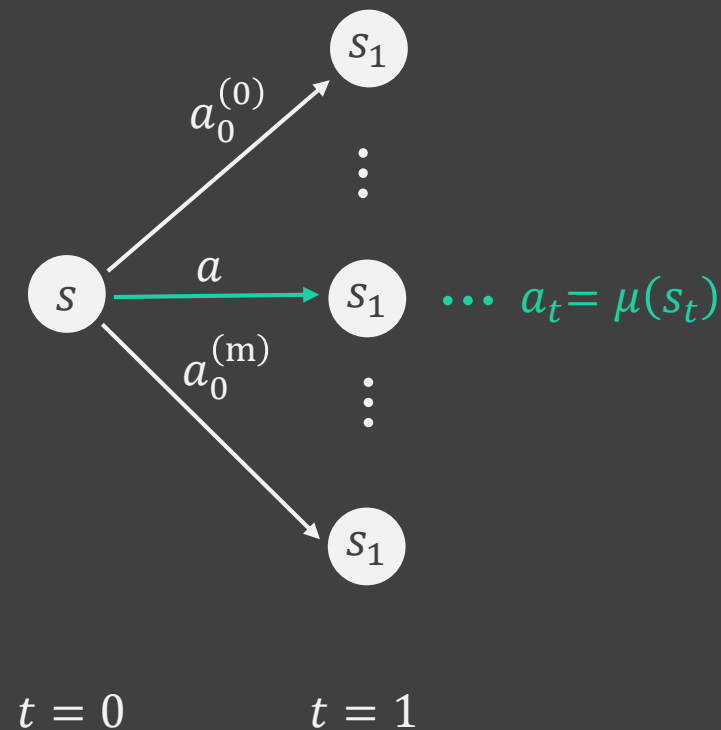
Quality Function:

$$\mathbb{E}_{a_t=\mu(s_t)}[R(\tau)]$$

$$Q^\mu(s,a) := \mathbb{E}_{s_0=s, a_0=a,\ a_t=\mu(s_t)}[R(\tau)]$$

Optimal Quality Function:

$$Q^*(s,a) := \max_\mu Q^\mu(s,a)$$

Optimal Action:

$$\mu^*(s) = \arg\max_a Q^*(s,a) = a^*$$



$$a_0^{(0)}$$

$$a$$

$$a_0^{(m)}$$

$$a_t = \mu(s_t)$$

$$t = 0 \qquad t = 1$$

# Important Functions

Quality · Bellman

Bellman Equation:

$$Q^*(s, a) = \mathbb{E}\left[R(s, a) + \gamma \max_{a'} Q^*(s', a')\right]$$

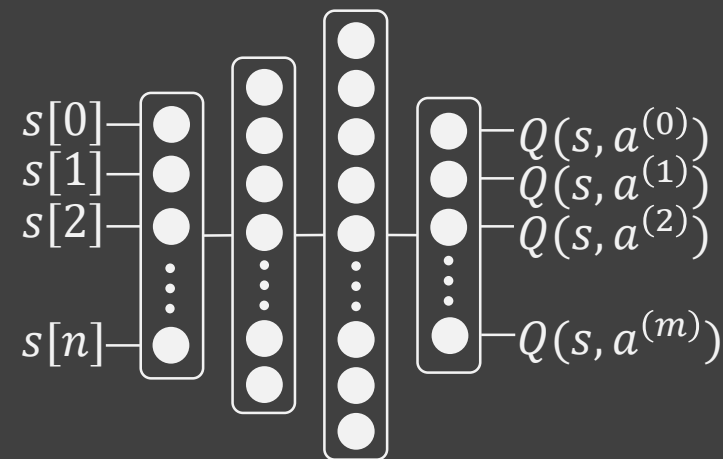# Algorithms Evolution

DQN · DPG · DDPG

Question:

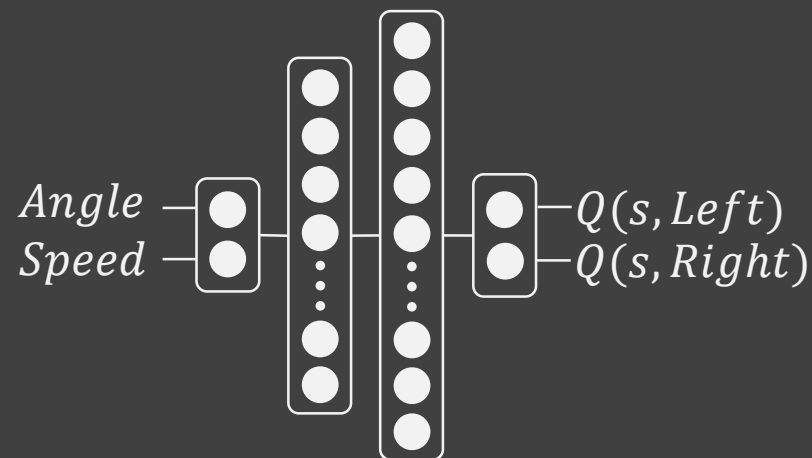How can we solve any (discrete and then continuous)

reinforcement learning problem?

Idea:

Approximate $Q^*(s, a)$ using a neural network.
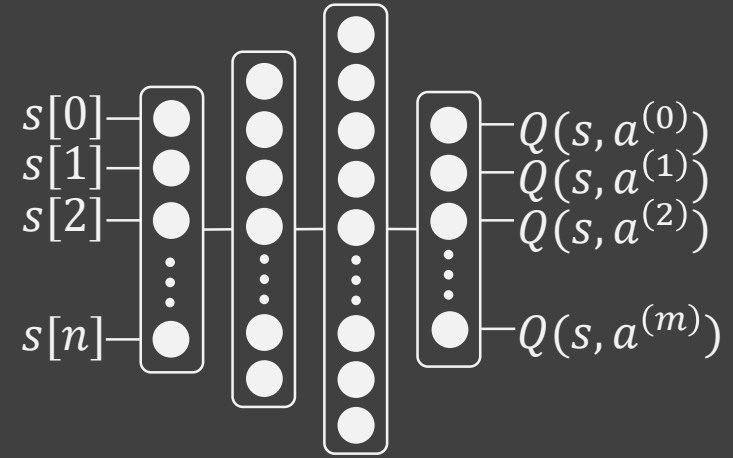
"Critic" Network · $Q^\mu(s, a \mid \theta^Q)$:

$$s[0] \quad s[1] \quad s[2] \quad \cdots \quad s[n]$$

$$Q(s, a^{(0)}) \quad Q(s, a^{(1)}) \quad Q(s, a^{(2)}) \quad \cdots \quad Q(s, a^{(m)})$$

Pendulum

$$Angle \quad Speed$$

$$Q(s, Left) \quad Q(s, Right)$$

# Algorithms Evolution

       *Deep Q Network*

1. Initialize randomly $Q^\mu(s, a \mid \theta^Q)$

2. Get initial state $s$

3. Repeat:

   a.   $a = \underset{a'}{argmax} \, Q^\mu(s, a' \mid \theta^Q)$

   b.   $a = a^{(i)}$ for $i \in \{0, \dots, m\}$ random with probability $\epsilon$

   c.   Execute $a$ and observe $r = R(s, a)$ and $s'$

   d.   $Q^T(s, a, r, s') := r + \gamma \underset{a'}{max} \, Q^\mu(s', a' \mid \theta^Q)$

   e.   $\theta^Q \leftarrow \theta^Q - \alpha \nabla_{\theta^Q} \left[ \left( Q^T(s, a, r, s') - Q^\mu(s, a \mid \theta^Q) \right)^2 \right]$

   f.   $s \leftarrow s'$

"Critic" Network · $Q^\mu(s, a \mid \theta^Q)$:



Bellman Equation

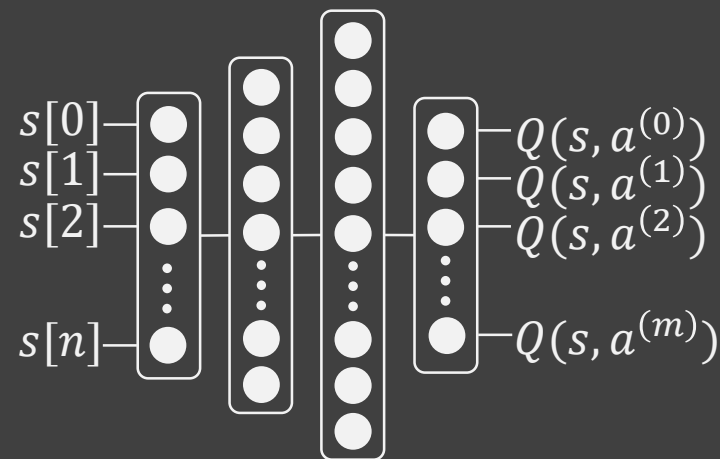$$Q^*(s, a) = \mathbb{E}\left[ R(s, a) + \gamma \, \underset{a'}{max} \, Q^*(s', a') \right]$$

# Algorithms Evolution

*Deep Q Network*

1. Initialize randomly $Q^\mu(s, a \mid \theta^Q)$

2. Get initial state $s$

3. Repeat:

   a. $a = \underset{a'}{argmax}\, Q^\mu(s, a' \mid \theta^Q)$

   b. $a = a^{(i)}$ for $i \in \{0, \dots, m\}$ random with probability $\epsilon$

   c. Execute $a$ and observe $r = R(s, a)$ and $s'$

   d. $Q^T(s, a, r, s') \coloneqq r + \gamma\, \underset{a'}{max}\, Q^\mu(s', a' \mid \theta^Q)$

   e. $\theta^Q \leftarrow \theta^Q - \alpha \nabla_{\theta^Q} \left[ \left( Q^T(s, a, r, s') - Q^\mu(s, a \mid \theta^Q) \right)^2 \right]$

   f. $s \leftarrow s'$

"Critic" Network · $Q^\mu(s, a \mid \theta^Q)$:



Bellman Equation

$$Q^*(s, a) = \mathbb{E}\left[ R(s, a) + \gamma\, \underset{a'}{max}\, Q^*(s', a') \right]$$
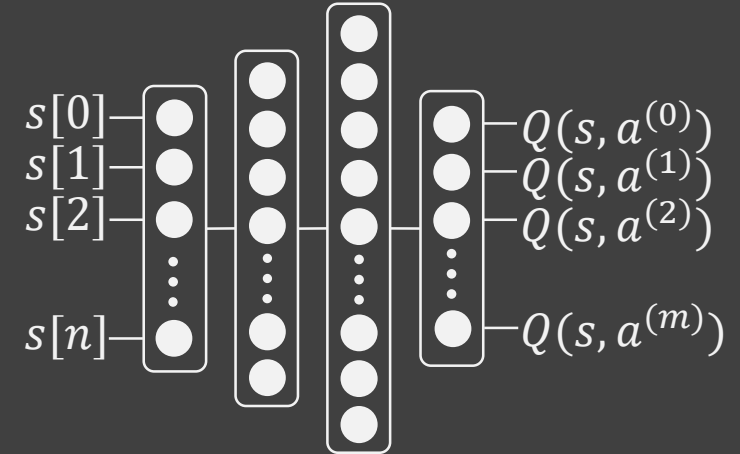
# Algorithms Evolution

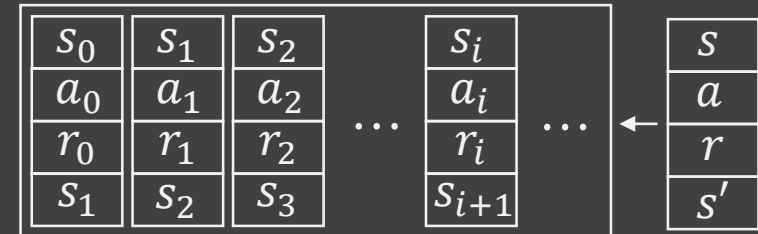DQN · DPG · DDPG                    *Deep Q Network*

Code Structure

1. *Initialize randomly* $Q^{\mu}\left(s, a \mid \theta^{Q}\right)$ *and* $Q^{\mu'}\left(s, a \mid \theta^{Q'}\right)$

2. *Initialize replay buffer* $\mathcal{R} = S \times A \times R \times S$

3. *Get initial state* $s$

4. *Repeat:*

   a. ■ · *Sample*

   b. ■ · *Train*

Target "Critic" Network · $Q^{\mu'}(s, a \mid \theta^{Q'})$ :



Replay Buffer · $\mathcal{R}$ :

| $s_0$ | $s_1$ | $s_2$ | ... | $s_i$ | ... | $s$ |
|-------|-------|-------|-----|-------|-----|-----|
| $a_0$ | $a_1$ | $a_2$ | | $a_i$ | | $a$ |
| $r_0$ | $r_1$ | $r_2$ | | $r_i$ | | $r$ |
| $s_1$ | $s_2$ | $s_3$ | | $s_{i+1}$ | | $s'$ |

# Algorithms Evolution
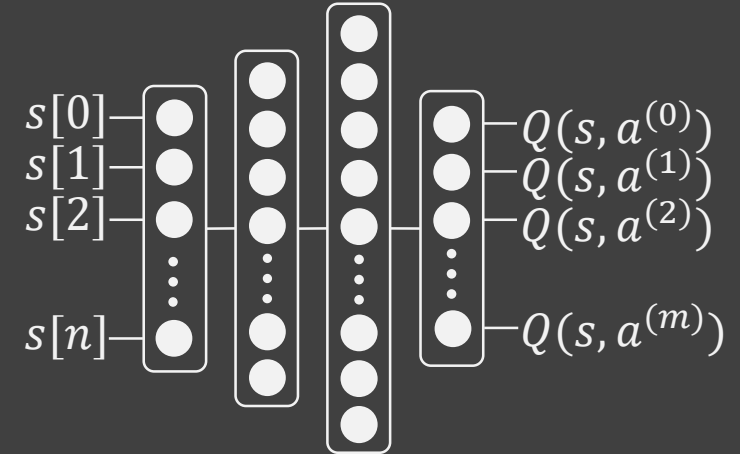
DQN · DPG · DDPG                    *Deep Q Network*
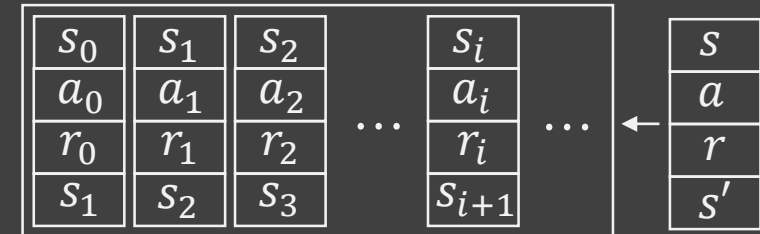
🟨 · *Sample*

i.   $a = \underset{a'}{argmax}\ Q^\mu(s, a' \mid \theta^Q)$

ii.  $a = a^{(i)}$ *for* $i \in \{0, \dots, m\}$ *random with probability* $\epsilon$

iii. *Execute* $a$ *and observe* $r$ *and* $s'$

iv.  *Store* $\mathcal{R} \leftarrow \mathcal{R} \cup (s, a, r, s')$

Target "Critic" Network · $Q^{\mu'}(s, a \mid \theta^{Q'})$ :

$s[0]$ — $Q(s, a^{(0)})$
$s[1]$ — $Q(s, a^{(1)})$
$s[2]$ — $Q(s, a^{(2)})$

$s[n]$ — $Q(s, a^{(m)})$

Replay Buffer · $\mathcal{R}$ :

| $s_0$ | $s_1$ | $s_2$ | | $s_i$ | | $s$ |
|-------|-------|-------|---|-------|---|-----|
| $a_0$ | $a_1$ | $a_2$ | $\dots$ | $a_i$ | $\dots$ | $a$ |
| $r_0$ | $r_1$ | $r_2$ | | $r_i$ | | $r$ |
| $s_1$ | $s_2$ | $s_3$ | | $s_{i+1}$ | | $s'$ |

# Algorithms Evolution

DQN · DPG · DDPG                    *Deep Q Network*

■ · *Train*

i.    *Sample random batch* $\mathcal{B} \subseteq \mathcal{R}$
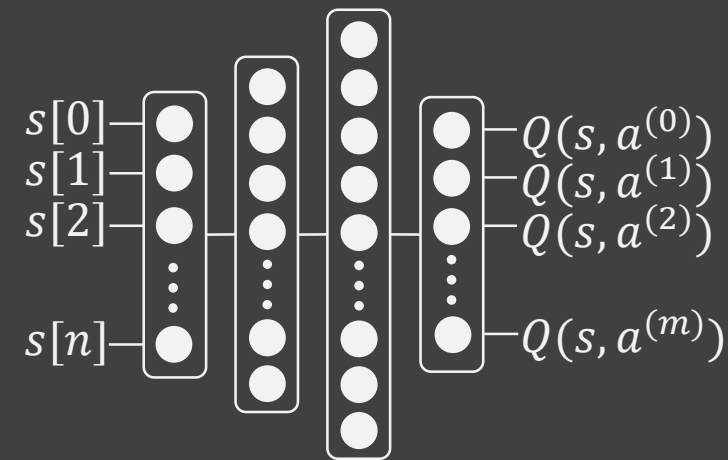
ii.   *For each* $(s, a, r, s') \in \mathcal{B}$

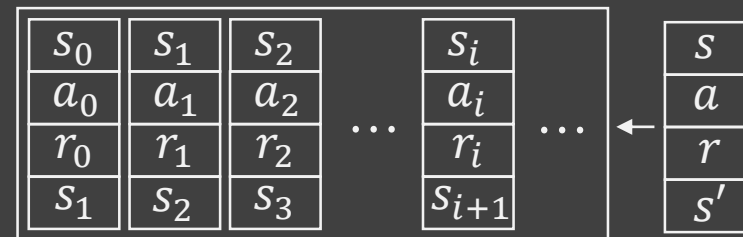  I.    $Q^T(s, a, r, s') := r + \gamma \max_{a'} Q^{\mu'}(s, a' \mid \theta^{Q'})$

  II.   $\theta^Q \leftarrow \theta^Q - \alpha \nabla_{\theta^Q} \left[ \left( Q^T(s, a, r, s') - Q^{\mu}(s, a \mid \theta^Q) \right)^2 \right]$

iii.  *Every* $C$ *steps reset* $\theta^{Q'} \leftarrow \theta^Q$

Target "Critic" Network · $Q^{\mu'}(s, a \mid \theta^{Q'})$ :



Replay Buffer · $\mathcal{R}$ :

# Algorithms Evolution

DQN · DPG · DDPG                    *Deep Q Network*

🟦 · *Train*

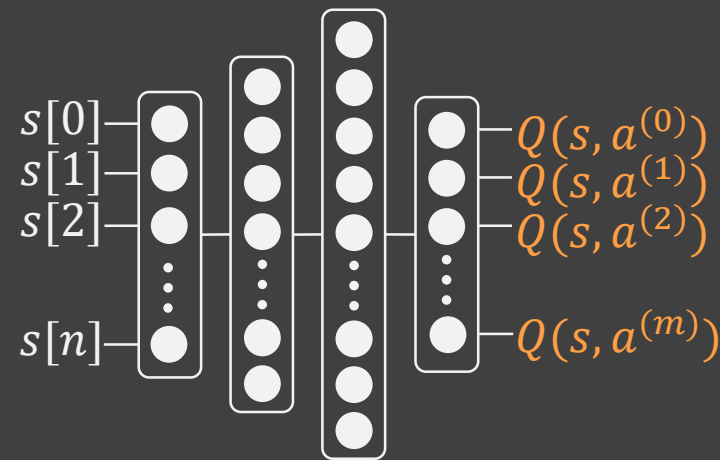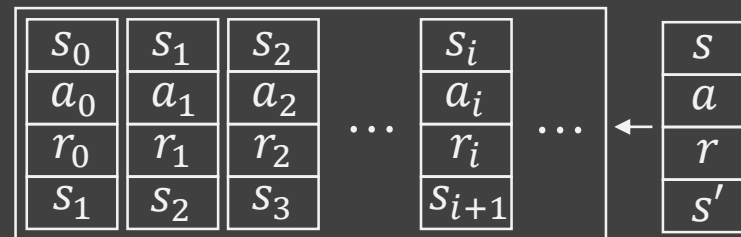i.  *Sample random batch* $\mathcal{B} \subseteq \mathcal{R}$

ii. *For each* $(s, a, r, s') \in \mathcal{B}$

    I.    $Q^T(s, a, r, s') := r + \gamma \max_{a'} Q^{\mu'}(s, a' \mid \theta^{Q'})$

    II.    $\theta^Q \leftarrow \theta^Q - \alpha \nabla_{\theta^Q} \left[ \left( Q^T(s, a, r, s') - Q^{\mu}(s, a \mid \theta^Q) \right)^2 \right]$

iii. *Every* $C$ *steps reset* $\theta^{Q'} \leftarrow \theta^Q$

Target "Critic" Network · $Q^{\mu'}(s, a \mid \theta^{Q'})$ :



Replay Buffer · $\mathcal{R}$ :

# Algorithms Evolution

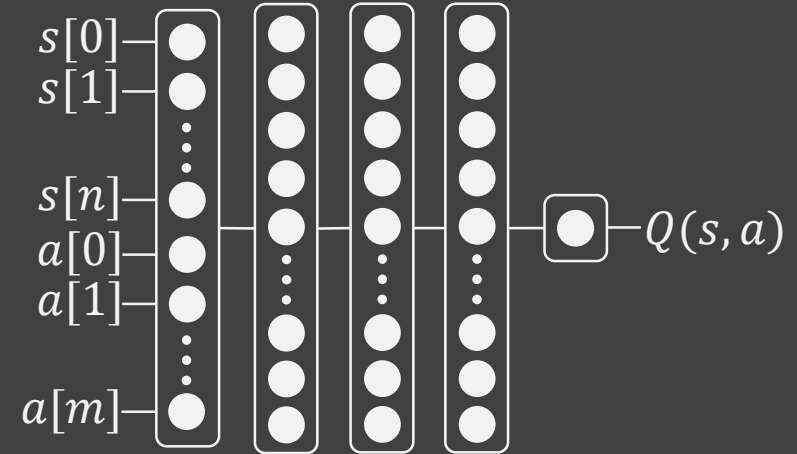DQN · **DPG** · DDPG        *Deterministic Policy Gradient*

Expected Return:

$$J(\mu_\theta) := \mathbb{E}_{a_t = \mu_\theta(s_t)}[R(\tau)]$$

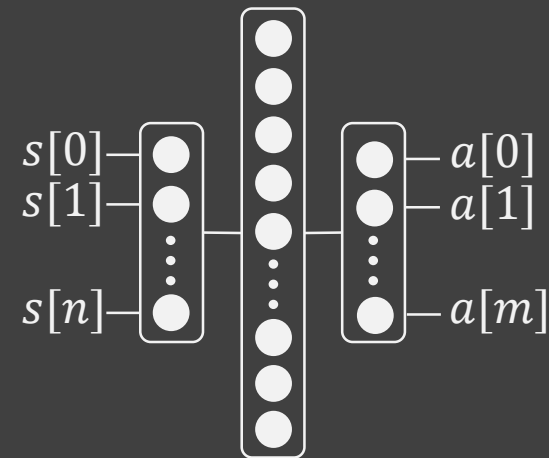Deterministic Policy Gradient Theorem:

$$\nabla_{\theta^\mu} J(\mu_{\theta^\mu}) = \mathbb{E}\left[\nabla_a Q^\mu(s, a \mid \theta^Q)\nabla_{\theta_\mu}\mu(s \mid \theta^\mu)\right]$$

"Critic" Network · $Q^\mu(s, a \mid \theta^Q)$:



"Actor" Network · $\mu(s \mid \theta^\mu)$:

# Algorithms Evolution
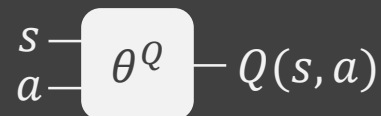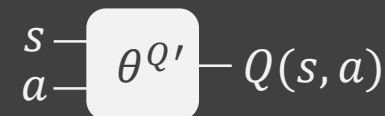
DQN · DPG · DDPG                    *Deep DPG*

Code Structure

1. *Initialize randomly* $Q^\mu(s, a \mid \theta^Q)$ *and* $\mu(s \mid \theta^\mu)$

2. *Initialize* $Q^{\mu\prime}(s, a \mid \theta^{Q\prime} \leftarrow \theta^Q)$ *and* $\mu'(s \mid \theta^{\mu\prime} \leftarrow \theta^\mu)$

3. *Initialize replay buffer* $\mathcal{R} = S \times A \times R \times S$

4. *Observe initial state* $s$
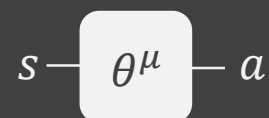
5. *Repeat:*

   a. ▇ · *Sample*

   b. ▇ · *Train*

"Critic" Network
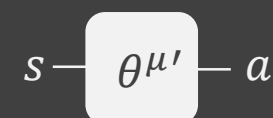$Q^\mu(s, a \mid \theta^Q)$

$$s,\ a \longrightarrow \boxed{\theta^Q} \longrightarrow Q(s, a)$$

Target "Critic" Network
$Q^{\mu\prime}(s, a \mid \theta^{Q\prime})$

$$s,\ a \longrightarrow \boxed{\theta^{Q\prime}} \longrightarrow Q(s, a)$$

"Actor" Network
$\mu(s \mid \theta^\mu)$

$$s \longrightarrow \boxed{\theta^\mu} \longrightarrow a$$

Target "Actor" Network
$\mu'(s \mid \theta^{\mu\prime})$

$$s \longrightarrow \boxed{\theta^{\mu\prime}} \longrightarrow a$$

Replay Buffer · $\mathcal{R}$

| $s_0$ | $s_1$ | $s_2$ | | $s_i$ | | $s$ |
|-------|-------|-------|---|-------|---|------|
| $a_0$ | $a_1$ | $a_2$ | $\cdots$ | $a_i$ | $\cdots$ $\leftarrow$ | $a$ |
| $r_0$ | $r_1$ | $r_2$ | | $r_i$ | | $r$ |
| $s_1$ | $s_2$ | $s_3$ | | $s_{i+1}$ | | $s'$ |

# Algorithms Evolution

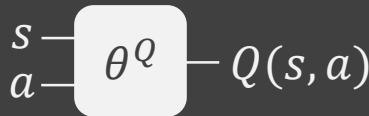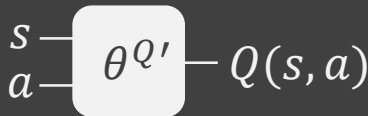DQN · DPG · DDPG                                   *Deep DPG*

🟨 · *Sample:*

i.   Execute $a = \mu(s \mid \theta^\mu) + \mathcal{N}$ and observe $r$ and $s'$

ii.  Store $\mathcal{R} \leftarrow \mathcal{R} \cup (s, a, r, s')$

"Critic" Network
$Q^\mu(s, a \mid \theta^Q)$

$$s \atop a \;-\; \boxed{\theta^Q} \;-\; Q(s, a)$$

Target "Critic" Network
$Q^{\mu'}(s, a \mid \theta^{Q'})$

$$s \atop a \;-\; \boxed{\theta^{Q'}} \;-\; Q(s, a)$$

"Actor" Network
$\mu(s \mid \theta^\mu)$

$$s \;-\; \boxed{\theta^\mu} \;-\; a$$

Target "Actor" Network
$\mu'(s \mid \theta^{\mu'})$

$$s \;-\; \boxed{\theta^{\mu'}} \;-\; a$$

Replay Buffer · $\mathcal{R}$

| $s_0$ | $s_1$ | $s_2$ | | $s_i$ | | | $s$ |
|-------|-------|-------|-----|-------|-----|---|-----|
| $a_0$ | $a_1$ | $a_2$ | $\cdots$ | $a_i$ | $\cdots$ | ← | $a$ |
| $r_0$ | $r_1$ | $r_2$ | | $r_i$ | | | $r$ |
| $s_1$ | $s_2$ | $s_3$ | | $s_{i+1}$ | | | $s'$ |

# Algorithms Evolution

DQN · DPG · DDPG                     *Deep DPG*

■ · *Train*

i.   *Sample random batch* $\mathcal{B} \subseteq \mathcal{R}$

ii.  $Q^T(s,a,r,s') := r + \gamma Q^{\mu'}\left(s', \overbrace{\mu'\left(s' \mid \theta^{\mu'}\right)}^{a'} \mid \theta^{Q'}\right)$

iii. $L^Q(\theta^Q) = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} \left[Q^T(s,a,r,s') - Q^{\mu}\left(s,a \mid \theta^Q\right)\right]^2$

iv.  $\theta^Q \leftarrow \theta^Q - \alpha \nabla_{\theta^Q} L(\theta^Q)$

v.   $L^{\mu}(\theta^{\mu}) = \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} \nabla_a Q^{\mu}\left(s,a \mid \theta^Q\right) \nabla_{\theta^{\mu}} \mu\left(s \mid \theta^{\mu}\right)$

vi.  $\theta^{\mu} \leftarrow \theta^{\mu} - \alpha \nabla_{\theta^{\mu}} L(\theta^{\mu})$
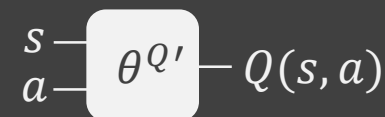
vii. $\theta^{Q'} \leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'}$ *where* $\tau \ll 1$

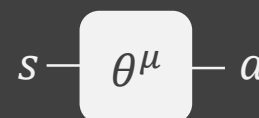viii. $\theta^{\mu'} \leftarrow \tau\theta^{\mu} + (1-\tau)\theta^{\mu'}$

"Critic" Network
$Q^{\mu}(s,a \mid \theta^Q)$

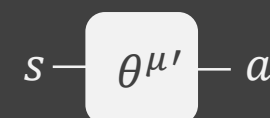$s$, $a$ → $\theta^Q$ → $Q(s,a)$

Target "Critic" Network
$Q^{\mu'}(s,a \mid \theta^{Q'})$

$s$, $a$ → $\theta^{Q'}$ → $Q(s,a)$

"Actor" Network
$\mu(s \mid \theta^{\mu})$

$s$ → $\theta^{\mu}$ → $a$

Target "Actor" Network
$\mu'(s \mid \theta^{\mu'})$

$s$ → $\theta^{\mu'}$ → $a$

Replay Buffer · $\mathcal{R}$

| $s_0$ | $s_1$ | $s_2$ | ⋯ | $s_i$ | ⋯ | $s$ |
|-------|-------|-------|---|-------|---|-----|
| $a_0$ | $a_1$ | $a_2$ | | $a_i$ | | $a$ |
| $r_0$ | $r_1$ | $r_2$ | | $r_i$ | | $r$ |
| $s_1$ | $s_2$ | $s_3$ | | $s_{i+1}$ | | $s'$ |

# DDPG Performance



Cart · Pendulum Swing-up · Cartpole Swing-up · Fixed Reacher · Monoped Balancing · Gripper · Blockworld · Puck Shooting · Cheetah · Moving Gripper

Normalized Reward
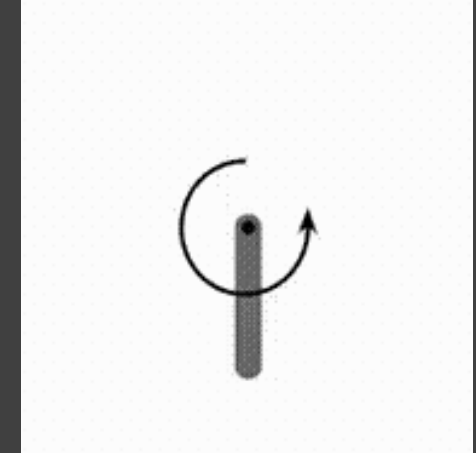
Million Steps

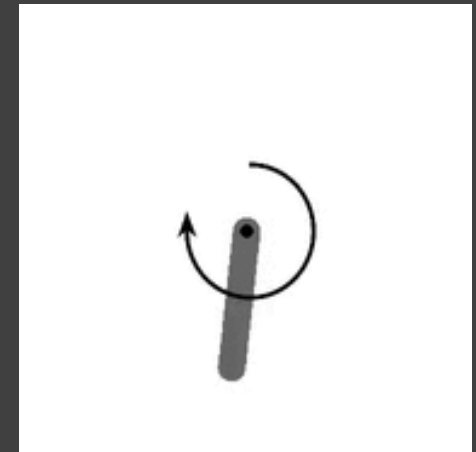· DPG   · DDPG   · DDPG + Batch Norm.   · DDPG + Batch Norm. + Raw pixels

# Conclusion

+ *Continuous control is much more stable.*

+ *Same algorithm/tuning solves many (20+) problems.*

+ *Pseudo-Code and hyperparameters are provided.*

- *Computationally expensive.*

- *Many hyperparameters must be tuned correctly.*

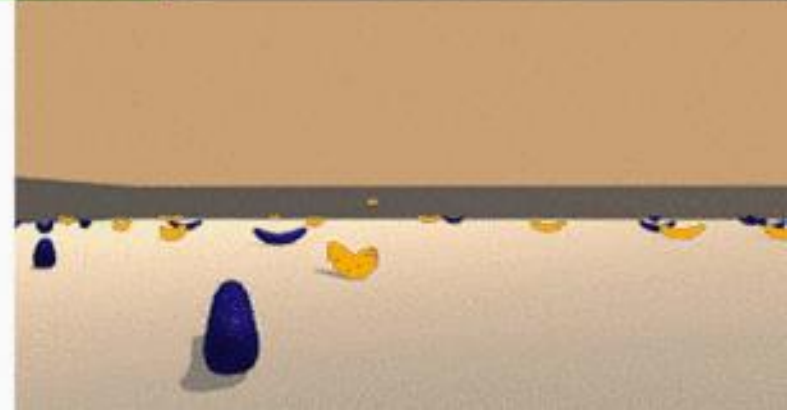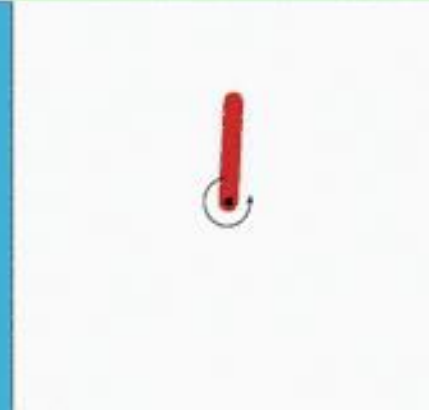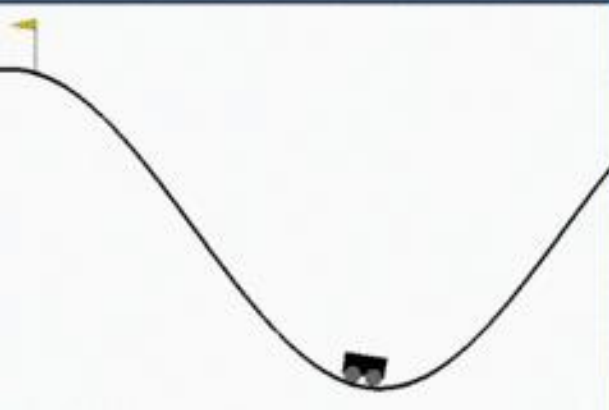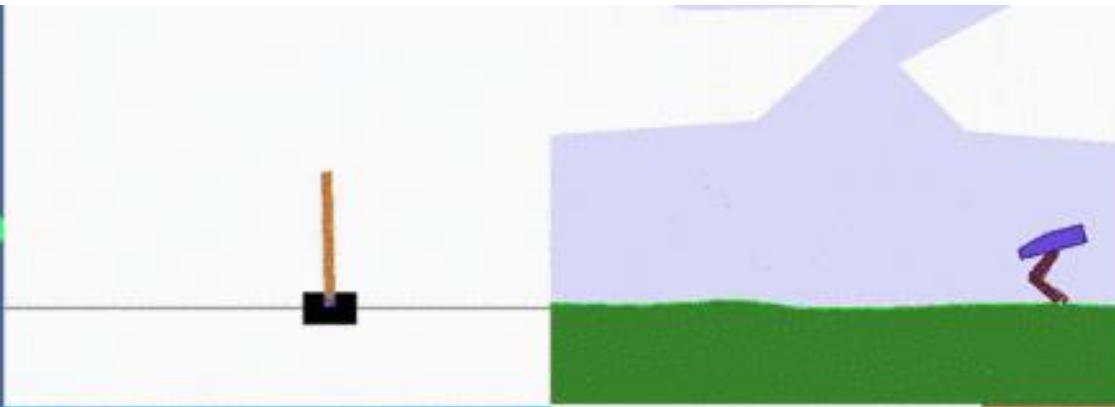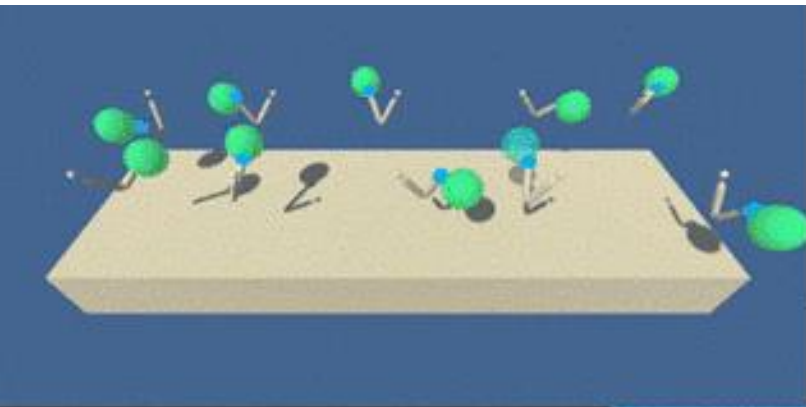- *Many decisions of the paper require lot of pre-knowledge.*

DQN · Discrete

DDPG · Continuous

# Any Questions?

# References

Presented Paper:
Lillicrap, Timothy P., et al. "Continuous control with deep reinforcement learning." arXiv preprint arXiv:1509.02971 (2015).

DQN Paper:
Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." *arXiv preprint arXiv:1312.5602* (2013).

DPG Paper:
Silver, David, et al. "Deterministic policy gradient algorithms." 2014.

Useful Websites:
- https://spinningup.openai.com/en/latest/spinningup/rl_intro.html
- https://towardsdatascience.com/deep-deterministic-policy-gradients-explained-2d94655a9b7b
- https://medium.com/@jonathan_hui/rl-dqn-deep-q-network-e207751f7ae4

# Hyperparameters

Learning Rates

$$\alpha_{actor} = 10^{-4}$$

$$\alpha_{critic} = 10^{-3}$$

Discount Factor

$$\gamma = 0.99$$

Target update

$$\tau = 0.001$$

Neural network

ReLU for all except last layer of the actor that uses tanh, 2 hidden layers of 300 respectively 400 neurons.

Buffer, Batch Size

$$10^6, 64$$

# Improvements after 2015

- *Paper uses batch normalization which has been found to improve only certain problems.*

- *An architecture with several actors can explore more of the environment.*