

# **CONTINUOUS CONTROL WITH DEEP REINFORCEMENT LEARNING (DDPG)**

SEONVIN CHO

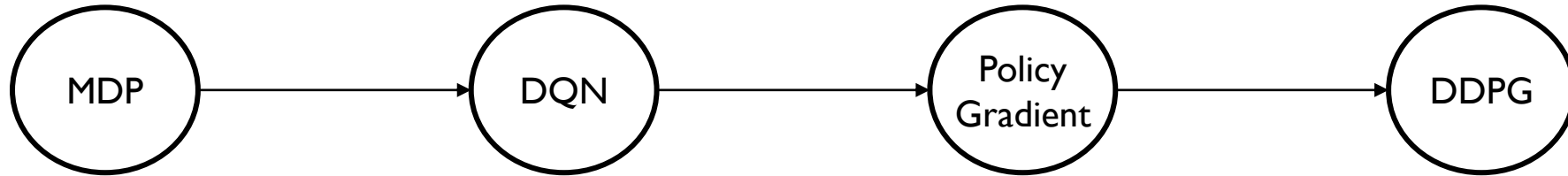
**INFORMATION AND INTELLIGENCE SYSTEMS LAB.**

ELECTRONIC ENGINEERING, HANYANG UNIVERSITY

**March 11, 2025**

# INTRO

---



- **Dynamic programming**

- instead of directly solving MDP is for the ease to solve difficult problems.
- splits and solves problems only when we know dynamics / transition probabilities.

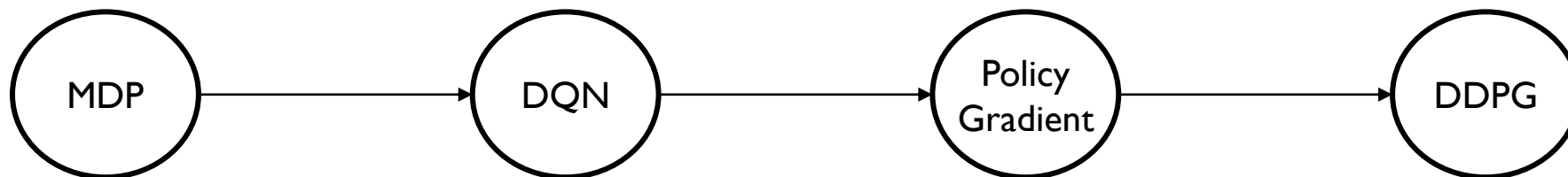
- **Monte Carlo**

- deals with all samples → unbiased sample & high variance.

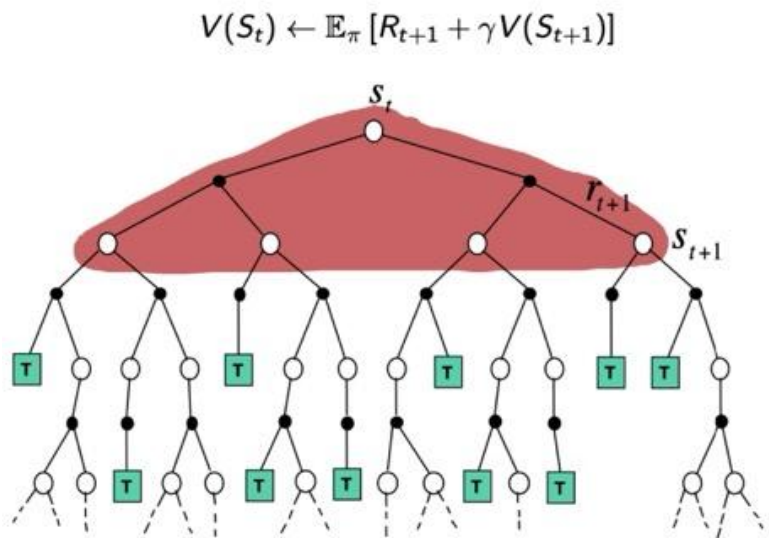
- **TD learning**

- uses only comparison with next one → biased sample & low variance.

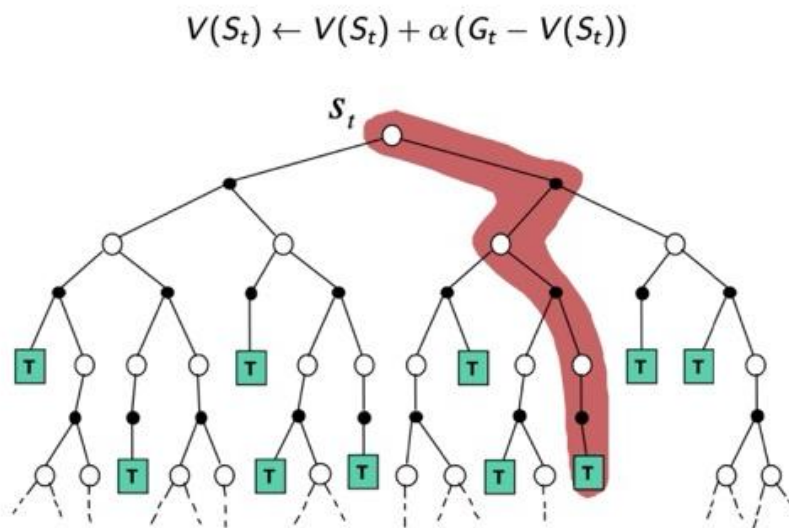
# INTRO



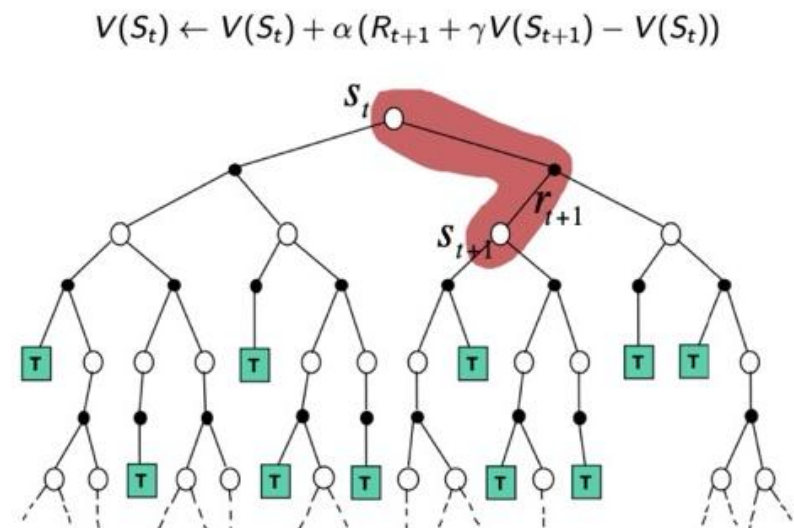
- **Dynamic Programming**



- **Monte Carlo**

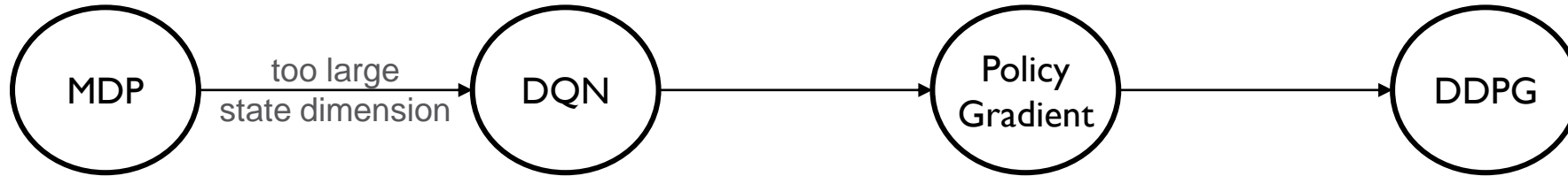


- **TD learning**



# INTRO

---

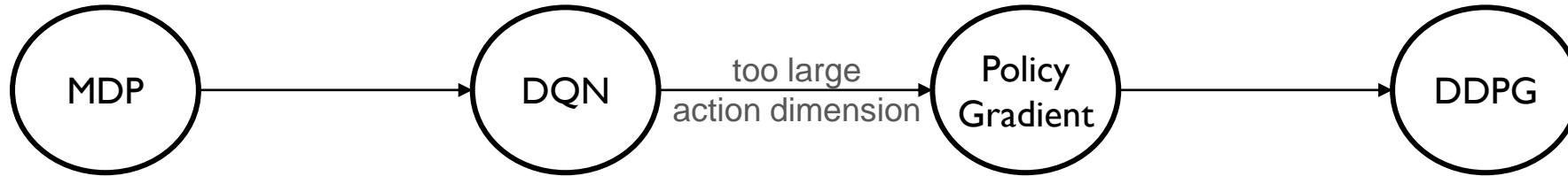


- **Deep Q-Network = Q-learning + Deep Neural Network**

- Estimates the action-value function = Value-based method
- Handles cases where the state dimension is large
  - › while dealing with discrete actions in continuous state spaces.
- Uses separate networks to mitigate value fluctuation.
- Still solves the "argmax Q" problem when updating weights
  - › Difficult to search for the optimal value if the function is nonlinear.

# INTRO

---

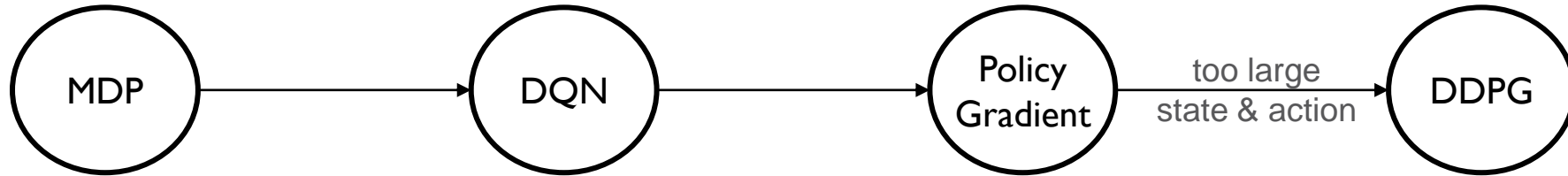


- **Policy Gradient**

- Estimates the parameterized policy = Policy-based method
- Handles cases where the action dimension is large
  - › while dealing with discrete states in continuous action spaces.
- Uses gradients to directly optimize a parameterized policy
  - › without considering whether the action space is continuous or discrete.
- Computes returns and gradient samples after all episodes have ended.
  - › unbiased and high variance

# INTRO

---



- **Deep Deterministic Policy Gradient**

**= DQN + Deterministic + Policy Gradient**

- Handles cases where both the state and action dimensions are large
  - › while dealing with continuous state and action spaces.
- Uses an Actor-Critic approach
  - › Employs a deterministic policy for continuous action spaces.
- Off-policy learning
  - › Utilize large replay buffers to leverage learning from uncorrelated transitions.

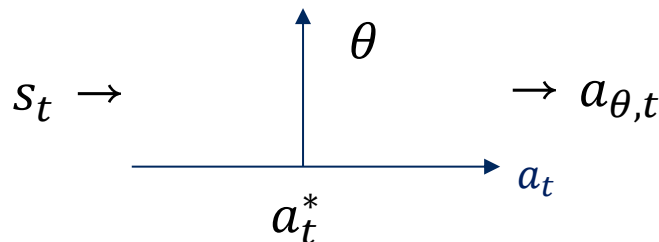
# DETERMINISTIC POLICY GRADIENT

- **Deterministic Policy**

- parameterized action

- $p(a_t|s_t) = \delta(a_t - a_t^*)$

- › delta function for only one optimal action

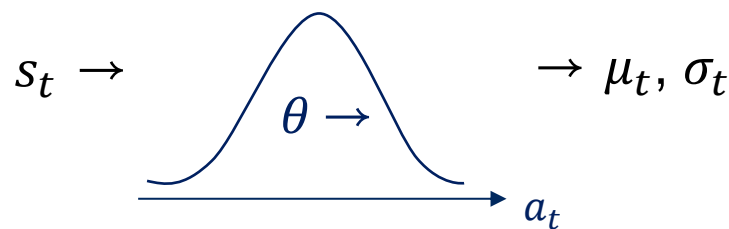


- **Stochastic Policy**

- parameterized policy

- $p_{\theta}(a_t|s_t)$

- › probability distribution for the action



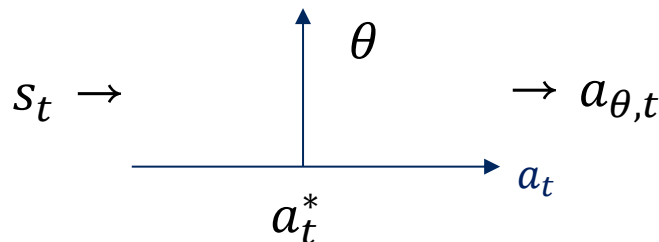
# DETERMINISTIC POLICY GRADIENT

- **Deterministic Policy**

- parameterized action

- $p(a_t|s_t) = \delta(a_t - a_t^*)$

- › delta function for only one optimal action



+



- **Policy Gradient Theorem**

- $\nabla_{\theta} J(\theta) = \mathbb{E}_{s,a \sim \pi} [\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot Q^{\pi}(s, a)]$

- **Deterministic Policy Gradient (DPG)**

- $\nabla_{\theta} J(\theta) = \mathbb{E}_{s,a \sim \pi} [\nabla_{\theta} \mu_{\theta}(s) \cdot \nabla_a Q^{\mu}(s, a)|_{a=\mu_{\theta}(s)}]$



# DETERMINISTIC POLICY GRADIENT (PROOF)

$$J_\theta = \mathbb{E}\{G_0\} = \int_{s_0: a_{0:00}} G_0 p(s_0: a_{0:00}) ds_0: a_{0:00} = \int_{s_0} \int_{a_0: a_{0:00}} G_0 P(a_0: a_{0:00} | s_0) da_0: a_{0:00} P(s_0) ds_0 = \int_{s_0} V(s_0) P(s_0) ds_0$$

$$V(s_0) = \int_{a_0: a_{0:00}} G_0 P(a_0: a_{0:00} | s_0) da_0: a_{0:00} = \int_{a_0} \int_{s_1: a_{0:00}} G_0 P(s_1: a_{0:00} | s_0, a_0) P(a_0 | s_0) ds_1: a_{0:00} da_0$$

$$= \int_{a_0} Q(s_0, a_0) P(a_0 | s_0) da_0 = Q(s_0, a_{\theta,0})$$

$= \int (a_0 - a_{\theta,0})$

$$Q(s_0, a_{\theta,0}) = \int_{s_1: a_{0:00}} (R_0 + \gamma G_1) p(s_1: a_{0:00}) ds_1: a_{0:00} = \int_{s_1, a_{0:1}} \int_{s_2: a_{0:00}} (R_0 + \gamma G_1) P(s_2: a_{0:00} | s_1, a_{0:1}) \underbrace{P(a_{0:1} | s_1)}_{= \int (a_{0:1} - a_{\theta,1})} P(s_1 | s_0, a_{\theta,0}) ds_2: a_{0:00} ds_1, a_{0:1}$$

$$= \int_{s_1, a_{0:1}} (R_0 + \gamma Q(s_1, a_{\theta,1})) P(a_{0:1} | s_1) P(s_1 | s_0, a_{\theta,0}) ds_1, a_{0:1}$$

$s_0, a_{\theta,0}$ 의 행동으로 정의  
action 주니 reward 받음 → next state  $s_0 \rightarrow a_0 \rightarrow r_0 \rightarrow s_1 \rightarrow \dots$

$$= R_0 + \gamma \int_{s_1} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1$$

$$\nabla_\theta J_\theta = \int_{s_0} \nabla_\theta V(s_0) P(s_0) ds_0 = \int_{s_0} \nabla_\theta Q(s_0, a_{\theta,0}) P(s_0) ds_0$$

$$\begin{aligned} \nabla_\theta Q(s_0, a_{\theta,0}) &= \nabla_{a_{\theta,0}} R_0 \cdot a'_{\theta,0} + \int_{s_1} \gamma Q(s_1, a_{\theta,1}) \nabla_\theta P(s_1 | s_0, a_{\theta,0}) a'_{\theta,0} ds_1 + \int_{s_1} \gamma \nabla_\theta Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 \\ &= a'_{\theta,0} \cdot \nabla_{a_{\theta,0}} (R_0 + \int_{s_1} \gamma Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1) + \int_{s_1} \gamma \nabla_\theta Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 \\ &= a'_{\theta,0} \cdot \nabla_{a_{\theta,0}} Q(s_0, a_{\theta,0}) + \int_{s_1} \gamma \nabla_\theta Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 \end{aligned}$$

$$\nabla_\theta J_\theta = \int_{s_0} a'_{\theta,0} \cdot \nabla_{a_{\theta,0}} Q(s_0, a_{\theta,0}) P(s_0) ds_0 + \int_{s_0} \int_{s_1} \gamma a_{\theta,1} \nabla_{a_{\theta,1}} Q(s_1, a_{\theta,1}) P(s_1 | s_0, a_{\theta,0}) ds_1 P(s_0) ds_0 + \dots$$

$$\approx \sum_{t=0}^{\infty} \int_{s_0: s_t} \underbrace{\delta^t a_{\theta,t} \nabla_{a_{\theta,t}} Q(s_t, a_{\theta,t})}_{\substack{t \uparrow \rightarrow \delta^t \approx 0 \rightarrow r^t \text{ 무시}}} P(s_0) \prod_{k=1}^t P(s_k | s_{k-1}, a_{\theta,k-1}) ds_0: s_t$$

replay buffer → sample → update      off-policy ∴ correlation ↓      actor-critic

$$= \sum_{t=0}^{\infty} \int_{s_0: s_t} \nabla_\theta Q_w(s_t, a_{\theta,t}) P(s_0) \prod_{k=1}^t P(s_k | s_{k-1}, a_{\theta,k-1}) ds_0: s_t$$

# DDPG

- **DQN + Deterministic Policy Gradient**

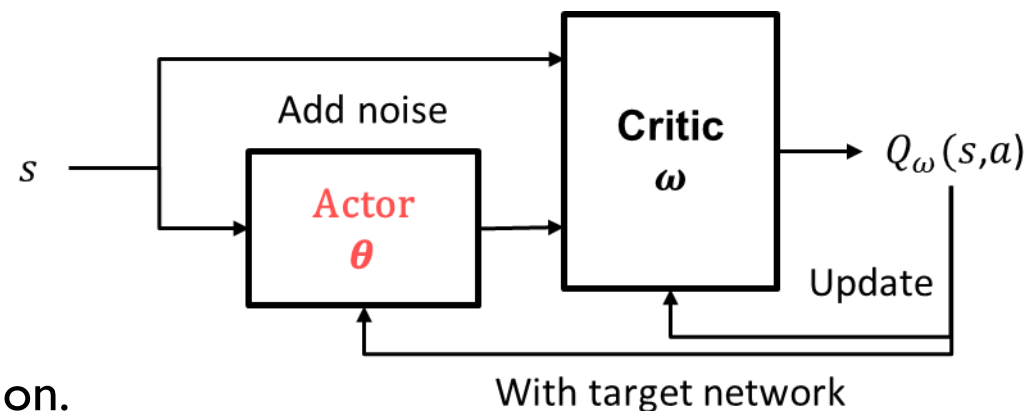
- parameterized action

- $p(a_t|s_t) = \delta(a_t - a_t^*)$

- › Suitable for continuous action spaces.

- › Limited exploration

- Requires Gaussian noise for better exploration.



- Actor update based on the Q-value gradient.

- Critic update using TD error.

# DDPG

---

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .

Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$

Initialize replay buffer  $R$

**for** episode = 1, M **do**

    Initialize a random process  $\mathcal{N}$  for action exploration

    Receive initial observation state  $s_1$

**for** t = 1, T **do** noise for exploration

        Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise

        Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$

        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$

        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$

        Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$  uniform random sample for replay buffer

        Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$

        Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$
uniform random sample for replay buffer

    Update the target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

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

**end for**

**end for**

---

# DDPG

---

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network  $Q(s, a|\theta^Q)$  and actor  $\mu(s|\theta^\mu)$  with weights  $\theta^Q$  and  $\theta^\mu$ .

Initialize target network  $Q'$  and  $\mu'$  with weights  $\theta^{Q'} \leftarrow \theta^Q, \theta^{\mu'} \leftarrow \theta^\mu$

Initialize replay buffer  $R$

**for** episode = 1, M **do**

    Initialize a random process  $\mathcal{N}$  for action exploration

    Receive initial observation state  $s_1$

**for** t = 1, T **do**

        Select action  $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$  according to the current policy and exploration noise

        Execute action  $a_t$  and observe reward  $r_t$  and observe new state  $s_{t+1}$

        Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$

        Sample a random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$

        Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$  target network

        Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$  main network

        Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

    Update the target networks:

$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

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

**end for**  
**end for**

---

separate networks  
similar to DQN!

# DDPG

- **DQN +Deterministic Policy Gradient**

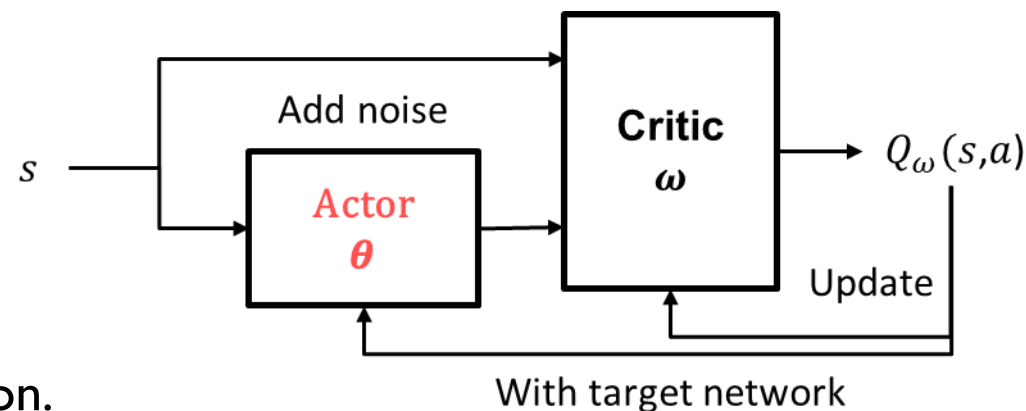
- parameterized action

- $p(a_t|s_t) = \delta(a_t - a_t^*)$

- › Suitable for continuous action spaces.

- › Limited exploration

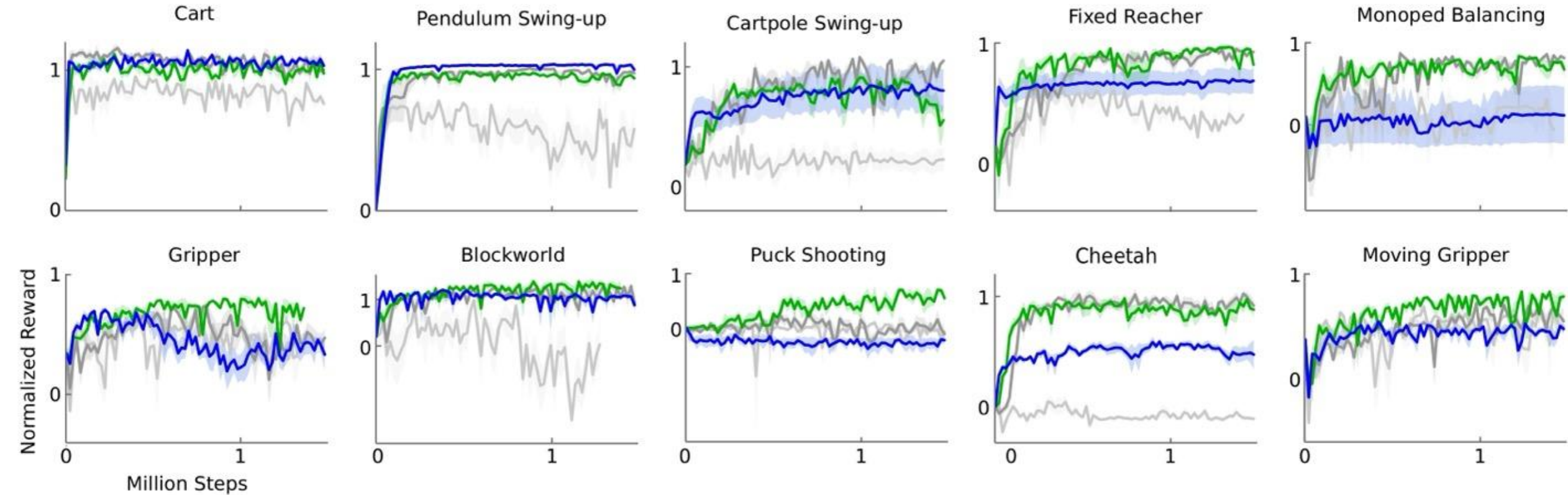
- Requires gaussian noise for better exploration.



- Actor update based on the Q-value gradient.

- Critic update using TD error.

# RESULTS



light grey = original DPG algorithm (minibatch NFQCA) with batch normalization  
dark grey = original DPG algorithm with target network  
green = original DPG algorithm with target networks and batch normalization  
blue = original DPG algorithm with target networks from pixel-only inputs

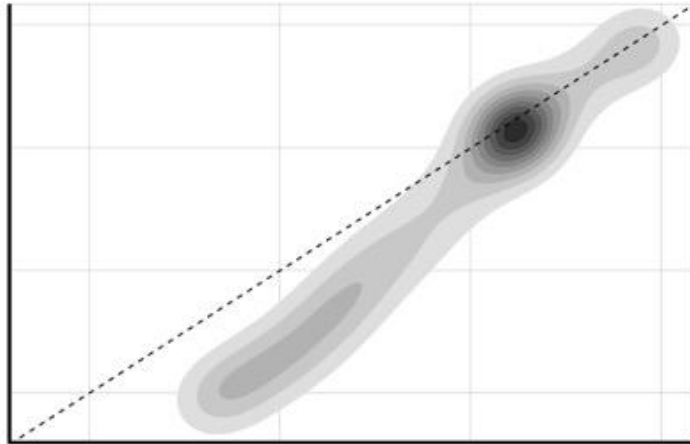
Target networks are crucial!

# RESULTS

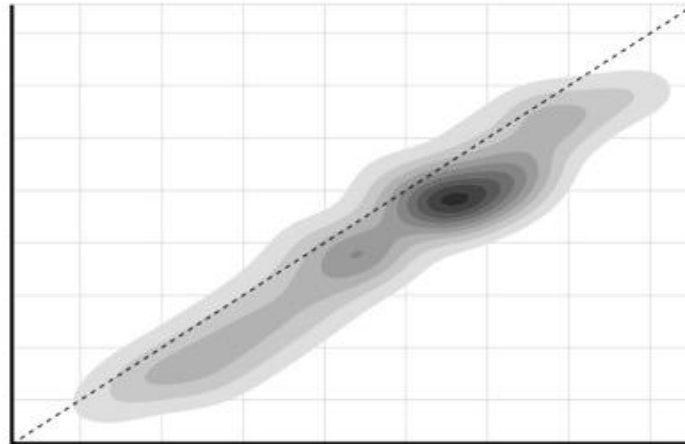
Q-value estimates deviate more from actual returns.  
Value estimation is more challenging in complex environments.

Pendulum

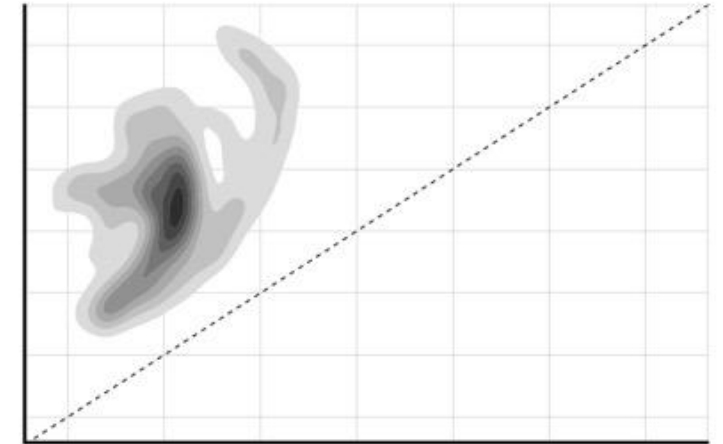
Estimated Q



Cartpole



Cheetah



Q-value estimates closely match actual returns.  
Accurate value estimation in simple environments.

Return

Dashed line: Ideal case where Q-values match actual returns  
Darker regions: More frequent Q-value and return pairs

# SUMMARY

---

## **DDPG = DQN + Deterministic Policy + Policy Gradient Theorem**

- DQN : off-policy & separate networks
  - › Off-policy : target policy  $\neq$  behavior policy
    - Uses a replay buffer to maximize sample efficiency.
  - › Separate networks
    - Uses target network and main network for stationarity.
- Actor-Critic approach
  - › Actor updates the policy network using the policy gradient.
  - › Critic updates the Q-value function using TD targets with a target network.
- Gaussian Noise for exploration
- Soft target update ( $\tau$ ) to slowly update the target network