

# 535514: Reinforcement Learning

## Lecture 22 — DQN, DDQN, and Distributional RL

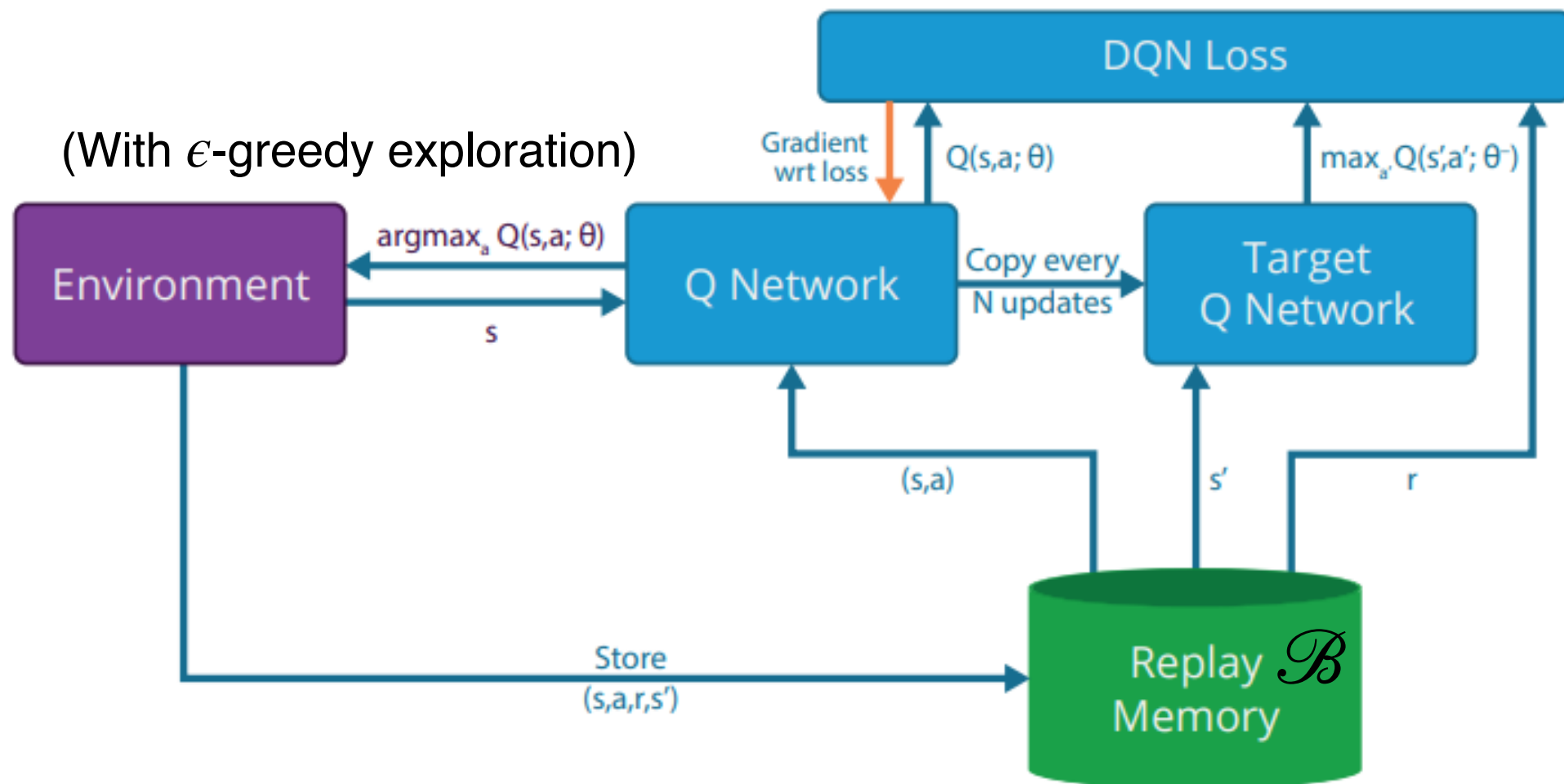
Ping-Chun Hsieh

May 9, 2024

# On-Policy vs Off-Policy Methods

	Policy Optimization	Value-Based	Model-Based	Imitation-Based
On-Policy	Exact PG REINFORCE (w/i baseline) A2C On-policy DAC TRPO Natural PG (NPG) PPO-KL & PPO-Clip RLHF by PPO-KL	Epsilon-Greedy MC Sarsa Expected Sarsa	Model-Predictive Control (MPC) PETS	IRL GAIL IQ-Learn
Off-Policy	Off-policy DPG & DDPG Twin Delayed DDPG (TD3)	Q-learning Double Q-learning DQN & DDQN C51 / QR-DQN / IQN Rainbow Soft Actor-Critic (SAC)		

# Review: Deep Q-Network



$$F(\mathbf{w}) := \frac{1}{2} \mathbb{E}_{(s,a,r,s') \sim \rho} \left[ \left( r + \gamma \max_{a' \in A} \underline{Q(s', a'; \bar{\mathbf{w}})} - Q(s, a; \mathbf{w}) \right)^2 \right]$$

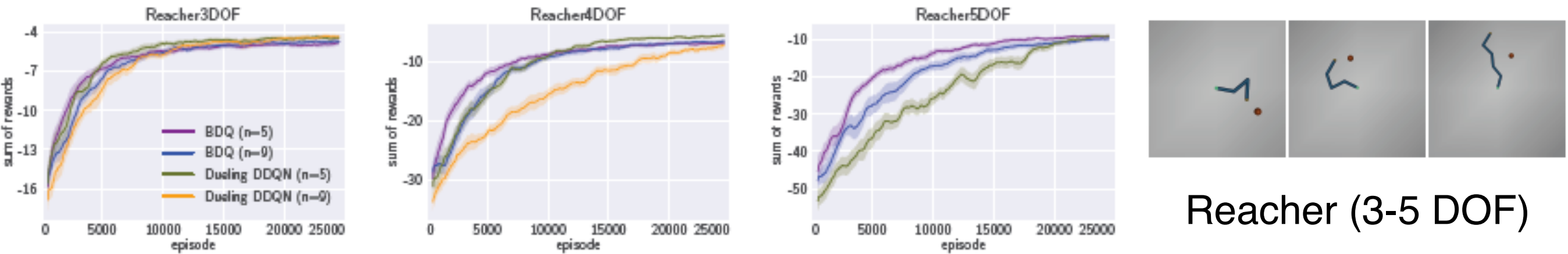
$$\approx \frac{1}{2} \sum_{\underline{(s,a,r,s') \in D \sim \mathcal{B}}} \left[ \left( r + \gamma \max_{a' \in A} Q(s', a'; \bar{\mathbf{w}}) - Q(s, a; \mathbf{w}) \right)^2 \right]$$

## Can DQN be Applied Under Continuous Actions?

The difficulty lies in the “ $\arg \max_{a \in \mathcal{A}} Q(s, a; \mathbf{w})$ ” operation

# Existing methods that adapts DQN to continuous actions

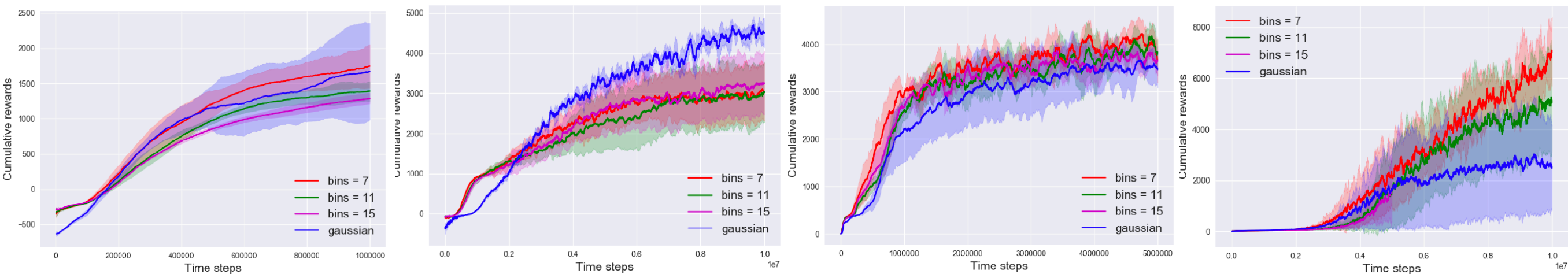
## 1. Naive Action Discretization [Tavakoli et al., AAAI 2020]



Issue: Naive discretization suffers from exponential growth of cardinality

## 2. Discretization + Factorization [Tang and Agrawal, AAAI 2020]

$$\pi(a|s) := \prod_{i=1}^d \pi_{\theta_i}(a_i|s) \quad \text{where } a = [a_0, a_1, \dots, a_{d-1}]^T$$



(a) **HalfCheetah + PPO**

(b) **Ant + PPO**

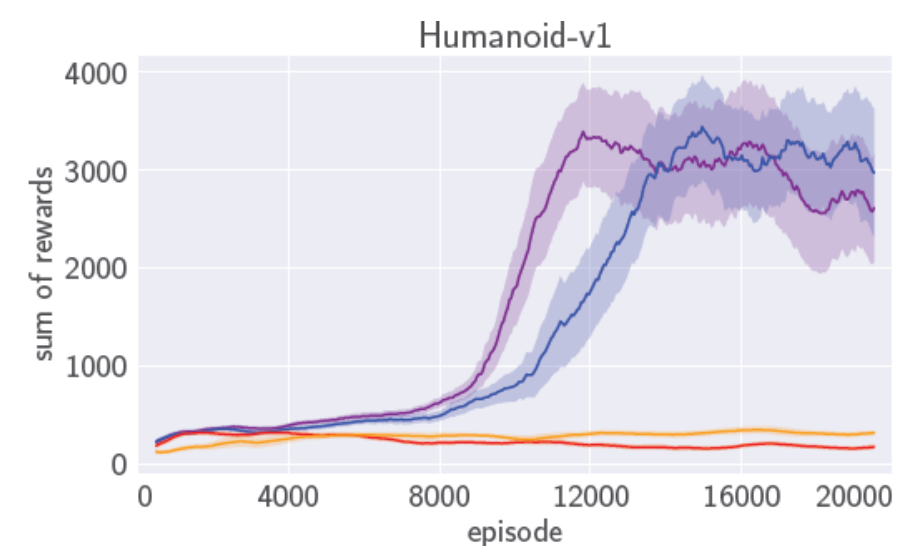
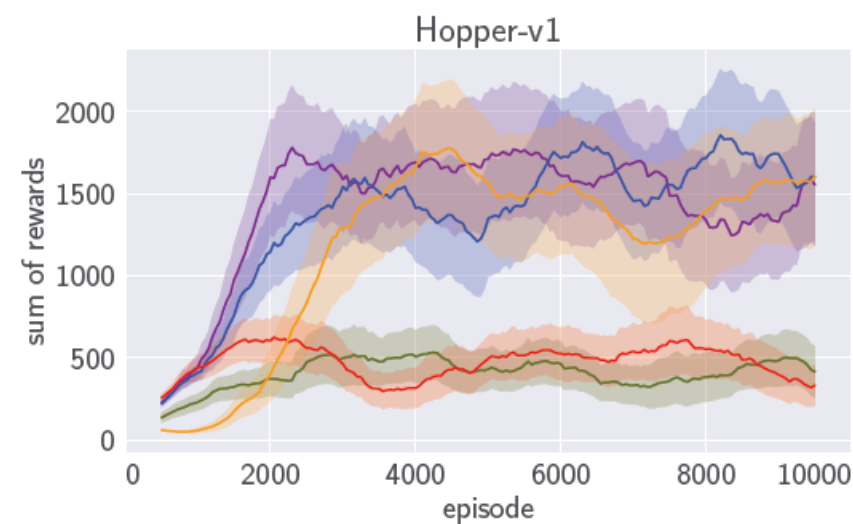
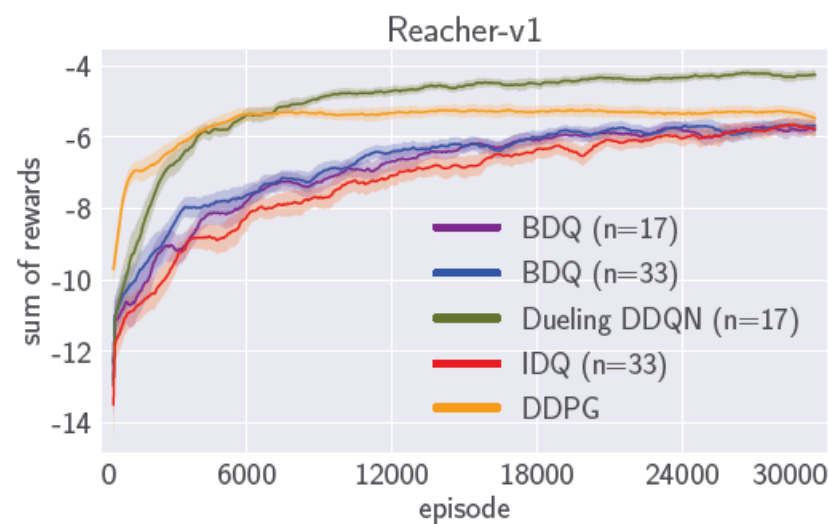
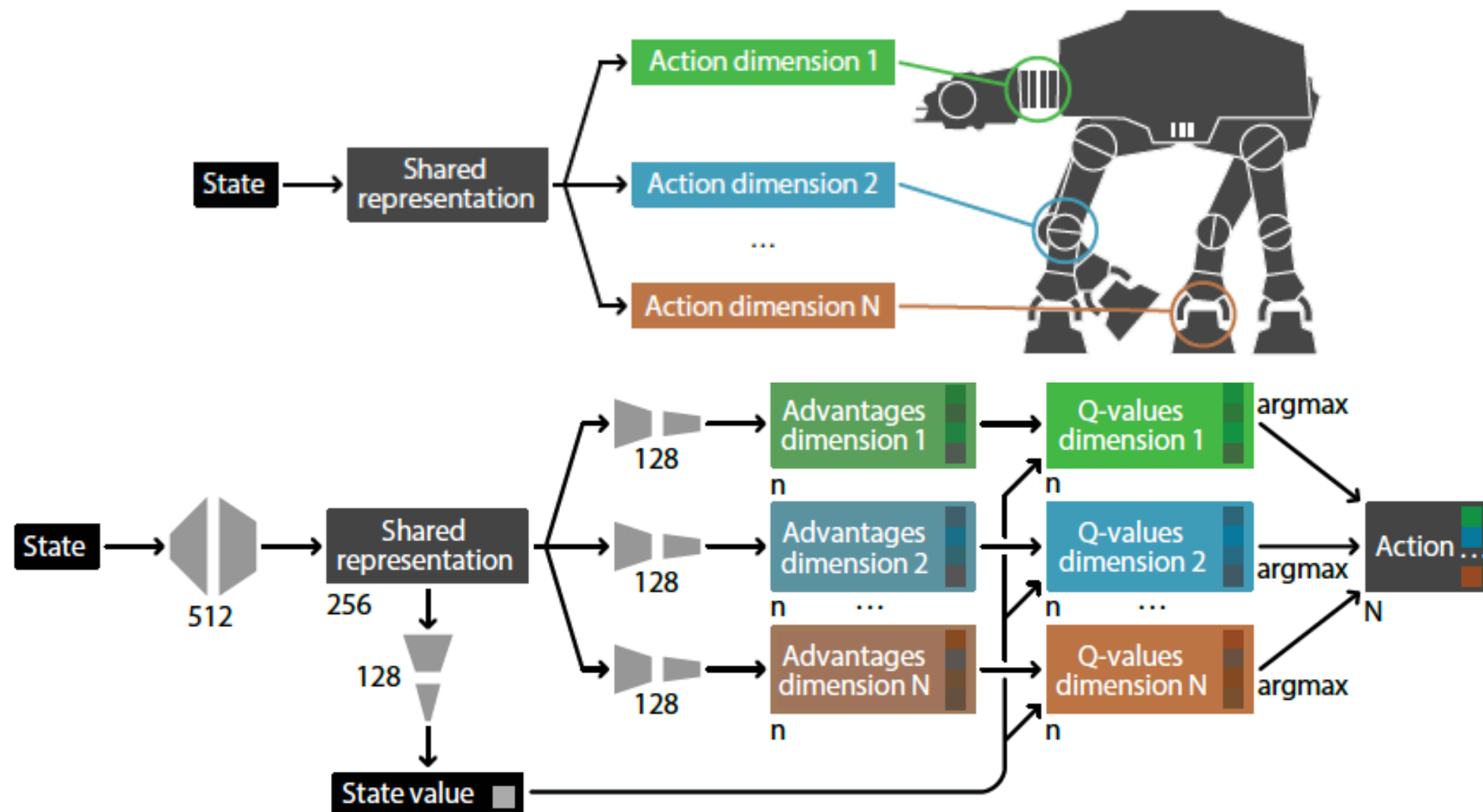
(c) **Walker + PPO**

(d) **Humanoid (R) + PPO**

Tavakoli et al., Action Branching Architectures for Deep Reinforcement Learning, AAAI 2018

Tang and Agrawal et al., Discretizing Continuous Action Space for On-Policy Optimization, AAAI 2020

### 3. Discretization + Branching [Tavakoli et al., AAAI 2018]



## 4. Normalized Advantage Functions (NAF): Quadratic Approximation!

Continuous Deep Q-Learning with Model-based Acceleration

[Gu et al., ICML 2016]

Shixiang Gu<sup>1 2 3</sup>

Timothy Lillicrap<sup>4</sup>

Ilya Sutskever<sup>3</sup>

Sergey Levine<sup>3</sup>

SG717@CAM.AC.UK

COUNTZERO@GOOGLE.COM

ILYASU@GOOGLE.COM

SLEVINE@GOOGLE.COM

<sup>1</sup>University of Cambridge <sup>2</sup>Max Planck Institute for Intelligent Systems <sup>3</sup>Google Brain <sup>4</sup>Google DeepMind

$$Q(s, a; \phi_A, \phi_V) = A(s, a; \phi_A) + V(s; \phi_V) \quad (P \text{ is state-dependent, positive definite})$$

$$A(s, a; \phi_A) := -\frac{1}{2}(a - \mu(s; \phi_\mu))^\top P(s; \phi_P)(a - \mu(s; \phi_\mu))$$

$$P(s; \phi_P) := L(s; \phi_P)L(s; \phi_P)^\top \quad (\text{This is known as the "Cholesky decomposition"})$$

( $L$  is a lower-triangular matrix)



## 5. Amortized Q-Learning (AQL): Sampling!

### Q-LEARNING IN ENORMOUS ACTION SPACES VIA AMORTIZED APPROXIMATE MAXIMIZATION

[NeurIPS 2018 Workshop]

Tom Van de Wiele\*, David Warde-Farley, Andriy Mnih & Volodymyr Mnih

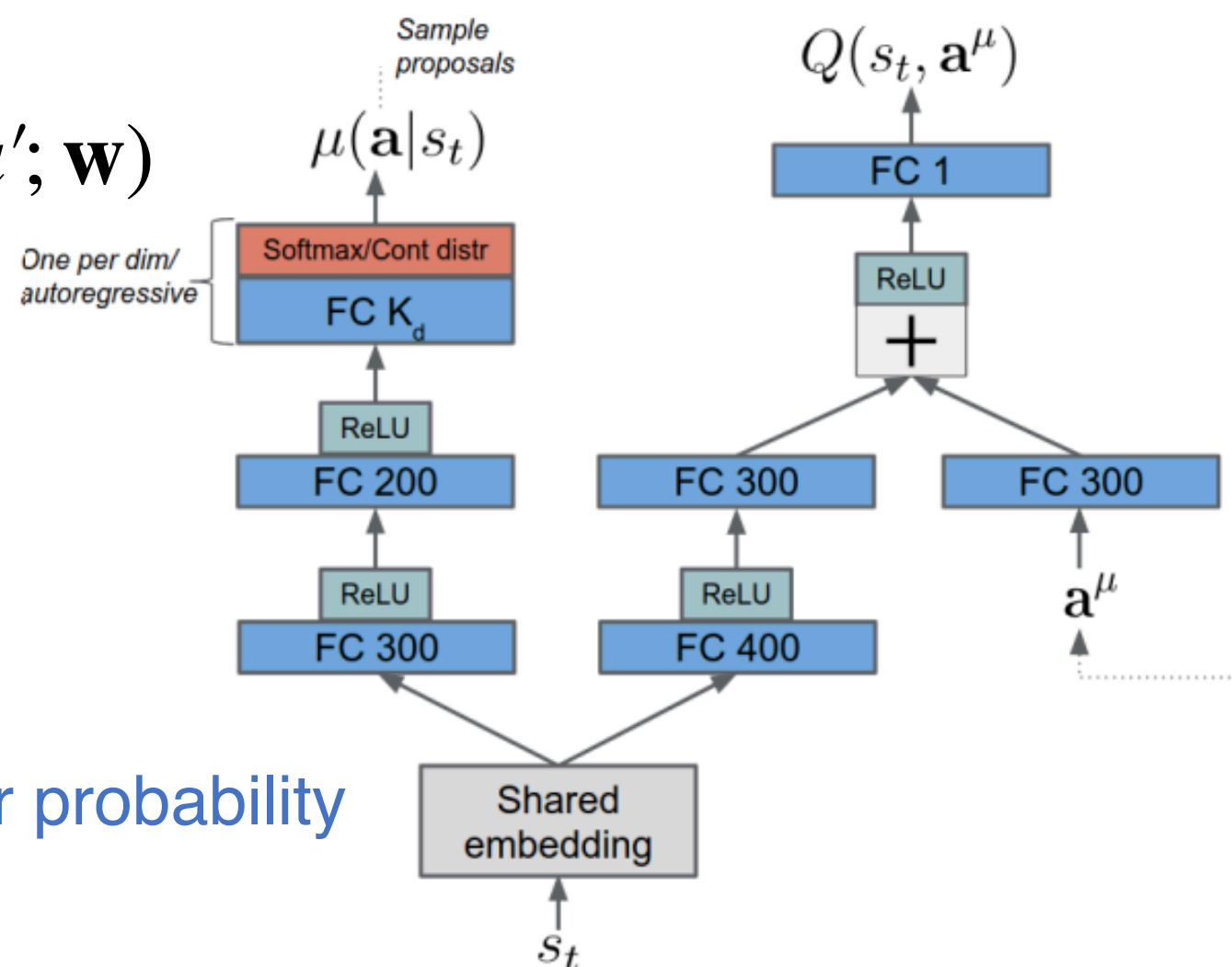
DeepMind

London, United Kingdom

tvdwiele@gmail.com, {dwf, amnih, vmnih}@google.com

$$\max_{a \in \mathcal{A}} Q(s, a; \mathbf{w}) \approx \max_{a' \in D \sim \mu(a)} Q(s, a'; \mathbf{w})$$

“proposal distribution”  
(Learned by maximum likelihood)



**Intuition:** Larger  $|D|$  induces a higher probability of seeing max-Q actions



## 6. DDPG: Reinterpret DDPG as an Adaptation of DQN for Continuous Actions!

(Quick Review)

Off-Policy Deterministic PG:  $\nabla_{\theta} J_{\beta}^{\pi_{\theta}} \approx \mathbb{E}_{s \sim d_{\mu}^{\beta}} \left[ \nabla_{\theta} \pi_{\theta}(s) \nabla_a Q^{\pi_{\theta}}(s, a) |_{a=\pi_{\theta}(s)} \right]$

- **Critic**: estimate  $Q_w \approx Q^{\pi_{\theta}}$  by bootstrapping
- **Actor**: updates policy parameters  $\theta$  by deterministic policy gradient

Step 1: Initialize  $\theta_0$ ,  $w_0$  and step sizes  $\alpha_{\theta}$ ,  $\alpha_w$

Step 2: Sample a trajectory  $\tau = (s_0, a_0, r_1, \dots) \sim P_{\mu}^{\beta}$

For each step of the current trajectory  $t = 0, 1, 2, \dots$

$$\Delta w_k \leftarrow \Delta w_k + \alpha_w (r_t + \gamma Q_{w_k}(s_{t+1}, \pi_{\theta}(s_{t+1})) - Q_{w_k}(s_t, a_t)) \nabla_w Q_w(s_t, a_t) |_{w=w_k}$$

$$\Delta \theta_k \leftarrow \Delta \theta_k + \alpha_{\theta} \gamma^t \left( \nabla_{\theta} \pi_{\theta}(s_t) \nabla_a Q_{w_k}(s_t, a) |_{a=\pi_{\theta}(s_t)} \right)$$

  $= \nabla_{\theta} Q_{w_k}(s_t, \pi_{\theta}(s_t)) |_{\theta=\theta_k}$

# Alternative Interpretation of DDPG: An Adaptation of DQN for Continuous Actions (Cont.)

- ▶ DDPG can be reinterpreted as DQN for continuous actions

1. **Deterministic policy:**  $\pi_{\theta}(s) \approx \arg \max_a Q_w(s, a)$

2. **How to find  $\theta$ :** solve  $\theta \leftarrow \arg \max_{\theta} Q_w(s, \pi_{\theta}(s))$  by SGD

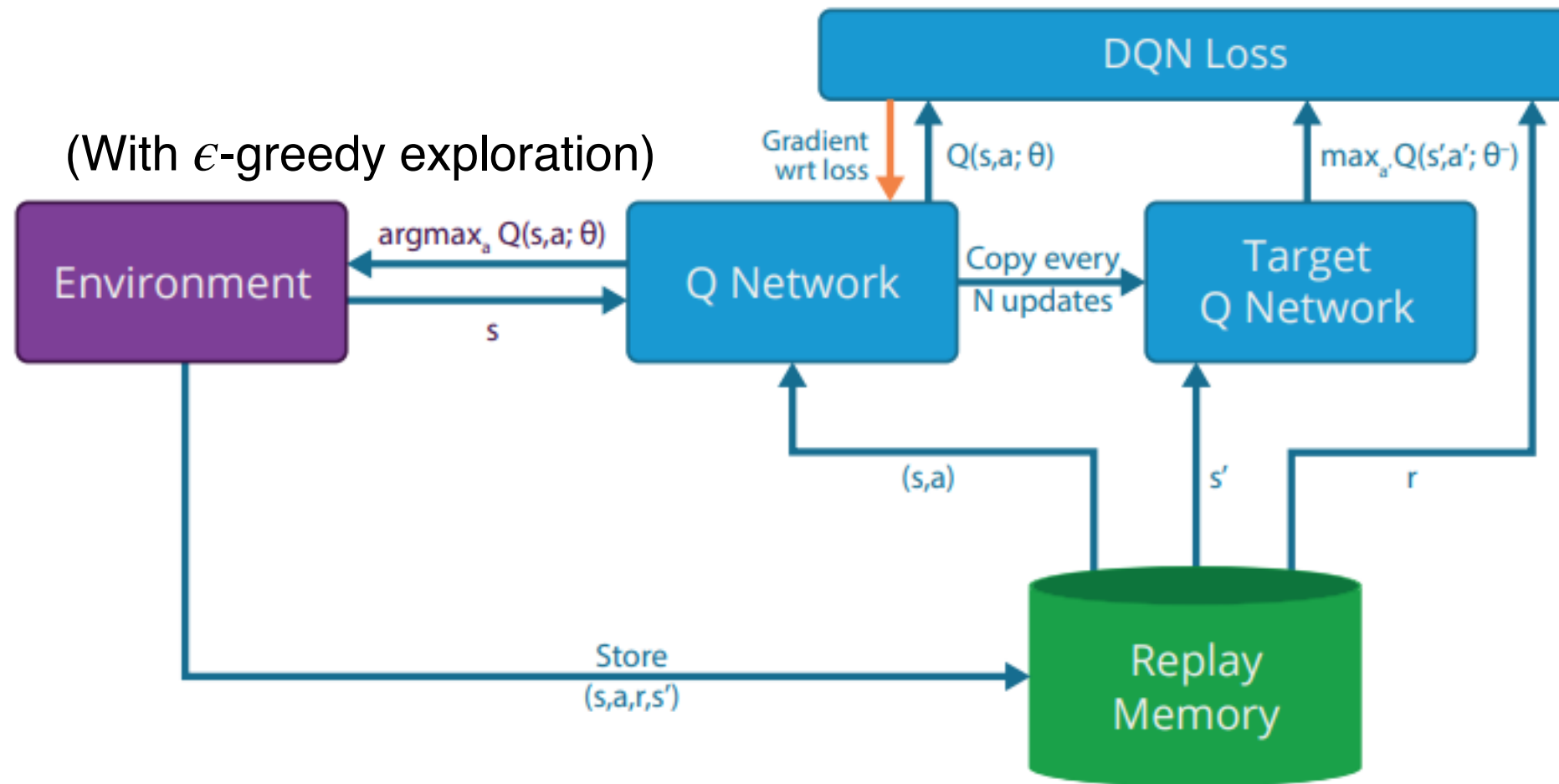
$$\Delta \theta_k \leftarrow \Delta \theta_k + \alpha_{\theta} \gamma^t \left( \nabla_{\theta} \pi_{\theta}(s_t) \nabla_a Q_{w_k}(s_t, a) \Big|_{a=\pi_{\theta}(s_t)} \right)$$

$\searrow = \nabla_{\theta} Q_{w_k}(s_t, \pi_{\theta}(s_t)) \Big|_{\theta=\theta_k}$

3. DQN and DDPG have a similar TD update scheme

$$\Delta w_k \leftarrow \Delta w_k + \alpha_w \left( r_t + \gamma Q_{w_k}(s_{t+1}, \pi_{\theta}(s_{t+1})) - Q_{w_k}(s_t, a_t) \right) \nabla_w Q_w(s_t, a_t) \Big|_{w=w_k}$$

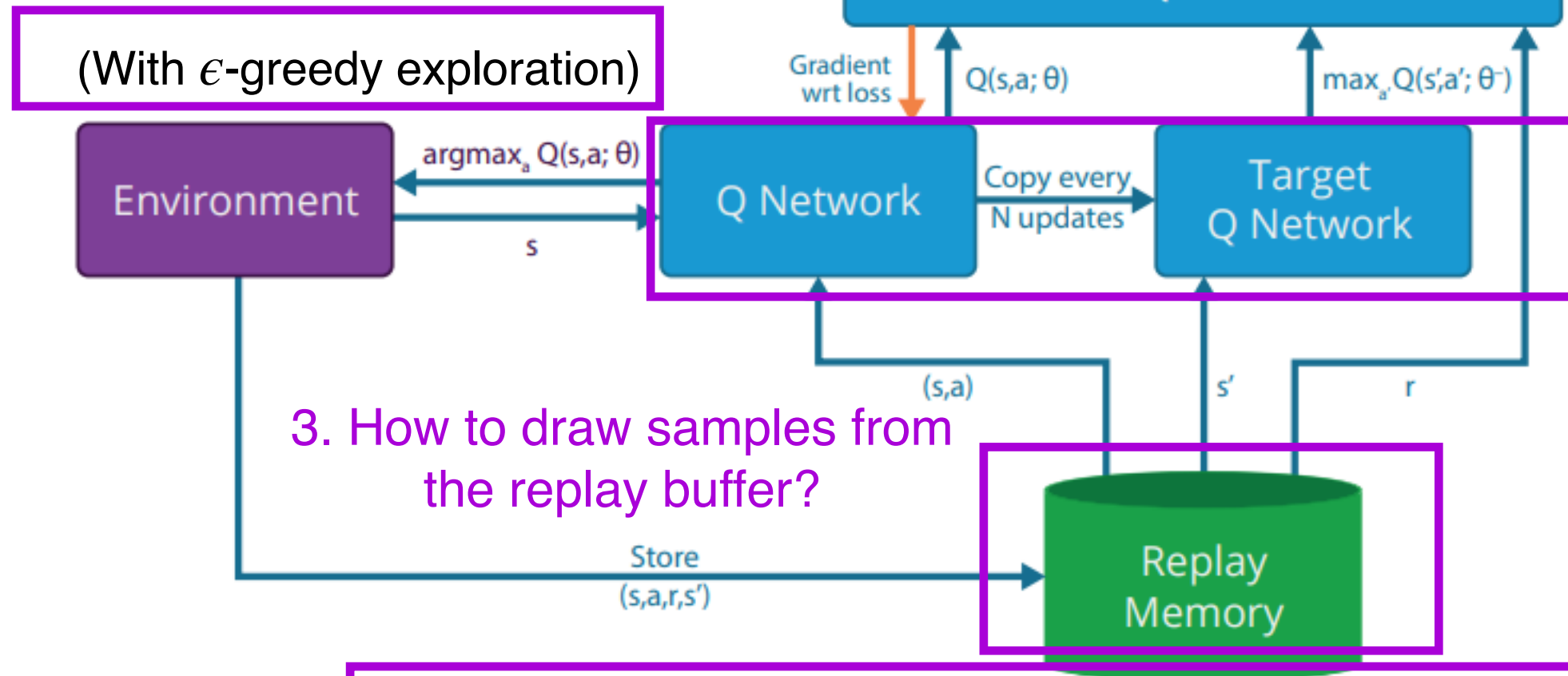
# Next Topic: What to Improve in Vanilla DQN?



$$L(\mathbf{w}) := \sum_{(s,a,r,s') \in D} \frac{1}{2} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \bar{\mathbf{w}}) - Q(s, a; \mathbf{w}) \right)^2 \right]$$

# Next Topic: What to Improve in Vanilla DQN?

4. A better exploration method?



3. How to draw samples from the replay buffer?

$$L(\mathbf{w}) := \sum_{(s,a,r,s') \in D} \frac{1}{2} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \bar{\mathbf{w}}) - Q(s, a; \mathbf{w}) \right)^2 \right]$$

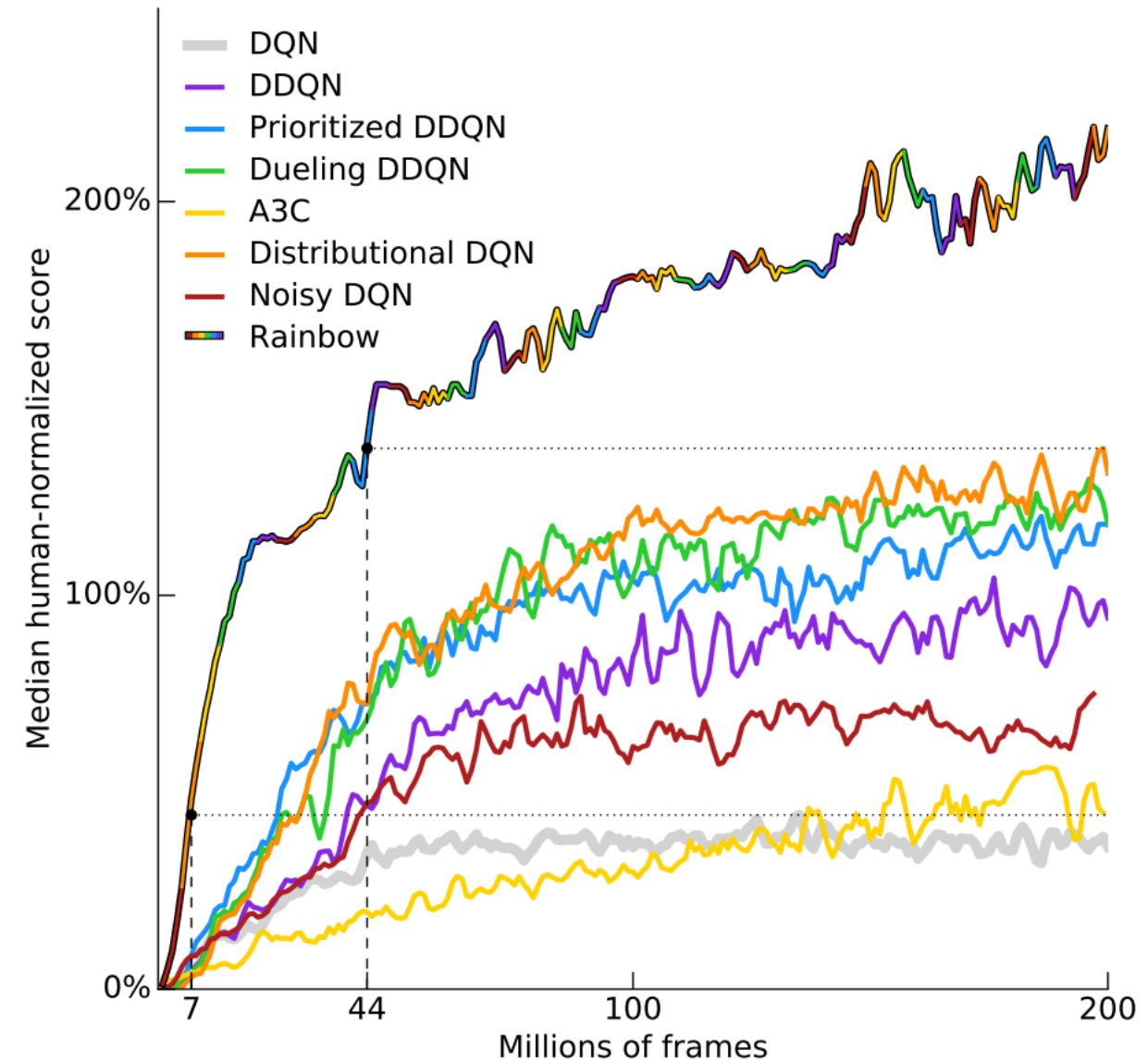
1. Overestimation Bias?

2. A better loss function?

5. A better way to represent Q function?

# Next Topic: Rainbow (= DQN with 6 Useful Tricks)

1. Double DQN (DDQN)
2. Distributional Q-learning
3. Prioritized experience replay (PER)
4. Dueling networks
5. Multi-step return in TD target
6. Noisy networks for exploration



# Double DQN (DDQN)

Hado van Hasselt, Arthur Guez, and David Silver,  
“Deep Reinforcement Learning with Double Q-learning,” AAAI 2016

# Recall: Double Q-Learning

Step 1: Initialize  $Q^A(s, a)$ ,  $Q^B(s, a)$  for all  $(s, a)$ , and initial state  $s_0$

Step 2: For each step  $t = 0, 1, 2, \dots$

Select  $a_t$  using  $\varepsilon$ -greedy w.r.t  $Q^A(s_t, a) + Q^B(s_t, a)$

Observe  $(r_{t+1}, s_{t+1})$

Choose one of the following updates uniformly at random

$$Q^A(s_t, a_t) \leftarrow Q^A(s_t, a_t) + \alpha(r_{t+1} + \gamma Q^B(s_{t+1}, \arg \max_a Q^A(s_{t+1}, a)) - Q^A(s_t, a_t))$$

$$Q^B(s_t, a_t) \leftarrow Q^B(s_t, a_t) + \alpha(r_{t+1} + \gamma Q^A(s_{t+1}, \arg \max_a Q^B(s_{t+1}, a)) - Q^B(s_t, a_t))$$

- ▶ **Key Idea:** Decouple “Q value” and “greedy action selection”
- ▶ **Question:** How to apply this to DQN?



# Double DQN

- Loss function of DQN:

$$F(\mathbf{w}) := \frac{1}{2} \mathbb{E}_{(s,a,r,s') \sim \rho} \left[ \left( r + \gamma \max_{a' \in A} Q(s', a'; \bar{\mathbf{w}}) - Q(s, a; \mathbf{w}) \right)^2 \right]$$
$$\approx \frac{1}{2} \sum_{(s,a,r,s') \in D} \left[ \left( r + \gamma \max_{a' \in A} Q(s', a'; \mathbf{w}) - Q(s, a; \mathbf{w}) \right)^2 \right]$$

- Loss function of Double DQN:

$$F(\mathbf{w}) := \frac{1}{2} \mathbb{E}_{(s,a,r,s') \sim \rho} \left[ \left( r + \gamma Q\left(s', \arg \max_{a' \in A} Q(s, a; \mathbf{w}); \bar{\mathbf{w}}\right) - Q(s, a; \mathbf{w}) \right)^2 \right]$$
$$\approx \frac{1}{2} \sum_{(s,a,r,s') \sim D} \left[ \left( r + \gamma Q\left(s', \arg \max_{a' \in A} Q(s, a; \mathbf{w}); \bar{\mathbf{w}}\right) - Q(s, a; \mathbf{w}) \right)^2 \right]$$

“We therefore propose to *evaluate the greedy policy according to the online network*, but using the target network to estimate its value.” —  
[van Hasselt et al., AAAI 2016]

# Distributional Q-Learning

(Learn value distribution  $Z(s, a)$  & use  $E[Z(s, a)]$  as  $Q(s, a)$  in Q-Learning)

# Why Shall We Consider “Value Distributions”?

- ▶ **Risky vs safe choices**
  - ▶ E.g., Same expected return but different variance
- ▶ **Good empirical performance** (despite that the underlying root cause is not fully known)
  - ▶ C51 [Belleware et al., ICML 2017]
  - ▶ QR-DQN [Dabney et al., AAAI 2018]
  - ▶ IQN [Dabney et al., ICML 2018]
- ▶ **New approaches for exploration**
  - ▶ Information-directed exploration [Nikolov et al., ICLR 2019]
  - ▶ Distributional RL for efficient exploration [Mavrin et al., ICML 2019]
- ▶ **Learn better critics**
  - ▶ Truncated Quantile Critics (TQC) [Kuznetsov et al., ICML 2020]

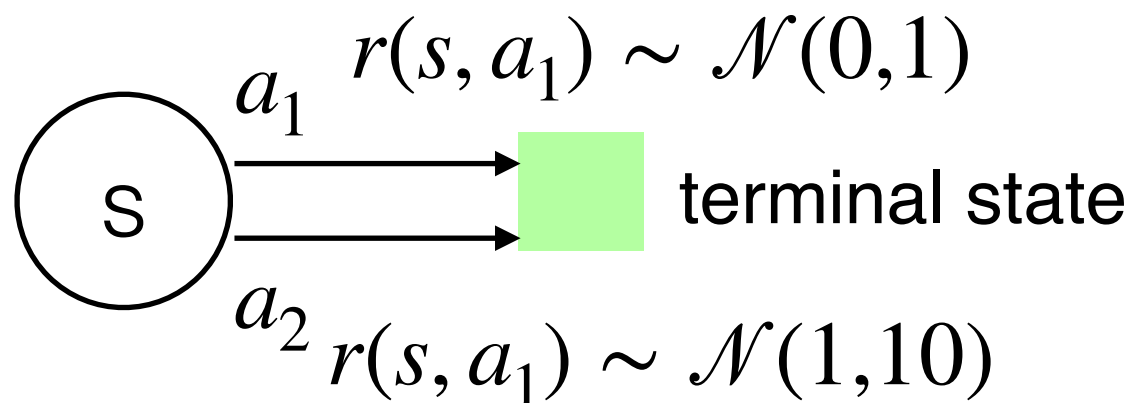
Question: How to learn the complete value distribution (instead of merely the expectation)?

# Sample Action-Value $Z^\pi(s, a)$

- ▶ Sample action-value  $Z^\pi(s, a)$ : sample return if we start from state  $s$  and take action  $a$ , and then follow policy  $\pi$

$$Q^\pi(s, a) = \mathbb{E}[Z^\pi(s, a)] = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)\right]$$

- ▶  $Z^\pi(s, a)$  is essentially a random variable (and hence follows some distribution)
- ▶ **Example**: 1-state MDP with 2 actions and  $\pi(s) = a_1$ 
  - ▶  $Q^\pi(s, a_1) = ?$   $Z^\pi(s, a_1) = ?$



- ▶  $Q^\pi(s, a_2) = ?$   $Z^\pi(s, a_2) = ?$

# Finding $Z^\pi$ via Distributional Bellman Equation

- ▶ **Mild assumption:**  $Z^\pi(s, a)$  has bounded moments
- ▶ Distributional Bellman equation for  $Z^\pi(s, a)$ :

$$Z^\pi(s, a) \stackrel{D}{=} r(s, a) + \gamma Z^\pi(s', a')$$

( $\stackrel{D}{=}$ : equal in distribution)

- ▶ **Question:** How to interpret this equation?
- ▶ **Question:** Are  $r(s, a)$  and  $Z^\pi(s', a')$  independent?
- ▶ **Question:** Is this consistent with Bellman expectation equation?

# Distributional Bellman Operator $B^\pi$

- ▶  $\mathcal{Z}$ : the space of all value distributions with bounded moments
- ▶ Transition operator  $P^\pi : \mathcal{Z} \rightarrow \mathcal{Z}$

$$P^\pi Z(s, a) \stackrel{D}{=} Z(s', a')$$

$$s' \sim P(\cdot | s, a), \quad a' \sim \pi(\cdot | s')$$

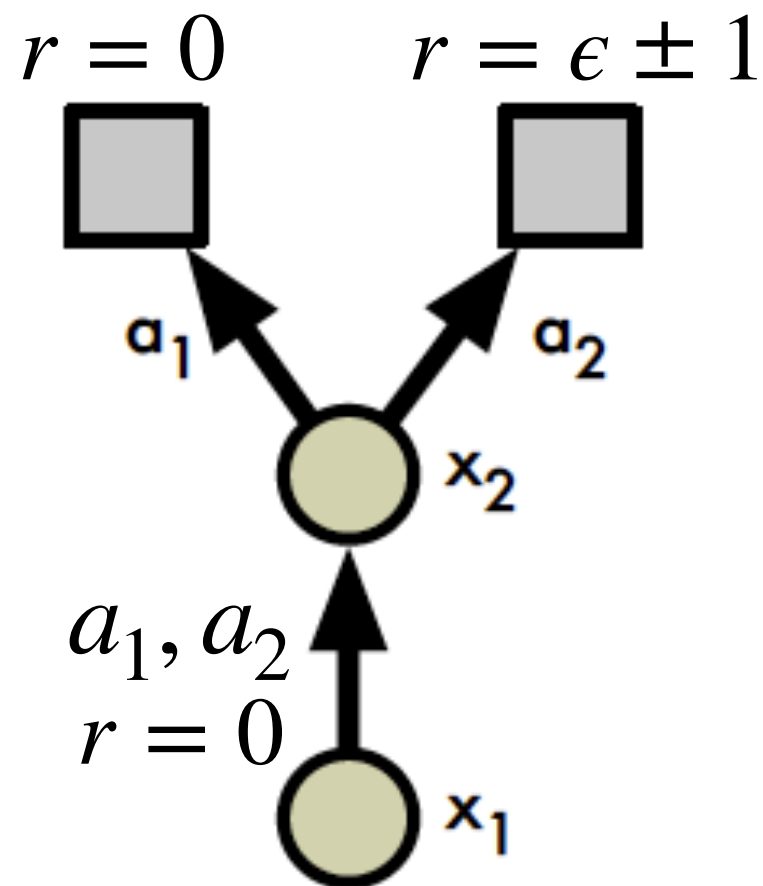
- ▶ Distributional Bellman operator  $B^\pi : \mathcal{Z} \rightarrow \mathcal{Z}$

$$B^\pi Z(s, a) \stackrel{D}{=} r(s, a) + \gamma P^\pi Z(s, a)$$



# An Example of Applying $B^\pi$

- ▶ **Example:** 2 states  $x_1, x_2$  and 2 actions  $a_1, a_2$
- ▶  $\pi(a_1 | x_2) = 0.3$ ,  $\pi(a_2 | x_2) = 0.7$ , and  $\gamma = 0.9$

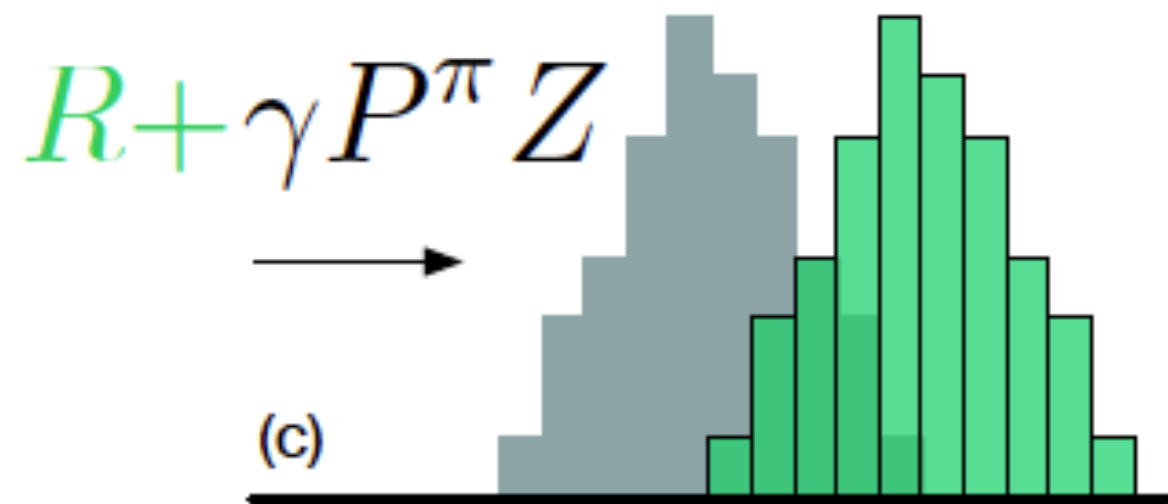
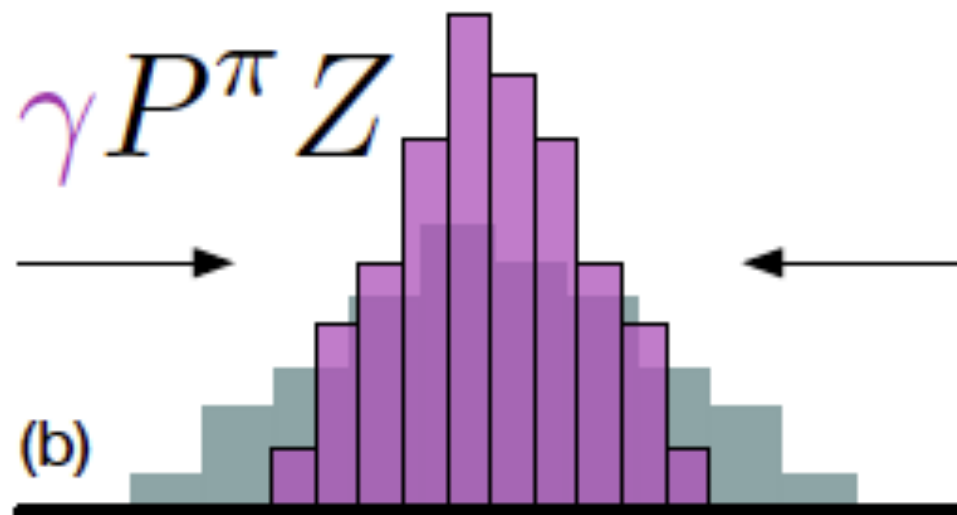
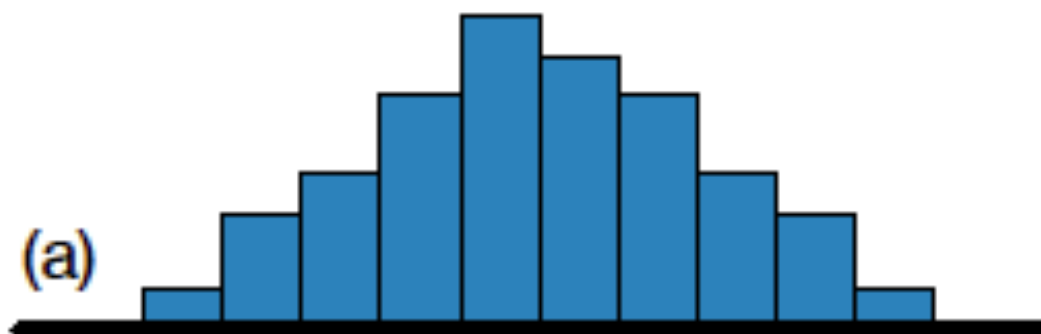


$$B^\pi Z(s, a) \stackrel{D}{=} r(s, a) + \gamma P^\pi Z(s, a)$$

- ▶ Suppose  $Z(x_1, a_1) = 0$ ,  $Z(x_2, a_1) = 0$  with probability 1 and  $Z(x_2, a_2) \sim \mathcal{N}(0, 1)$
- ▶ **Question:**  $B^\pi Z(x_2, a_2) = ?$   $B^\pi Z(x_1, a_1) = ?$

# Visualization of Distributional Bellman Operator

$$P^\pi Z$$



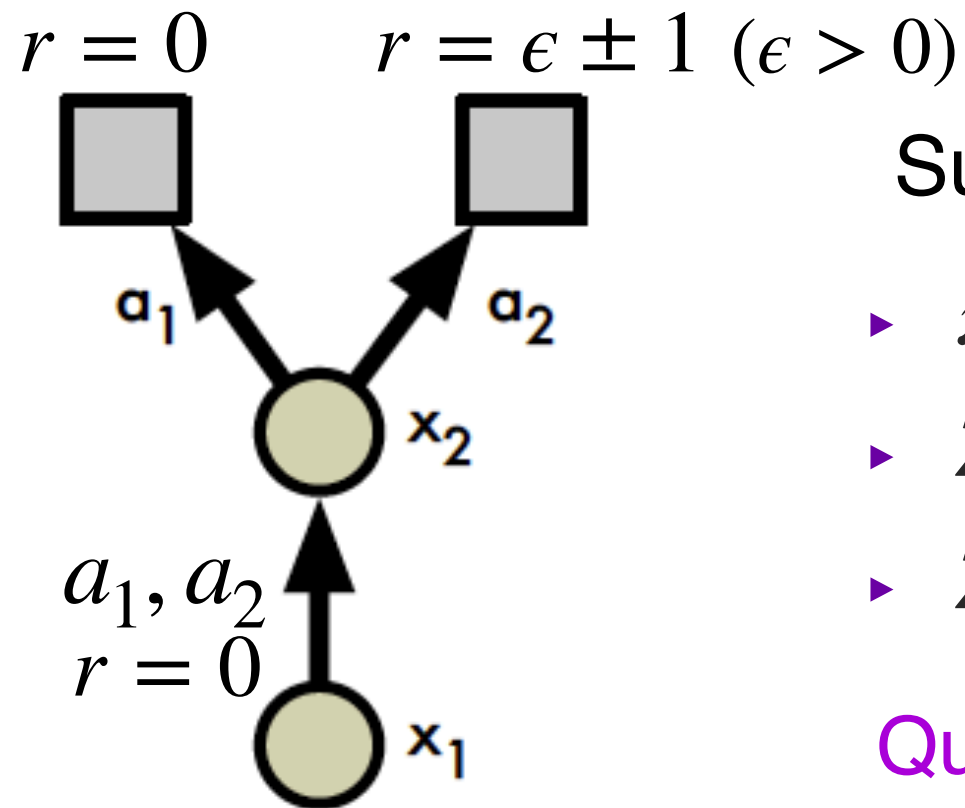
# Distributional “Optimality” Operator

Recall— Distributional Bellman operator  $B^\pi : \mathcal{Z} \rightarrow \mathcal{Z}$

$$B^\pi Z(s, a) \stackrel{D}{=} r(s, a) + \gamma P^\pi Z(s, a)$$

- Distributional optimality operator  $B^*$ : The  $B^\pi$  resulting from a greedy policy  $\pi$  (what does “greedy” mean here?)

# An Example of $B^*$



Suppose we have the following:

- ▶  $\pi(a_1 | x_2) = 0.3$ ,  $\pi(a_2 | x_2) = 0.7$ , and  $\gamma = 1$
- ▶  $Z(x_1, a_1) = 0$ ,  $Z(x_2, a_1) = 0$  with probability 1
- ▶  $Z(x_2, a_2) \sim \mathcal{N}(0, 1)$

**Question:** What's the PDF of  $B^*Z(x_1, a_1) = ?$

$$B^*Z(s, a) \stackrel{D}{=} r(s, a) + \gamma P^{\pi_{\text{greedy}}} Z(s, a)$$