

Deep Reinforcement Learning and Control


# Determinist PG, Re-parametrized PG

Fall 2020, CMU 10-703

Katerina Fragkiadaki



# Advantage Actor-Critic

- 
0. Initialize policy parameters  $\theta$  and critic parameters  $\phi$ .
  1. Sample trajectories  $\{\tau_i = \{s_t^i, a_t^i\}_{t=0}^T\}$  by deploying the current policy  $\pi_\theta(a_t | s_t)$ .
  2. Fit value function  $V_\phi^\pi(s)$  by MC or TD estimation (update  $\phi$ )
  3. Compute action advantages  $A^\pi(s_t^i, a_t^i) = R(s_t^i, a_t^i) + \gamma V_\phi^\pi(s_{t+1}^i) - V_\phi^\pi(s_t^i)$
  4.  $\nabla_\theta U(\theta) \approx \hat{g} = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t^i | s_t^i) A^\pi(s_t^i, a_t^i)$
  5.  $\theta \leftarrow \theta + \alpha \nabla_\theta U(\theta)$

# Policy gradients so far

Policy objective:

$$\max_{\theta} . \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)]$$

Advantage actor critic policy gradient:

$$\mathbb{E}_{s \sim d^{\pi_{\theta}}(s), a \sim \pi_{\theta}(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) [A(s, a; \phi)]$$

# Another policy objective

Previous policy objective:

$$\max_{\theta} . \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)]$$

New policy objective:

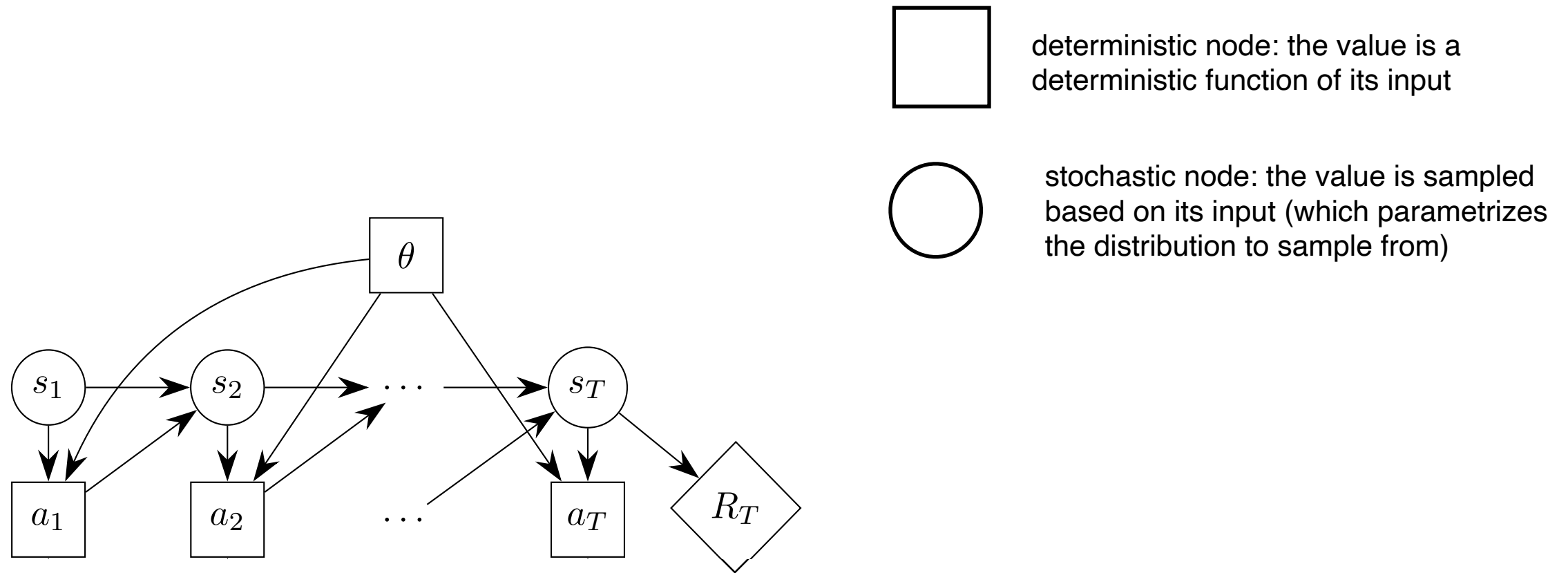
$$\max_{\theta} . \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[ \sum_{t=1}^T Q(s_t, a_t) \right]$$

Qs:

- Can we backpropagate through the Q function approximator?
- If the policy is deterministic we already know how to do it via the chain rule!

$$\mathbb{E} \sum_t \frac{dQ(s_t, a_t)}{d\theta} = \mathbb{E} \sum_t \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta}$$

# Deep Deterministic Policy Gradients

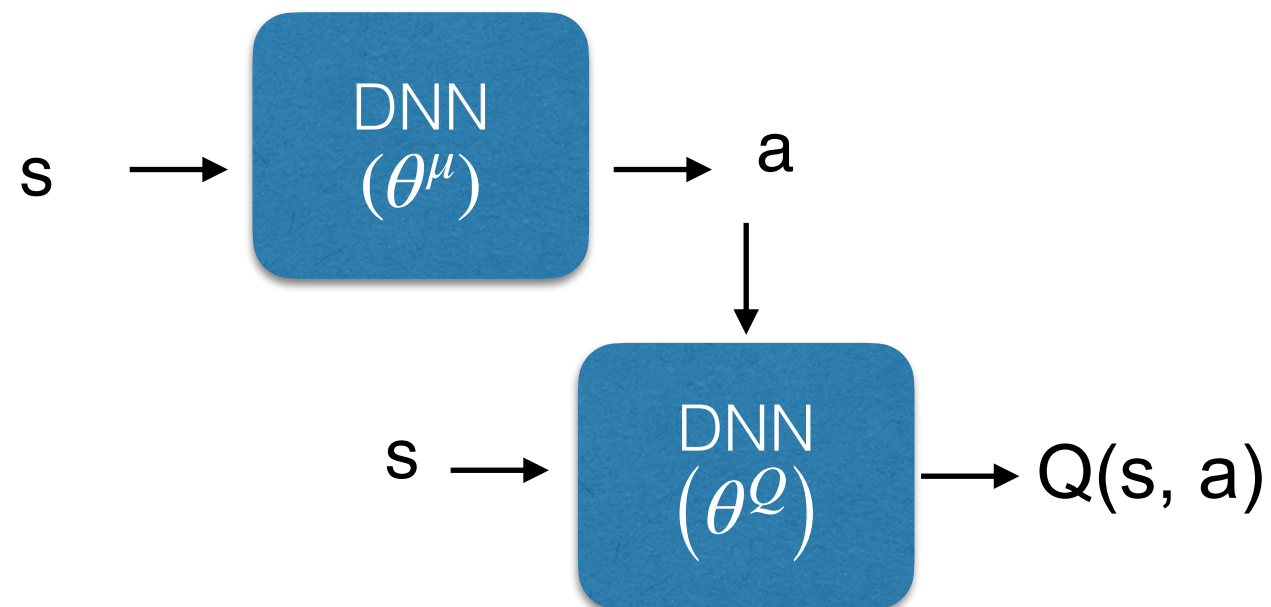


$$a = \pi_{\theta}(s)$$

$$\mathbb{E} \sum_t \frac{dQ(s_t, a_t)}{d\theta} = \mathbb{E} \sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta}$$

# Deep Deterministic Policy Gradients

The computational graph:



We are following a stochastic behavior policy to collect data.  
DDPG : Deep Q learning for continuous actions

# Deep Deterministic Policy Gradients

---

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

        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:

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

    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**

---

Fitting the Q function

# Another policy objective

Previous policy objective:

$$\max_{\theta} . \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)]$$

New policy objective:

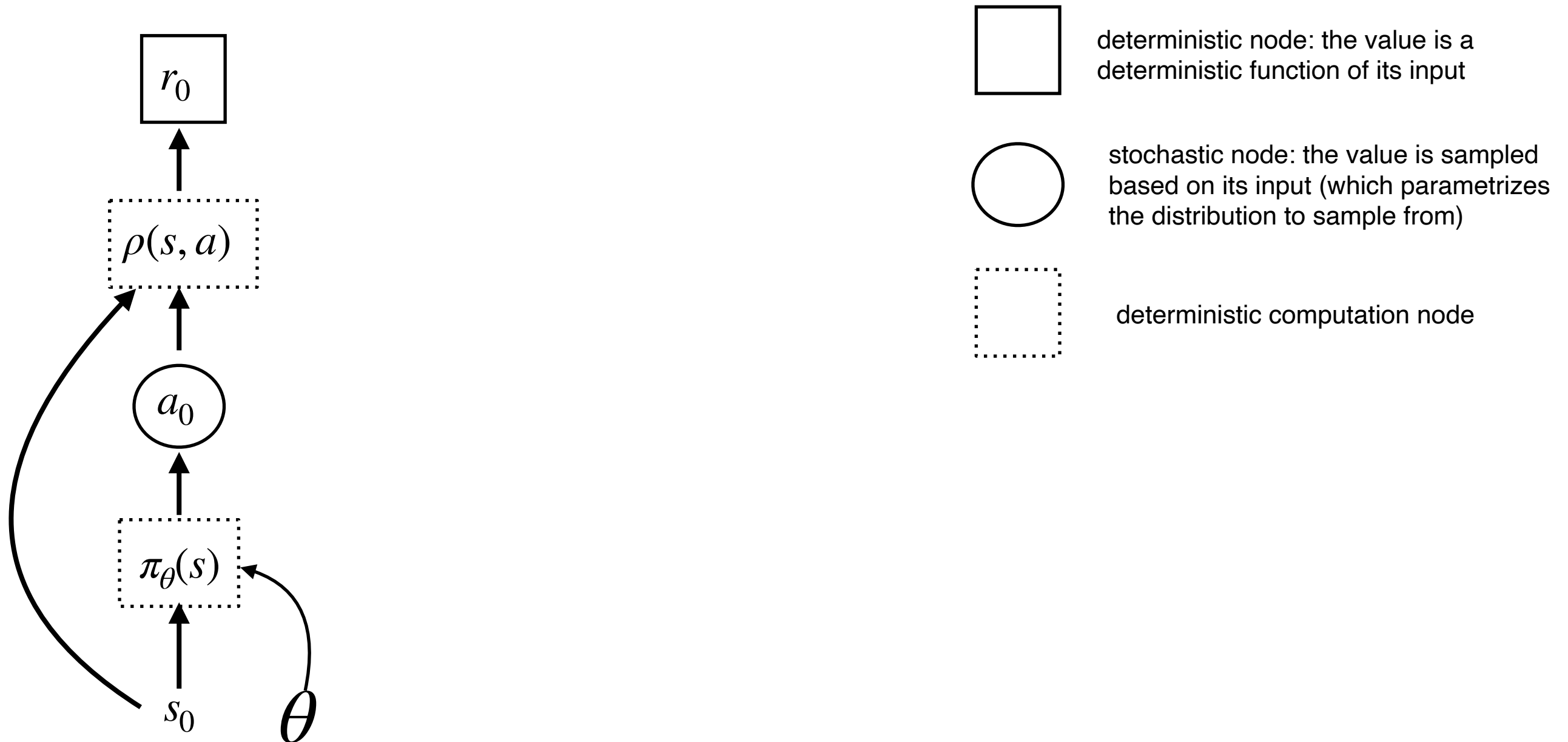
$$\max_{\theta} . \mathbb{E}_{\tau \sim P_{\theta}(\tau)} \left[ \sum_{t=1}^T Q(s_t, a_t) \right]$$

Qs:

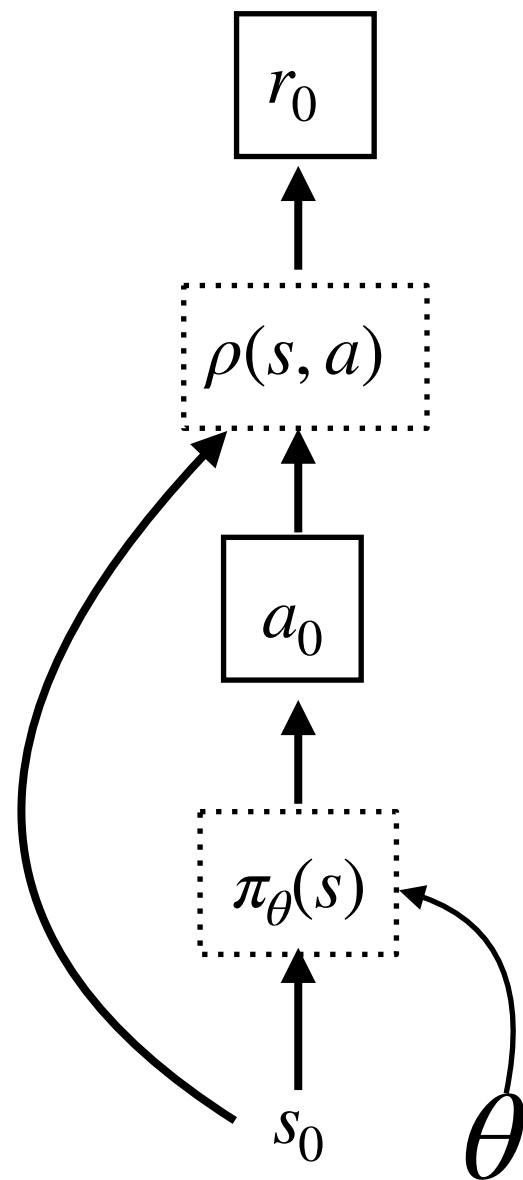
- Can we backpropagate through the Q function approximator?
- If the policy is deterministic we already know how to do it via the chain rule!
- What if the policy is a parametrized Gaussian distribution?



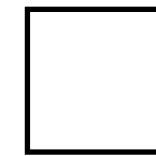
# Imagine we knew the reward function $\rho(s, a)$



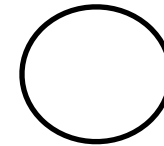
# Deterministic policy



$$a = \pi_\theta(s)$$



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

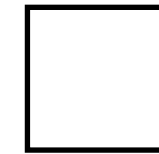
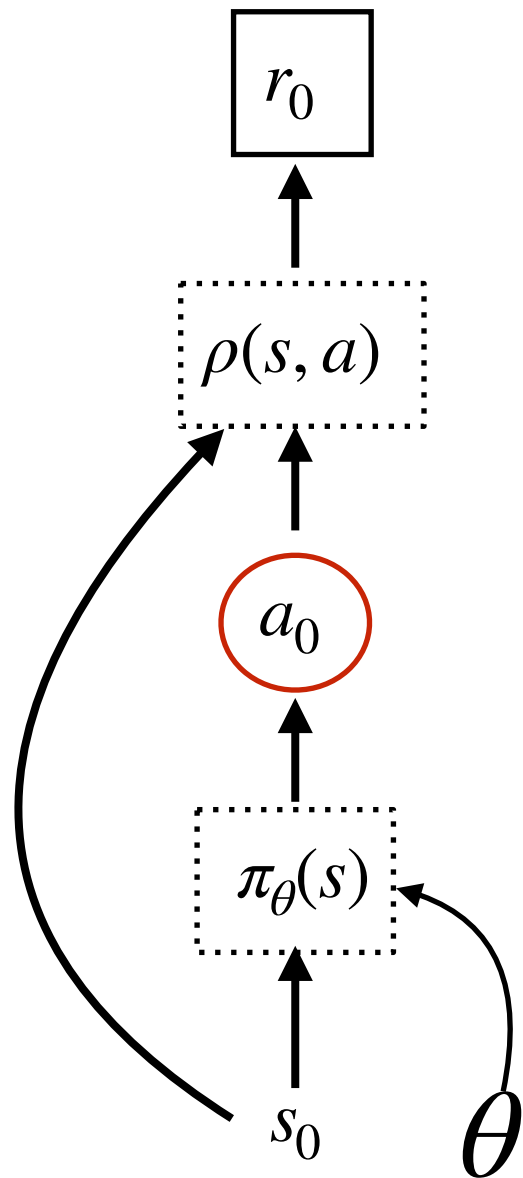
I want to learn  $\theta$  to maximize the average reward obtained.

$$\max_{\theta} \rho(s_0, a)$$

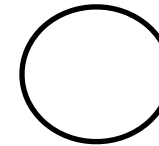
I can compute the gradient with the chain rule.

$$\nabla_{\theta} \rho(s, a) = \frac{d\rho}{da} \frac{da}{d\theta}$$

# Stochastic policy



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



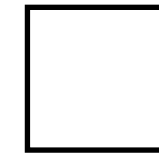
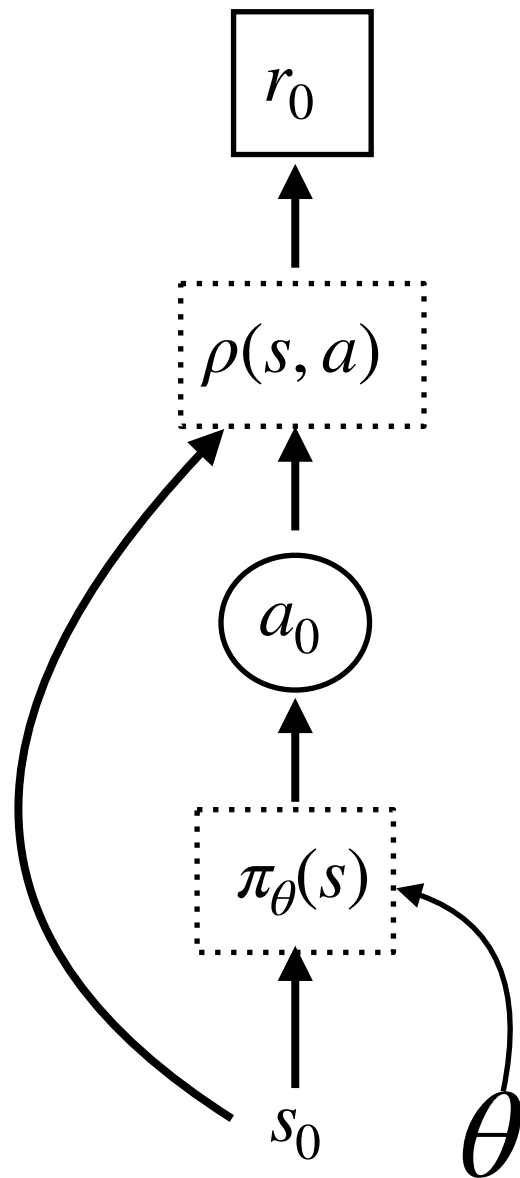
deterministic computation node

I want to learn  $\theta$  to maximize the average reward obtained.

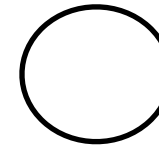
$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

$$\nabla_{\theta} \mathbb{E}_a \rho(s_0, a)$$

# Stochastic policy



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

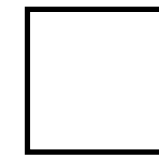
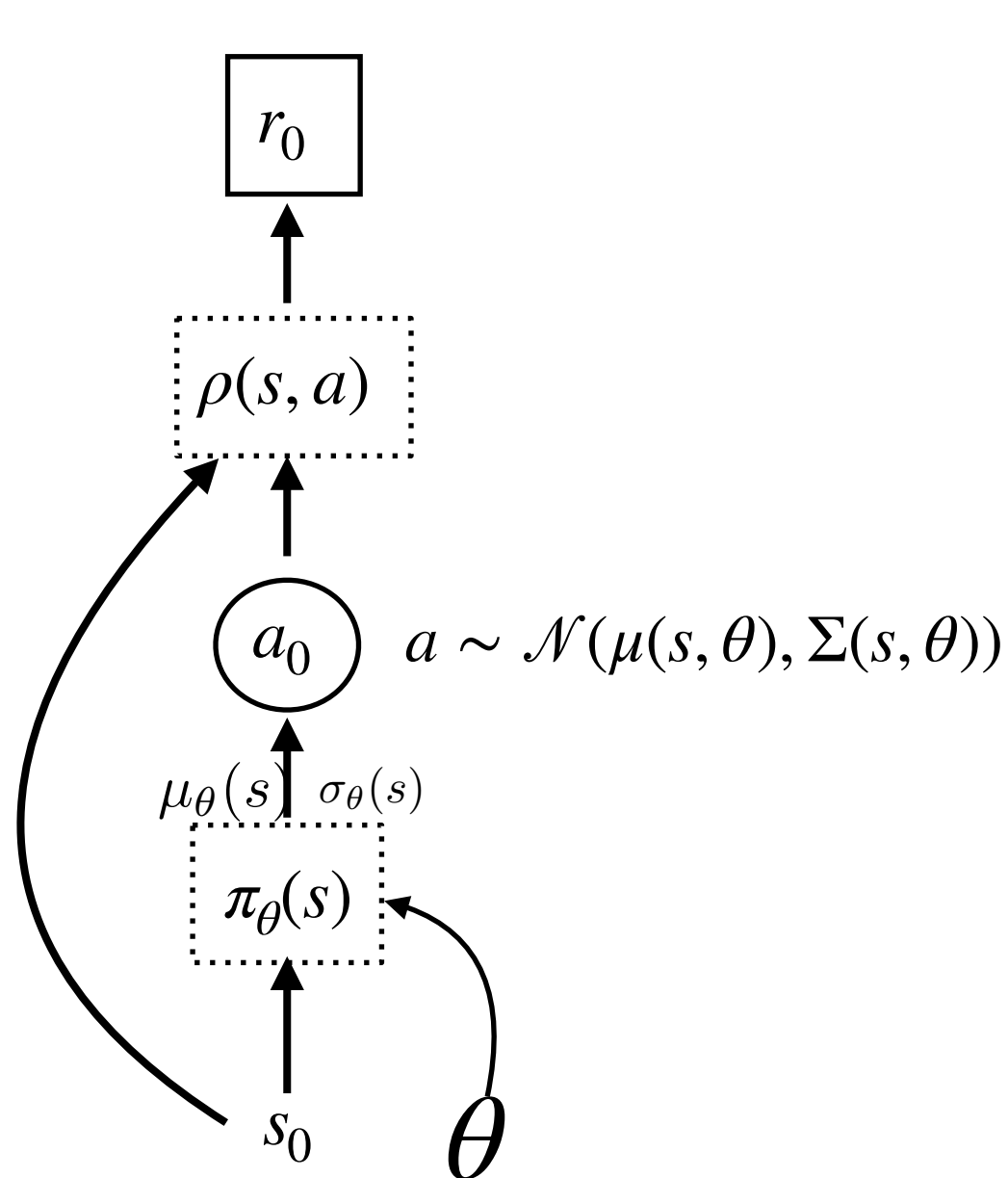
I want to learn  $\theta$  to maximize the average reward obtained.

$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

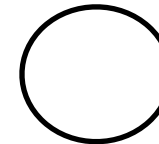
Likelihood ratio estimator, works for both continuous and discrete actions

$$\mathbb{E}_a \nabla_{\theta} \log \pi_{\theta}(s) \rho(s_0, a)$$

# Example: Gaussian policy



deterministic node: the value is a deterministic function of its input



stochastic node: the value is sampled based on its input (which parametrizes the distribution to sample from)



deterministic computation node

I want to learn  $\theta$  to maximize the average reward obtained.

$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$

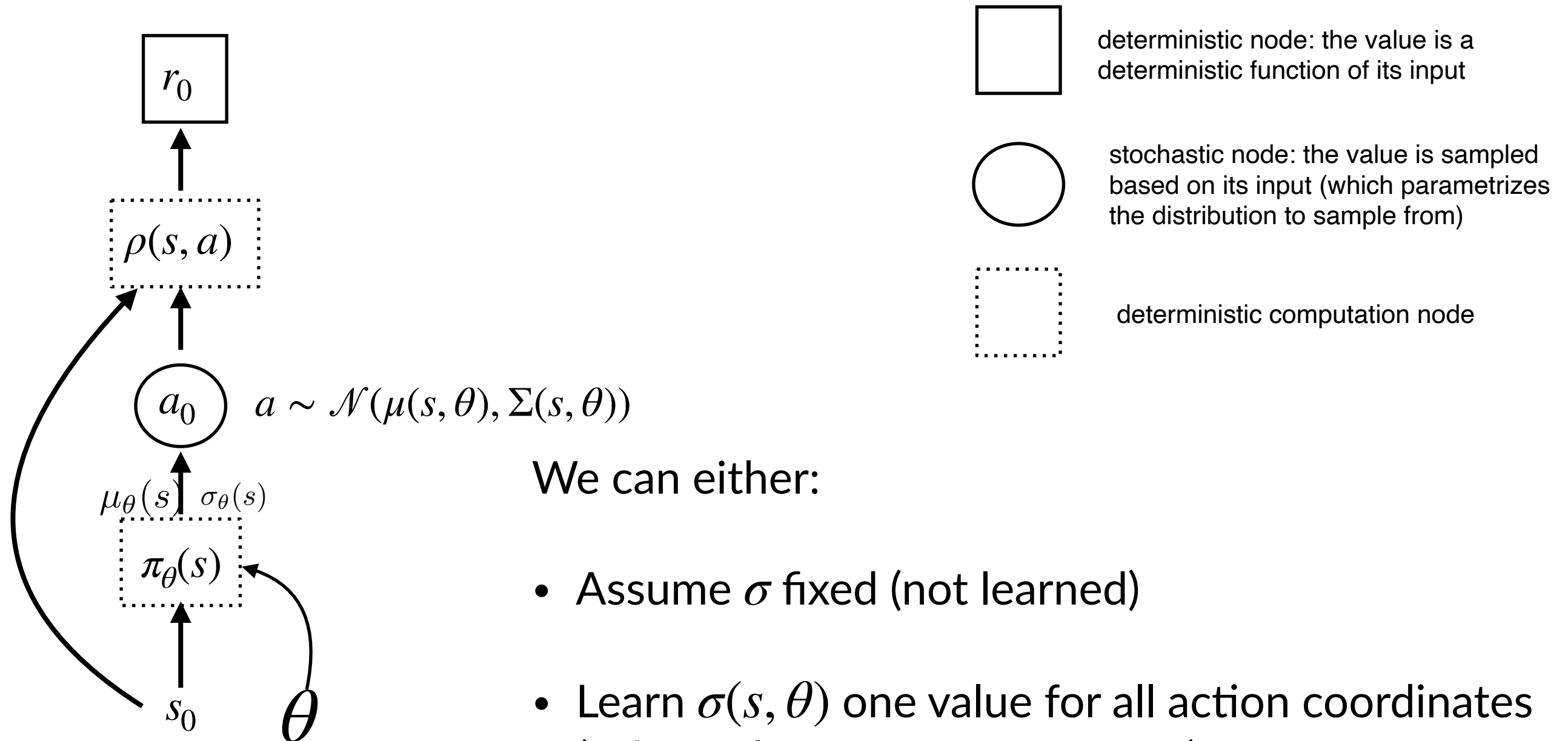
Likelihood ratio estimator, works for both continuous and discrete actions

$$\mathbb{E}_a \nabla_{\theta} \log \pi_{\theta}(s) \rho(s_0, a)$$

If  $\sigma^2$  is constant:

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s; \theta)) \frac{\partial \mu(s; \theta)}{\partial \theta}}{\sigma^2}$$

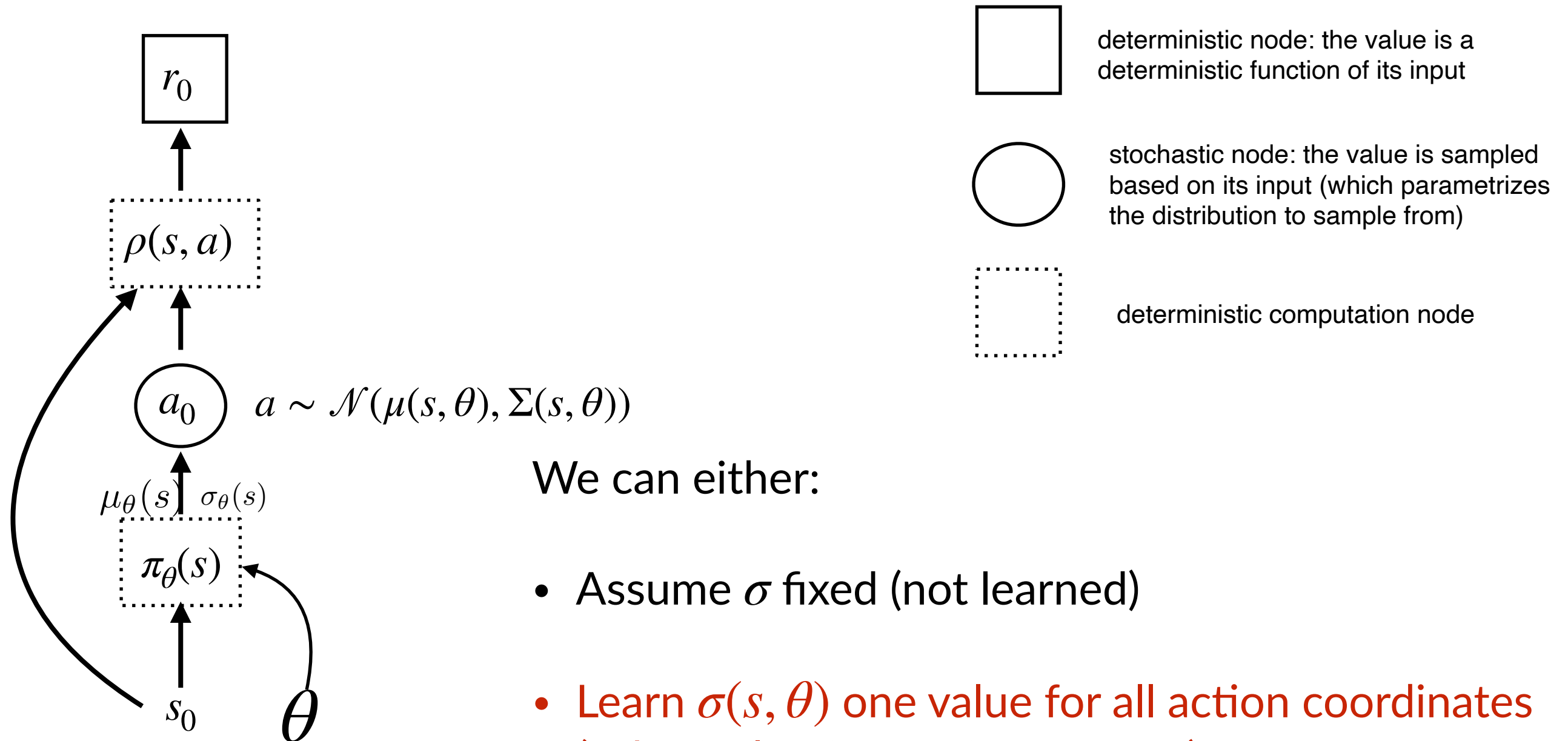
# Example: Gaussian policy



We can either:

- Assume  $\sigma$  fixed (not learned)
- Learn  $\sigma(s, \theta)$  one value for all action coordinates (spherical or isotropic Gaussian)
- Learn  $\sigma^i(s, \theta), i = 1 \cdots n$  (diagonal covariance)
- Learn a full covariance matrix  $\Sigma(s, \theta)$

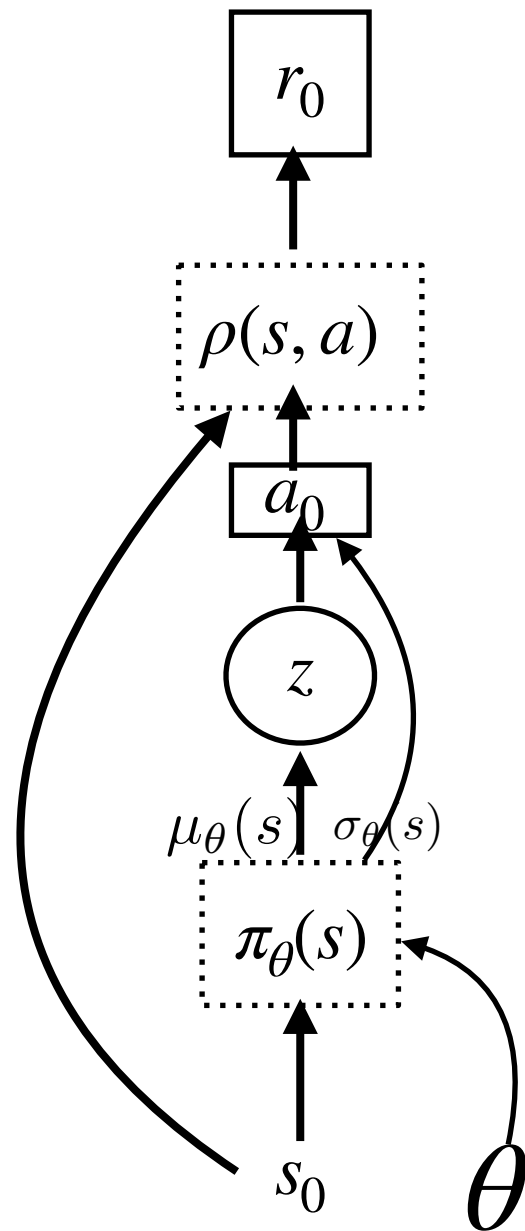
# Example: Gaussian policy



We can either:

- Assume  $\sigma$  fixed (not learned)
- Learn  $\sigma(s, \theta)$  one value for all action coordinates (spherical or isotropic Gaussian)
- Learn  $\sigma^i(s, \theta), i = 1 \cdots n$  (diagonal covariance)
- Learn a full covariance matrix  $\Sigma(s, \theta)$

# Re-parametrization for Gaussian



Instead of:  $a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$

We can write:  $a = \mu(s, \theta) + z\sigma(s, \theta) \quad z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$

Because:  $\mathbb{E}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \mu(s, \theta)$   
 $\text{Var}_z(\mu(s, \theta) + z\sigma(s, \theta)) = \sigma(s, \theta)^2 \mathbf{I}_{n \times n}$

$$\max_{\theta} \mathbb{E}_a \rho(s_0, a)$$



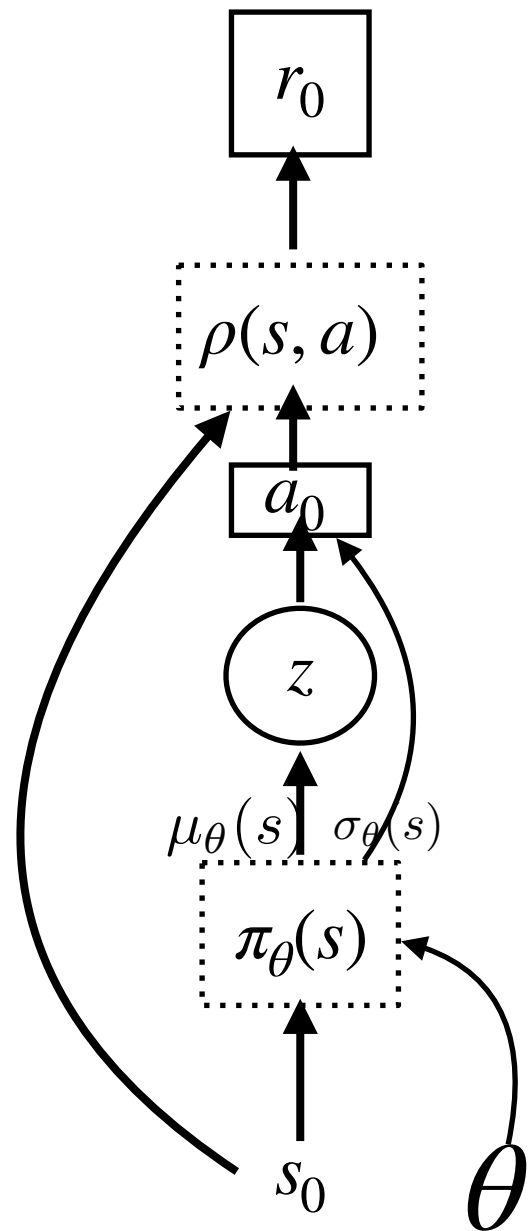
$$\max_{\theta} \mathbb{E}_z \rho(s_0, a(z))$$

Qs:

- Does  $a$  depend on  $\theta$  ?
- Does  $z$  depend on  $\theta$  ?



# Re-parametrization for Gaussian



Instead of:  $a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$

We can write:  $a = \mu(s, \theta) + z\sigma(s, \theta) \quad z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$

What do we gain?

$$\nabla_{\theta} \mathbb{E}_z \left[ \rho(a(\theta, z), s) \right] = \mathbb{E}_z \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta}$$

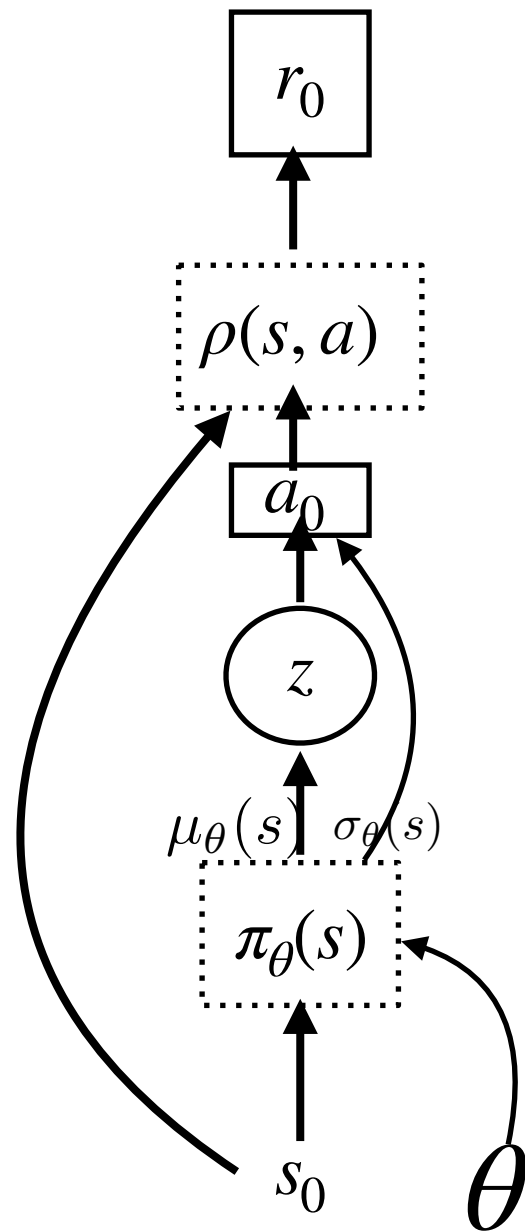
$$\frac{da(\theta, z)}{d\theta} = \frac{d\mu(s, \theta)}{d\theta} + z \frac{d\sigma(s, \theta)}{d\theta}$$

$$\max_{\theta} \cdot \mathbb{E}_a \rho(s_0, a)$$



$$\max_{\theta} \cdot \mathbb{E}_z \rho(s_0, a(z))$$

# Re-parametrization for Gaussian



Instead of:  $a \sim \mathcal{N}(\mu(s, \theta), \Sigma(s, \theta))$

We can write:  $a = \mu(s, \theta) + z\sigma(s, \theta) \quad z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n})$

What do we gain?

$$\nabla_{\theta} \mathbb{E}_z \left[ \rho(a(\theta, z), s) \right] = \mathbb{E}_z \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta}$$

$$\frac{da(\theta, z)}{d\theta} = \frac{d\mu(s, \theta)}{d\theta} + z \frac{d\sigma(s, \theta)}{d\theta}$$

$$\max_{\theta} \cdot \mathbb{E}_a \rho(s_0, a)$$

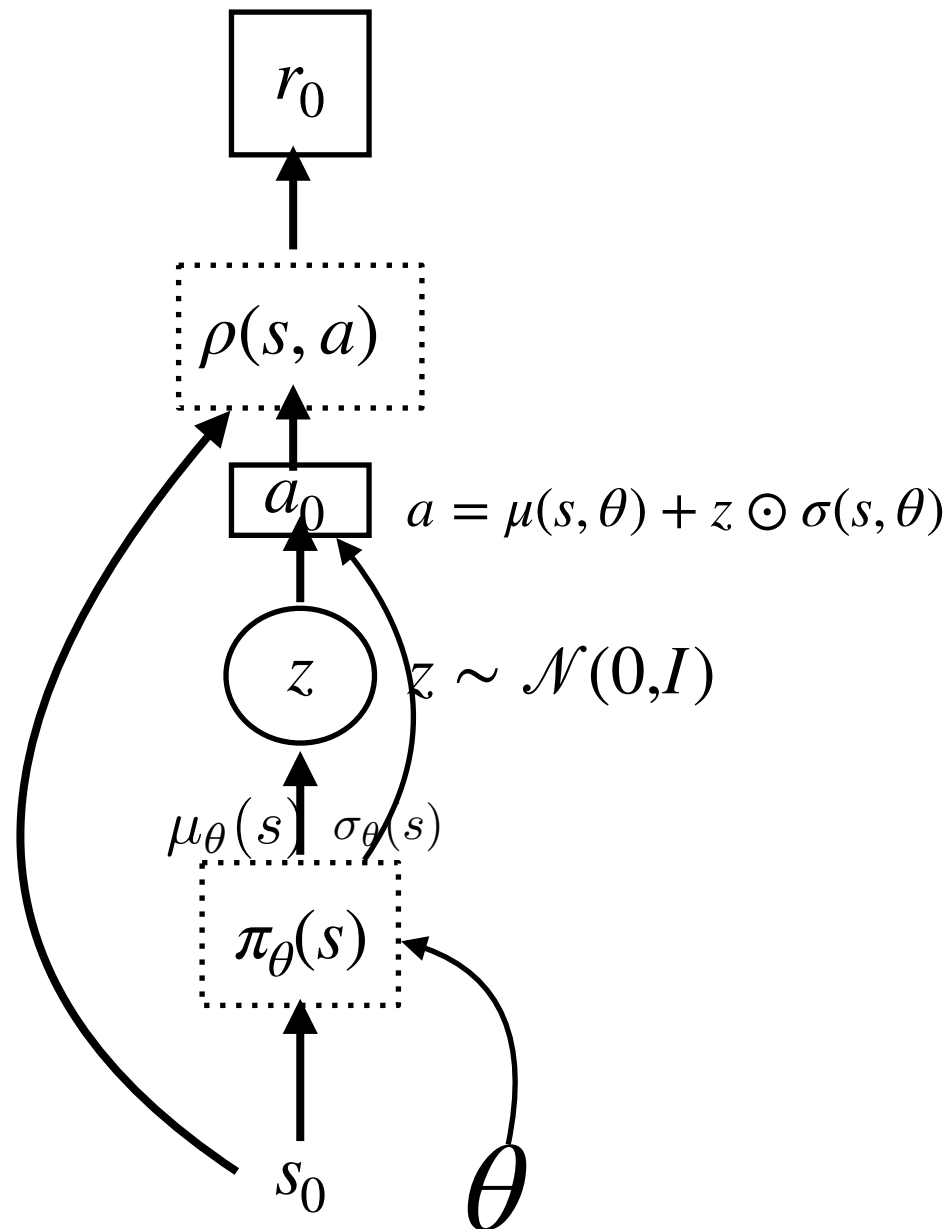


$$\max_{\theta} \cdot \mathbb{E}_z \rho(s_0, a(z))$$

Sample estimate:

$$\nabla_{\theta} \frac{1}{N} \sum_{i=1}^N \left[ \rho(a(\theta, z_i), s) \right] = \frac{1}{N} \sum_{i=1}^N \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta} \Big|_{z=z_i}$$

# Re-parametrization for Gaussian



Likelihood ratio grad estimator:

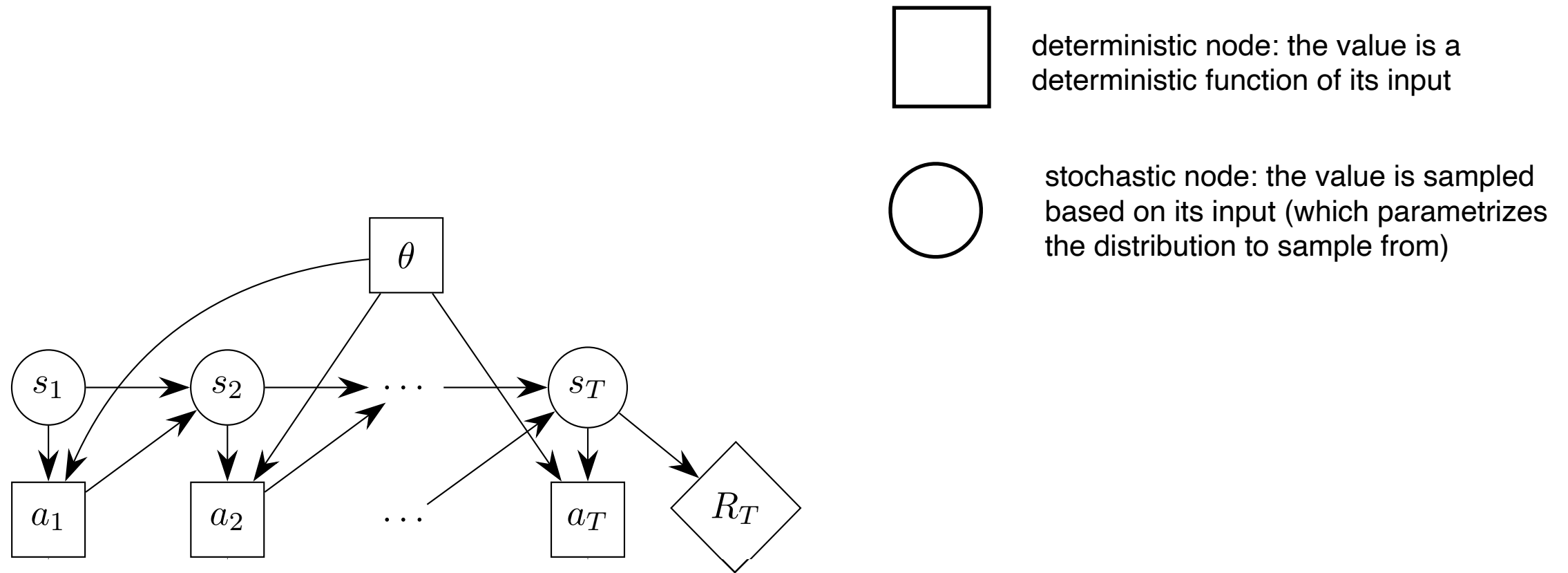
$$\mathbb{E}_a \nabla_\theta \log \pi_\theta(s, a) \rho(s, a)$$

Pathwise derivative:

$$\mathbb{E}_z \frac{d\rho(a(\theta, z), s)}{da} \frac{da(\theta, z)}{d\theta}$$

The pathwise derivative uses the derivative of the reward w.r.t. the action!

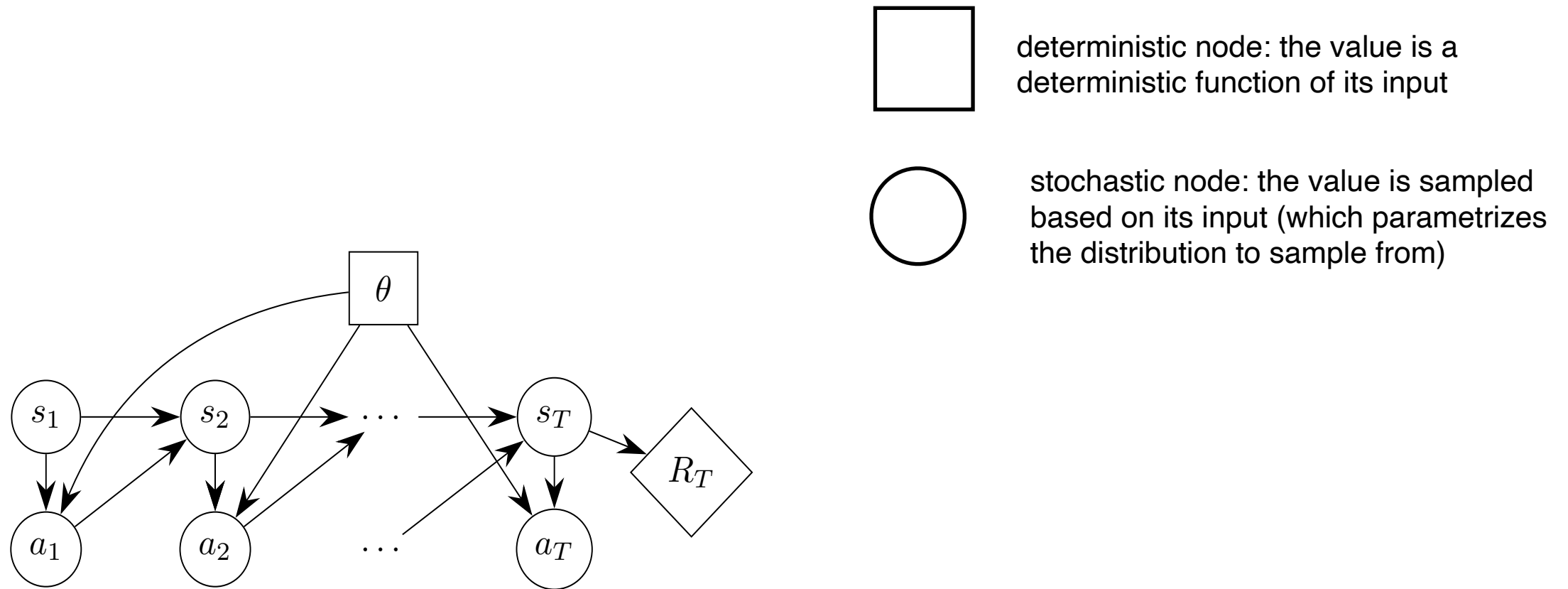
# Deep Deterministic Policy Gradients



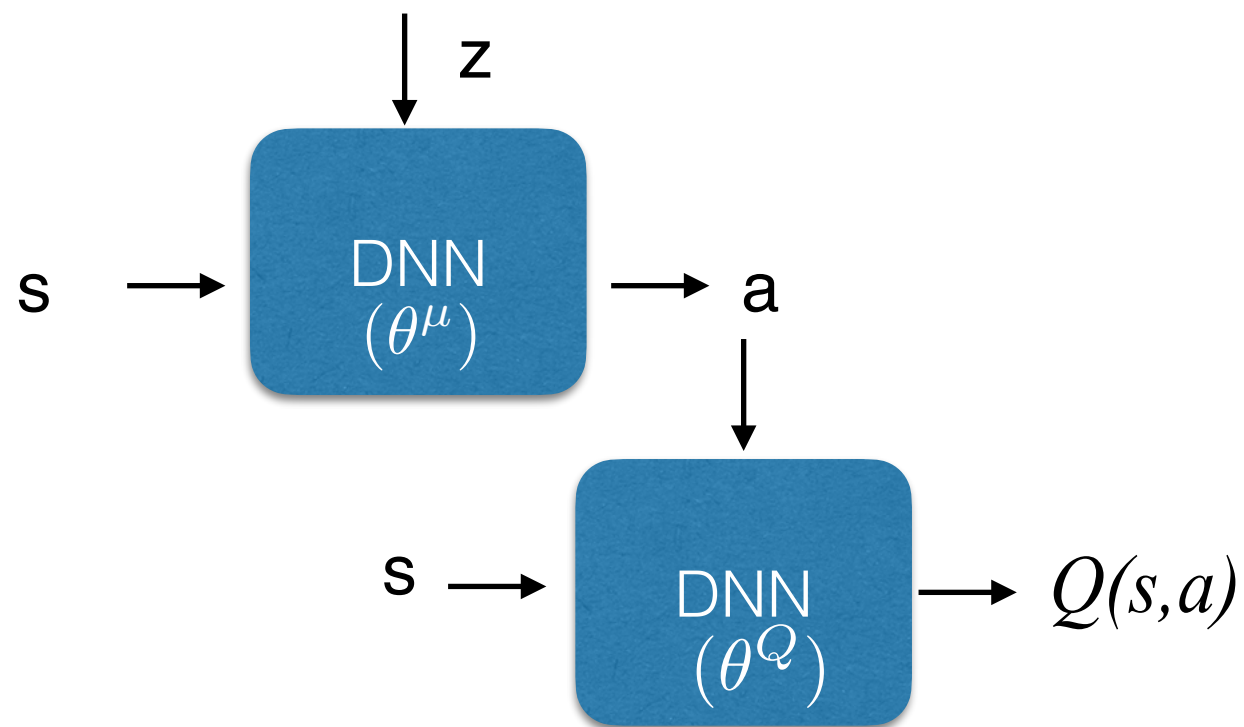
$$a = \pi_{\theta}(s)$$

$$\mathbb{E} \sum_t \frac{dQ(s_t, a_t)}{d\theta} = \mathbb{E} \sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta}$$

# Re-parametrized Policy Gradients

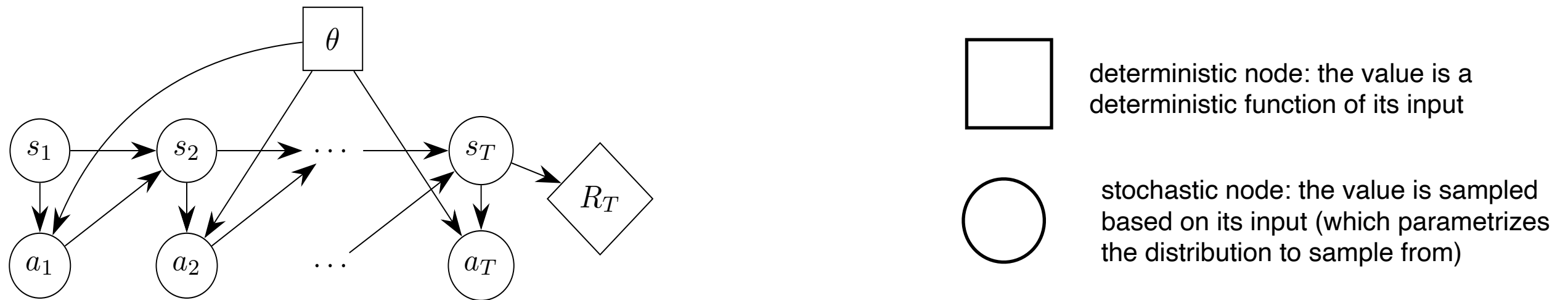


$$z \sim \mathcal{N}(0, 1)$$

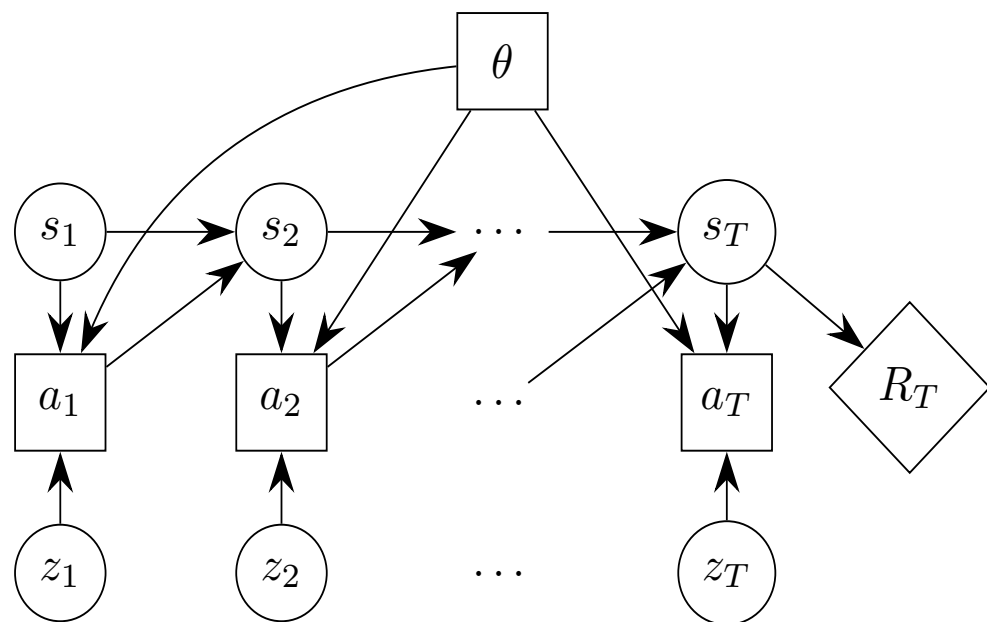


$$a = \mu(s; \theta) + z\sigma(s; \theta)$$

# Re-parametrized Policy Gradients



- Reparameterize:  $a_t = \pi(s_t, z_t, \theta)$ .  $z_t$  is noise from fixed distribution



$$\mathbb{E} \sum_t \frac{dQ(s_t, a_t)}{d\theta} = \mathbb{E} \sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta} = \mathbb{E} \sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \left( \frac{d\mu(s_t; \theta)}{d\theta} + z_t \frac{d\sigma(s_t; \theta)}{d\theta} \right)$$

# Stochastic Value Gradients V0

**for** iteration=1, 2, ... **do**

Execute policy  $\pi_\theta$  to collect  $T$  timesteps of data

Update  $\pi_\theta$  using  $g \propto \nabla_\theta \sum_{t=1}^T Q(s_t, \pi(s_t, z_t; \theta))$

Update  $Q_\phi$  using  $g \propto \nabla_\phi \sum_{t=1}^T (Q_\phi(s_t, a_t) - \hat{Q}_t)^2$ , e.g. with TD( $\lambda$ )

**end for**



# Computing Gradients of Expectations

- When the variable w.r.t. which we are differentiating appears inside the expectation:

$$\nabla_{\theta} \mathbb{E}_{x \sim P(x)} f(x(\theta)) = \mathbb{E}_{x \sim P(x)} \nabla_{\theta} f(x(\theta)) = \mathbb{E}_{x \sim P(x)} \frac{df(x(\theta))}{dx} \frac{dx}{d\theta}$$

- When the variable w.r.t. which we are differentiating appears in the distribution:  $\nabla_{\theta} \mathbb{E}_{x \sim P_{\theta}(x)} f(x)$

Likelihood ratio gradient estimator:

$$\mathbb{E}_{x \sim P_{\theta}(x)} \nabla_{\theta} \log P_{\theta}(x) f(x)$$

Re-parametrized gradient for Gaussian distributions:

$$\nabla_{\theta} \mathbb{E}_{z \sim \mathcal{N}(0, I)} f(x(z, \theta)) = \mathbb{E}_{z \sim \mathcal{N}(0, I)} \frac{df}{dx} \left( \frac{d\mu(\theta)}{d\theta} + z \frac{d\sigma(\theta)}{d\theta} \right)$$

Deep Reinforcement Learning and Control

# Goal Relabeling

Fall 2020, CMU 10-703

Katerina Fragkiadaki



# Universal value function Approximators

$$V(s; \theta) \rightarrow V(s, g; \theta)$$

$$\pi(s; \theta) \rightarrow \pi(s, g; \theta)$$

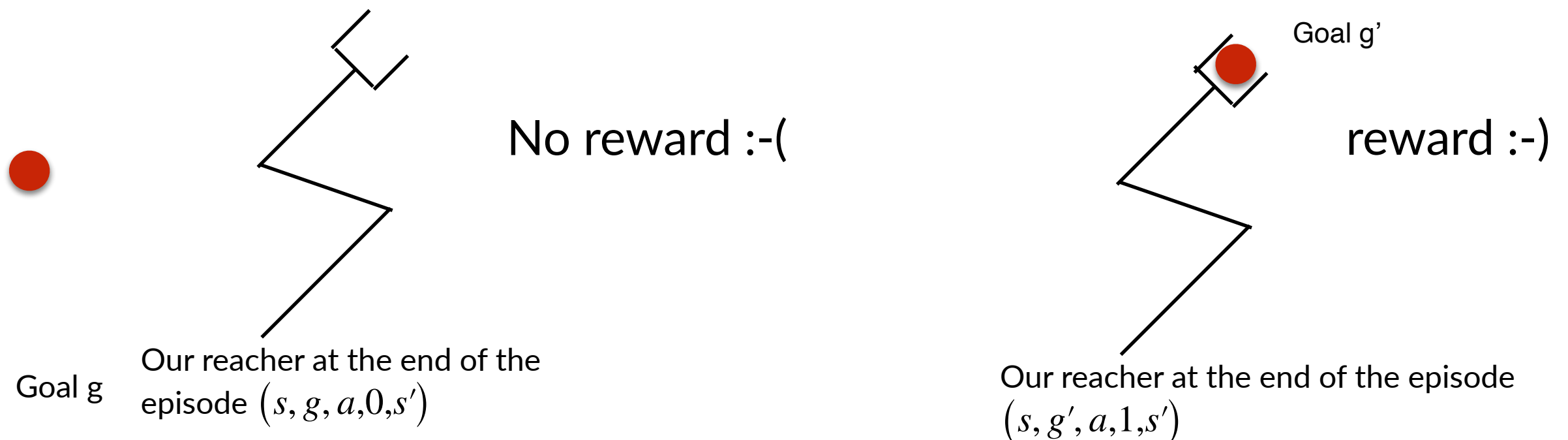
- All methods we have learnt so far can be used.
- At the beginning of an episode, we sample not only a start state but also a goal  $g$ , which stays constant throughout the episode
- The experience tuples should contain the goal.

$$(s, a, r, s') \rightarrow (s, g, a, r, s')$$

# Hindsight Experience Replay

**Marcin Andrychowicz<sup>\*</sup>, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel<sup>†</sup>, Wojciech Zaremba<sup>†</sup>**  
OpenAI

**Main idea:** use failed executions under one goal  $g$ , as successful executions under an alternative goal  $g'$  (which is where we ended at the end of the episode).



# Hindsight Experience Replay

---

**Algorithm 1** Hindsight Experience Replay (HER)

---

**Given:**

- an off-policy RL algorithm  $\mathbb{A}$ ,
- a strategy  $\mathbb{S}$  for sampling goals for replay,
- a reward function  $r : \mathcal{S} \times \mathcal{A} \times \mathcal{G} \rightarrow \mathbb{R}$ .

▷ e.g. DQN, DDPG, NAF, SDQN

▷ e.g.  $\mathbb{S}(s_0, \dots, s_T) = m(s_T)$

▷ e.g.  $r(s, a, g) = -[f_g(s) = 0]$

▷ e.g. initialize neural networks

Initialize  $\mathbb{A}$

Initialize replay buffer  $R$

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

    Sample a goal  $g$  and an initial state  $s_0$ .

**for**  $t = 0, T - 1$  **do**

        Sample an action  $a_t$  using the behavioral policy from  $\mathbb{A}$ :

$$a_t \leftarrow \pi_b(s_t || g)$$

▷  $||$  denotes concatenation

        Execute the action  $a_t$  and observe a new state  $s_{t+1}$

**end for**

**for**  $t = 0, T - 1$  **do**

$$r_t := r(s_t, a_t, g)$$

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

▷ standard experience replay

    Sample a set of additional goals for replay  $G := \mathbb{S}(\text{current episode})$

**for**  $g' \in G$  **do**

$$r' := r(s_t, a_t, g')$$

        Store the transition  $(s_t || g', a_t, r', s_{t+1} || g')$  in  $R$

←  $G$  : the states of the current episode

▷ HER

**end for**

**end for**

**for**  $t = 1, N$  **do**

        Sample a minibatch  $B$  from the replay buffer  $R$

        Perform one step of optimization using  $\mathbb{A}$  and minibatch  $B$

**end for**

**end for**

---

← Usually as additional goal we pick the goal that this episode achieved, and the reward becomes non zero

# Hindsight Experience Replay

--- DDPG    — DDPG+count-based exploration    — DDPG+HER    — DDPG+HER (version from Sec. 4.5)

