

10703 Deep Reinforcement Learning and Control

Russ Salakhutdinov

Machine Learning Department

rsalakhu@cs.cmu.edu

Slides developed and borrowed from
Katerina Fragkiadaki

End-to-end Model Based
Reinforcement Learning

Today's Lecture

End-to-end policy optimization
through back-propagation

Last week: Trajectory optimization



\mathbf{x}^0



target


Last week: Trajectory optimization

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$


 \mathbf{x}^0

 target

Last week: Trajectory optimization

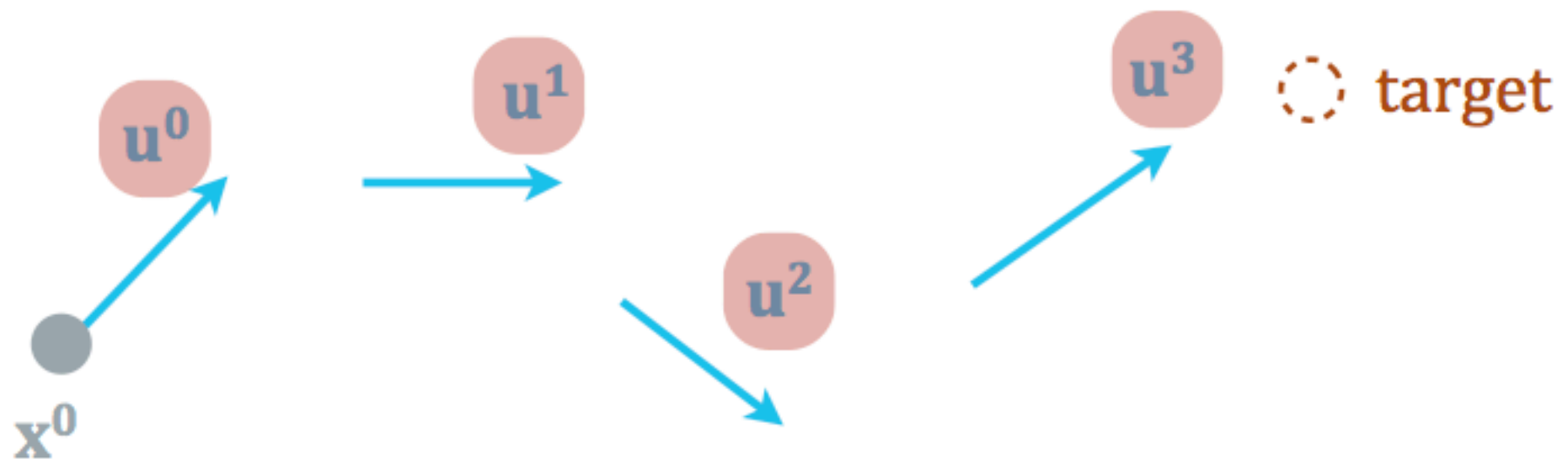

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$

○ target

●
 \mathbf{x}^0

Last week: Trajectory optimization

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$

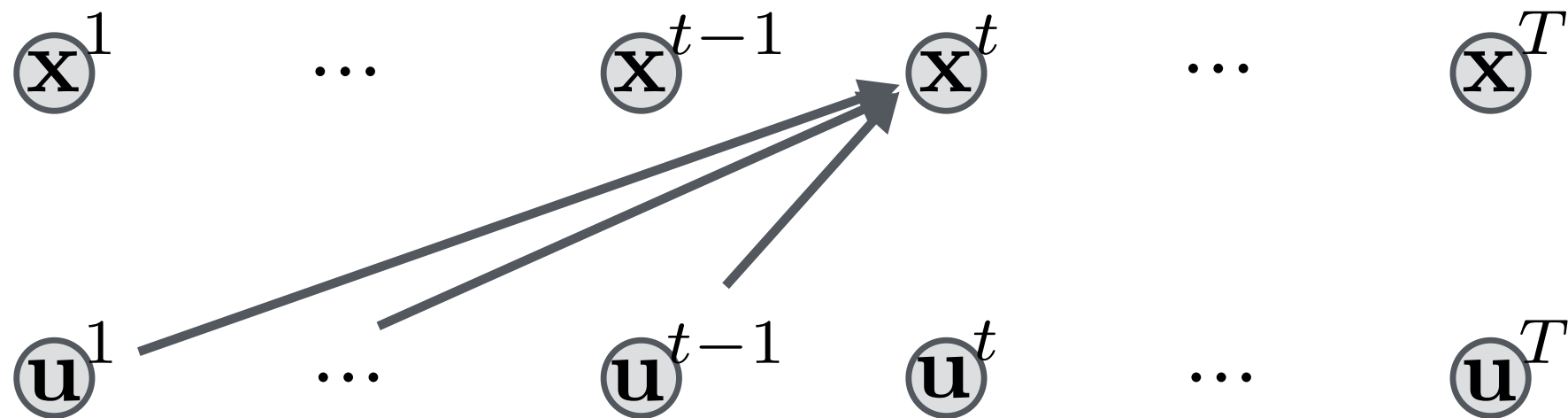


Poor Conditioning

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots \\ \dots + c(f(f(\dots) \dots), \mathbf{u}_T)$$

Poor Conditioning

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_T} c(\mathbf{x}_1, \mathbf{u}_1) + c(f(\mathbf{x}_1, \mathbf{u}_1), \mathbf{u}_2) + \dots \\ \dots + c(f(f(\dots) \dots), \mathbf{u}_T)$$



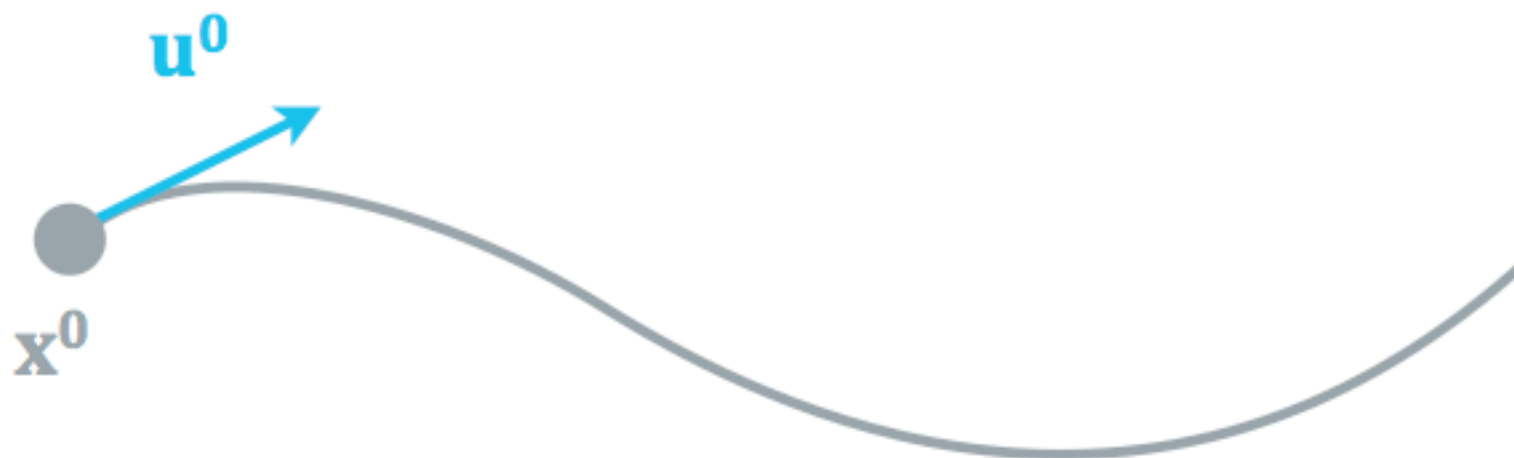
Poor Conditioning

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$



Poor Conditioning

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$



Poor Conditioning

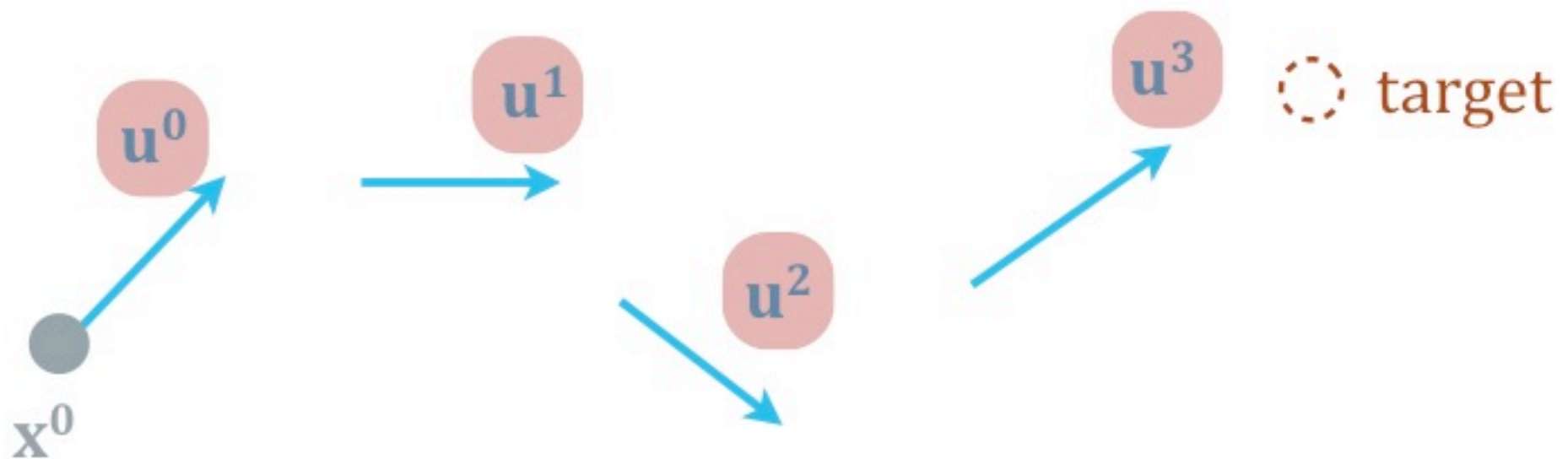
$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$



Trajectory Optimization

- Consider the special case of quadratic c and linear f

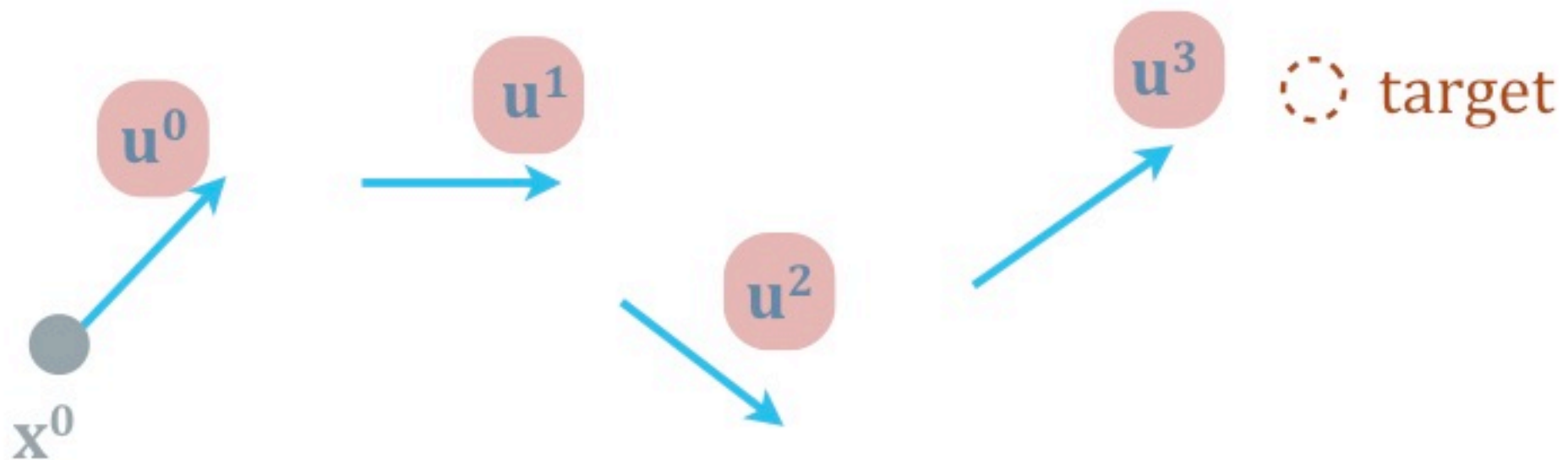
$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1} c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$



Trajectory Optimization

- Consider the special case of quadratic c and linear f

$$\min_{\mathbf{u}^0 \dots \mathbf{u}^T} \sum_{t=1} c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \mathbf{u}^t)$$



Solve it using **dynamic programming** :

- write u_t^* as function of the state x_t at each $t = T, \dots, 1$
- substitute x_0 (known)
- for $t = 1, \dots, T$ substitute x_t into u_t^* and fire the dynamics forward to obtain next state state ($x_{t+1} = f(x_t, u_t^*)$)

Learning Control Policies

$$\pi_{\theta} : \mathbf{x} \mapsto \mathbf{u}$$


 \mathbf{x}^0



Learning Control Policies

$$\pi_{\theta} : \mathbf{x} \mapsto \mathbf{u}$$

$$\min_{\theta} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \pi_{\theta}(\mathbf{x}^t))$$

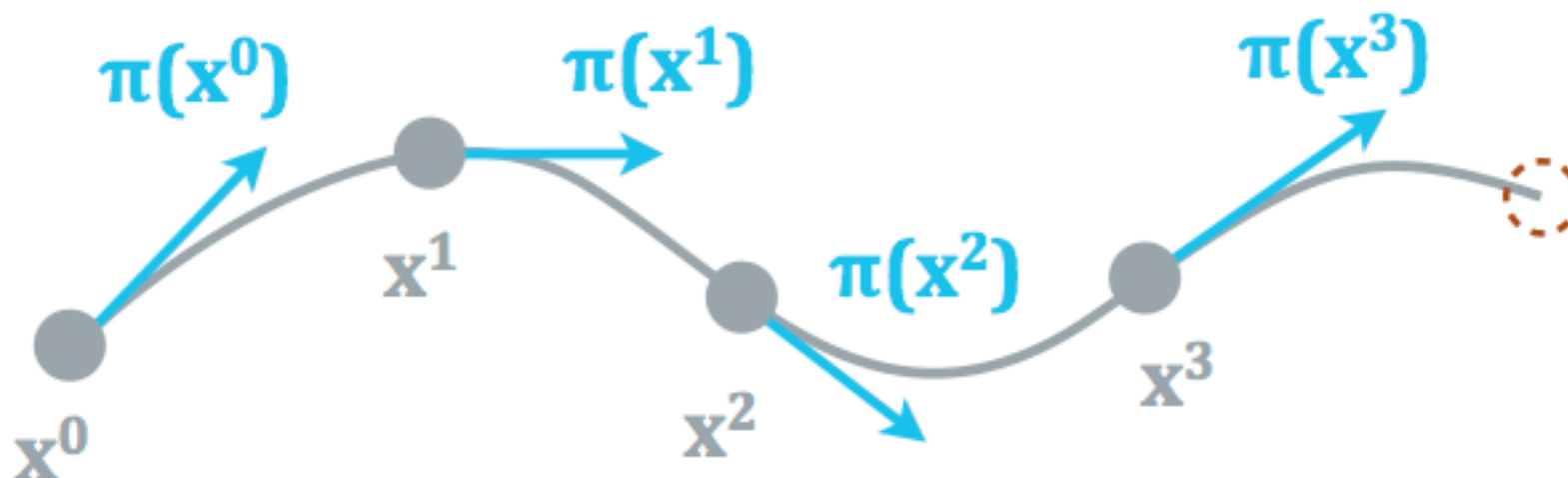

 \mathbf{x}^0



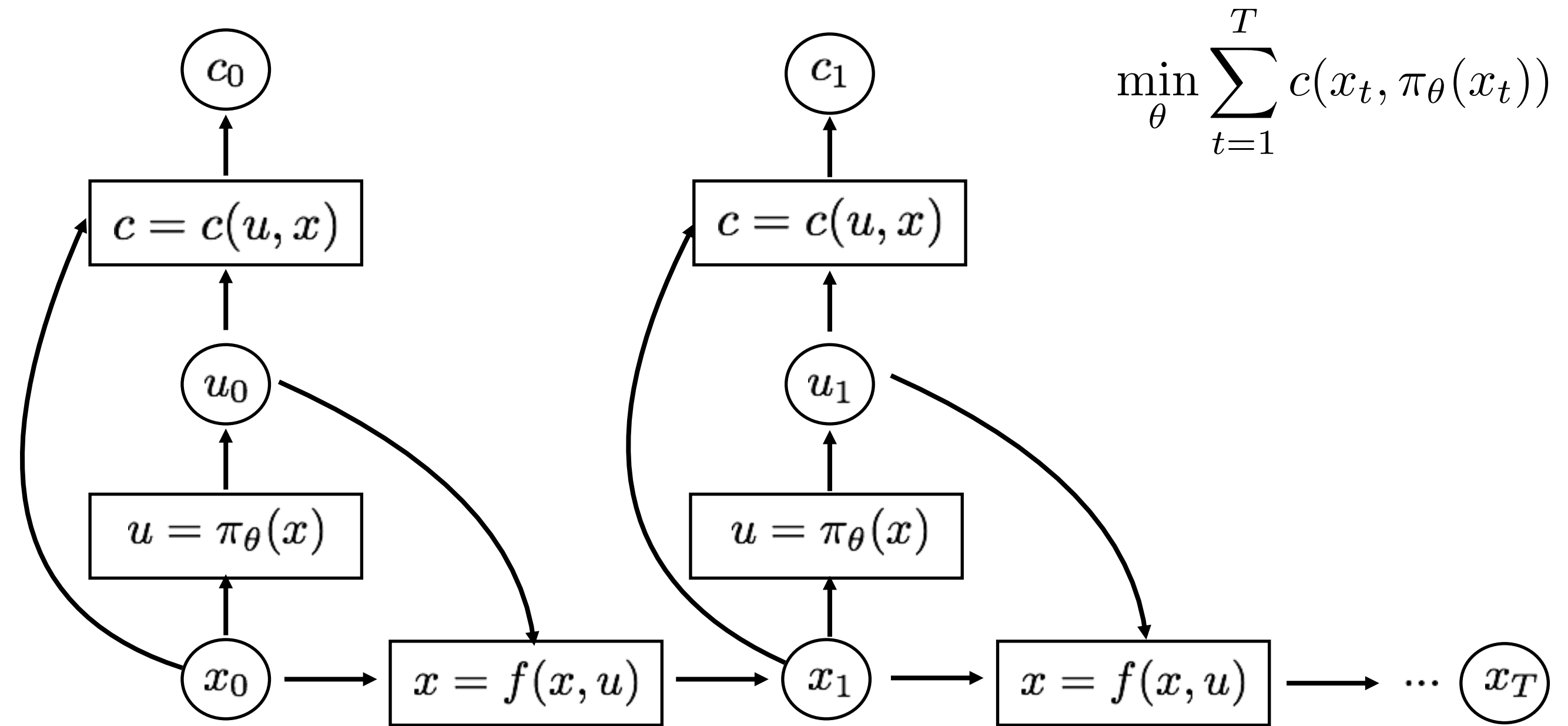
Learning Control Policies

$$\pi_{\theta} : \mathbf{x} \mapsto \mathbf{u}$$

$$\min_{\theta} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \quad \mathbf{x}^{t+1} = f(\mathbf{x}^t, \pi_{\theta}(\mathbf{x}^t))$$

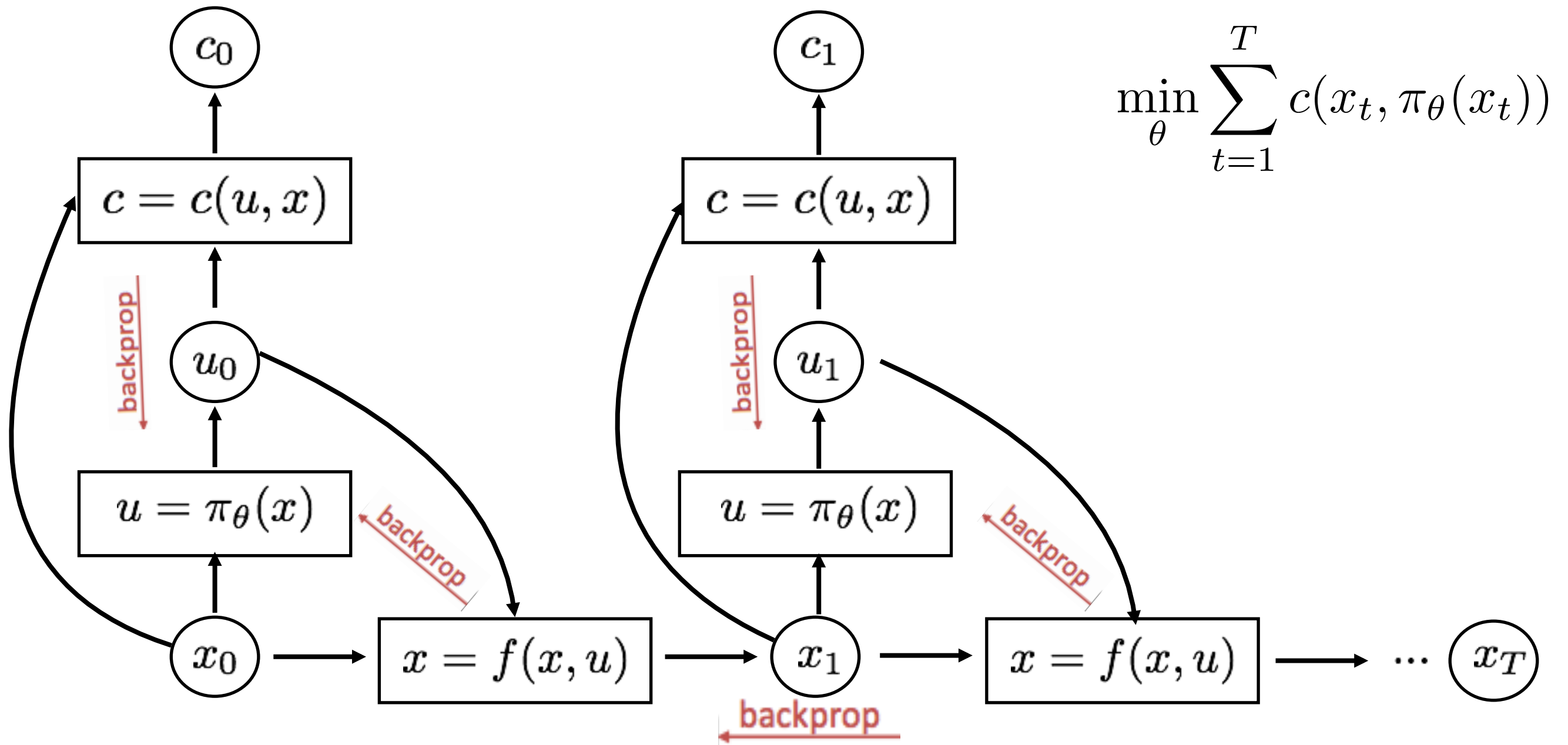


Learning Control Policies through Backpropagation



- So far, dynamics are assumed **known and deterministic**.
- Policy is assumed deterministic.
- We solve for policy parameters θ using back-propagating (through time).

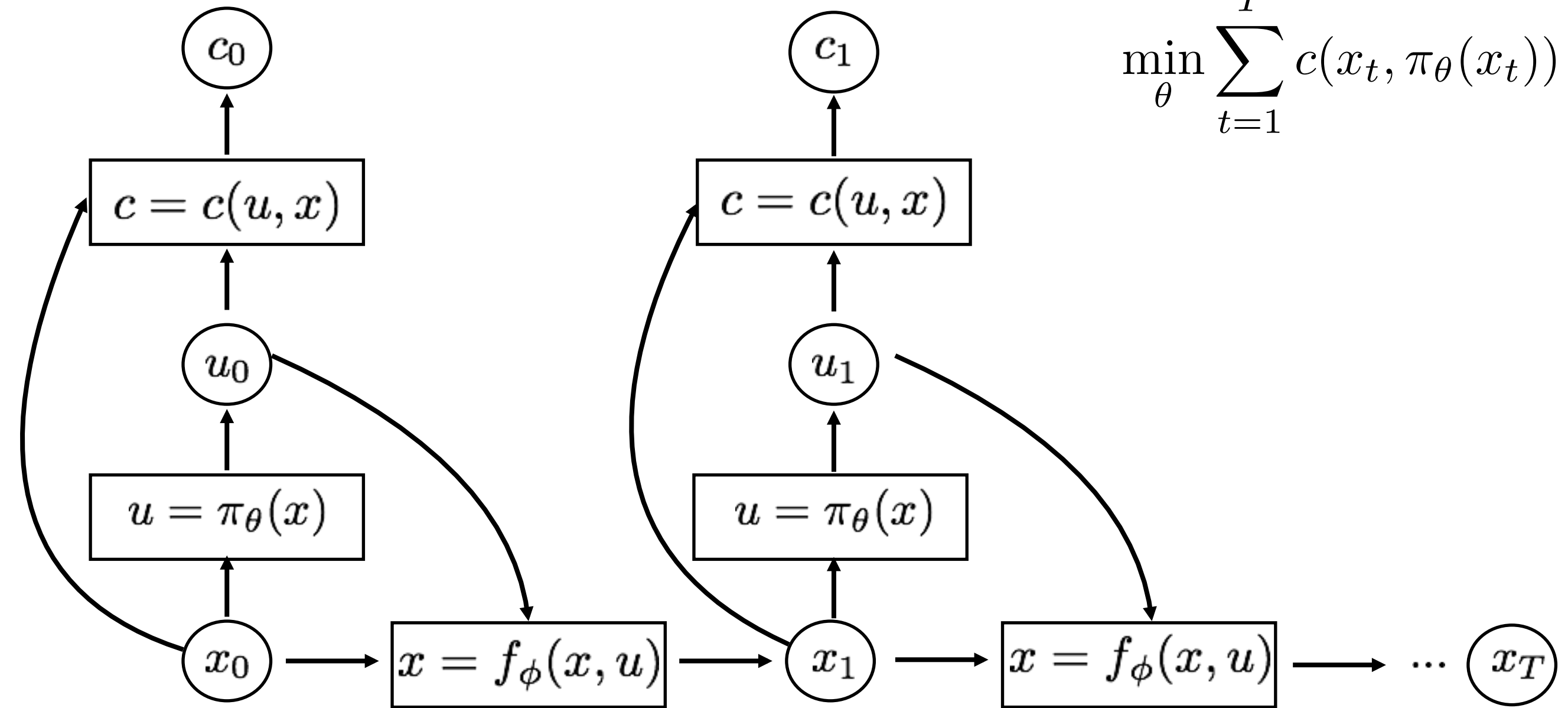
Learning Control Policies through Backpropagation



- So far, dynamics are assumed **known and deterministic**.
- Policy is assumed deterministic.
- We solve for policy parameters θ using back-propagating (through time).

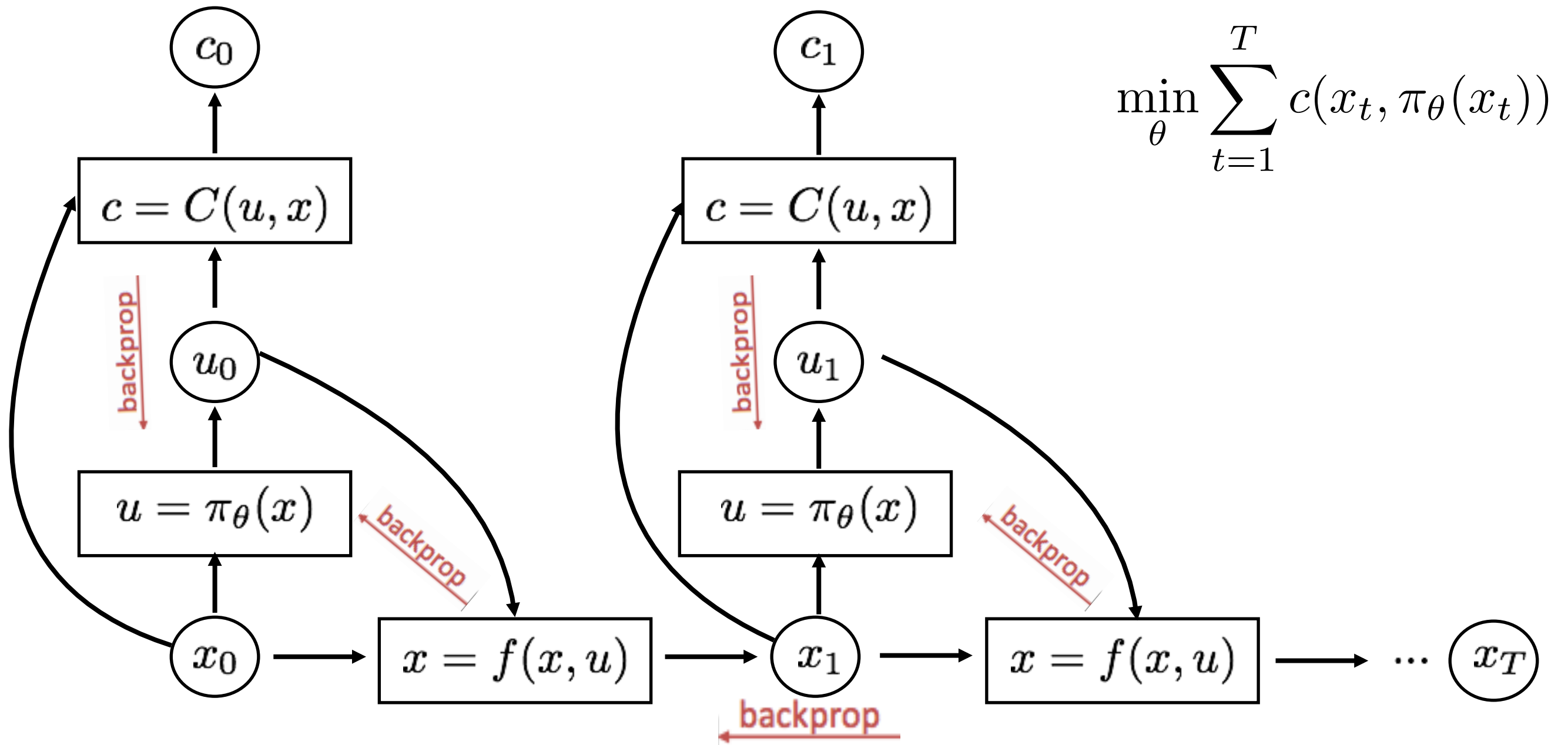
Learning Control Policies through Backpropagation

$$\min_{\theta} \sum_{t=1}^T c(x_t, \pi_{\theta}(x_t))$$



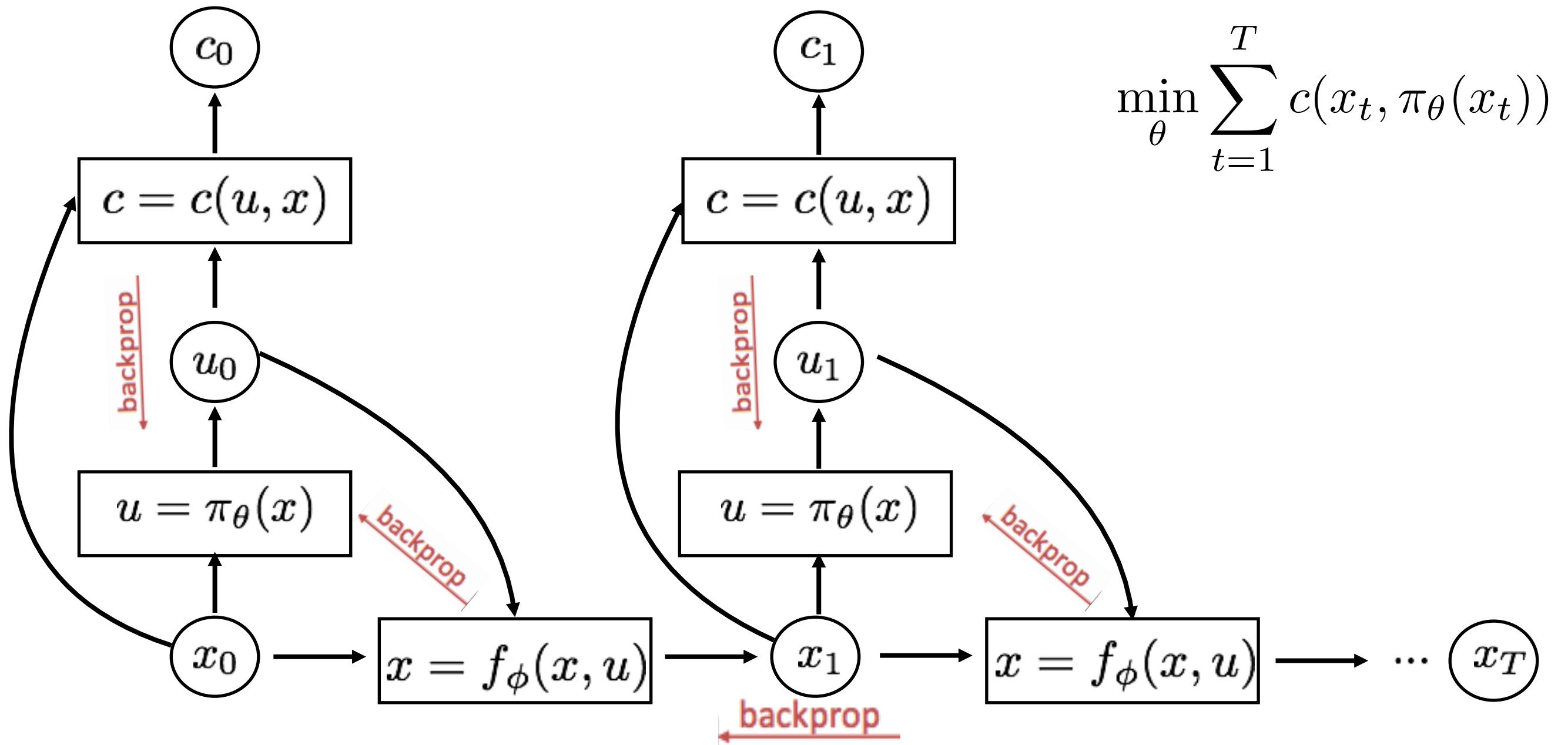
- Dynamics are **unknown**.
- Policy is assumed deterministic.
- We alternate solving for dynamics parameters ϕ (standard regression) and solving for policy parameters θ using back-propagating (through time).

Learning Control Policies through Backpropagation



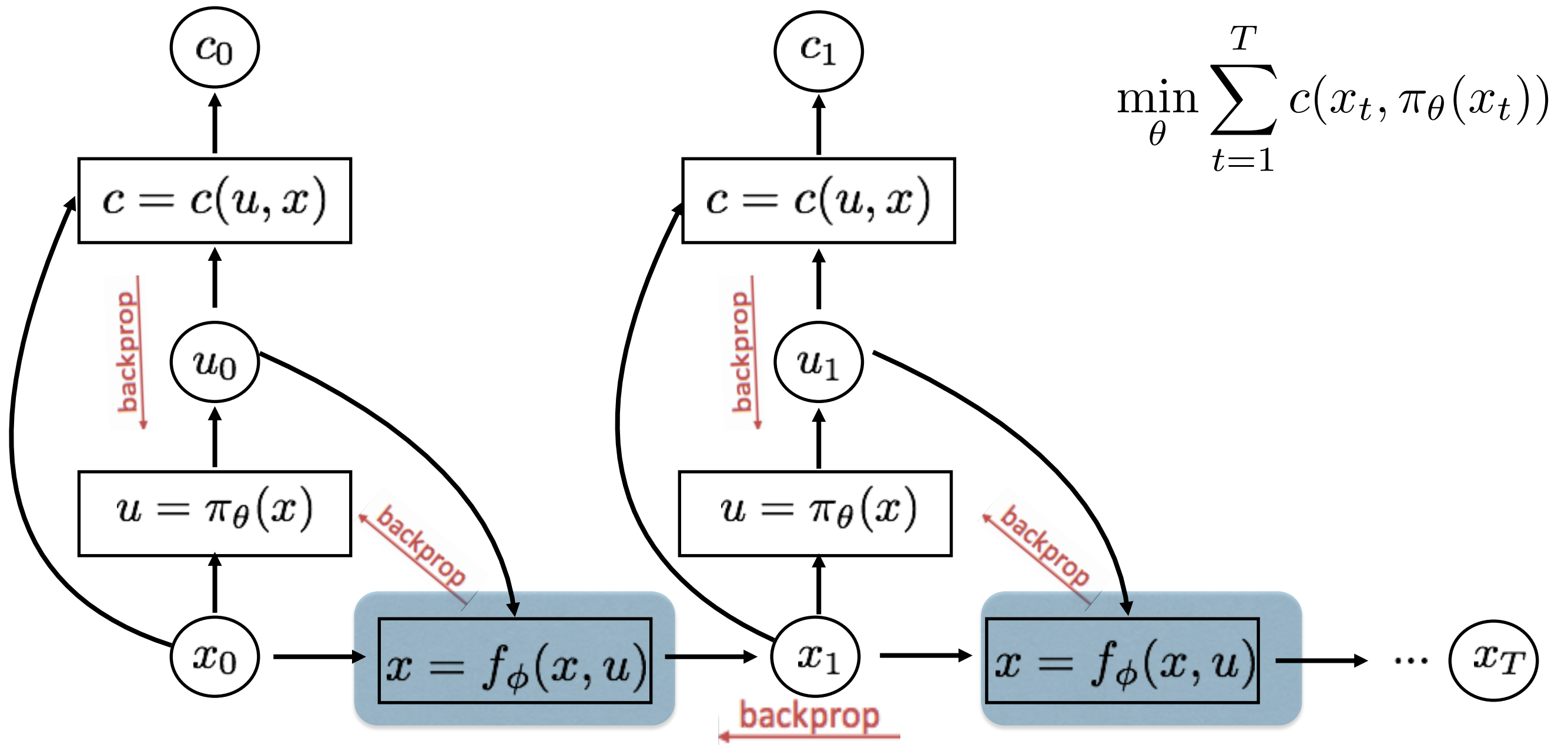
1. run base policy $\pi_0(\mathbf{u}_t|\mathbf{x}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

Learning Control Policies through Backpropagation



1. run base policy $\pi_0(\mathbf{u}_t|\mathbf{x}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$
2. learn dynamics model $f_\phi(x, u)$ to minimize $\sum_i \|f_\phi(x_i, u_i) - x'_i\|^2$

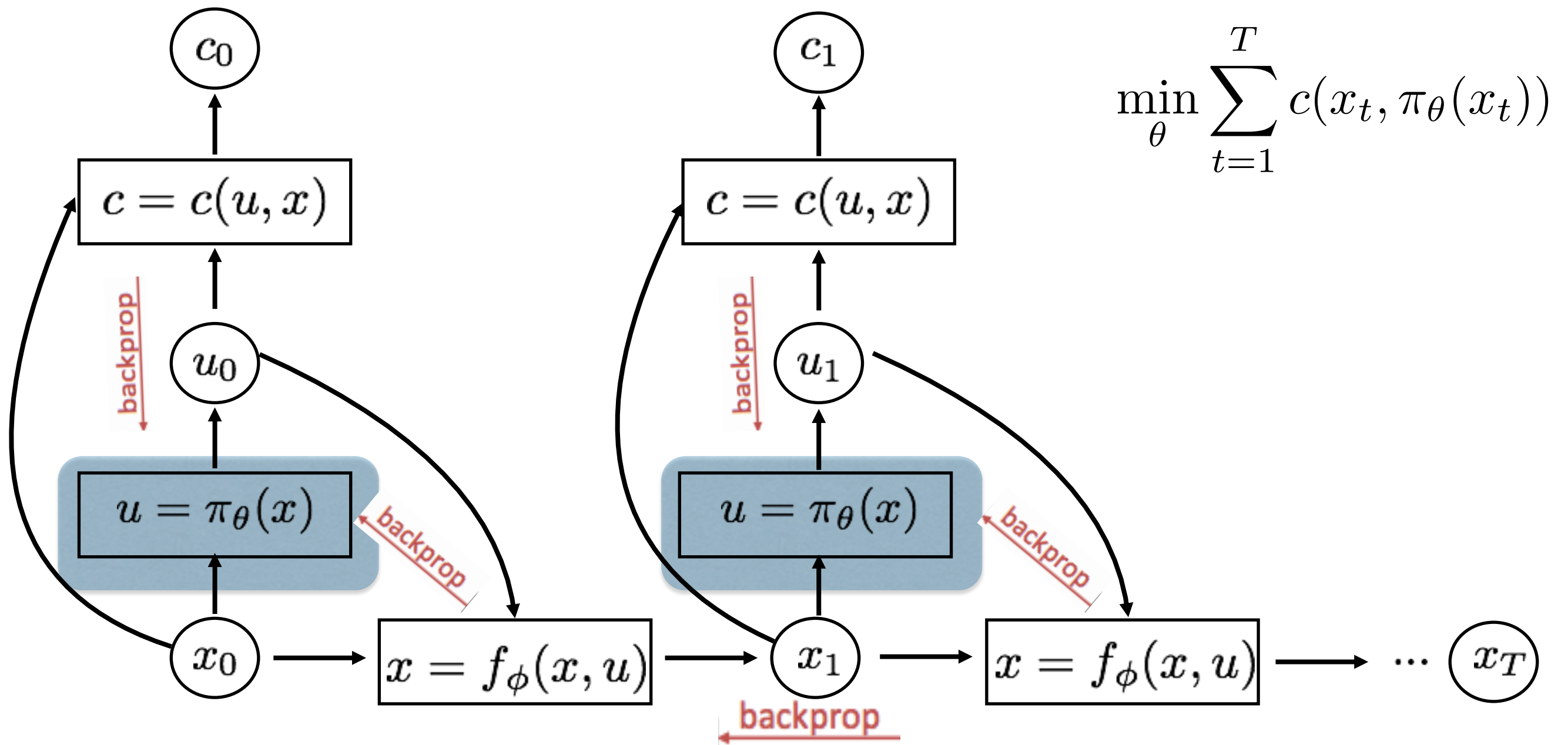
Learning Control Policies through Backpropagation



1. run base policy $\pi_0(\mathbf{u}_t|\mathbf{x}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

2. learn dynamics model $f_\phi(x, u)$ to minimize $\sum_i \|f_\phi(x_i, u_i) - x'_i\|^2$

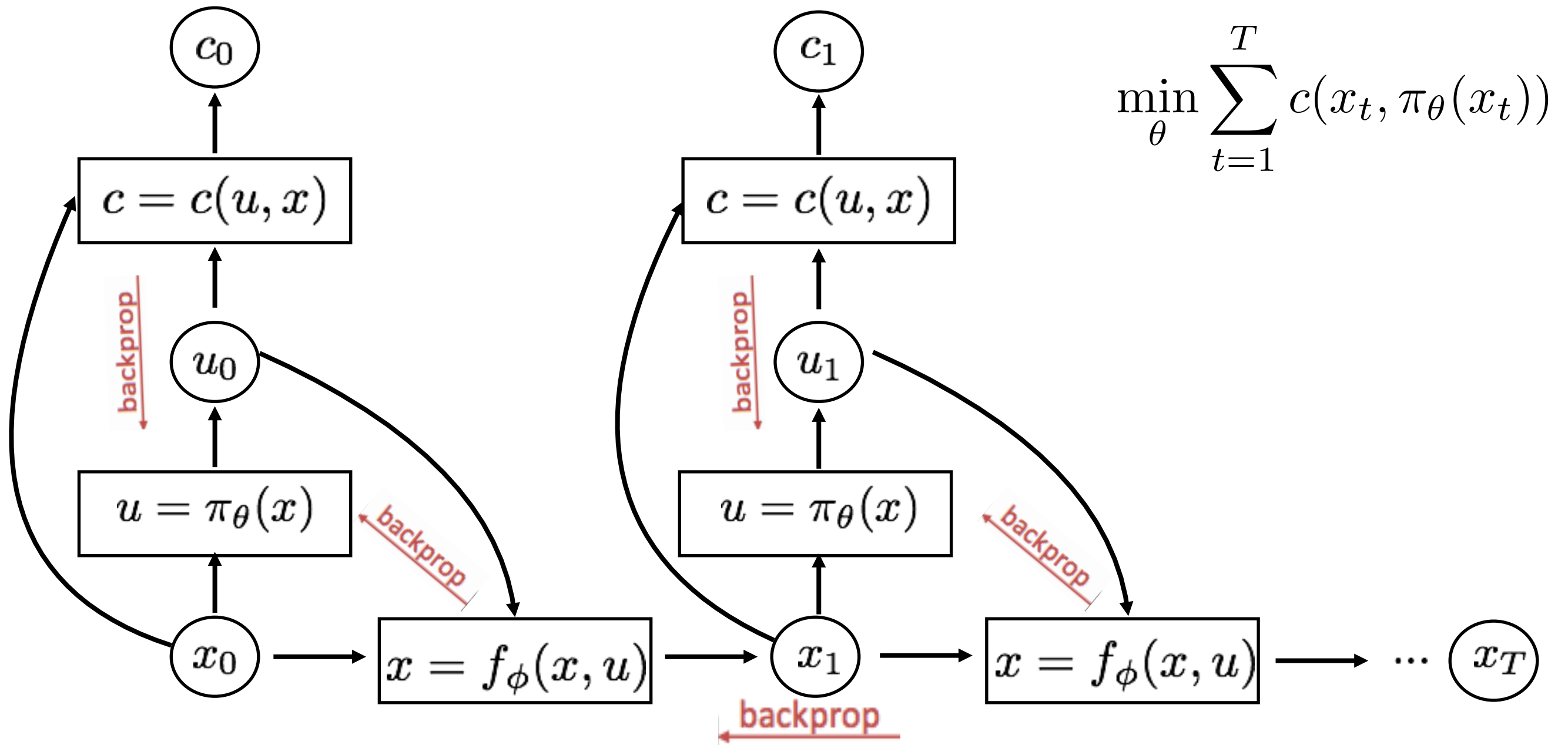
Learning Control Policies through Backpropagation



1. run base policy $\pi_0(\mathbf{u}_t|\mathbf{x}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$
2. learn dynamics model $f_\phi(x, u)$ to minimize $\sum_i ||f_\phi(x_i, u_i) - x'_i||^2$
3. backpropagate through $f_\phi(x, u)$ into the policy to optimize $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$

while dynamics are frozen

Learning Control Policies through Backpropagation



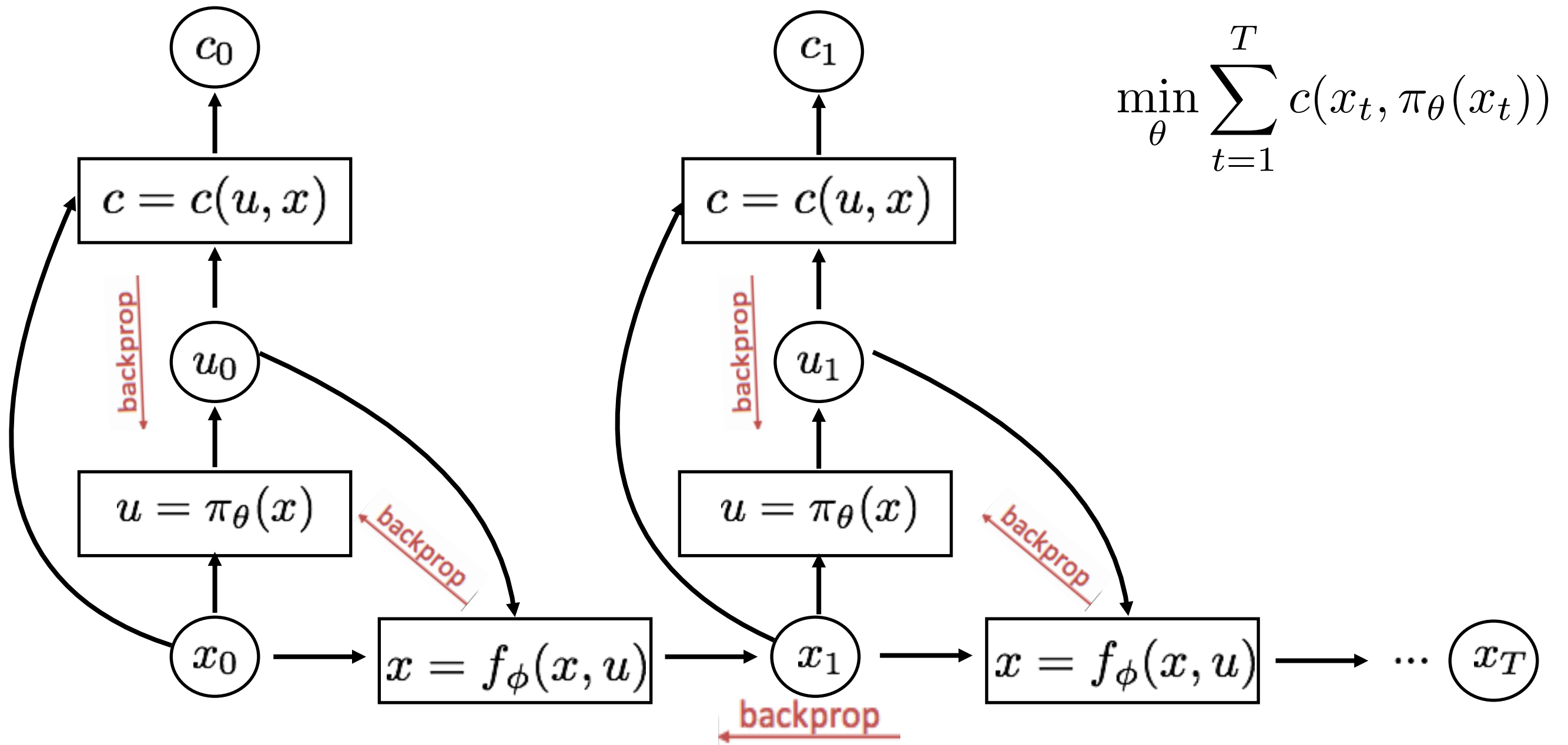
1. run base policy $\pi_0(\mathbf{u}_t|\mathbf{x}_t)$ (e.g., random policy) to collect $\mathcal{D} = \{(\mathbf{x}, \mathbf{u}, \mathbf{x}')_i\}$

2. learn dynamics model $f_\phi(x, u)$ to minimize $\sum_i \|f_\phi(x_i, u_i) - x'_i\|^2$

3. backpropagate through $f_\phi(x, u)$ into the policy to optimize $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$

4. run $\pi_\theta(\mathbf{u}_t|\mathbf{x}_t)$, appending the visited tuples $(\mathbf{x}, \mathbf{u}, \mathbf{x}')$ to \mathcal{D}

Learning Control Policies through Backpropagation



Challenges:

- Poor conditioning
- θ couples actions across all steps \rightarrow no DP
- Chaining inaccurate dynamics naturally leads to errors

Learning Control Policies through Imitation

$$\pi_{\theta} : \mathbf{x} \mapsto \mathbf{u}$$

Learning from Demonstrations:

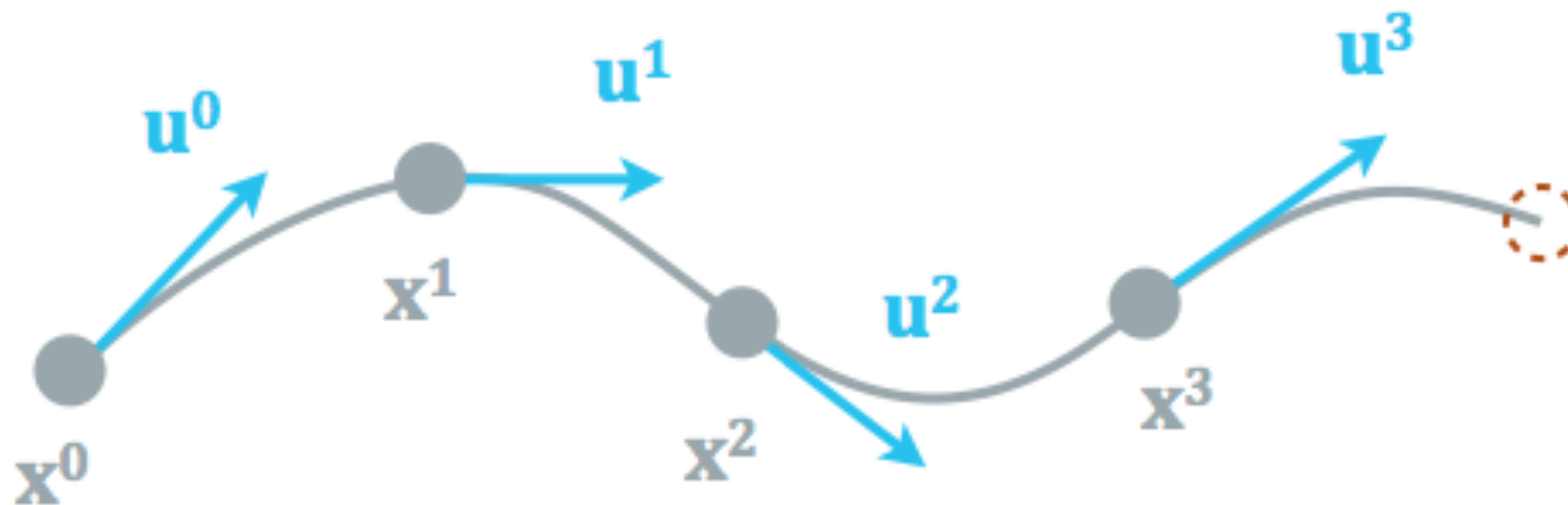
supervised learning

$$\min_{\theta} \sum_i ||\pi_{\theta}(\mathbf{x}^i) - \mathbf{u}^i||^2$$

Training Data

input: \mathbf{x}^i

output: \mathbf{u}^i



Learning Control Policies through Imitation

$$\pi_{\theta} : \mathbf{x} \mapsto \mathbf{u}$$

Learning from Demonstrations:

supervised learning

$$\min_{\theta} \sum_i ||\pi_{\theta}(\mathbf{x}^i) - \mathbf{u}^i||^2$$

Training Data

input: \mathbf{x}^i

output: \mathbf{u}^i

Where does training data come from?

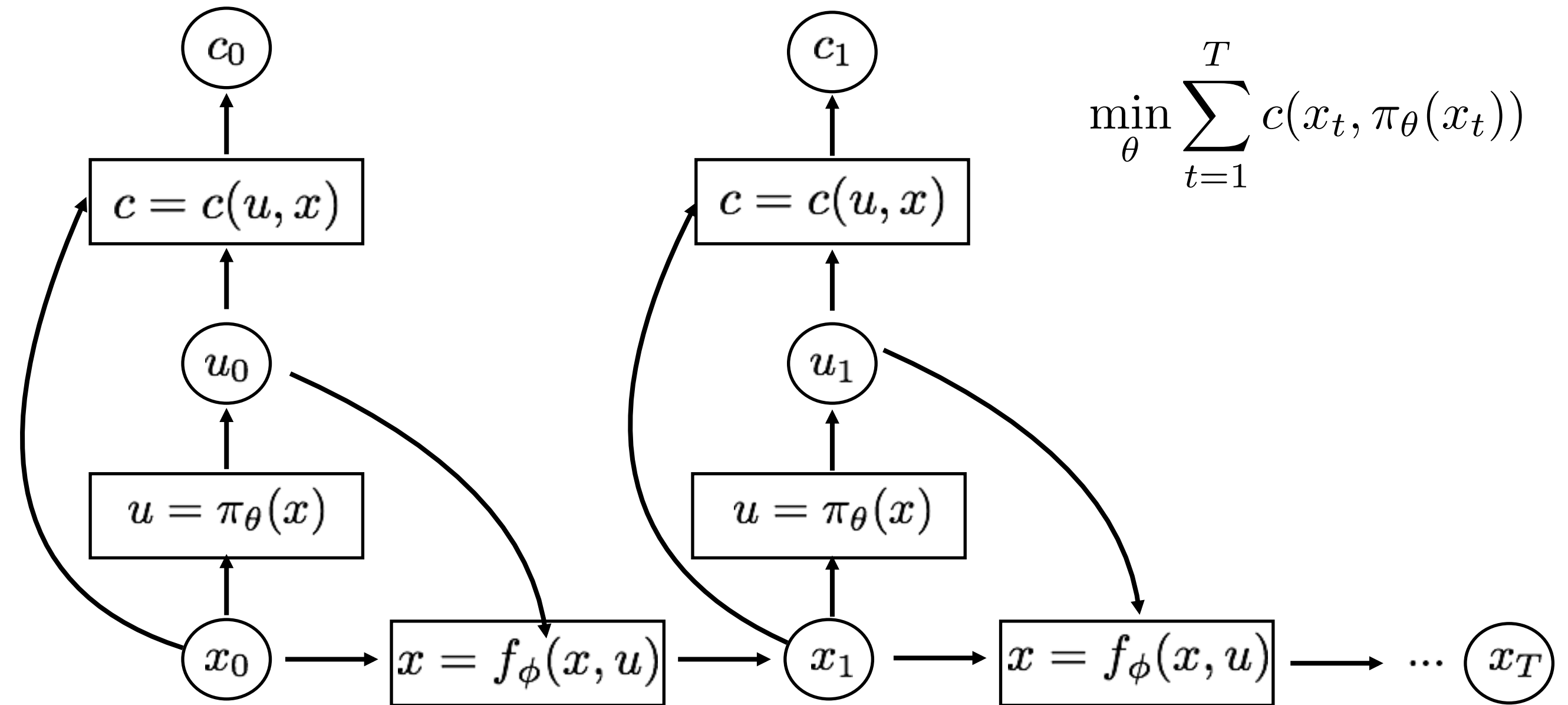
Optimal controllers trained with trajectory optimization

Learning Control Policies through Imitation

- Joint trajectory and policy optimization (last week):

$$\min_{\mathbf{u}_1, \dots, \mathbf{u}_T, \mathbf{x}_1, \dots, \mathbf{x}_T, \theta} \sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t) \text{ s.t. } \mathbf{x}_t = f(\mathbf{x}_{t-1}, \mathbf{u}_{t-1})$$
$$\text{s.t. } \mathbf{u}_t = \pi_{\theta}(\mathbf{x}_t)$$

Learning Control Policies through Backpropagation



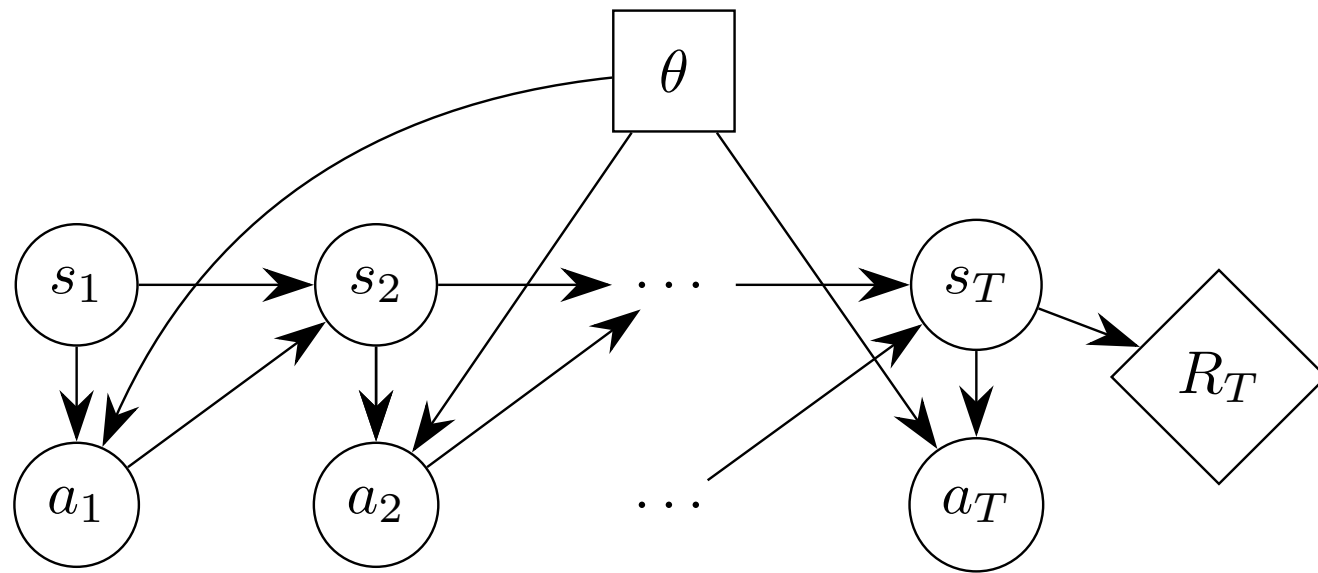
Challenges:

- Backproping through **stochastic** policies and **stochastic** environments
- Avoiding error accumulation through chaining of one step dynamics

This Lecture

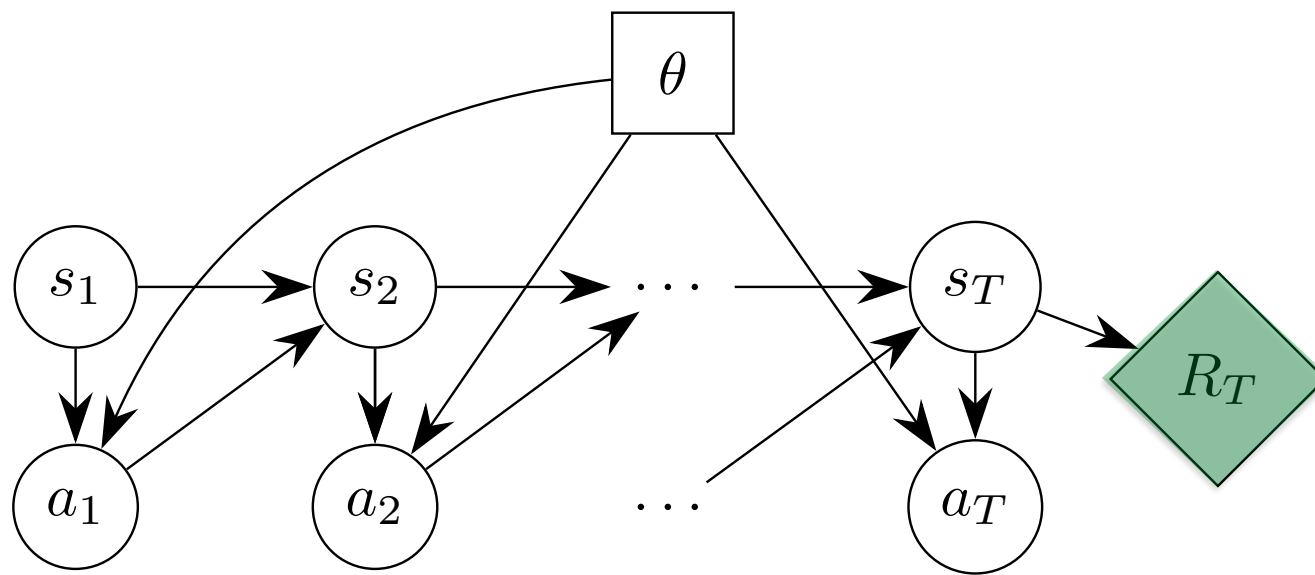
- Backproping through **stochastic** policies and **stochastic** environments
 - **Re-parametrization trick**
- Avoiding error accumulation through chaining of one step dynamics
 - Use function approximation for action and state value functions to predict future returns so that you do not rely on your model for long chaining

Policy optimization

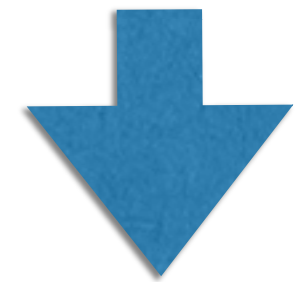


- You get a reward when the jeannie appears

Policy optimization



- You get a reward when the jeannie appears



Policy Optimization

$$\underset{\pi}{\text{maximize}} \mathbb{E}_{\pi} [\textit{expression}]$$

Fixed-horizon episodic: $\sum_{t=0}^{T-1} r_t$

Average-cost: $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} r_t$

Infinite-horizon discounted: $\sum_{t=0}^{\infty} \gamma^t r_t$

Variable-length undiscounted: $\sum_{t=0}^{T_{\text{terminal}}-1} r_t$

Infinite-horizon undiscounted: $\sum_{t=0}^{\infty} r_t$

Episodic Settings

$$s_0 \sim \mu(s_0)$$

$$a_0 \sim \pi(a_0 \mid s_0)$$

$$s_1, r_0 \sim P(s_1, r_0 \mid s_0, a_0)$$

$$a_1 \sim \pi(a_1 \mid s_1)$$

$$s_2, r_1 \sim P(s_2, r_1 \mid s_1, a_1)$$

...

$$a_{T-1} \sim \pi(a_{T-1} \mid s_{T-1})$$

$$s_T, r_{T-1} \sim P(s_T \mid s_{T-1}, a_{T-1})$$

Objective:

maximize $\eta(\pi)$, where

$$\eta(\pi) = E[r_0 + r_1 + \cdots + r_{T-1} \mid \pi]$$

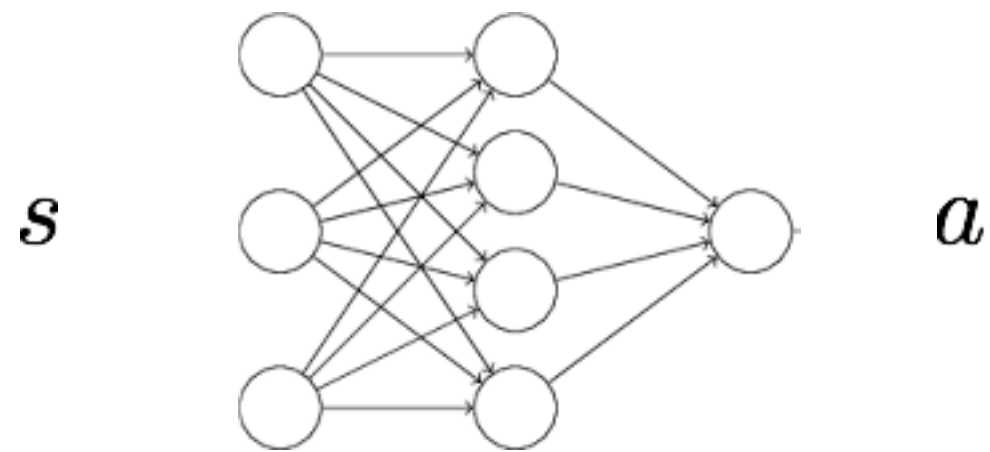
\mathcal{T} : trajectory, a sequence of action states

Parameterized Policies

- A family of policies indexed by parameter vector $\theta \in \mathbb{R}^d$
 - **Deterministic:** $a = \pi(s, \theta)$
 - **Stochastic:** $\pi(a|s, \theta)$
- Analogous to classification or regression with input s , output a
 - Discrete action space: network outputs vector of probabilities
 - Continuous actions space: network outputs mean and diagonal covariance of Gaussian

Parametrized policies

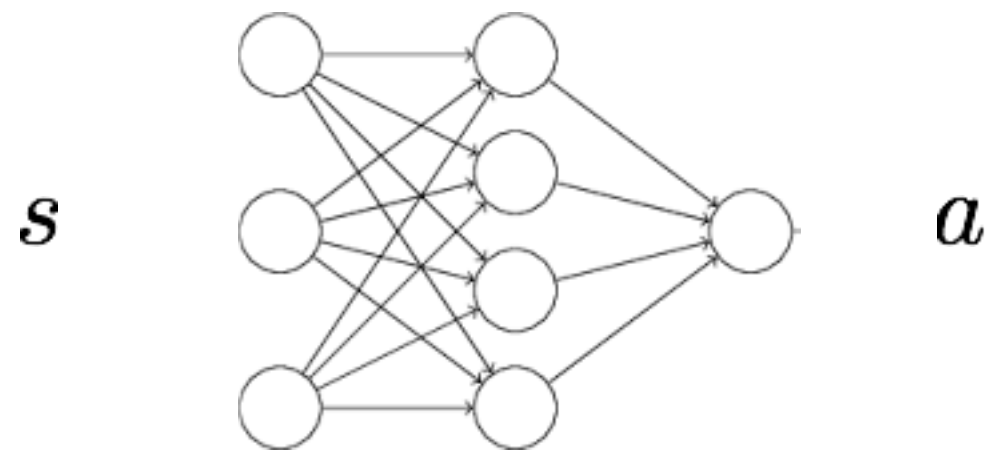
deterministic policy



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

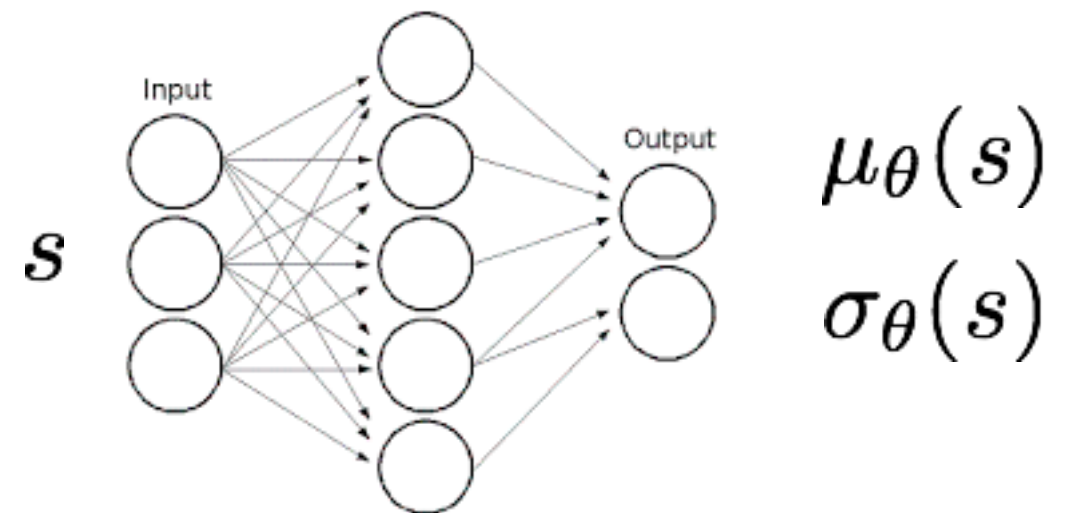
Parametrized policies

deterministic policy



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

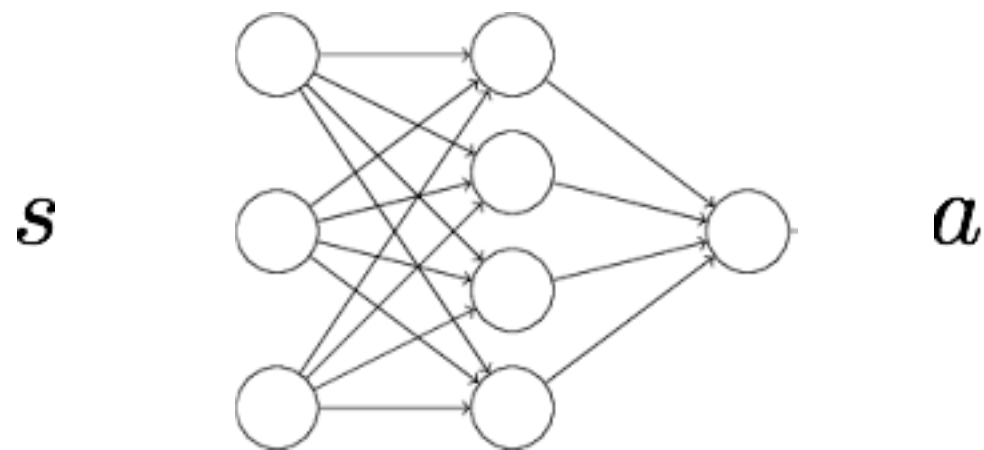
stochastic continuous policy:
usually unimodal gaussian



$$a \sim \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}^2(s))$$

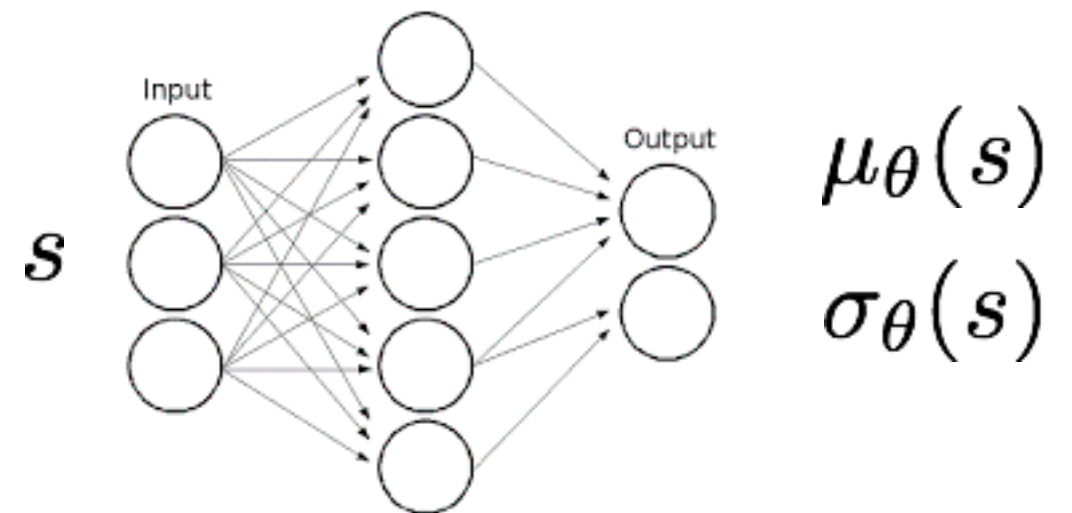
Parametrized policies

deterministic policy



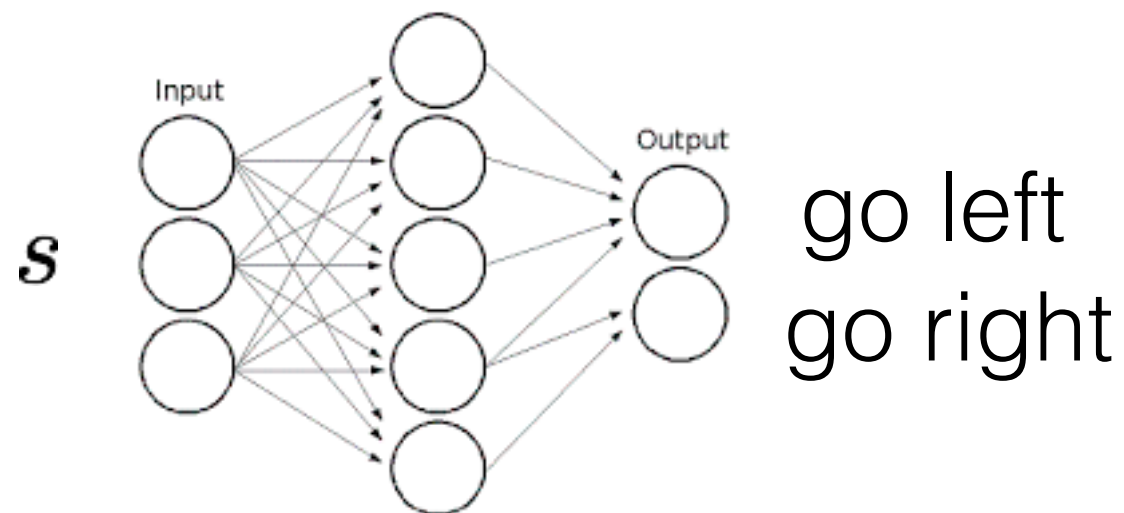
$$a = \pi_{\theta}(s)$$

stochastic continuous policy:
usually unimodal gaussian



$$a \sim \mathcal{N}(\mu_{\theta}(s), \sigma_{\theta}^2(s))$$

discrete action space



How do we compute gradients?

$$\nabla_{\theta} \mathbb{E}[R_{\tau}]$$

- Numerically: finite differencing



How do we compute gradients?

$$\nabla_{\theta} \mathbb{E}[R_{\tau}]$$

- Numerically: finite differencing
- Score function gradient estimator (a.k.a. likelihood ratio gradient estimator)



Likelihood Ratio Policy Gradient

$$J(\theta) = \sum_{\tau} P[\tau; \theta] R(\tau)$$

Taking the gradient w.r.t. θ gives

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \sum_{\tau} P[\tau; \theta] R(\tau)$$

Likelihood Ratio Policy Gradient

$$J(\theta) = \sum_{\tau} P[\tau; \theta] R(\tau)$$

Taking the gradient w.r.t. θ gives

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P[\tau; \theta] R(\tau)\end{aligned}$$

Likelihood Ratio Policy Gradient

$$J(\theta) = \sum_{\tau} P[\tau; \theta] R(\tau)$$

Taking the gradient w.r.t. θ gives

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P[\tau; \theta]} R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)\end{aligned}$$

Likelihood Ratio Policy Gradient

$$J(\theta) = \sum_{\tau} P[\tau; \theta] R(\tau)$$

Taking the gradient w.r.t. θ gives

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P[\tau; \theta] R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \frac{\nabla_{\theta} P(\tau; \theta)}{P[\tau; \theta]} R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) \nabla_{\theta} \log P(\tau; \theta) R(\tau)\end{aligned}$$

Approximate with the empirical estimate for m sample paths under policy

$$\nabla_{\theta} J(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

Decompose Path into States and Actions

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \nabla_{\theta} \log \left[\prod_{t=0}^H \underbrace{P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)})}_{\text{dynamics model}} \cdot \underbrace{\pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{policy}} \right]$$

Decompose Path into States and Actions

$$\begin{aligned}\nabla_{\theta} \log P(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[\prod_{t=0}^H \underbrace{P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)})}_{\text{dynamics model}} \cdot \underbrace{\pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{policy}} \right] \\ &= \nabla_{\theta} \left[\sum_{t=0}^H \log P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right]\end{aligned}$$

Decompose Path into States and Actions

$$\begin{aligned}\nabla_{\theta} \log P(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[\prod_{t=0}^H \underbrace{P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)})}_{\text{dynamics model}} \cdot \underbrace{\pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{policy}} \right] \\ &= \nabla_{\theta} \left[\sum_{t=0}^H \log P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right] \\ &= \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})\end{aligned}$$

Decompose Path into States and Actions

$$\begin{aligned}\nabla_{\theta} \log P(\tau^{(i)}; \theta) &= \nabla_{\theta} \log \left[\prod_{t=0}^H \underbrace{P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)})}_{\text{dynamics model}} \cdot \underbrace{\pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{policy}} \right] \\&= \nabla_{\theta} \left[\sum_{t=0}^H \log P(s_{t+1}^{(i)} | s_t^{(i)}, u_t^{(i)}) + \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \right] \\&= \nabla_{\theta} \sum_{t=0}^H \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)}) \\&= \sum_{t=0}^H \underbrace{\nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{no dynamics model required!!}}\end{aligned}$$

Gaussian Policy

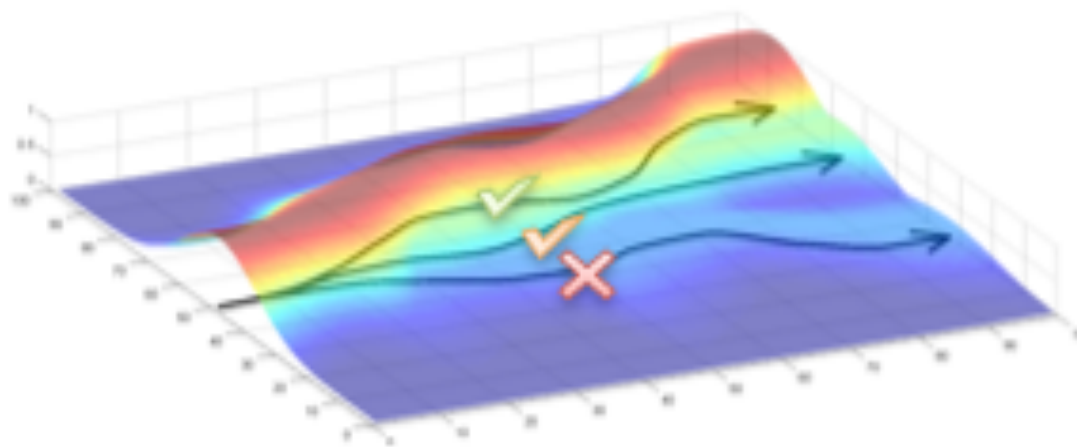
- Variance may be fixed σ^2 , or can also be parametrized
- Policy is Gaussian, $a \sim \mathcal{N}(\mu(s; \theta), \sigma^2)$
- The score function is

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

Likelihood Ratio Gradient: Intuition

$$\nabla_{\theta} J(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

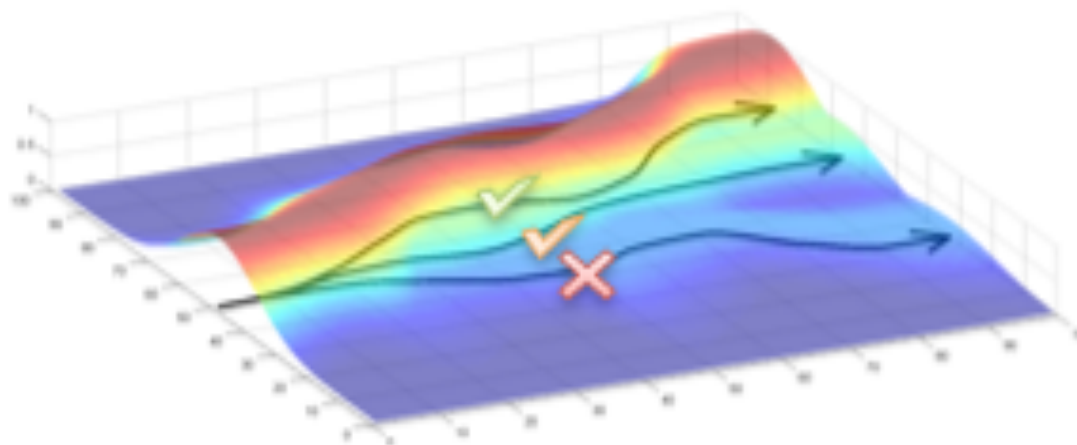
- Gradient tries to:
 - **Increase** probability of paths with positive R
 - **Decrease** probability of paths with negative R



Likelihood Ratio Gradient: Intuition

$$\nabla_{\theta} J(\theta) \approx \hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

- Gradient tries to:
 - **Increase** probability of paths with positive R
 - **Decrease** probability of paths with negative R



- The reward function is a black box and dynamics are not used anywhere
- The world is a black box.

Likelihood Ratio Gradient Estimate

$$\hat{g} = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} \log P(\tau^{(i)}; \theta) R(\tau^{(i)})$$

Here:

$$\nabla_{\theta} \log P(\tau^{(i)}; \theta) = \sum_{t=0}^H \underbrace{\nabla_{\theta} \log \pi_{\theta}(u_t^{(i)} | s_t^{(i)})}_{\text{no dynamics model required!!}}$$

Unbiased means:

$$\mathbb{E}[\hat{g}] = \nabla_{\theta} U(\theta)$$

Reduce Variance using a Critic

- A **critic** provides an estimate of the expectation of the future reward as opposed to a single return sample.
- Use a function approximator for Q function or the advantage function.

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\tau} [R] &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) Q^{\pi}(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t, \theta) A^{\pi}(s_t, a_t) \right]\end{aligned}$$

Reduce Variance using a Critic

Q-function is state-action-value function:

$$Q^{\pi,\gamma}(s, a) = \mathbb{E}_{\pi} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s, a_0 = a]$$

State-value function:

$$\begin{aligned} V^{\pi,\gamma}(s) &= \mathbb{E}_{\pi} [r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \mid s_0 = s] \\ &= \mathbb{E}_{a \sim \pi} [Q^{\pi,\gamma}(s, a)] \end{aligned}$$

Advantage function:

$$A^{\pi,\gamma}(s, a) = Q^{\pi,\gamma}(s, a) - V^{\pi,\gamma}(s)$$

Q Actor-Critic

function QAC

 Initialise s, θ

 Sample $a \sim \pi_\theta$

for each step **do**

 Sample reward $r = \mathcal{R}_s^a$; sample transition $s' \sim \mathcal{P}_{s,\cdot}^a$.

 Sample action $a' \sim \pi_\theta(s', a')$

$\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$

$\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$

$w \leftarrow w + \beta \delta \phi(s, a)$

$a \leftarrow a', s \leftarrow s'$

end for

end function

How do we compute gradients?

$$\nabla_{\theta} \mathbb{E} [R_T]$$

- Numerically: finite differencing
- Score function estimator
 - **Problems:** high variance! In particular, as the policy becomes more and more deterministic, the variance explodes.

Variance in the Gaussian case

$$x \sim \mathcal{N}(\mu(\theta), \sigma(\theta))$$

$$g = \nabla_{\theta} \log p_{\theta}(x) f(x)$$

$$= \frac{(x - \mu(\theta))\mu'(\theta)}{\sigma^2} f(x)$$

Variance in the Gaussian case

$$x \sim \mathcal{N}(\mu(\theta), \sigma(\theta))$$

$$g = \nabla_{\theta} \log p_{\theta}(x) f(x)$$

$$= \frac{(x - \mu(\theta))\mu'(\theta)}{\sigma^2} f(x)$$

Sample x : $x = \mu(\theta) + z\sigma, z \sim \mathcal{N}(0, 1)$

$$\hat{g} = \frac{z\sigma\mu'(\theta)}{\sigma^2} f(\mu(\theta) + z\sigma)$$

$$= \frac{z\mu'(\theta)}{\sigma} f(\mu(\theta) + z\sigma)$$

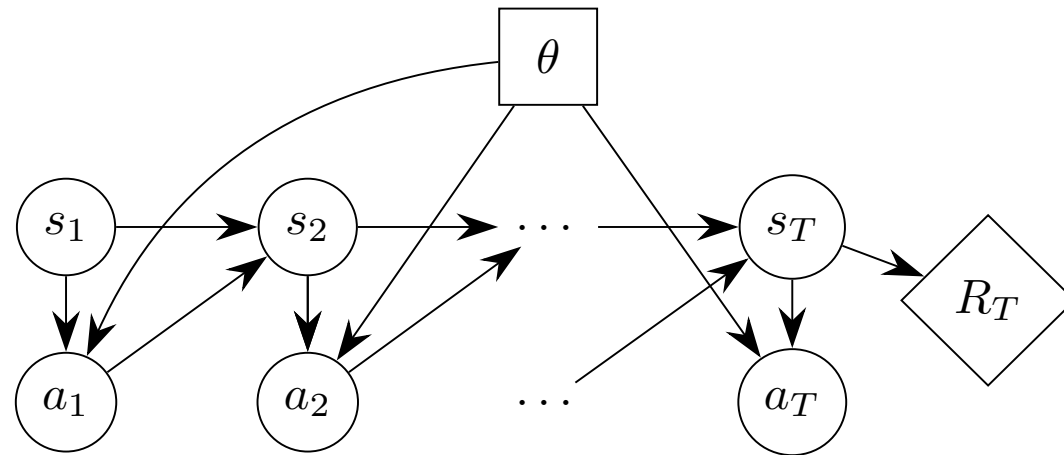
$$\mathbb{V}(\hat{g}) = \mathbb{E}[\hat{g} - g]^2 = \mathbb{E}\left[\frac{z\mu'(\theta)}{\sigma} f(\mu(\theta) + z\sigma) - g\right]^2$$

How do we compute gradients?

$$\nabla_{\theta} \mathbb{E}[R_T]$$

- Numerically: finite differencing
- Score function estimator
- Deep deterministic policy gradients: giving up stochastic policies

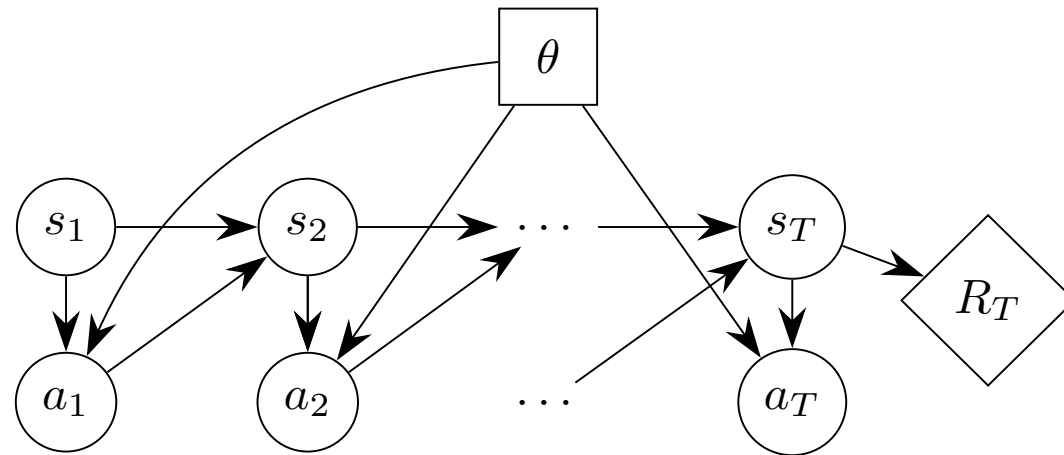
Deep Deterministic Policy Gradients



R_T : the return of a trajectory

$$\frac{d}{d\theta} \mathbb{E} [R_T] = \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right]$$

Deep Deterministic Policy Gradients



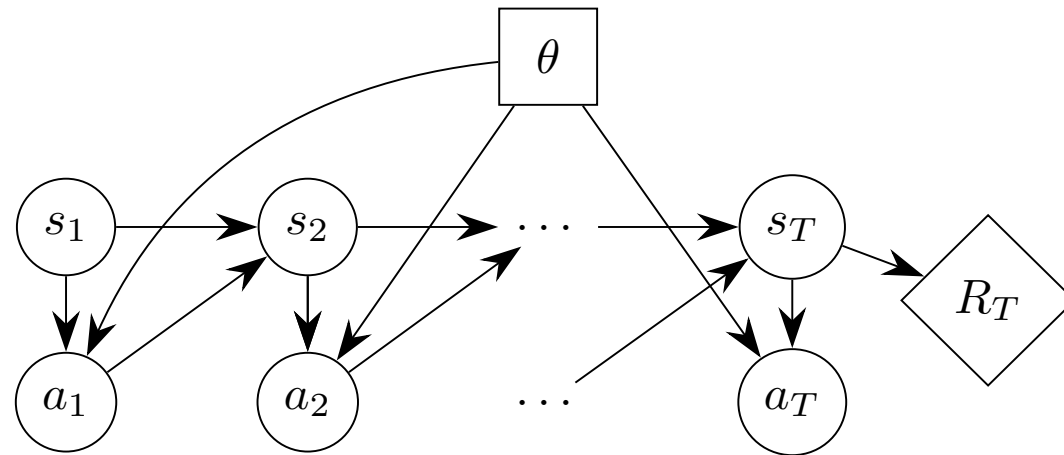
R_T : the return of a trajectory

This expectation refers to the actions after time t

$$\frac{d}{d\theta} \mathbb{E} [R_T] = \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{da_t} \mathbb{E} [R_T | a_t] \frac{da_t}{d\theta} \right]$$

An arrow points from the text "This expectation refers to the actions after time t " to the term $\mathbb{E} [R_T | a_t]$ in the equation.

Deep Deterministic Policy Gradients



R_T : the return of a trajectory

$$\begin{aligned} \frac{d}{d\theta} \mathbb{E} [R_T] &= \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{da_t} \mathbb{E} [R_T | a_t] \frac{da_t}{d\theta} \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{d\theta} Q(s_t, \pi(s_t, z_t; \theta)) \right] \end{aligned}$$

Remember: Q learning

- Definition

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_{i \geq t}, s_{i > t} \sim E, a_{i > t} \sim \pi} [R_t | s_t, a_t]$$

Remember: Q learning

- Definition

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_{i \geq t}, s_{i > t} \sim E, a_{i > t} \sim \pi} [R_t | s_t, a_t]$$

- Bellman equation

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi} [Q^\pi(s_{t+1}, a_{t+1})]]$$

- Using a deterministic policy

$$Q^\mu(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma Q^\mu(s_{t+1}, \mu(s_{t+1}))]$$

Remember: Q learning

- Definition

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_{i \geq t}, s_{i > t} \sim E, a_{i > t} \sim \pi} [R_t | s_t, a_t]$$

- Bellman equation

$$Q^\pi(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma \mathbb{E}_{a_{t+1} \sim \pi} [Q^\pi(s_{t+1}, a_{t+1})]]$$

- Using a deterministic policy

$$Q^\mu(s_t, a_t) = \mathbb{E}_{r_t, s_{t+1} \sim E} [r(s_t, a_t) + \gamma Q^\mu(s_{t+1}, \mu(s_{t+1}))]$$

- Deep Q learning:

$$L(\theta^Q) = \mathbb{E}_{s_t \sim \rho^\beta, a_t \sim \beta, r_t \sim E} [(Q(s_t, a_t | \theta^Q) - y_t)^2]$$

$$y_t = r(s_t, a_t) + \gamma Q(s_{t+1}, \mu(s_{t+1}) | \theta^Q)$$

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

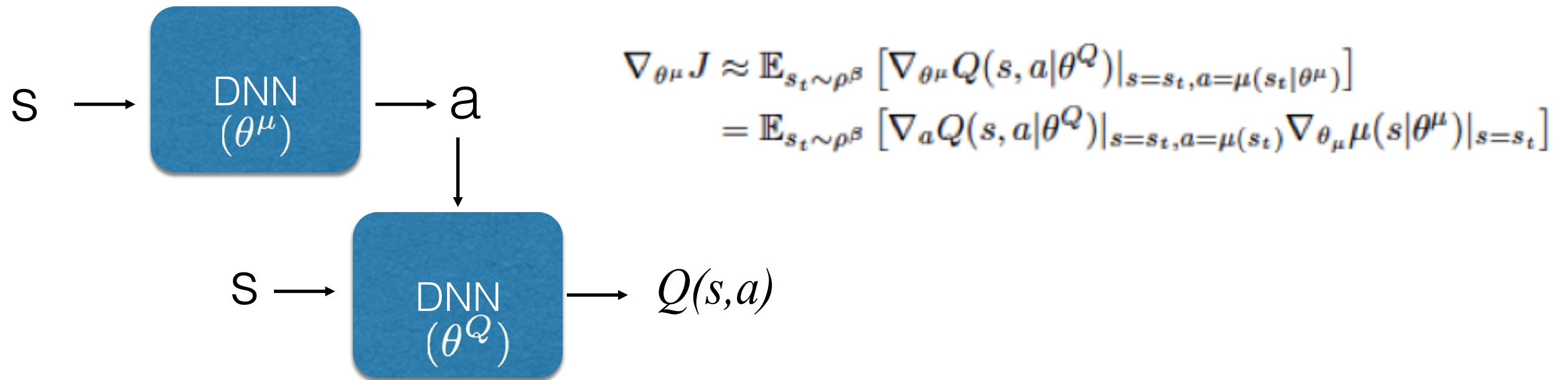
Q learning in continuous action space

- This optimization takes too long for continuous action spaces and has to be performed at every iteration:

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

Q learning in continuous action space

- Instead of parametrizing Q let's also parametrize the policy $a = \mu(\theta)$



How do we compute gradients?

$$\nabla_{\theta} \mathbb{E}[R_T]$$

- Numerically: finite differencing
- Score function estimator
- Deep deterministic policy gradients: giving up stochastic policies
- Pathwise derivatives

Pathwise derivatives for Gaussian samples

- Consider normally distributed variable y

$$p(y|x) = \mathcal{N}(y|\mu(x), \sigma^2(x))$$

Pathwise derivatives for Gaussian samples

- Consider normally distributed variable y

$$p(y|x) = \mathcal{N}(y|\mu(x), \sigma^2(x))$$

- Reparametrization:

$$y = \mu(x) + \sigma(x)\xi, \text{ where } \xi \sim \mathcal{N}(0, 1)$$

Pathwise derivatives for Gaussian samples

- Consider normally distributed variable y

$$p(y|x) = \mathcal{N}(y|\mu(x), \sigma^2(x))$$

- Reparametrization:

$$y = \mu(x) + \sigma(x)\xi, \text{ where } \xi \sim \mathcal{N}(0, 1)$$

- Sampling: Sample ξ and then deterministically generate $\mathbf{y} = \mathbf{f}(\mathbf{x}, \xi)$

$$\mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \int \mathbf{g}(\mathbf{f}(\mathbf{x}, \xi)) \rho(\xi) d\xi$$

$$\nabla_{\mathbf{x}} \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \mathbb{E}_{\rho(\xi)} \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \approx \frac{1}{M} \sum_{i=1}^M \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \big|_{\xi=\xi_i}$$

Pathwise derivatives for Gaussian samples

- Consider normally distributed variable y

$$p(y|x) = \mathcal{N}(y|\mu(x), \sigma^2(x))$$

- Reparametrization:

$$y = \mu(x) + \sigma(x)\xi, \text{ where } \xi \sim \mathcal{N}(0, 1)$$

- Sampling: Sample ξ and then deterministically generate $\mathbf{y} = \mathbf{f}(\mathbf{x}, \xi)$

$$\mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \int \mathbf{g}(\mathbf{f}(\mathbf{x}, \xi)) \rho(\xi) d\xi$$

$$\nabla_{\mathbf{x}} \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \mathbb{E}_{\rho(\xi)} \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \approx \frac{1}{M} \sum_{i=1}^M \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \big|_{\xi=\xi_i}$$

- Compare to the score function gradient estimator: $\frac{1}{M} \sum_{i=1}^M \nabla_{\mathbf{x}} \log(p(\mathbf{y}|\mathbf{x})) \mathbf{g}(\mathbf{y})$

Pathwise derivatives for Gaussian samples

Sampling: sample ξ and then deterministically generate y : $\mathbf{y} = \mathbf{f}(\mathbf{x}, \xi)$

$$\mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \int \mathbf{g}(\mathbf{f}(\mathbf{x}, \xi)) \rho(\xi) d\xi$$

$$\nabla_{\mathbf{x}} \mathbb{E}_{p(\mathbf{y}|\mathbf{x})} \mathbf{g}(\mathbf{y}) = \mathbb{E}_{\rho(\xi)} \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \approx \frac{1}{M} \sum_{i=1}^M \mathbf{g}_{\mathbf{y}} \mathbf{f}_{\mathbf{x}} \big|_{\xi=\xi_i}$$

Compare to the score
function gradient estimator:

$$\frac{1}{M} \sum_{i=1}^M \nabla_{\mathbf{x}} \log(p(\mathbf{y}|\mathbf{x})) \mathbf{g}(\mathbf{y})$$

- The pathwise derivative **makes use of the gradient of g** .
- Of course, that assumes we know the function g (our reward function) and how it is related to our actions.

Pathwise derivative for Gaussian Policies

- Gaussian Policies:

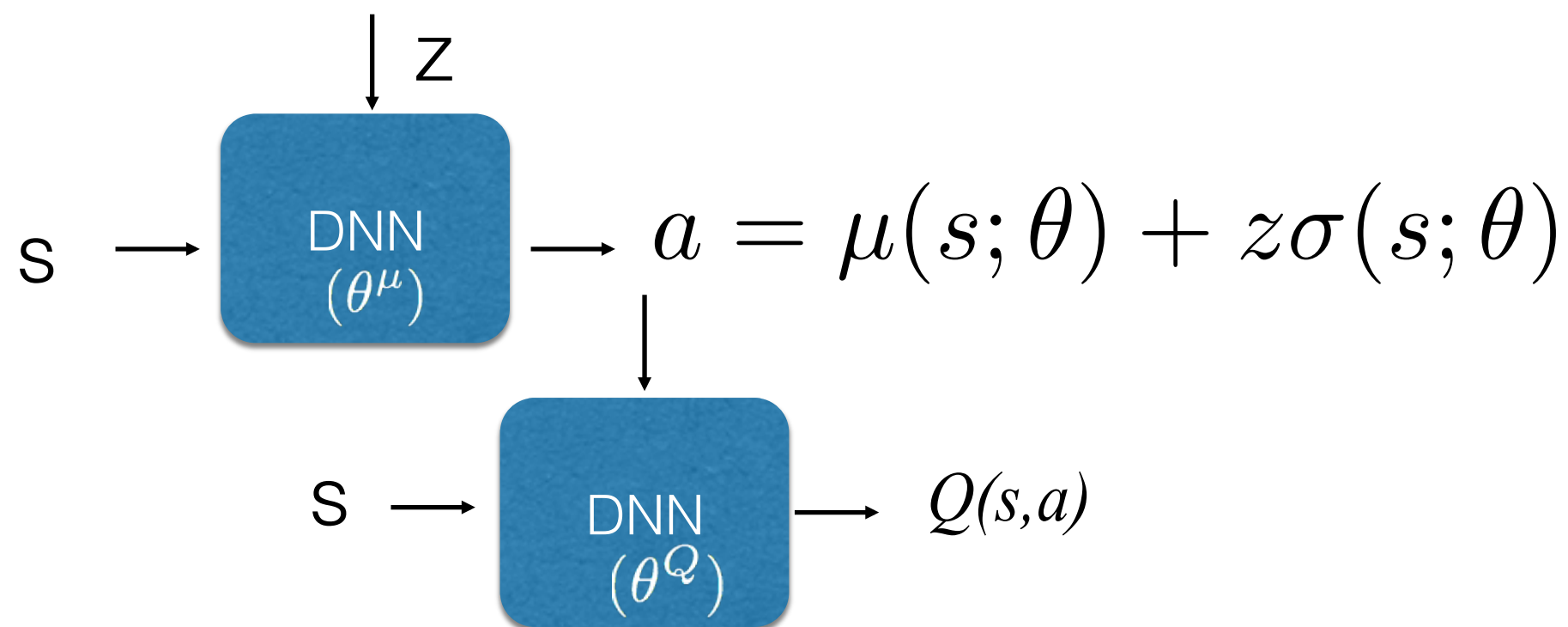
$$a = \mu(s, \theta) + z^* \sigma(s, \theta)$$
$$\frac{da}{d\theta} = \frac{d\mu(s, \theta)}{d\theta} + z \frac{d\sigma(s, \theta)}{d\theta}$$
$$\nabla_{\theta} \mathbb{E}_z (R(a(\theta, z)))$$
$$\mathbb{E}_z \left(R'(a(\theta, z)) \frac{da(\theta, z)}{d\theta} \right)$$

- R should be **known and differentiable**
- To propagate for more than 1 step, dynamics should be known and differentiable

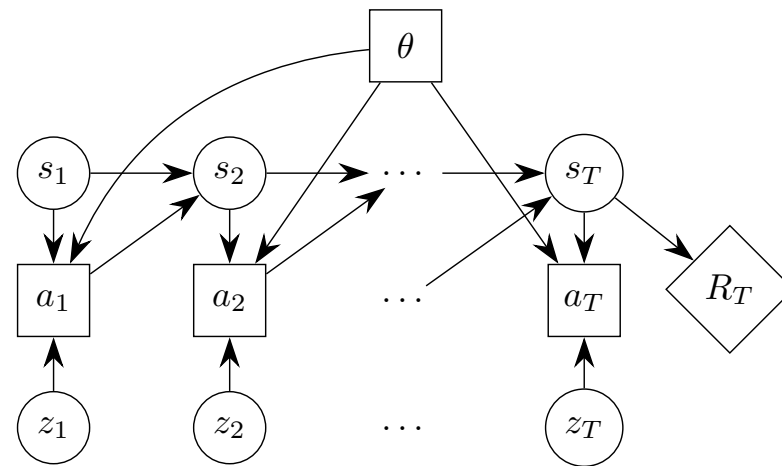
Q learning in continuous **stochastic** action space

- Now use **stochastic policies** using the reparametrization trick!

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



Q learning in continuous **stochastic** action space



$$\begin{aligned} \frac{d}{d\theta} \mathbb{E} [R_T] &= \mathbb{E} \left[\sum_{t=1}^T \frac{dR_T}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{da_t} \mathbb{E} [R_T \mid a_t] \frac{da_t}{d\theta} \right] \\ &= \mathbb{E} \left[\sum_{t=1}^T \frac{dQ(s_t, a_t)}{da_t} \frac{da_t}{d\theta} \right] = \mathbb{E} \left[\sum_{t=1}^T \frac{d}{d\theta} Q(s_t, \pi(s_t, z_t; \theta)) \right] \end{aligned}$$

Q learning in continuous **stochastic** action space

SVG(0)

Learn Q_ϕ to approximate $Q^{\pi,\gamma}$, and use it to compute gradient estimates.

Pseudocode:

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

 Execute policy π_θ to collect T timesteps of data

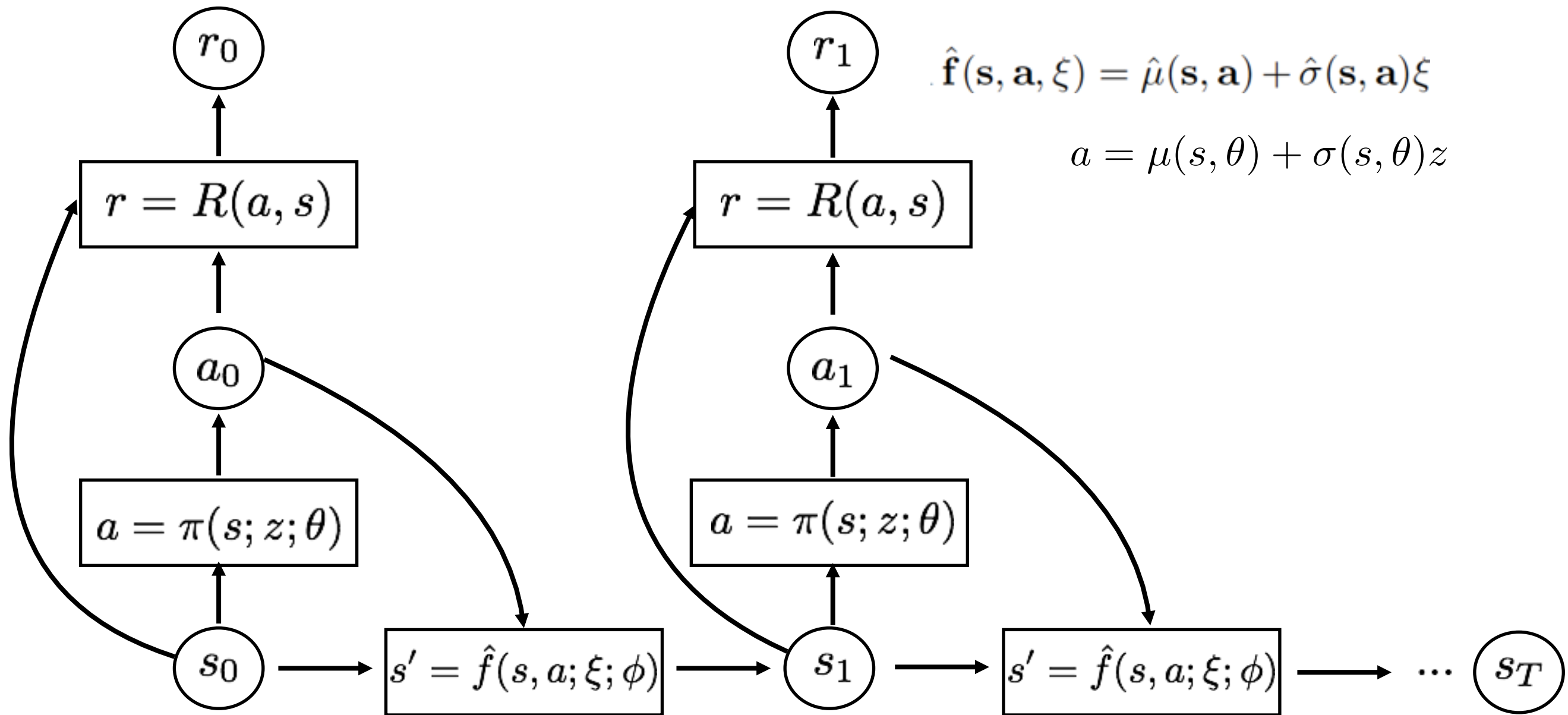
 Update π_θ using $g \propto \nabla_\theta \sum_{t=1}^T Q(s_t, \pi(s_t, z_t; \theta))$

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

end for

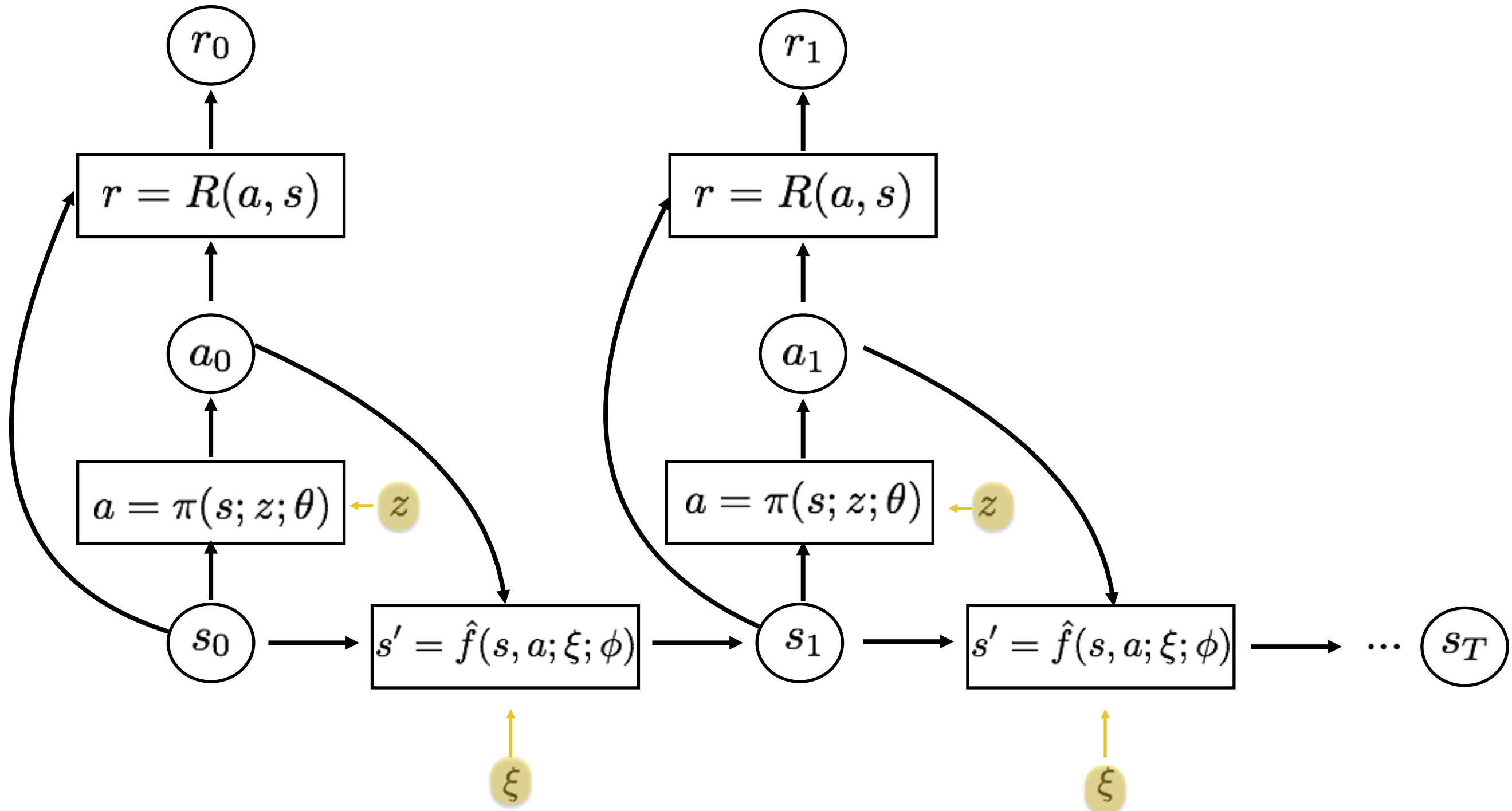
End-to-end model based RL

- Reparametrization trick for both policies and dynamics



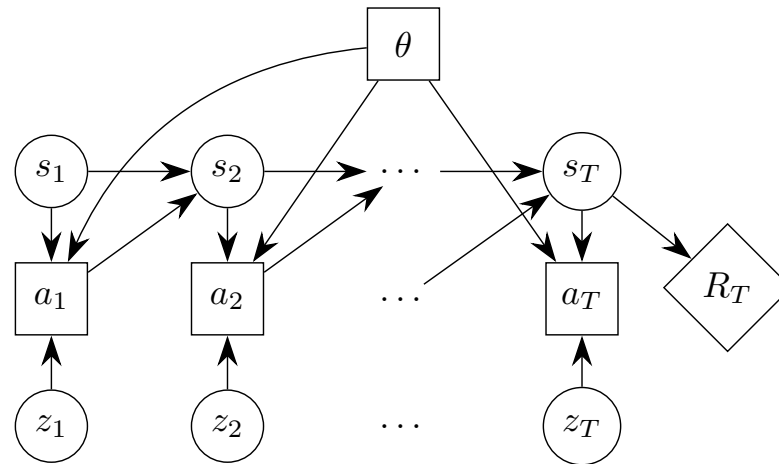
End-to-end model based RL

- Reparametrization trick for both policies and dynamics



End-to-end model based RL

SVG(∞)

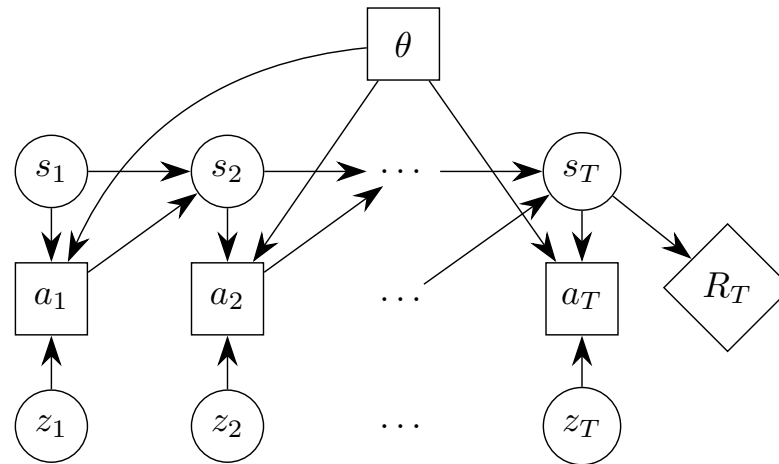


- Just learn dynamics model f
- Given whole trajectory, infer all noise variables
- Given transition (s_t, a_t, s_{t+1}) , infer $\zeta_t = s_{t+1} - f(s_t, a_t)$
- Freeze all policy and dynamics noise, differentiate through entire deterministic computation graph

R should be **known and differentiable**

End-to-end model based RL

SVG(1)



- Instead of learning Q , we learn
 - State-value function $V \approx V^{\pi, \gamma}$
 - Dynamics model f , approximating $s_{t+1} = f(s_t, a_t) + \zeta_t$
- Given transition (s_t, a_t, s_{t+1}) infer $\zeta_t = s_{t+1} - f(s_t, a_t)$
- $Q(s_t, a_t) = \mathbb{E}[r_t + \gamma V(s_{t+1})] = \mathbb{E}[r_t + \gamma V(f(s_t, a_t) + \zeta_t)]$
 $a_t = \pi(s_t, \theta, \zeta_t)$

Re-parametrization trick for categorical distributions

- Consider variable y following the K categorical distribution:

$$y_k \sim \frac{\exp((\log p_k)/\tau)}{\sum_{j=0}^K \exp((\log p_j)/\tau)}$$

Re-parametrization trick for categorical distributions

- Consider variable y following the K categorical distribution:

$$y_k \sim \frac{\exp((\log p_k)/\tau)}{\sum_{j=0}^K \exp((\log p_j)/\tau)}$$

- Reparametrization:

$$y_k \sim G(\log p) = \frac{\exp((\log p_k + \varepsilon)/\tau)}{\sum_{j=0}^K \exp((\log p_j + \varepsilon)/\tau)}, \quad \varepsilon = -\log(-\log(u)), u \sim \mathcal{U}[0, 1]$$

- Sampling: sample u and then sample from $G(\log p)$ to generate y_k

Re-parametrization trick for categorical distributions

- Consider variable y following the K categorical distribution:

$$y_k \sim \frac{\exp((\log p_k)/\tau)}{\sum_{j=0}^K \exp((\log p_j)/\tau)}$$

- Reparametrization:**

$$y_k \sim G(\log p) = \frac{\exp((\log p_k + \varepsilon)/\tau)}{\sum_{j=0}^K \exp((\log p_j + \varepsilon)/\tau)}, \quad \varepsilon = -\log(-\log(u)), u \sim \mathcal{U}[0, 1]$$

- Sampling: sample u and then sample from $G(\log p)$ to generate y_k
- In the forward pass you sample from the parametrized distribution

$$c \sim G(\log p)$$

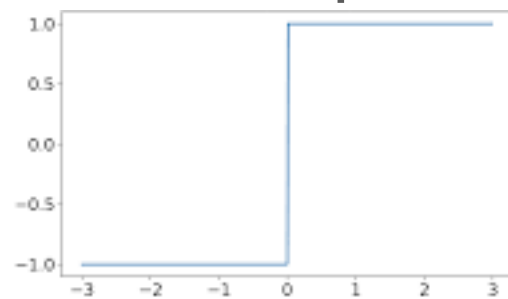
- In the backward pass you use the soft distribution:

$$\frac{dc}{d\theta} = \frac{dG}{dp} \frac{dp}{d\theta}$$

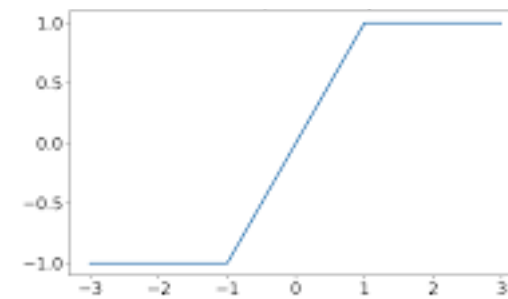
Bacproping through discrete variables

For binary neurons:

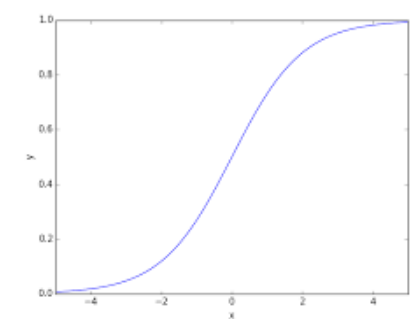
forward pass



backward pass



Straight-through



sigmoidal

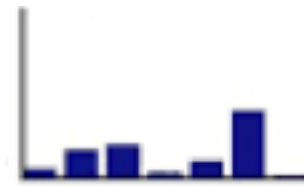
Bacproping through discrete variables

For categorically distributed neurons:

forward pass



backward pass



Summary

- Recap of estimating gradients
- Backpropagating through sampling using the reparametrization trick.