

# The Importance of Non-Markovianity in Maximum State Entropy Exploration

Mirco Mutti\*

Politecnico di Milano  
Università di Bologna

 @mirco\_mutti

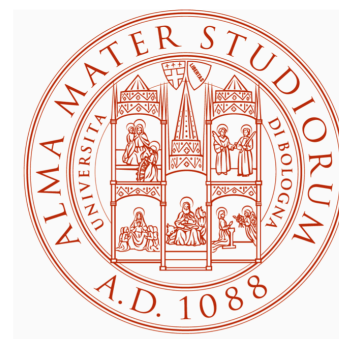
Riccardo De Santi\*

ETH Zürich

 @desariky

Marcello Restelli

Politecnico di Milano



**ETH** zürich

\* equal contribution

# Motivation

# Motivation

Markovian

$$\pi(a|s)$$

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$

# Motivation

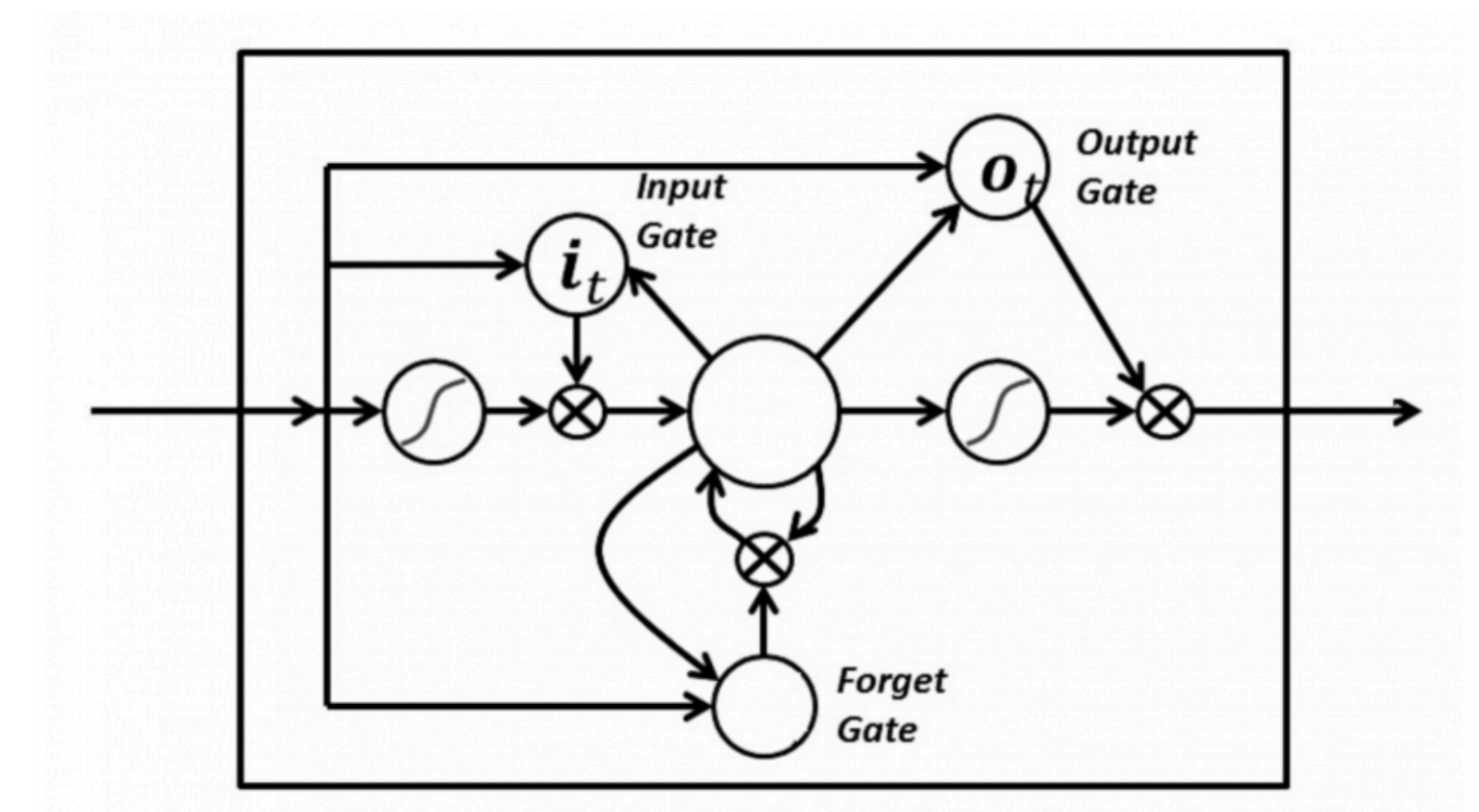
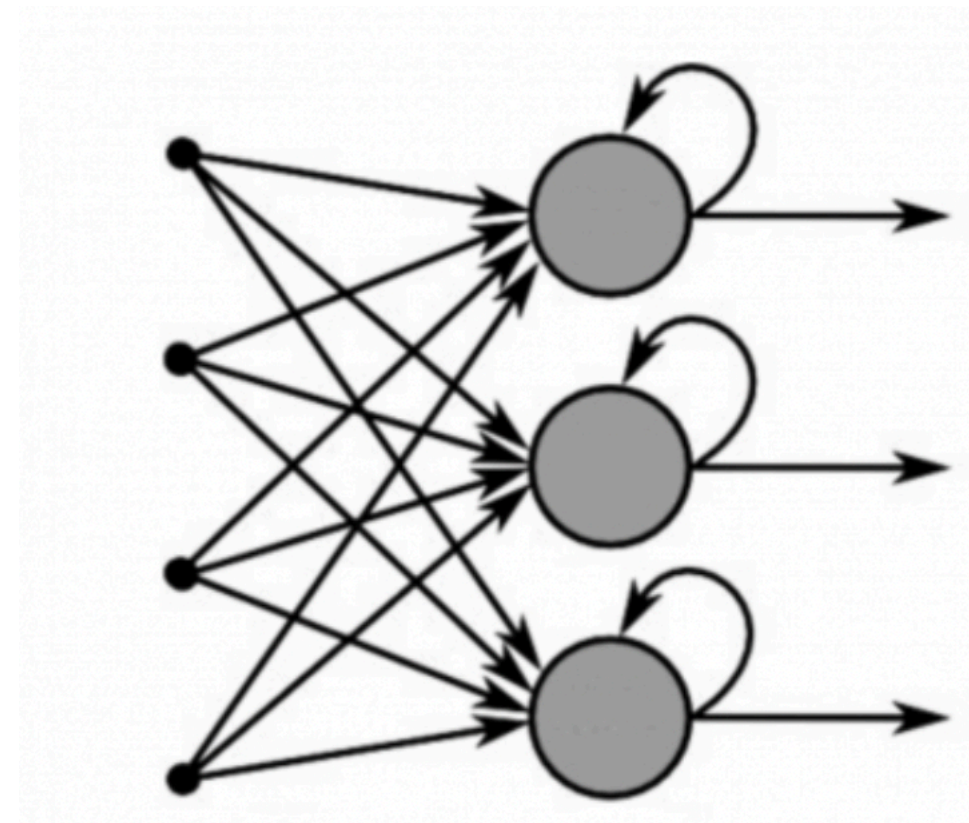
Markovian

$$\pi(a|s)$$

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$



(Williams & Zipser, 1989), (Hochreiter & Schmidhuber, 1997)



# Motivation

Markovian

$$\pi(a|s)$$

$$s_0 \rightarrow a_1$$

$$s_1 \rightarrow a_3$$

$$s_2 \rightarrow a_0$$

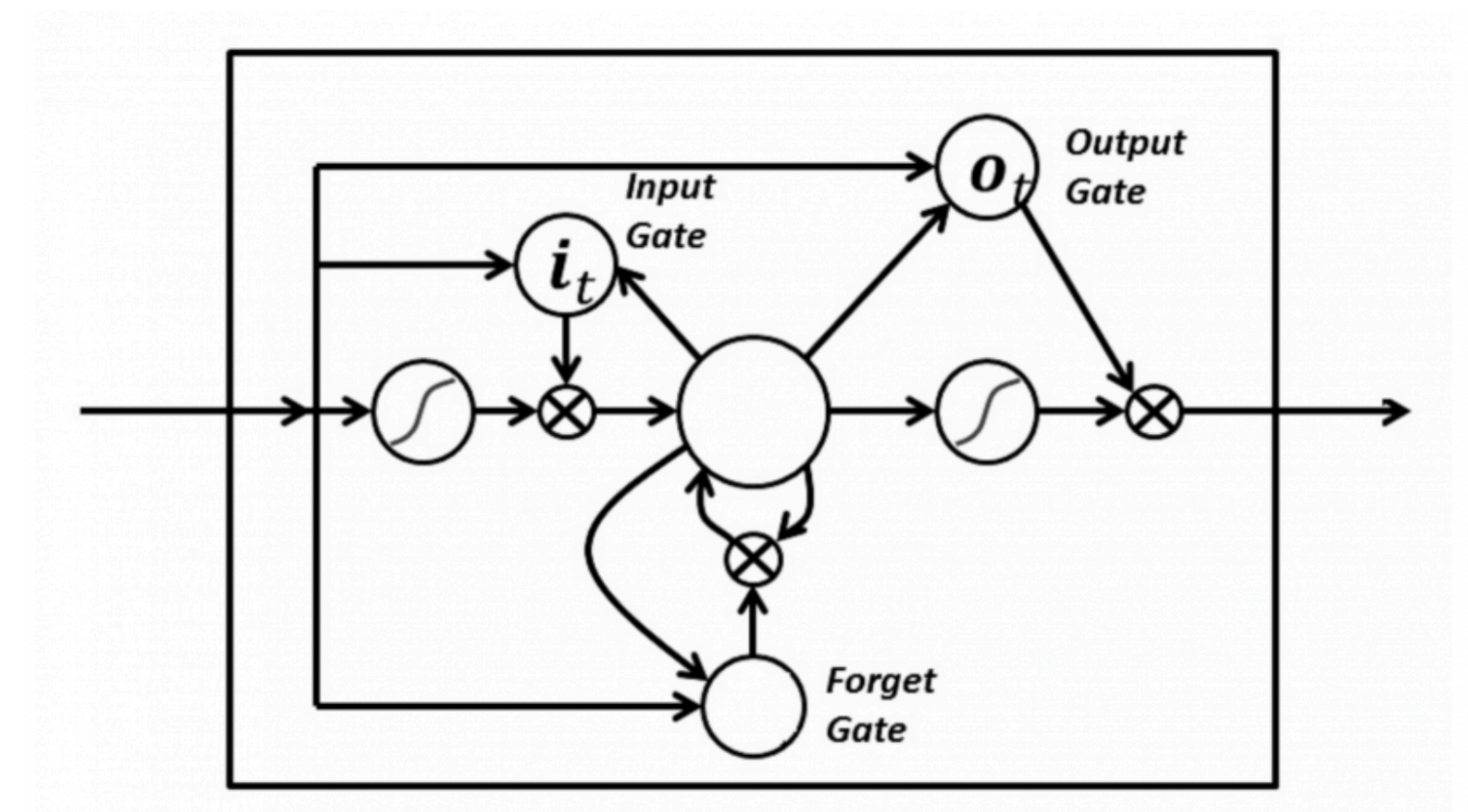
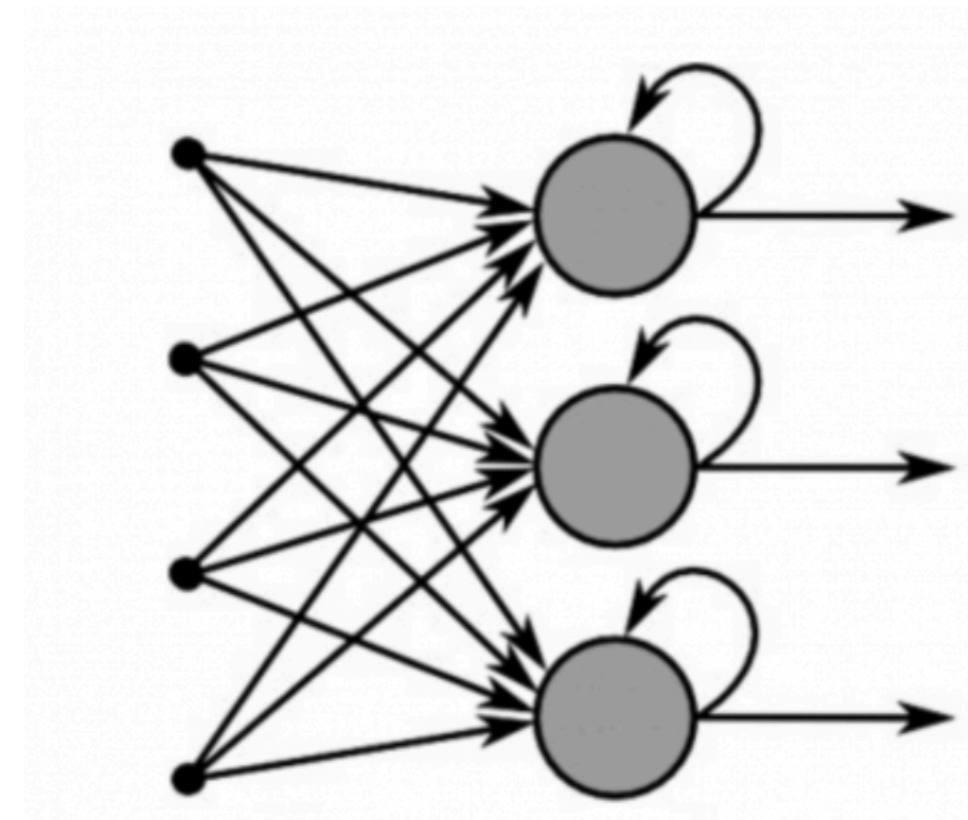
...

existence of an optimal  
deterministic Markovian policy<sup>1</sup>

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$



(Williams & Zipser, 1989), (Hochreiter & Schmidhuber, 1997), <sup>1</sup>(Proposition 4.4.3, Puterman, 2014)

# Motivation

Markovian

$$\pi(a|s)$$

$$s_0 \rightarrow a_1$$

$$s_1 \rightarrow a_3$$

$$s_2 \rightarrow a_0$$

...

existence of an optimal  
deterministic Markovian policy<sup>1</sup>

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$

Who cares about non-Markovian policies?

<sup>1</sup>(Proposition 4.4.3, Puterman, 2014)

# Motivation

Markovian

$$\pi(a|s)$$

$$s_0 \rightarrow a_1$$

$$s_1 \rightarrow a_3$$

$$s_2 \rightarrow a_0$$

...

existence of an optimal  
deterministic Markovian policy<sup>1</sup>

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$

Who cares about non-Markovian policies?

Partial observability

<sup>1</sup>(Proposition 4.4.3, Puterman, 2014)



# Motivation

Markovian

$$\pi(a|s)$$

$$s_0 \rightarrow a_1$$

$$s_1 \rightarrow a_3$$

$$s_2 \rightarrow a_0$$

...

existence of an optimal  
deterministic Markovian policy<sup>1</sup>

Non-Markovian

$$\pi(a|h)$$

$$h = (s_0, a_0, s_1, a_1, \dots, s)$$

Who cares about non-Markovian policies?

Partial observability

Imitation learning, risk-aversion, pure exploration, ... ?

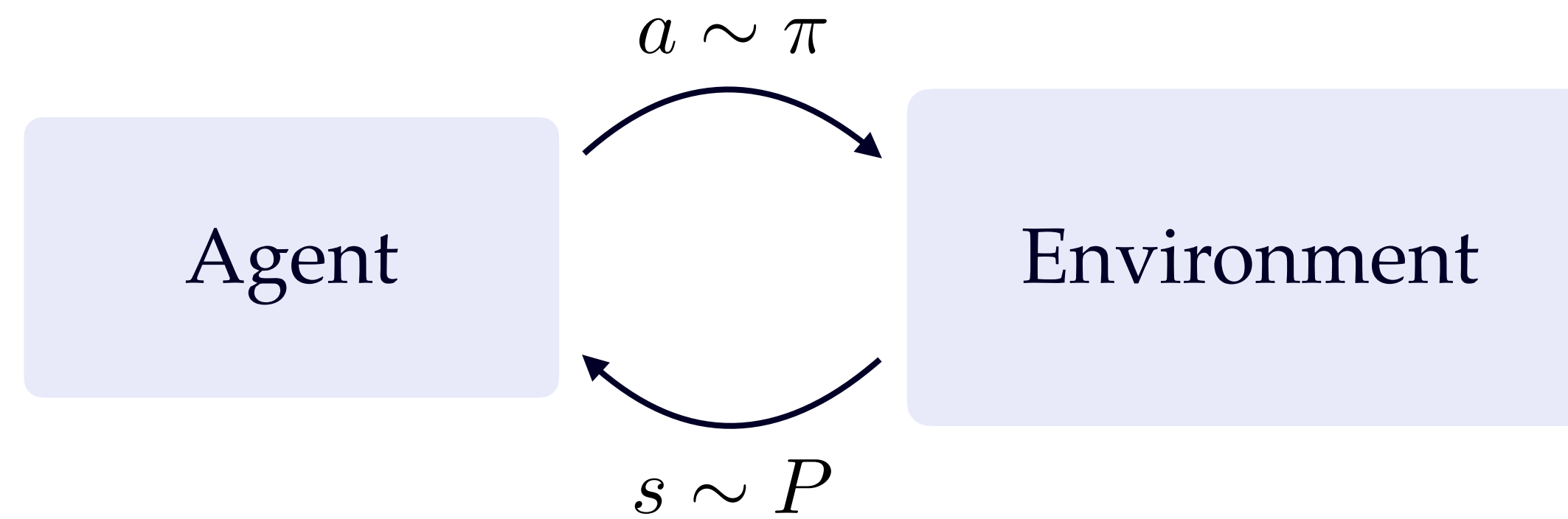
<sup>1</sup>(Proposition 4.4.3, Puterman, 2014)



# Problem Setting

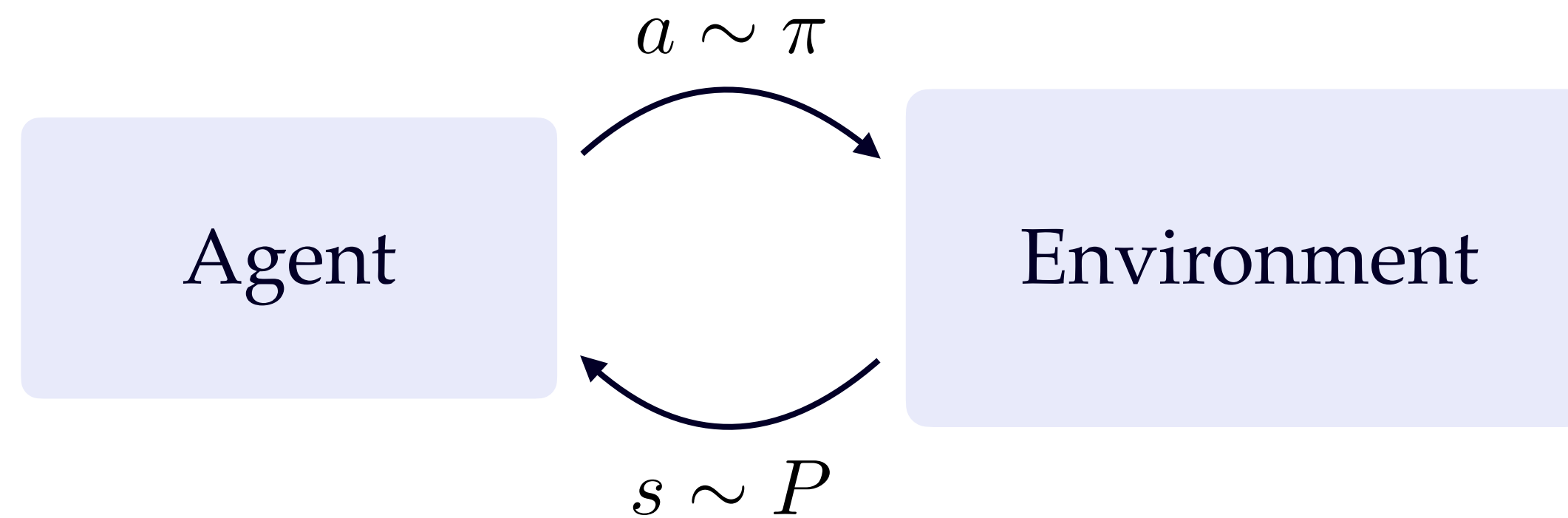
# Problem Setting

Controlled Markov Process (CMP)



# Problem Setting

Controlled Markov Process (CMP)



$\mathcal{S}$  discrete set of states

$\mathcal{A}$  discrete set of actions

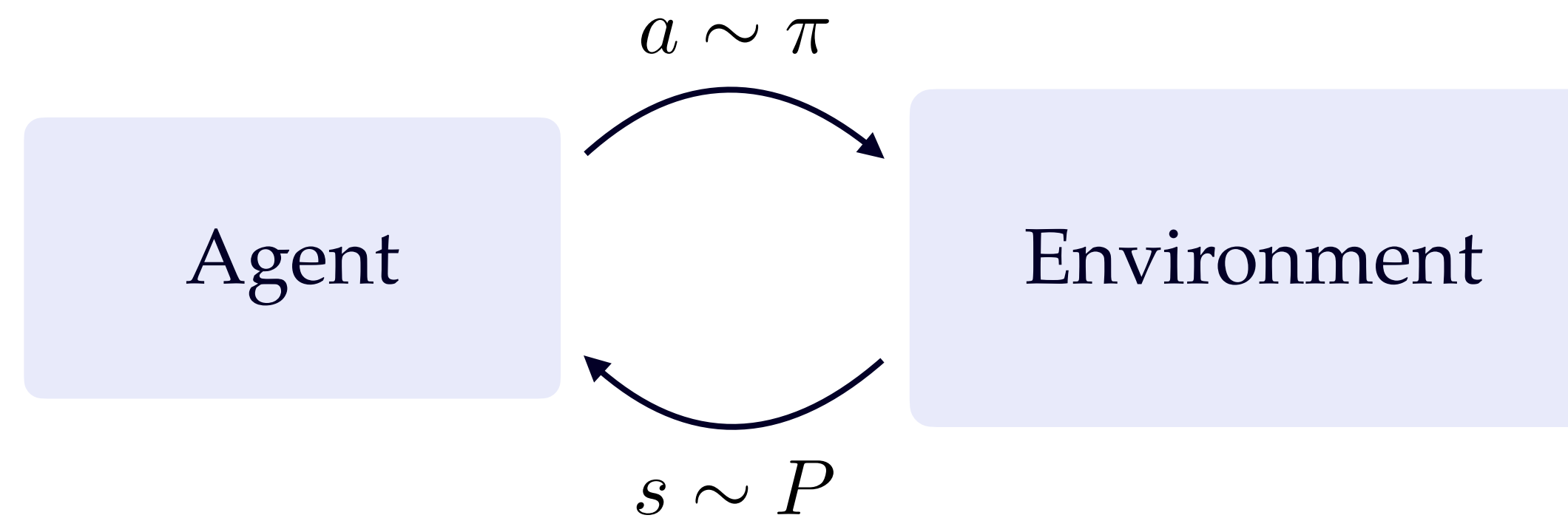
$P$  transition matrix

$\mu$  initial state distribution

$T$  episode horizon

# Problem Setting

Controlled Markov Process (CMP)



$\Pi_{\text{NM}}$  set of **non-Markovian** policies

$$\pi : \mathcal{H} \rightarrow \Delta(\mathcal{A})$$

$$\mathcal{H} := \mathcal{S} \times \mathcal{S} \times \dots$$

$\Pi_{\text{M}}$  set of **Markovian** policies

$$\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$$

$\mathcal{S}$  discrete set of states

$\mathcal{A}$  discrete set of actions

$P$  transition matrix

$\mu$  initial state distribution

$T$  episode horizon



# Problem Setting

policy  $\pi$  + CMP



marginal state distribution

$$d^\pi(s) = \frac{1}{T} \sum_{t \in [T]} \Pr(s_t = s)$$

# Problem Setting

policy  $\pi$  + CMP



marginal state distribution

$$d^\pi(s) = \frac{1}{T} \sum_{t \in [T]} Pr(s_t = s)$$

**Reinforcement Learning (RL)**

$$\mathcal{J}(\pi) = d^\pi \cdot R$$

# Problem Setting

policy  $\pi$  + CMP



marginal state distribution

$$d^\pi(s) = \frac{1}{T} \sum_{t \in [T]} Pr(s_t = s)$$

**Reinforcement Learning (RL)**

$$\mathcal{J}(\pi) = d^\pi \cdot R$$

**Convex Reinforcement Learning (CRL)<sup>1,2</sup>**

$$\mathcal{J}(\pi) = \mathcal{F}(d^\pi)$$

$\mathcal{F}$  is a convex/concave function

<sup>1</sup>(Zhang et al., 2020), <sup>2</sup>(Zahavy et al., 2021)

# Problem Setting

policy  $\pi$  + CMP



marginal state distribution

$$d^\pi(s) = \frac{1}{T} \sum_{t \in [T]} \Pr(s_t = s)$$

CFOL Workshop @ ICML

Mutti et al. "Challenging Common Assumptions in Convex Reinforcement Learning". 2022.

Convex Reinforcement Learning (CRL)<sup>1,2</sup>

$$\mathcal{J}(\pi) = \mathcal{F}(d^\pi)$$

$\mathcal{F}$  is a convex/concave function

<sup>1</sup>(Zhang et al., 2020), <sup>2</sup>(Zahavy et al., 2021)



# Problem Setting

policy  $\pi$  + CMP



marginal state distribution

$$d^\pi(s) = \frac{1}{T} \sum_{t \in [T]} \Pr(s_t = s)$$

CFOL Workshop @ ICML

Mutti et al. "Challenging Common Assumptions in Convex Reinforcement Learning". 2022.

this paper

**Maximum State Entropy (MSE)<sup>3</sup>**

$$\mathcal{E}(\pi) = H(d^\pi) = d^\pi \log d^\pi$$



**Convex Reinforcement Learning (CRL)<sup>1,2</sup>**

$$\mathcal{J}(\pi) = \mathcal{F}(d^\pi)$$

$\mathcal{F}$  is a convex/concave function

<sup>1</sup>(Zhang et al., 2020), <sup>2</sup>(Zahavy et al., 2021), <sup>3</sup>(Hazan et al., 2019)

# Does Non-Markovianity Matter?

$$\mathcal{E}(\pi) := H(d^\pi)$$




Markovian policies  
are sufficient<sup>1</sup>

<sup>1</sup>[Puterman, 2014]

# Does Non-Markovianity Matter?

state visitation frequency  $d(s) = \frac{1}{T} \sum_{t \in [T]} \mathbb{1}(s_t = s)$

$$\mathcal{E}_{\infty}(\pi) := H(d^{\pi}) = H\left(\mathbb{E}_{d \sim p^{\pi}}[d]\right)$$


↓  
Markovian policies  
are sufficient<sup>1</sup>

<sup>1</sup>[Puterman, 2014]

# Does Non-Markovianity Matter?

$$\mathcal{E}_\infty(\pi) := H(d^\pi) = H\left(\mathbb{E}_{d \sim p^\pi}[d]\right) \geq \mathbb{E}_{d \sim p^\pi} [H(d)]$$



Markovian policies  
are sufficient<sup>1</sup>

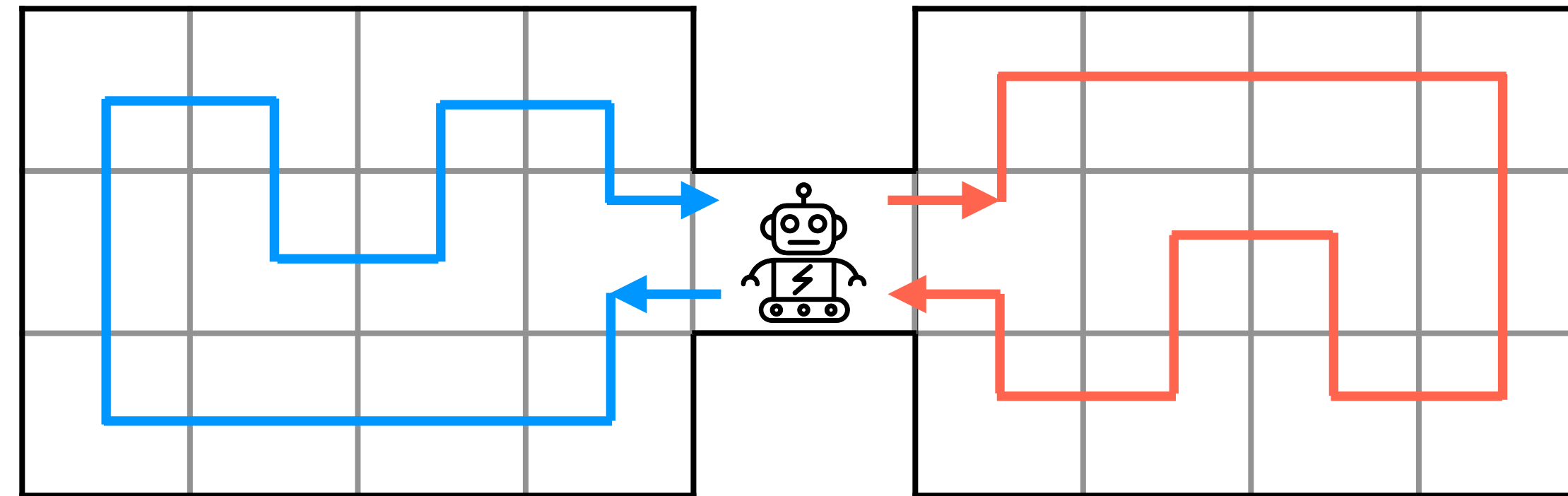
(through Jensen's)

<sup>1</sup>[Puterman, 2014]



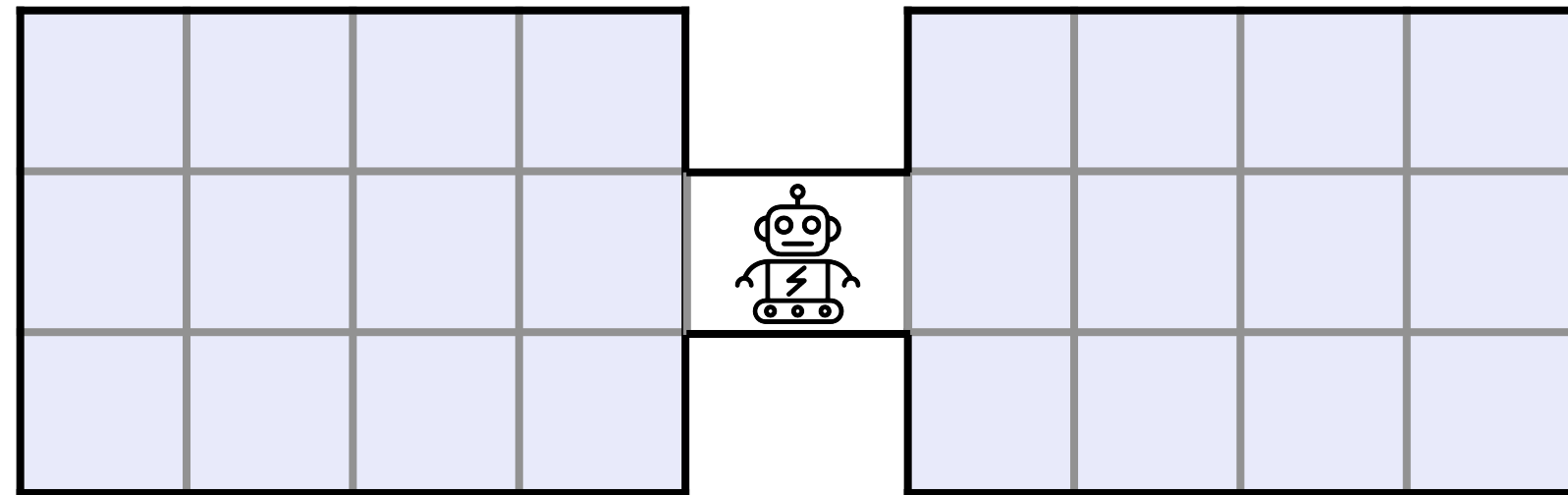


# Illustrative Example

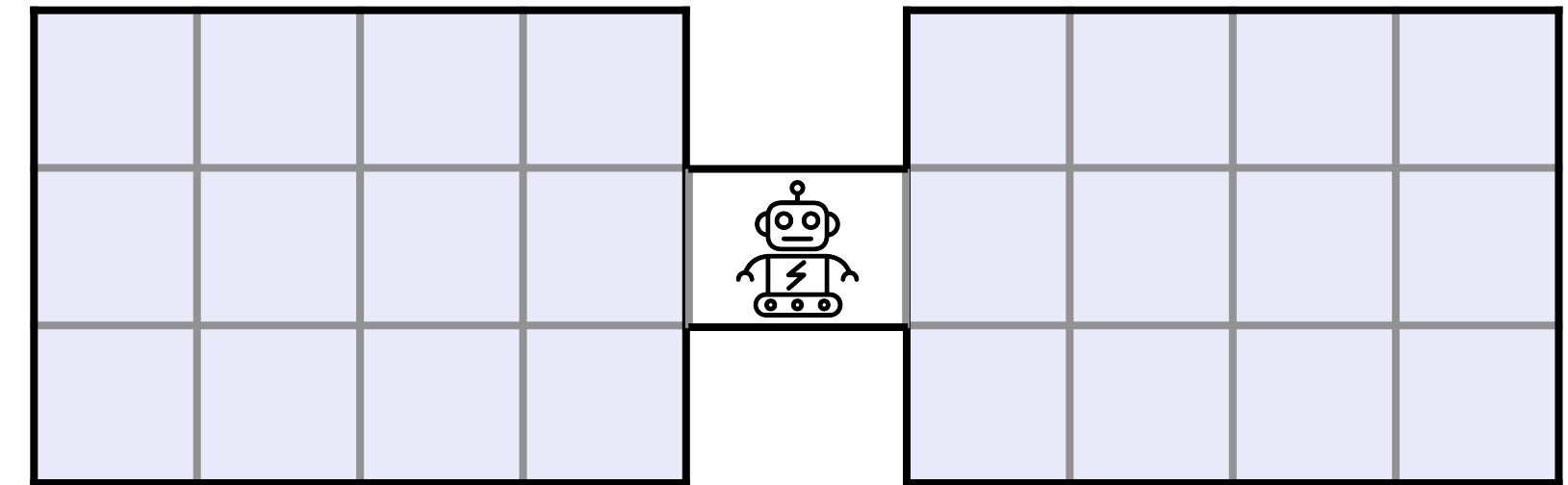


# Illustrative Example

Best Markovian Policy

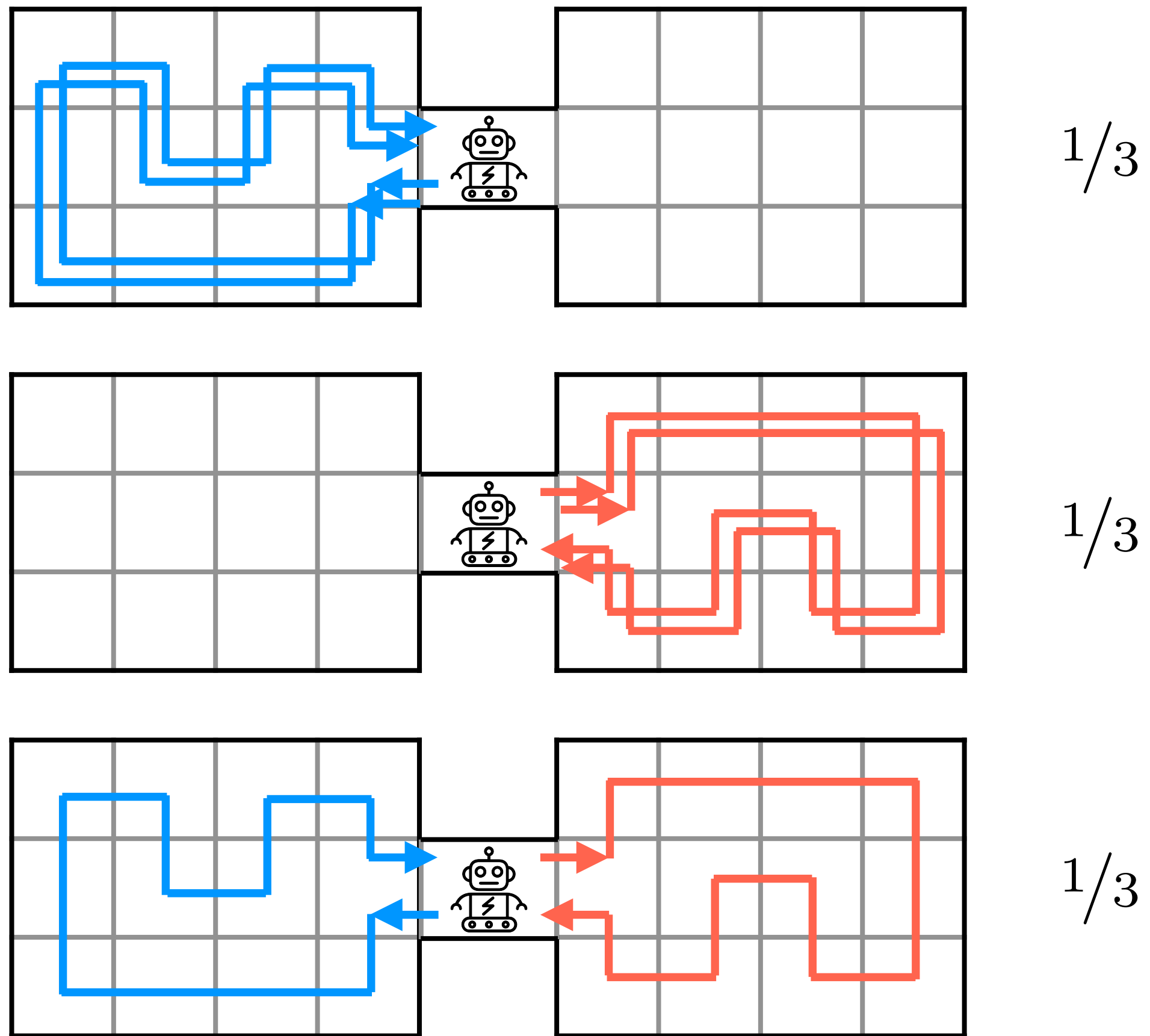


Best Non-Markovian Policy



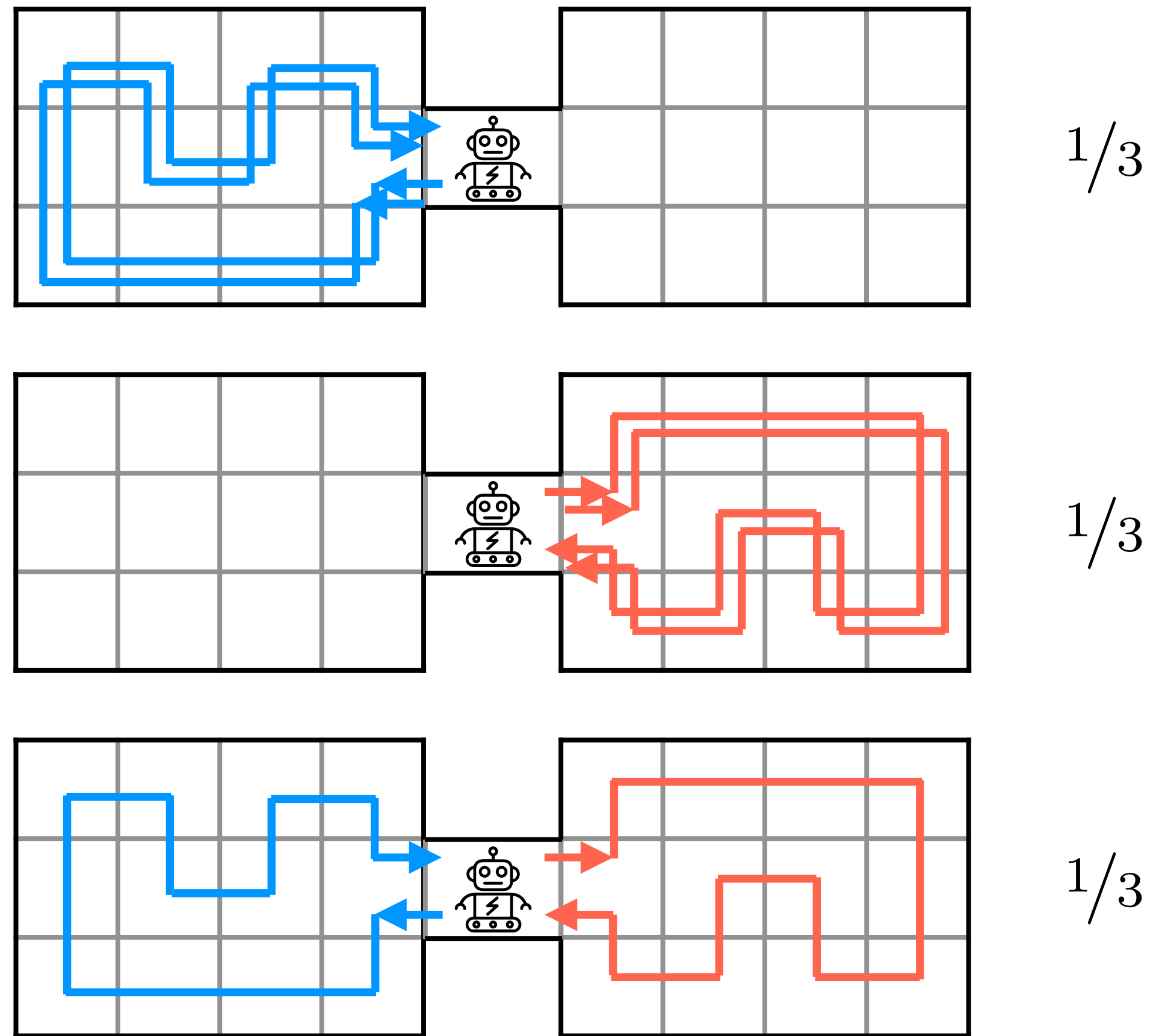
# Illustrative Example

## Best Markovian Policy

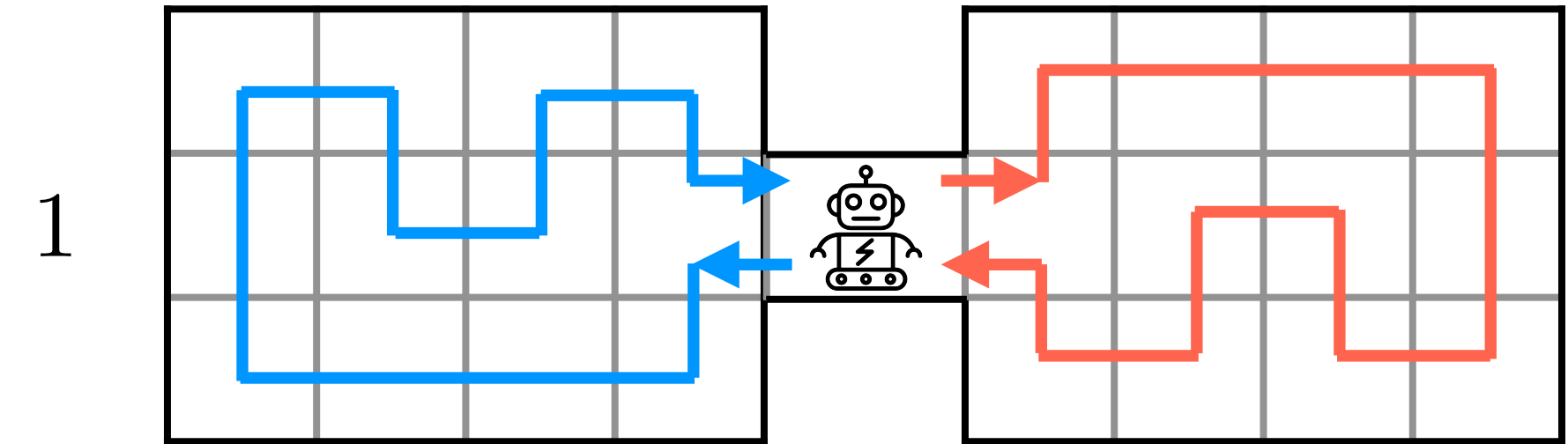


# Illustrative Example

Best Markovian Policy



Best Non-Markovian Policy



# The Importance of Non-Markovianity

Finite-Sample Maximum State Entropy

$$\mathcal{E}_1(\pi) = \mathbb{E}_{d \sim p^\pi} [H(d)]$$



# The Importance of Non-Markovianity

Finite-Sample Maximum State Entropy

$$\mathcal{E}_1(\pi) = \mathbb{E}_{d \sim p^\pi} [H(d)]$$

A tool to compare Markovian and non-Markovian policies?

# The Importance of Non-Markovianity

**Finite-Sample** Maximum State Entropy

$$\mathcal{E}_1(\pi) = \mathbb{E}_{d \sim p^\pi} [H(d)]$$

A tool to compare Markovian and non-Markovian policies?

The entropy is non-additive, standard regret cannot be used

# Regret

Definition (Expected Regret-to-go). Let  $\pi$  be a policy interacting with an MDP over  $T - t$  steps starting from trajectory  $h_t$ . We define the expected regret-to-go as

$$R_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p^\pi} [H(d_{h_t \oplus h_{T-t}})]$$

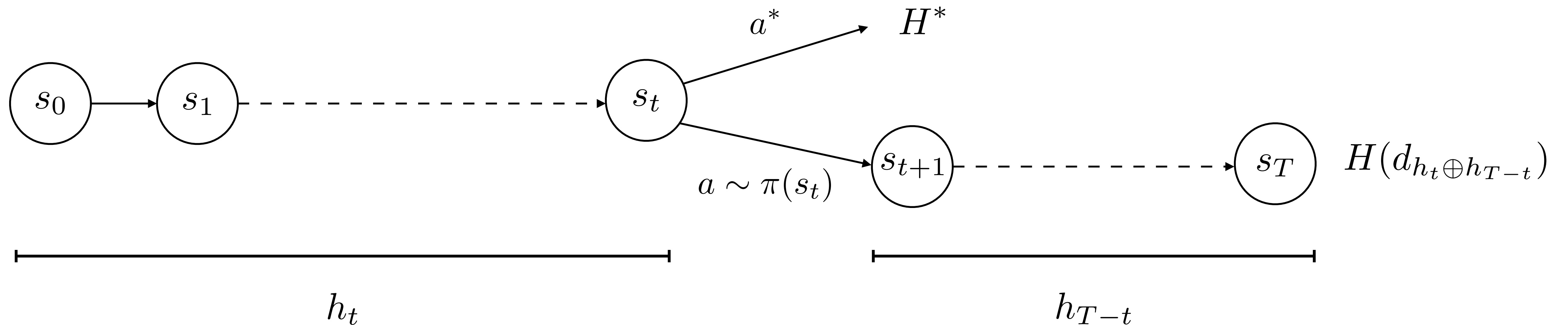
where  $H^*$  is the maximum entropy that can be achieved by any policy starting from  $h_t$ .

# Regret

Definition (Expected Regret-to-go). Let  $\pi$  be a policy interacting with an MDP over  $T - t$  steps starting from trajectory  $h_t$ . We define the expected regret-to-go as

$$R_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p^\pi} [H(d_{h_t \oplus h_{T-t}})]$$

where  $H^*$  is the maximum entropy that can be achieved by any policy starting from  $h_t$ .



# Regret Bounds

Definition (Expected Regret-to-go). Let  $\pi$  be a policy interacting with an MDP over  $T - t$  steps starting from trajectory  $h_t$ . We define the expected regret-to-go as

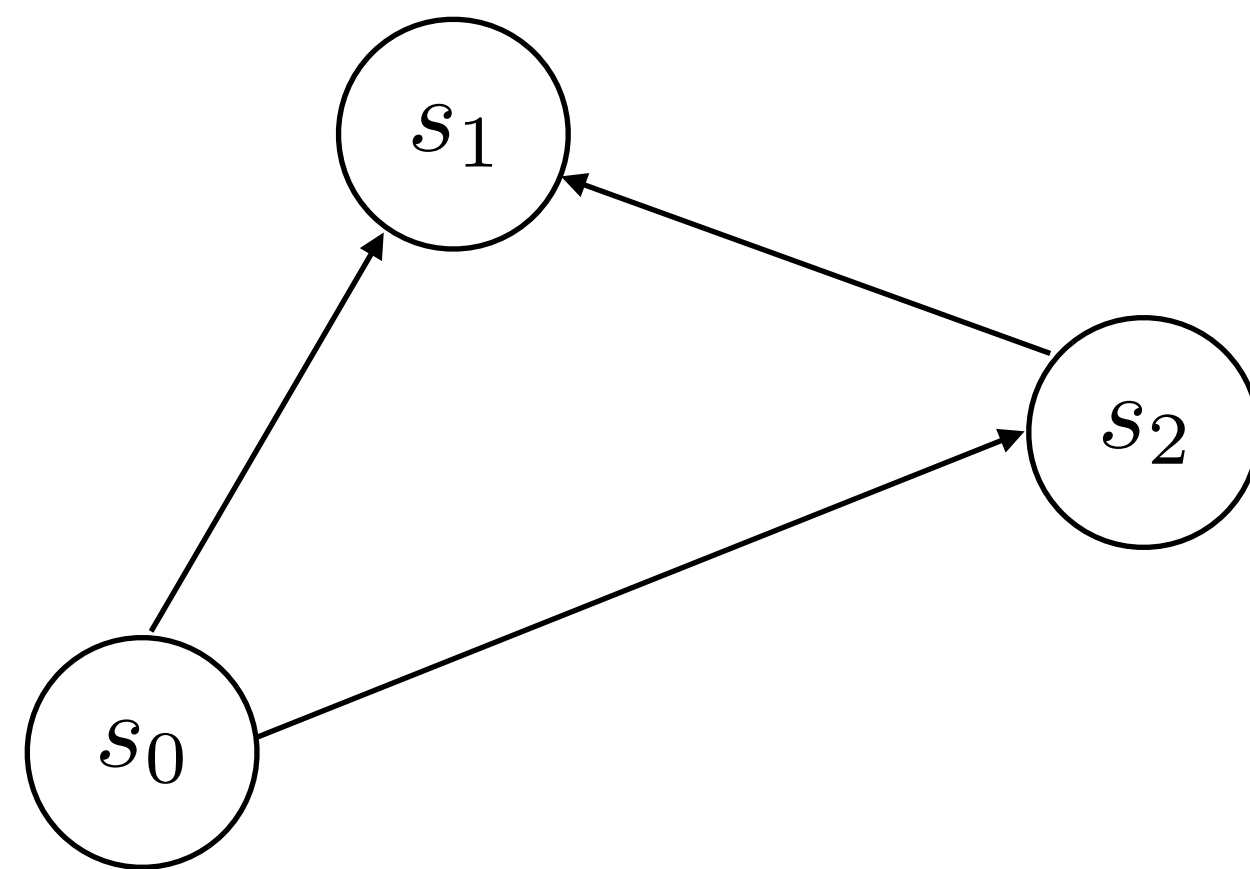
$$R_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p^\pi} [H(d_{h_t \oplus h_{T-t}})]$$

where  $H^*$  is the maximum entropy that can be achieved by any policy starting from  $h_t$ .

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$

# Regret Bounds

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$

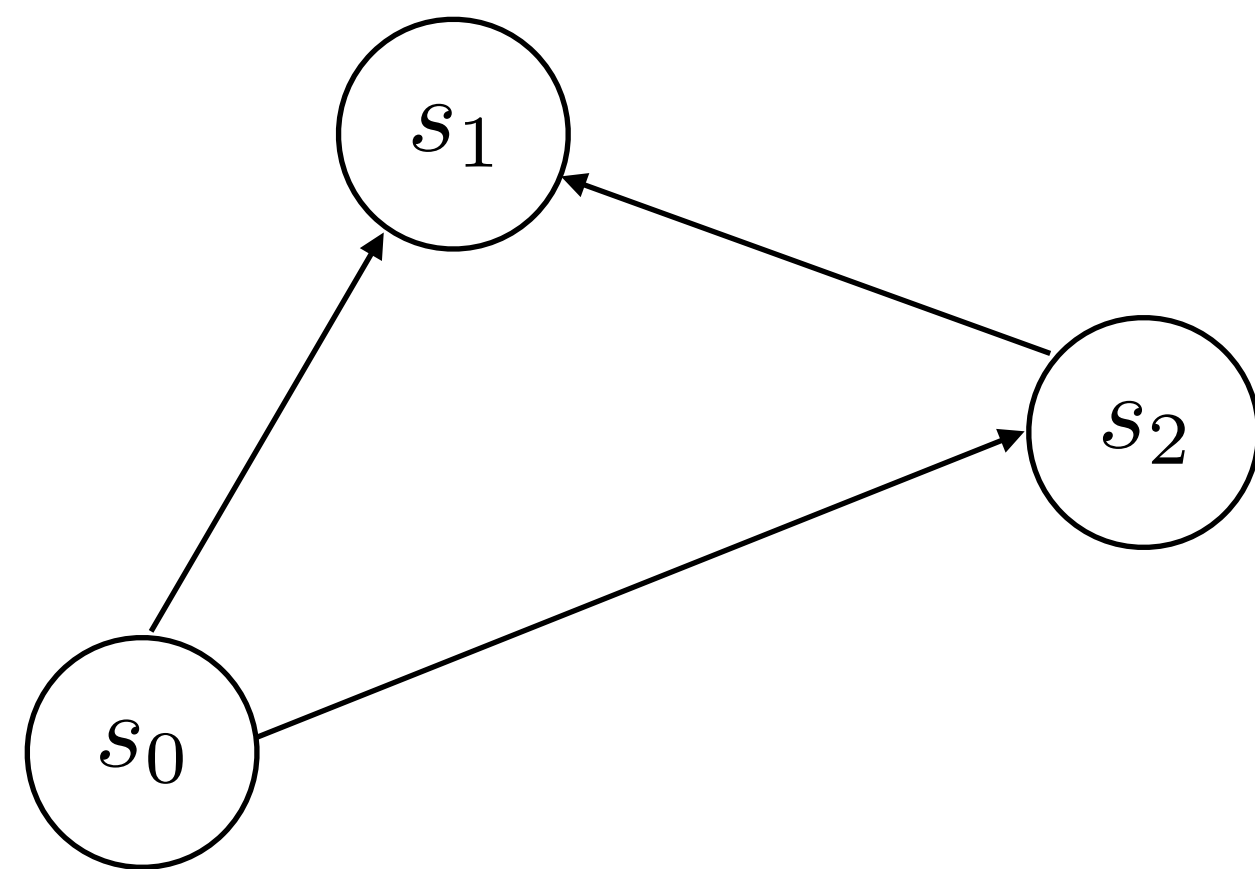


**CMP** with **MSE** objective

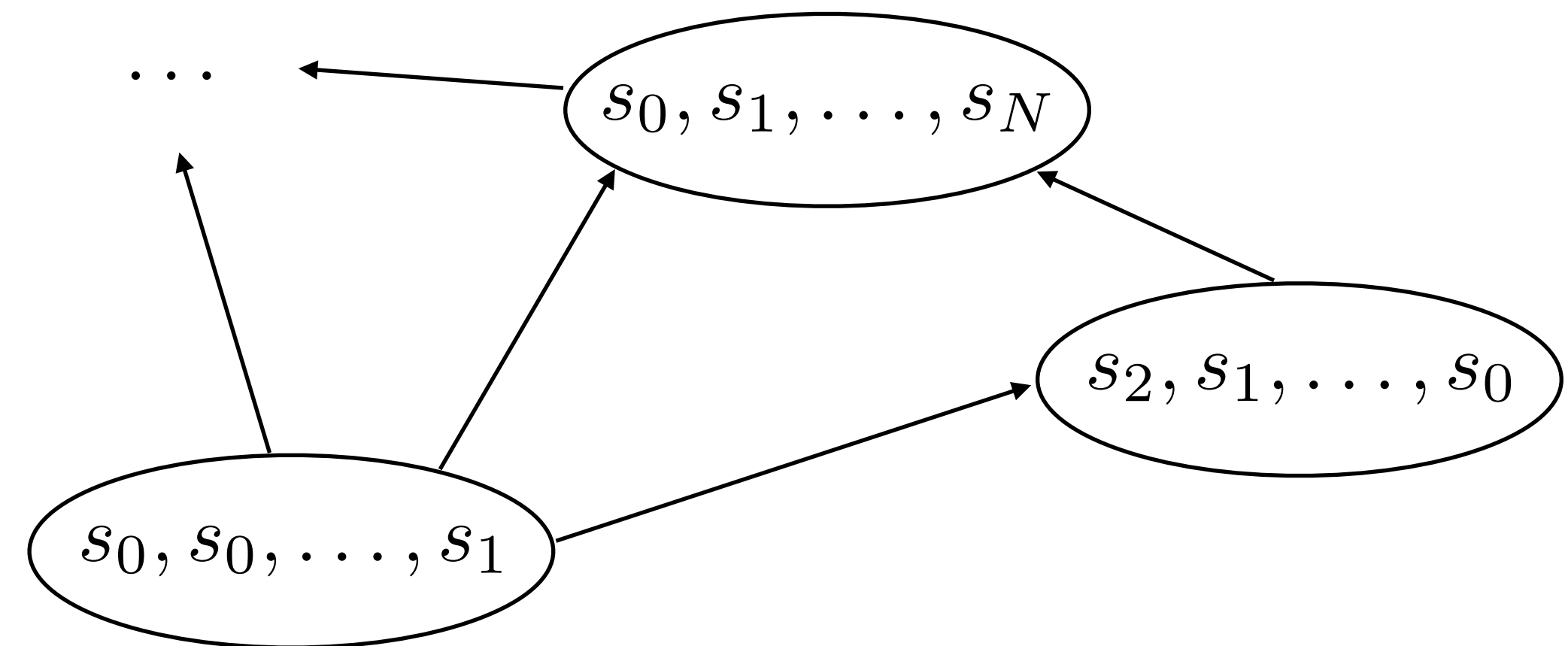


# Regret Bounds

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$



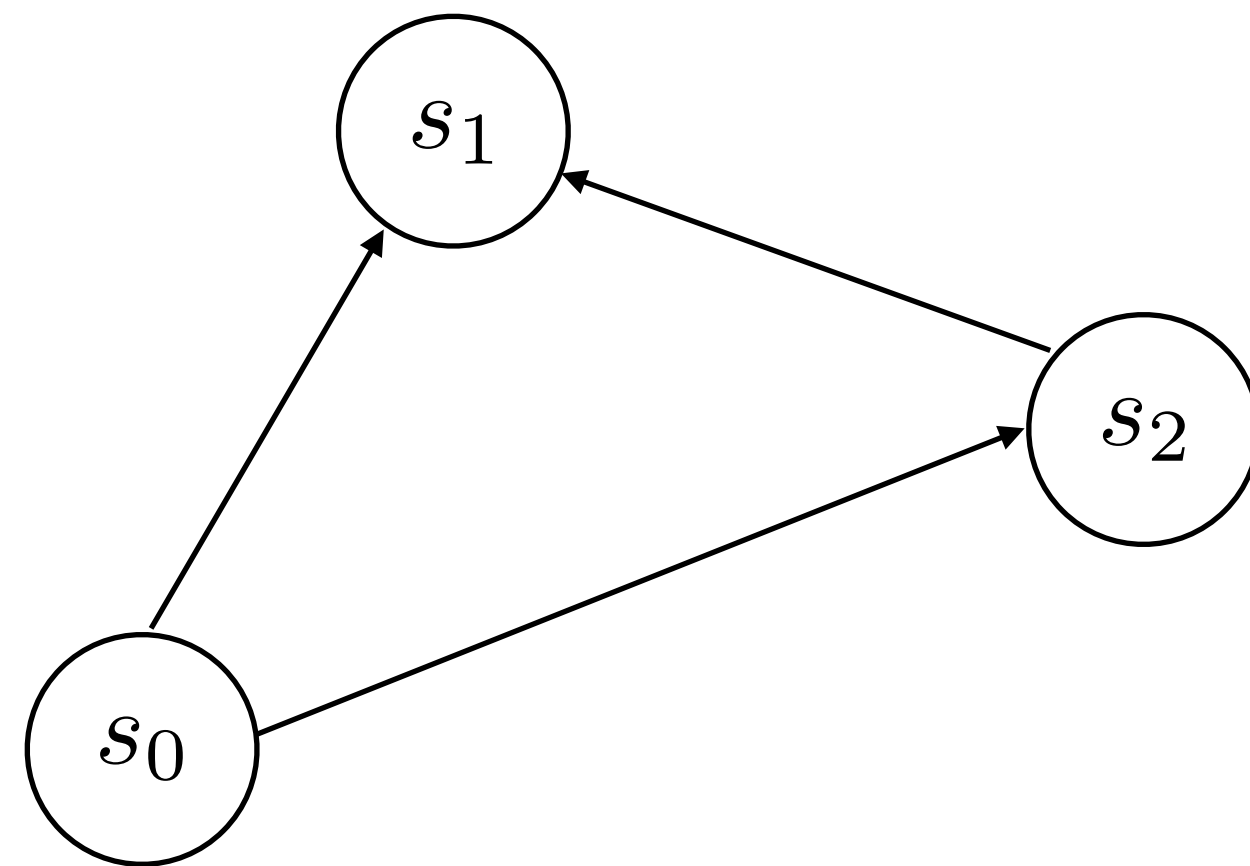
**CMP** with MSE objective



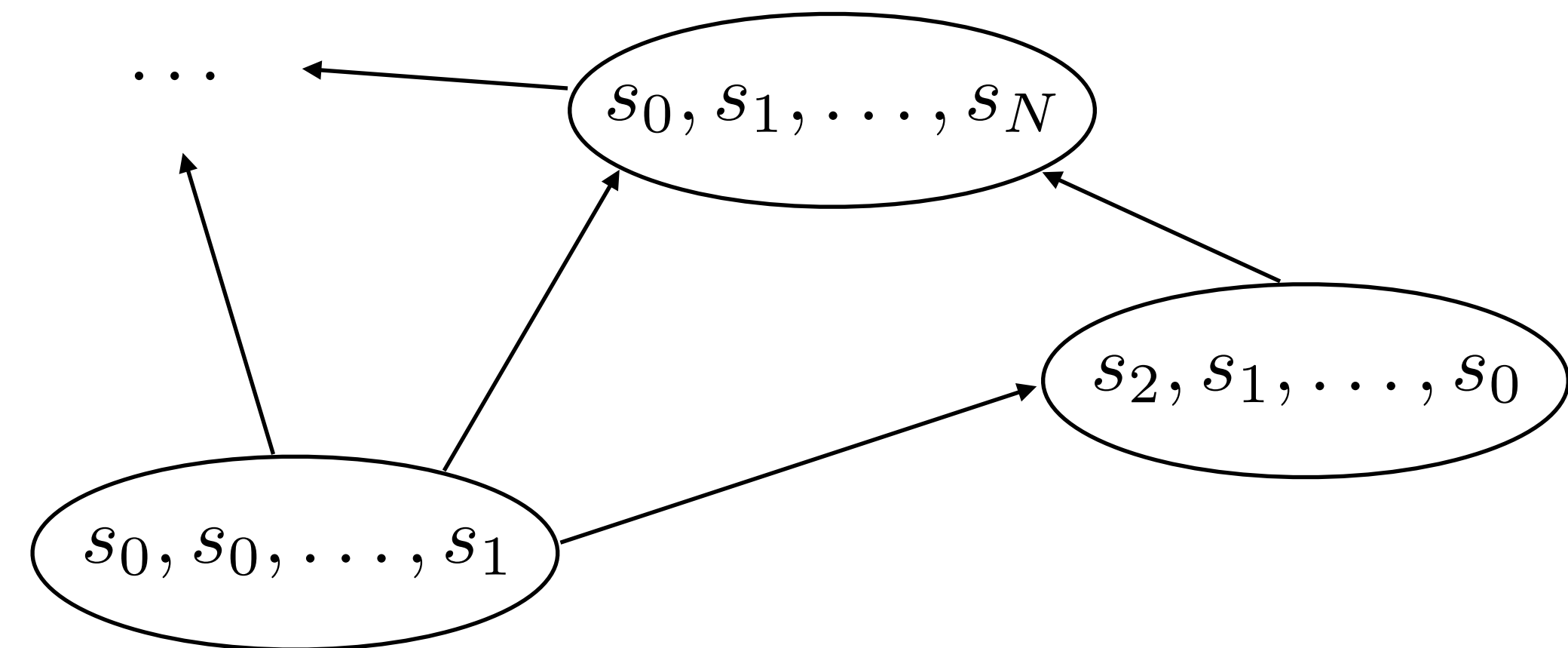
**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

# Regret Bounds

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$



**CMP** with **MSE** objective



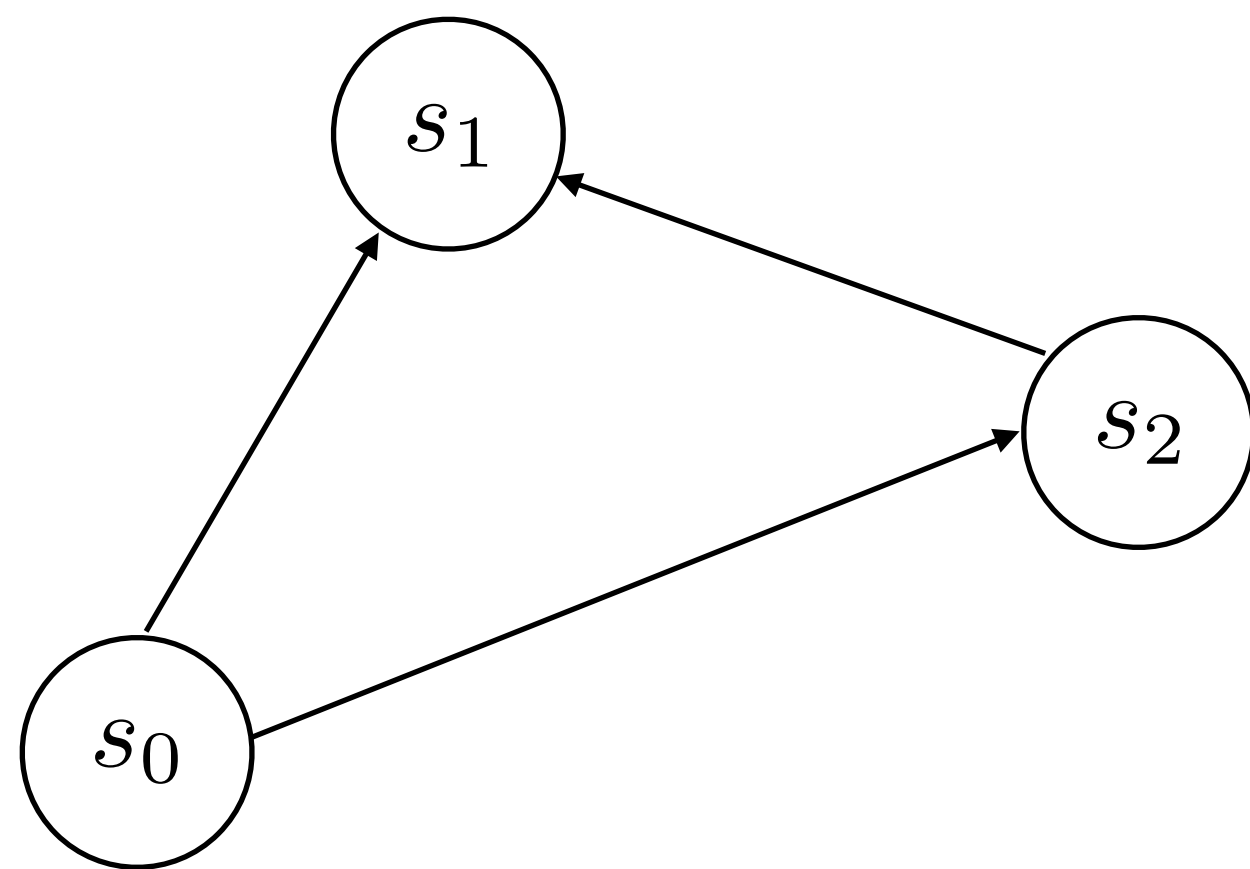
**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

optimal deterministic Markovian policy<sup>1</sup>

<sup>1</sup>(Proposition 4.4.3, Puterman, 2014)

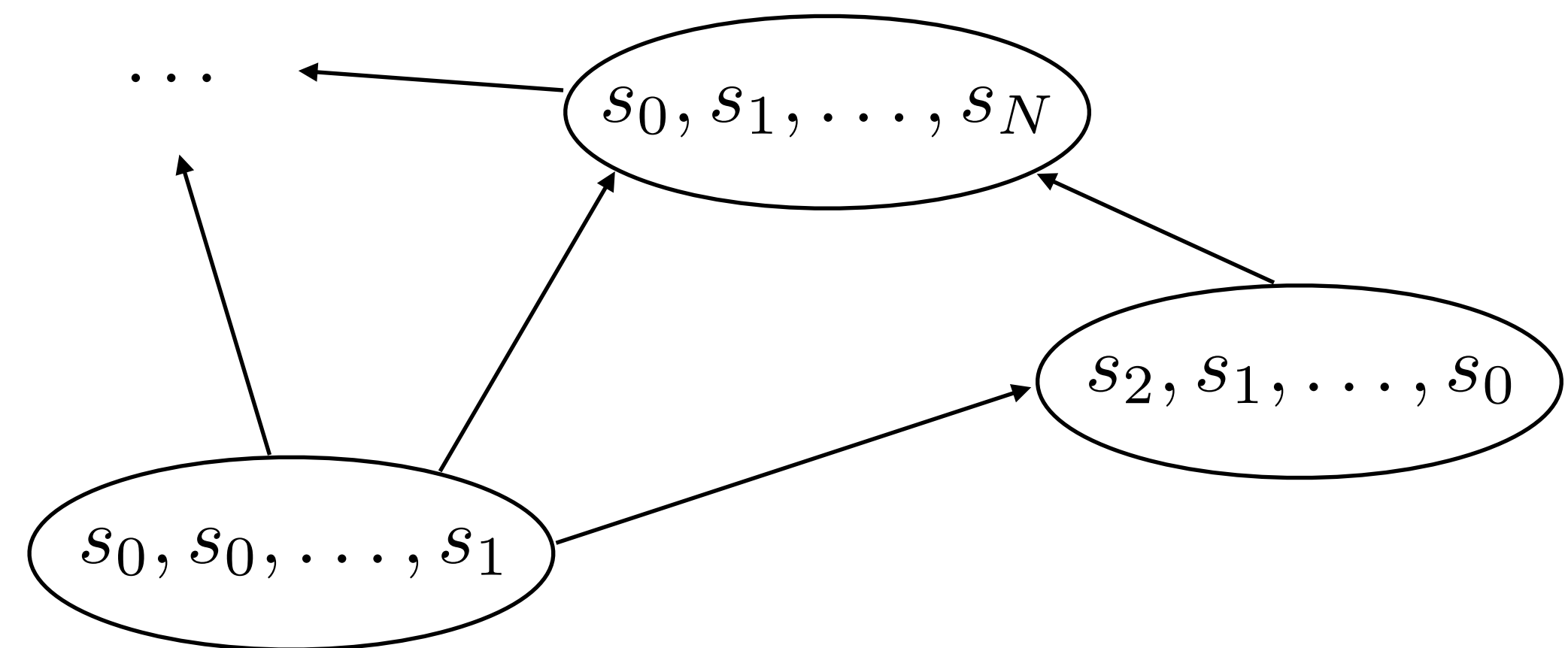
# Regret Bounds

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$



**CMP with MSE objective**

optimal deterministic non-Markovian policy



**Extended CMP with reward  $R(h_T) = H(d_{h_T})$**

$\Leftarrow$  optimal deterministic Markovian policy<sup>1</sup>

<sup>1</sup>(Proposition 4.4.3, Puterman, 2014)

# Regret Bounds

Definition (Expected Regret-to-go). Let  $\pi$  be a policy interacting with an MDP over  $T - t$  steps starting from trajectory  $h_t$ . We define the expected regret-to-go as

$$R_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p^\pi} [H(d_{h_t \oplus h_{T-t}})]$$

where  $H^*$  is the maximum entropy that can be achieved by any policy starting from  $h_t$ .

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$

**Second step.** The optimal Markovian policy  $\pi_{\text{M}}$  is randomized

# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] =$$

# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$

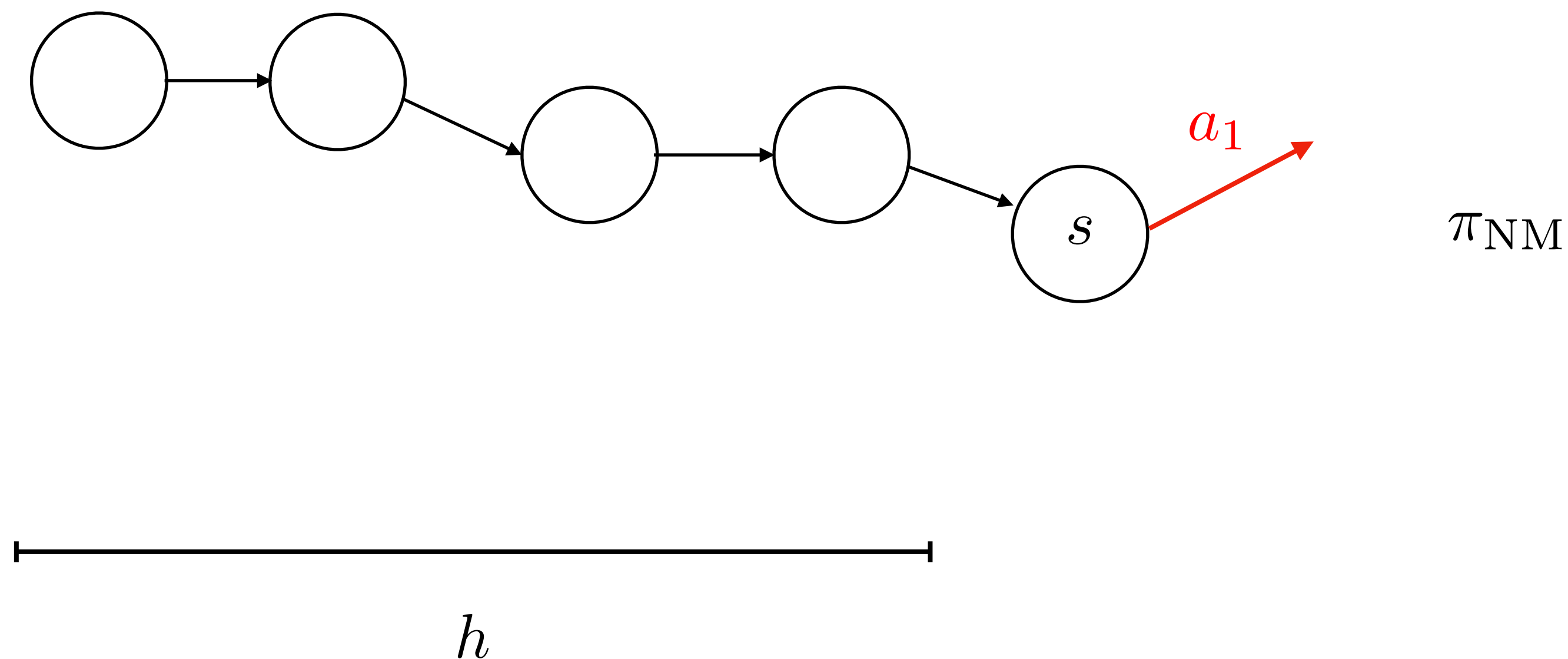
(through the Law of Total Variance and the determinism of  $\pi_{NM}$ )



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

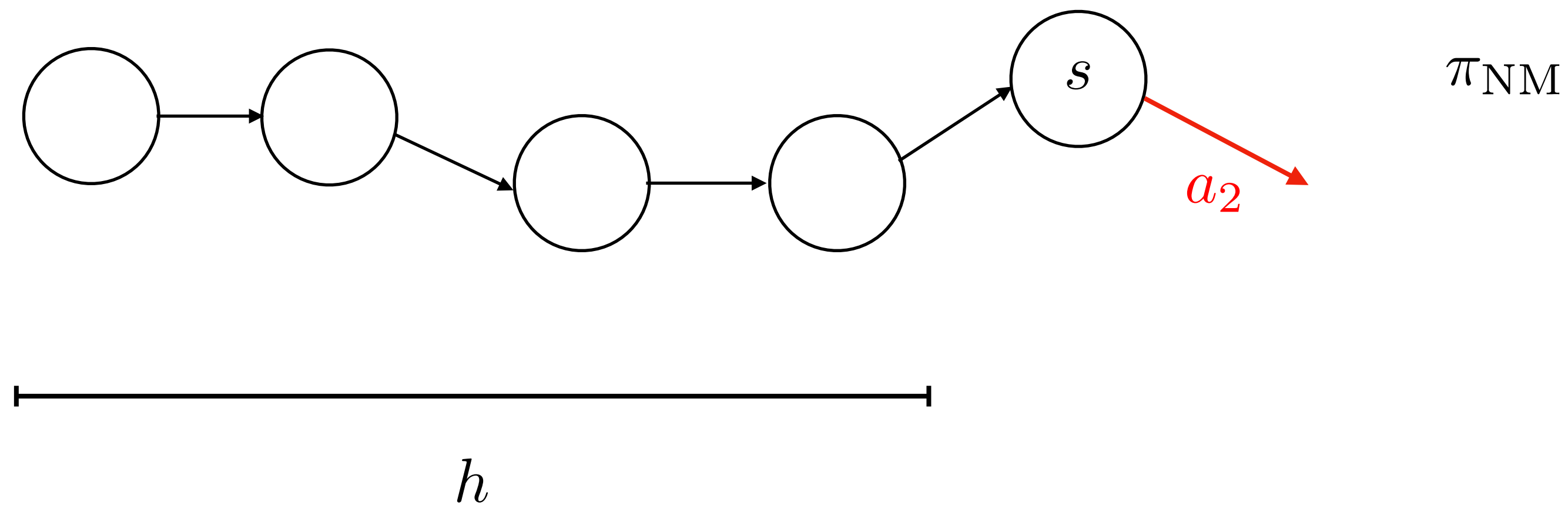
$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

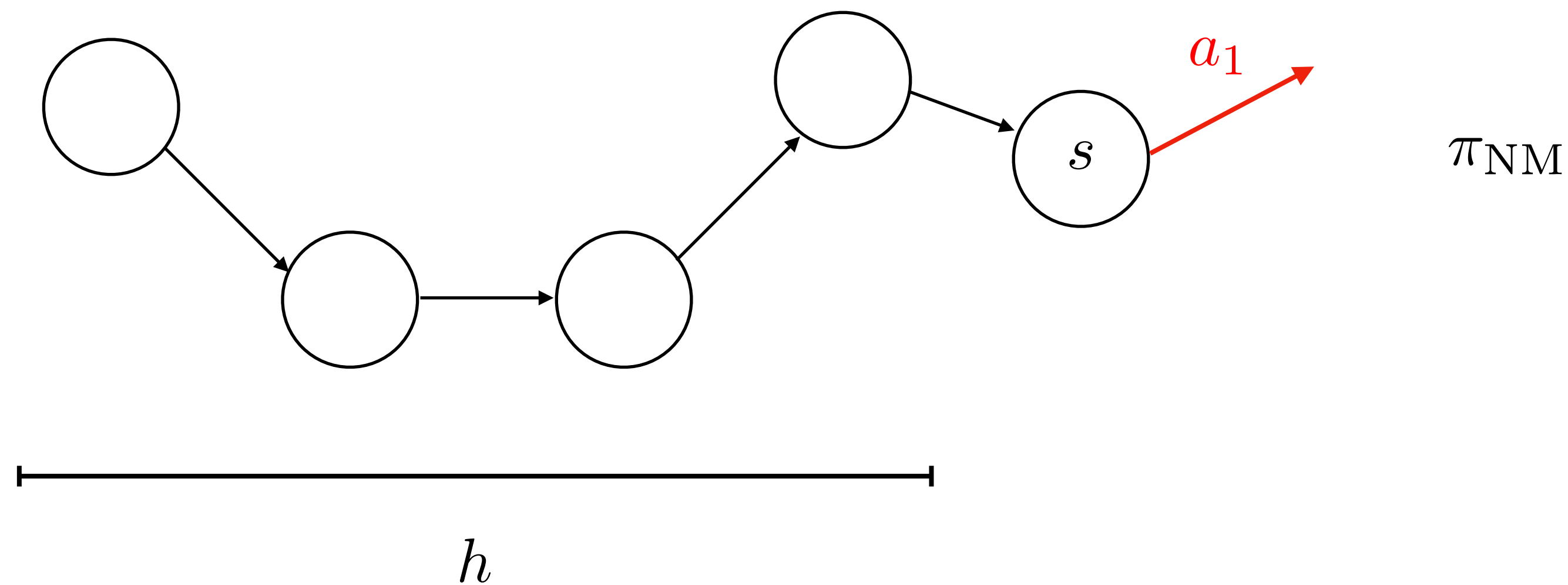
$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

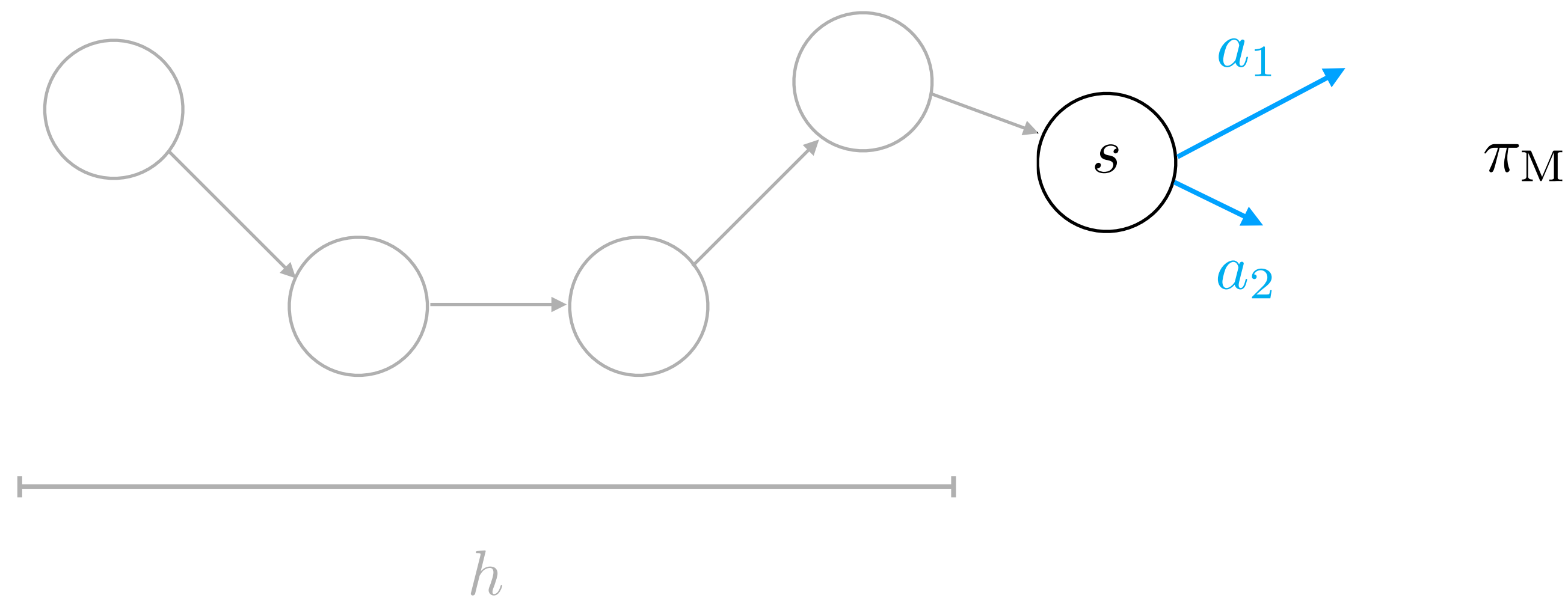
$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

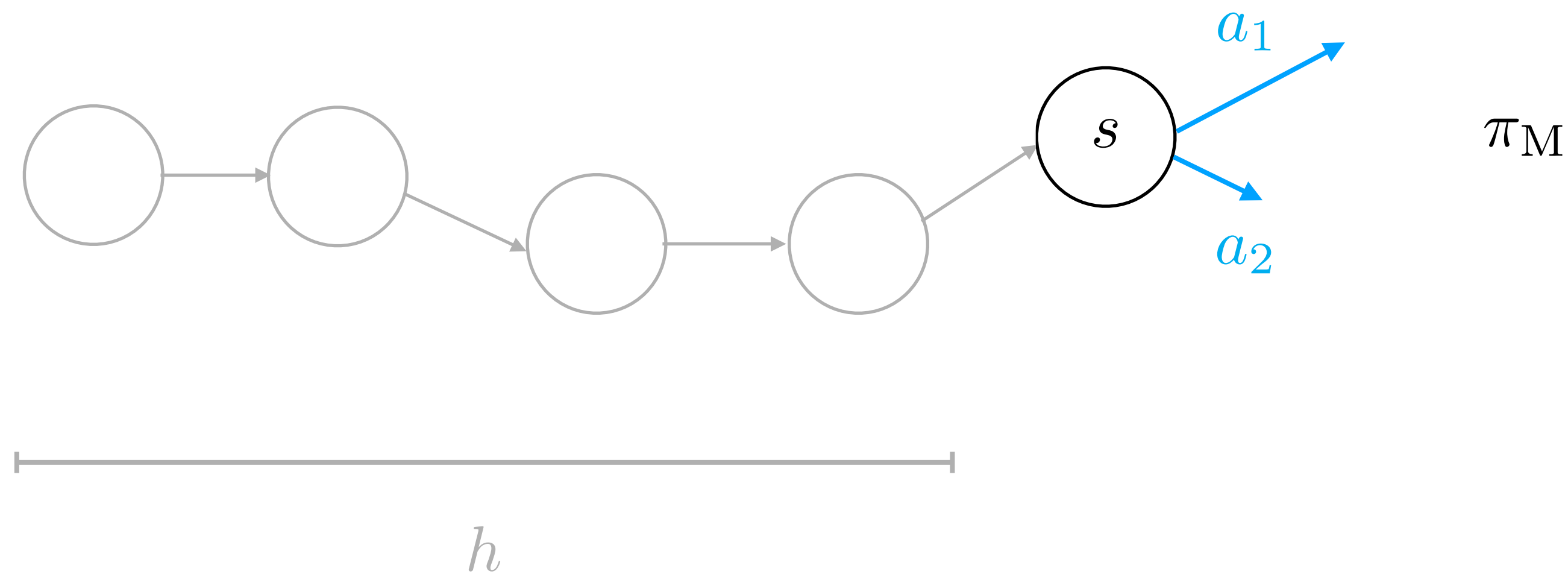
$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

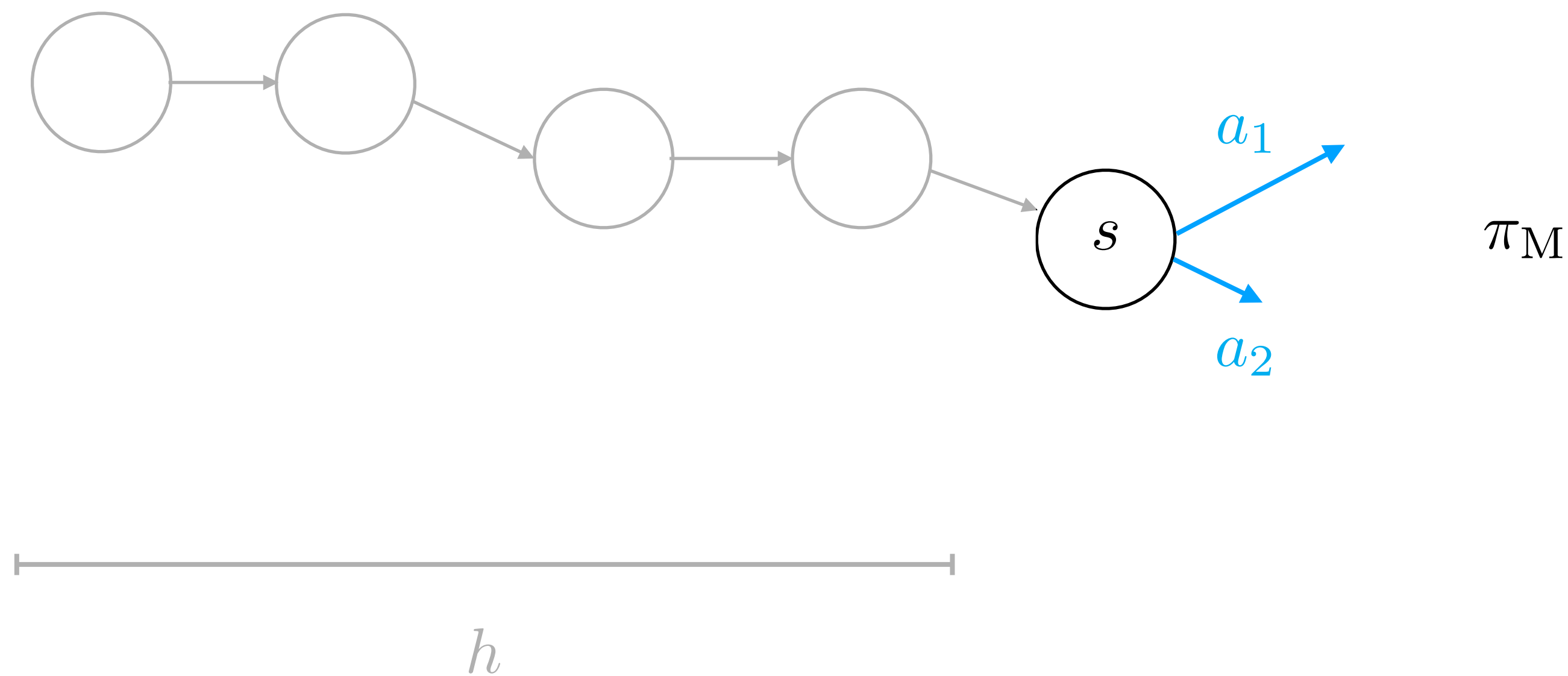
$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$



# Regret Bounds

**Second step.** The optimal Markovian policy  $\pi_M$  is randomized

$$\text{Var} [\text{Ber}(\pi_M(a^* | s, t))] = \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$





# Regret Bounds

Definition (Expected Regret-to-go). Let  $\pi$  be a policy interacting with an MDP over  $T - t$  steps starting from trajectory  $h_t$ . We define the expected regret-to-go as

$$R_{T-t}(\pi, h_t) = H^* - \mathbb{E}_{h_{T-t} \sim p^\pi} [H(d_{h_t \oplus h_{T-t}})]$$

where  $H^*$  is the maximum entropy that can be achieved by any policy starting from  $h_t$ .

**First step.** There exists a deterministic optimal non-Markovian policy  $\pi_{\text{NM}}$

**Second step.** The optimal Markovian policy  $\pi_{\text{M}}$  is randomized

**Result.** The optimal Markovian policy suffers positive regret

# Regret Bounds

**Result.** The optimal Markovian policy suffers positive regret

$$\mathcal{R}_{T-t}(\pi_M, h_t) \propto \text{Var}[\text{Ber}(\pi_M(a^* | s))]$$

# Regret Bounds

**Result.** The optimal Markovian policy suffers positive regret

$$\mathcal{R}_{T-t}(\pi_M, h_t) \propto \text{Var}[\text{Ber}(\pi_M(a^* | s))]$$

$$\propto \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$

# Regret Bounds

**Result.** The optimal Markovian policy suffers positive regret

$$\mathcal{R}_{T-t}(\pi_M, h_t) \propto \text{Var}[\text{Ber}(\pi_M(a^* | s))]$$

$$\propto \text{Var}_{hs \sim p_t^{\pi_{NM}}} [\mathbb{E} [\text{Ber}(\pi_{NM}(a^* | hs))]]$$

$\implies$  **non-Markovianity** matters in **finite-sample** MSE

# Computational Tractability

Learning the optimal Markovian policy for MSE is known to be provably efficient<sup>1,2</sup>

<sup>1</sup>(Hazan et al., 2019), <sup>2</sup>(Zhang et al., 2020)

# Computational Tractability

Learning the optimal Markovian policy for MSE is known to be provably efficient<sup>1,2</sup>

Is computing the optimal non-Markovian policy for the finite-sample MSE even tractable?

<sup>1</sup>(Hazan et al., 2019), <sup>2</sup>(Zhang et al., 2020)

# Computational Tractability

Learning the optimal Markovian policy for MSE is known to be provably efficient<sup>1,2</sup>

Is computing the optimal non-Markovian policy for the finite-sample MSE even tractable?

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*

<sup>1</sup>(Hazan et al., 2019), <sup>2</sup>(Zhang et al., 2020)



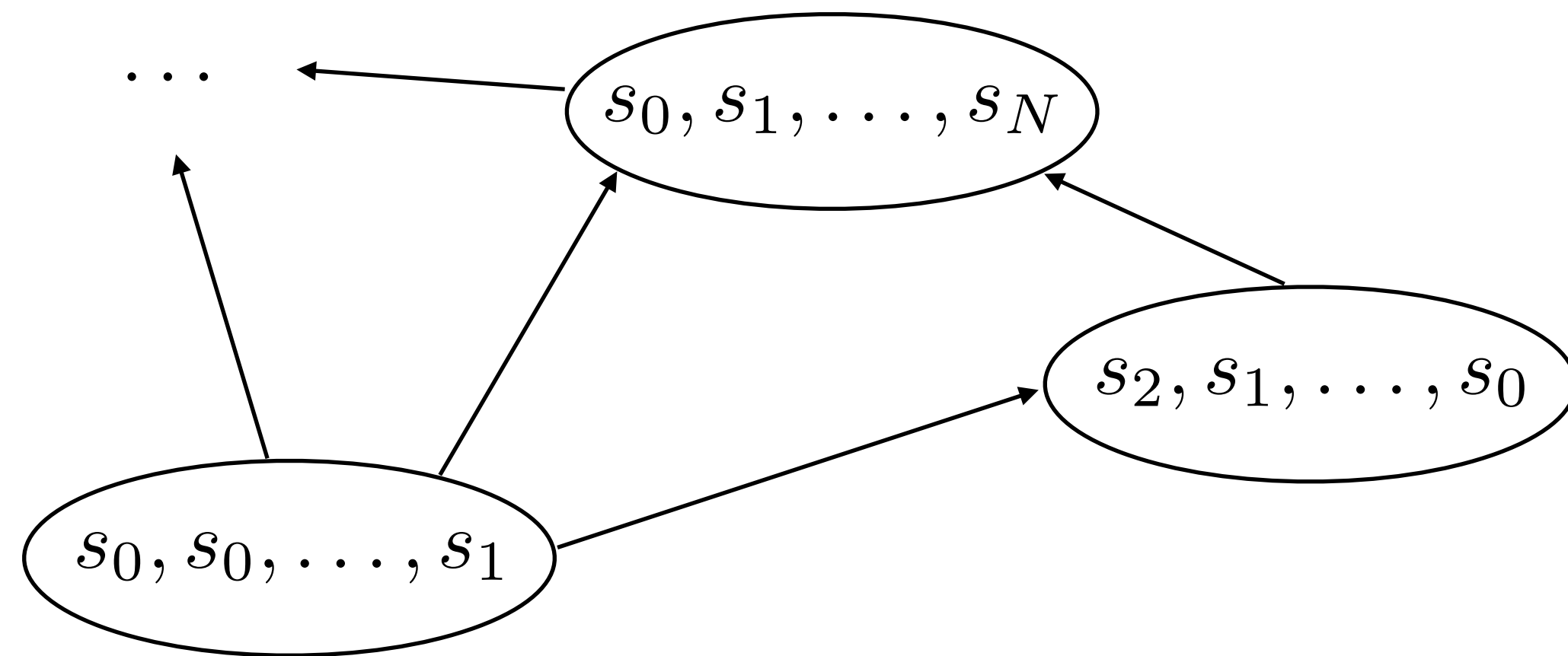
# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*

through reduction to POMDP

# Computational Tractability

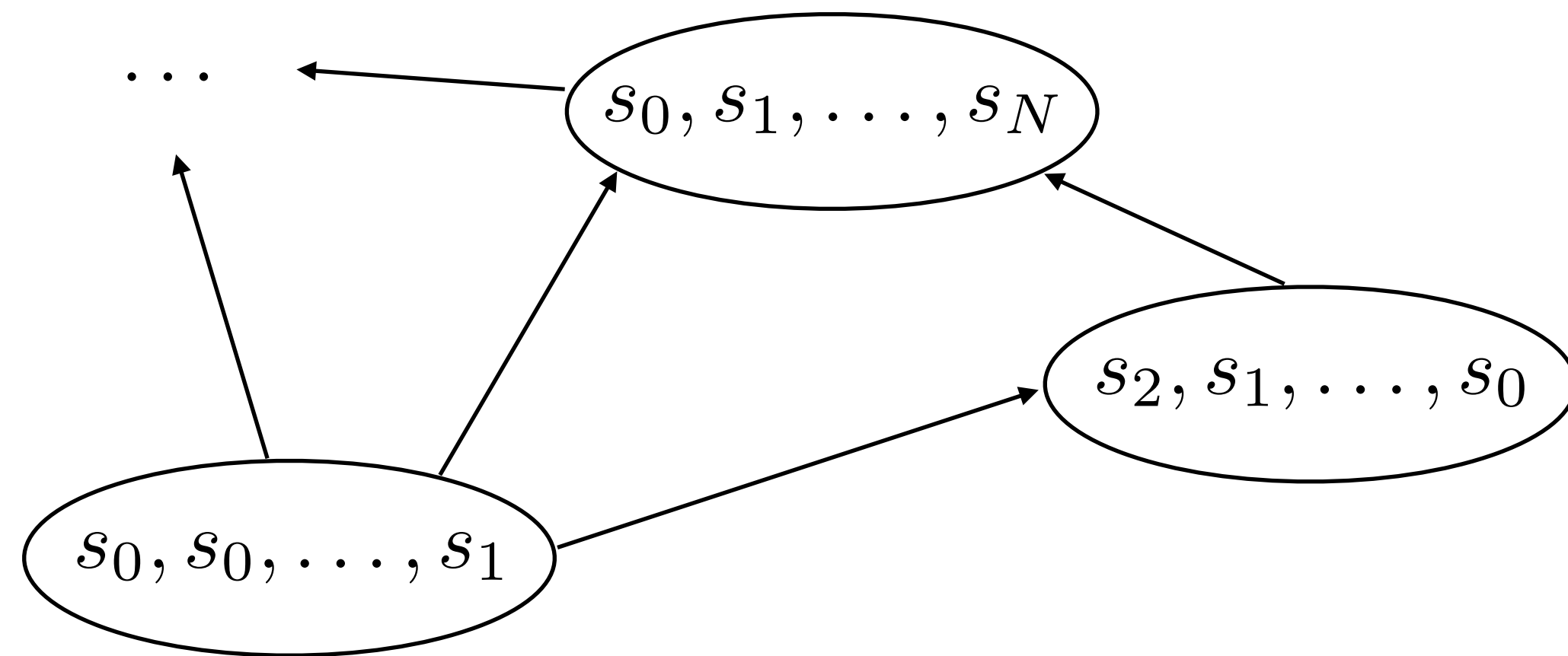
Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is NP-hard.*



**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*

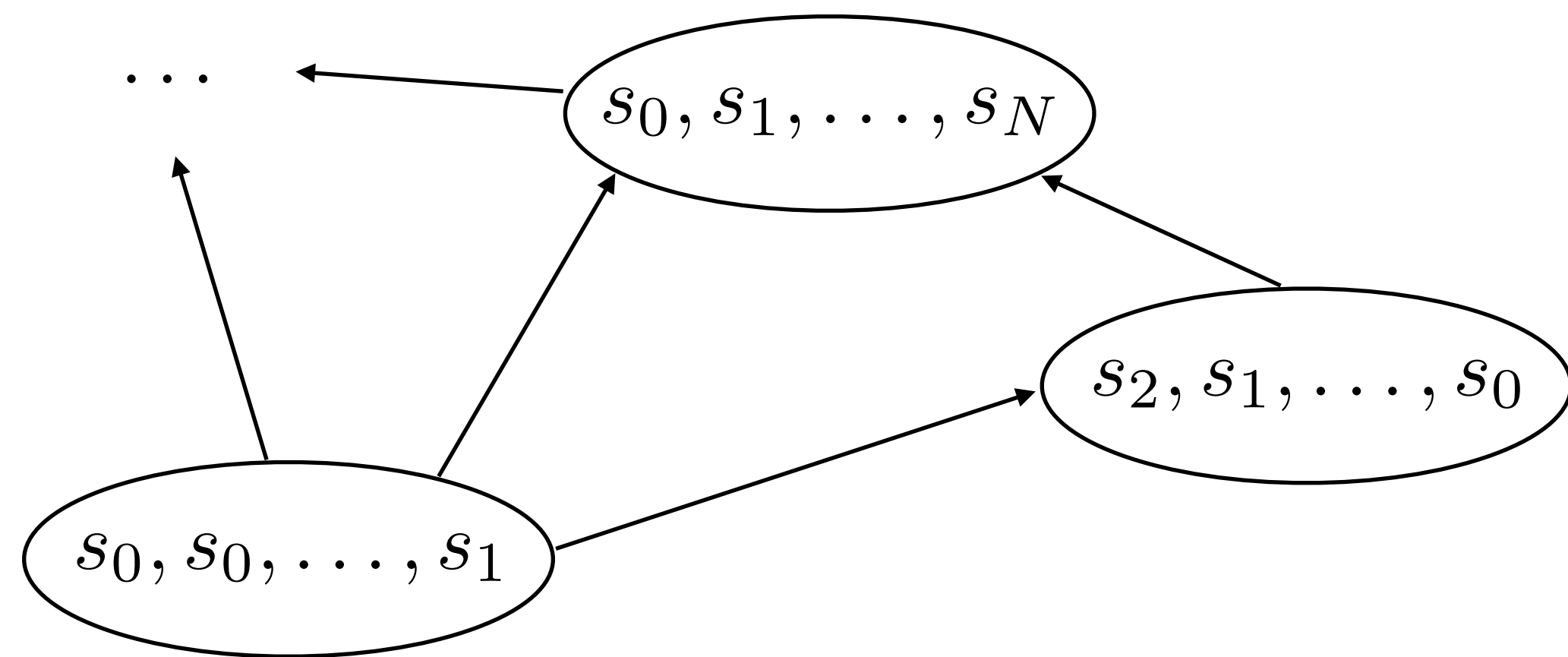


**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

exponential blowup with the horizon

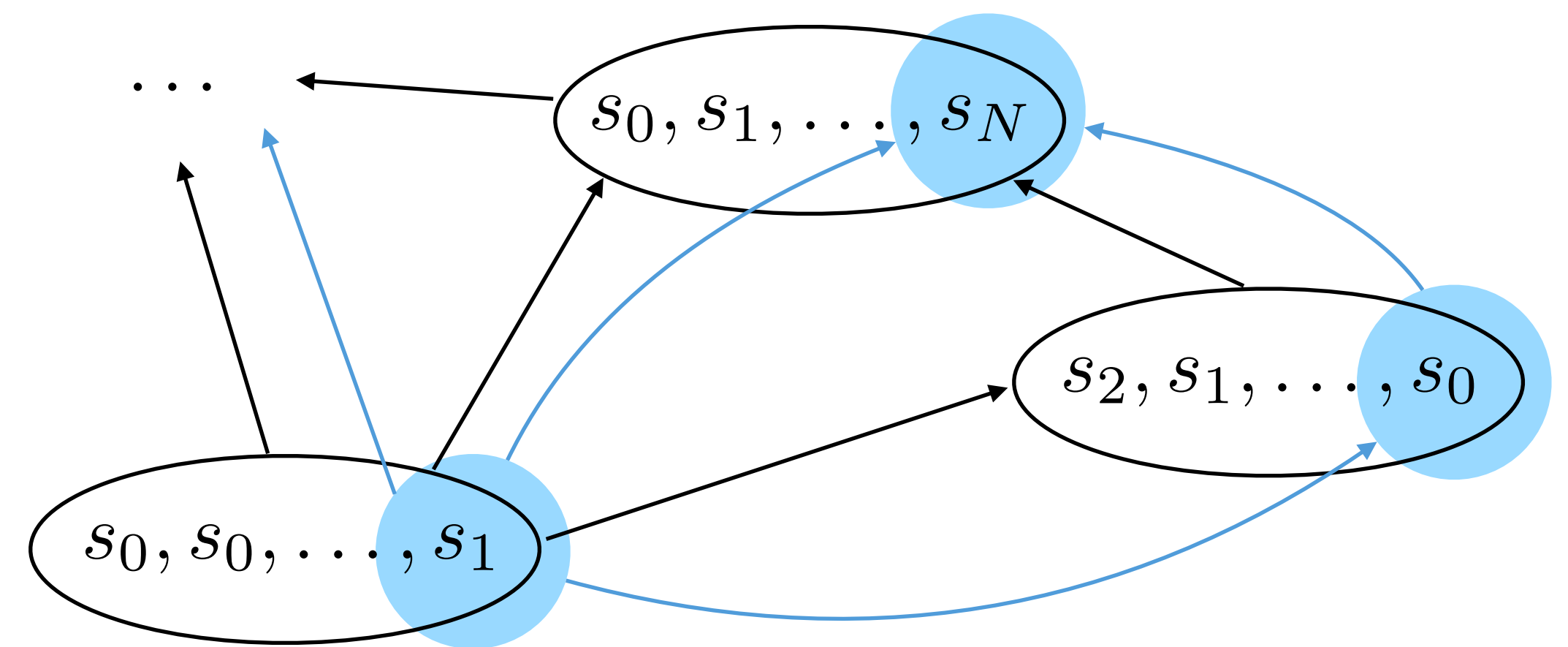
# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*



**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

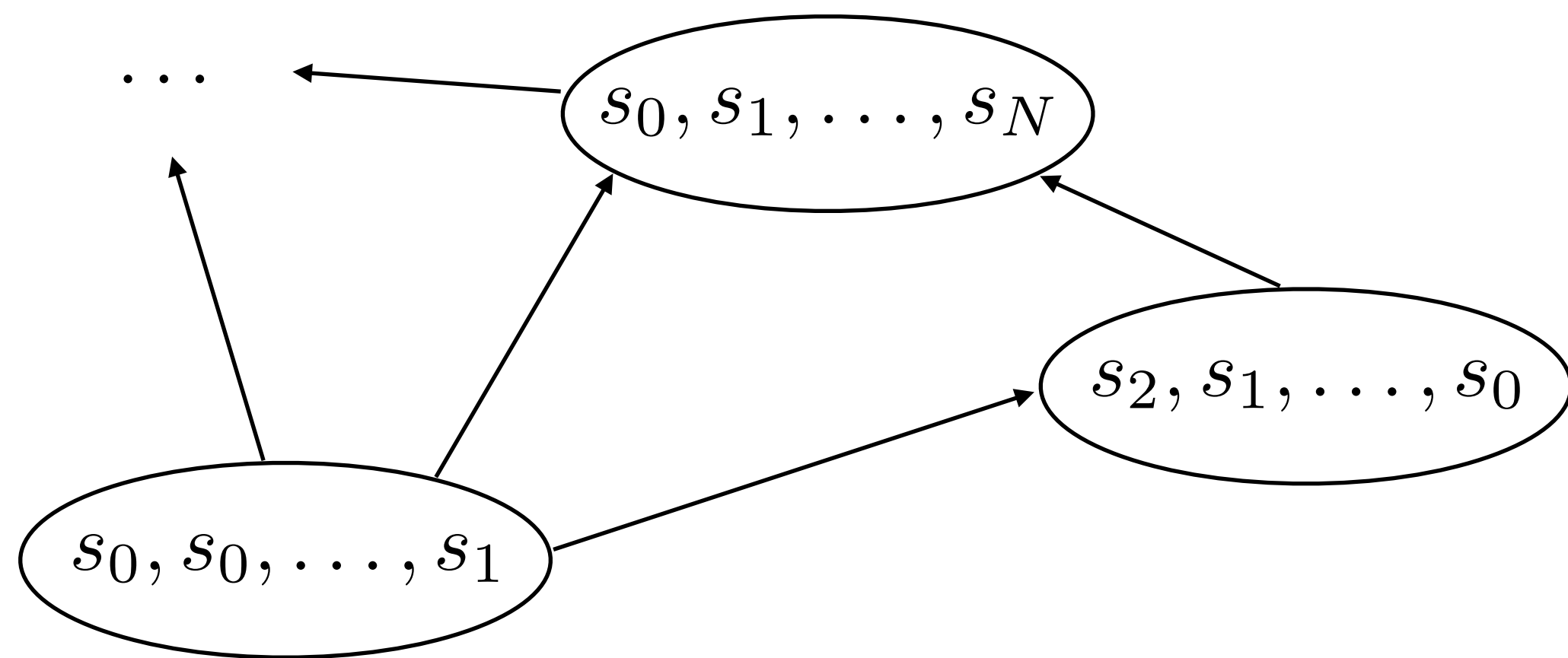
exponential blowup with the horizon



Reduction to a class of **POMDPs**

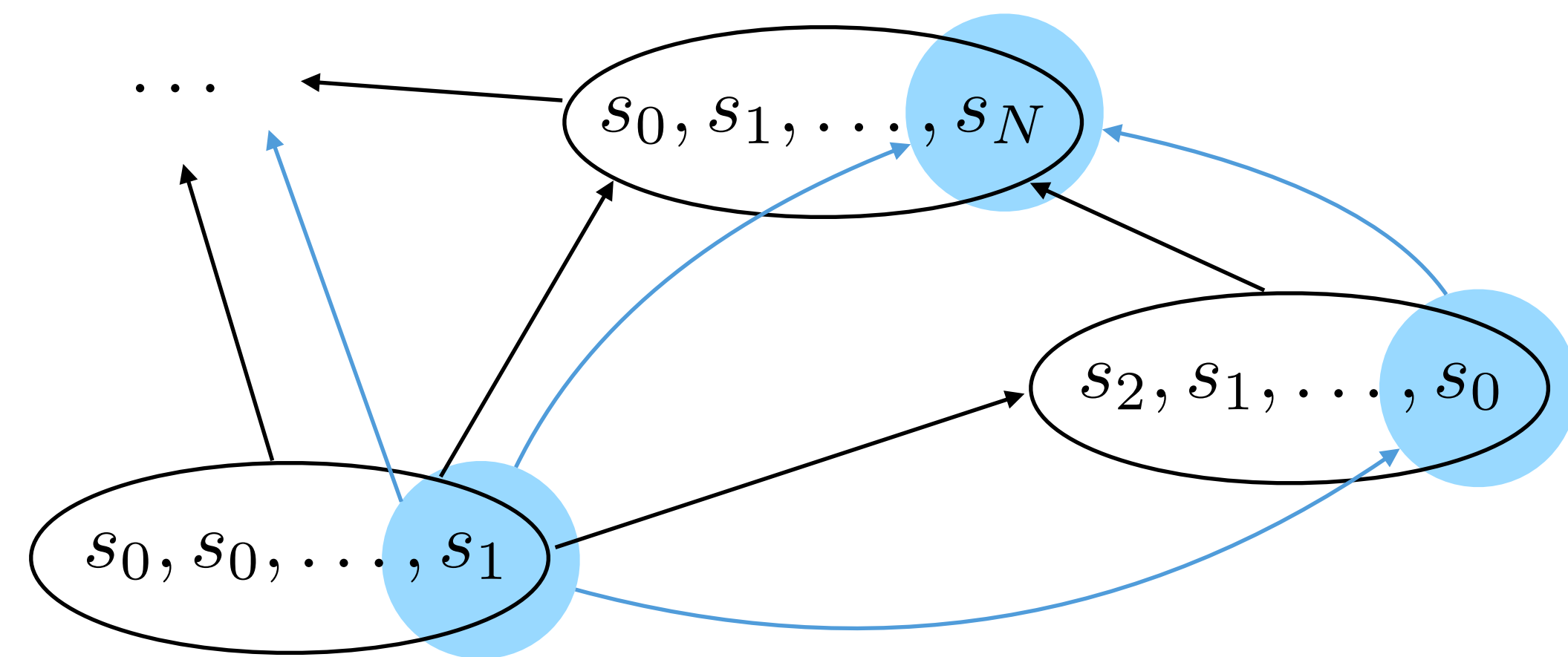
# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*



**Extended CMP** with reward  $R(h_T) = H(d_{h_T})$

exponential blowup with the horizon



Reduction to a class of **POMDPs**  $\geq_p$  **3SAT**

NP-hard problem<sup>1</sup>

<sup>1</sup>(Mundhenk et al., 2000)

# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*

Can we solve the problem with function approximation?



# Computational Tractability

Theorem (Computational Complexity). *Optimizing the finite-sample MSE within the space of non-Markovian policies is **NP-hard**.*

Can we solve the problem with function approximation?

DARL Workshop @ ICML & Pre-Training Workshop @ ICML

Mutti et al. "**Non-Markovian Policies for Unsupervised Reinforcement Learning in Multiple Environments**". 2022.



# Take Home

**Non-Markovian** policies are better for **finite-sample** convex objectives

**Optimizing** non-Markovian policies exactly is often **intractable**

# Take Home

**Non-Markovian** policies are better for **finite-sample** convex objectives

**Optimizing** non-Markovian policies exactly is often **intractable**

## What Is Next?

**Approximate methods** to optimize non-Markovian policies for convex objectives

**Applications:** When is it critical to consider a finite-sample objective?

# References

- (Williams & Zipser, 1989) A learning algorithm for continually running fully recurrent neural networks. *Neural computation*, 1989.
- (Hochreiter & Schmidhuber, 1997) Long short-term memory. *Neural computation*, 1997.
- (Puterman, 2014) *Markov decision processes: Discrete stochastic dynamic programming*. 2014.
- (Zhang et al., 2020) Variational policy gradient method for reinforcement learning with general utilities. *NeurIPS*, 2020.
- (Zahavy et al., 2021) Reward is enough for convex MDPs. *NeurIPS*, 2021.
- (Hazan et al., 2019) Provably efficient maximum entropy exploration. *ICML 2019*.
- (Mundhenk et al., 2000) Complexity of finite-horizon Markov decision process problems. *Journal of the ACM*, 2020.

# Questions?

## Poster

Hall E #838

## Paper

