

On the Curses of Horizon in Off-policy Evaluation in non-Markov Environments

Nan Jiang

University of Illinois at Urbana-Champaign

September 2024

@Simons Institute

Based on: (1) Uehara et al. NeurIPS 2023. <https://arxiv.org/pdf/2207.13081.pdf>
(2) Zhang and Jiang. 2024. <https://arxiv.org/pdf/2402.14703.pdf>



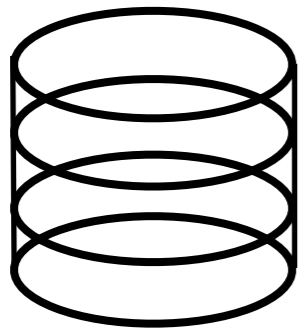
Masatoshi
Uehara



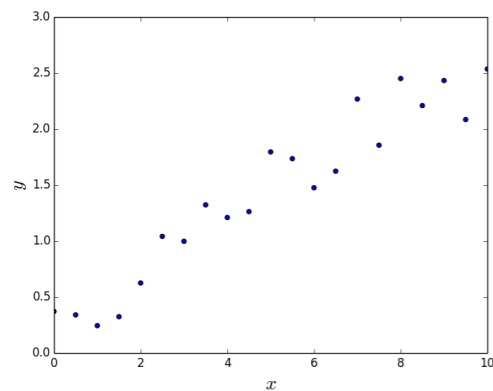
Yuheng
Zhang

Generalization in Prediction

- How do we know if an algorithm generalizes?

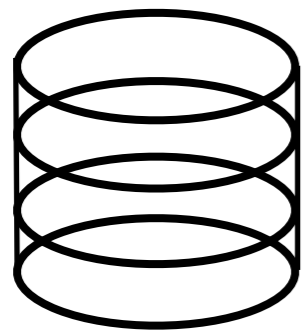


data

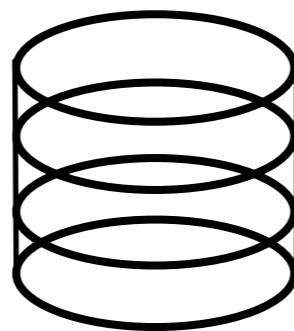
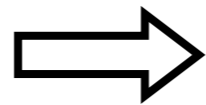


Generalization in Prediction

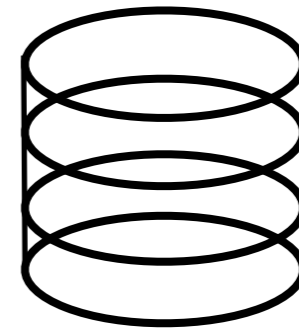
- How do we know if an algorithm generalizes?



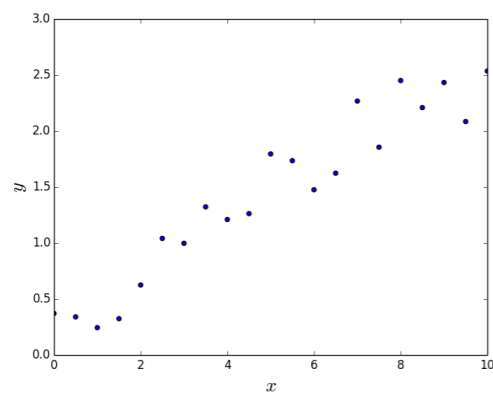
data



training

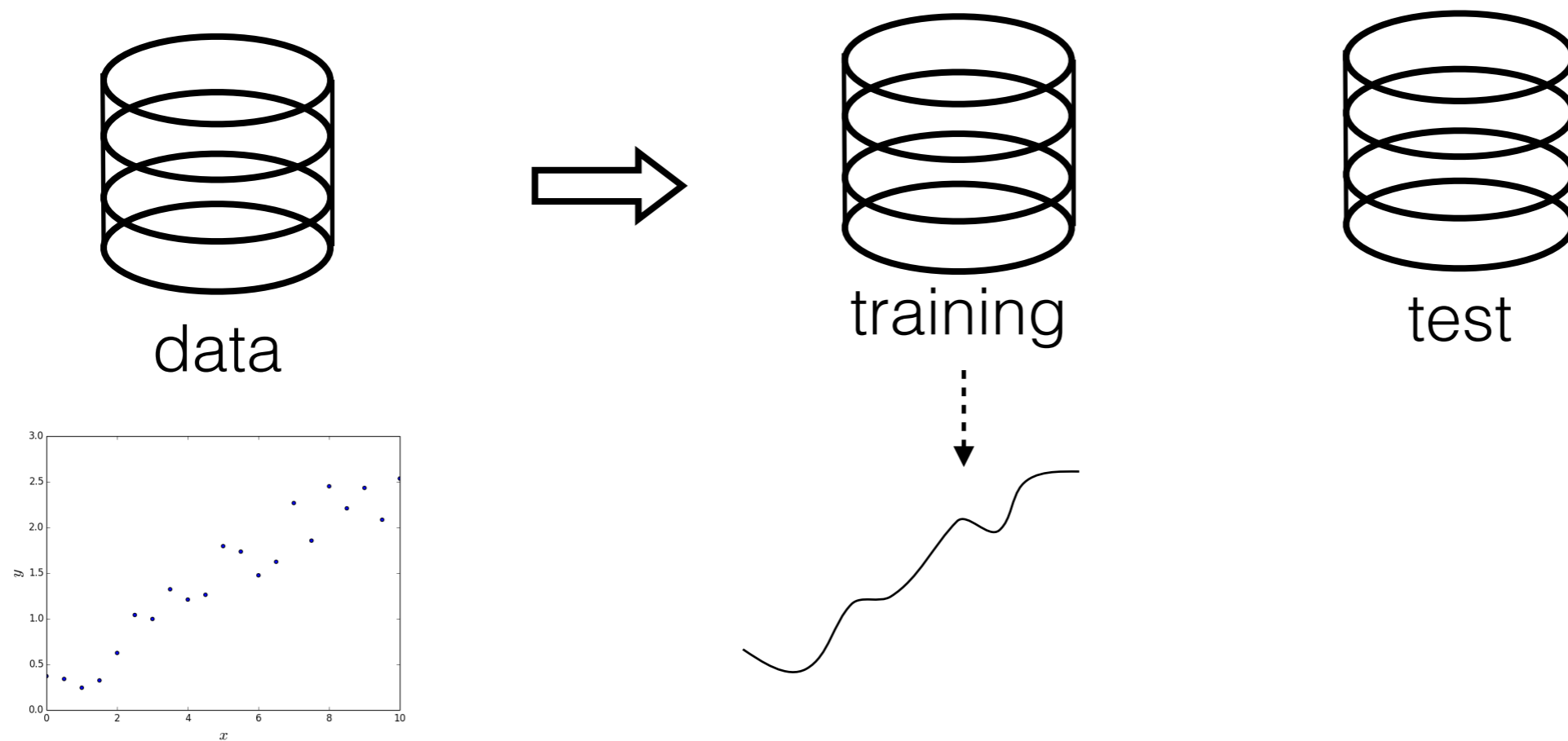


test



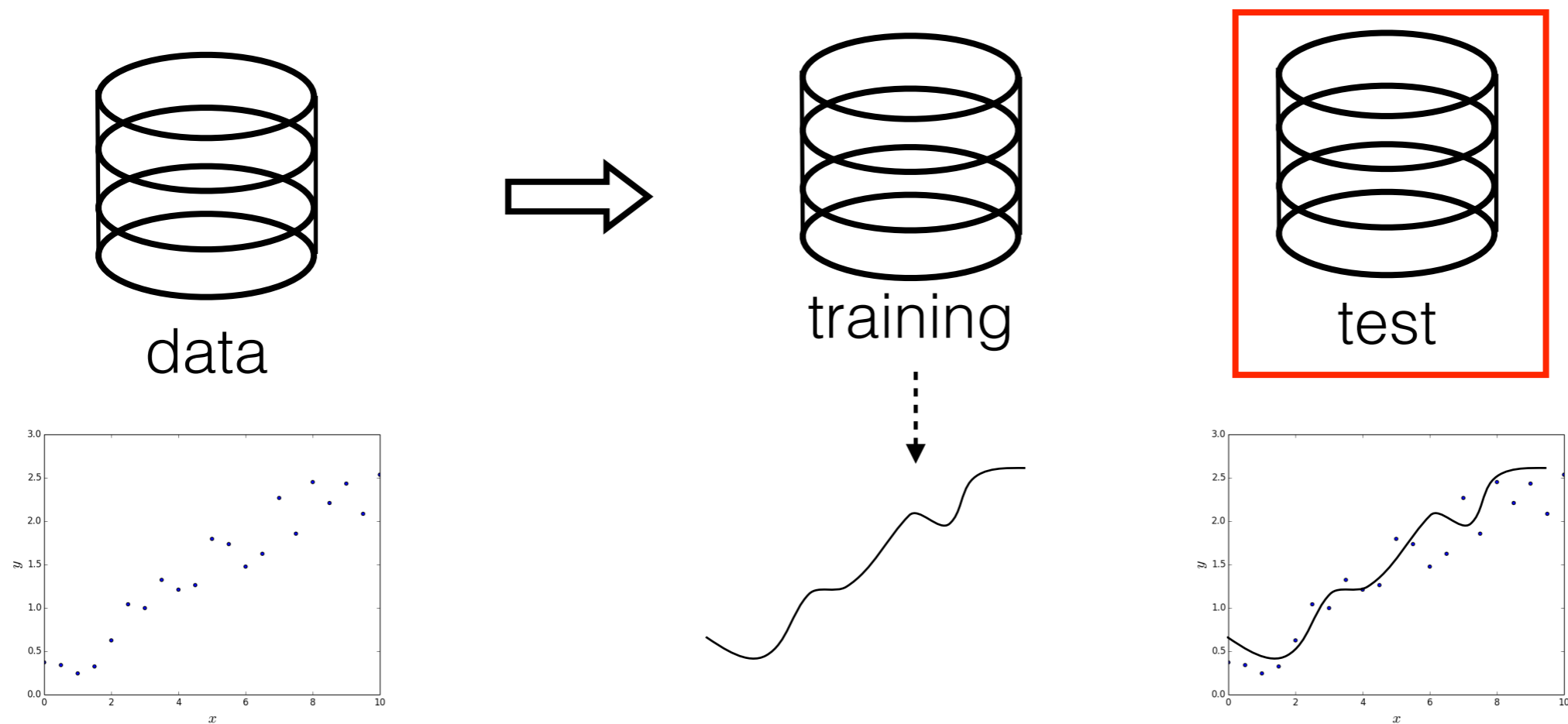
Generalization in Prediction

- How do we know if an algorithm generalizes?



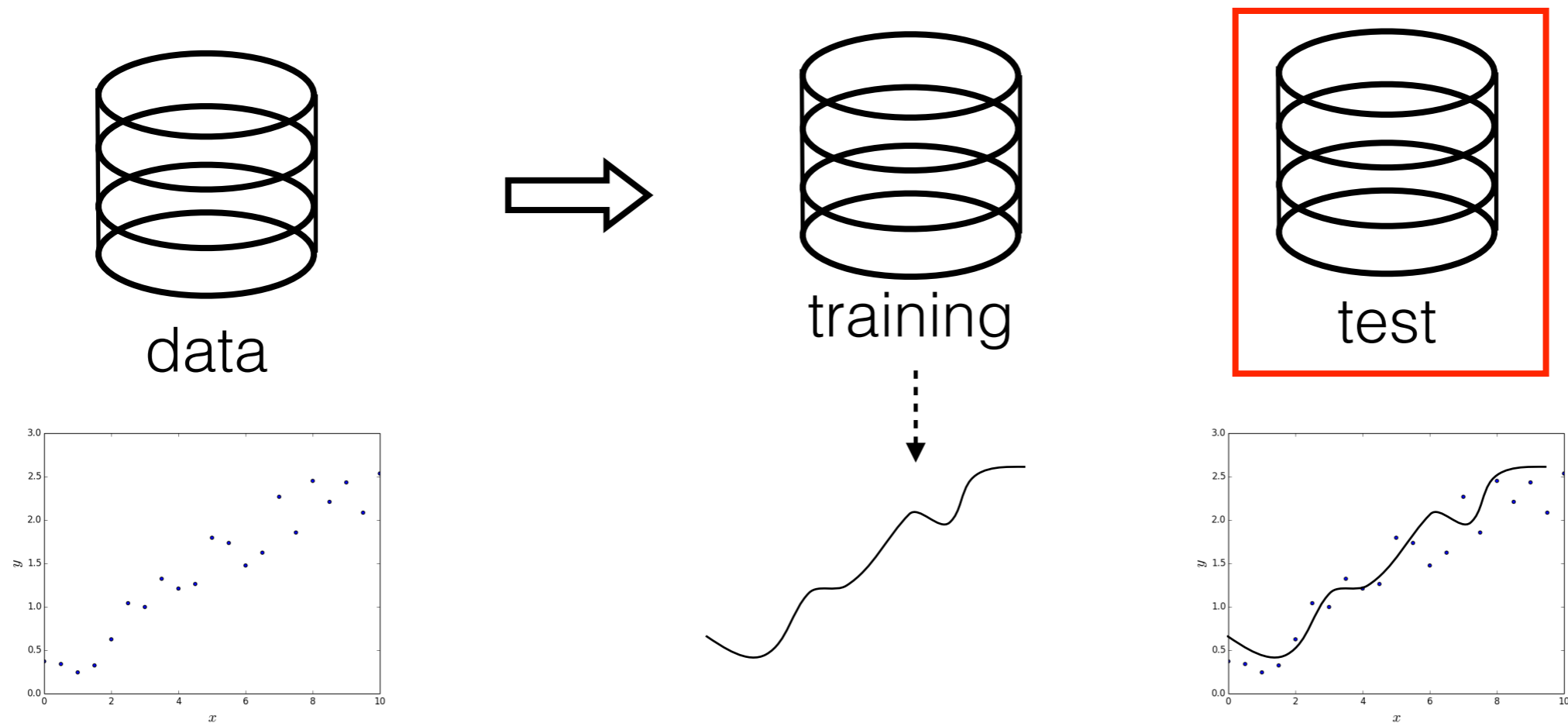
Generalization in Prediction

- How do we know if an algorithm generalizes?



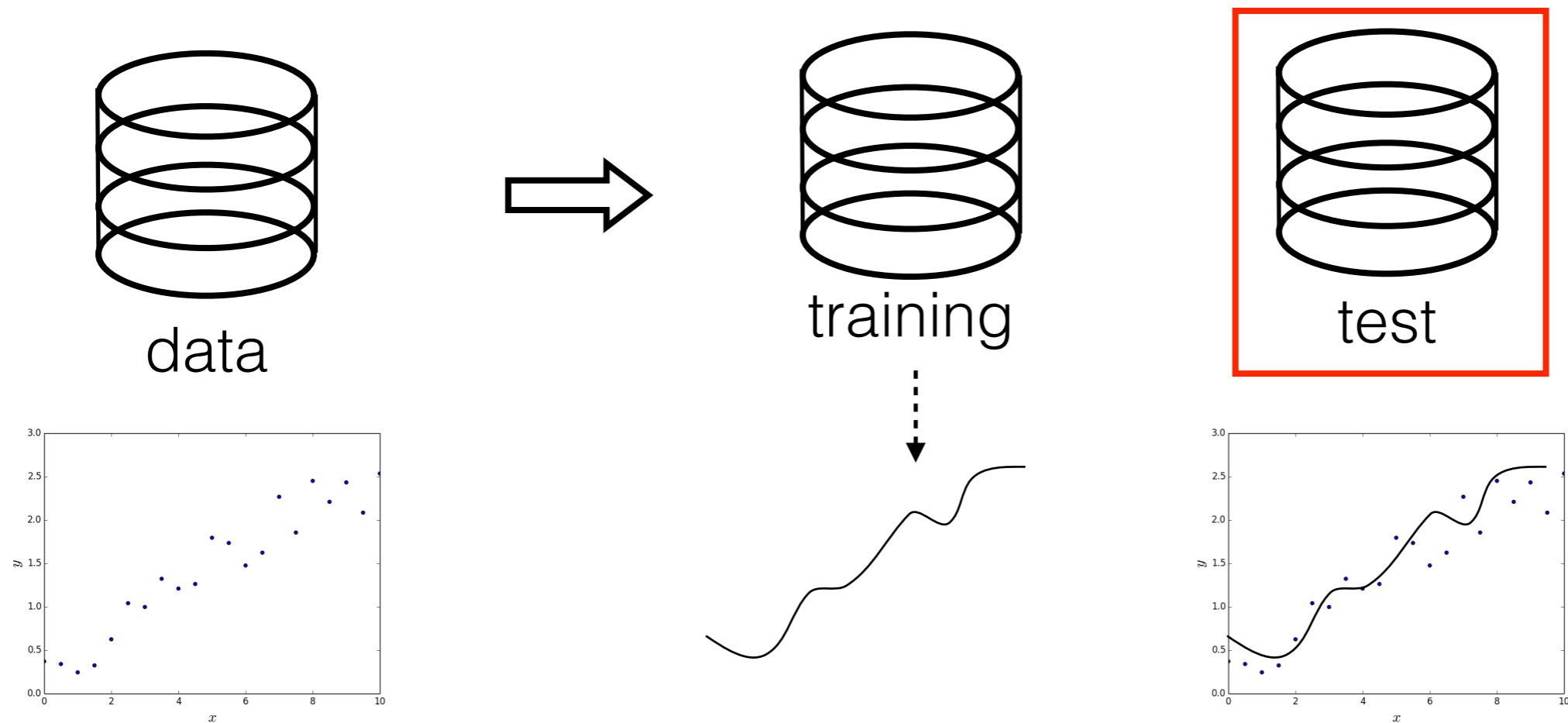
Generalization in Prediction

- How do we know if an algorithm generalizes?
 - Training error \approx test error (calculated from test data)



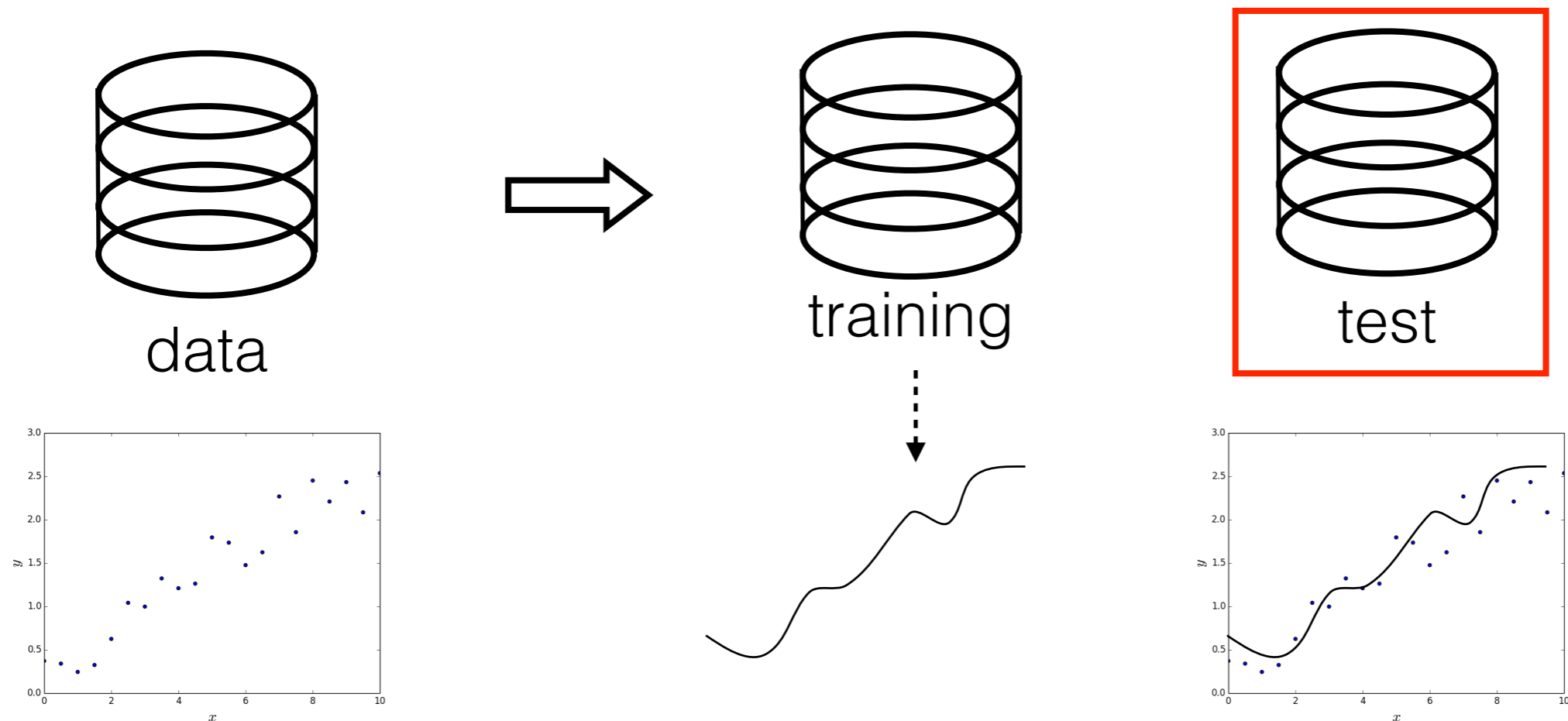
Generalization in Prediction

- How do we know if an algorithm generalizes?
 - Training error \approx test error (calculated from test data)
- When generalization happens?

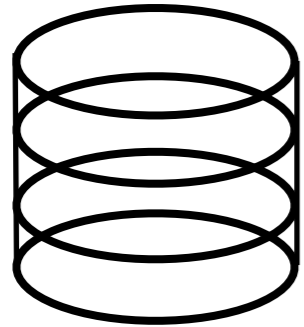


Generalization in Prediction

- How do we know if an algorithm generalizes?
 - Training error \approx test error (calculated from test data)
- When generalization happens?
 - Sufficient training data ($>$ capacity of function class)

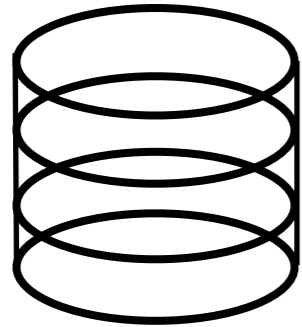


Generalization in Decision-Making

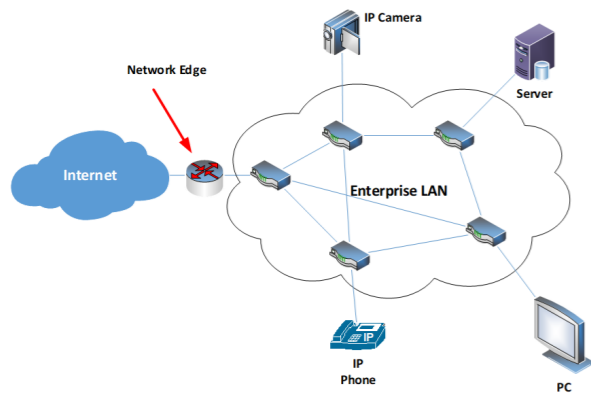
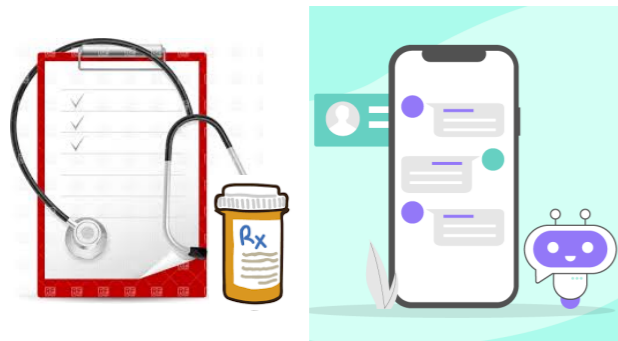


data

Generalization in Decision-Making

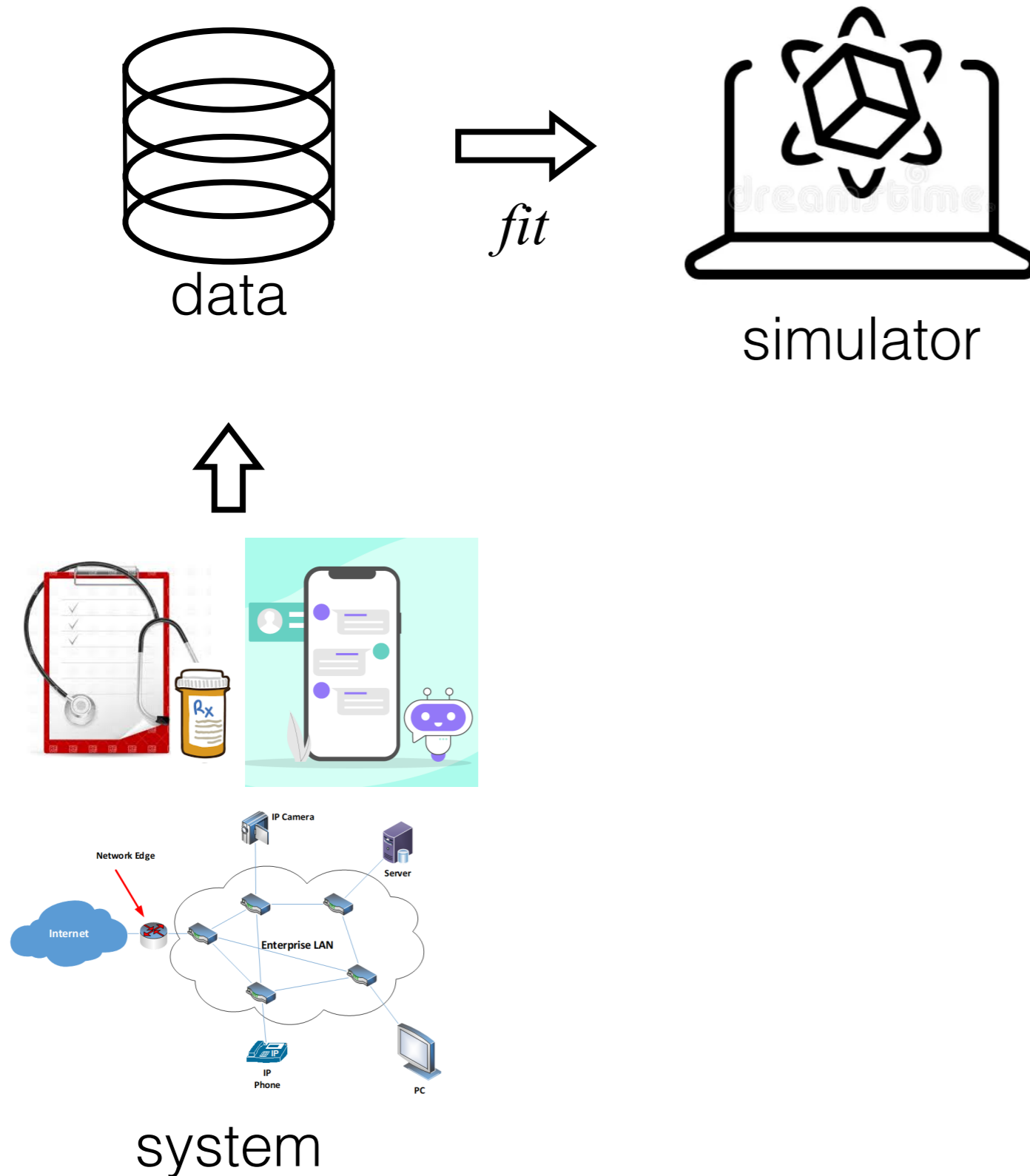


data

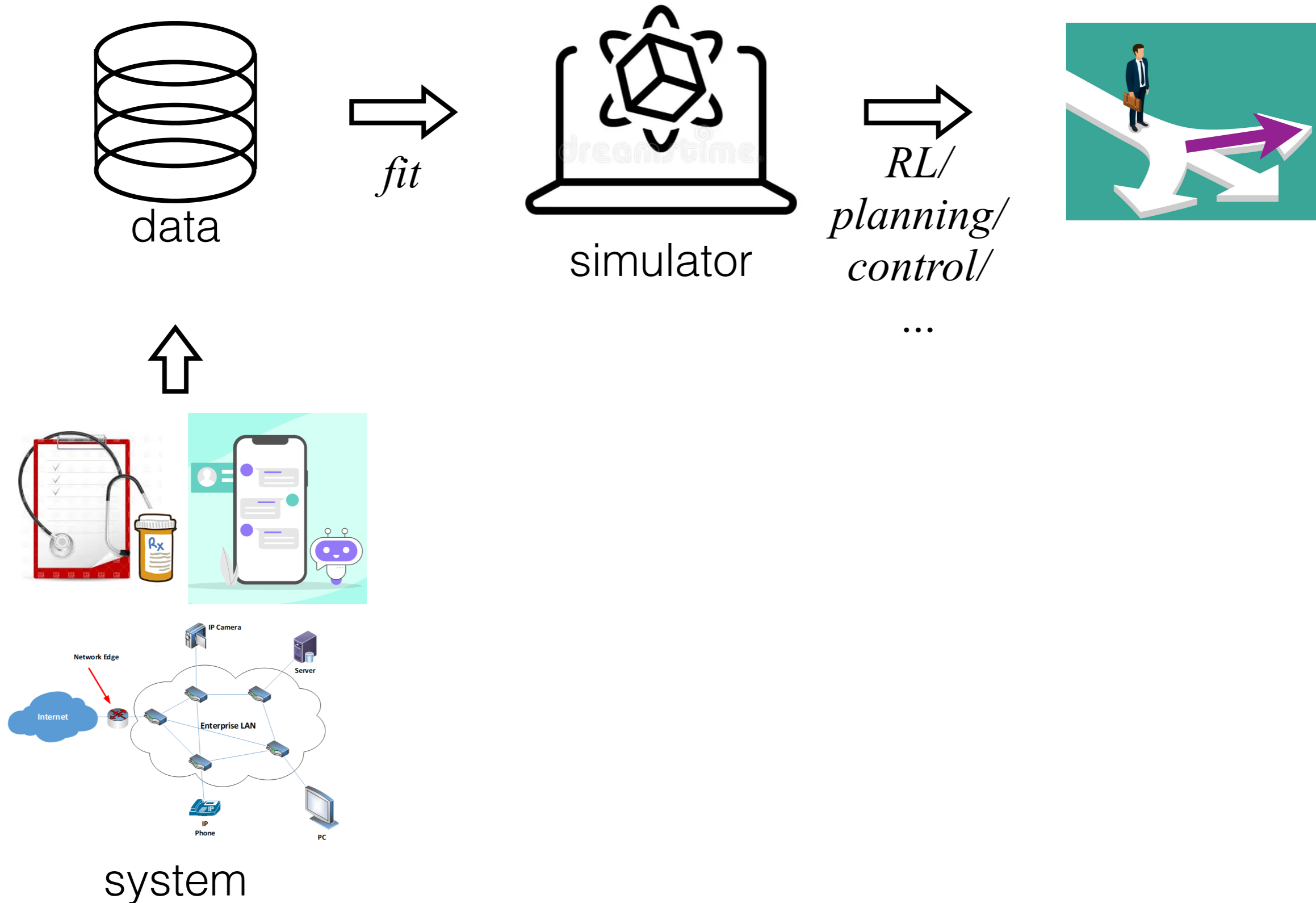


system

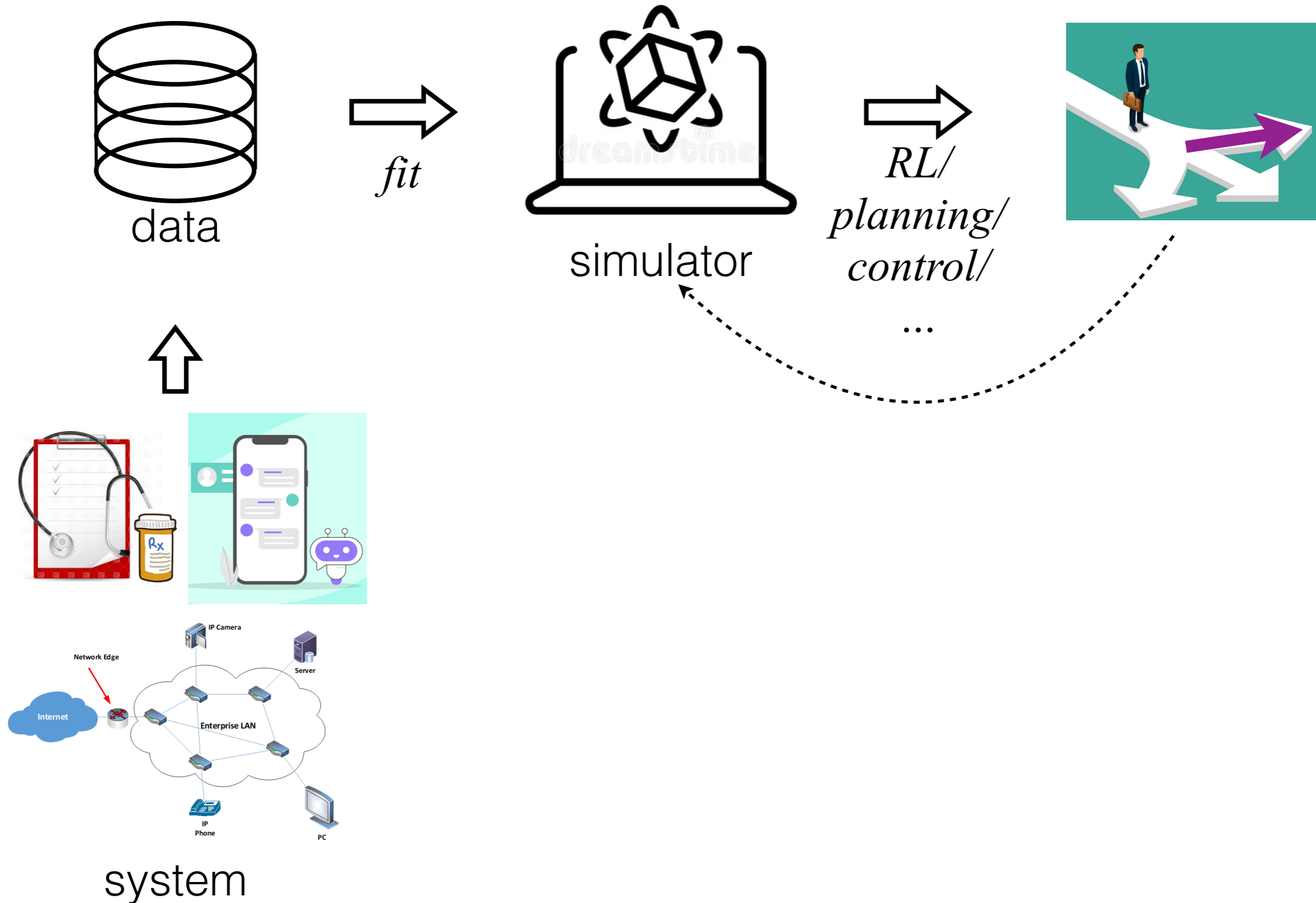
Generalization in Decision-Making



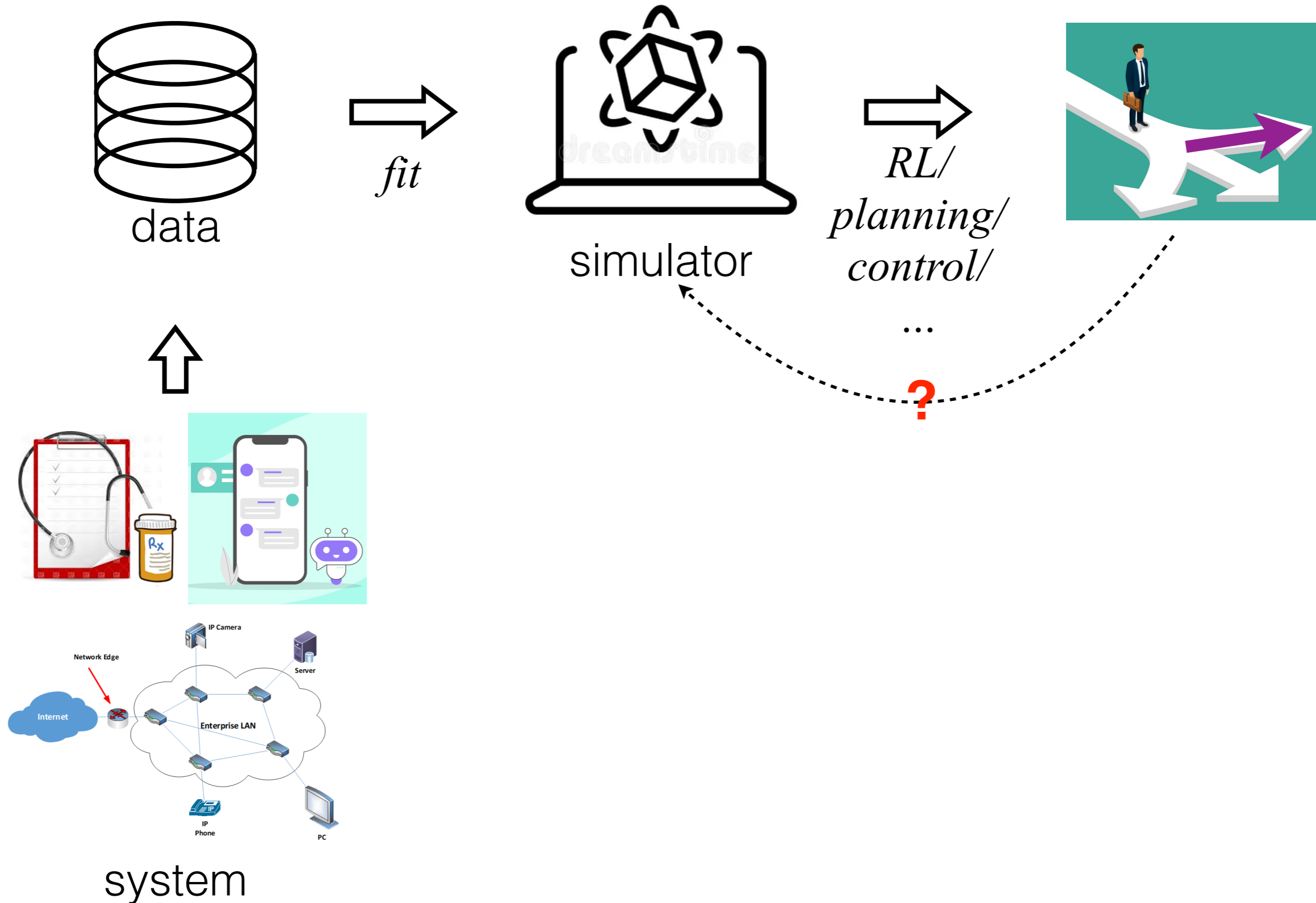
Generalization in Decision-Making



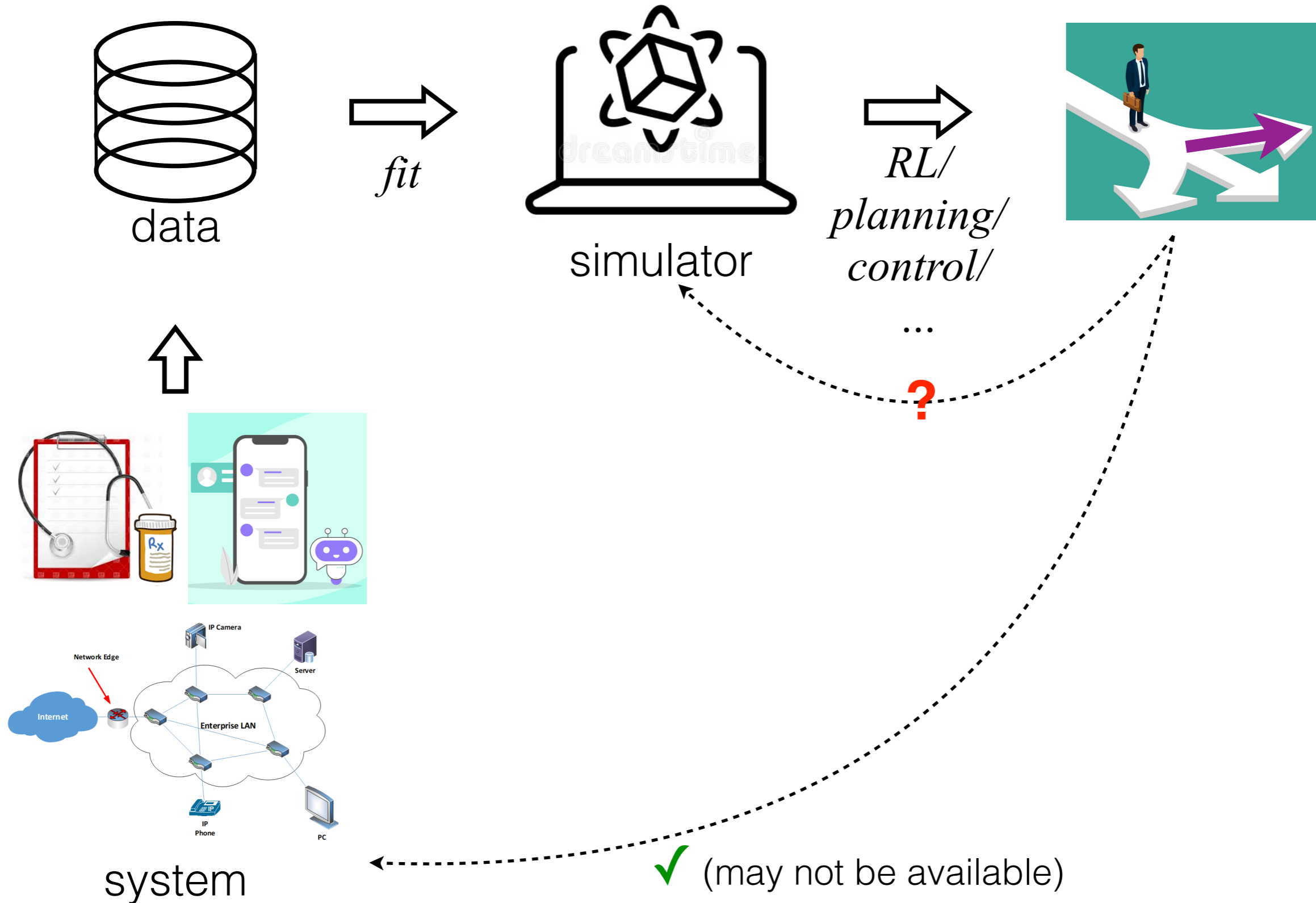
Generalization in Decision-Making



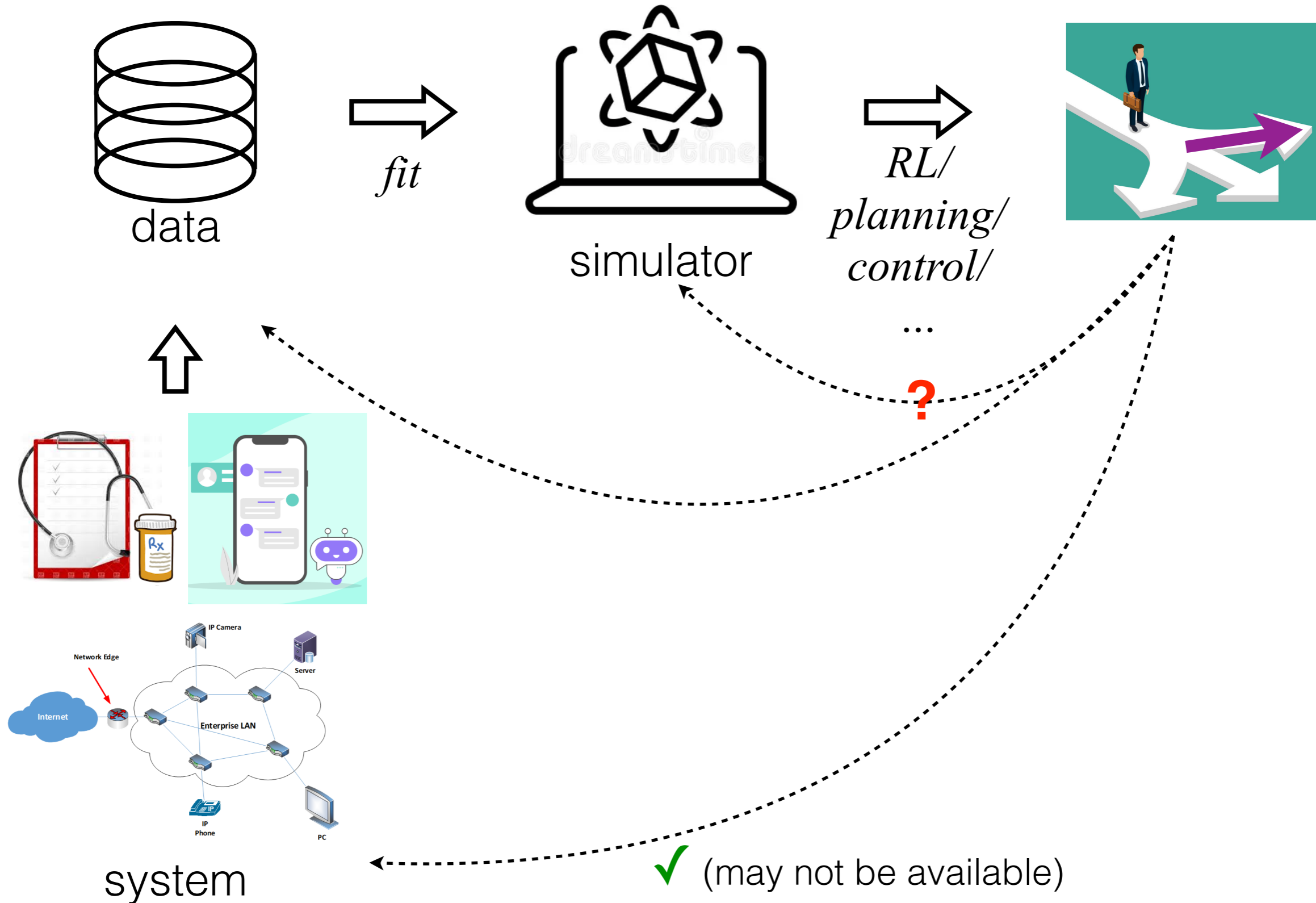
Generalization in Decision-Making



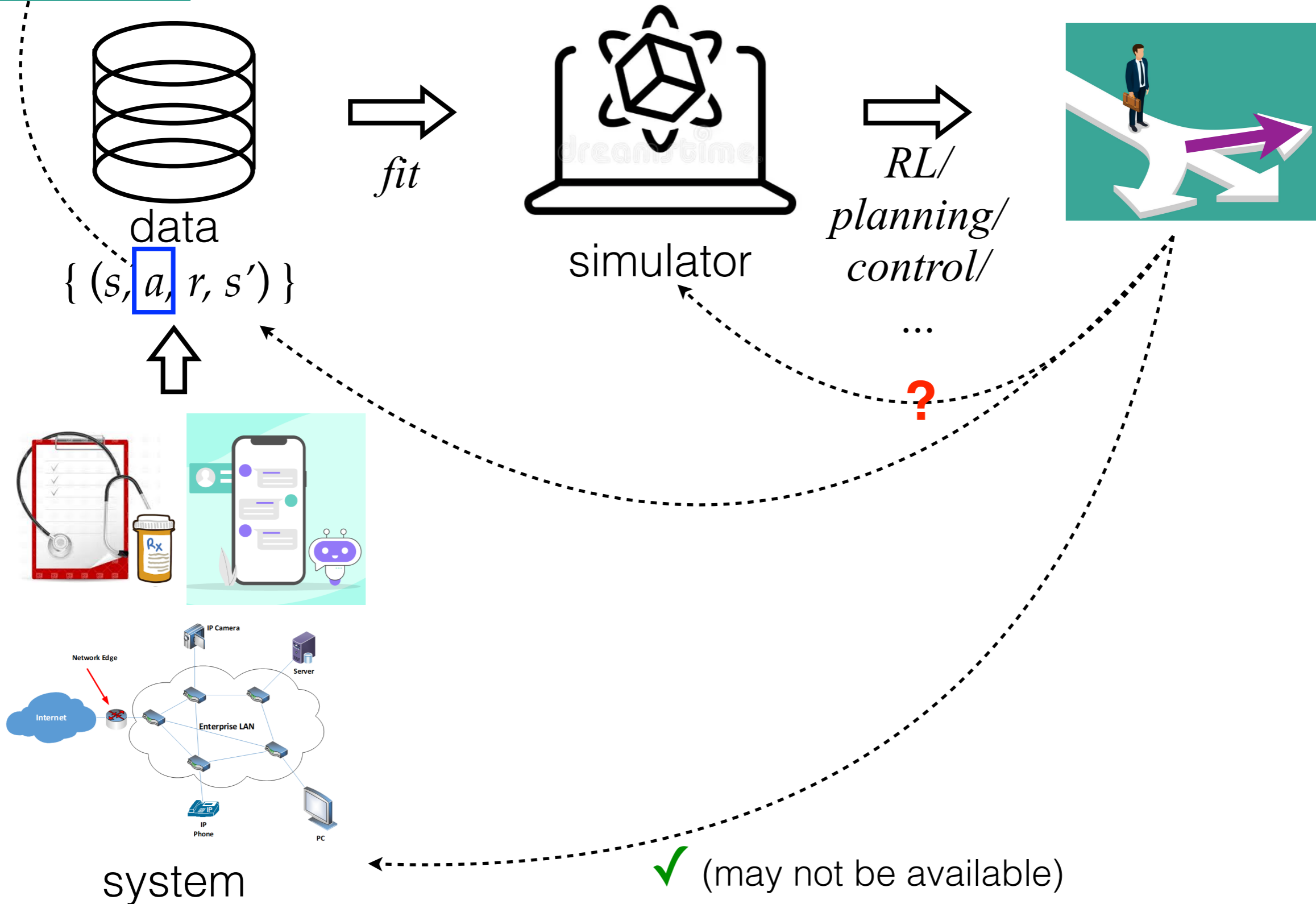
Generalization in Decision-Making



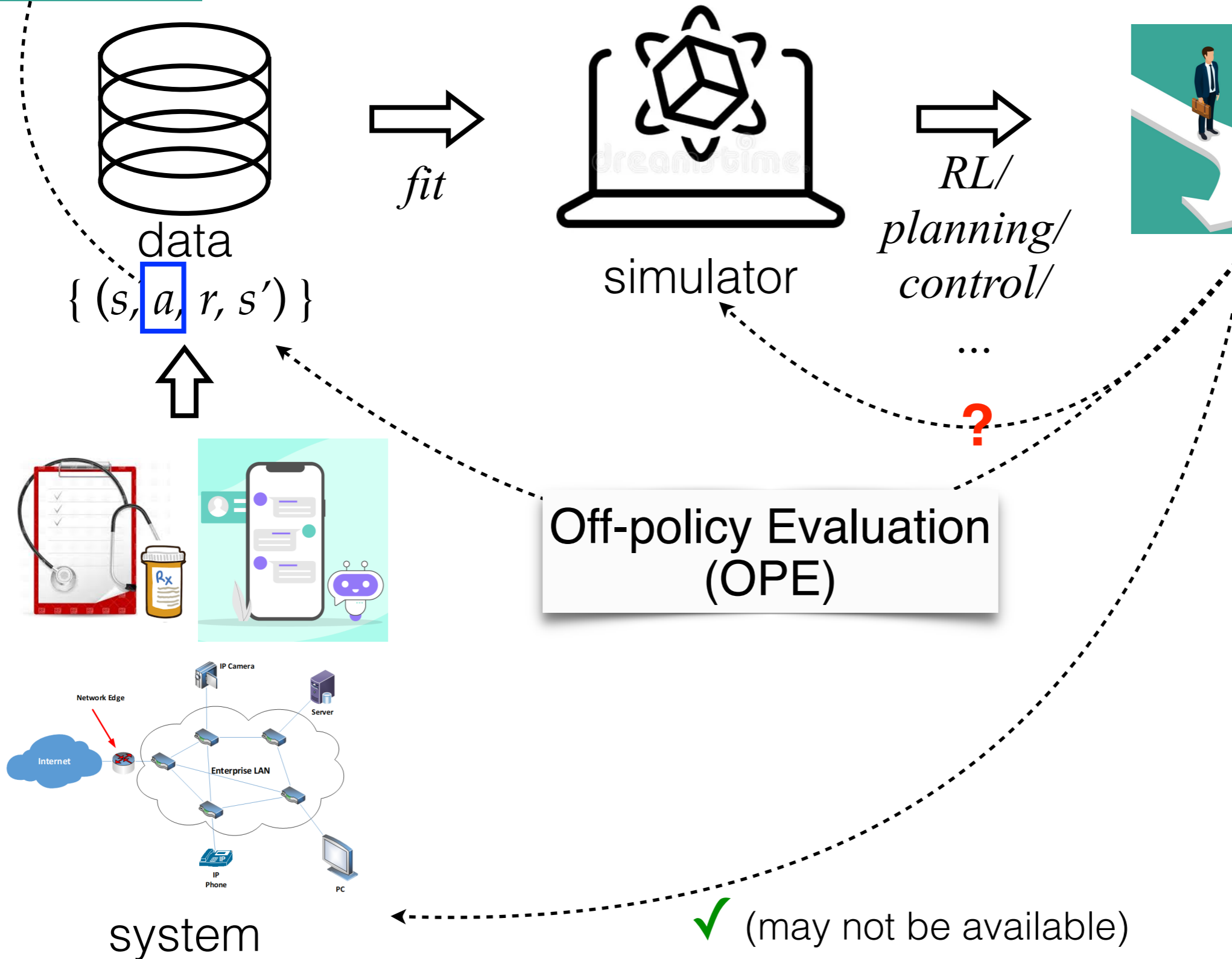
Generalization in Decision-Making



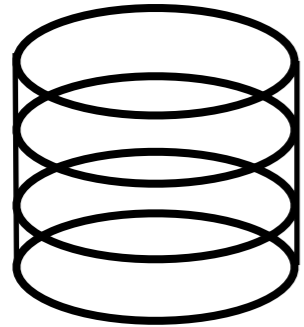
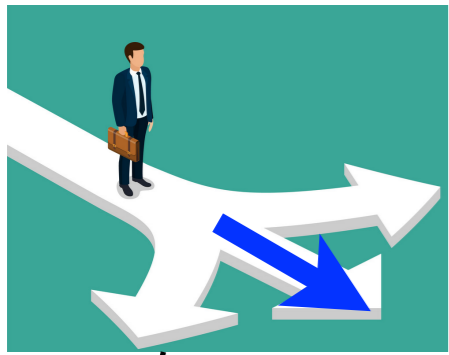
Generalization in Decision-Making



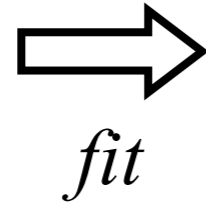
Generalization in Decision-Making



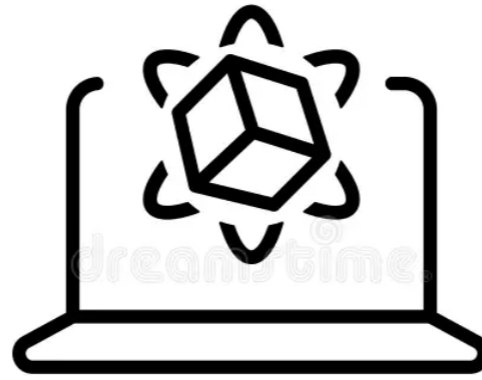
Example: RLHF in LLMs



data



fit



simulator

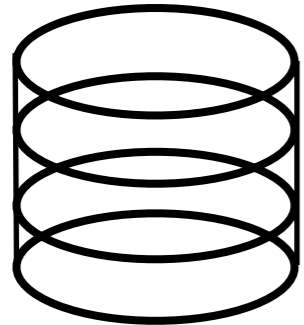


*RL/
planning/
control/*

...

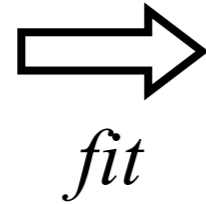


Example: RLHF in LLMs

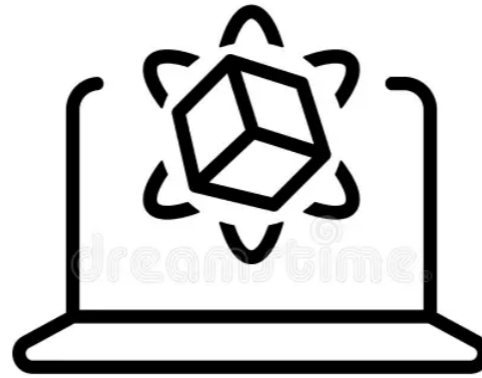


data

{prompt,
response1,
response2,
winner}



fit



simulator

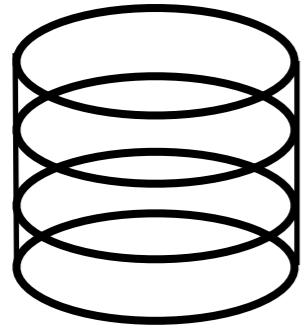
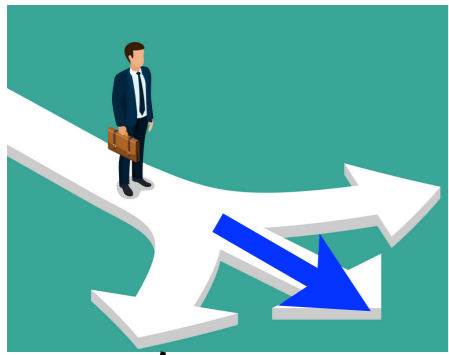


*RL/
planning/
control/*

...

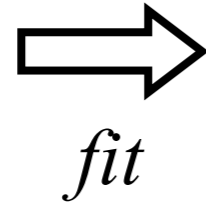


Example: RLHF in LLMs

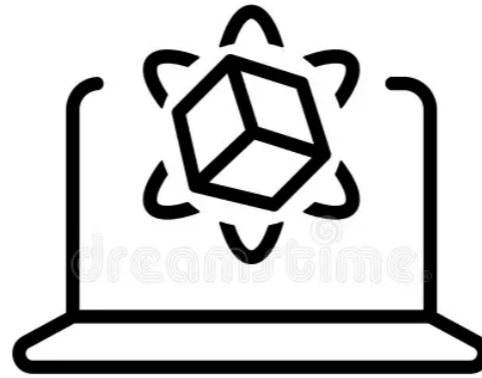


data

{prompt,
response,
reward}



fit



simulator

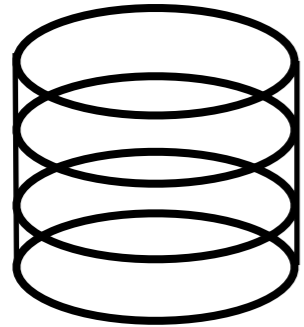


*RL/
planning/
control/*

...

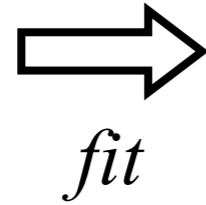


Example: RLHF in LLMs

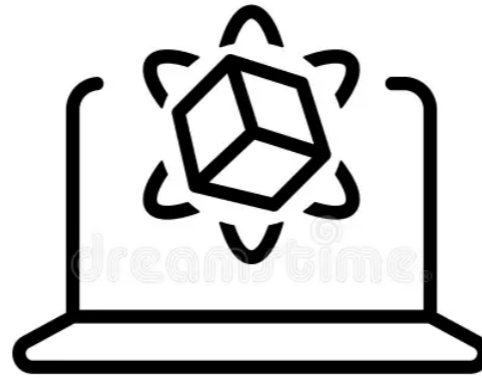


data

{prompt,
response,
reward}



fit



simulator

reward model

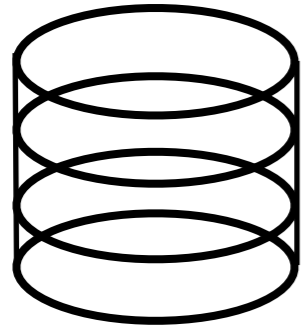


*RL/
planning/
control/*

...

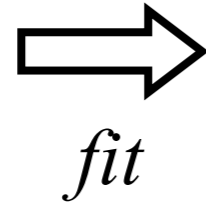


Example: RLHF in LLMs

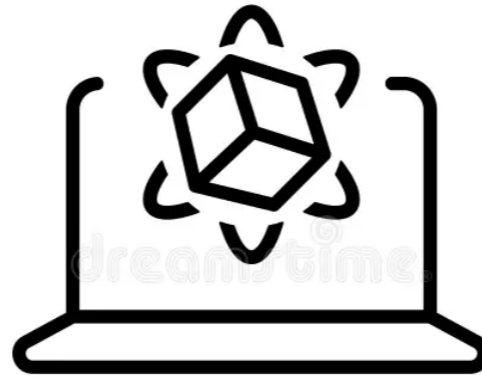


data

{prompt,
response,
reward}



fit



simulator

reward model



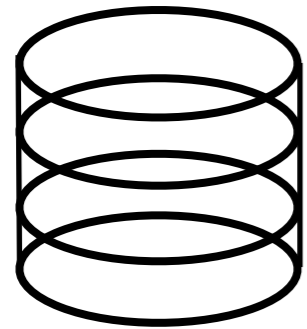
*RL/
planning/
control/*

...



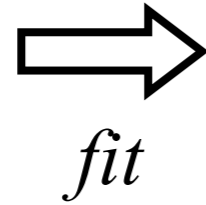
PPO-trained policy

Example: RLHF in LLMs

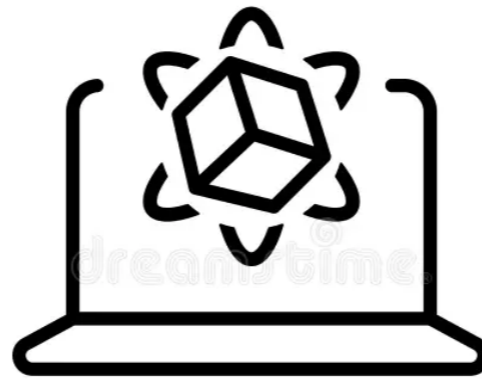


data

{prompt,
response,
reward}



fit



simulator

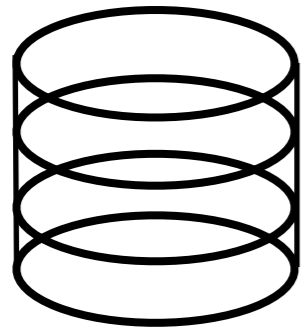
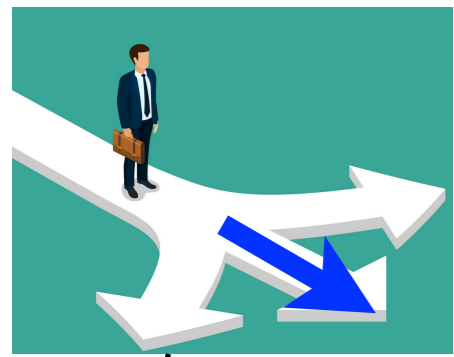
reward model

*RL/
planning/
control/
...*



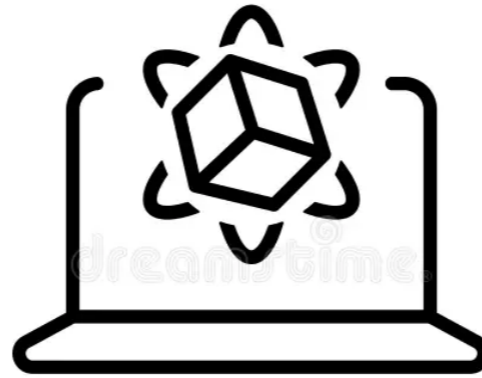
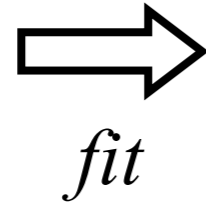
PPO-trained policy

Example: RLHF in LLMs



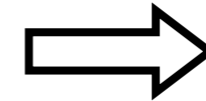
data

{prompt,
response,
reward}



simulator

reward model

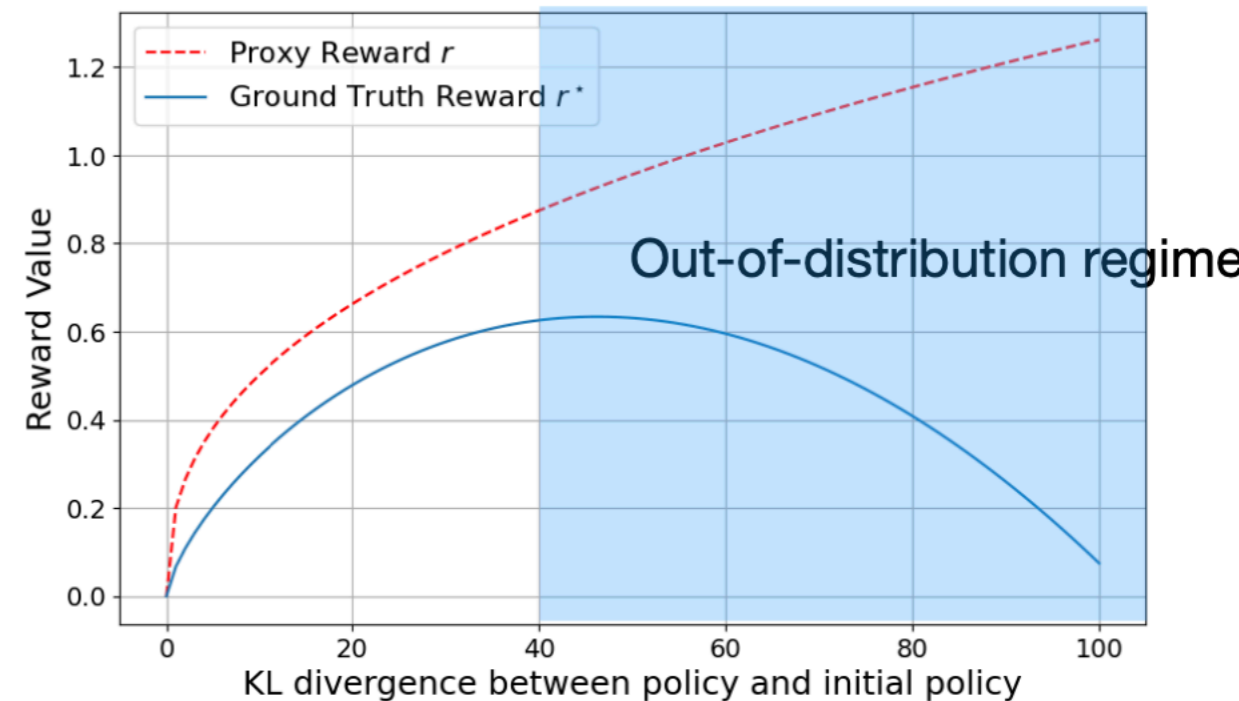


RL/
planning/
control/

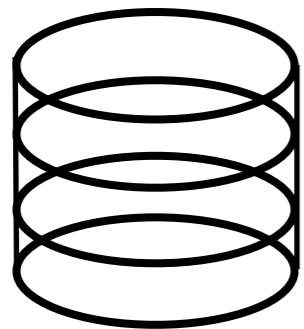
...



PPO-trained policy

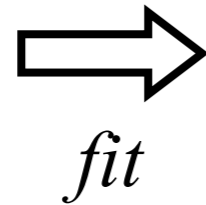


Example: RLHF in LLMs

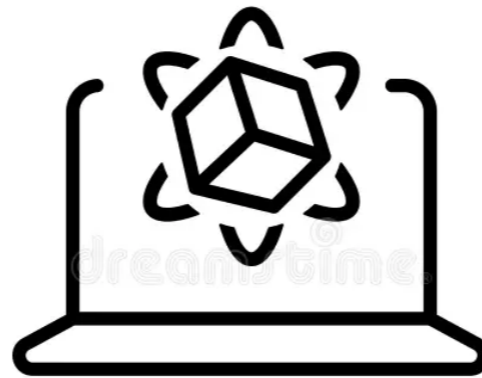


data

{prompt,
response,
reward}



fit



simulator

reward model



*RL/
planning/
control/*

...



PPO-trained policy

- Generalization only happens to policies “*covered by*” behavior policy
- How to define *coverage*, and what’s its interplay with algorithms?

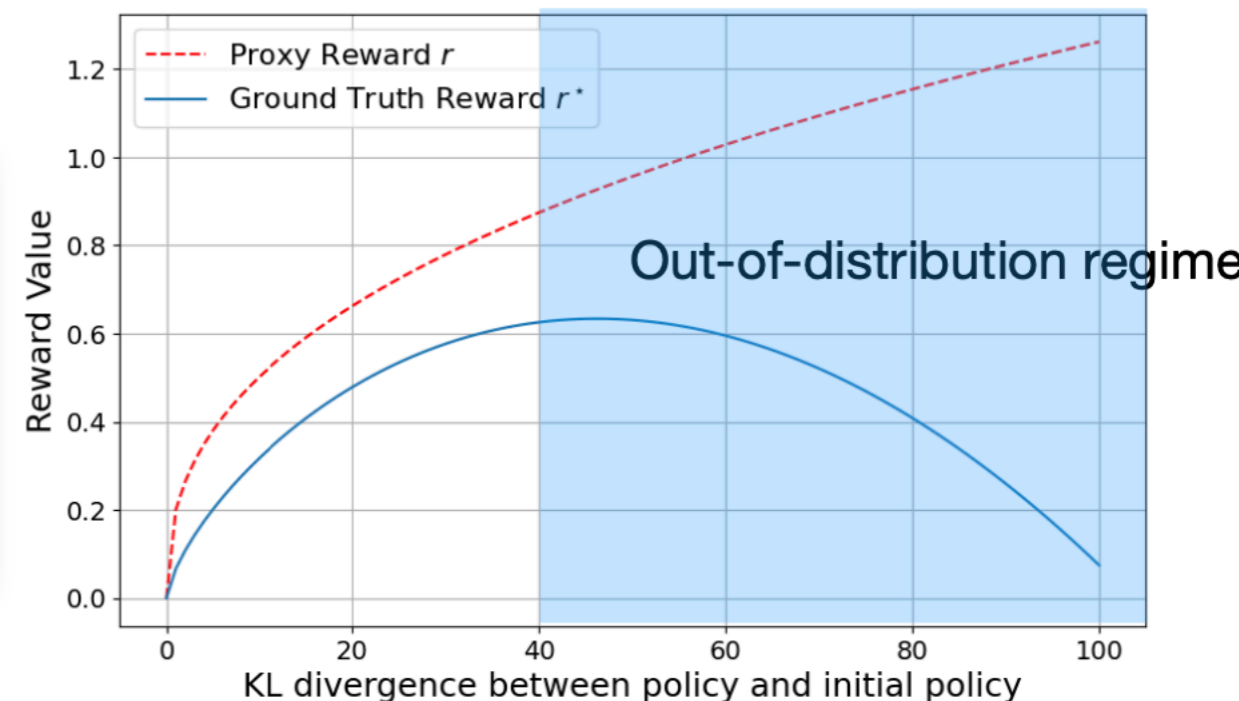


Figure from Gao et al. 2023

Framework for decision-making

- Episodic RL:

$$O_1, a_1, r_1, \dots, O_H, a_H, r_H$$

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF
 - o_1 is prompt, each a_h is a token

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF

- o_1 is prompt, each a_h is a token
- $r_h = 0$ except for $h=H$
- o_h is constant symbol (nothing happens between tokens)

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF

- o_1 is prompt, each a_h is a token
- $r_h = 0$ except for $h=H$
- o_h is constant symbol (nothing happens between tokens)
- alternatively, bandit ($H=1$) with combinatorial action space

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF

- o_1 is prompt, each a_h is a token
 - $r_h = 0$ except for $h=H$
 - o_h is constant symbol (nothing happens between tokens)
 - alternatively, bandit ($H=1$) with combinatorial action space
- A policy π maps history

$$\tau_h = (o_1, a_1, r_1, \dots, o_h)$$

to action distribution: $a_h \sim \pi(\cdot \mid \tau_h)$

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF

- o_1 is prompt, each a_h is a token
- $r_h = 0$ except for $h=H$
- o_h is constant symbol (nothing happens between tokens)
- alternatively, bandit ($H=1$) with combinatorial action space

- A policy π maps history

$$\tau_h = (o_1, a_1, r_1, \dots, o_h)$$

to action distribution: $a_h \sim \pi(\cdot \mid \tau_h)$

- Performance measured by $J(\pi) := \mathbb{E}_\pi \left[\sum_{h=1}^H r_h \right]$

Framework for decision-making

- Episodic RL:

$$o_1, a_1, r_1, \dots, o_H, a_H, r_H$$

- RLHF

- o_1 is prompt, each a_h is a token
- $r_h = 0$ except for $h=H$
- o_h is constant symbol (nothing happens between tokens)
- alternatively, bandit ($H=1$) with combinatorial action space

- A policy π maps history

$$\tau_h = (o_1, a_1, r_1, \dots, o_h)$$

to action distribution: $a_h \sim \pi(\cdot \mid \tau_h)$

- Performance measured by $J(\pi) := \mathbb{E}_\pi \left[\sum_{h=1}^H r_h \right]$
- OPE: estimate $J(\pi)$ using data episodes collected with π_b

Unbiased OPE

Importance sampling (IS) [Precup'00]

Behavior



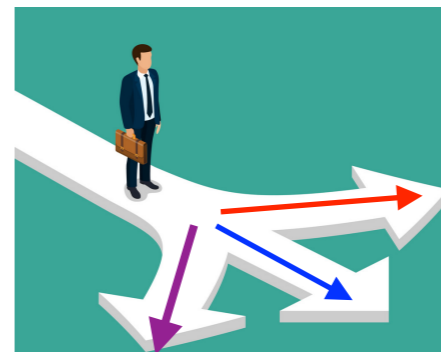
Target



Unbiased OPE

Importance sampling (IS) [Precup'00]

Behavior



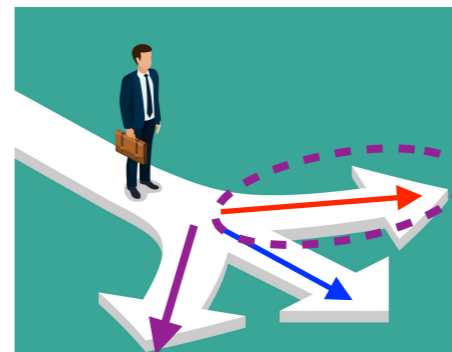
Target



Unbiased OPE

Importance sampling (IS) [Precup'00]

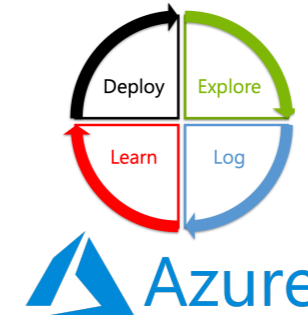
Behavior



Target



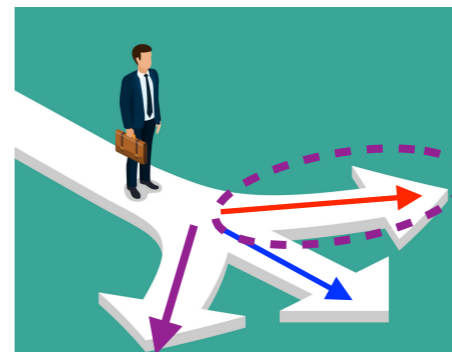
Unbiased OPE



Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓

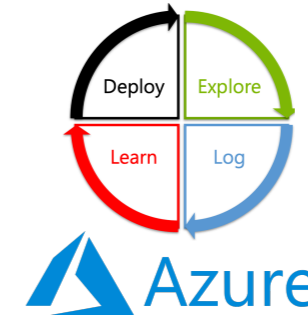
Behavior



Target

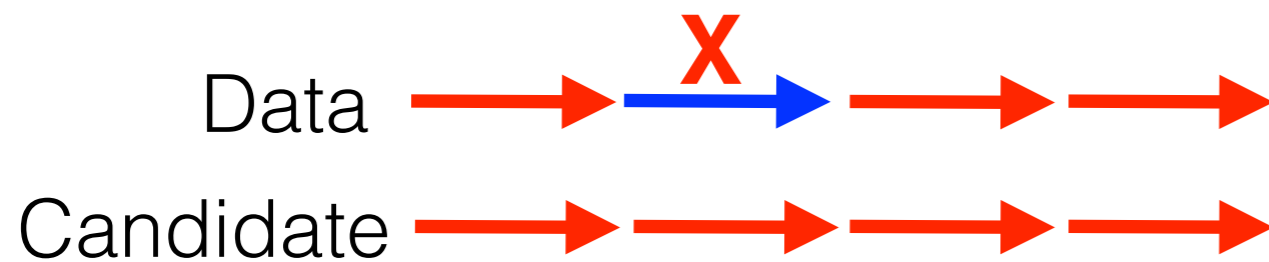


Unbiased OPE

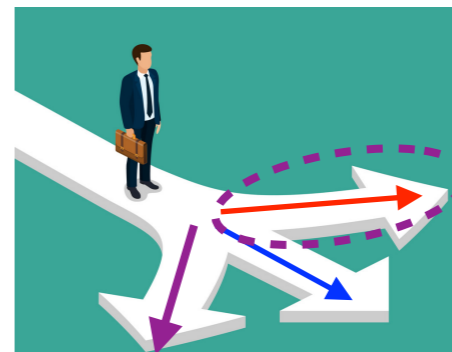


Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



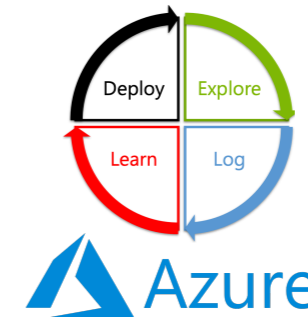
Behavior



Target

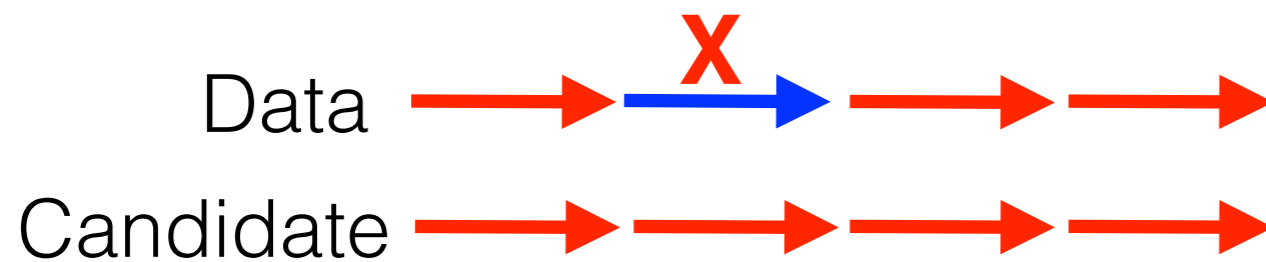


Unbiased OPE

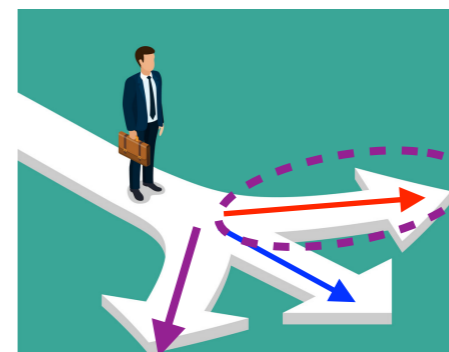


Importance sampling (IS) [Precup'00]

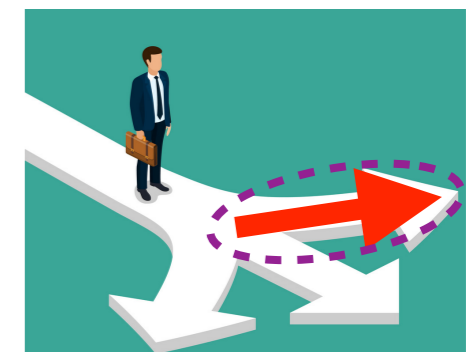
- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



Behavior

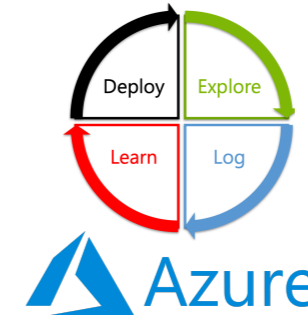


Target



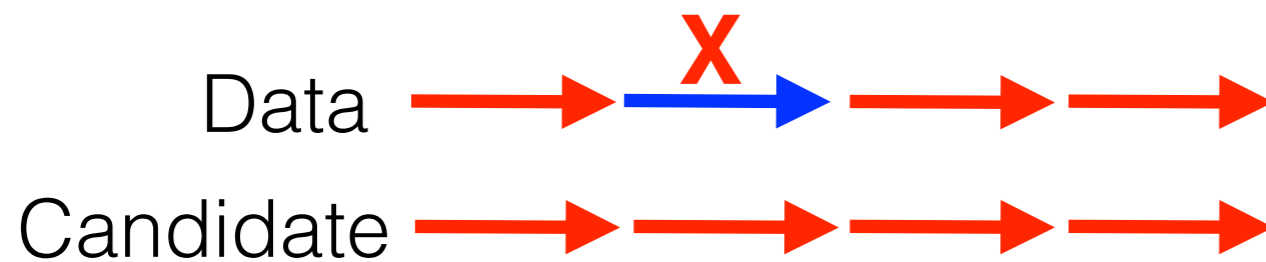
$$o_1, a_1, r_1, \dots, o_H, a_H, r_H \Rightarrow \left(\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)} \right) \left(\sum_{h=1}^H r_h \right)$$

Unbiased OPE

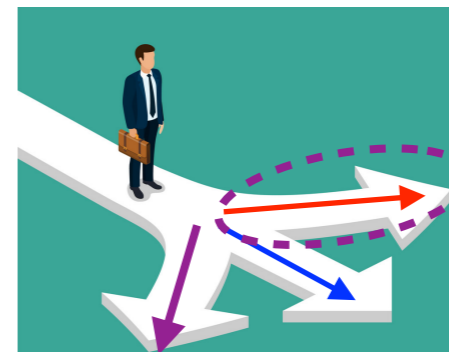


Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



Behavior

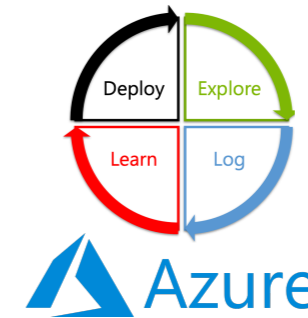


Target



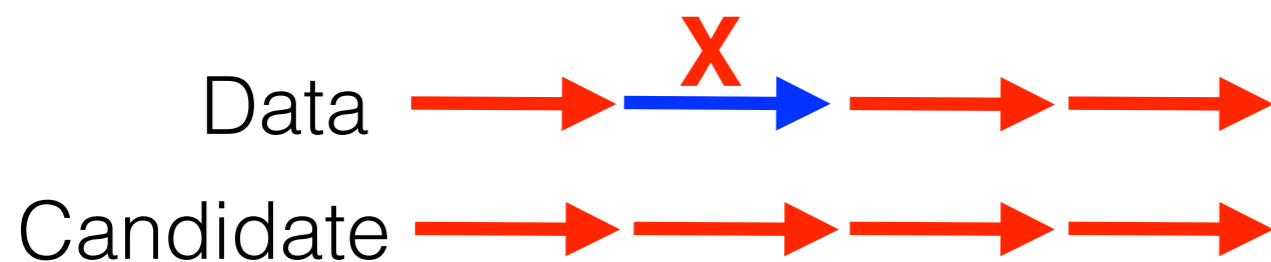
$$o_1, a_1, r_1, \dots, o_H, a_H, r_H \Rightarrow \left(\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)} \right) \left(\sum_{h=1}^H r_h \right)$$

Unbiased OPE

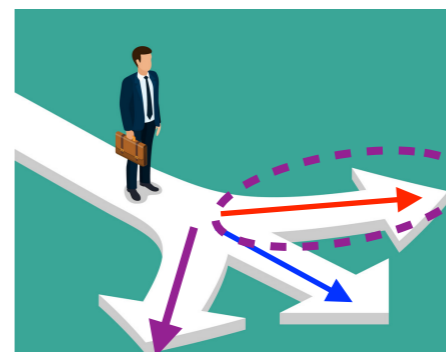


Importance sampling (IS) [Precup'00]

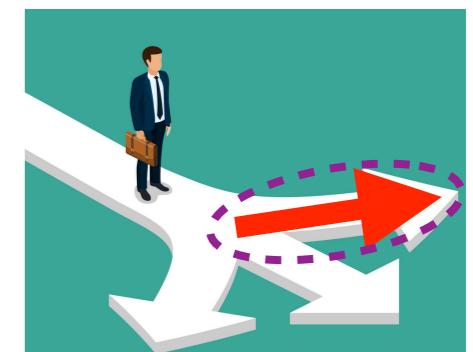
- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



Behavior



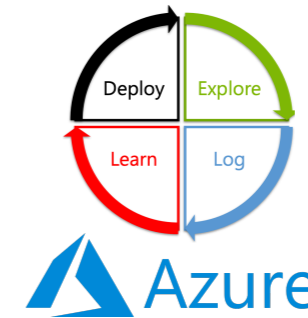
Target



$$o_1, a_1, r_1, \dots, o_H, a_H, r_H \Rightarrow \left(\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)} \right) \left(\sum_{h=1}^H r_h \right)$$

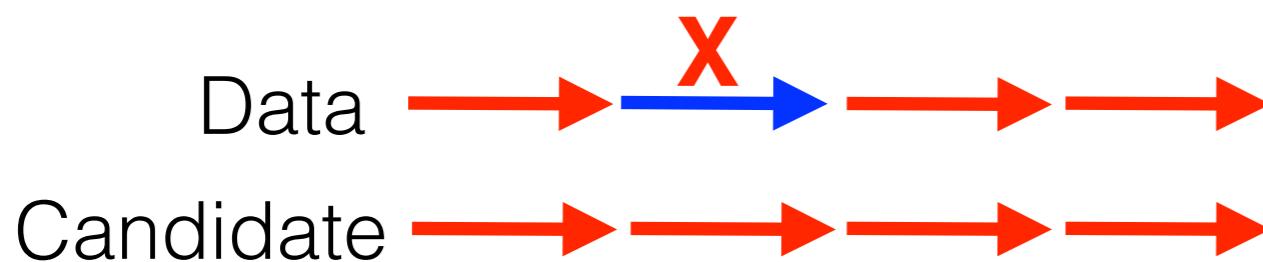
- Or, can only evaluate π when $\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)}$ small

Unbiased OPE

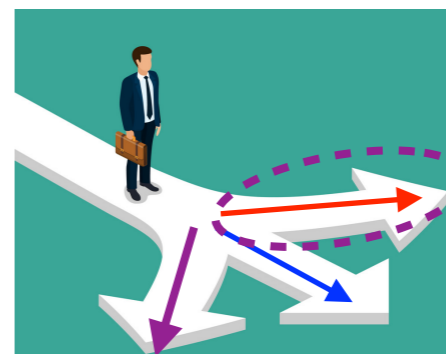


Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



Behavior



Target



$$o_1, a_1, r_1, \dots, o_H, a_H, r_H \Rightarrow \left(\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)} \right) \left(\sum_{h=1}^H r_h \right)$$

- Or, can only evaluate π when $\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)}$ small

IS' measure of *coverage*

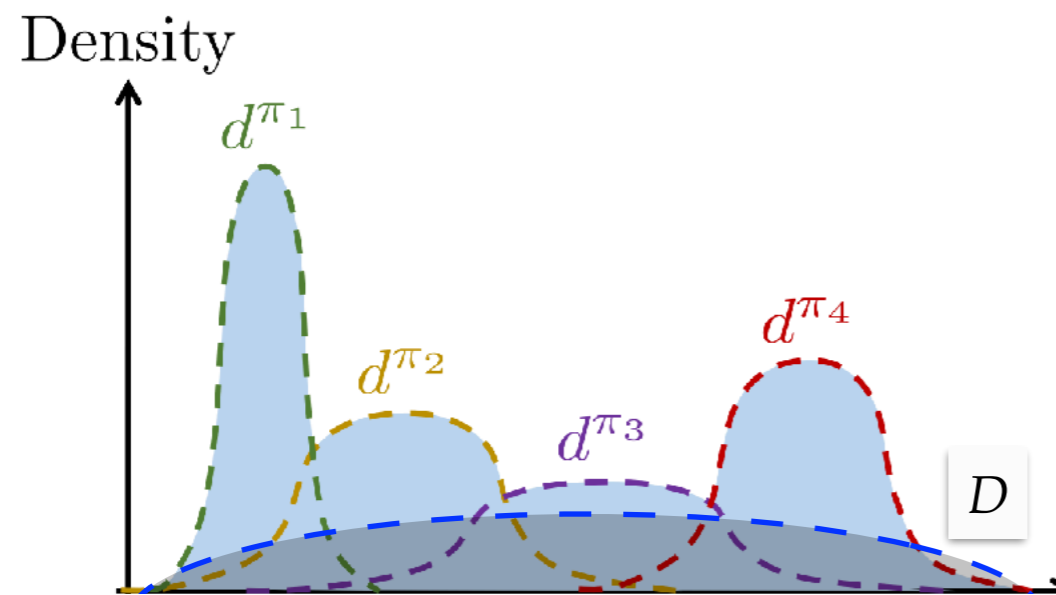
Why OPE & Coverage?

Evaluation is the basis of optimization

Why OPE & Coverage?

Evaluation is the basis of optimization

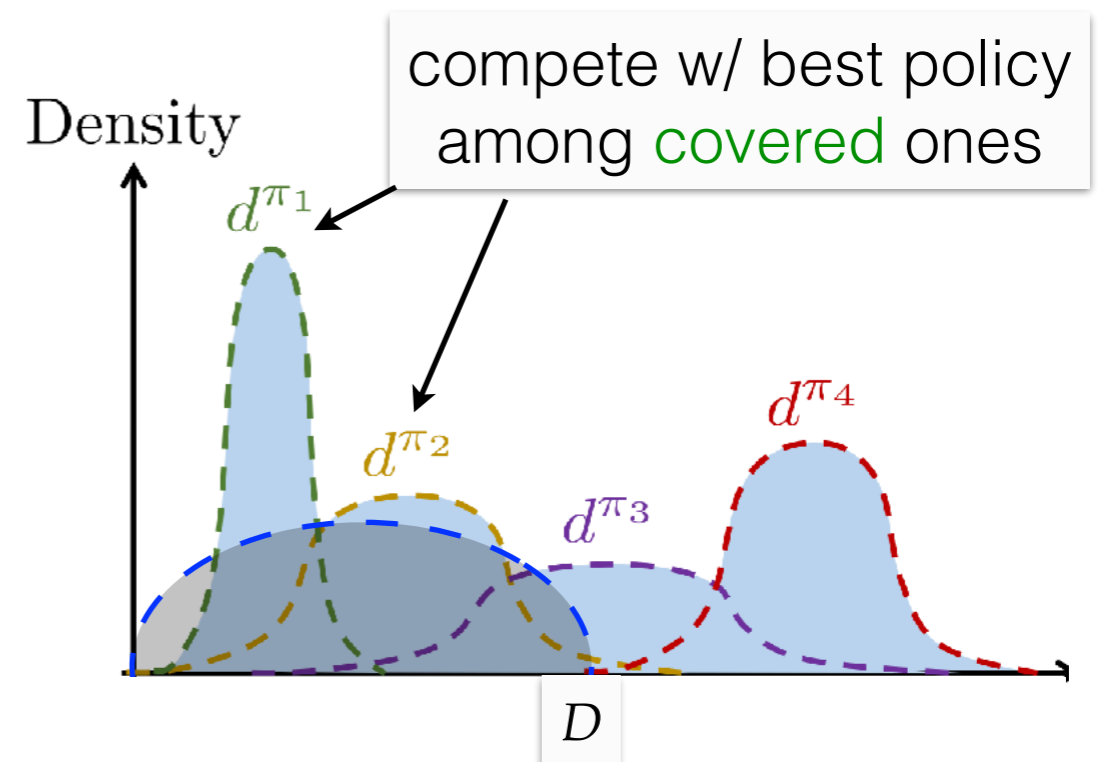
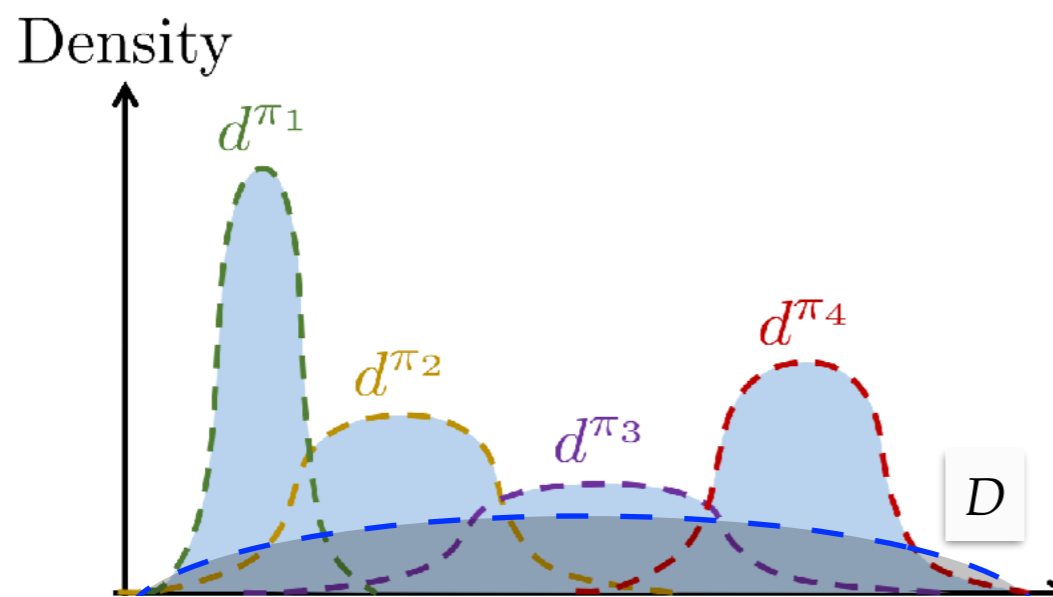
- If all policies are covered & accurately evaluated, we can pick the approximate best



Why OPE & Coverage?

Evaluation is the basis of optimization

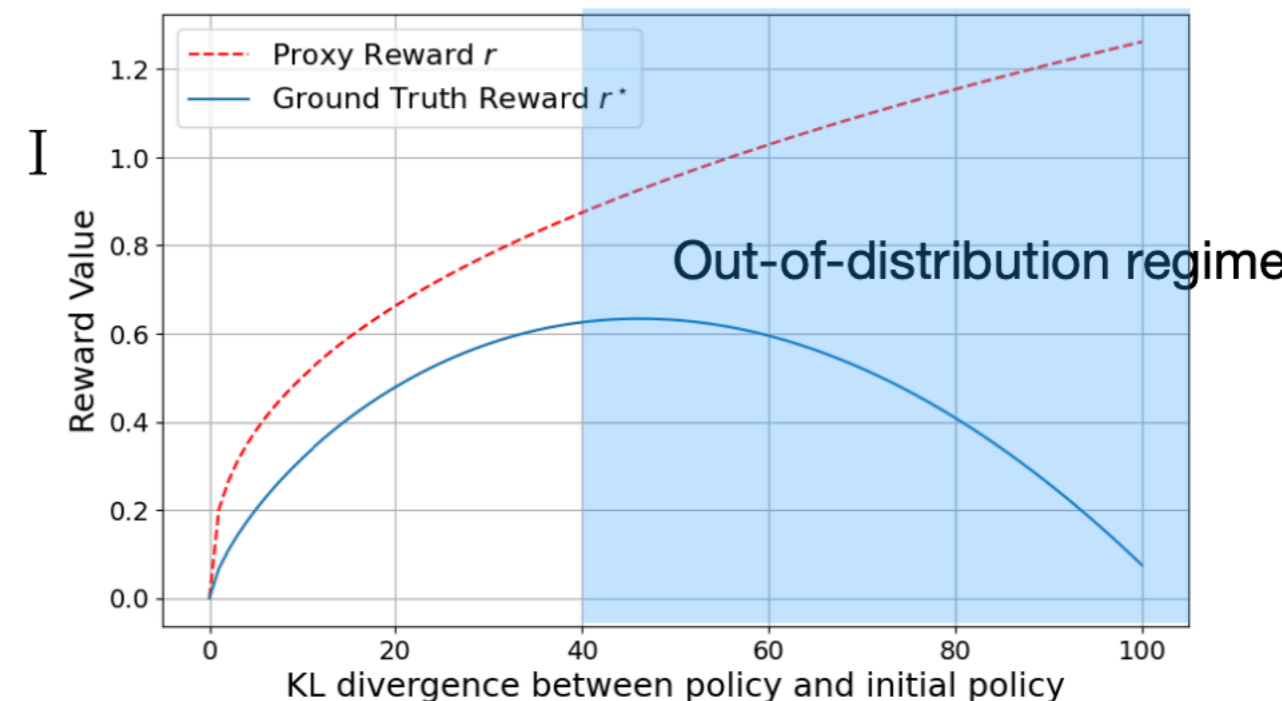
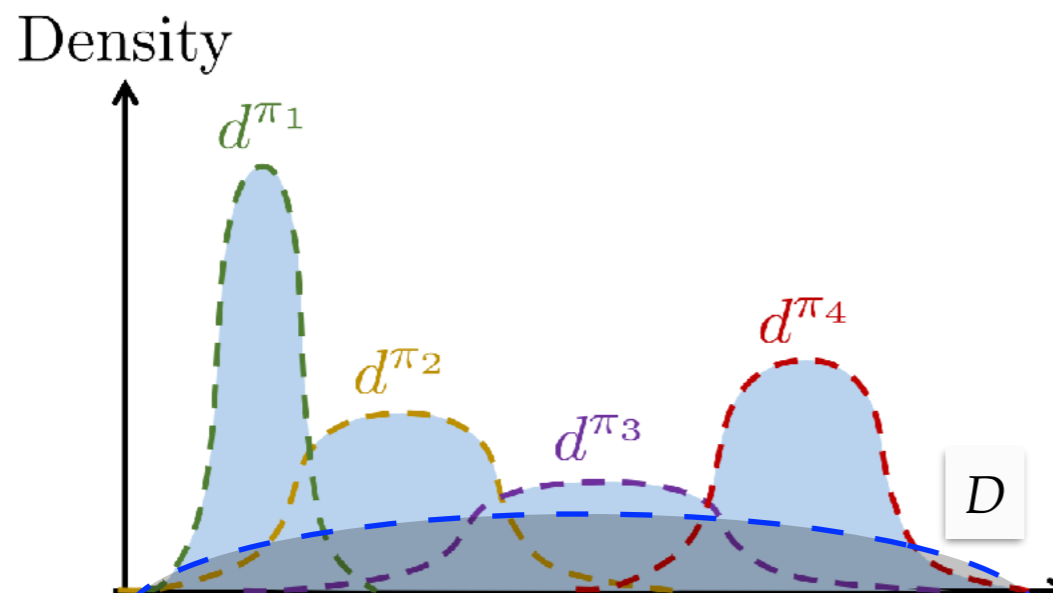
- If all policies are covered & accurately evaluated, we can pick the approximate best
- Doesn't have to! Can **constrain** learned policy to be **covered** => compete with best covered policy



Why OPE & Coverage?

Evaluation is the basis of optimization

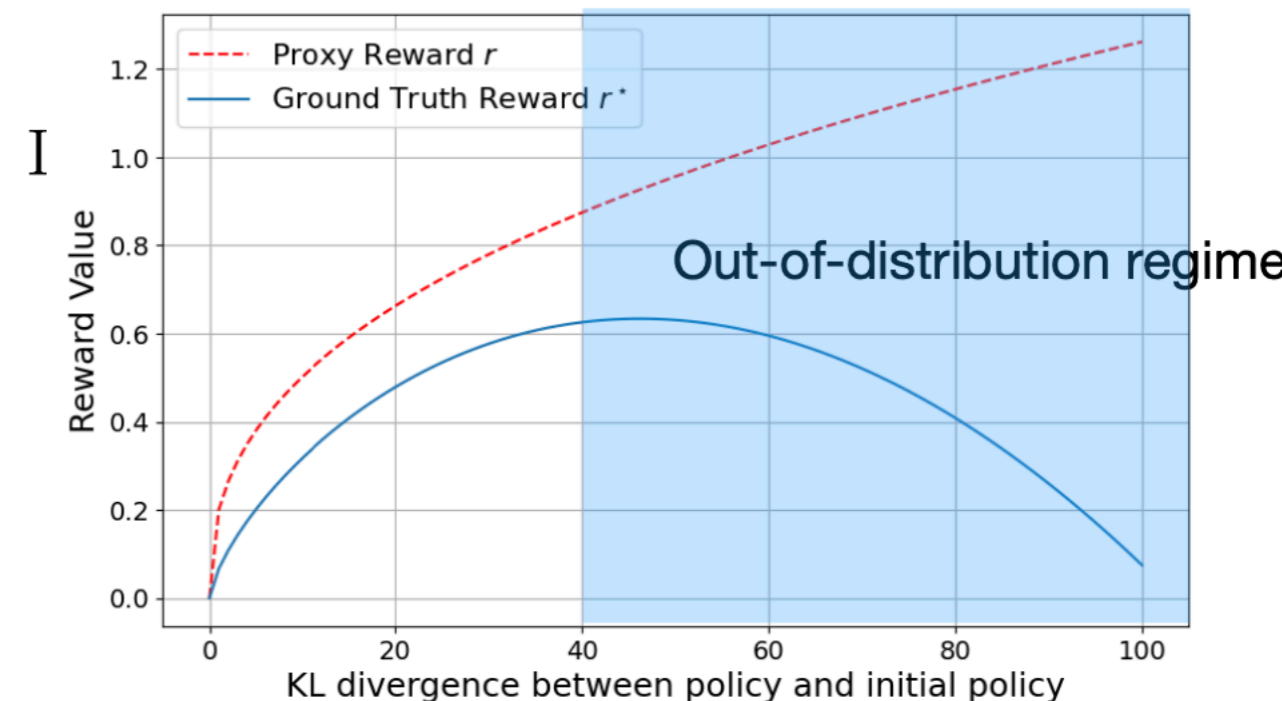
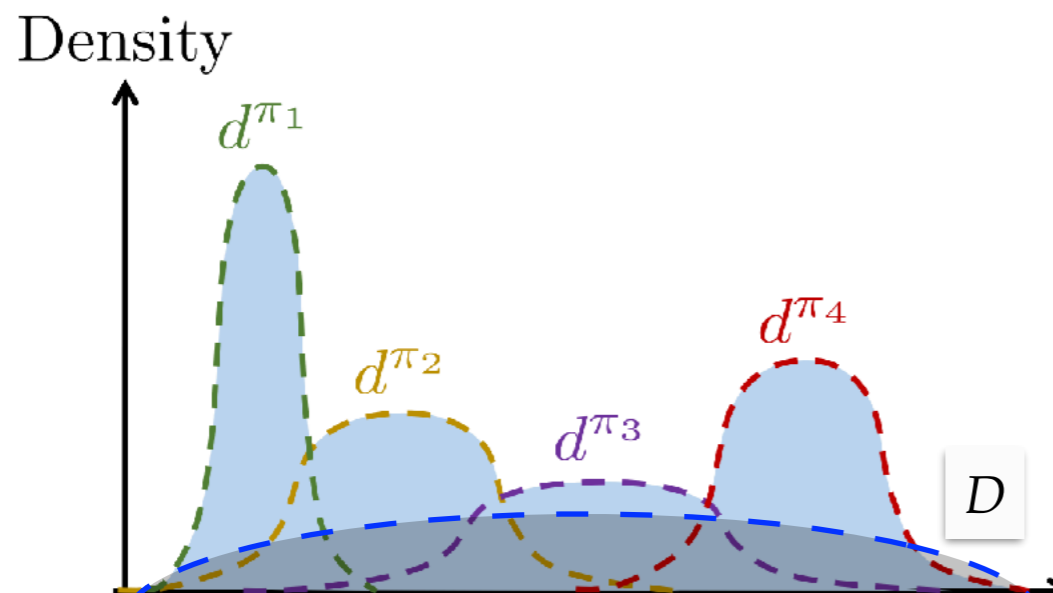
- If all policies are covered & accurately evaluated, we can pick the approximate best
- Doesn't have to! Can **constrain** learned policy to be **covered** => compete with best covered policy
 - behavior regularization, pessimism in face of uncertainty, etc.



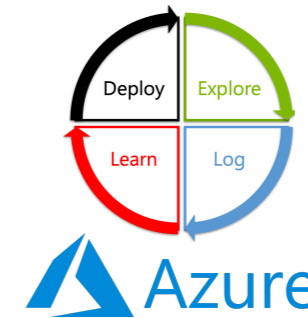
Why OPE & Coverage?

Evaluation is the basis of optimization

- If all policies are covered & accurately evaluated, we can pick the approximate best
- Doesn't have to! Can **constrain** learned policy to be **covered** => compete with best covered policy
 - behavior regularization, pessimism in face of uncertainty, etc.
- More lenient **coverage** => stronger competing target

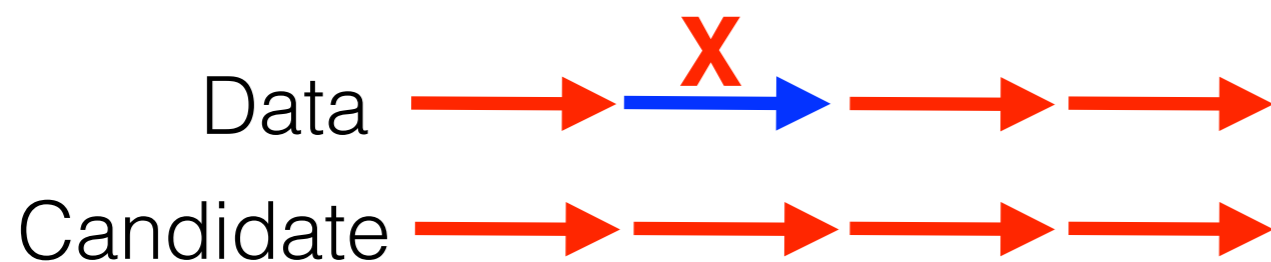


Unbiased OPE

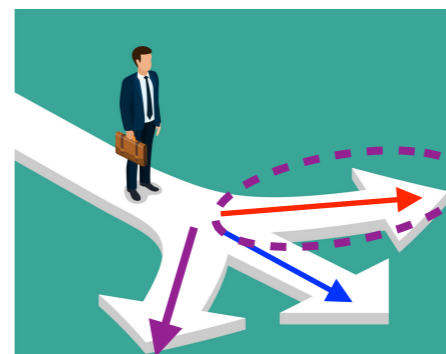


Importance sampling (IS) [Precup'00]

- Industry deployment (ctx. bandit, horizon=1)
- No Markovianity required ✓
- **Exponential-in-horizon** variance!



Behavior



Target




$$o_1, a_1, r_1, \dots, o_H, a_H, r_H \Rightarrow \left(\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)} \right) \left(\sum_{h=1}^H r_h \right)$$

- Or, can only evaluate π when $\prod_{h=1}^H \frac{\pi(a_h | \tau_h)}{\pi_b(a_h | \tau_h)}$ small

IS' measure of *coverage*

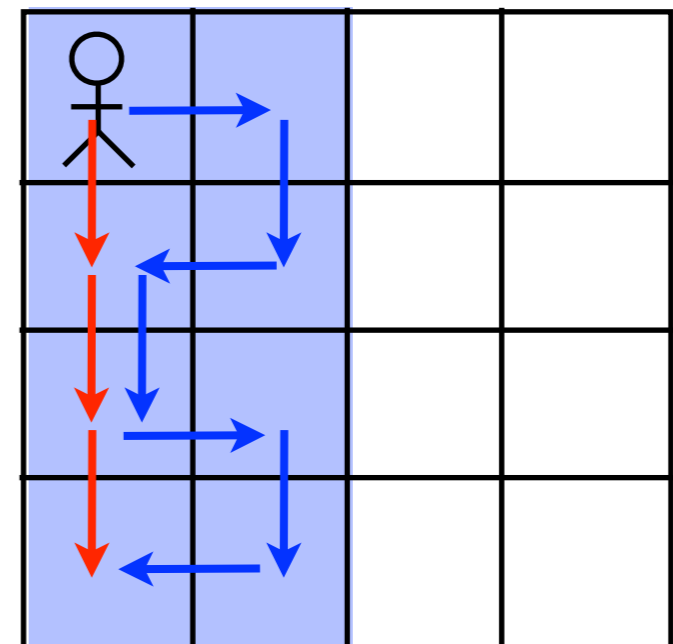
Better coverage?

Better coverage?

FQE [Munos, Szepesvari... CJ'19, ...] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, ...]

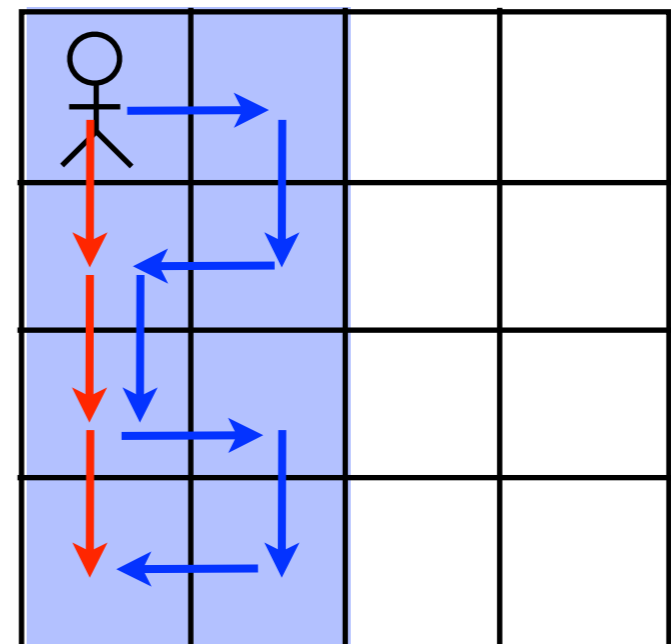
- Assume MDPs: $o_1 (= s_1), a_1, r_1, \dots, o_H (= s_H), a_H, r_H$



Better coverage?

FQE [Munos, Szepesvari... CJ'19, ...] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, ...]

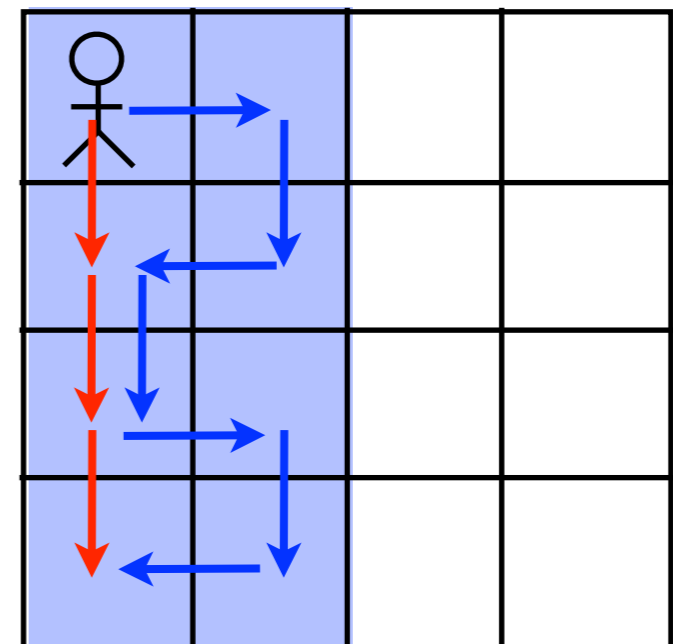
- Assume MDPs: $o_1 (= s_1), a_1, r_1, \dots, o_H (= s_H), a_H, r_H$
- Learn **value functions**: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^H r_{h'} | s_h]$
 - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.



Better coverage?

FQE [Munos, Szepesvari... CJ'19, ...] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, ...]

- Assume MDPs: $o_1 (= s_1), a_1, r_1, \dots, o_H (= s_H), a_H, r_H$
- Learn **value functions**: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^H r_{h'} | s_h]$
 - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.
 - Estimate from class \mathcal{V} by $\arg \min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b} [(V - \mathcal{T}^\pi V)^2]$
- **Coverage**: *marginal* state distribution
 - require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small



* Also needs Bellman completeness

Better coverage?

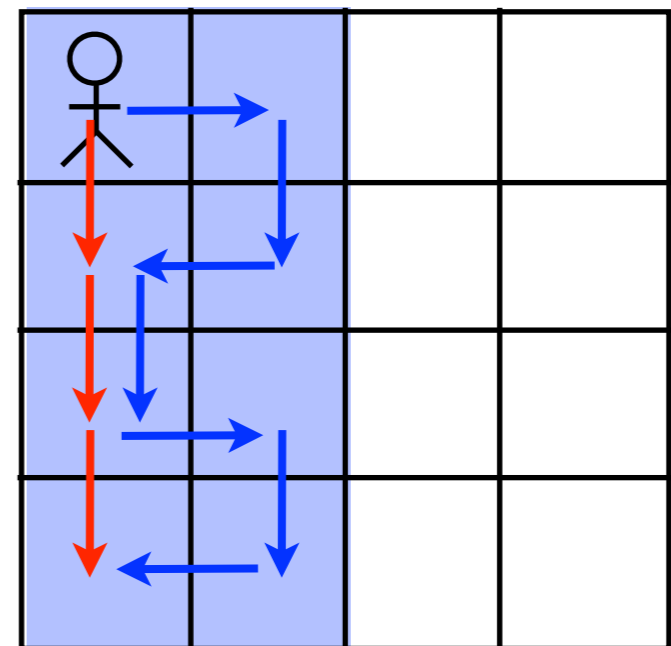
FQE [Munos, Szepesvari... CJ'19, ...] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, ...]

- Assume MDPs: $o_1 (= s_1), a_1, r_1, \dots, o_H (= s_H), a_H, r_H$
- Learn **value functions**: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^H r_{h'} | s_h]$
 - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.
 - Estimate from class \mathcal{V} by $\arg \min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b} [(V - \mathcal{T}^\pi V)^2]$
- **Coverage**: *marginal* state distribution

- require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small

- Can further improve to e.g.,

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi [V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b} [(V - \mathcal{T}^\pi V)^2]}}$$



* Also needs Bellman completeness

Better coverage?

FQE [Munos, Szepesvari... CJ'19, ...] / **MIS** [Liu et al'18, Nachum et al'19, UHJ'20, ...]

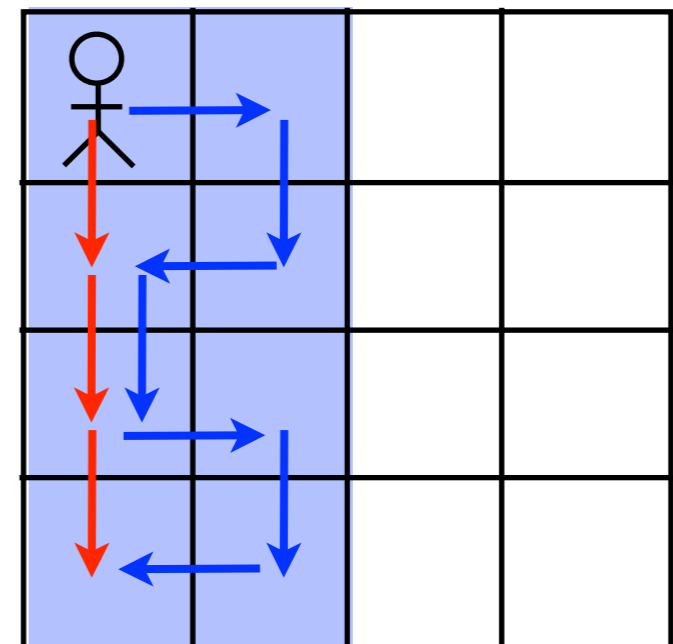
- Assume MDPs: $o_1 (= s_1), a_1, r_1, \dots, o_H (= s_H), a_H, r_H$
- Learn **value functions**: $V^\pi(s_h) := \mathbb{E}_\pi[\sum_{h'=h}^H r_{h'} | s_h]$
 - Satisfies $V^\pi(s_h) = (\mathcal{T}^\pi V^\pi)(s_h) := R(s_h, \pi) + \mathbb{E}_{s_{h+1} \sim P(s_h, \pi)}[V^\pi(s_{h+1})]$.
 - Estimate from class \mathcal{V} by $\arg \min_{V \in \mathcal{V}} \mathbb{E}_{\pi_b} [(V - \mathcal{T}^\pi V)^2]$
- **Coverage**: *marginal* state distribution

- require $\frac{d^\pi(s_h)}{d^{\pi_b}(s_h)}$ and $\frac{\pi(a_h|o_h)}{\pi_b(a_h|o_h)}$ small

- Can further improve to e.g.,

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}}$$

- Fundamental to offline training
& online exploration

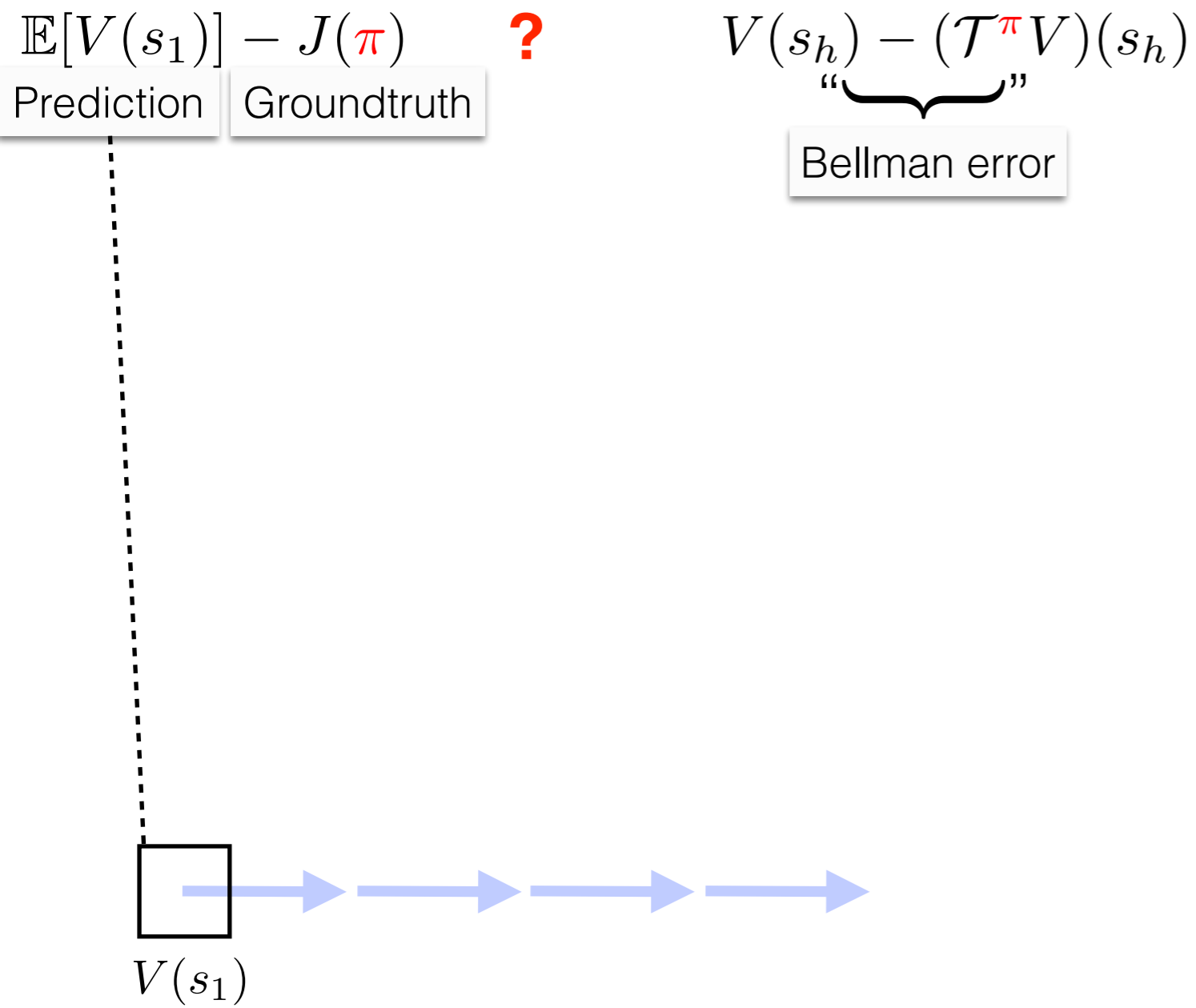


* Also needs Bellman completeness

How do value functions help in MDPs?

$$V(s_h) - \underbrace{(\mathcal{T}^\pi V)}_{\text{Bellman error}}(s_h)$$

How do value functions help in MDPs?

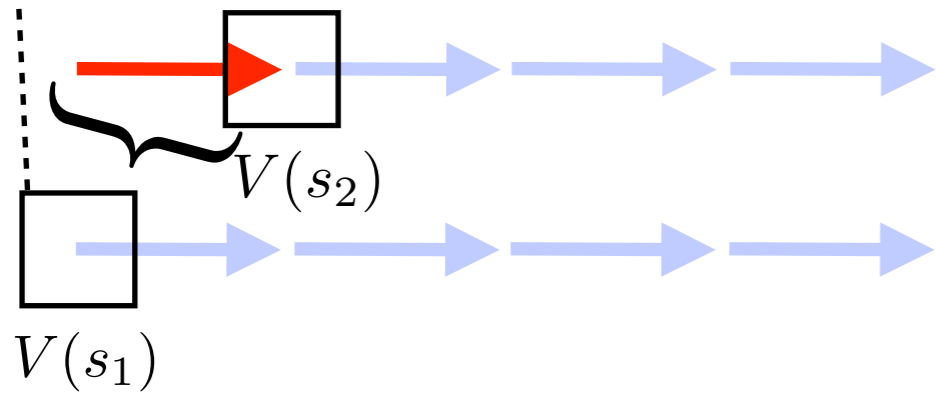


How do value functions help in MDPs?

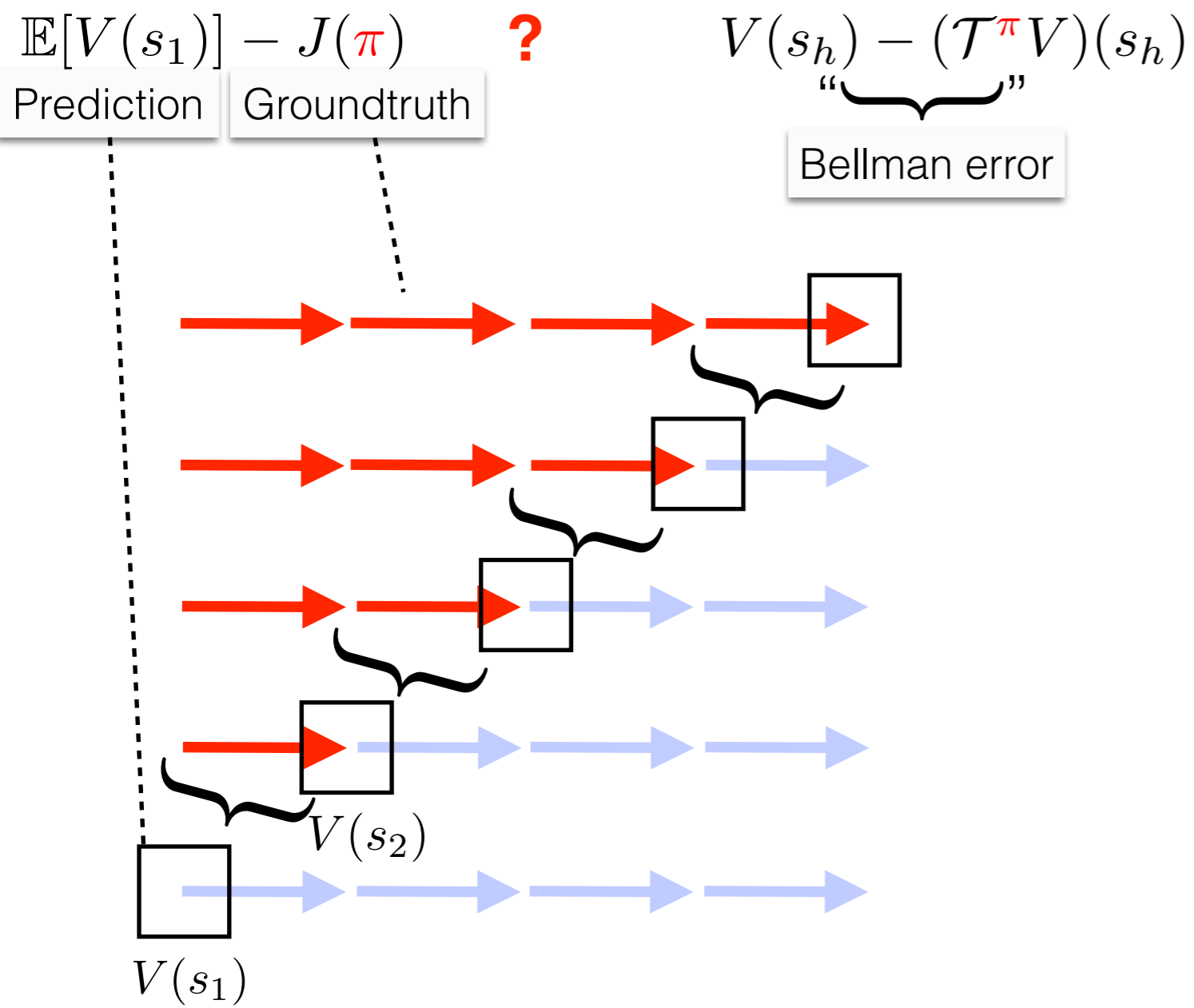
$$\mathbb{E}[V(s_1)] - J(\pi) \quad ?$$

Prediction Groundtruth

$$V(s_h) - \underbrace{(\mathcal{T}^\pi V)(s_h)}_{\text{Bellman error}}$$



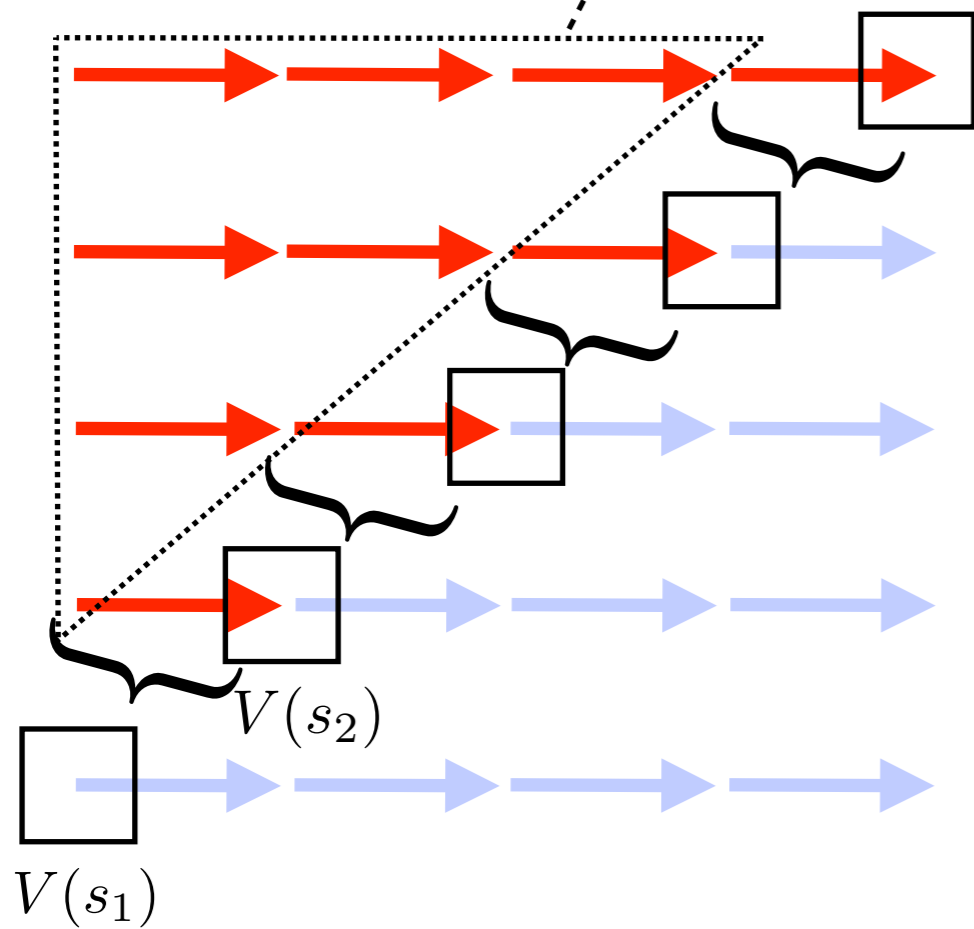
How do value functions help in MDPs?



How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)(s_h)}_{\text{Bellman error}}]$$

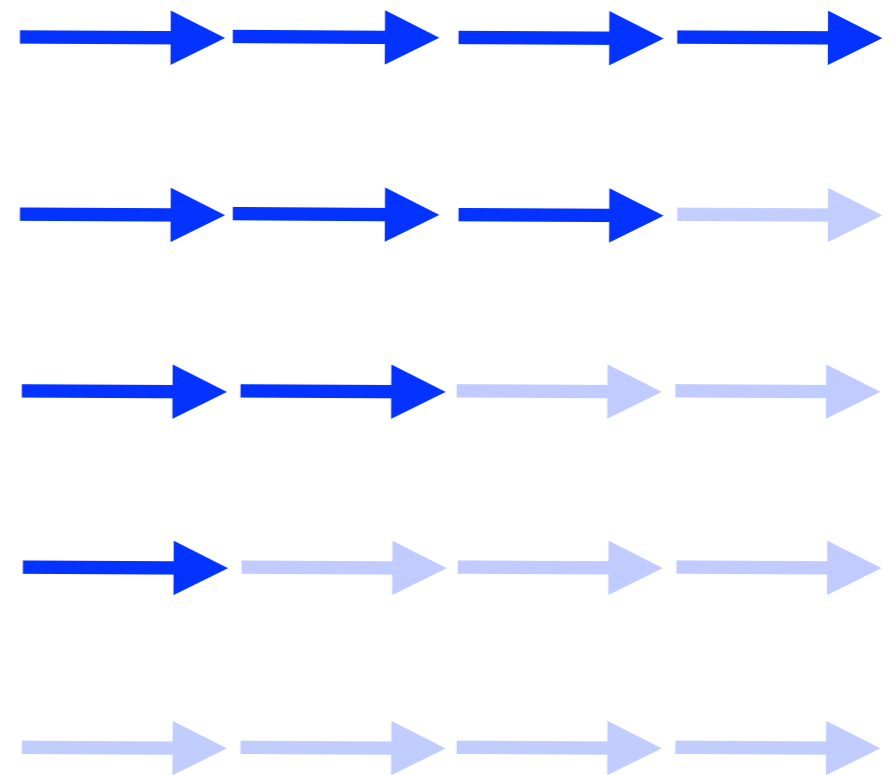
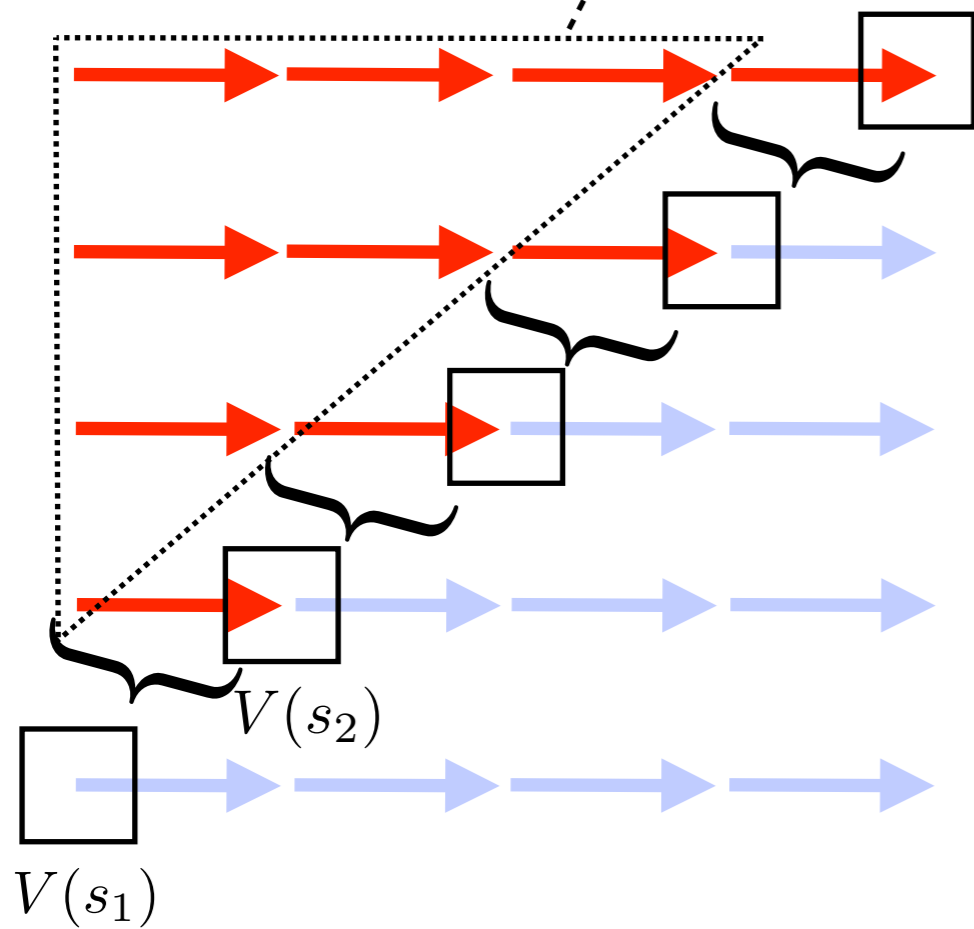
Prediction Groundtruth



How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)}_{\text{Bellman error}}(s_h)]$$

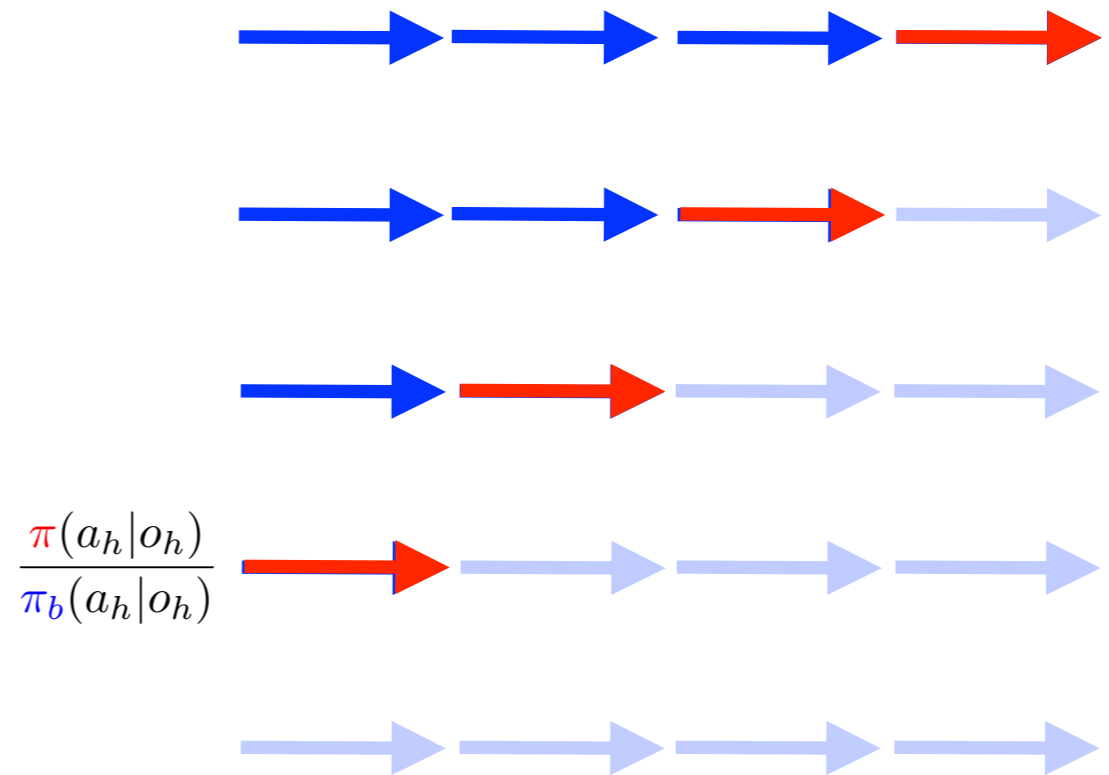
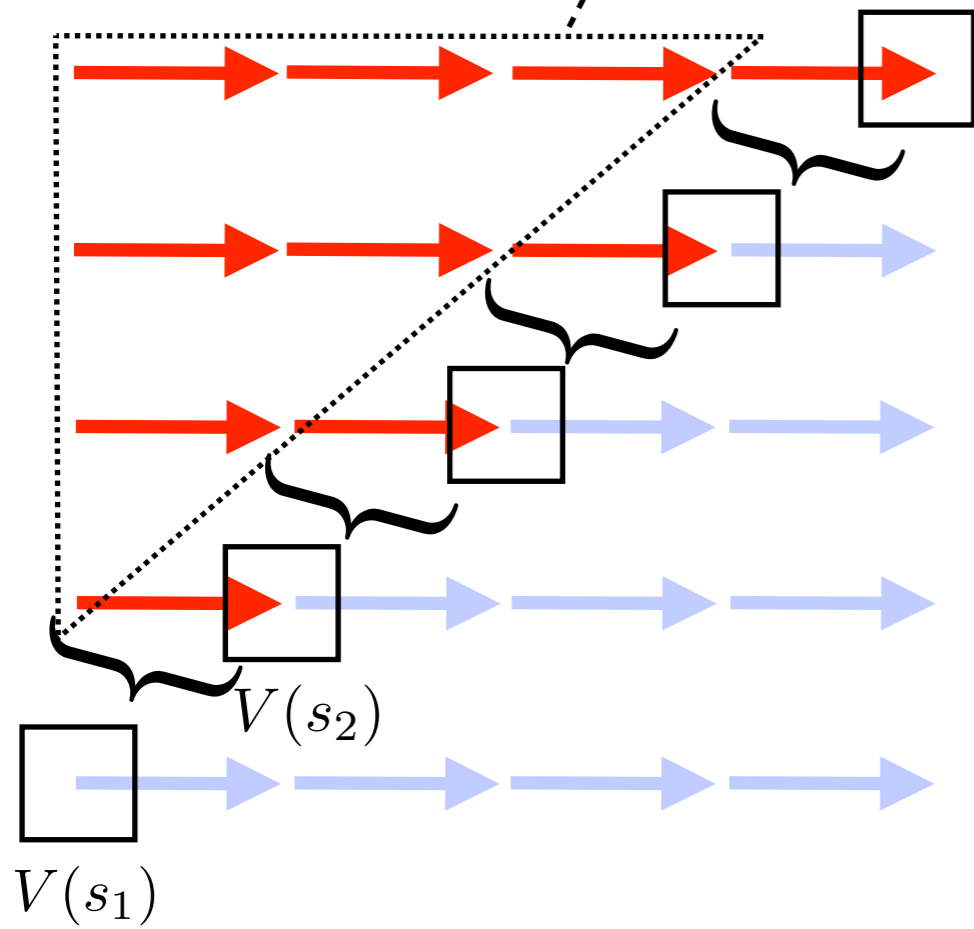
Prediction Groundtruth



How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)}_{\text{Bellman error}}(s_h)]$$

Prediction Groundtruth

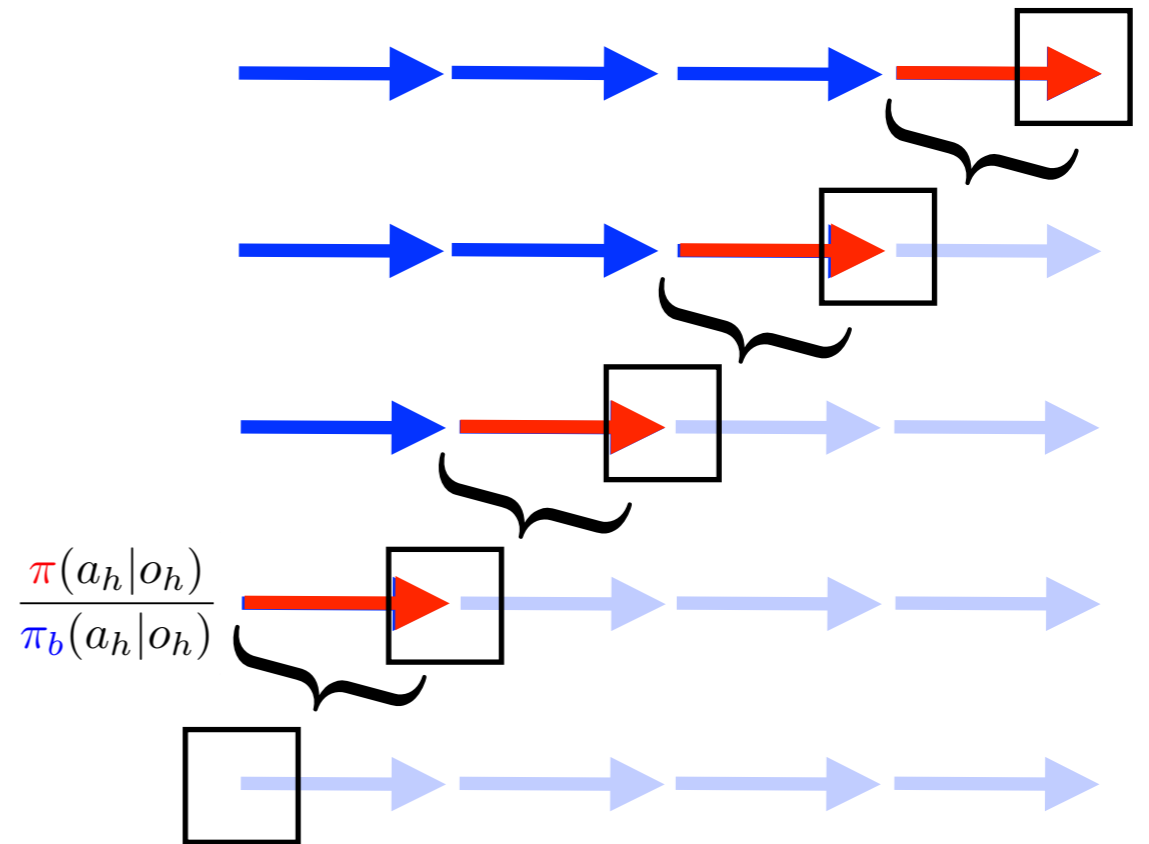
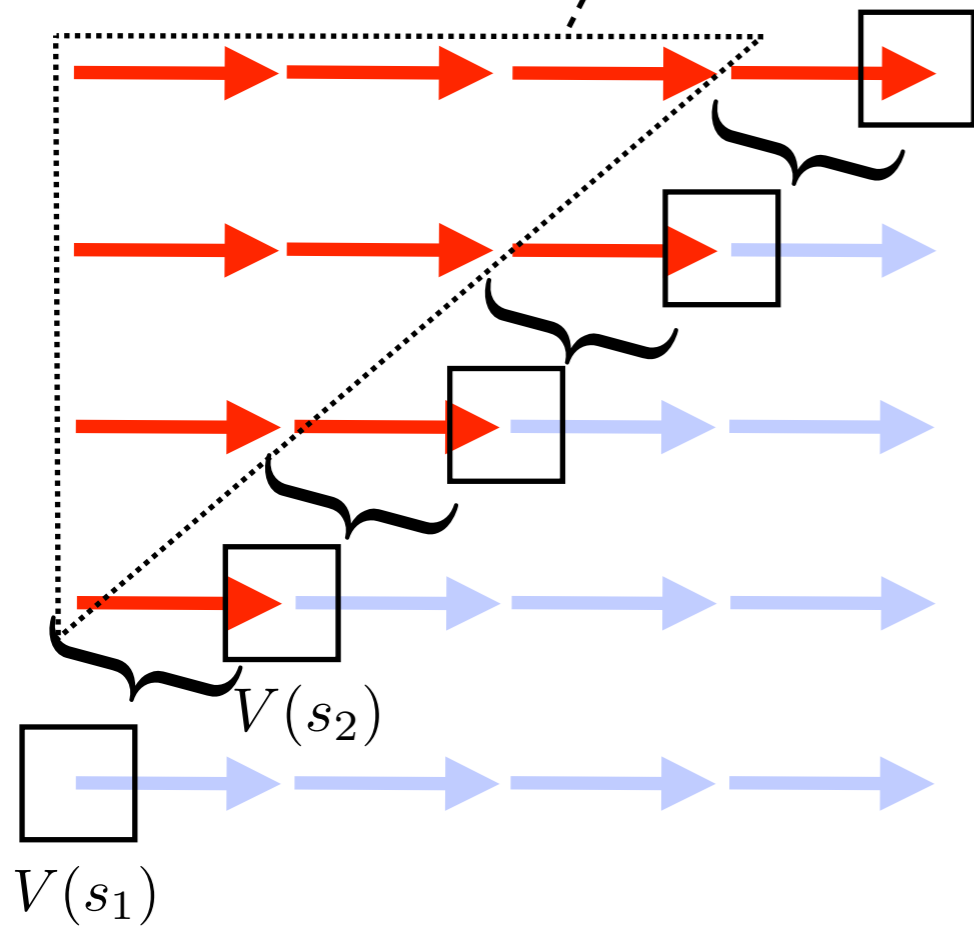


How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)(s_h)}_{\text{Bellman error}}]$$

Prediction Groundtruth

$$\arg \min_V \sum_{h=1}^H \mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$

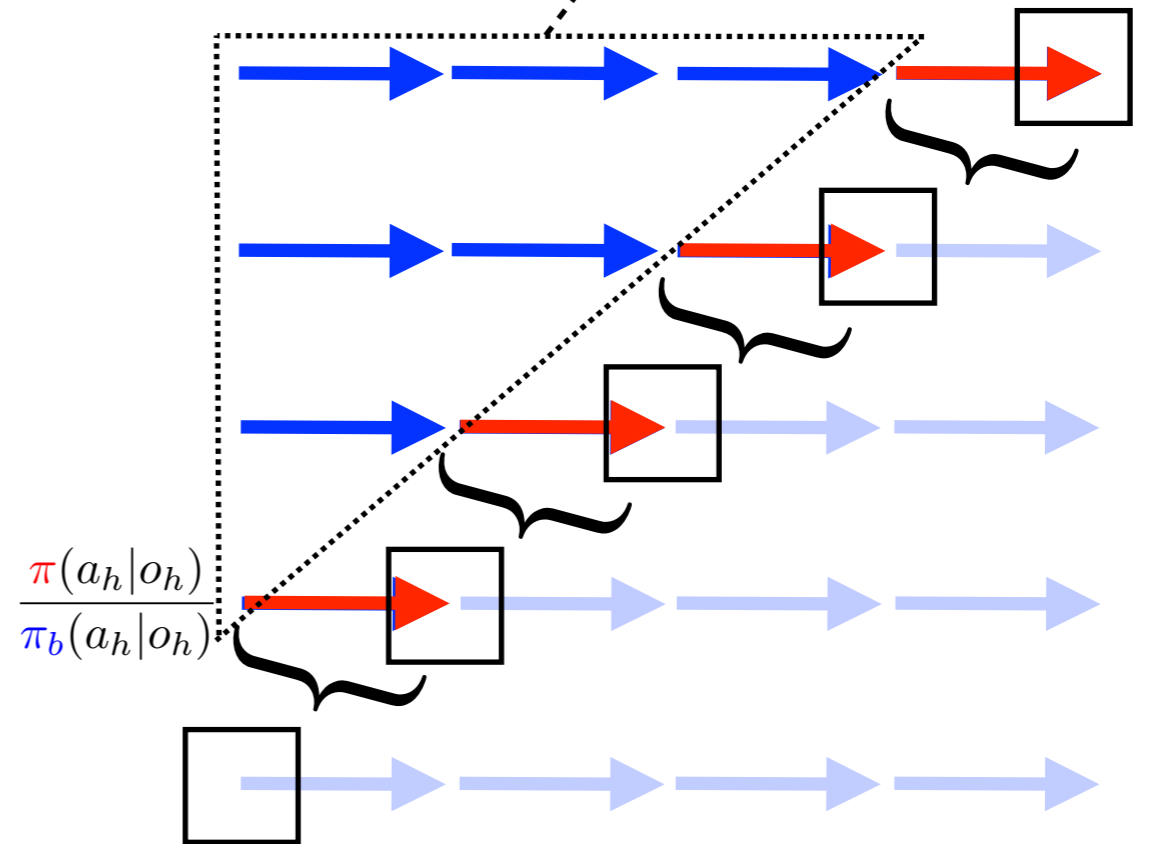
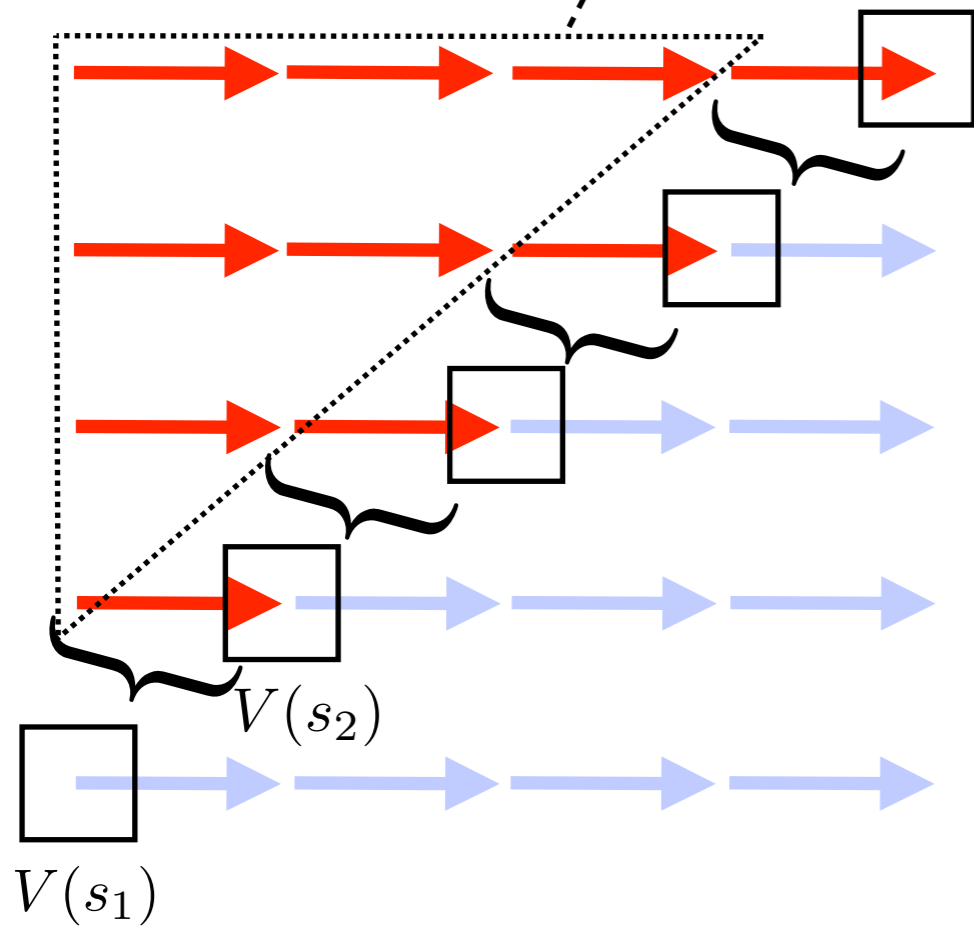


How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)}_{\text{Bellman error}}(s_h)]$$

Prediction Groundtruth

$$\arg \min_V \sum_{h=1}^H \mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$



How do value functions help in MDPs?

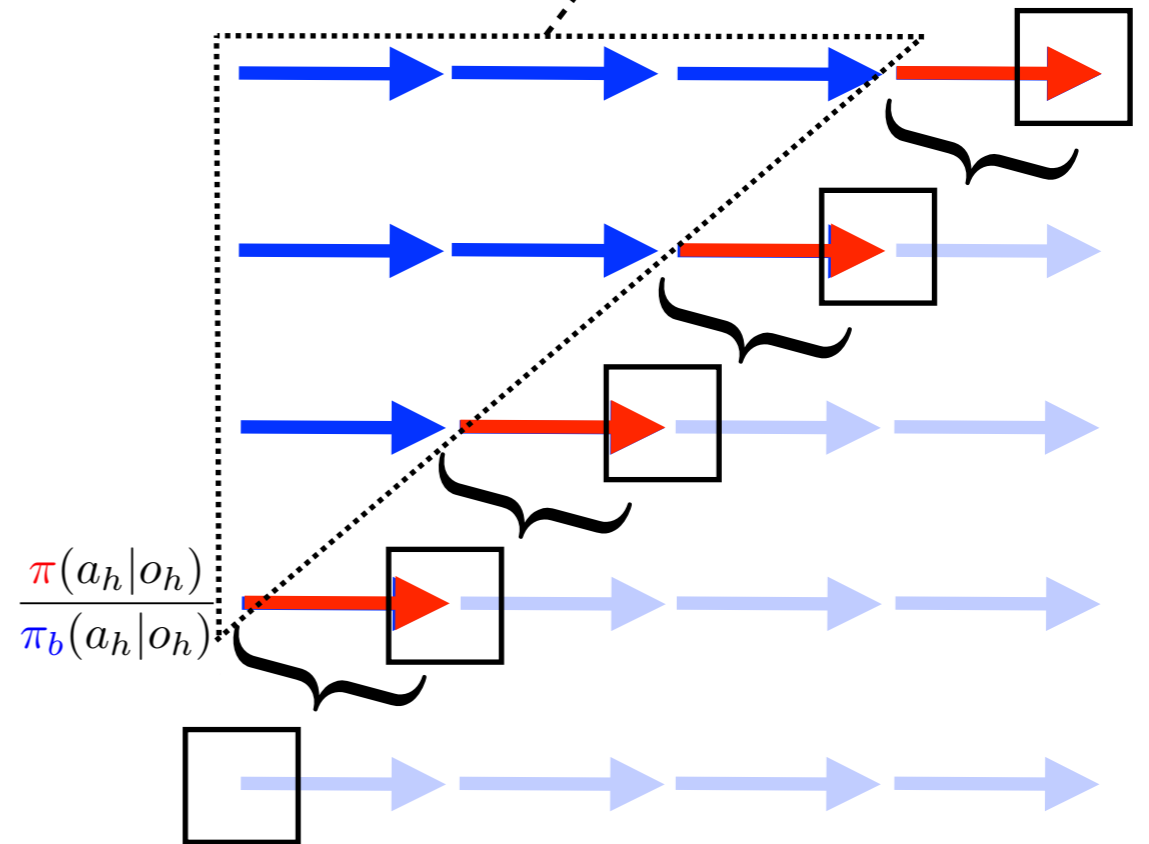
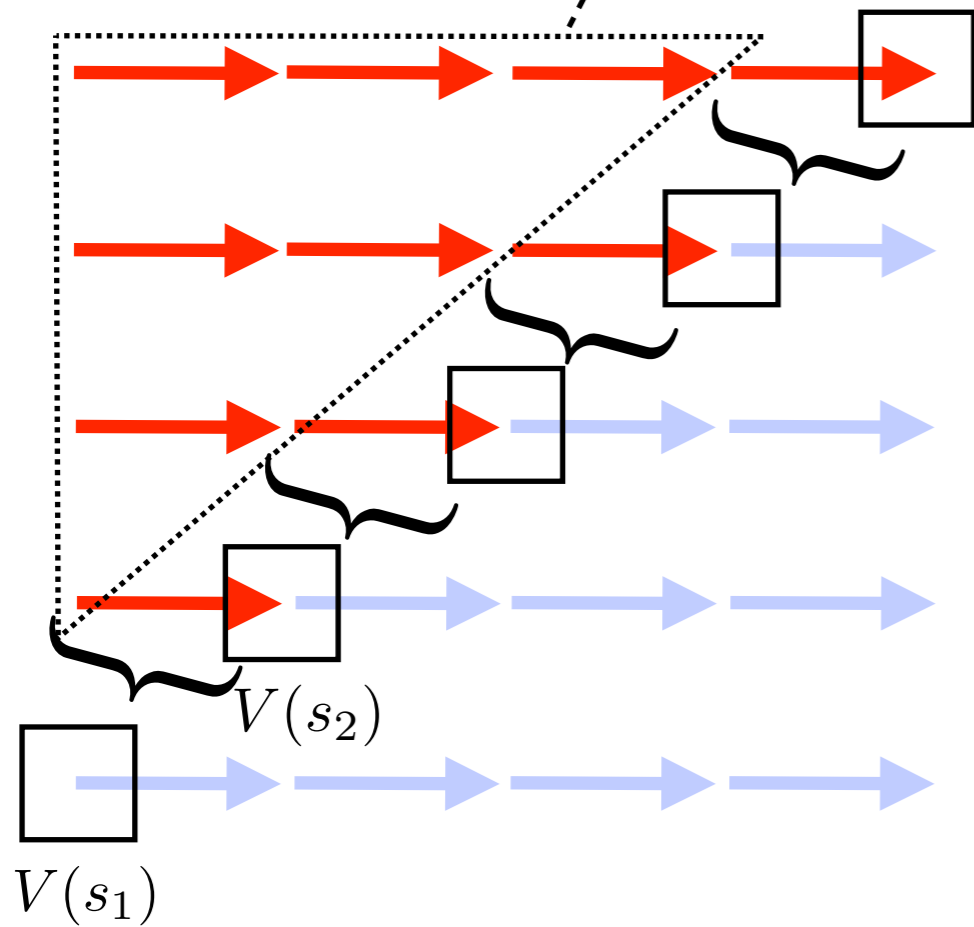
$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi} V)(s_h)]$$

Prediction Groundtruth

Bellman error

$$\mathbb{E}_p[g^2] \leq \|p/q\|_{\infty} \cdot \mathbb{E}_q[g^2]$$

$$\arg \min_V \sum_{h=1}^H \mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$



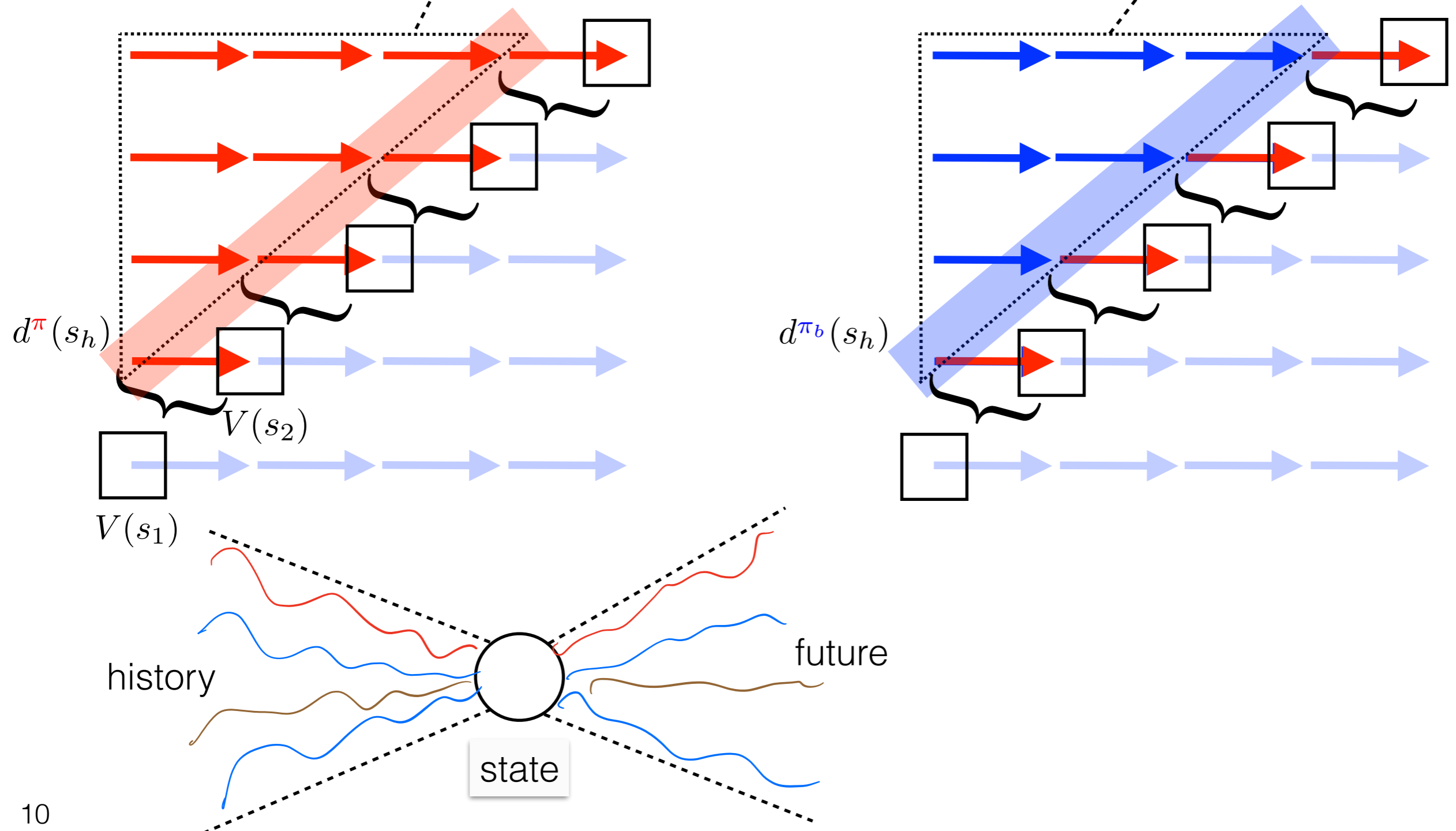
How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)}_{\text{Bellman error}}(s_h)]$$

Prediction Groundtruth

$$\mathbb{E}_p[g^2] \leq \|p/q\|_{\infty} \cdot \mathbb{E}_q[g^2]$$

$$\arg \min_V \sum_{h=1}^H \mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$



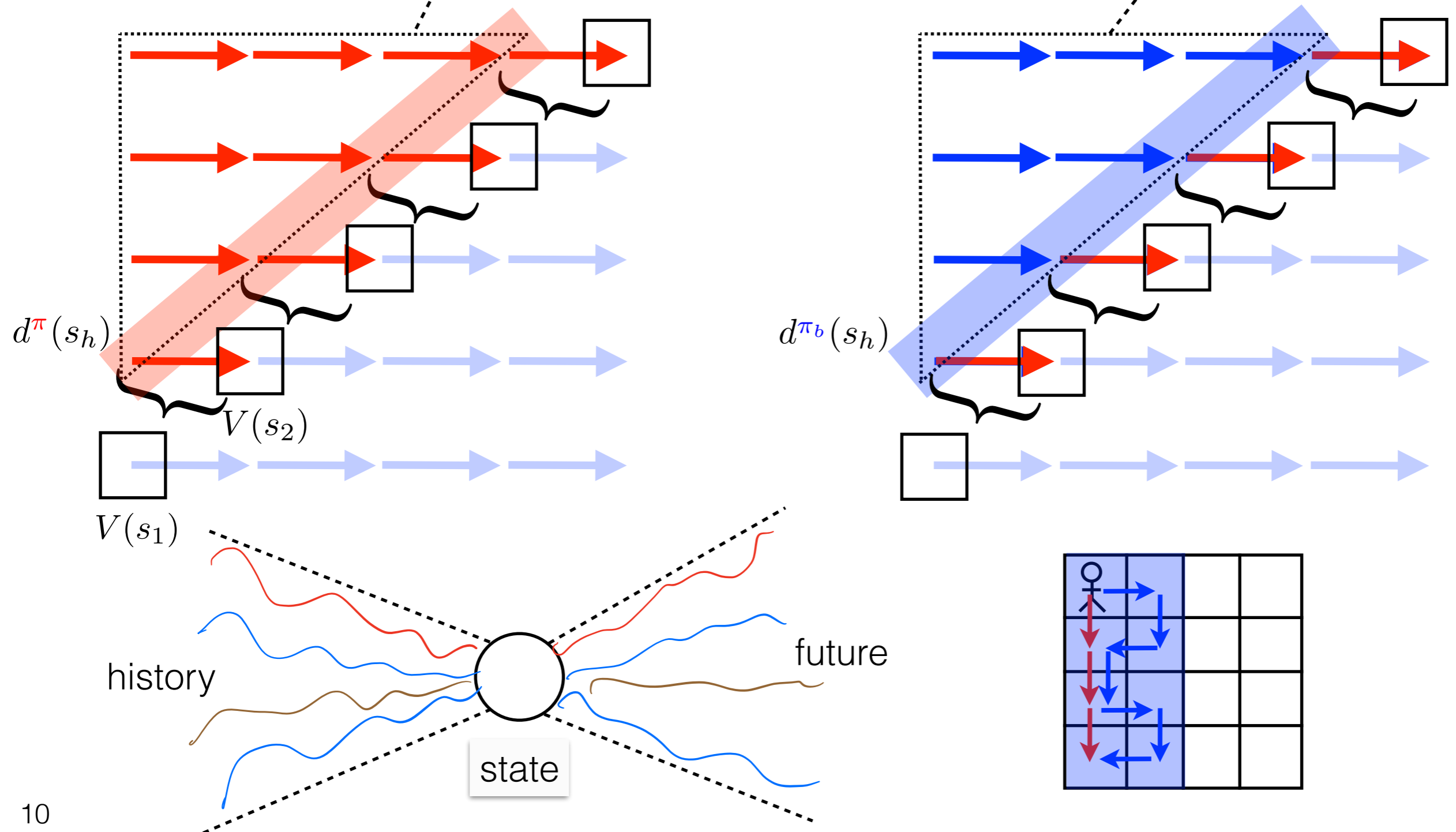
How do value functions help in MDPs?

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - \underbrace{(\mathcal{T}^{\pi} V)}_{\text{Bellman error}}(s_h)]$$

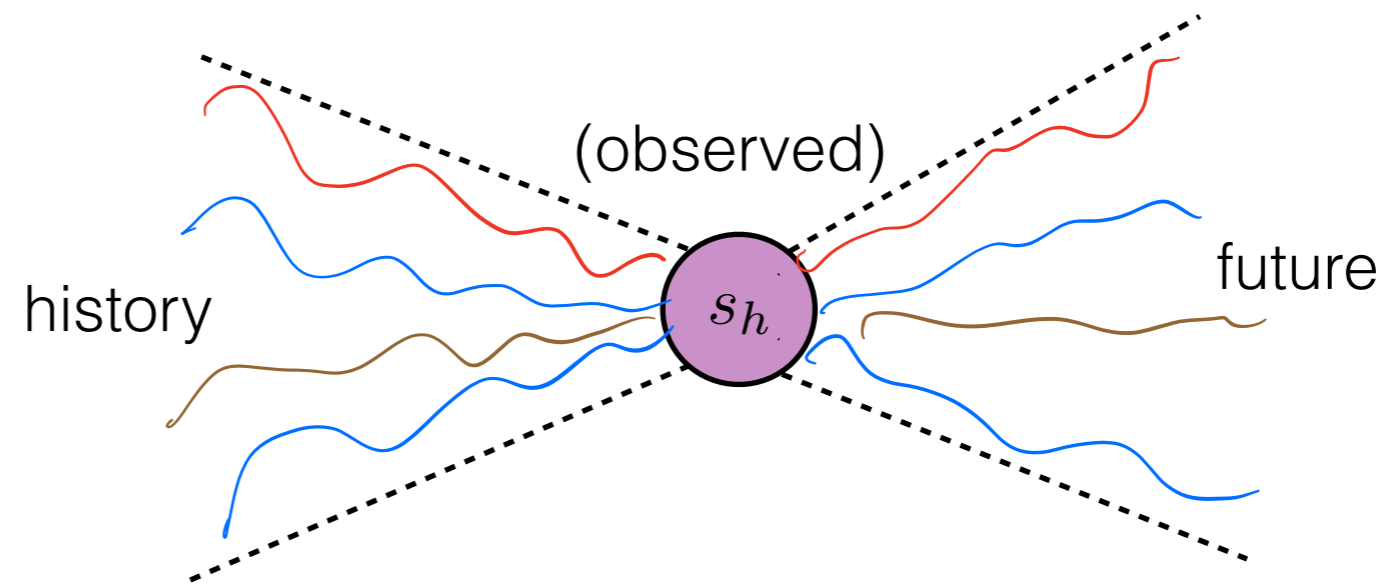
Prediction Groundtruth

$$\mathbb{E}_p[g^2] \leq \|p/q\|_{\infty} \cdot \mathbb{E}_q[g^2]$$

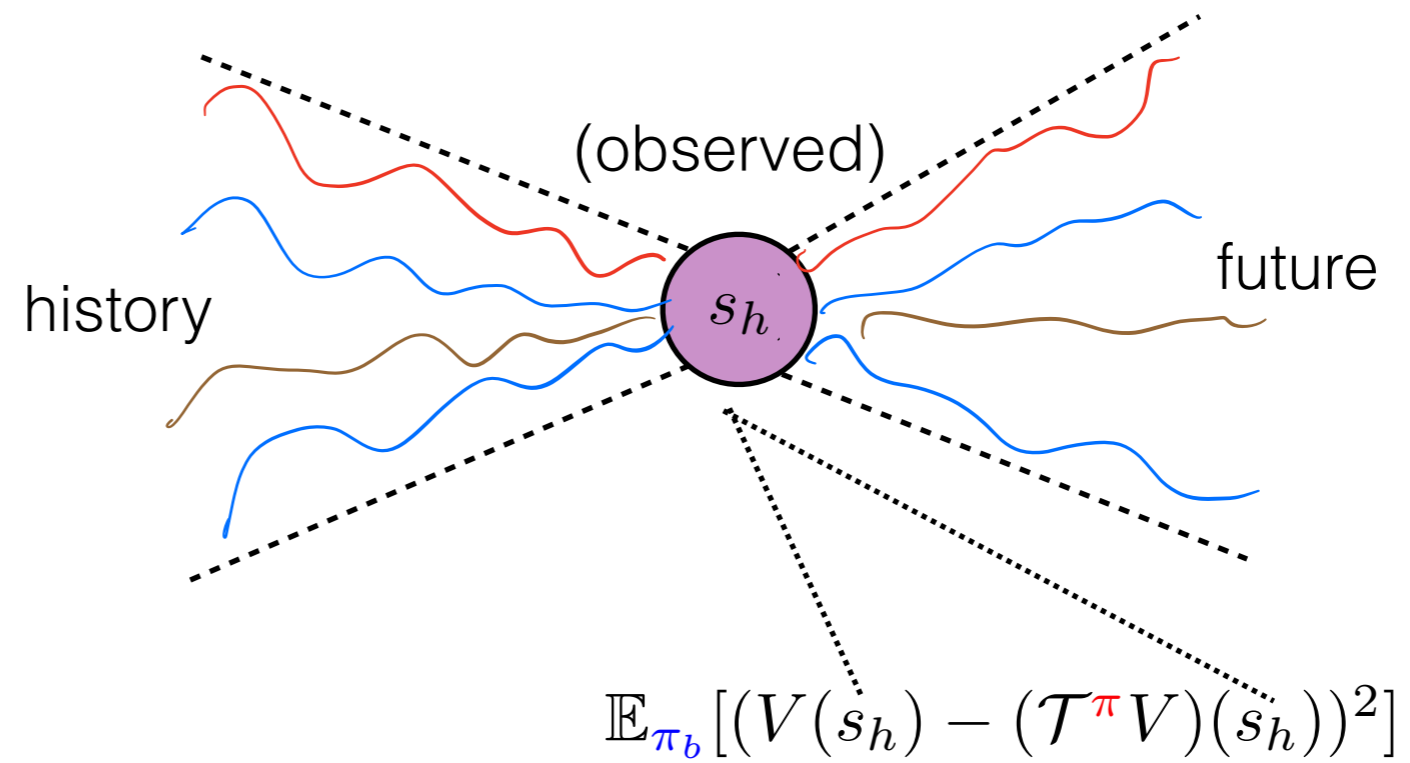
$$\arg \min_V \sum_{h=1}^H \mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$



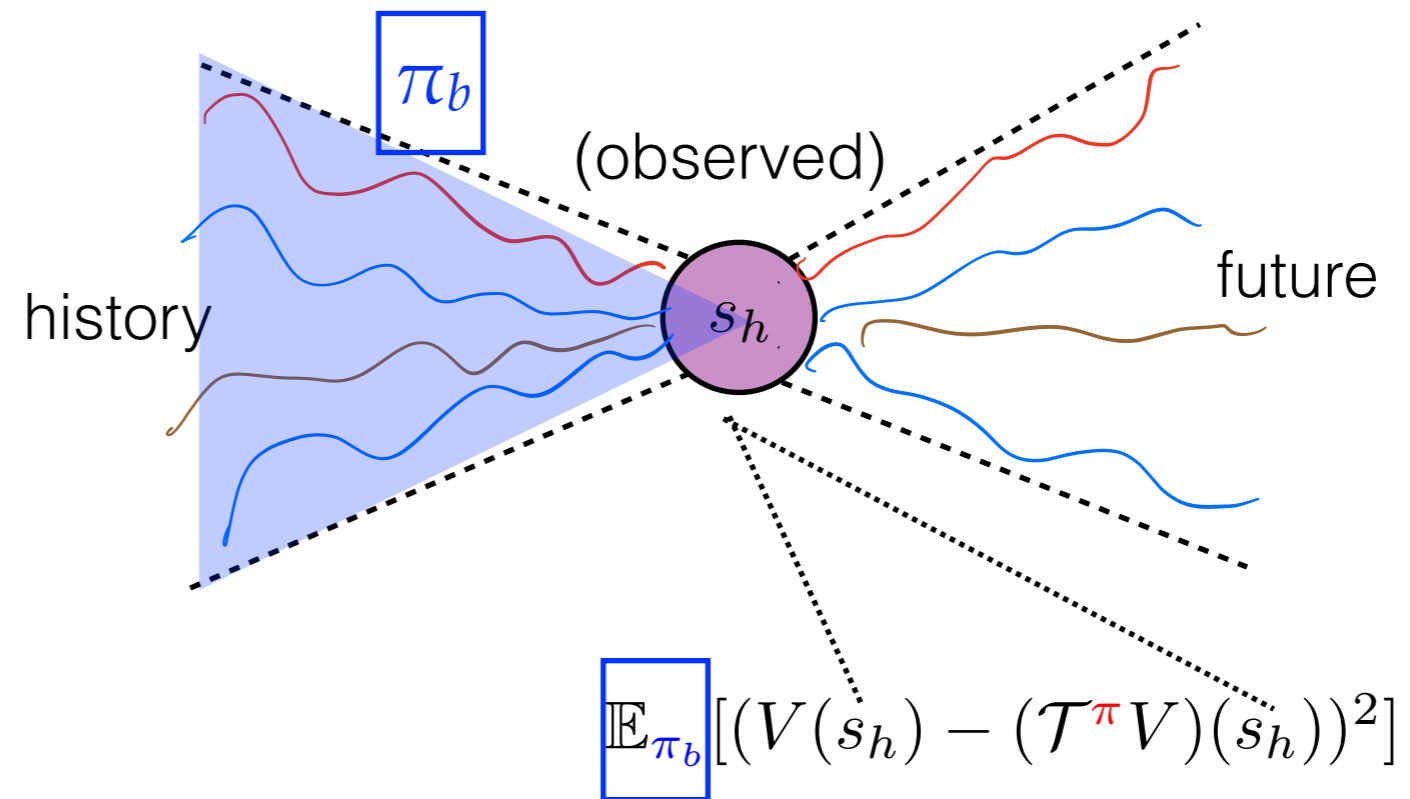
How do value functions help in MDPs?



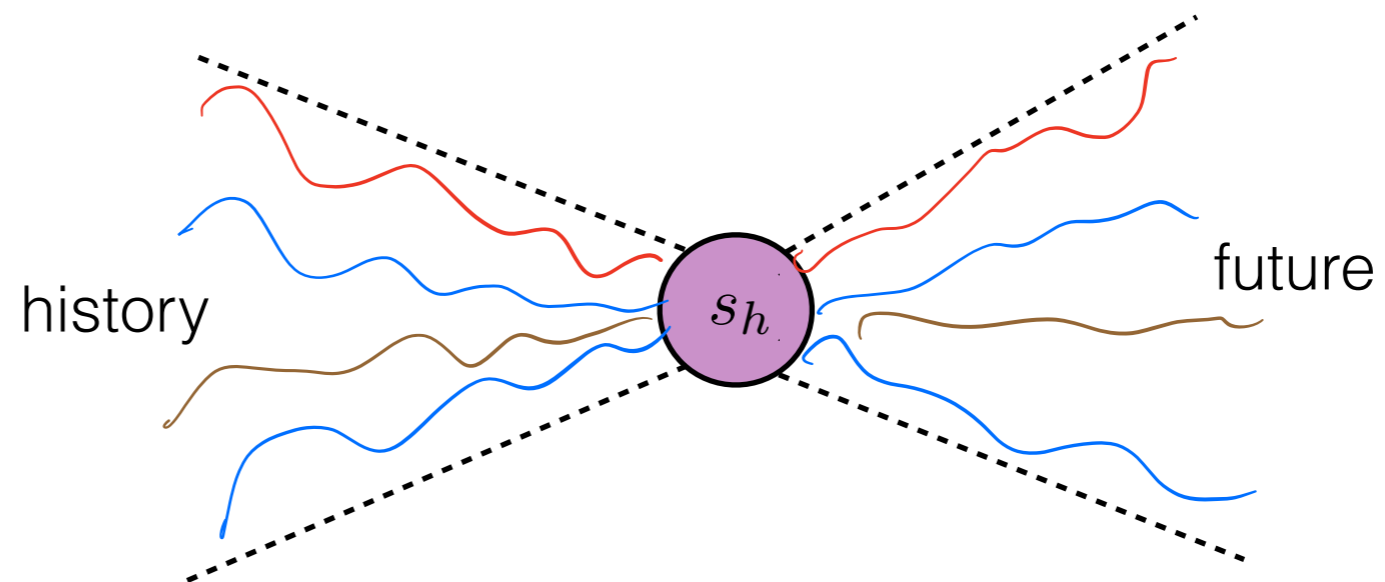
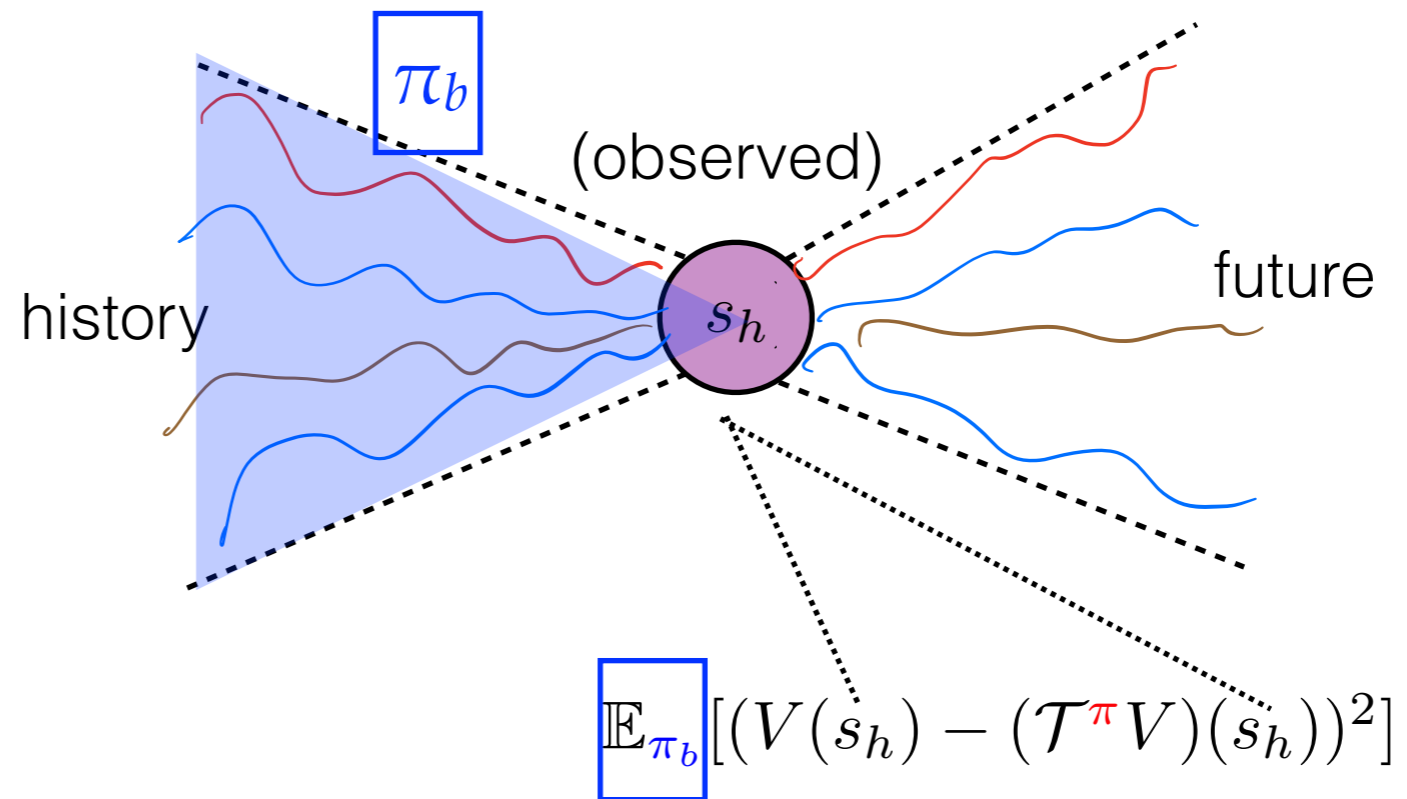
How do value functions help in MDPs?



How do value functions help in MDPs?



How do value functions help in MDPs?

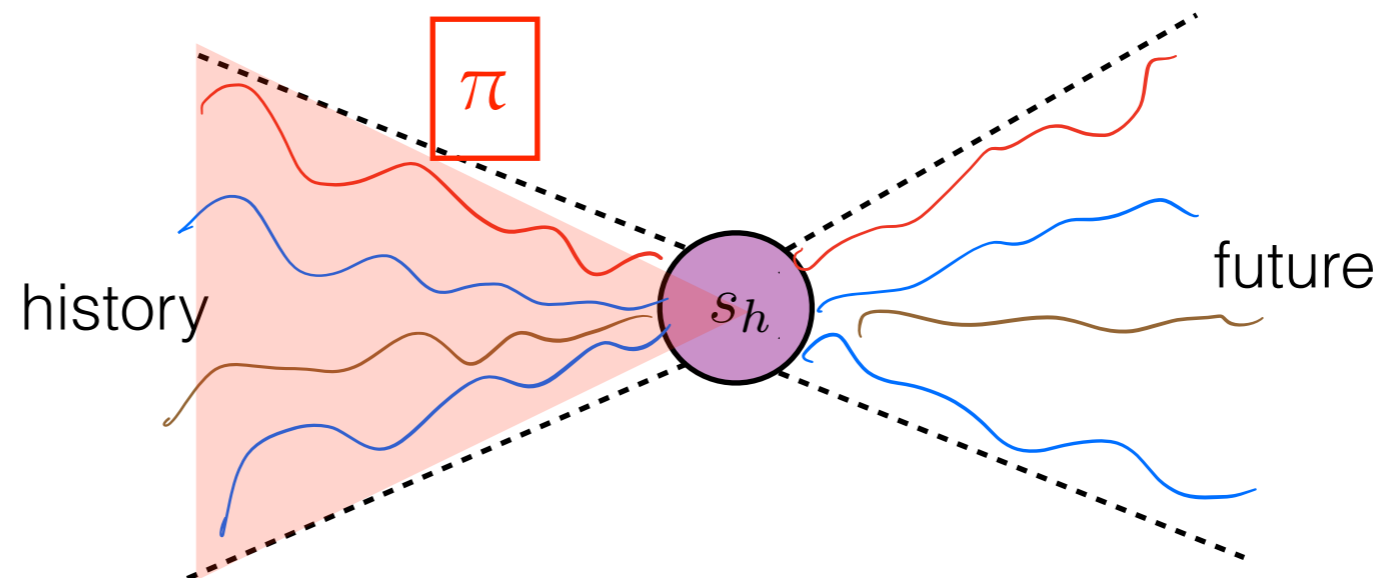
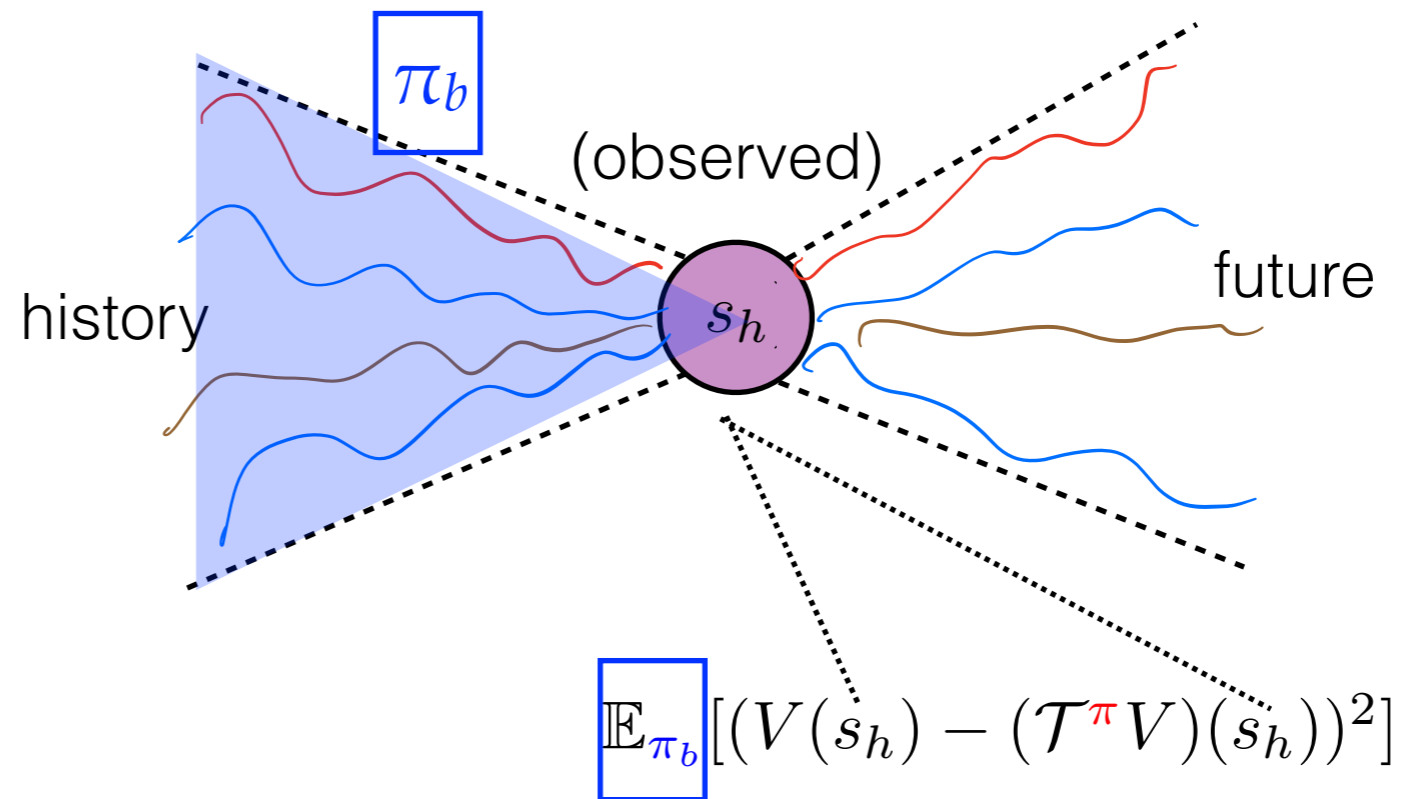


Prediction

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi}V)(s_h)]$$

Groundtruth

How do value functions help in MDPs?

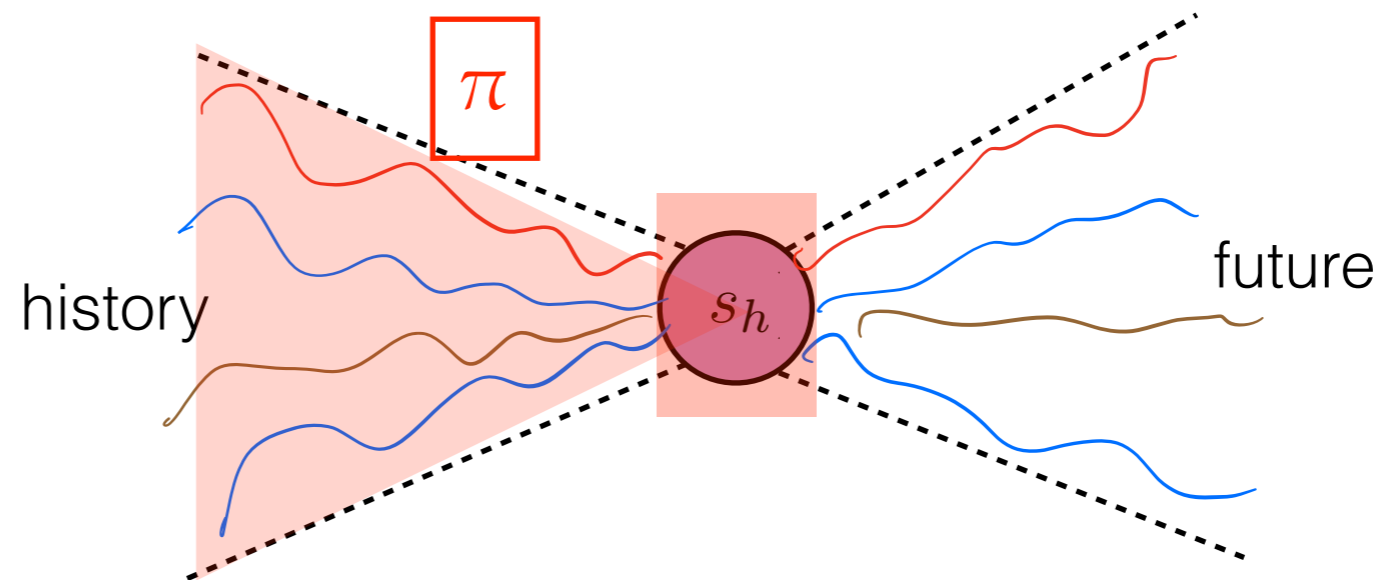
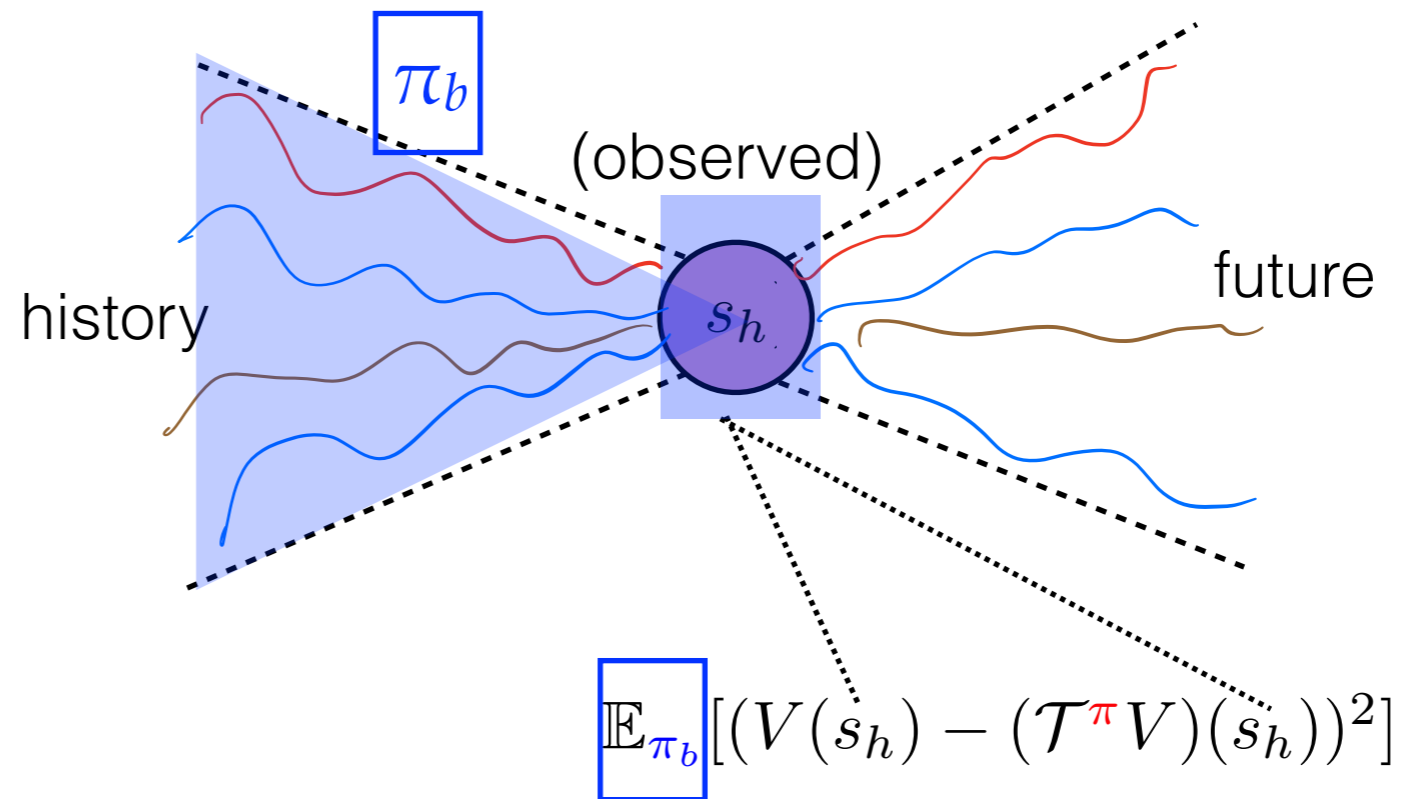


Prediction

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi} V)(s_h)]$$

Groundtruth

How do value functions help in MDPs?

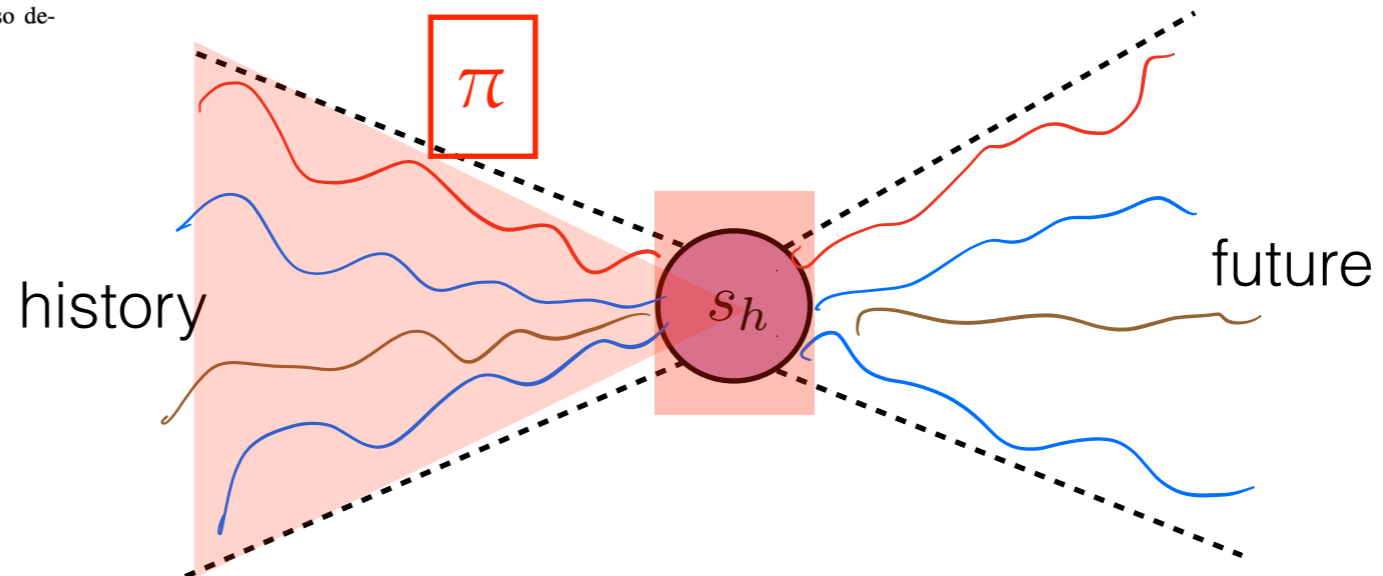
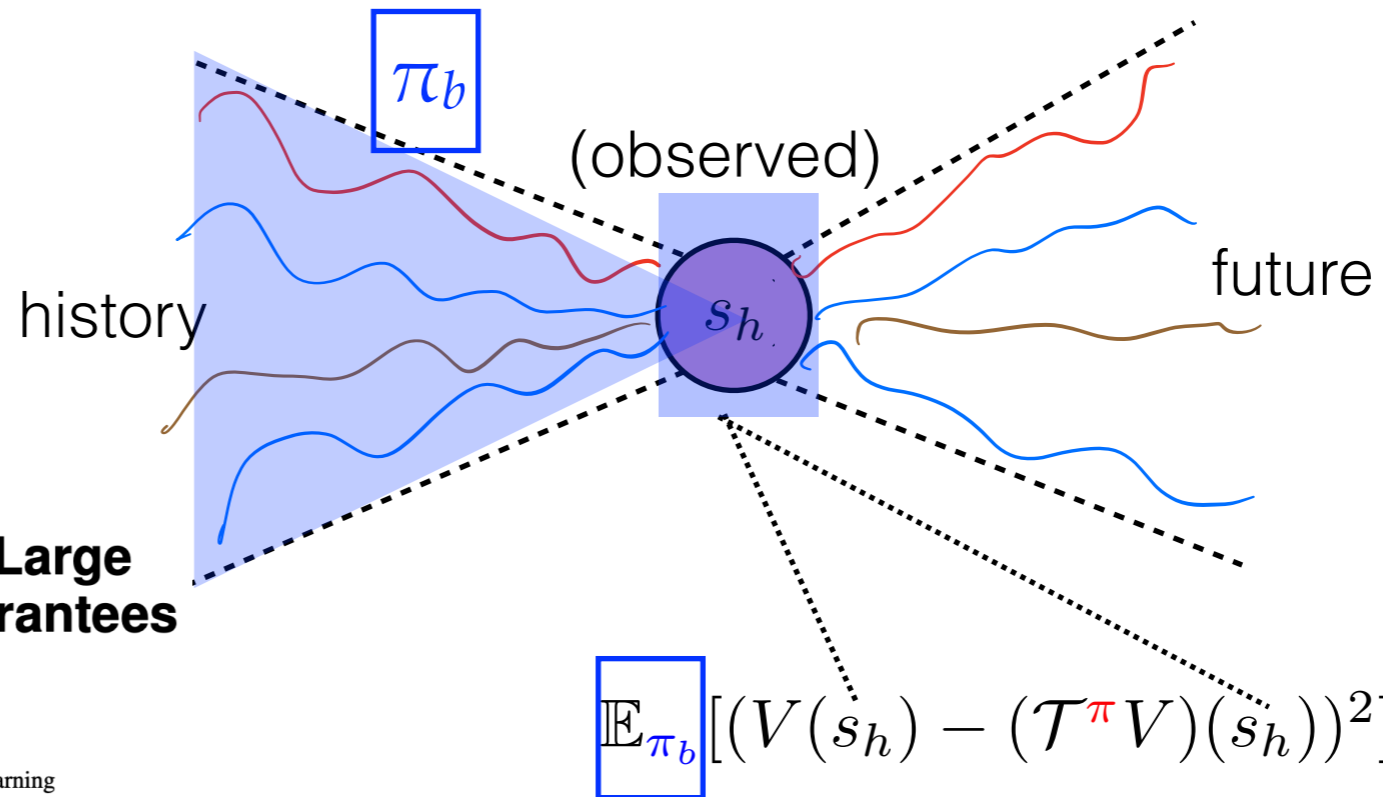


Prediction

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi} V)(s_h)]$$

Groundtruth

How do value functions help in MDPs?



Prediction

$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi} V)(s_h)]$$

Groundtruth

Submitted to Statistical Science

Offline Reinforcement Learning in Large State Spaces: Algorithms and Guarantees

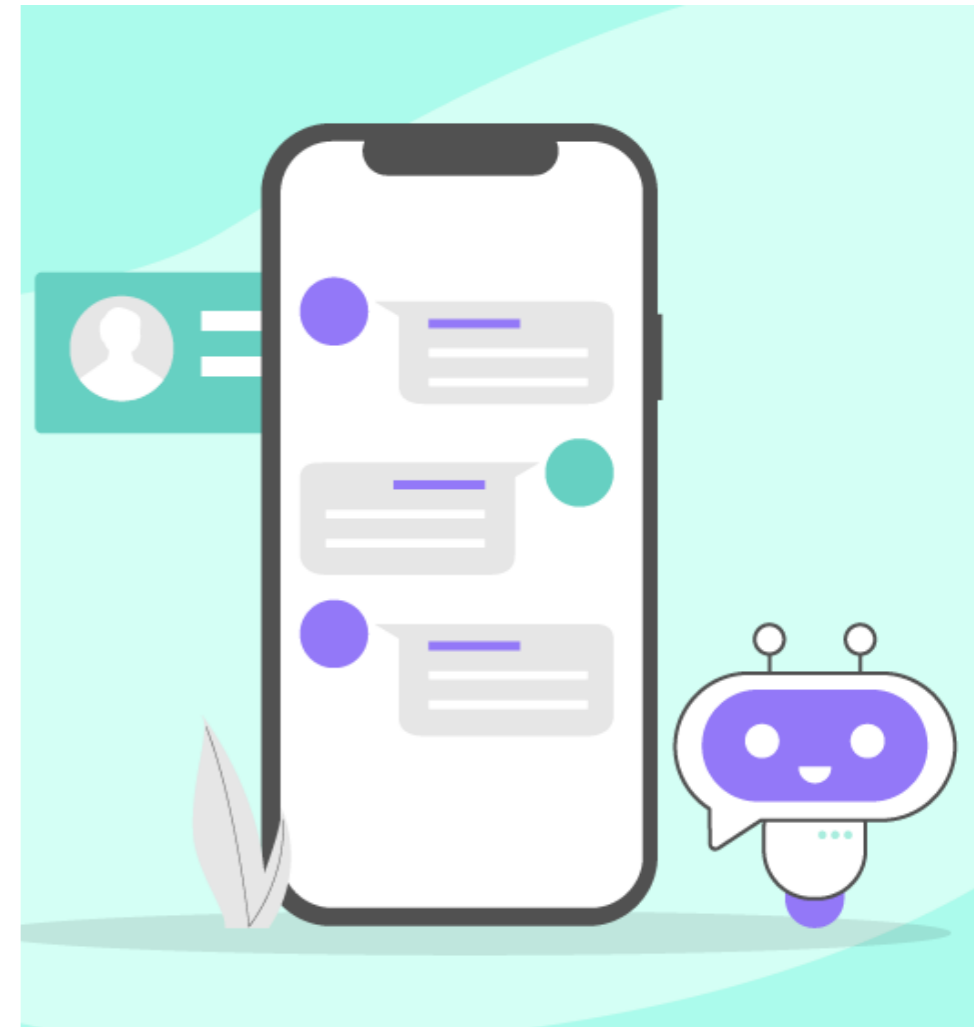
Nan Jiang and Tengyang Xie (draft version; under review)

Abstract. This article introduces the theory of offline reinforcement learning in large state spaces, where good policies are learned from historical data without online interactions with the environment. Key concepts introduced include expressivity assumptions on function approximation (e.g., Bellman-completeness vs. realizability) and data coverage (e.g., all-policy vs. single-policy coverage), and a rich landscape of algorithms and results is described, depending on the assumptions one is willing to make and the sample and computational complexity guarantees one wishes to achieve. We also describe open questions and connections to adjacent areas.

Key words and phrases: offline reinforcement learning.

Partially Observed (non-Markov) Problems

$$O_1, a_1, r_1, \dots, O_h, a_h, r_h, \dots, O_H, a_H, r_H$$

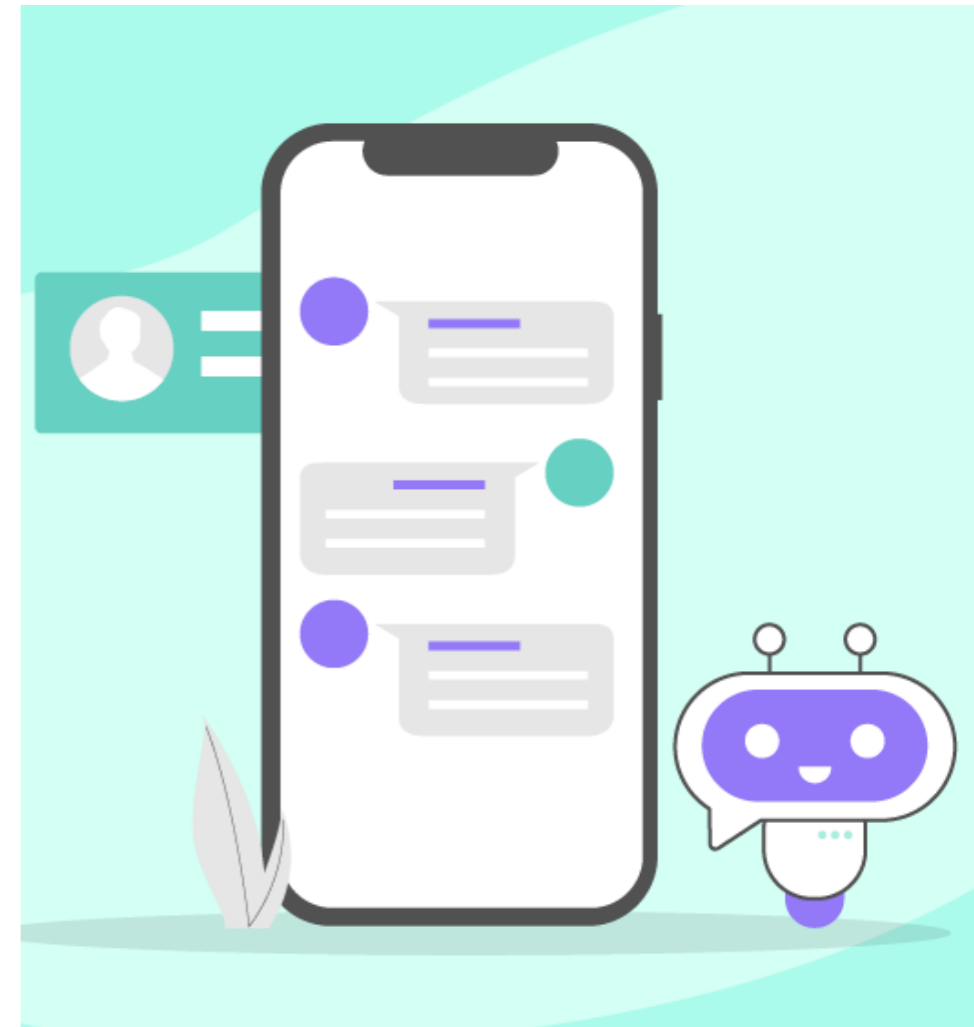


Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H$$

- Define new state τ_h . Problem solved?



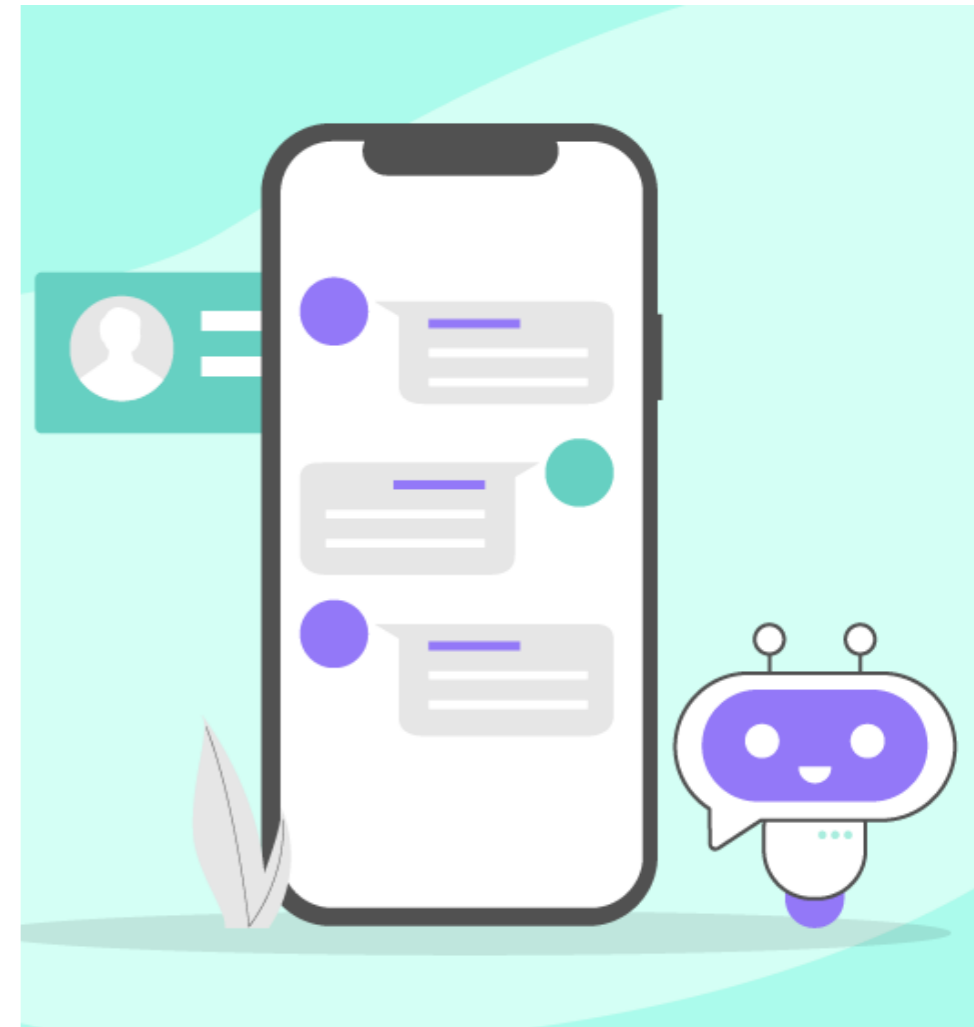
Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H$$


- Define new state τ_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)}$



Partially Observed (non-Markov) Problems

- Can always convert to MDP


$$o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H$$


- Define new state τ_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'} | o_{h'})}{\pi_b(a_{h'} | o_{h'})} \quad !!$

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H$$


- Define new state τ_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'} | o_{h'})}{\pi_b(a_{h'} | o_{h'})}$!!

- Value in RLHF: $> \exp(25)$!!!

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H}$$

- Define new state τ_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \quad !!$

- Structure can help: e.g., if \mathcal{V} is linear in feature $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H}$$

- Define new state \mathcal{T}_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \quad !!$

- Structure can help: e.g., if \mathcal{V} is linear in feature $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H}$$

- Define new state \mathcal{T}_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \quad !!$

- Structure can help: e.g., if \mathcal{V} is linear in feature $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

- this assumes given low-dim linear feature to encode history...

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H}$$

- Define new state \mathcal{T}_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \quad !!$

- Structure can help: e.g., if \mathcal{V} is linear in feature $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$

$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

- this assumes given low-dim linear feature to encode history...
- side q: what structure in \mathcal{V} balances expressivity and coverage

Partially Observed (non-Markov) Problems

- Can always convert to MDP

$$\underbrace{o_1, a_1, r_1, \dots, o_h, a_h, r_h, \dots, o_H, a_H, r_H}$$

- Define new state \mathcal{T}_h . Problem solved?

- State density ratio: $\frac{d^\pi(o_1, a_1, \dots, o_h)}{d^{\pi_b}(o_1, a_1, \dots, o_h)} = \prod_{h'=1}^{h-1} \frac{\pi(a_{h'}|o_{h'})}{\pi_b(a_{h'}|o_{h'})} \quad !!$

- Structure can help: e.g., if \mathcal{V} is linear in feature $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$

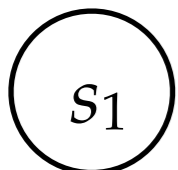
$$\sup_{V \in \mathcal{V}} \frac{|\mathbb{E}_\pi[V - \mathcal{T}^\pi V]|}{\sqrt{\mathbb{E}_{\pi_b}[(V - \mathcal{T}^\pi V)^2]}} \leq \mathbb{E}_\pi[\phi]^\top \mathbb{E}_{\pi_b}[\phi\phi^\top]^{-1} \mathbb{E}_\pi[\phi]$$

- this assumes given low-dim linear feature to encode history...

Can we avoid the **exponentials** in OPE in **PO** settings, without relying on structured function classes?

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in S_h$ (small?)



$h=1$

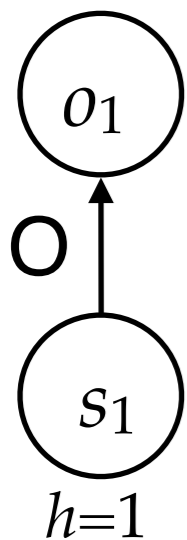
$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$

emission process
 $\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$



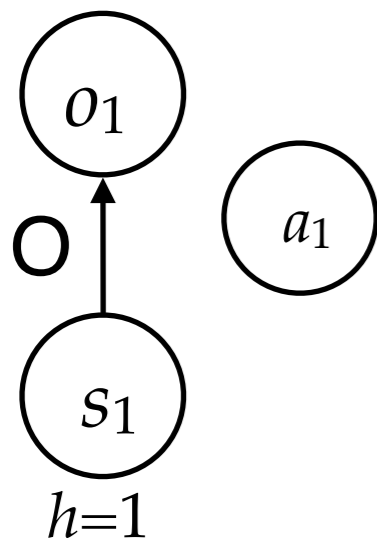
$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)

emission process
 $\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

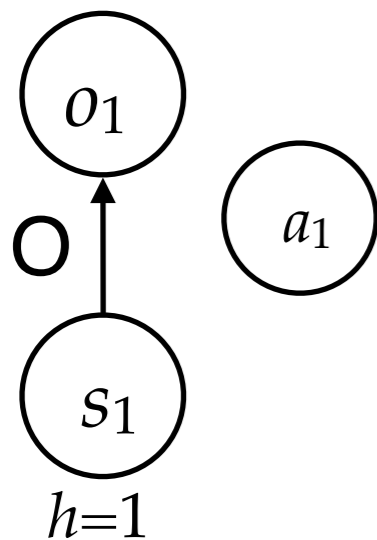
- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$

emission process

$$\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$$

reward function

$$R: \mathcal{S} \times A \rightarrow [0, 1]$$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$

transition dynamics

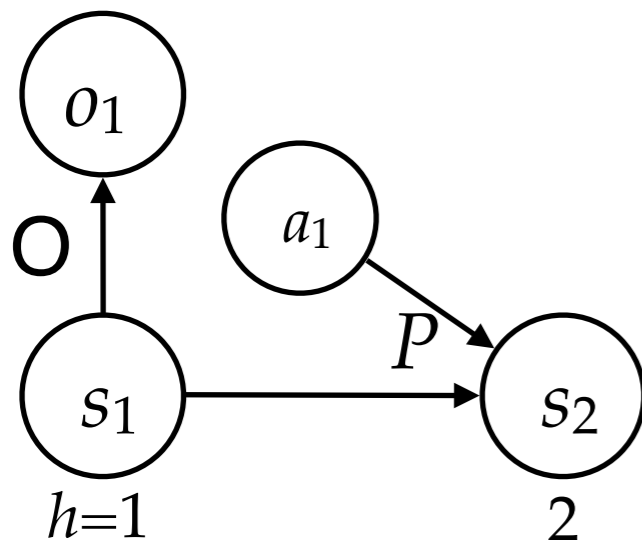
$$P: S \times A \rightarrow \Delta(S)$$

emission process

$$\mathbf{O}: S \rightarrow \Delta(\mathcal{O})$$

reward function

$$R: S \times A \rightarrow [0,1]$$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$

transition dynamics

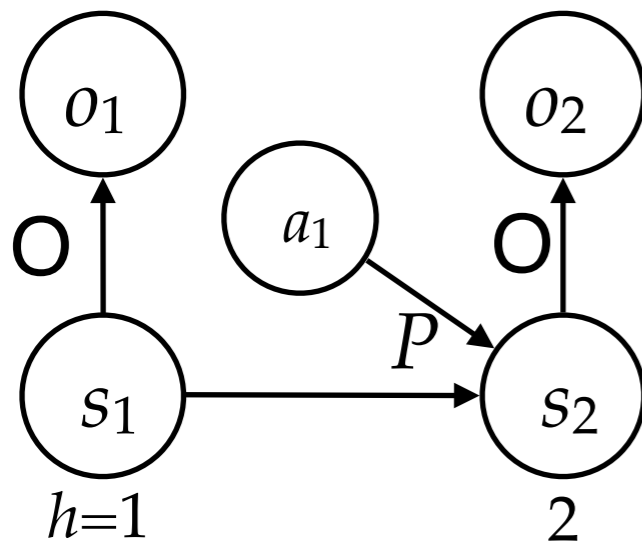
$$P: \mathcal{S} \times A \rightarrow \Delta(\mathcal{S})$$

emission process

$$\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$$

reward function

$$R: \mathcal{S} \times A \rightarrow [0,1]$$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$

transition dynamics

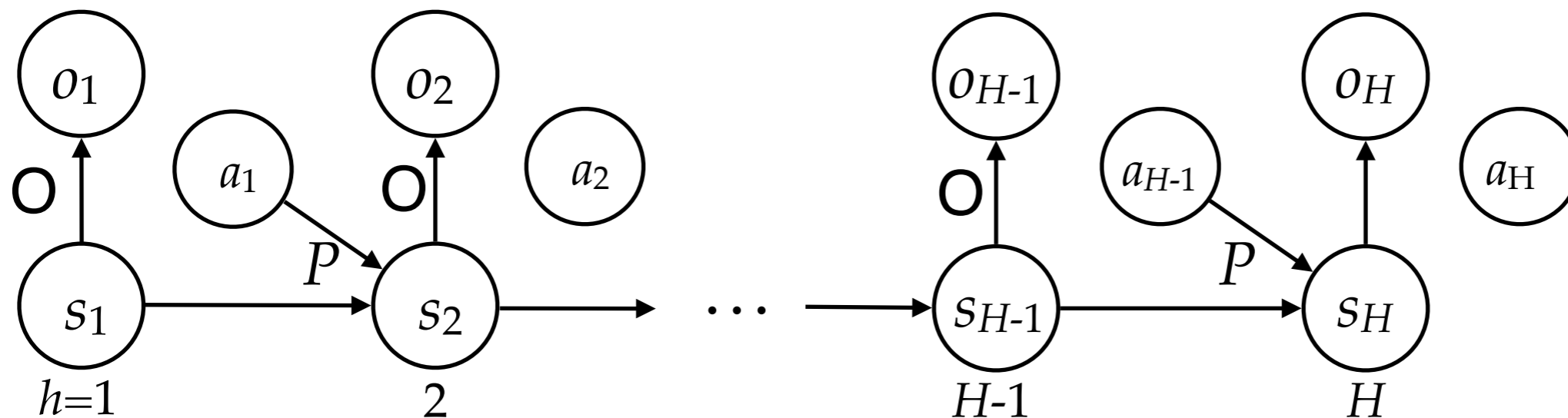
$$P: \mathcal{S} \times A \rightarrow \Delta(\mathcal{S})$$

emission process

$$\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$$

reward function

$$R: \mathcal{S} \times A \rightarrow [0, 1]$$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in \mathcal{A}$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$
- *Memoryless* policies $\pi : \mathcal{O} \rightarrow \Delta(\mathcal{A})$

transition dynamics

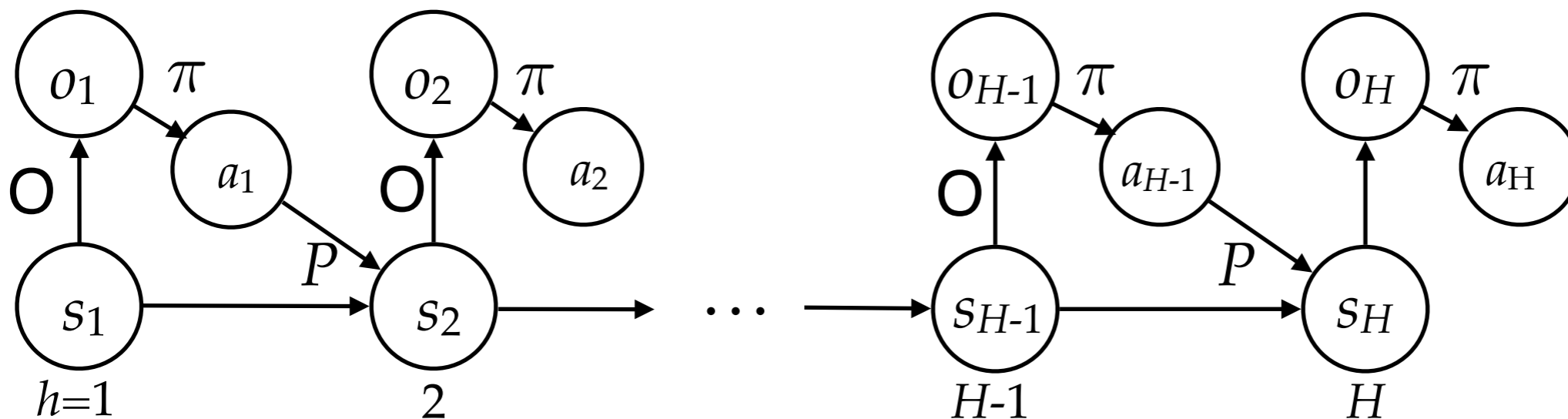
$$P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$$

emission process

$$\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$$

reward function

$$R: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

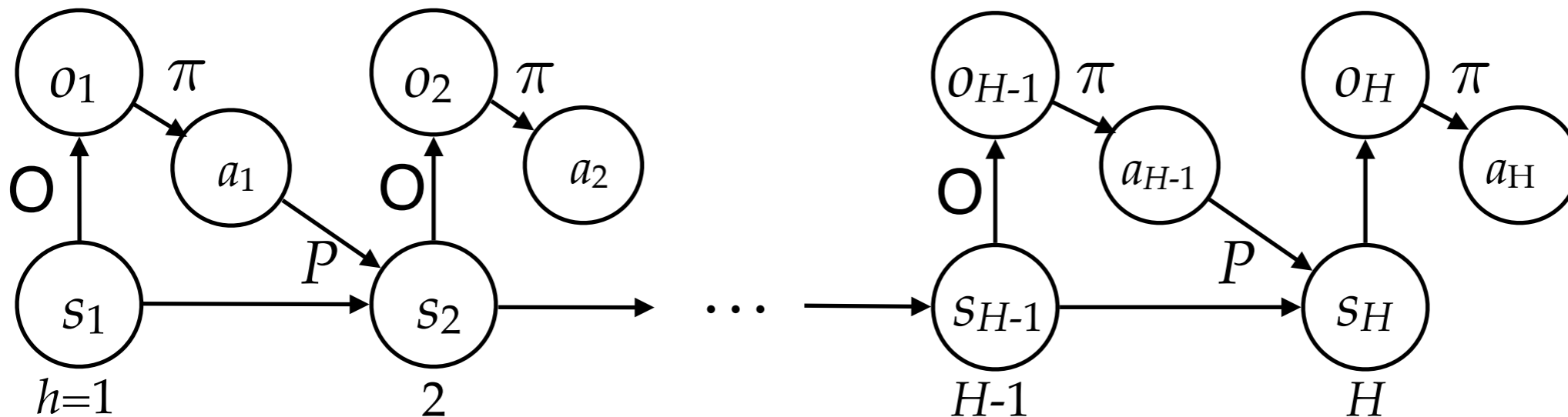
Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in A$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$
- Memoryless* policies $\pi : \mathcal{O} \rightarrow \Delta(\mathcal{A})$
- Goal: estimate $J(\pi) := \mathbb{E}_\pi [\sum_{h=1}^H r_h]$ using episodes collected by π_b

transition dynamics
 $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$

emission process
 $\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$

reward function
 $R: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$



$$\mathcal{S} = \bigcup_h \mathcal{S}_h$$

$$\mathcal{O} = \bigcup_h \mathcal{O}_h$$

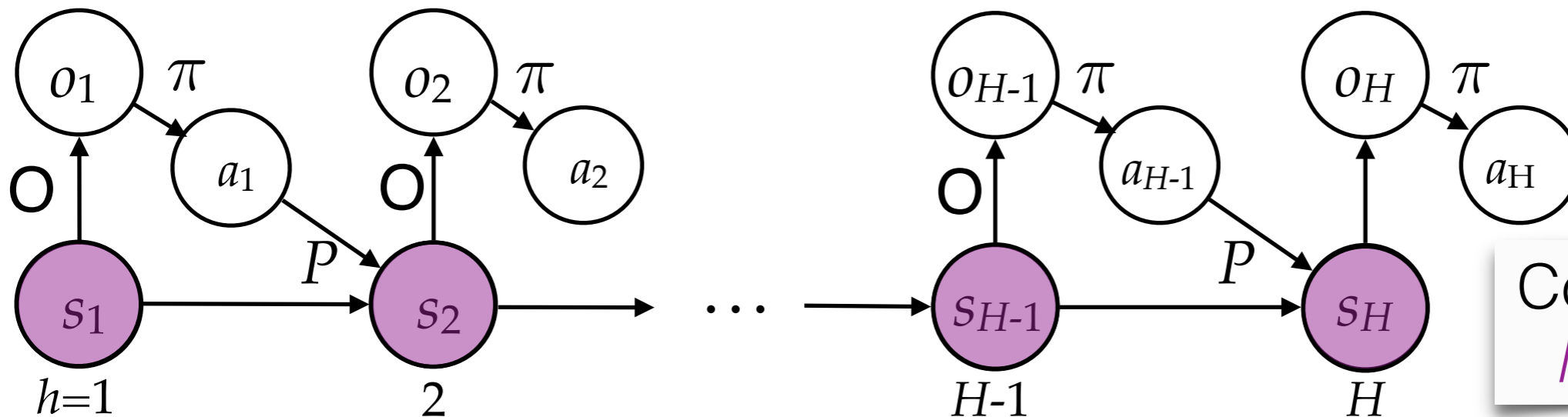
Partially Observable MDPs (POMDPs)

- For $h = 1, 2, \dots, H$,
 - nature generates *latent state* $s_h \in \mathcal{S}_h$ (small?)
 - $s_h \sim P(\cdot | s_{h-1}, a_{h-1})$ for $h \geq 2$
 - agent *observes* $o_h \in \mathcal{O}_h$ (large), $o_h \sim \mathbf{O}(\cdot | s_h)$
 - chooses *action* $a_h \in \mathcal{A}$ (small)
 - receives *reward* $r_h = R(o_h, a_h)$
- Memoryless* policies $\pi : \mathcal{O} \rightarrow \Delta(\mathcal{A})$
- Goal: estimate $J(\pi) := \mathbb{E}_\pi [\sum_{h=1}^H r_h]$ using episodes collected by π_b

transition dynamics
 $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$

emission process
 $\mathbf{O}: \mathcal{S} \rightarrow \Delta(\mathcal{O})$

reward function
 $R: \mathcal{S} \times \mathcal{A} \rightarrow [0,1]$

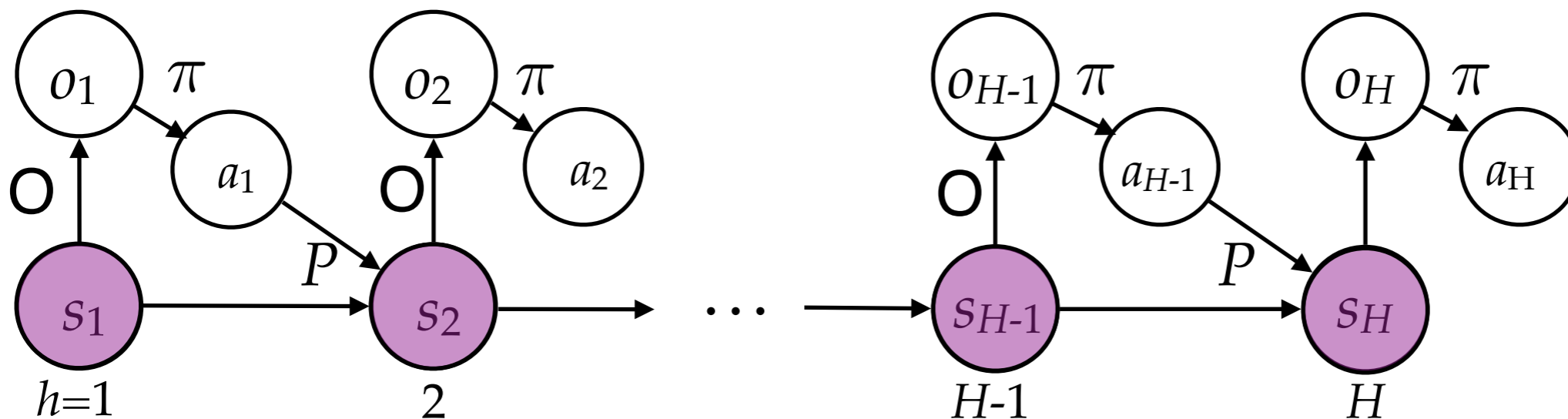


Coverage over
latent state?

Future-Dependent Value Function

- Define: value function of latent state

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

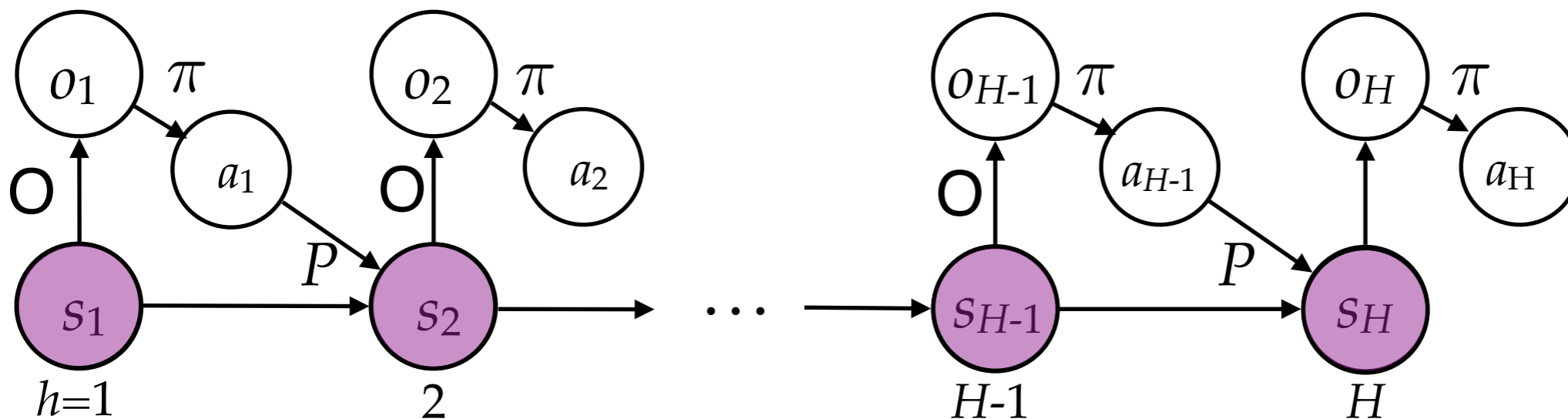


Future-Dependent Value Function

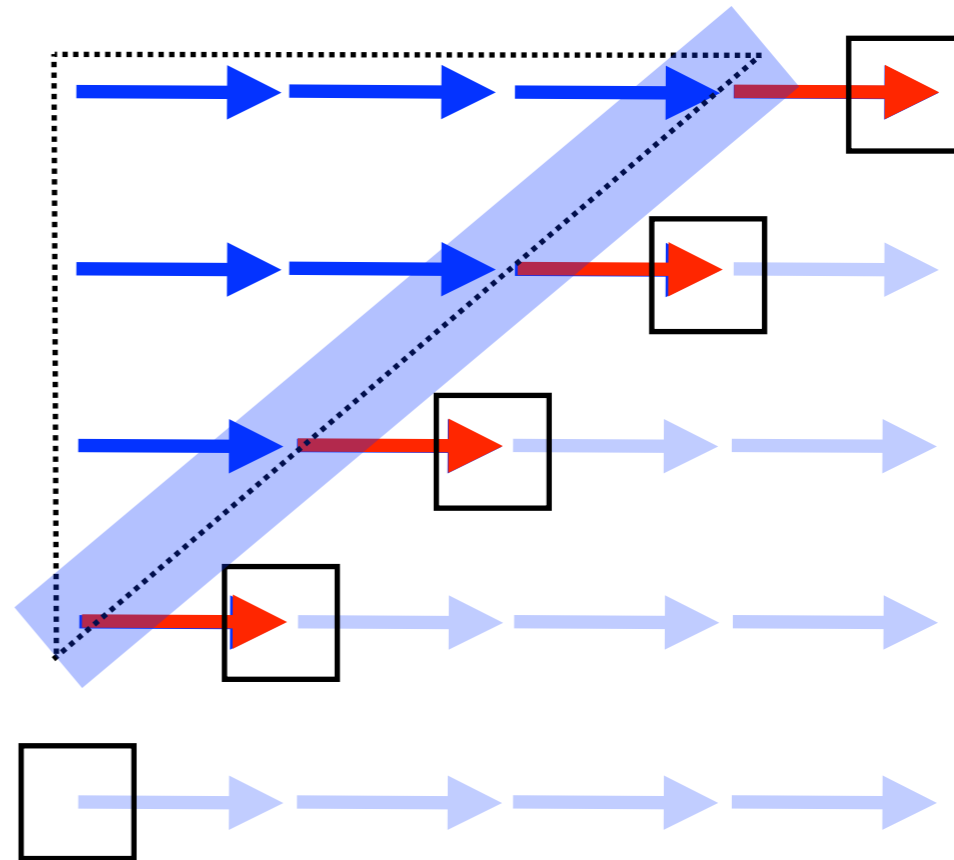
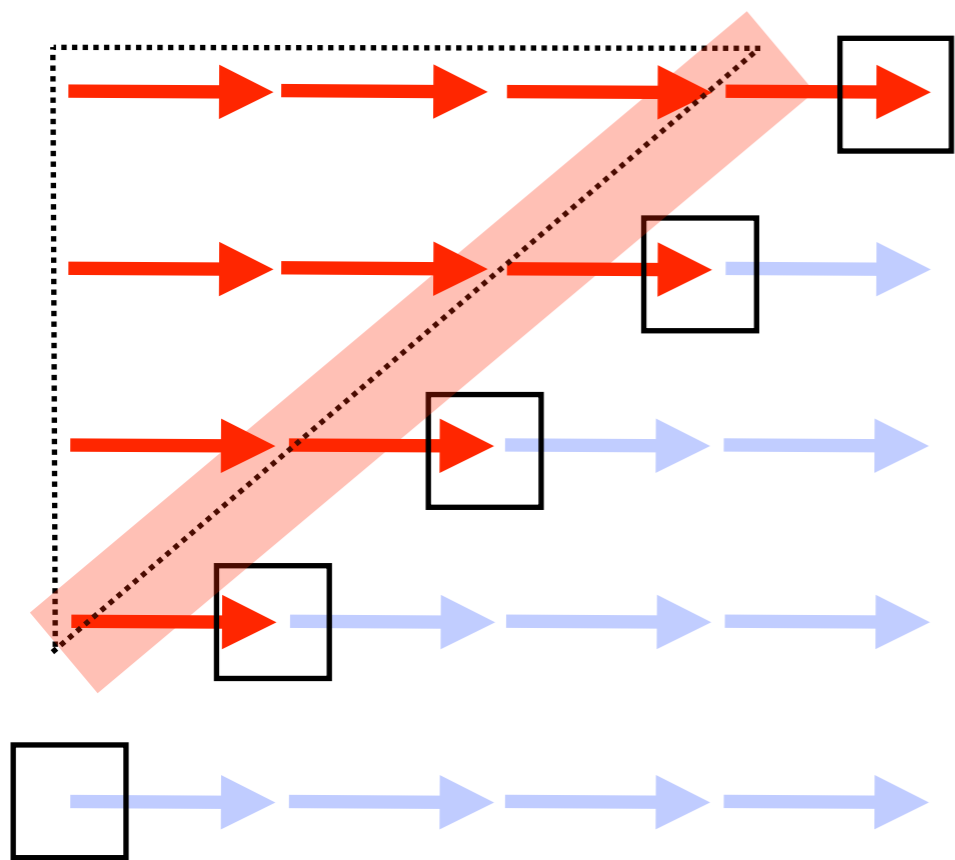
- Define: **value function** of **latent state**

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

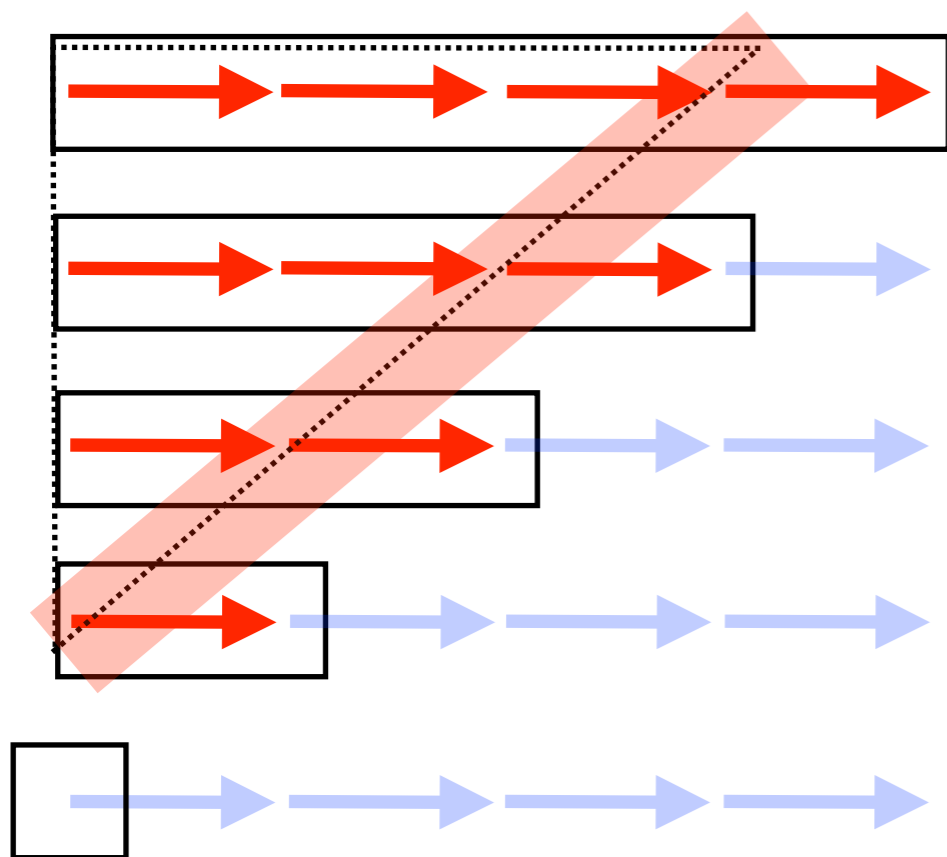
- Problem: s_h is latent — can't even use this function!



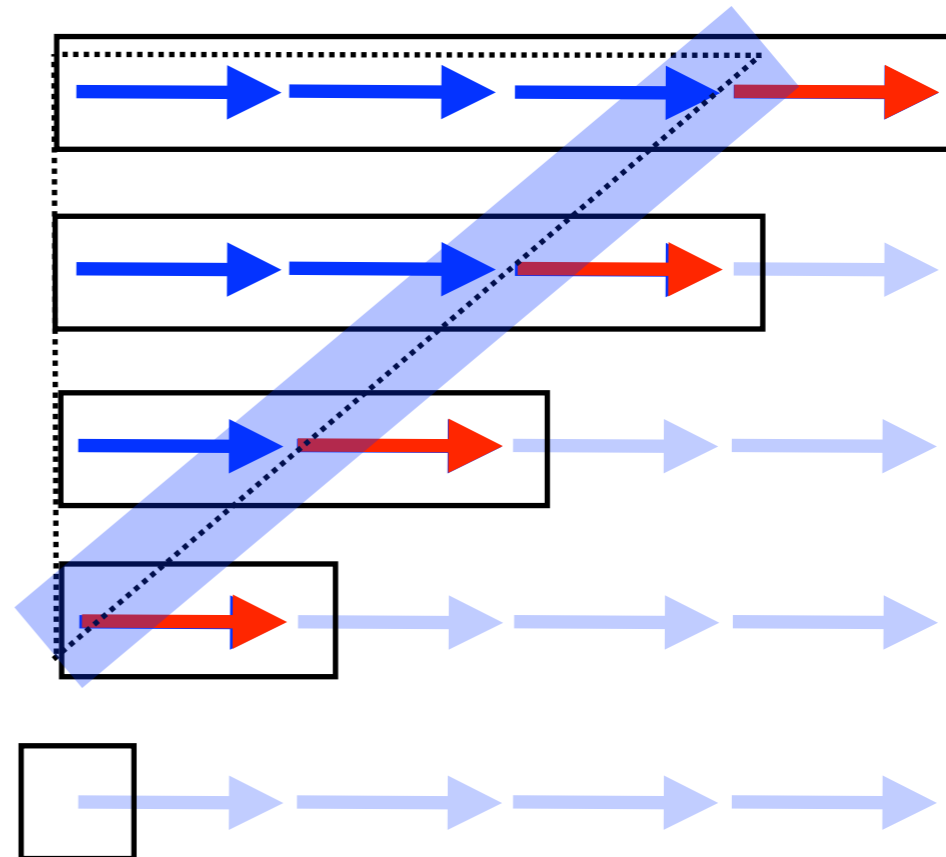
Value function?



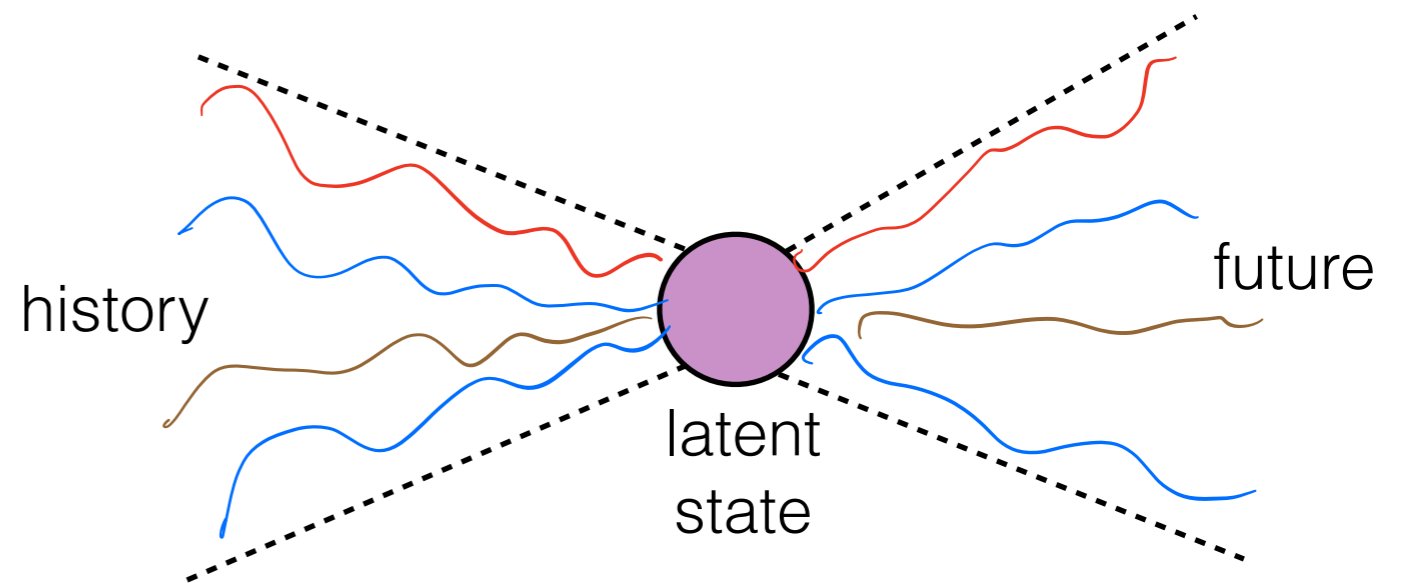
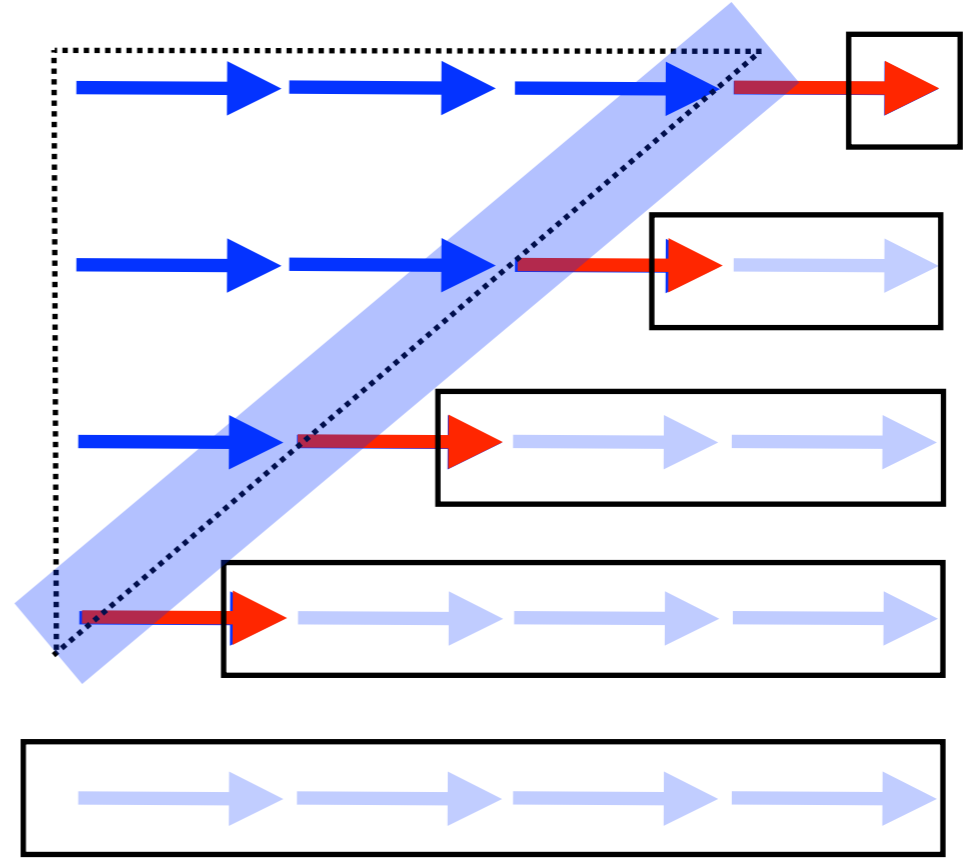
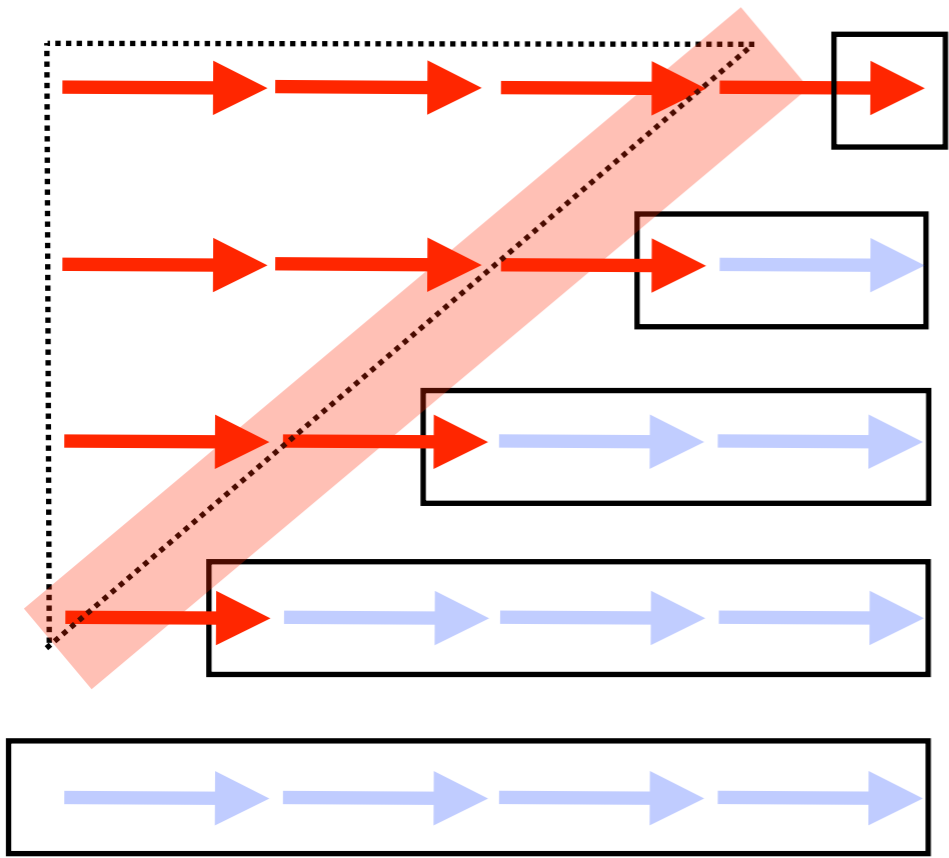
Value function?



X



Value function?

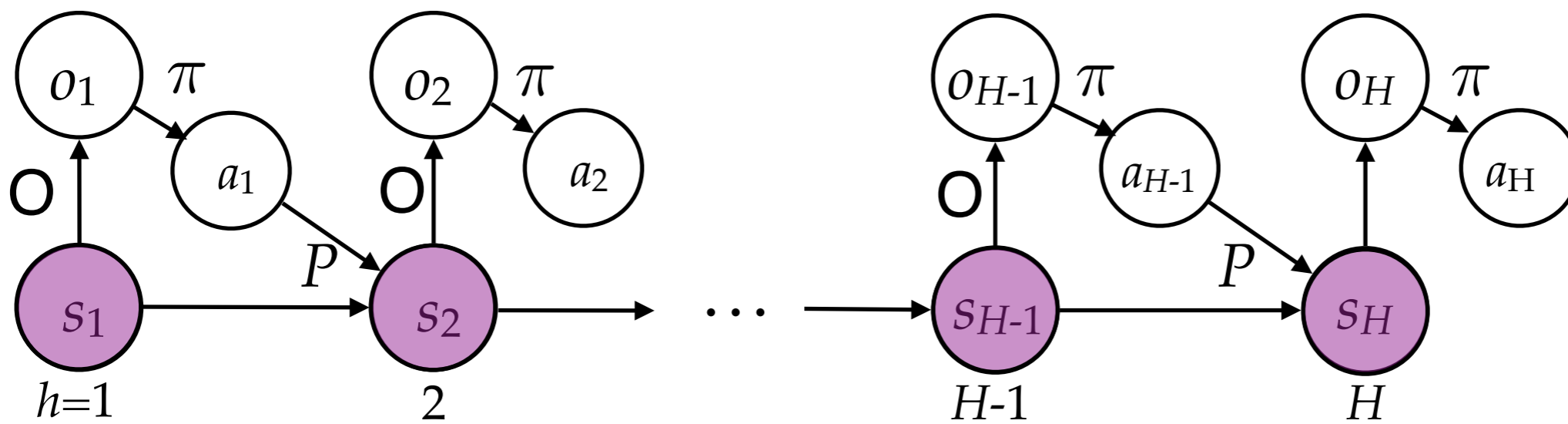


Future-Dependent Value Function

- Define: **value function** of **latent state**

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

- Problem: s_h is latent — can't even use this function!
- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using **future** as input!

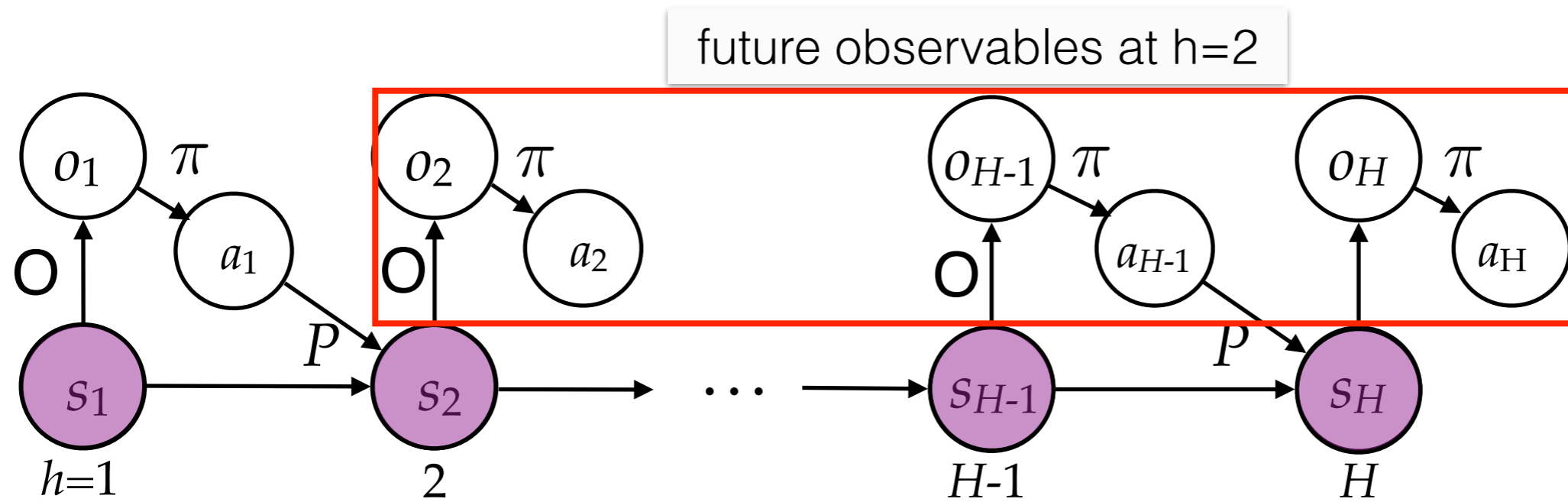


Future-Dependent Value Function

- Define: **value function** of **latent state**

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

- Problem: s_h is latent — can't even use this function!
- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using **future** as input!
 - $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) \mid s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

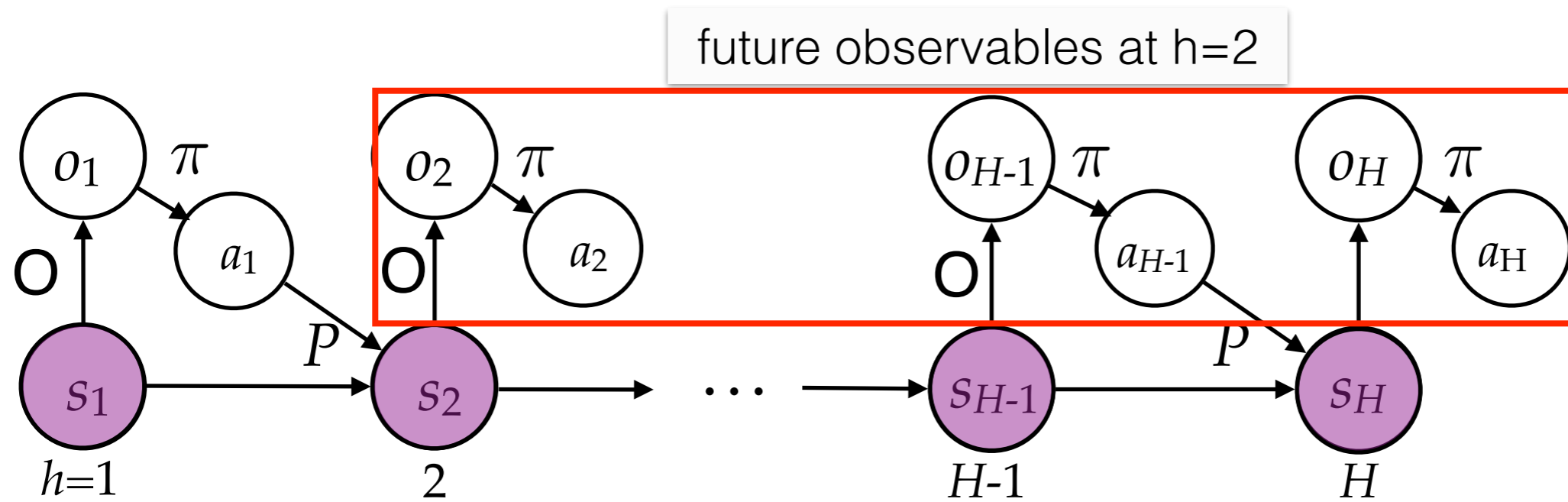


Future-Dependent Value Function

- Define: **value function** of **latent state**

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

- Problem: s_h is latent — can't even use this function!
- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using **future** as input!
 - $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) \mid s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
 - Long history for using future (prediction) as state
 - OOM, SMA, PSR, etc - see Thon & Jaeger'15

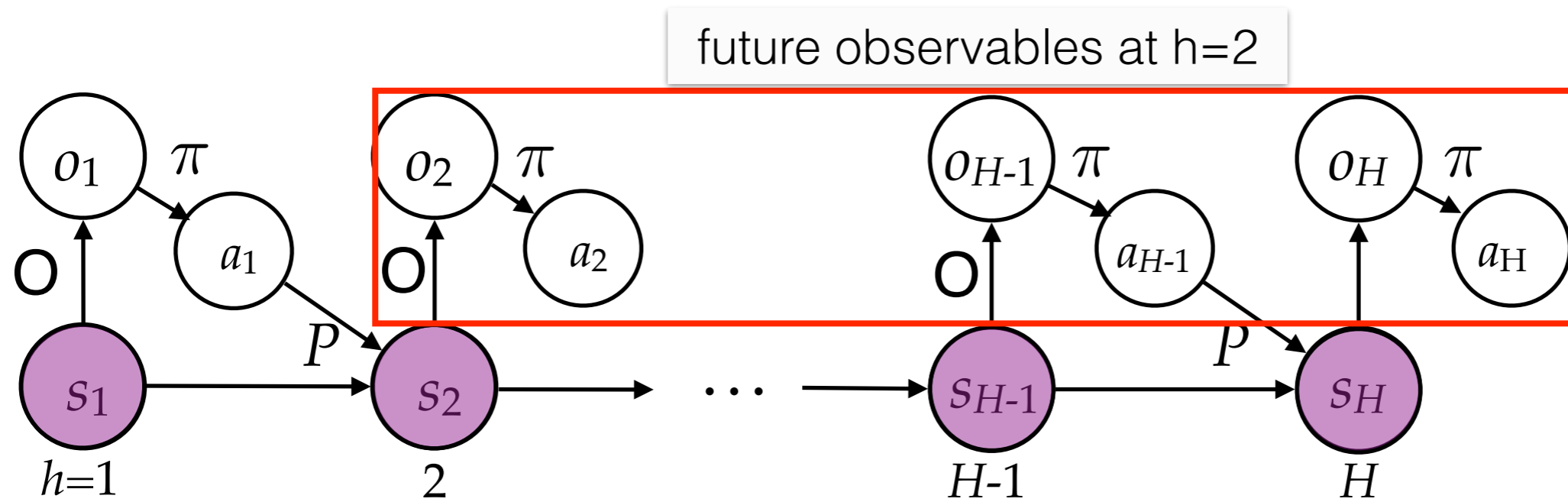


Future-Dependent Value Function

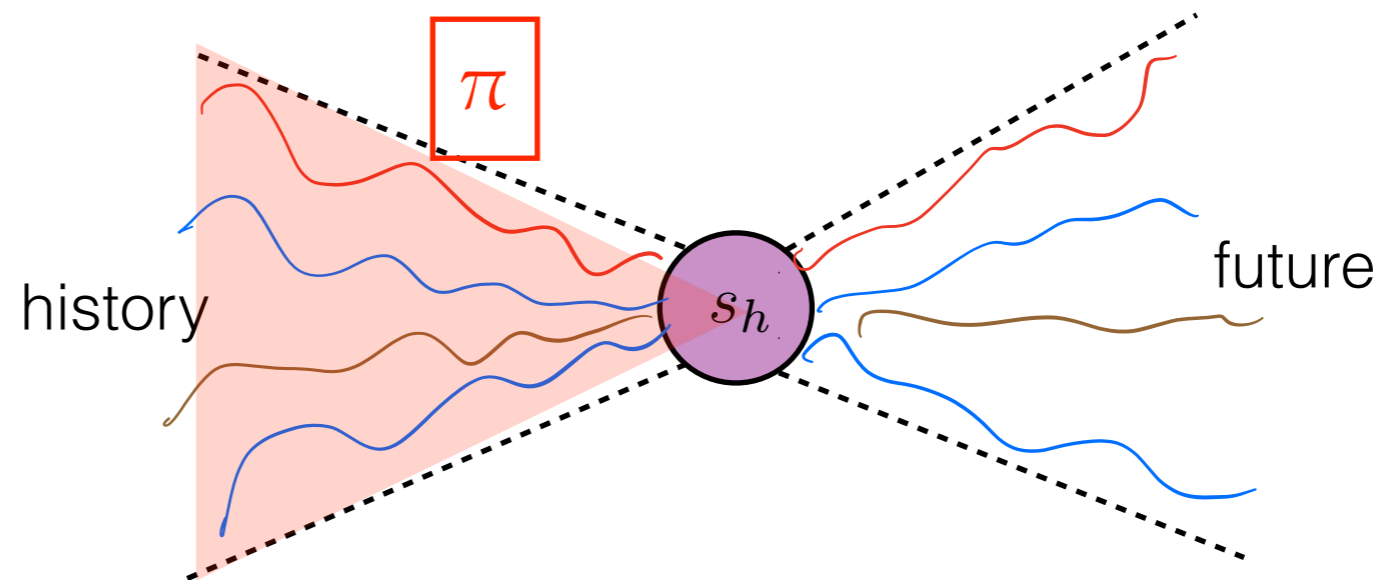
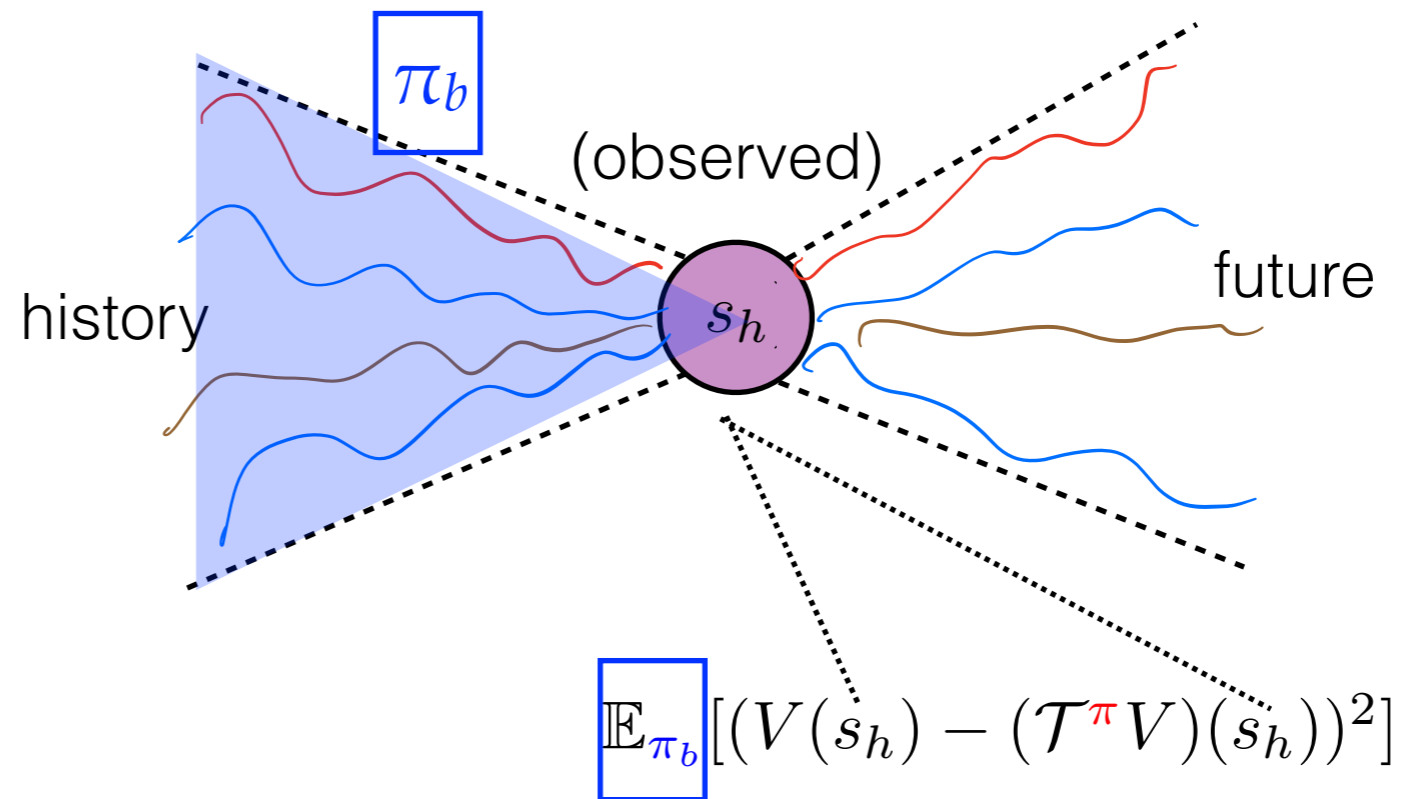
- Define: **value function** of **latent state**

$$V_{\mathcal{S}}^{\pi}(s_h) := \mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'} \mid s_h \right] \in [0, H]$$

- Problem: s_h is latent — can't even use this function!
- Solution: $V_{\mathcal{F}}^{\pi}$ as proxy of $V_{\mathcal{S}}^{\pi}$, using **future** as input!
 - $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) \mid s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
 - Long history for using future (prediction) as state
 - OOM, SMA, PSR, etc - see Thon & Jaeger'15
 - Distribution of future observables is low-rank ($\leq |S|$)

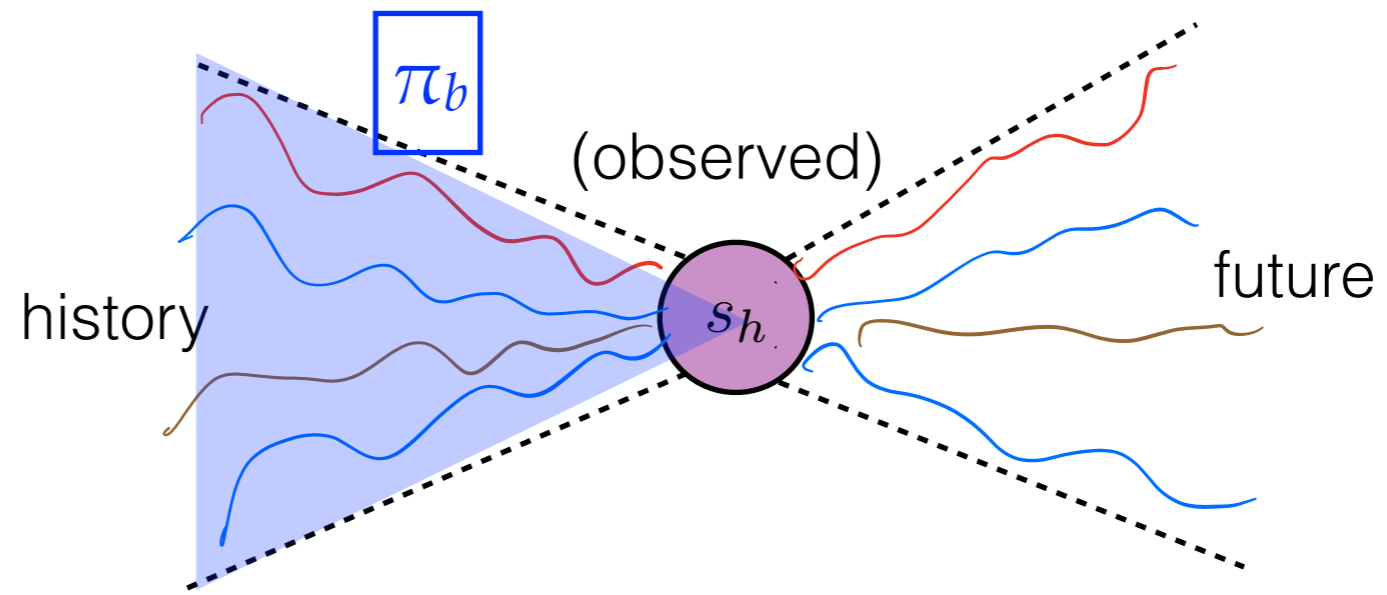


Markov case

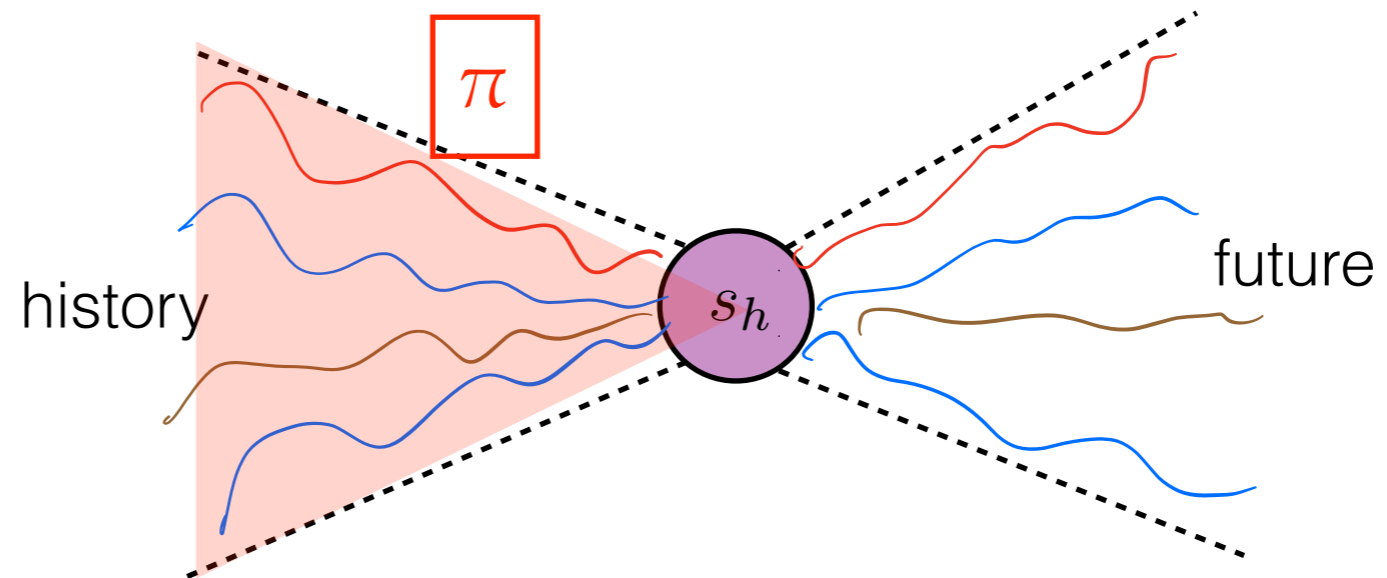


$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi}V)(s_h)]$$

Markov case

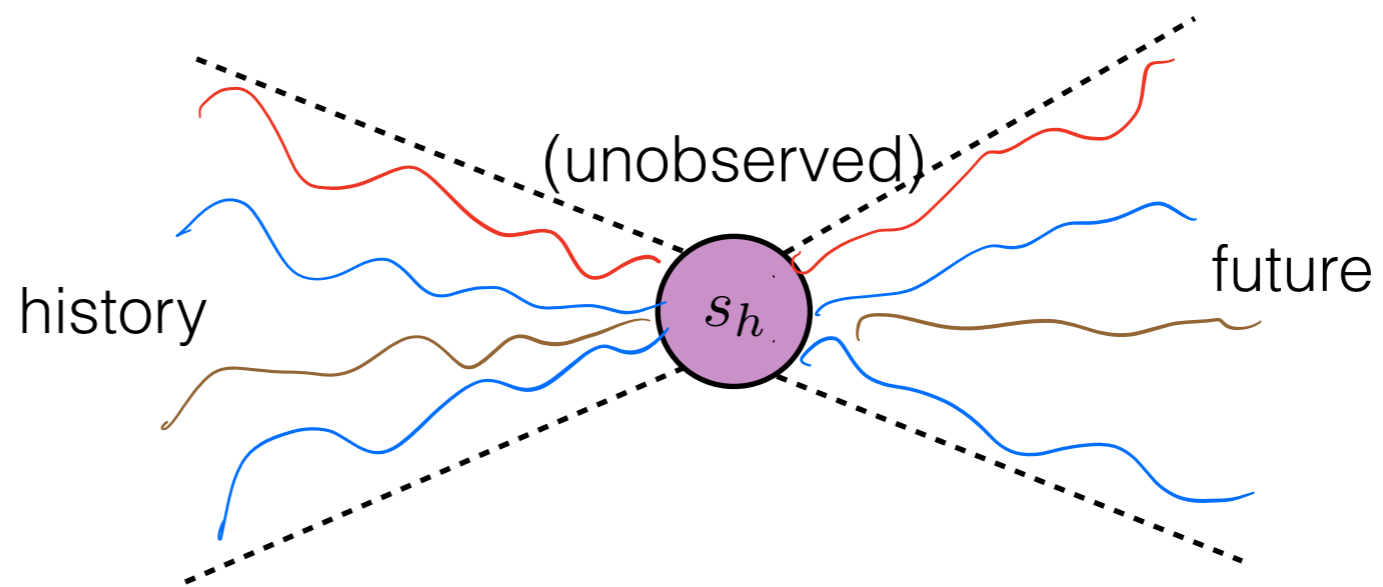


$$\mathbb{E}_{\pi_b} [(V(s_h) - (\mathcal{T}^{\pi} V)(s_h))^2]$$

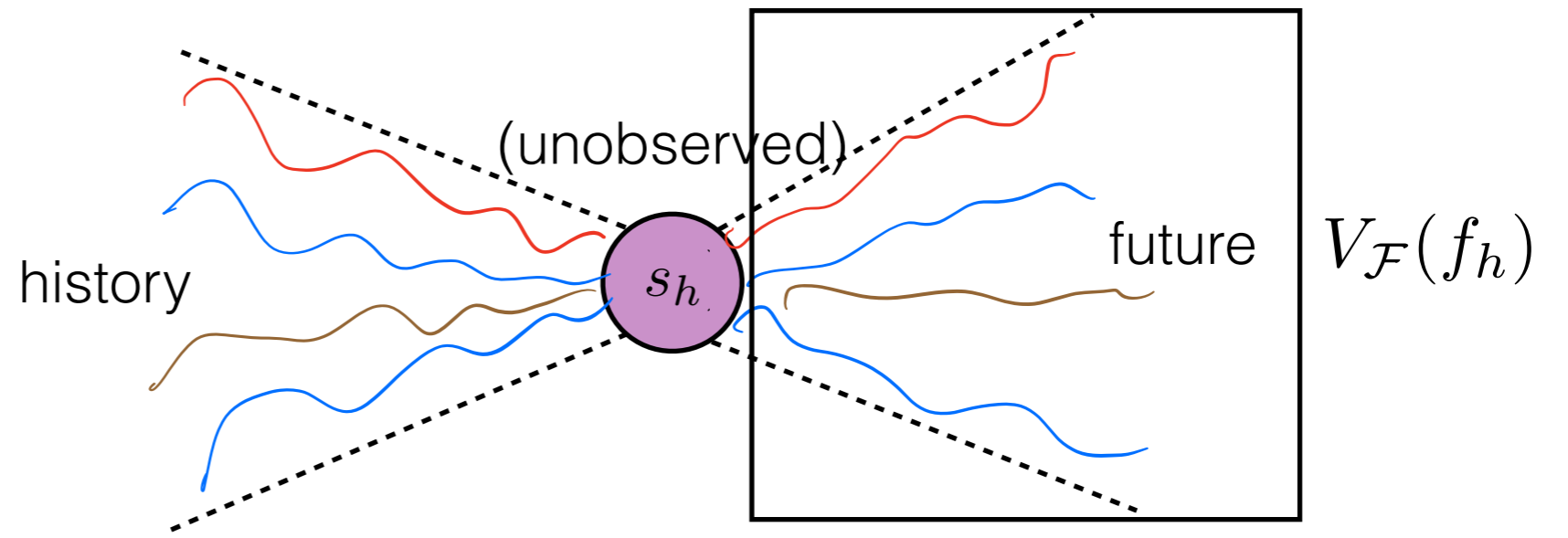


$$\mathbb{E}[V(s_1)] - J(\pi) = \sum_{h=1}^H \mathbb{E}_{\pi} [V(s_h) - (\mathcal{T}^{\pi} V)(s_h)]$$

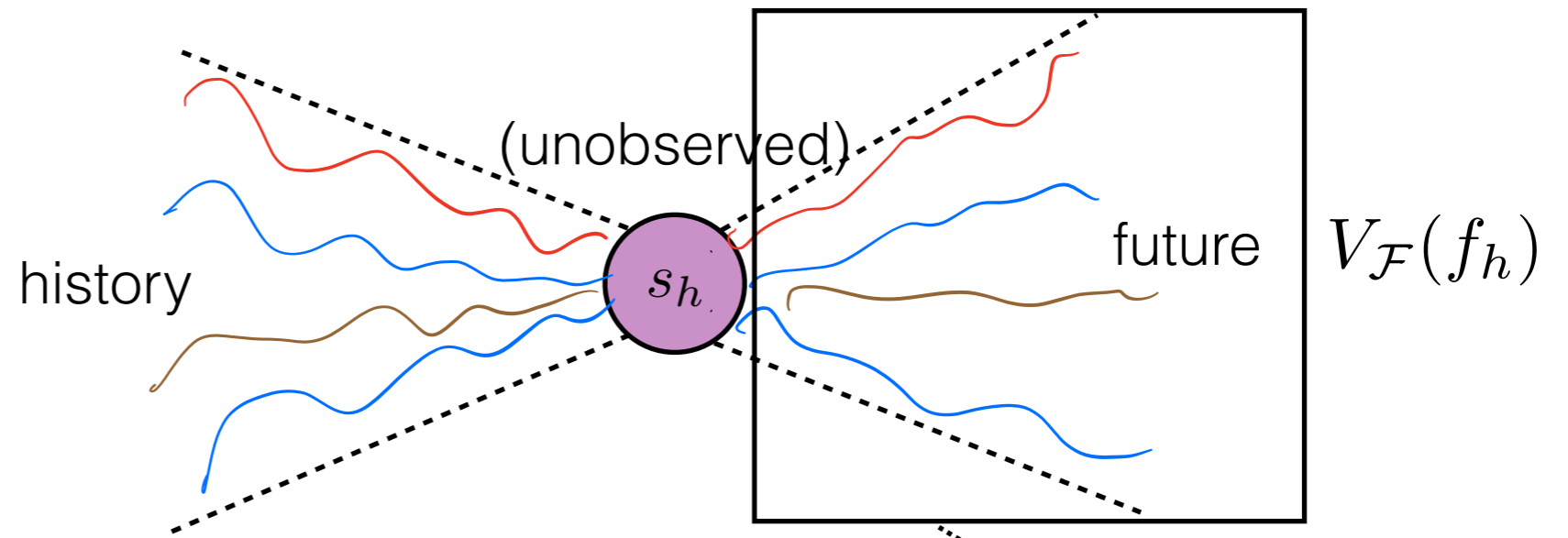
POMDP case



POMDP case



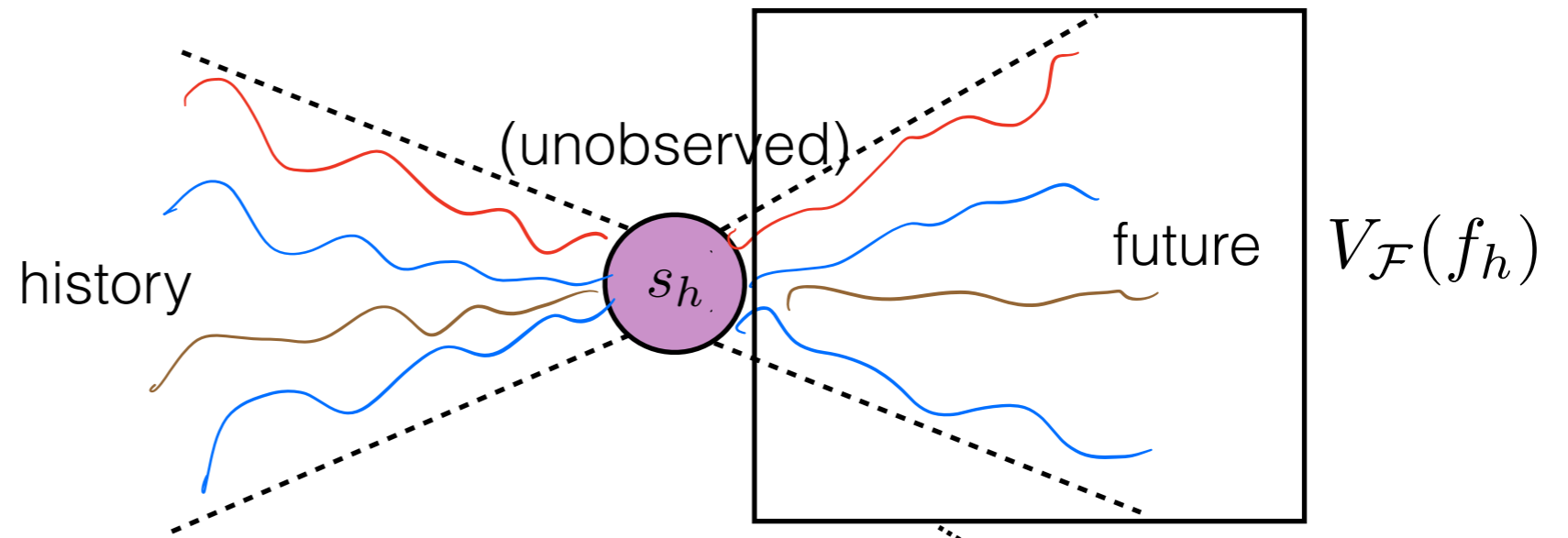
POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right]$$

POMDP case

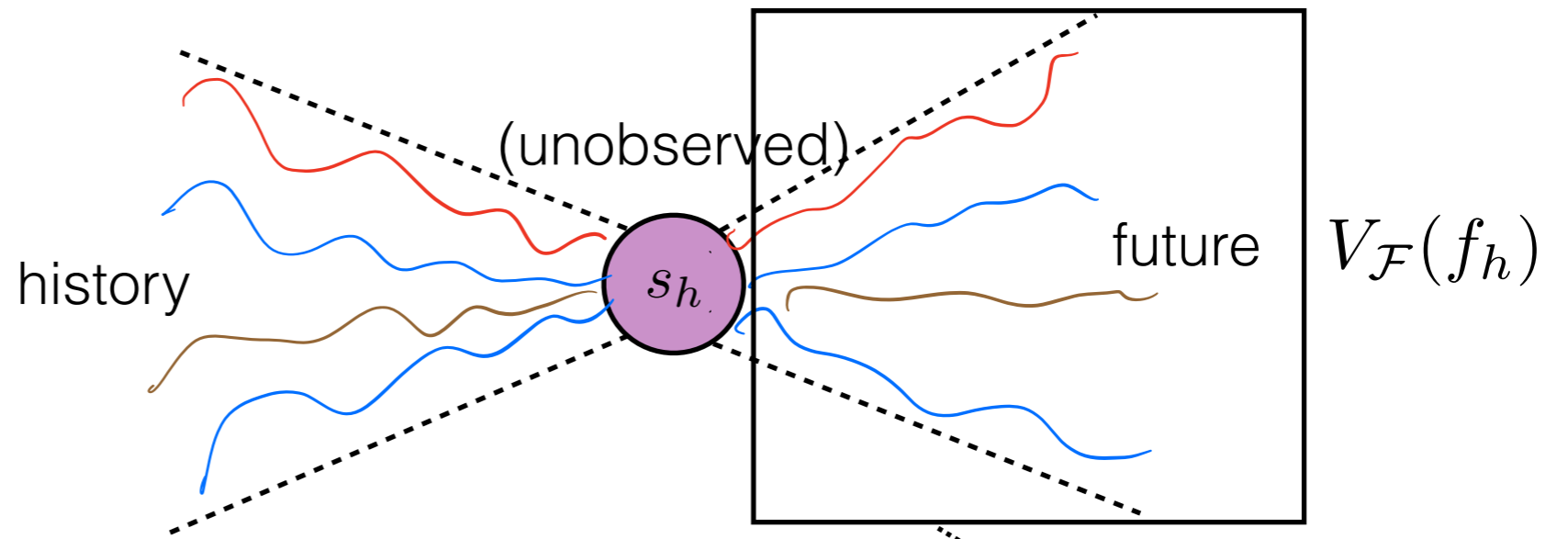


Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}}(s_h)}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right]$$

$$\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$$

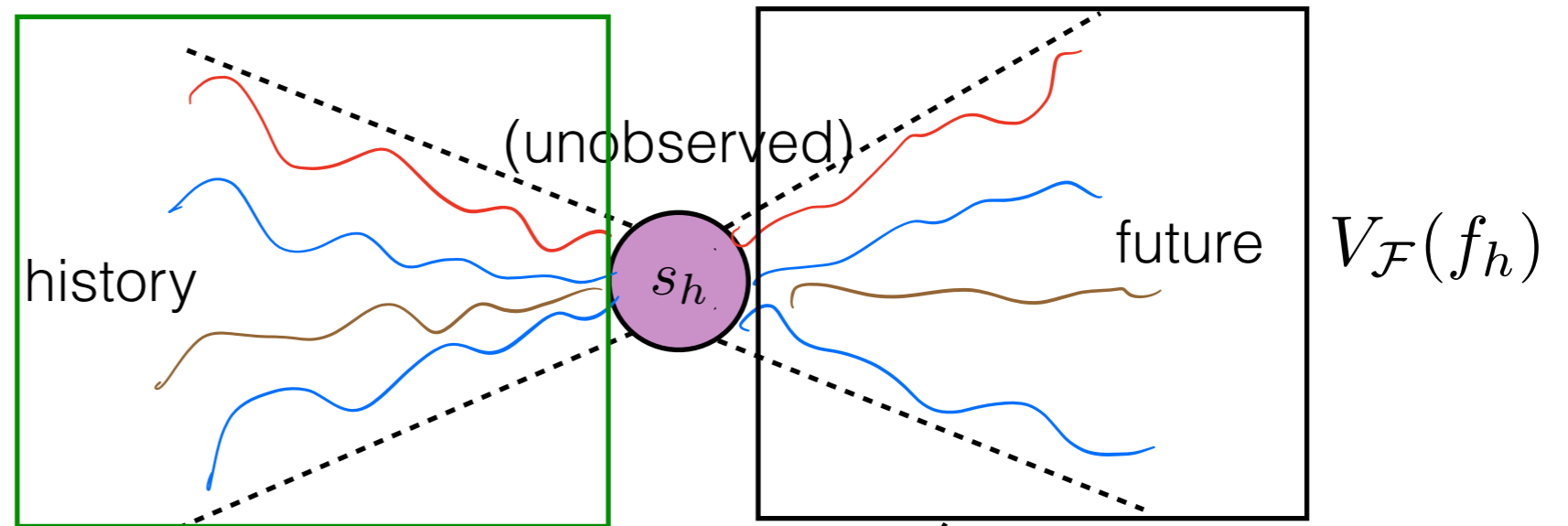
POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

POMDP case



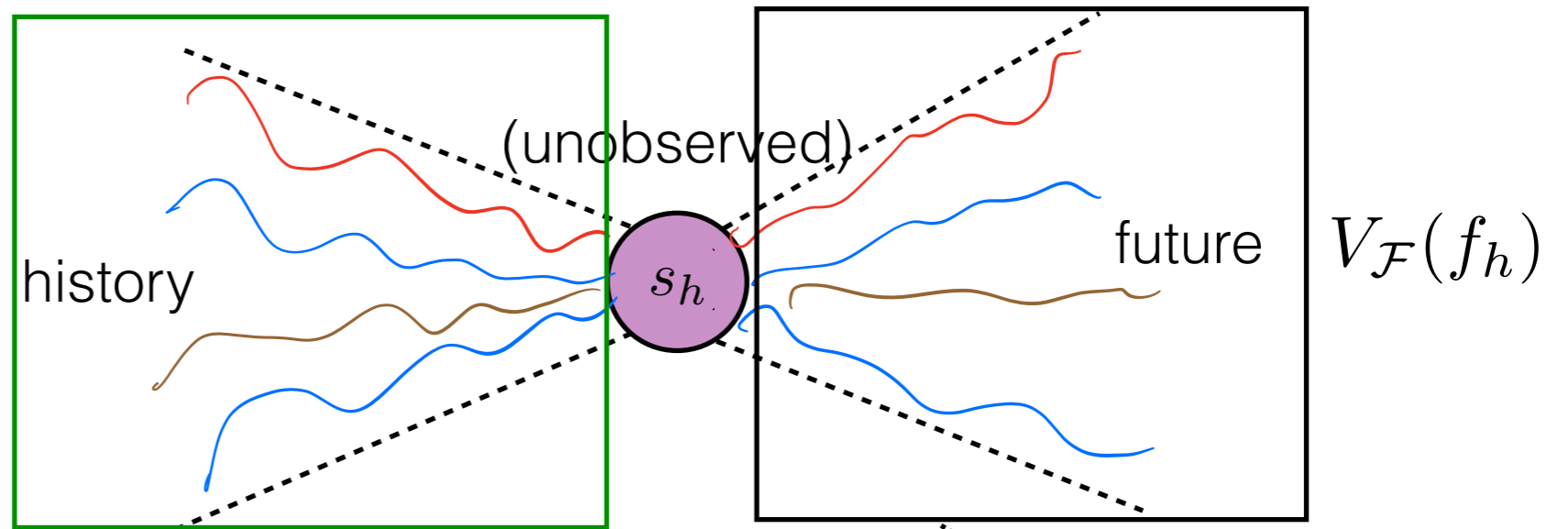
Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \tau_h \right]^2 \right]$$

POMDP case



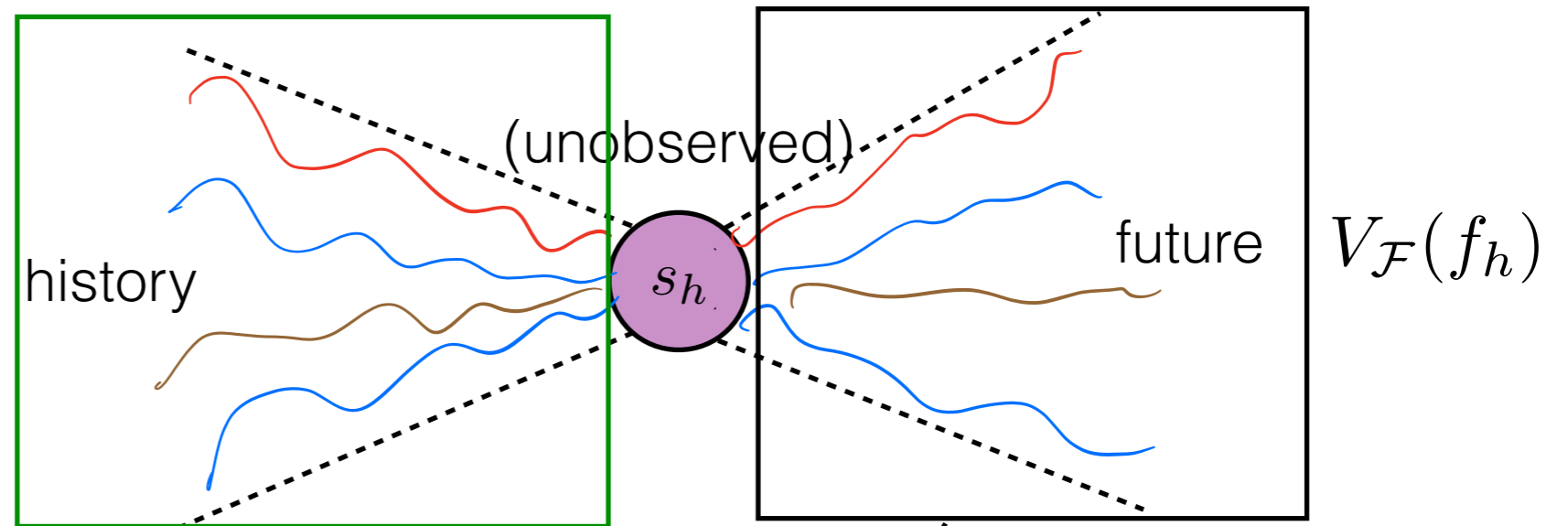
Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \tau_h \right]^2 \right] \quad \checkmark$$

POMDP case



Ideal loss:

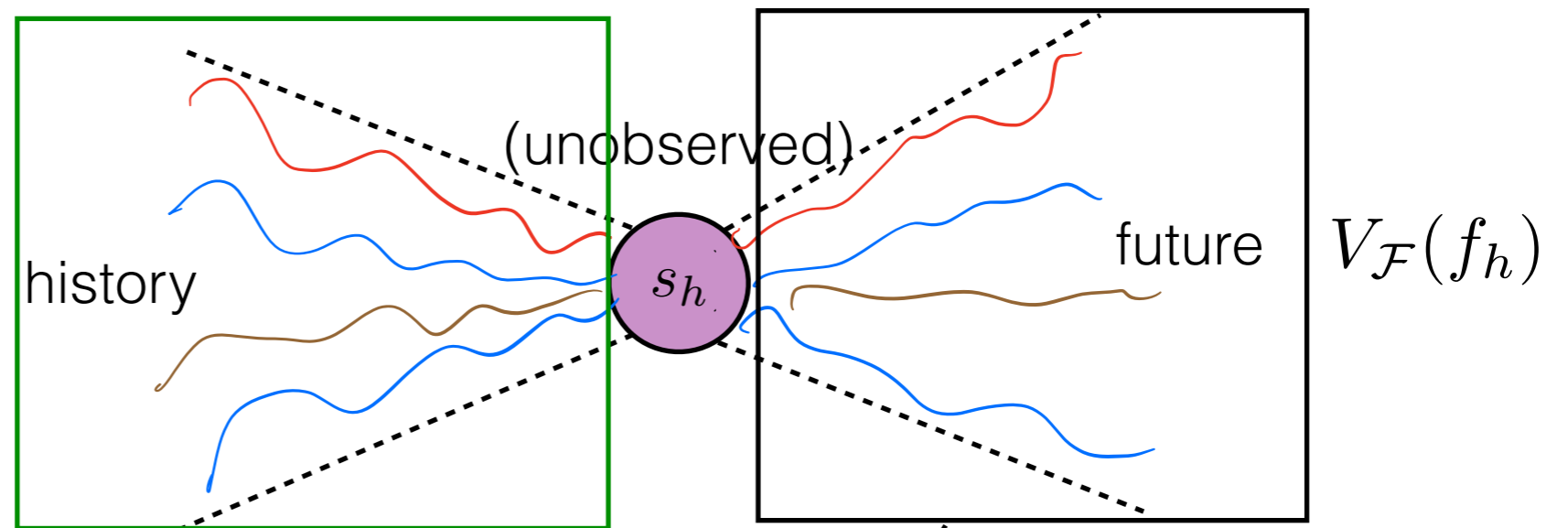
$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | s_h], \mathbf{b}(\cdot | \tau_h) \right\rangle^2 \right] \quad \parallel$$

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} [\Delta_h V_{\mathcal{F}} | \tau_h]^2 \right] \quad \checkmark$$

POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{\text{(unobserved)}} \right]^2 \right] \quad \mathbf{X}$$

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$

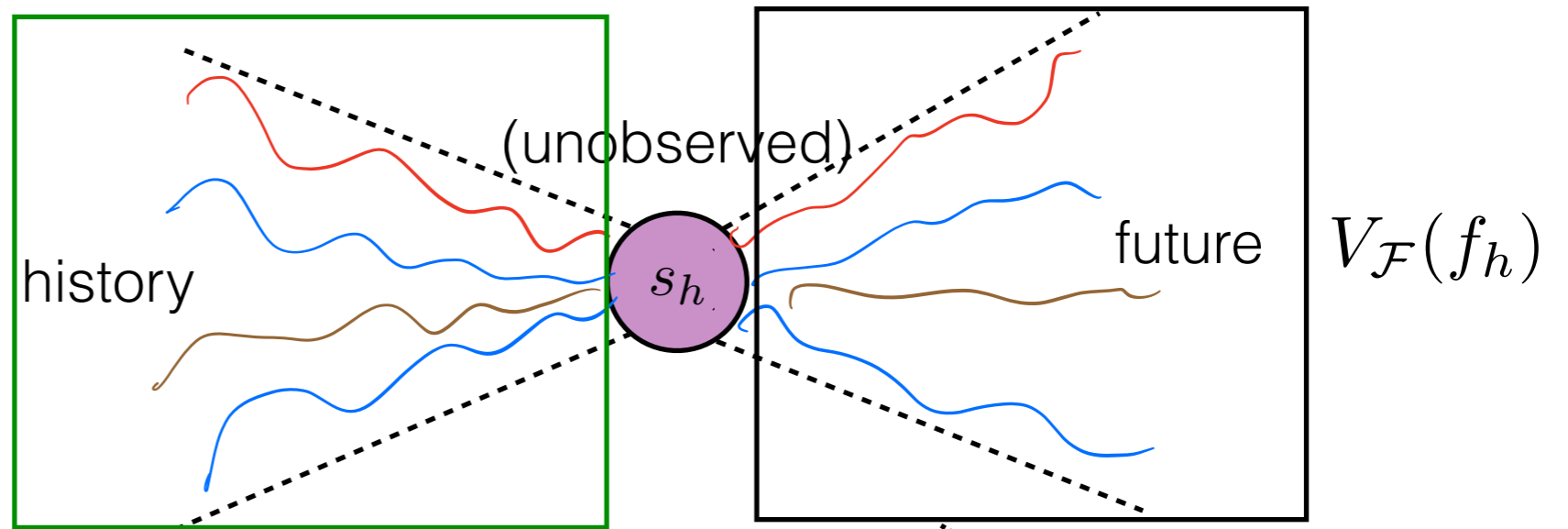
$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \tau_h) \right\rangle \right]^2$$

||

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \tau_h \right]^2 \right] \quad \checkmark$$

POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{\text{(unobserved)}} \right]^2 \right] \quad \mathbf{X}$$

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$

belief state $\Pr[s_h | \tau_h]$

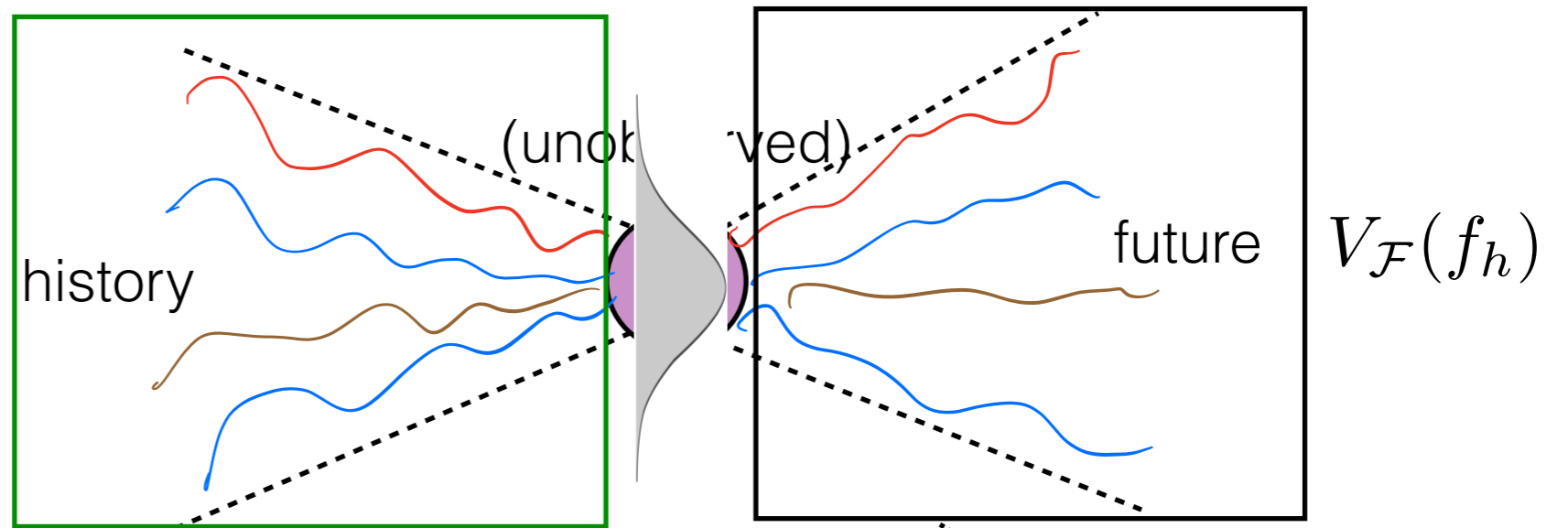
$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \tau_h) \right\rangle \right]^2$$

||

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \tau_h \right]^2 \right] \quad \checkmark$$

POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \tau_h) \right\rangle^2 \right]$$

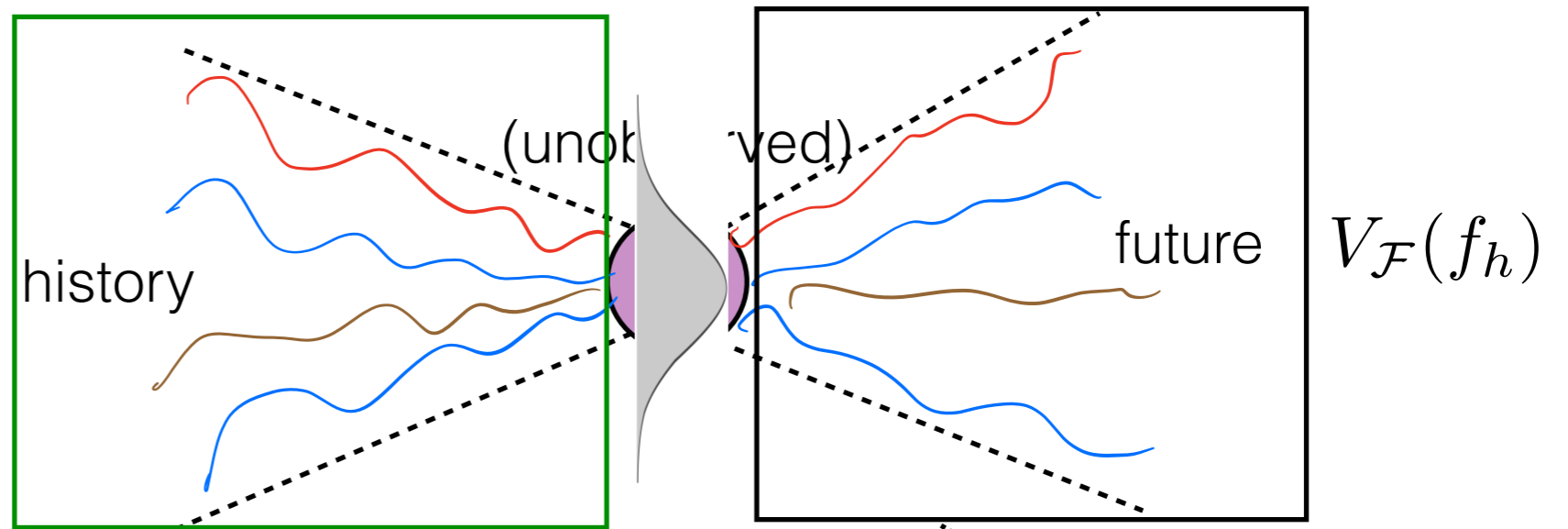
belief state $\Pr[s_h | \tau_h]$

||

Actual loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \tau_h \right]^2 \right] \quad \checkmark$$

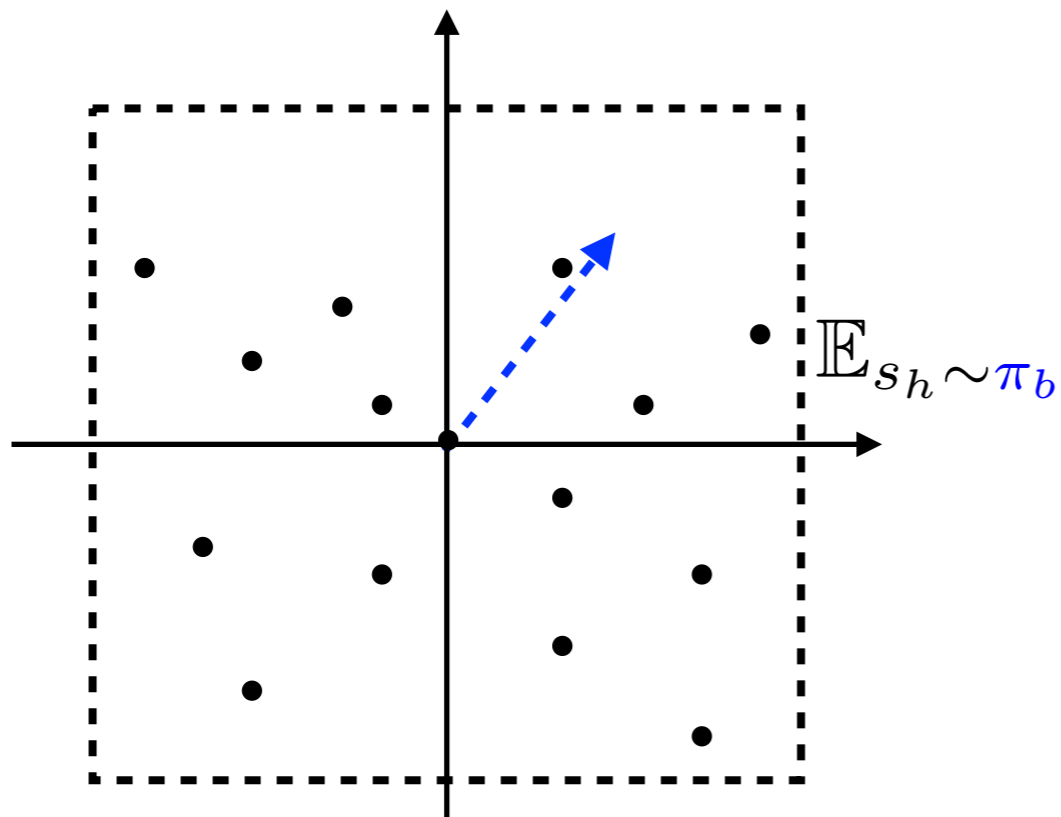
POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{\text{unknown}} \right]^2 \right] \quad \mathbf{X}$$

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$



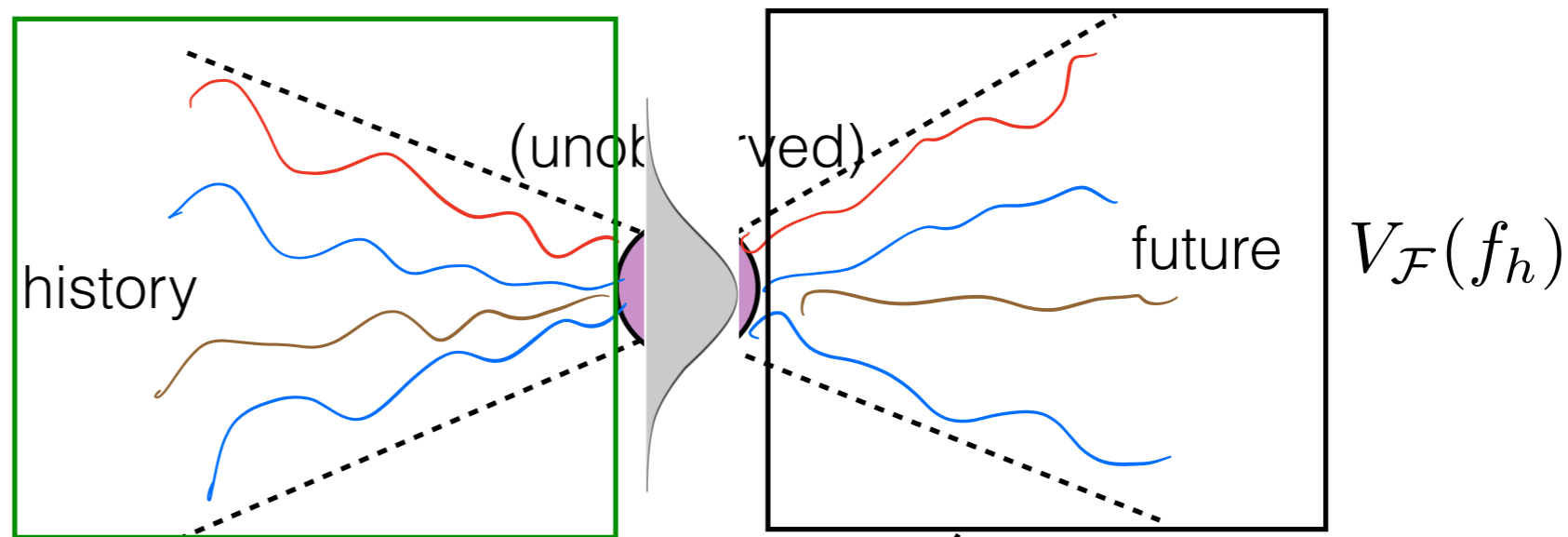
$S = 2$

belief state $\Pr[s_h | \mathcal{T}_h]$

$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \mathcal{T}_h) \right\rangle^2 \right]$$

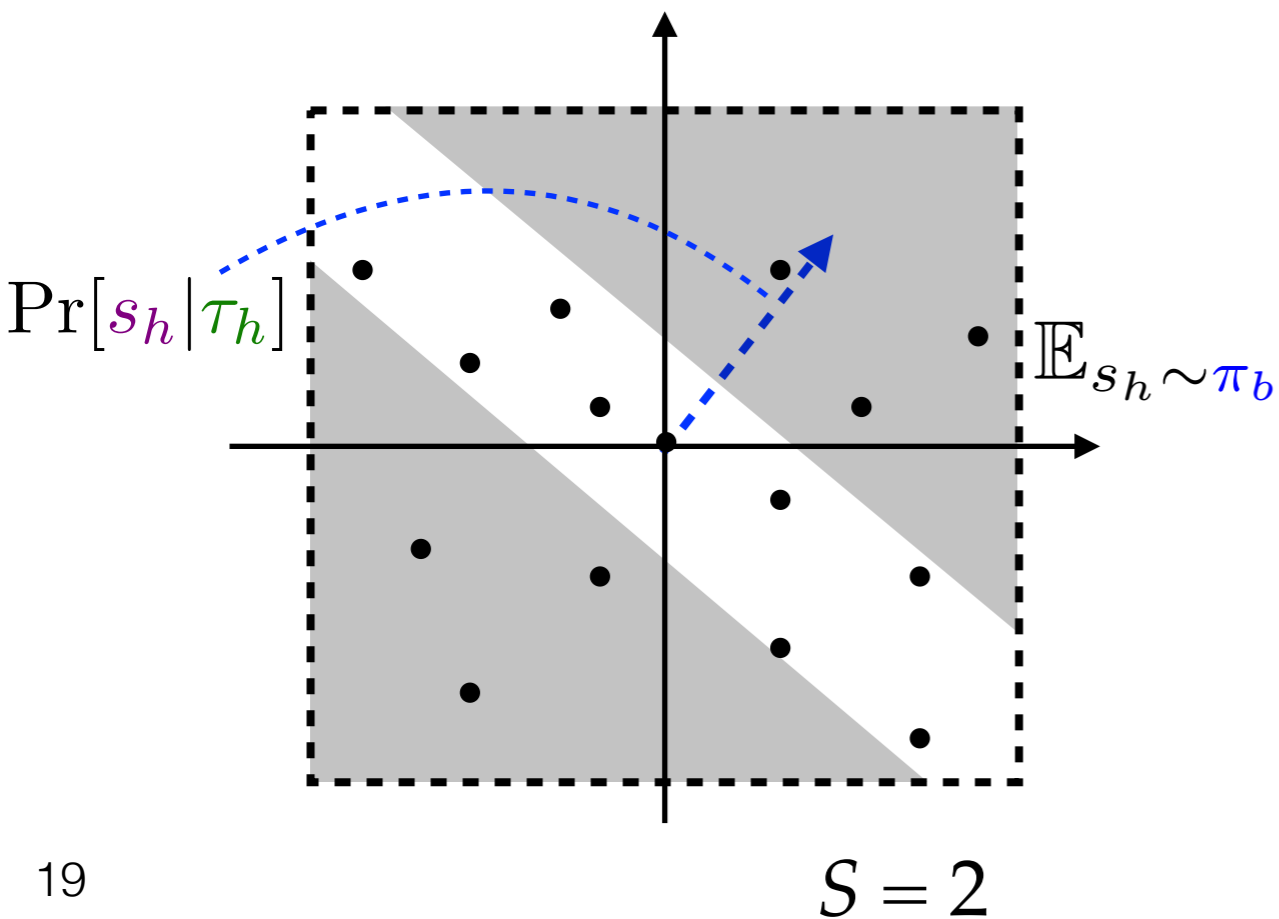
$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \mathcal{T}_h \right]^2 \right] \quad \checkmark$$

POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})} \right]^2 \right] \quad \mathbf{X}$$

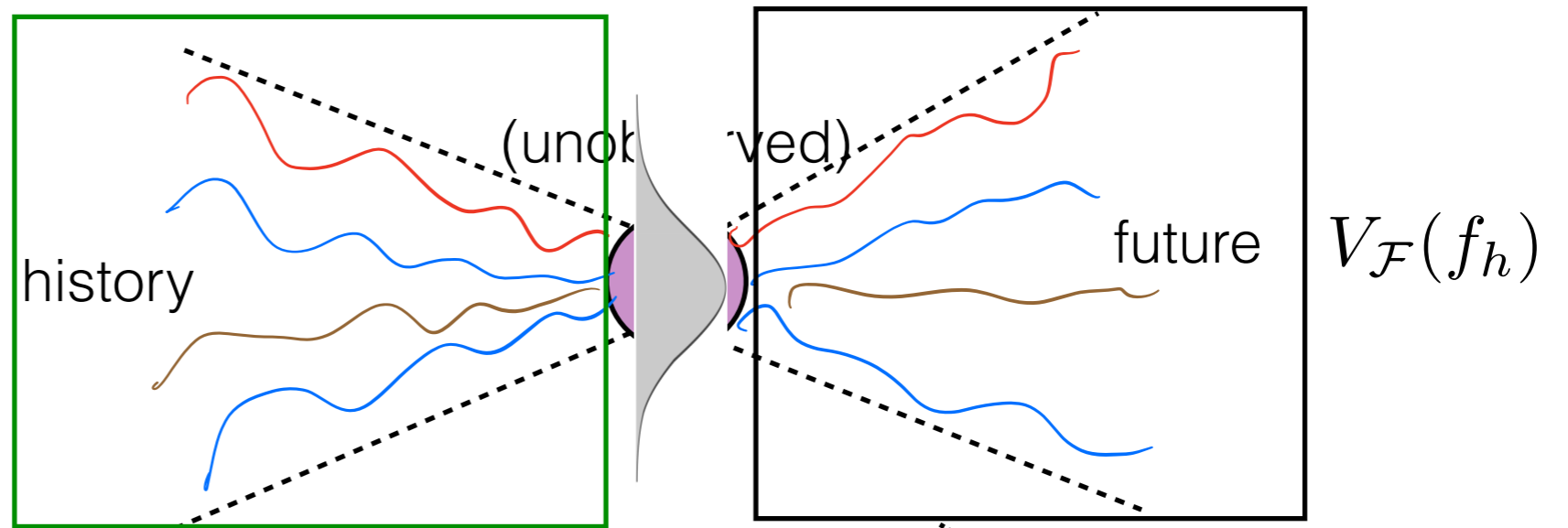


$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \mathcal{T}_h) \right\rangle^2 \right]$$

||

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \mathcal{T}_h \right]^2 \right] \quad \checkmark$$

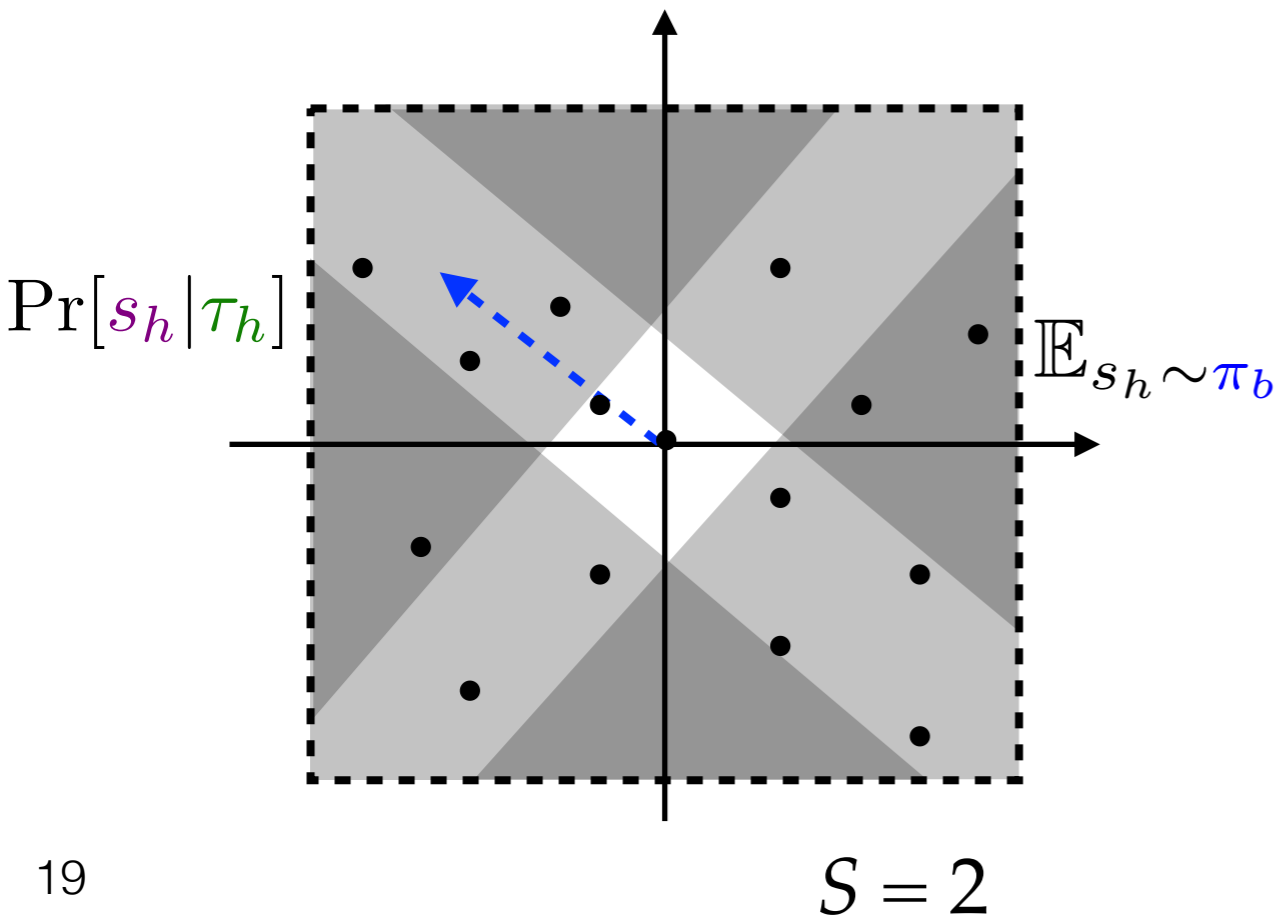
POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{\text{observed}} \right]^2 \right] \quad \mathbf{X}$$

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$

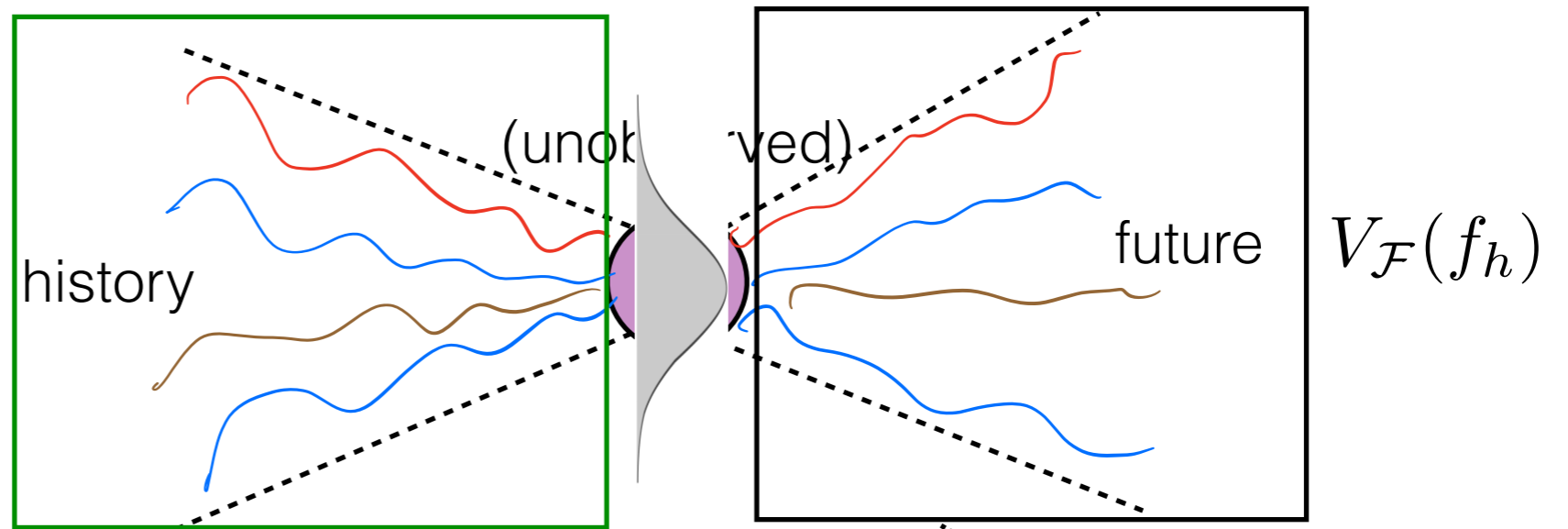


belief state $\Pr[s_h | \mathcal{T}_h]$

$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \mathcal{T}_h) \right\rangle^2 \right]$$

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \mathcal{T}_h \right]^2 \right] \quad \checkmark$$

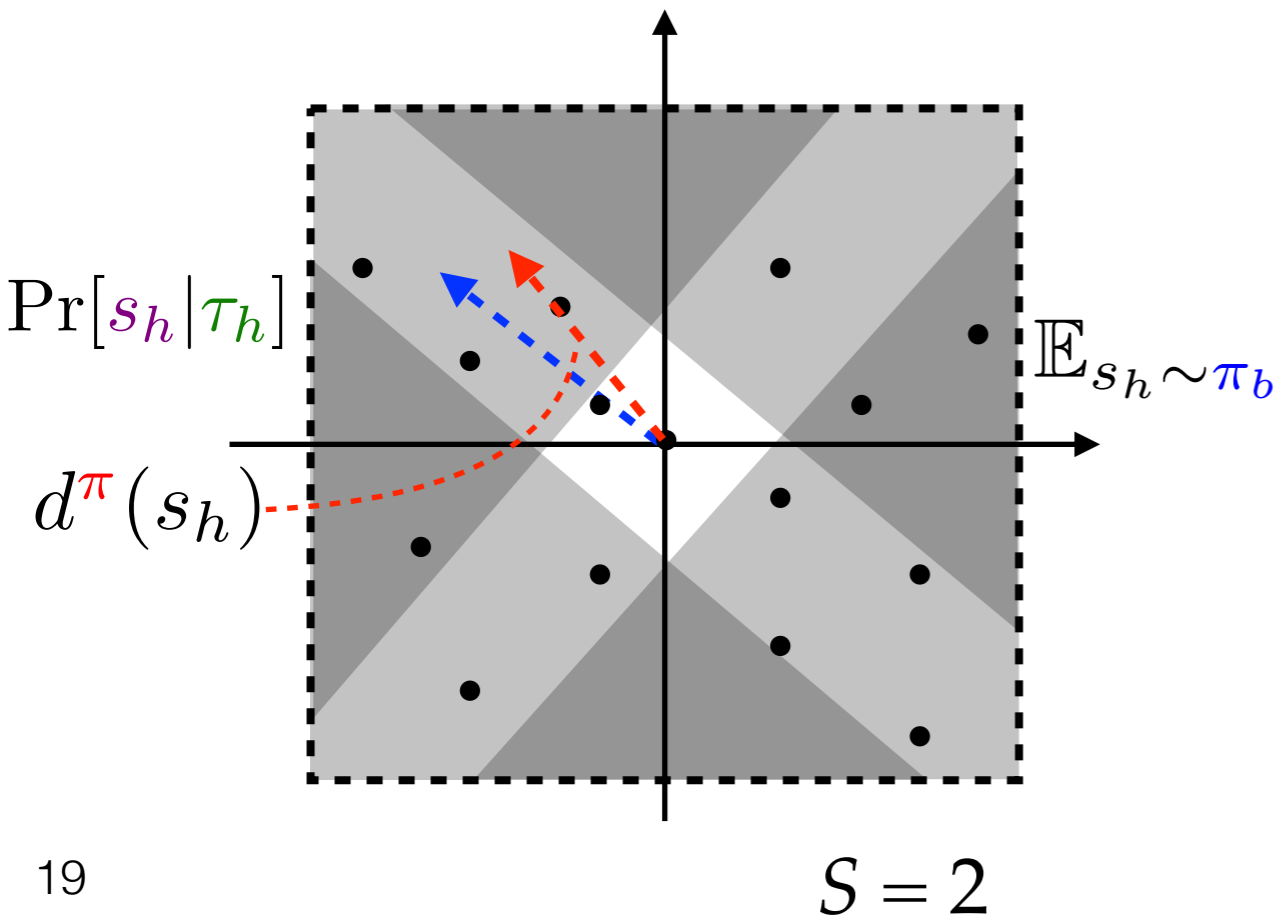
POMDP case



Ideal loss:

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\underbrace{\Delta_h V_{\mathcal{F}} | s_h}_{\text{unknown}} \right]^2 \right] \quad \mathbf{X}$$

$$V_{\mathcal{F}}(f_h) - r_h - V_{\mathcal{F}}(f_{h+1})$$



$$\mathbb{E}_{s_h \sim \pi_b} \left[\left\langle \mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | s_h \right], \mathbf{b}(\cdot | \mathcal{T}_h) \right\rangle^2 \right]$$

||

$$\mathbb{E}_{s_h \sim \pi_b} \left[\mathbb{E}_{\substack{a_h \sim \pi \\ a_{h+1:H} \sim \pi_b}} \left[\Delta_h V_{\mathcal{F}} | \mathcal{T}_h \right]^2 \right] \quad \checkmark$$

Guarantee under **belief coverage**

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- **Coverage** parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$
- Ranges & complexities of function classes (e.g., that of \mathcal{V})

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- **Coverage** parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$
- Ranges & complexities of function classes (e.g., that of \mathcal{V})

- Similar to $\mathbb{E}_{\pi} [\phi]^\top \mathbb{E}_{\pi_b} [\phi \phi^\top]^{-1} \mathbb{E}_{\pi} [\phi]$

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- **Coverage** parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$
- Ranges & complexities of function classes (e.g., that of \mathcal{V})

- Similar to $\mathbb{E}_{\pi} [\phi]^\top \mathbb{E}_{\pi_b} [\phi \phi^\top]^{-1} \mathbb{E}_{\pi} [\phi]$
- $\pi = \pi_b : C_{\mathcal{H}} = 1$

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- **Coverage** parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$
- Ranges & complexities of function classes (e.g., that of \mathcal{V})

- Similar to $\mathbb{E}_{\pi} [\phi]^\top \mathbb{E}_{\pi_b} [\phi \phi^\top]^{-1} \mathbb{E}_{\pi} [\phi]$
- $\pi = \pi_b : C_{\mathcal{H}} = 1$
- 1-hot $\mathbf{b}(\tau_h)$: $\mathbb{E}_{\pi_b} [(d_h^\pi / d_h^{\pi_b})^2]$

Guarantee under **belief coverage**

Theorem (Informal): Assume

$$(d_h^\pi)^\top \mathbb{E}_{\tau_h \sim \pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi \leq C_{\mathcal{H}}$$

and standard representation assumptions (realizability & Bellman-completeness), the sample complexity of OPE is poly in

- **Coverage** parameters: $C_{\mathcal{H}}$ and $C_{\mathcal{A}} := \max_{s_h, a_h} \frac{\pi(a_h|s_h)}{\pi_b(a_h|s_h)}$
- Ranges & complexities of function classes (e.g., that of \mathcal{V})

- Similar to $\mathbb{E}_{\pi} [\phi]^\top \mathbb{E}_{\pi_b} [\phi \phi^\top]^{-1} \mathbb{E}_{\pi} [\phi]$
- $\pi = \pi_b : C_{\mathcal{H}} = 1$
- 1-hot $\mathbf{b}(\tau_h)$: $\mathbb{E}_{\pi_b} [(d_h^\pi / d_h^{\pi_b})^2]$
- $\|d_h^\pi / d_h^{\pi_b}\|_\infty \Rightarrow \|\mathbb{E}_{\pi_b} [\mathbf{b}(\tau_h) \mathbf{b}(\tau_h)^\top]^{-1} d_h^\pi\|_\infty$

Does FDVF exist...?

- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$

Does FDVF exist...?

- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
- Equation has no solution (!) if
 - two latent states have the same future distribution under π_b
 - ... but they have different values in $V_{\mathcal{S}}^{\pi}(s_h)$

Does FDVF exist...?

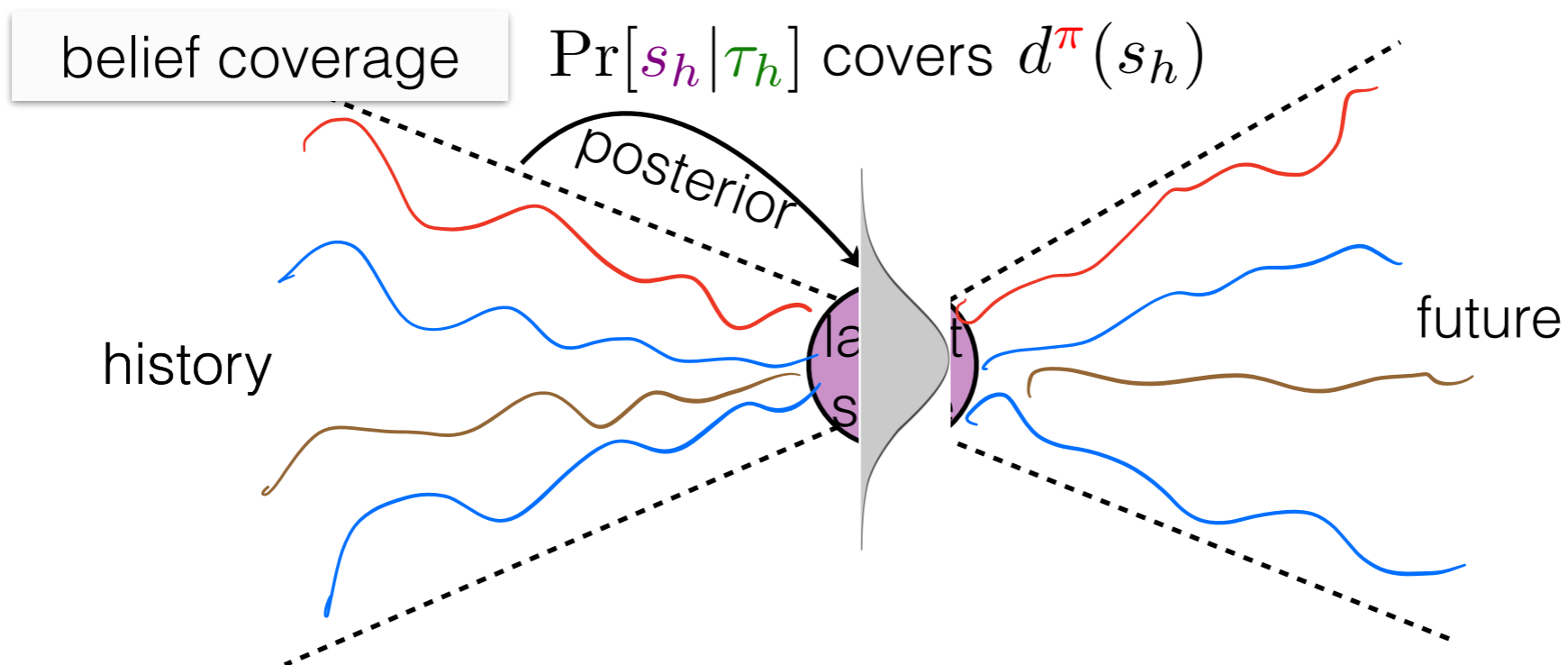
- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
- Equation has no solution (!) if
 - two latent states have the same future distribution under π_b
 - ... but they have different values in $V_{\mathcal{S}}^{\pi}(s_h)$
- New condition: **outcome coverage**
 - Let $\Sigma_{\mathcal{F}}$ be the confusion matrix of predicting s_h from future

Does FDVF exist...?

- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
- Equation has no solution (!) if
 - two latent states have the same future distribution under π_b
 - ... but they have different values in $V_{\mathcal{S}}^{\pi}(s_h)$
- New condition: **outcome coverage**
 - Let $\Sigma_{\mathcal{F}}$ be the confusion matrix of predicting s_h from future
 - Finite & bounded $\|\Sigma_{\mathcal{F}}^{-1} V_{\mathcal{S}}^{\pi}\|_{\infty}$

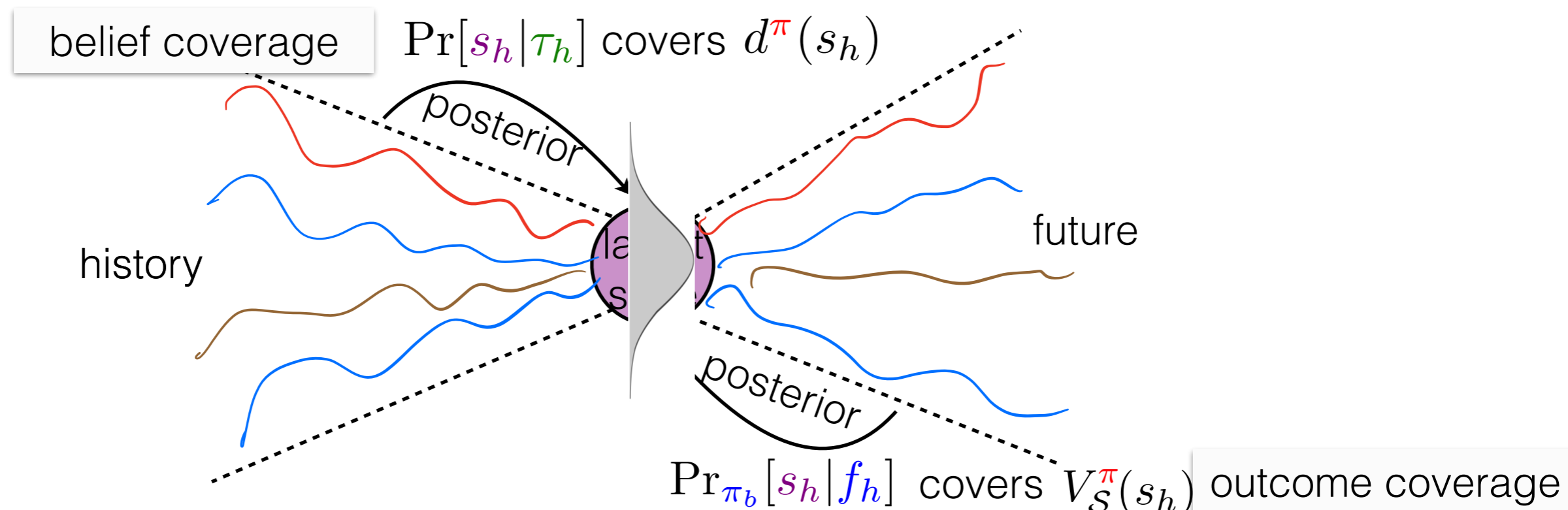
Does FDVF exist...?

- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
- Equation has no solution (!) if
 - two latent states have the same future distribution under π_b
 - ... but they have different values in $V_{\mathcal{S}}^{\pi}(s_h)$
- New condition: **outcome coverage**
 - Let $\Sigma_{\mathcal{F}}$ be the confusion matrix of predicting s_h from future
 - Finite & bounded $\|\Sigma_{\mathcal{F}}^{-1} V_{\mathcal{S}}^{\pi}\|_{\infty}$



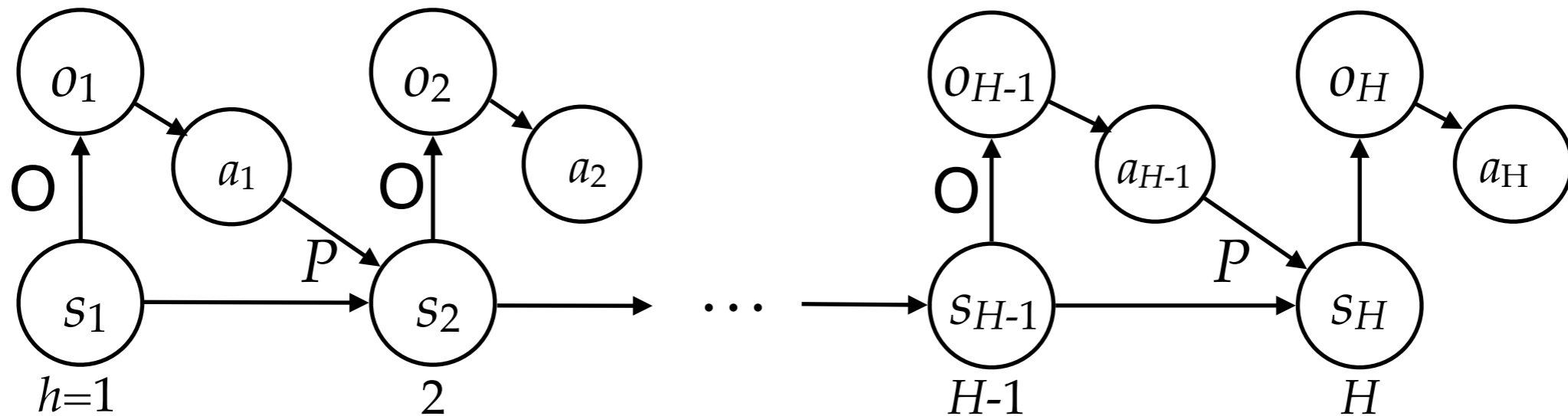
Does FDVF exist...?

- $\mathbb{E}_{\pi_b} [V_{\mathcal{F}}^{\pi}(f_h) | s_h] = V_{\mathcal{S}}^{\pi}(s_h)$
- Equation has no solution (!) if
 - two latent states have the same future distribution under π_b
 - ... but they have different values in $V_{\mathcal{S}}^{\pi}(s_h)$
- New condition: **outcome coverage**
 - Let $\Sigma_{\mathcal{F}}$ be the confusion matrix of predicting s_h from future
 - Finite & bounded $\|\Sigma_{\mathcal{F}}^{-1} V_{\mathcal{S}}^{\pi}\|_{\infty}$



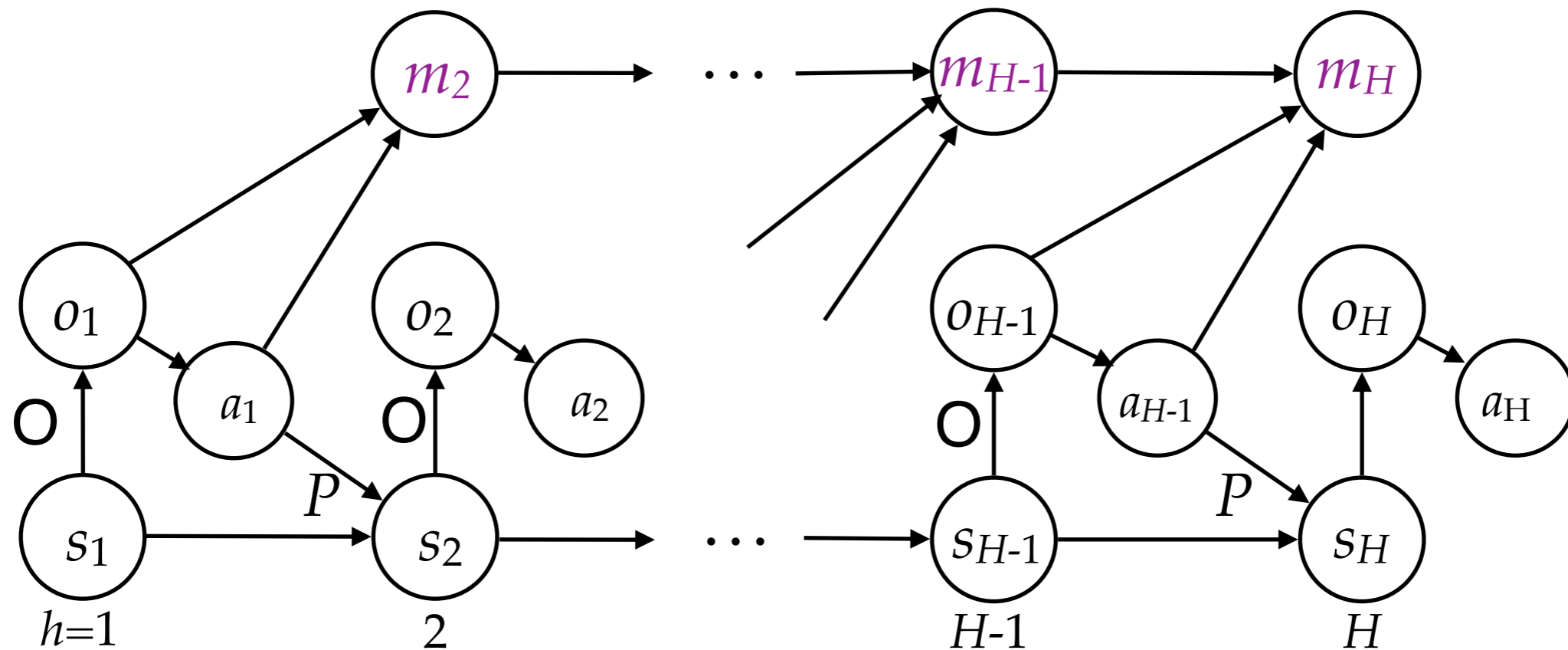
Reduction from FSM to Memoryless Policies

- Recurrent memory: $m_h = \text{update}(m_{h-1}, o_{h-1}, a_{h-1})$



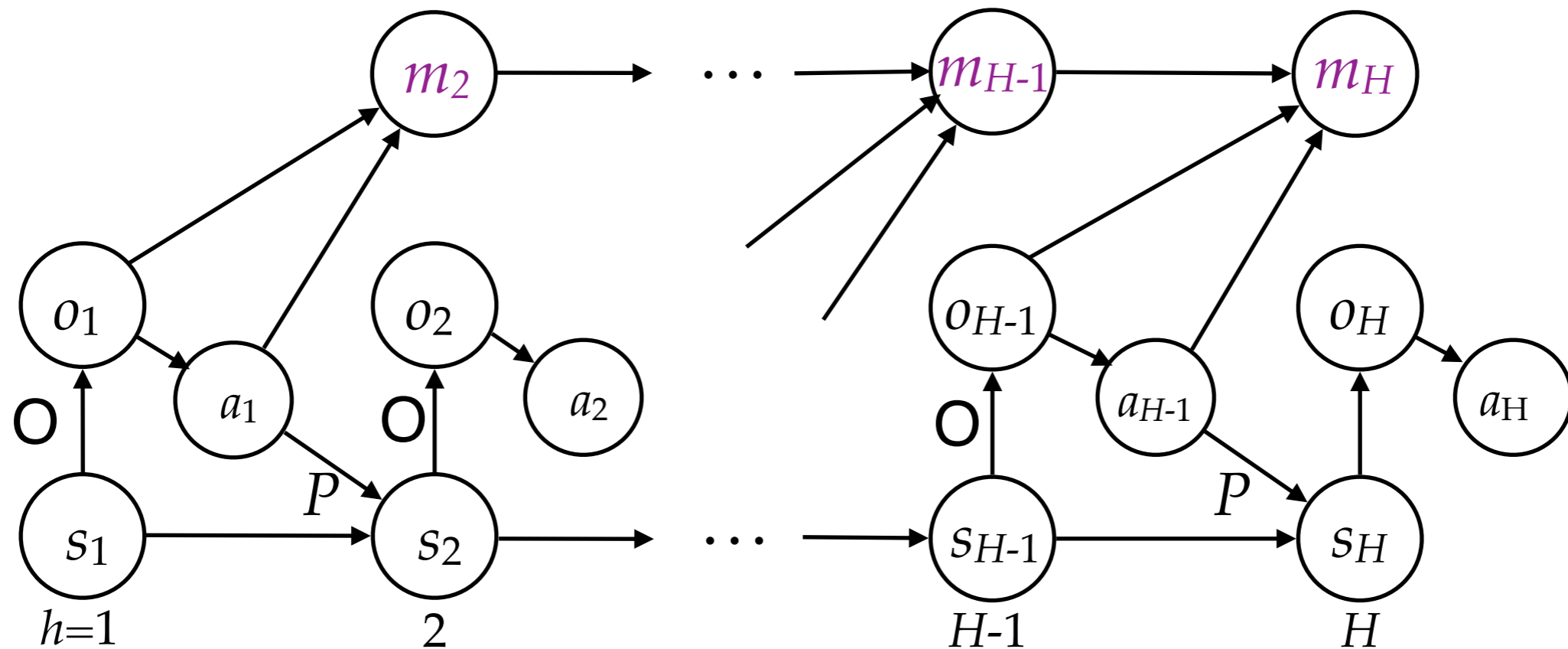
Reduction from FSM to Memoryless Policies

- Recurrent memory: $m_h = \text{update}(m_{h-1}, o_{h-1}, a_{h-1})$



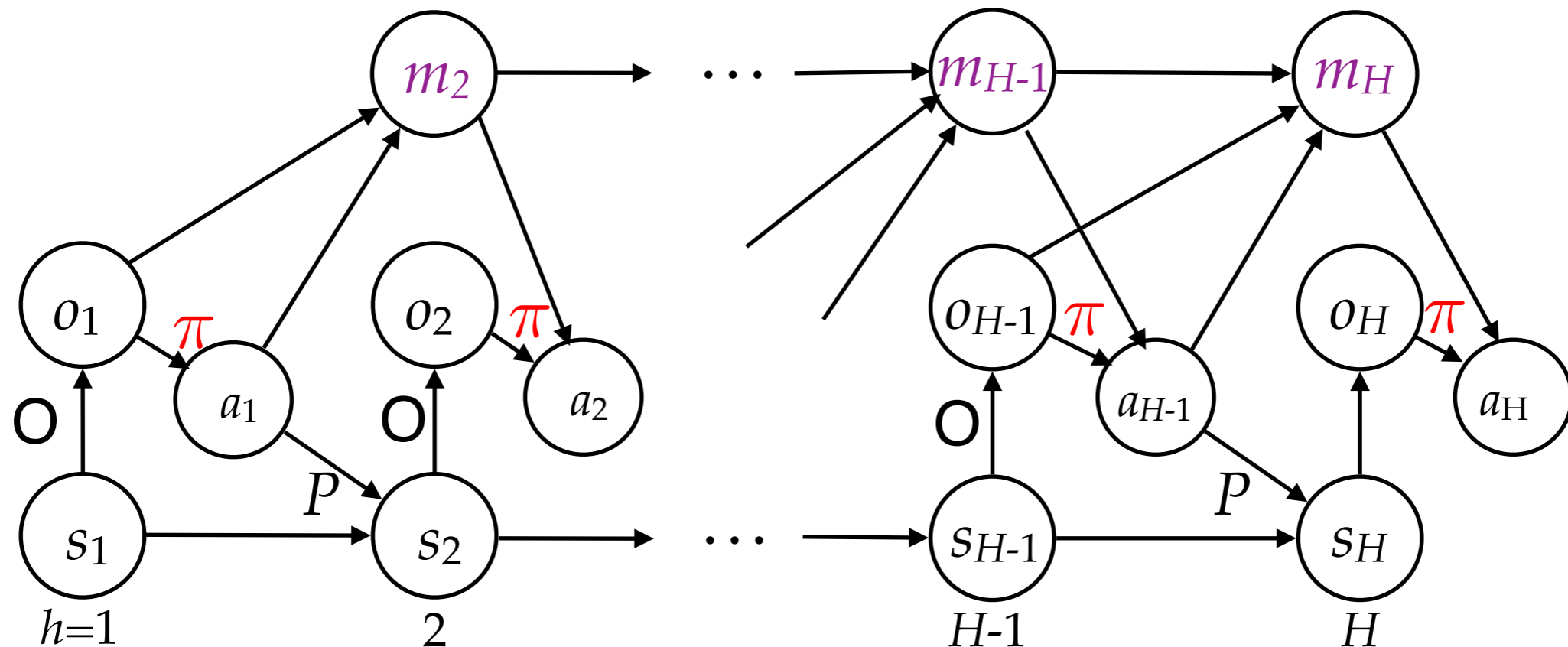
Reduction from FSM to Memoryless Policies

- Recurrent memory: $m_h = \text{update}(m_{h-1}, o_{h-1}, a_{h-1})$
 - Subsume $m_h = \tau_h$ as a special case



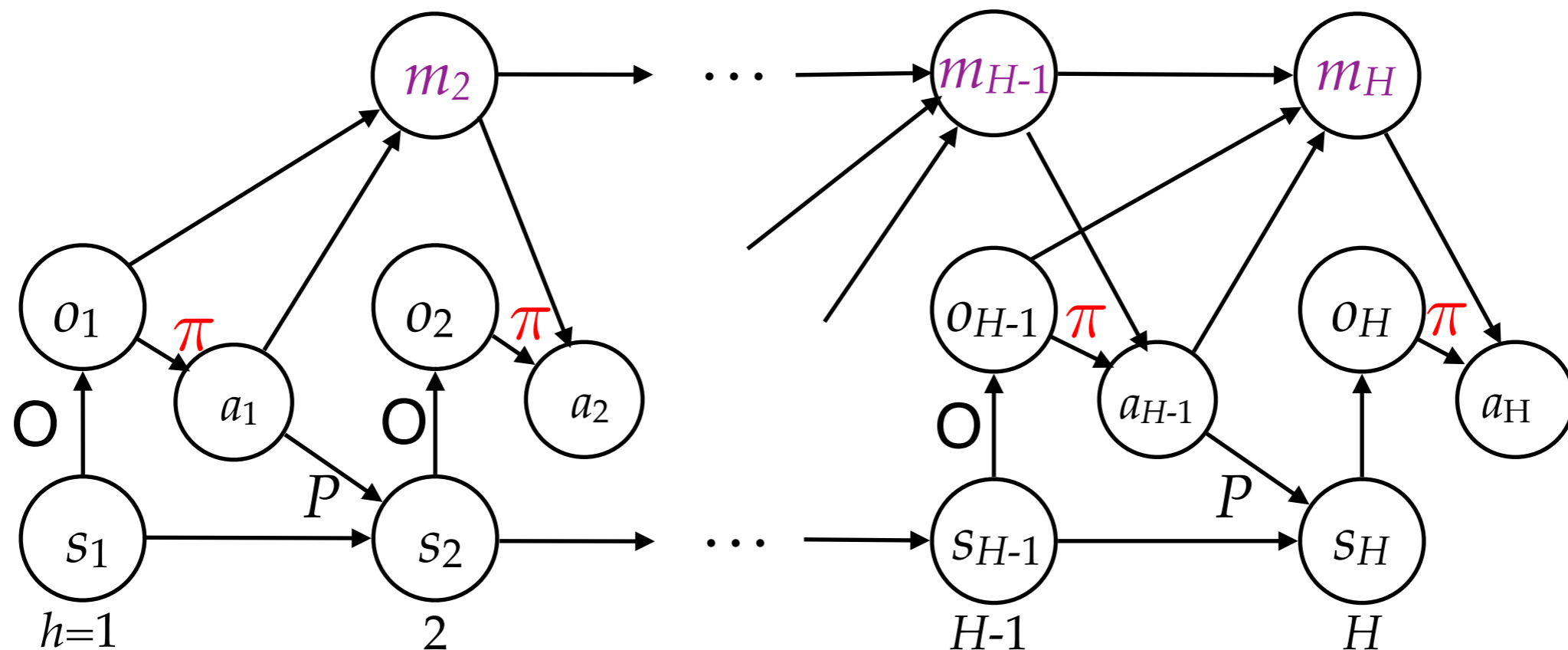
Reduction from FSM to Memoryless Policies

- Recurrent memory: $m_h = \text{update}(m_{h-1}, o_{h-1}, a_{h-1})$
 - Subsume $m_h = \tau_h$ as a special case
- FSM policy: $\pi(a_h | o_h, m_h)$



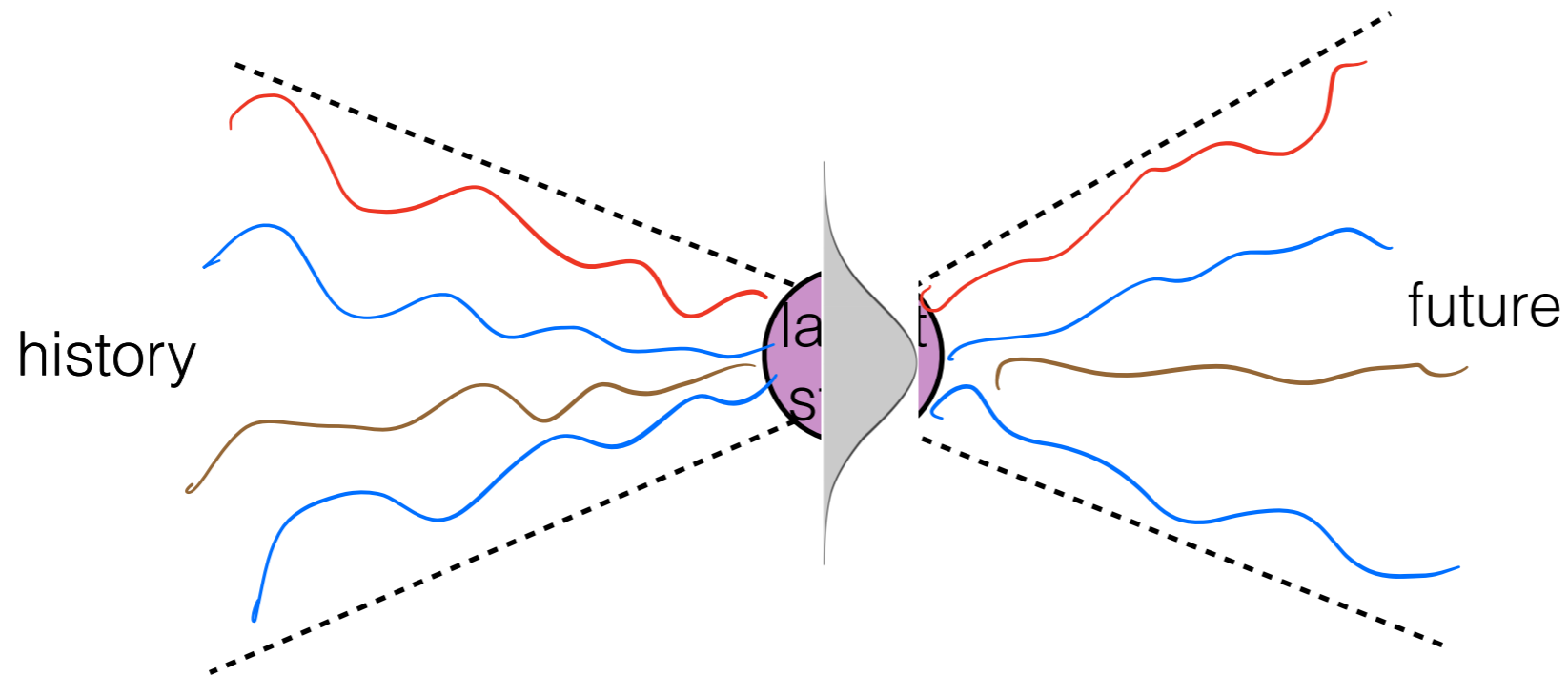
Reduction from FSM to Memoryless Policies

- Recurrent memory: $\mathcal{m}_h = \text{update}(\mathcal{m}_{h-1}, o_{h-1}, a_{h-1})$
 - Subsume $\mathcal{m}_h = \tau_h$ as a special case
- FSM policy: $\pi(a_h | o_h, \mathcal{m}_h)$
- Reduction: define an augmented POMDP
 - latent state $\tilde{s}_h = (s_h, \mathcal{m}_h)$
 - observation $\tilde{o}_h = (o_h, \mathcal{m}_h)$



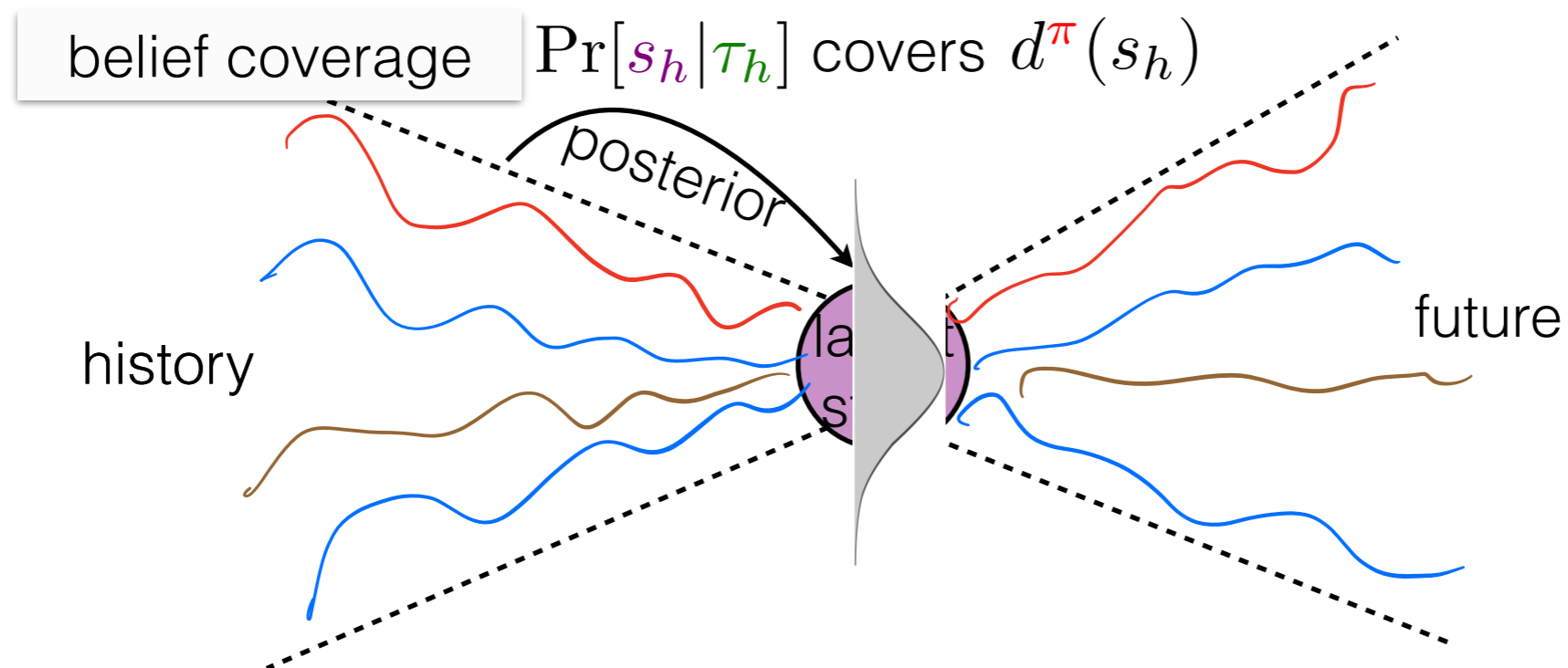
Conclusion

- **Problem:** OPE in POMDPs
- New **framework:** future-dependent value function $V_{\mathcal{F}}^{\pi}$



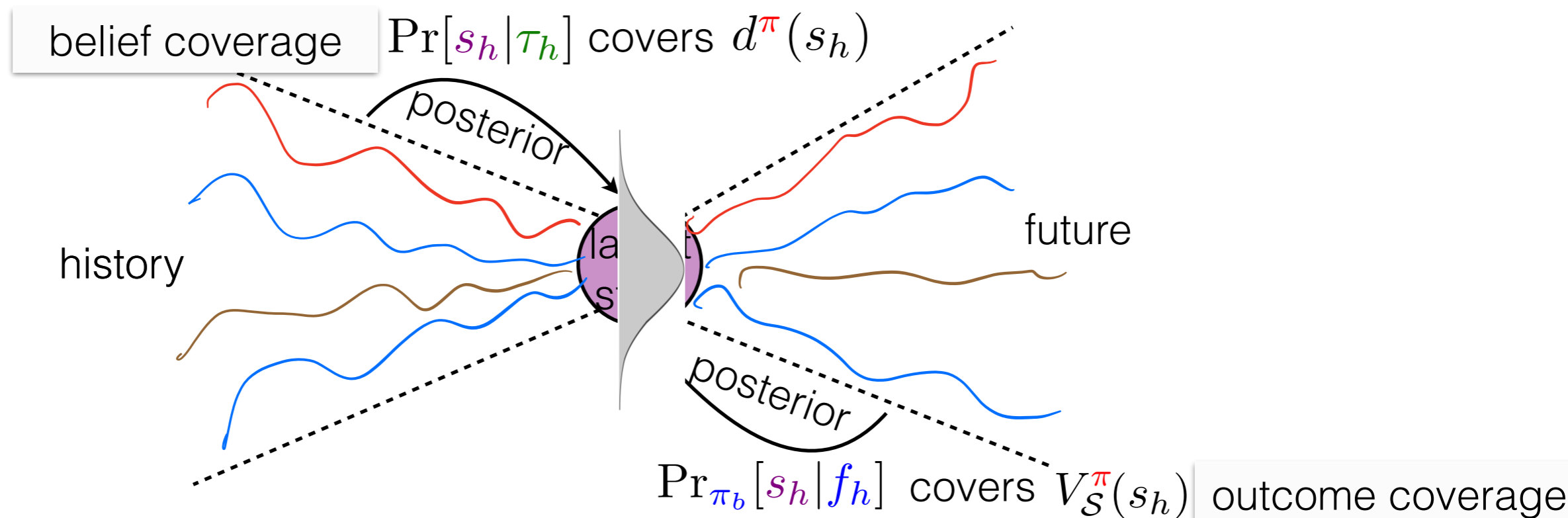
Conclusion

- **Problem:** OPE in POMDPs
- New **framework:** future-dependent value function $V_{\mathcal{F}}^{\pi}$
- New **assumptions:**
 - **Belief coverage** \Rightarrow error transfer



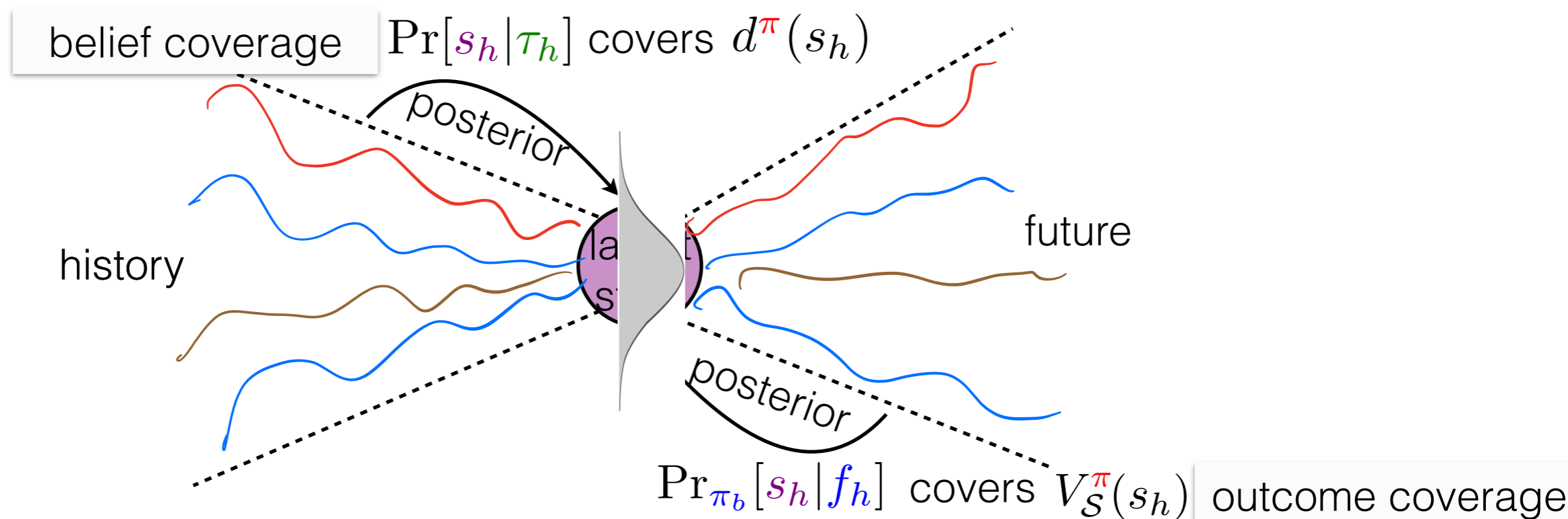
Conclusion

- **Problem:** OPE in POMDPs
- New **framework:** future-dependent value function $V_{\mathcal{F}}^{\pi}$
- New **assumptions:**
 - **Belief coverage** \Rightarrow error transfer
 - **Outcome coverage** \Rightarrow bounded $V_{\mathcal{F}}^{\pi}$



Conclusion

- **Problem:** OPE in POMDPs
- New **framework:** future-dependent value function $V_{\mathcal{F}}^{\pi}$
- New **assumptions:**
 - **Belief coverage** \Rightarrow error transfer
 - **Outcome coverage** \Rightarrow bounded $V_{\mathcal{F}}^{\pi}$
- Open **question:** beyond memoryless & FSM policies



Conclusion



Masatoshi Uehara



Yuheng Zhang

- **Problem:** OPE in POMDPs
- New **framework:** future-dependent value function $V_{\mathcal{F}}^{\pi}$
- New **assumptions:**
 - **Belief coverage** \Rightarrow error transfer
 - **Outcome coverage** \Rightarrow bounded $V_{\mathcal{F}}^{\pi}$
- Open **question:** beyond memoryless & FSM policies

