

# Model-based RL in Contextual Decision Processes: PAC Bounds and Exponential Improvements over Model-free Approaches

**Wen Sun**

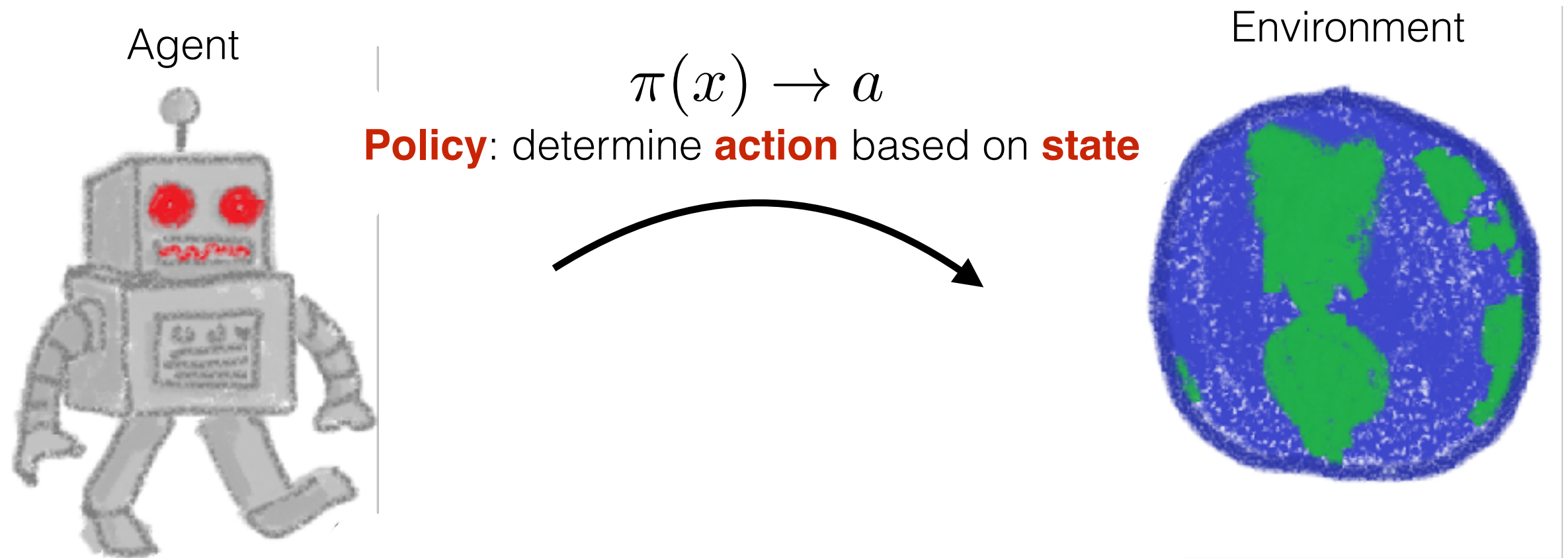
CMU -> MSR NYC

Joint work with Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford



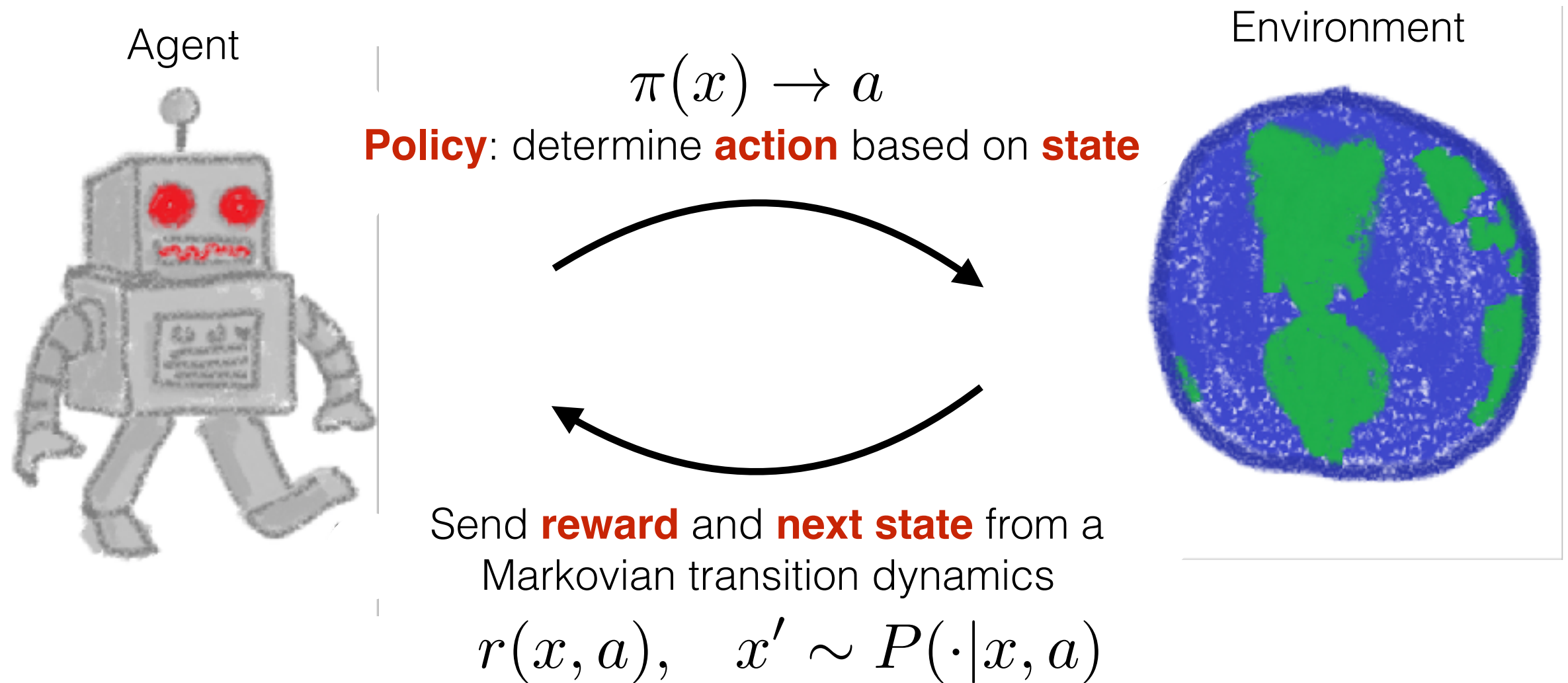
# Reinforcement Learning

## Markov Decision Process



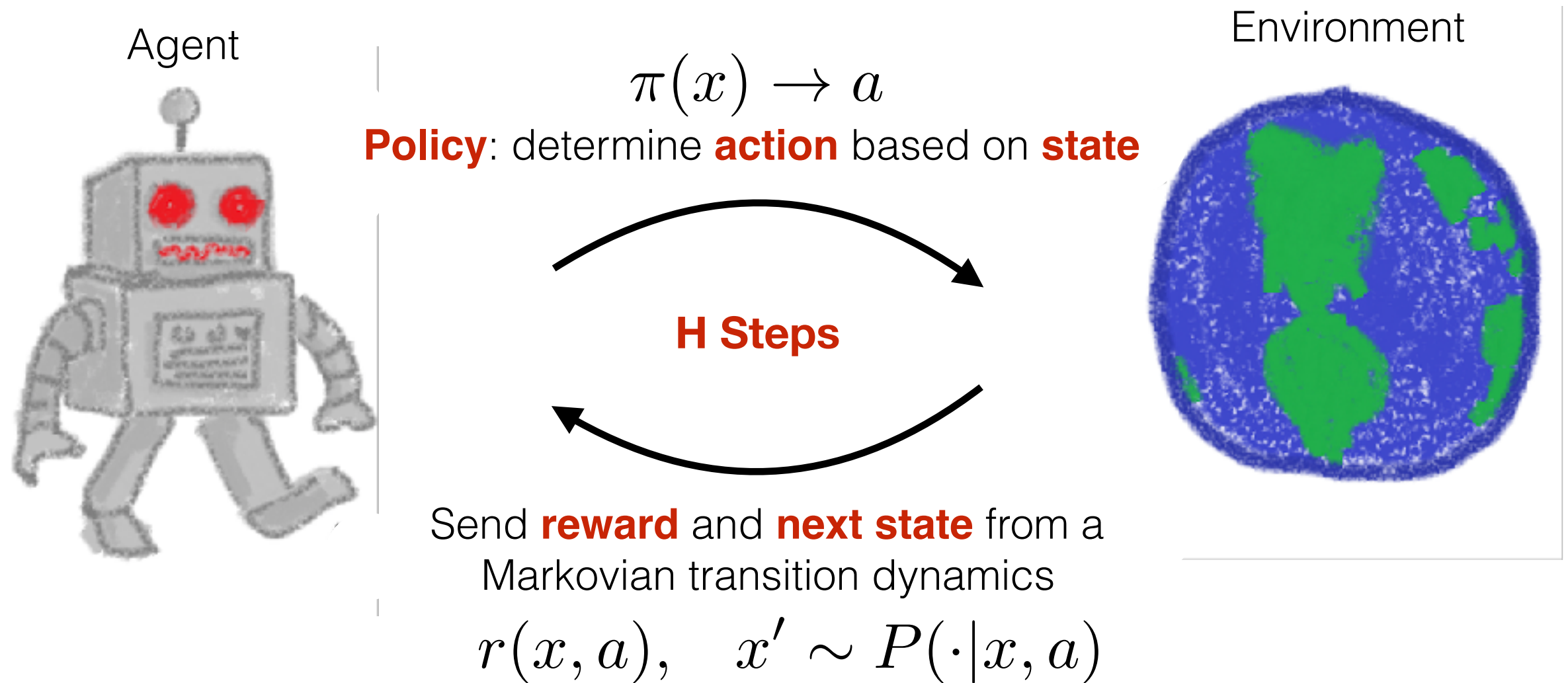
# Reinforcement Learning

## Markov Decision Process



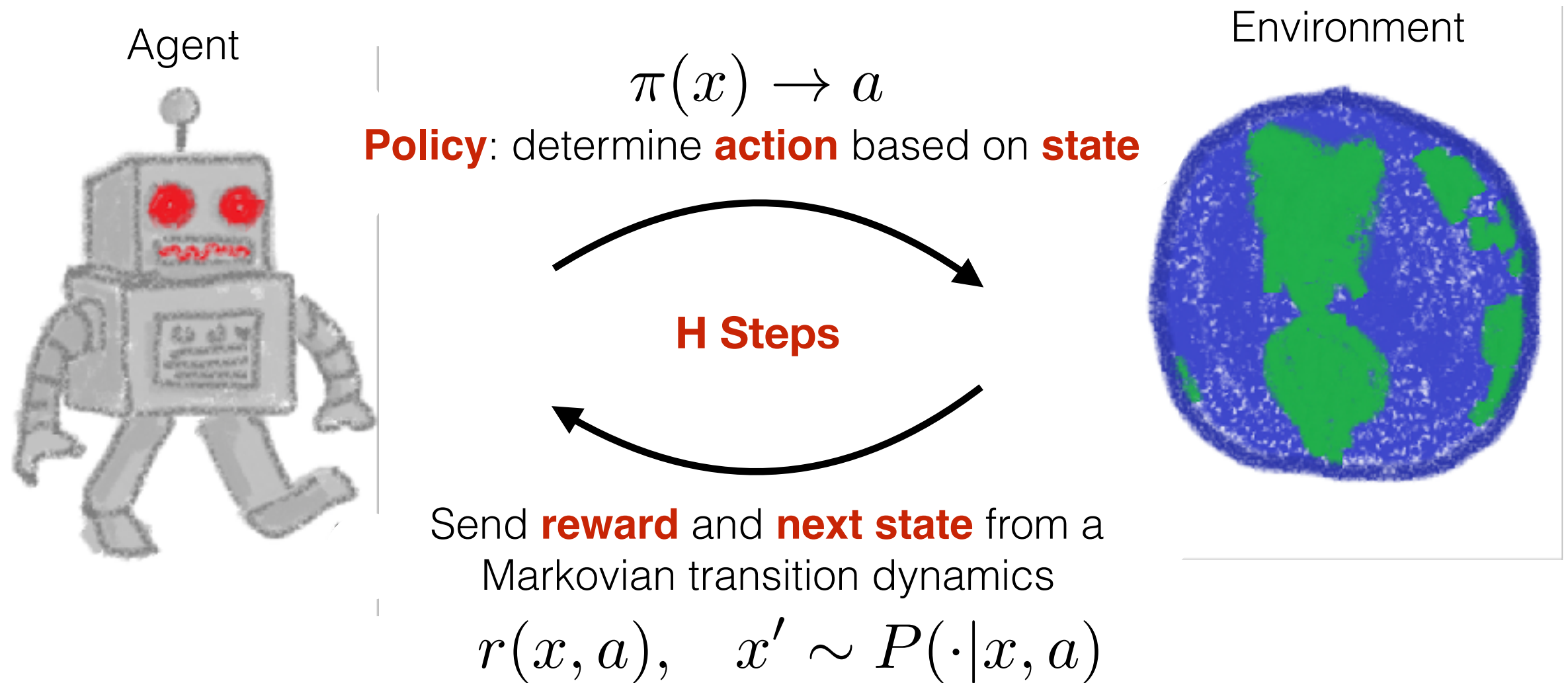
# Reinforcement Learning

## Markov Decision Process



# Reinforcement Learning

## Markov Decision Process

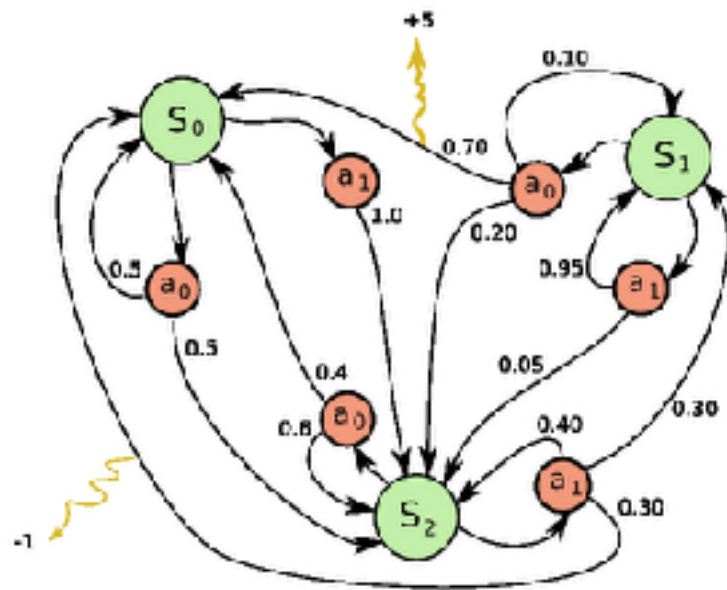


Maximize expected total reward:

$$J(\pi) = \mathbb{E}[r_1 + r_2 + \cdots + r_H | \pi]$$

# Progress of RL in Theory

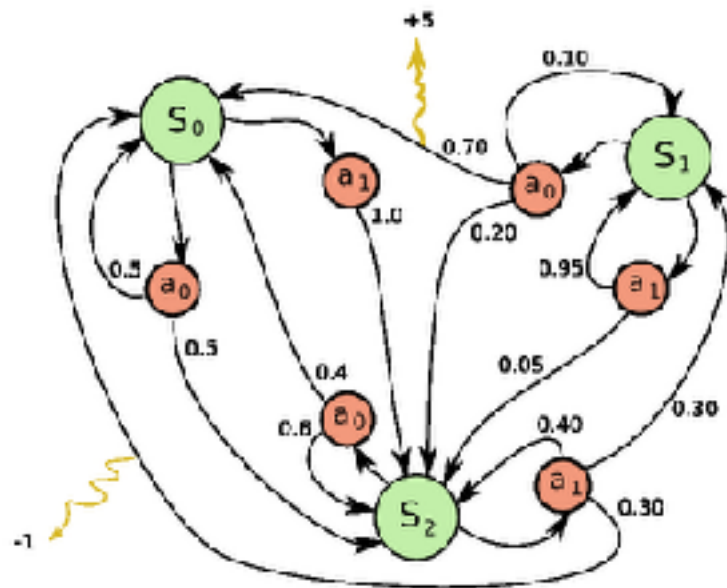
## Sample Efficiency in Small Discrete MDPs





# Progress of RL in Theory

## Sample Efficiency in Small Discrete MDPs



## Sample Complexity:

To achieve  $\epsilon$  near-optimal policy,  
need at most

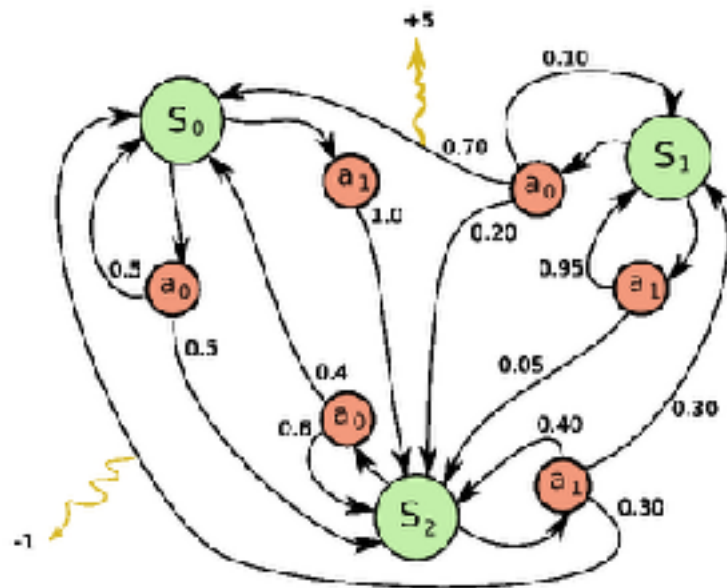
$\text{poly}(\# \text{ of states, } \# \text{ of actions, Horizon, } 1/\epsilon)$

many interactions

[e.g., Kearns & Singh, 02, Dann & Brunskill, 15, Azar et.al, 17]

# Progress of RL in Theory

## Sample Efficiency in Small Discrete MDPs



## Sample Complexity:

To achieve  $\epsilon$  near-optimal policy,  
need at most

$\text{poly}(\# \text{ of states}, \# \text{ of actions}, \text{Horizon}, 1/\epsilon)$

many interactions

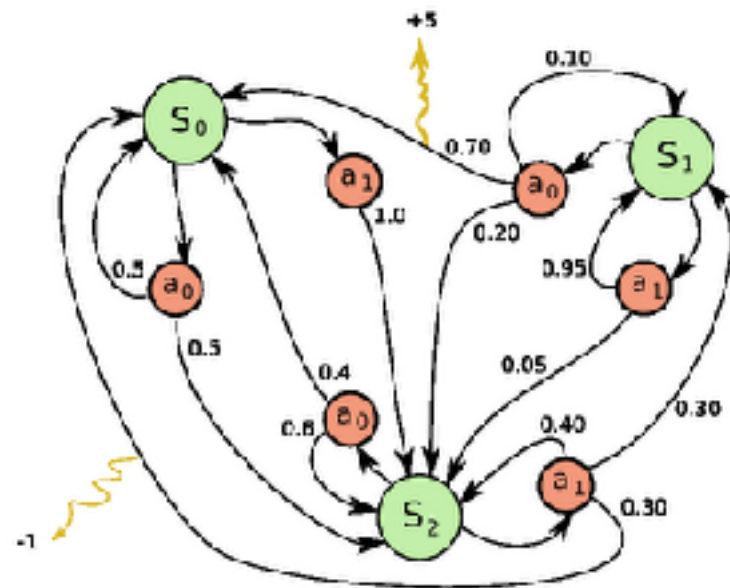
[e.g., Kearns & Singh, 02, Dann & Brunskill, 15, Azar et.al, 17]



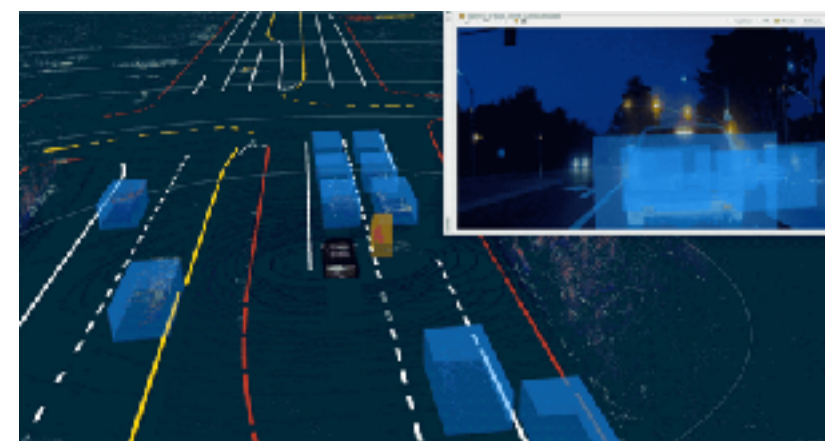
# Progress of RL in Theory

## Large-Scale Decision Making Problems

### Sample Efficiency in Small Discrete MDPs



$\neq$



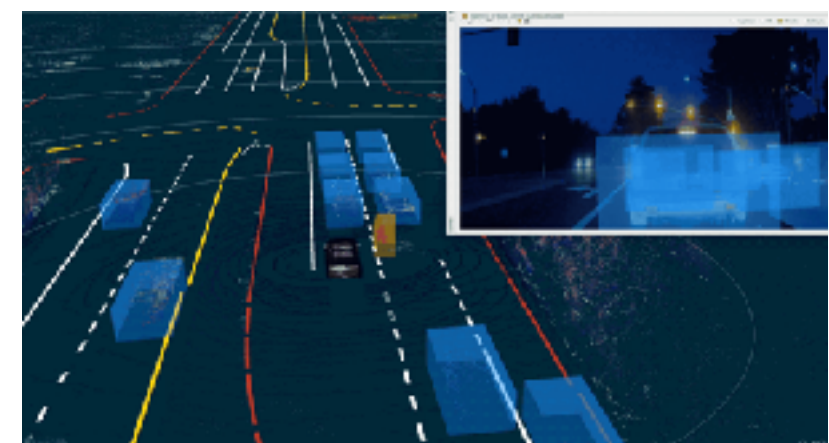
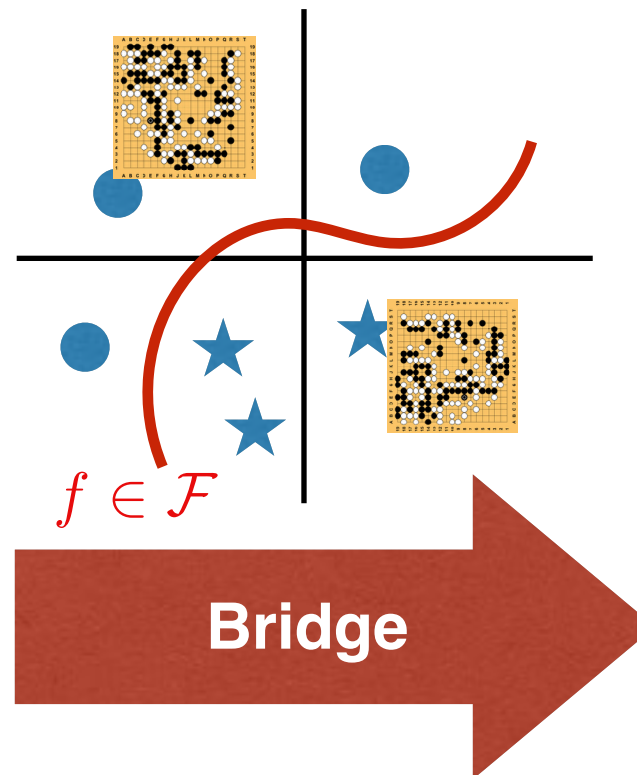
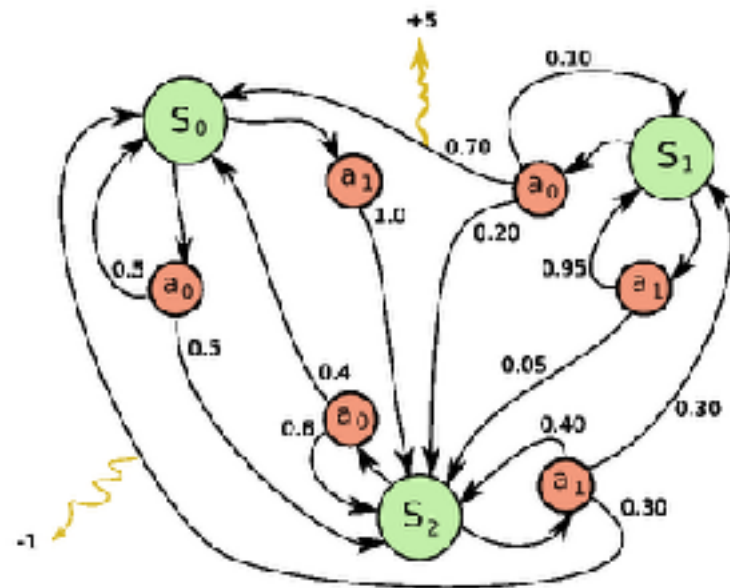
**Sample Complexity:**  
To achieve  $\epsilon$  near-optimal policy,  
need at most  
 $\text{poly}(\# \text{ of states}, \# \text{ of actions}, \text{Horizon}, 1/\epsilon)$   
many interactions

[e.g., Kearns & Singh, 02, Dann & Brunskill, 15, Azar et.al, 17]

# Progress of RL in Theory

## Large-Scale Decision Making Problems

### Sample Efficiency in Small Discrete MDPs



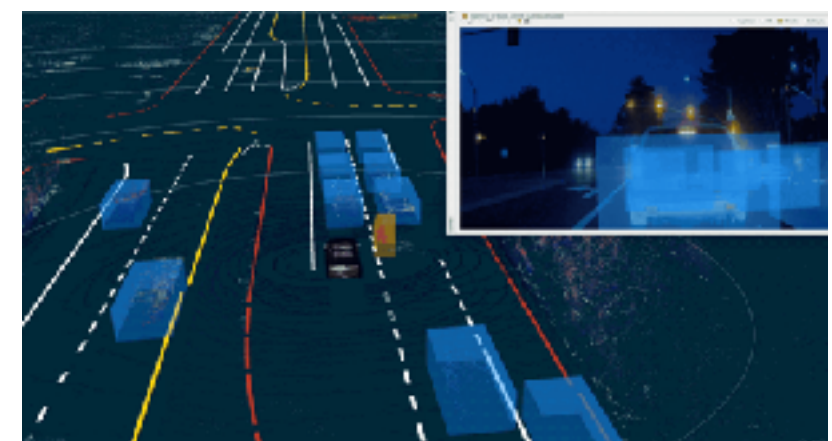
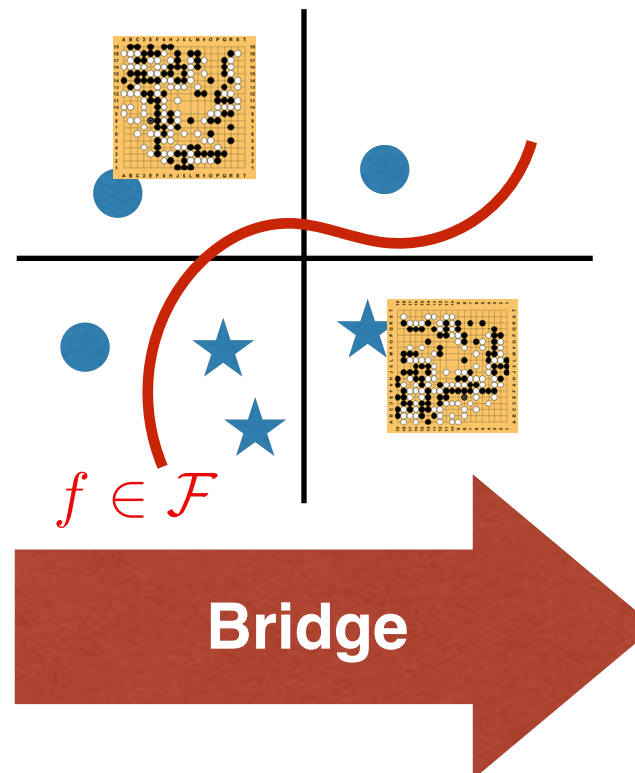
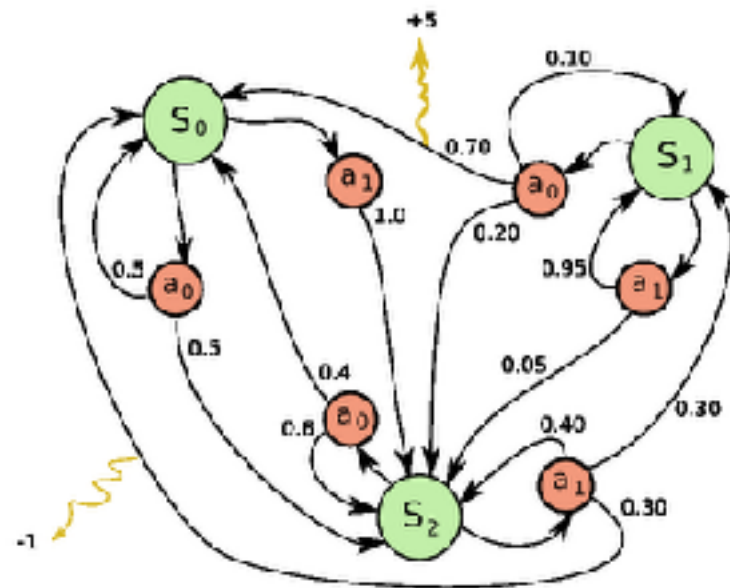
**Sample Complexity:**  
To achieve  $\epsilon$  near-optimal policy,  
need at most  
 $\text{poly}(\# \text{ of states}, \# \text{ of actions}, \text{Horizon}, 1/\epsilon)$   
many interactions

[e.g., Kearns & Singh, 02, Dann & Brunskill, 15, Azar et.al, 17]

# Progress of RL in Theory

## Large-Scale Decision Making Problems

### Sample Efficiency in Small Discrete MDPs



**Sample Complexity:**  
To achieve  $\epsilon$  near-optimal policy,  
need at most

$\text{poly}(\# \text{ of states, } \# \text{ of actions, Horizon, } 1/\epsilon)$

many interactions

[e.g., Kearns & Singh, 02, Dann & Brunskill, 15, Azar et.al, 17]

VC-dim

# **Previous Works on PAC RL w/ General Function Approximation**

# Previous Works on PAC RL w/ General Function Approximation

Contextual Bandits (horizon=1)

(e.g., Auer et al., 02, Langford & Zhang, 07)



# Previous Works on PAC RL w/ General Function Approximation

## Contextual Bandits (horizon=1)

(e.g., Auer et al., 02, Langford & Zhang, 07)

## Contextual Decision Process

(Krishnamurthy et al., 16, Jiang et al., 17, Dann et al., 18)



# Previous Works on PAC RL w/ General Function Approximation

## Contextual Bandits (horizon=1)

(e.g., Auer et al., 02, Langford & Zhang, 07)

## Contextual Decision Process

(Krishnamurthy et al., 16, Jiang et al., 17, Dann et al., 18)

## Model-based vs Model-free

# **Previous Works on PAC RL w/ General Function Approximation**

Contextual Bandits (horizon=1)

(e.g., Auer et al., 02, Langford & Zhang, 07)

Contextual Decision Process

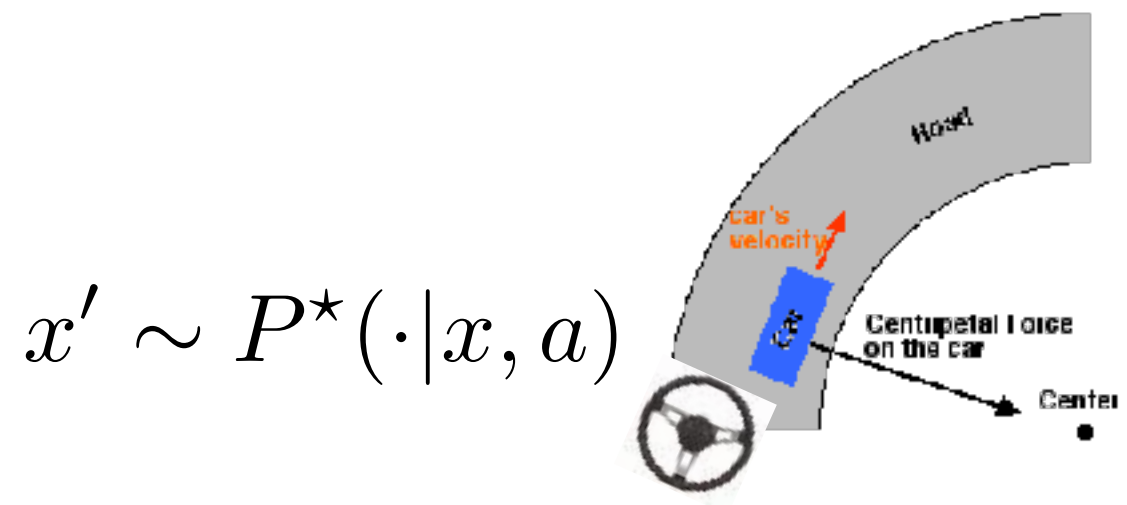
(Krishnamurthy et al., 16, Jiang et al., 17, Dann et al., 18)

## **Model-based vs Model-free**

## **A PAC model-based Algorithm**

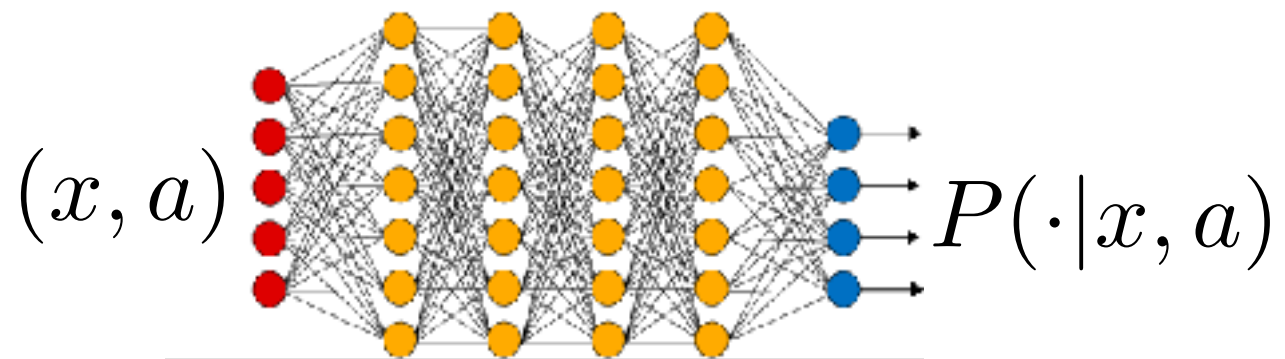
# Setup of Model-Based RL

# Setup of Model-Based RL

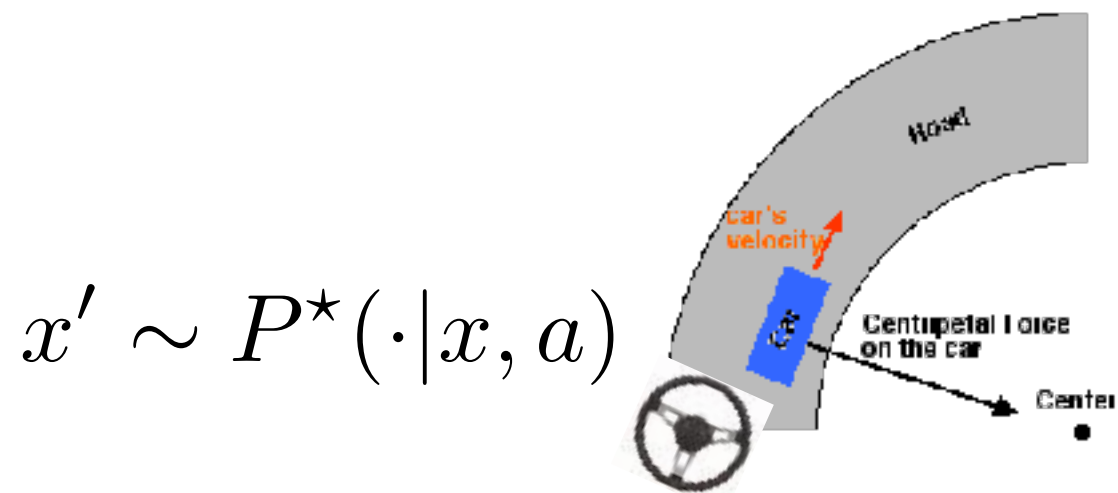


# Setup of Model-Based RL

## Function Approximators

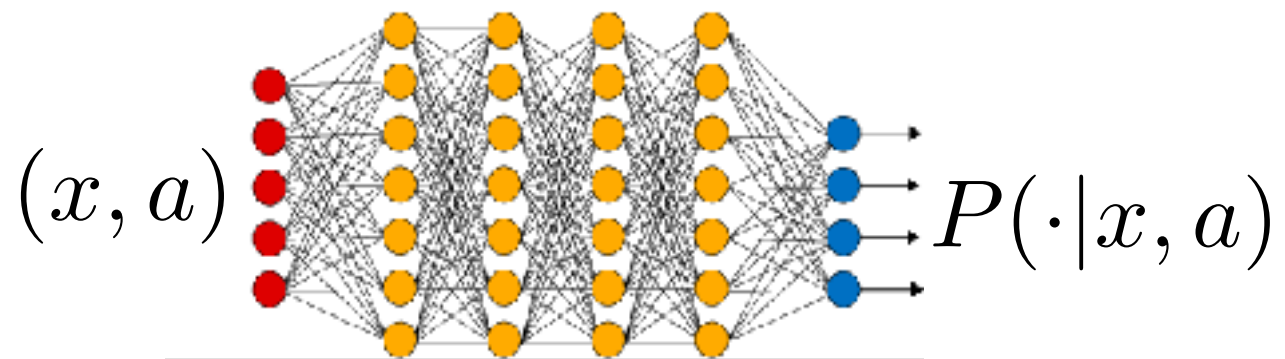


$$\mathcal{P} = \{P : \mathcal{X} \times \mathcal{A} \rightarrow \Delta(\mathcal{X})\}$$



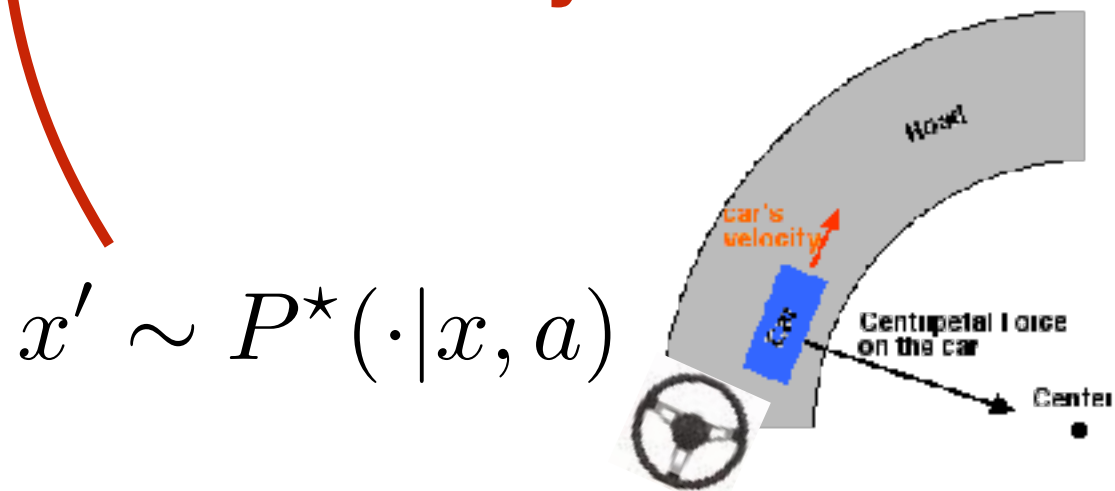
# Setup of Model-Based RL

## Function Approximators



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

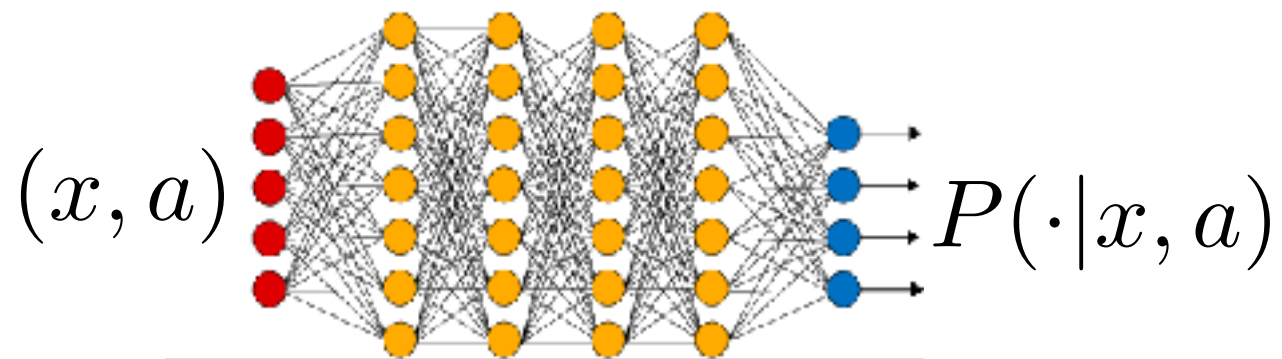
**1. Realizability:**  $P^* \in \mathcal{P}$





# Setup of Model-Based RL

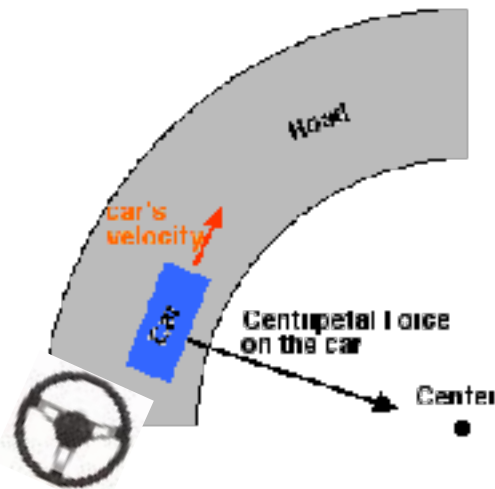
## Function Approximators



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

**1. Realizability:**  $P^* \in \mathcal{P}$

$$x' \sim P^*(\cdot|x, a)$$



**2. Access to Optimal Planner (OP)**

$$OP(P, r) \Rightarrow \pi_P$$

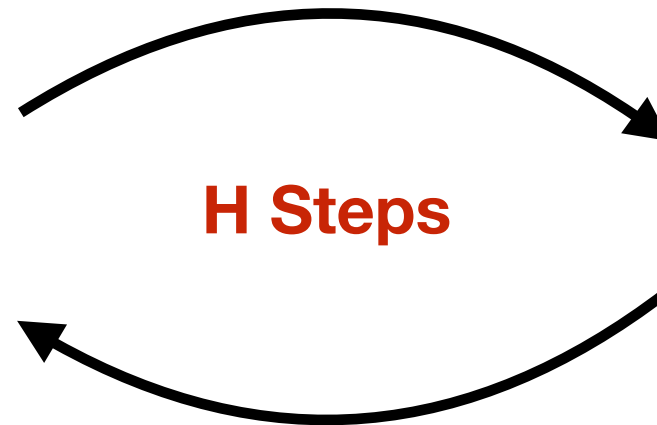
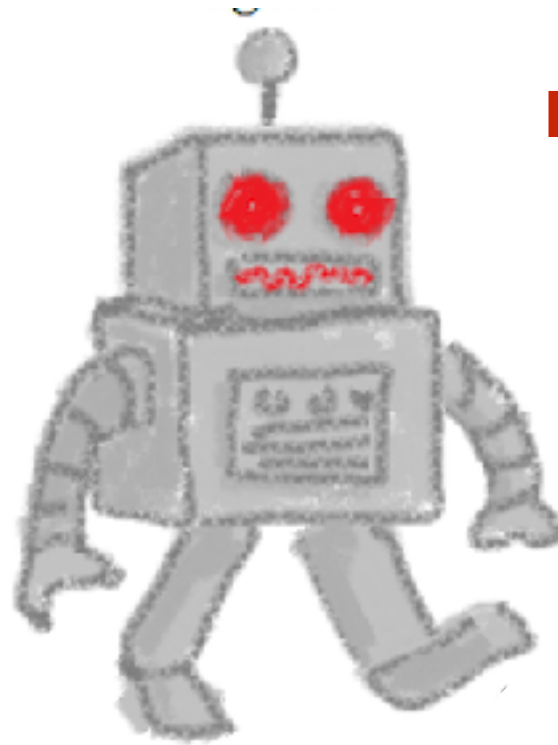
# Definition of Model-free RL

**Input:**  $\mathcal{Q} \triangleq \{Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}\}$

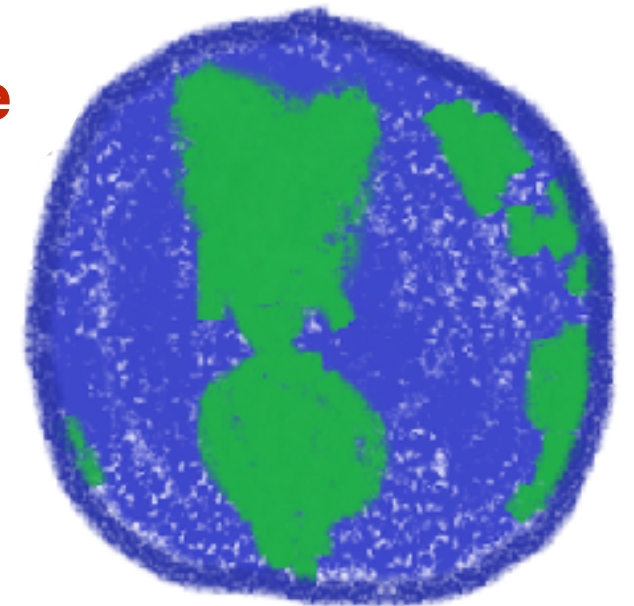
# Definition of Model-free RL

**Input:**  $Q \triangleq \{Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}\}$

**Policy:** determine **action** based on **state**



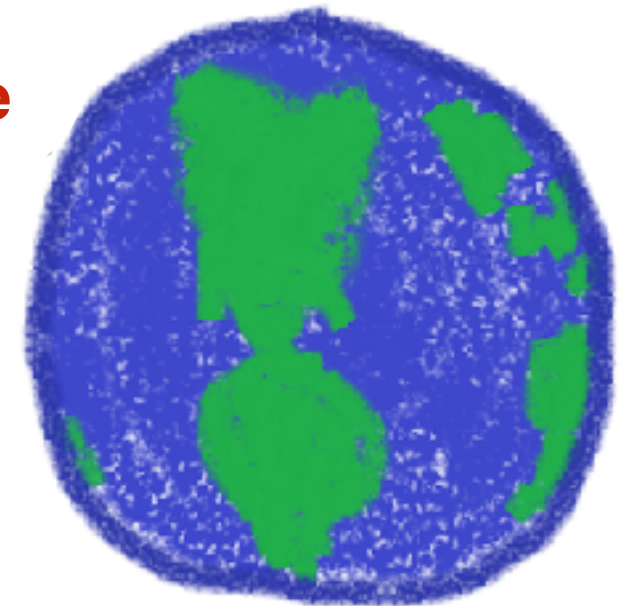
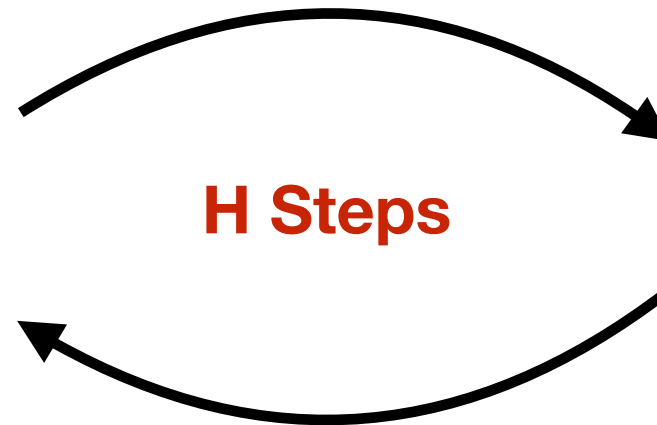
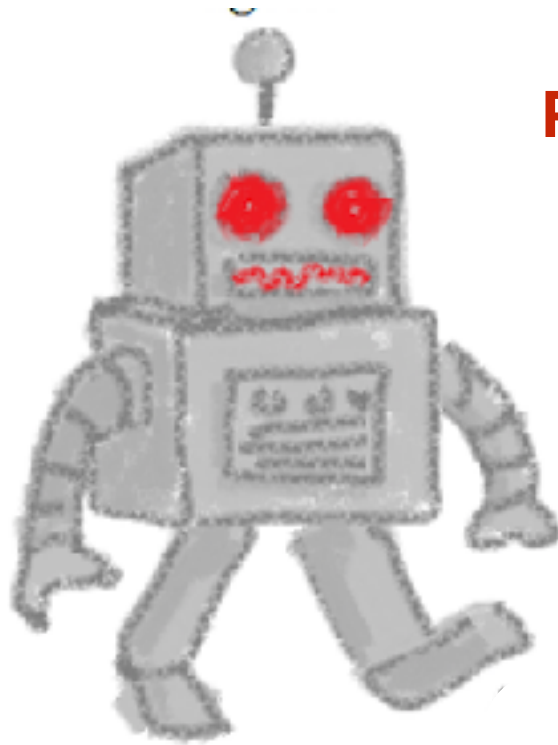
**Env:** reveal next state  $\mathcal{X}$



# Definition of Model-free RL

**Input:**  $\mathcal{Q} \triangleq \{Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}\}$

**Policy:** determine **action** based on **state**



**Q-profile:**

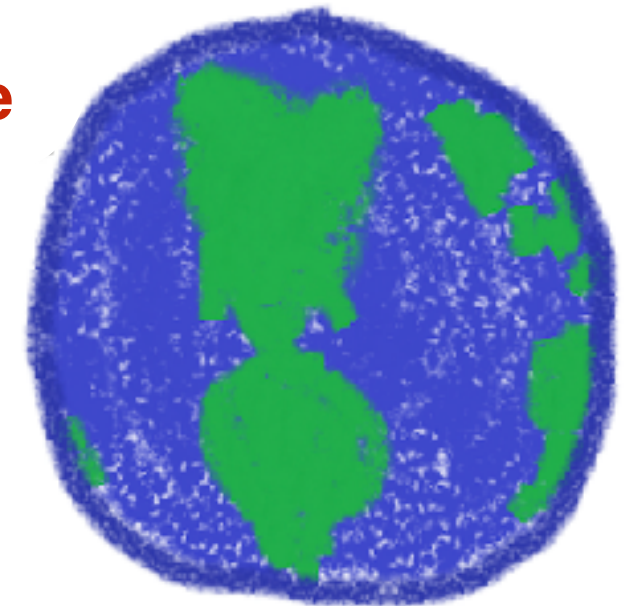
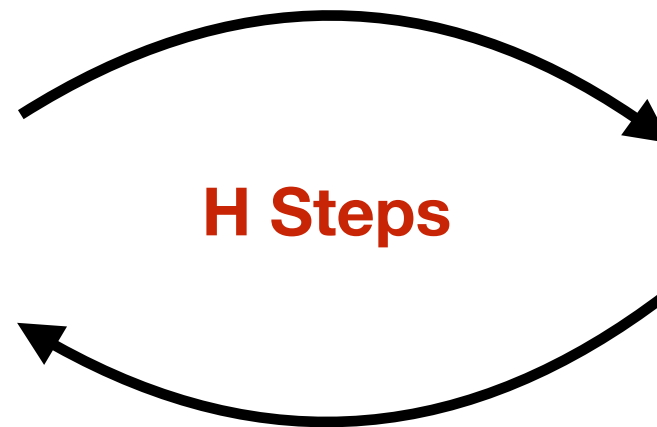
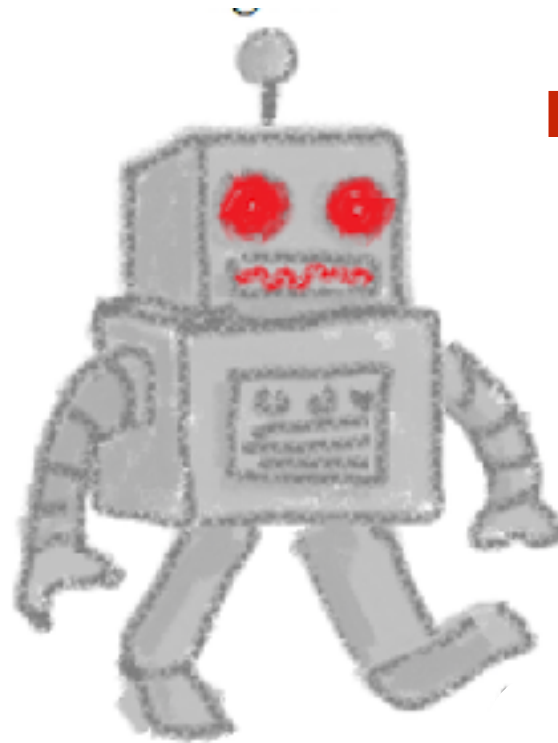
$$[Q(x, a)]_{Q \in \mathcal{Q}, a \in \mathcal{A}}$$

all possible Q values evaluated at state  $x$

# Definition of Model-free RL

**Input:**  $\mathcal{Q} \triangleq \{Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}\}$

**Policy:** determine **action** based on **state**



**Q-profile:**

$$[Q(x, a)]_{Q \in \mathcal{Q}, a \in \mathcal{A}}$$

all possible Q values evaluated at state  $x$

Efficient Q-learning (Jin et.al, 18)

Fitted Q-Iteration (Ernst et.al., 05)

OLIVE (Jiang et.al, 17)

Policy Gradient (Williams 92)

# **An Exponential Improvement over Model-free RL**



# An Exponential Improvement over Model-free RL

There exists MDPs (e.g., Factored MDPs), s.t., to learn near optimal policy,

Model-Based RL:

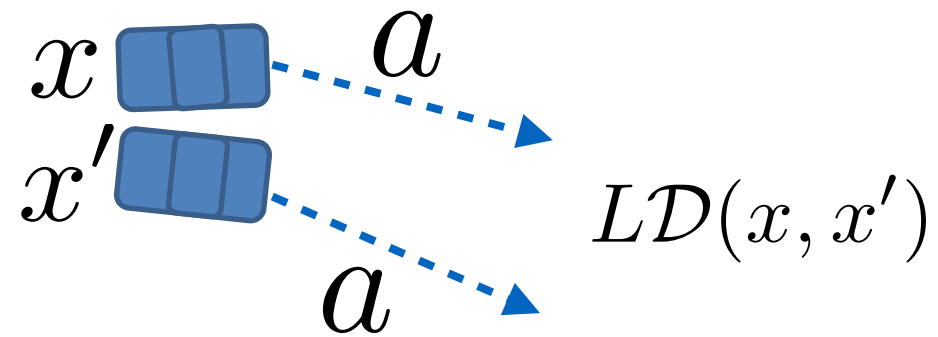
**Polynomial** Sample  
Complexity

**VS**

**ANY** Model-Free RL:

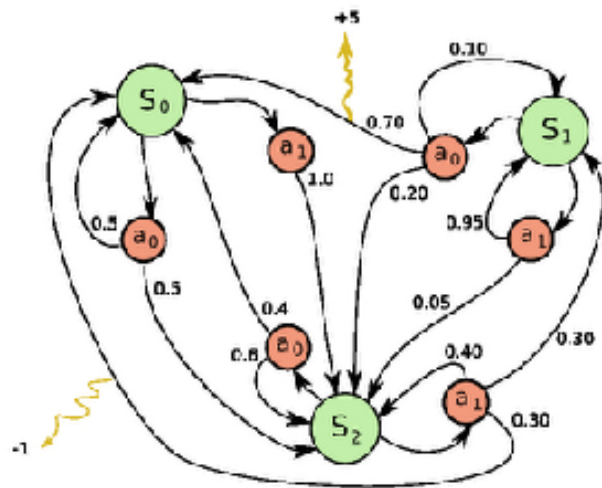
$\Omega(\exp(H))$

# We have been studying model-based RL, BUT...



## Lipschitz MDPs

[Kearns, Langford, Kakade, 03]



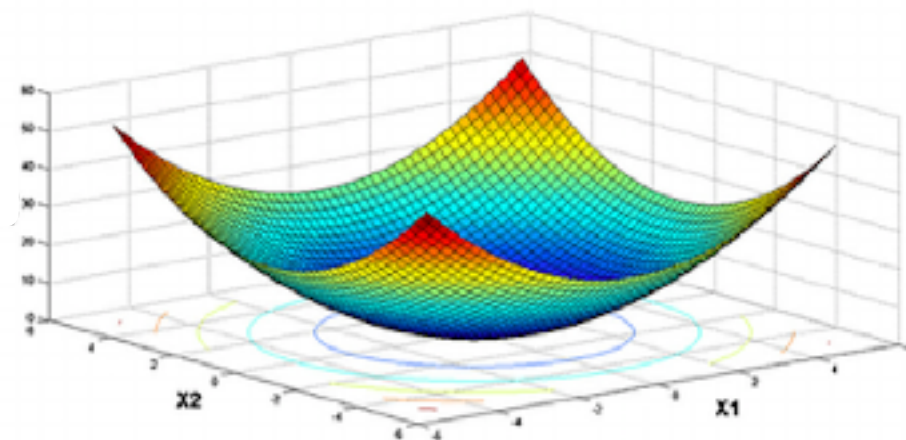
## Small Tabular MDP

[Kearns & Singh, 02]



## Factored MDPs

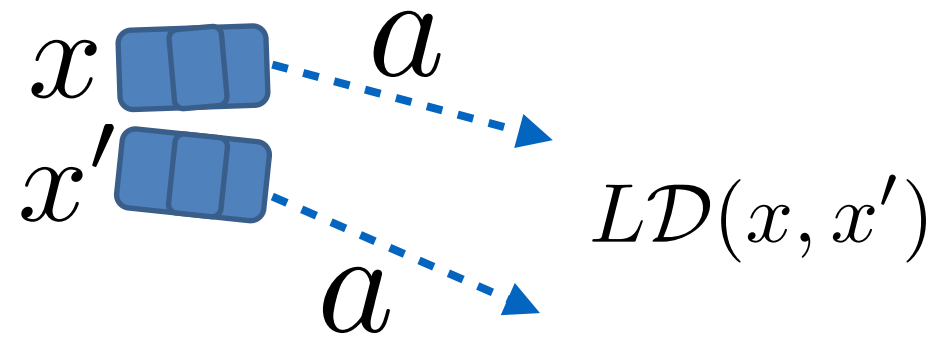
[Guestrin et.al, 03; Osband & Van Roy, 13]



## Linear Quadratic Regulator (LQR)

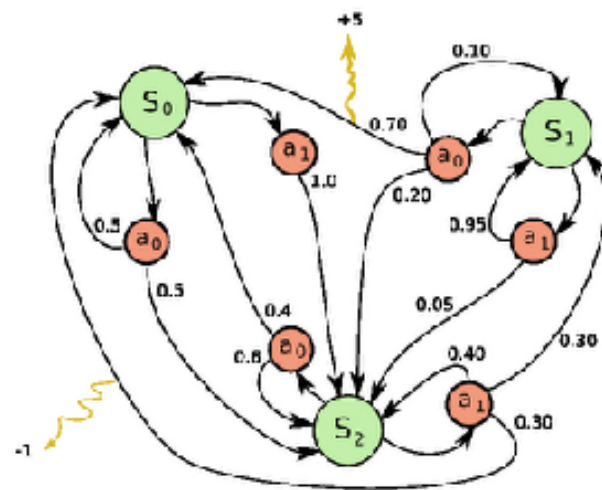
[Dean et.al, 18]

# We have been studying model-based RL, BUT...



Lipschitz MDPs

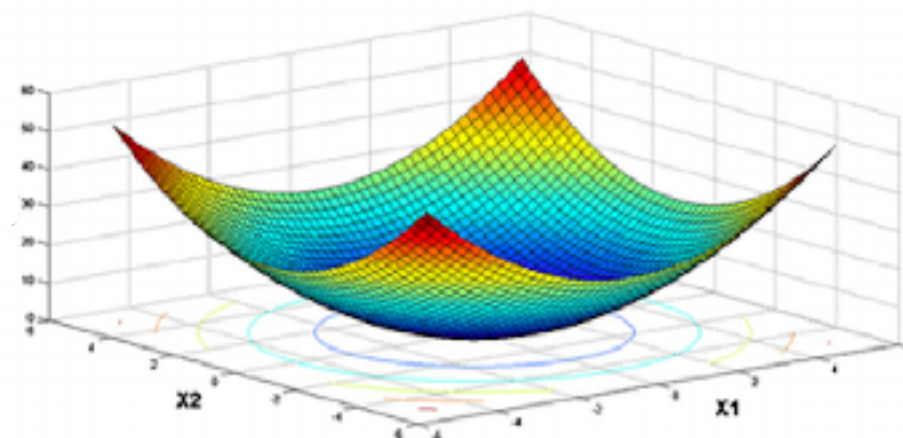
[Kearns, Langford, Kakade, 03]



Small Tabular MDP

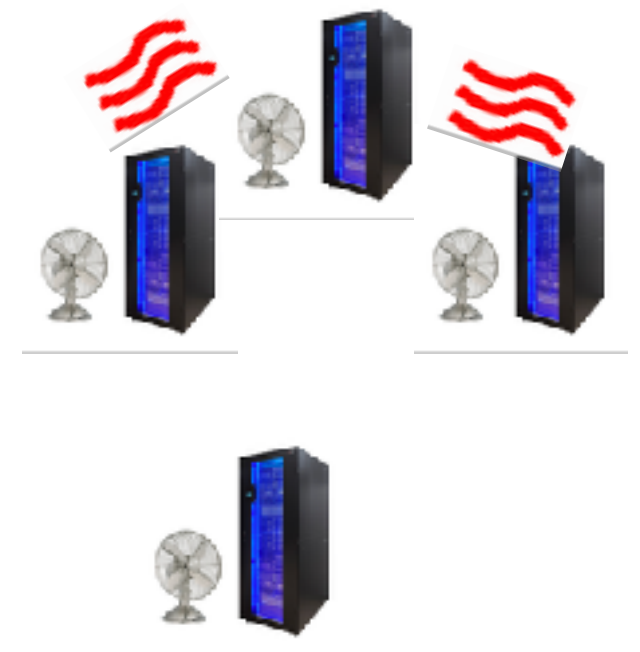
[Kearns & Singh, 02]

**A Unified  
Algorithm?**



Linear Quadratic Regulator (LQR)

[Dean et.al, 18]

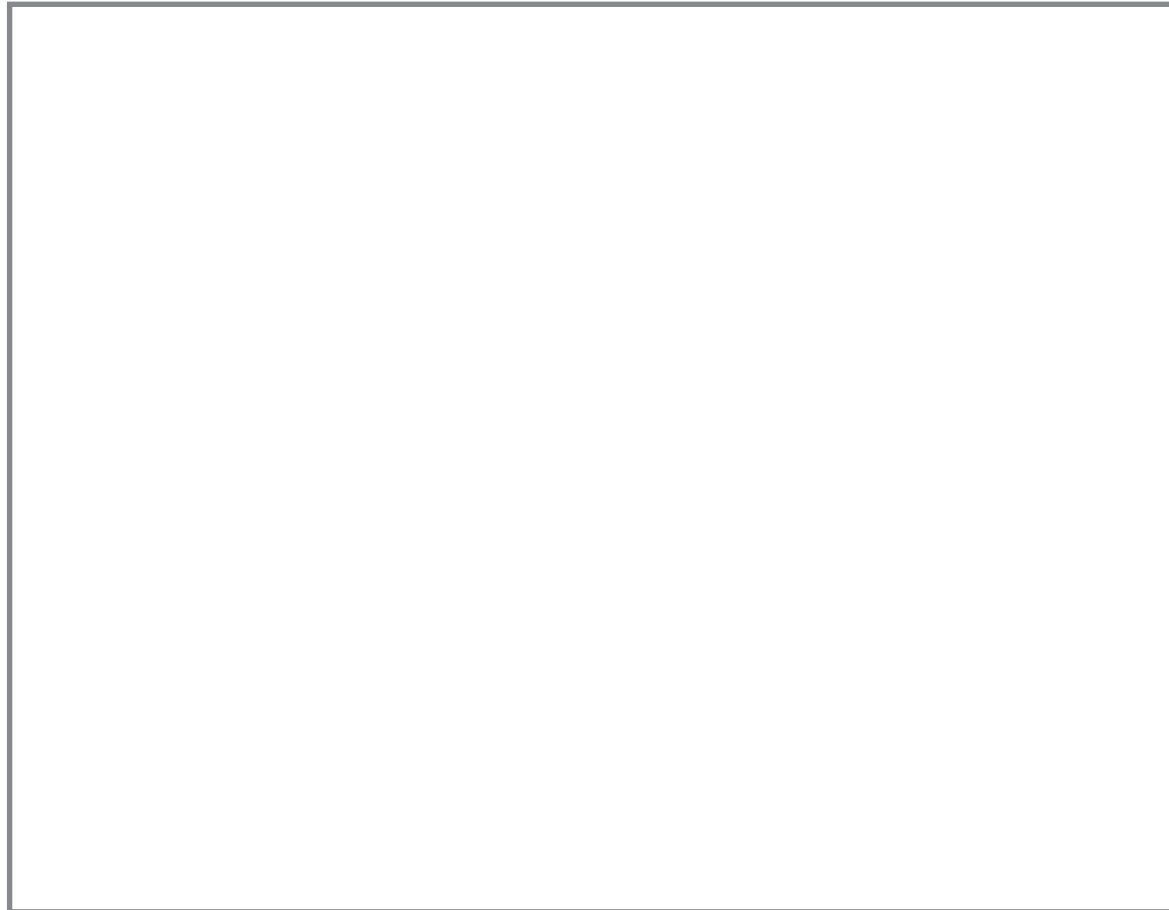


Factored MDPs

[Guestrin et.al, 03; Osband & Van Roy, 13]

# A Unified Measure—*Witness Rank*

**Misfit Matrix:**

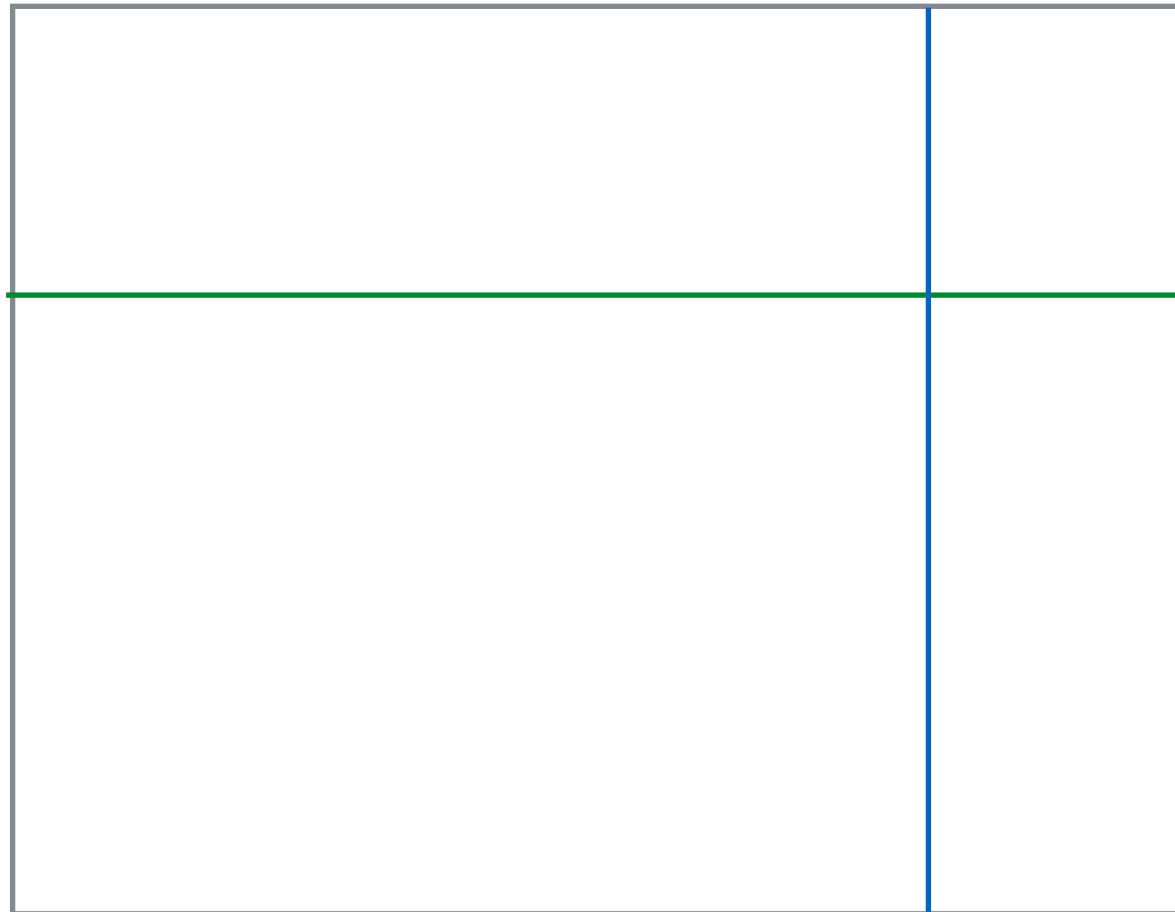


$$\in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$$

# A Unified Measure—*Witness Rank*

**Misfit Matrix:**

$P_r$

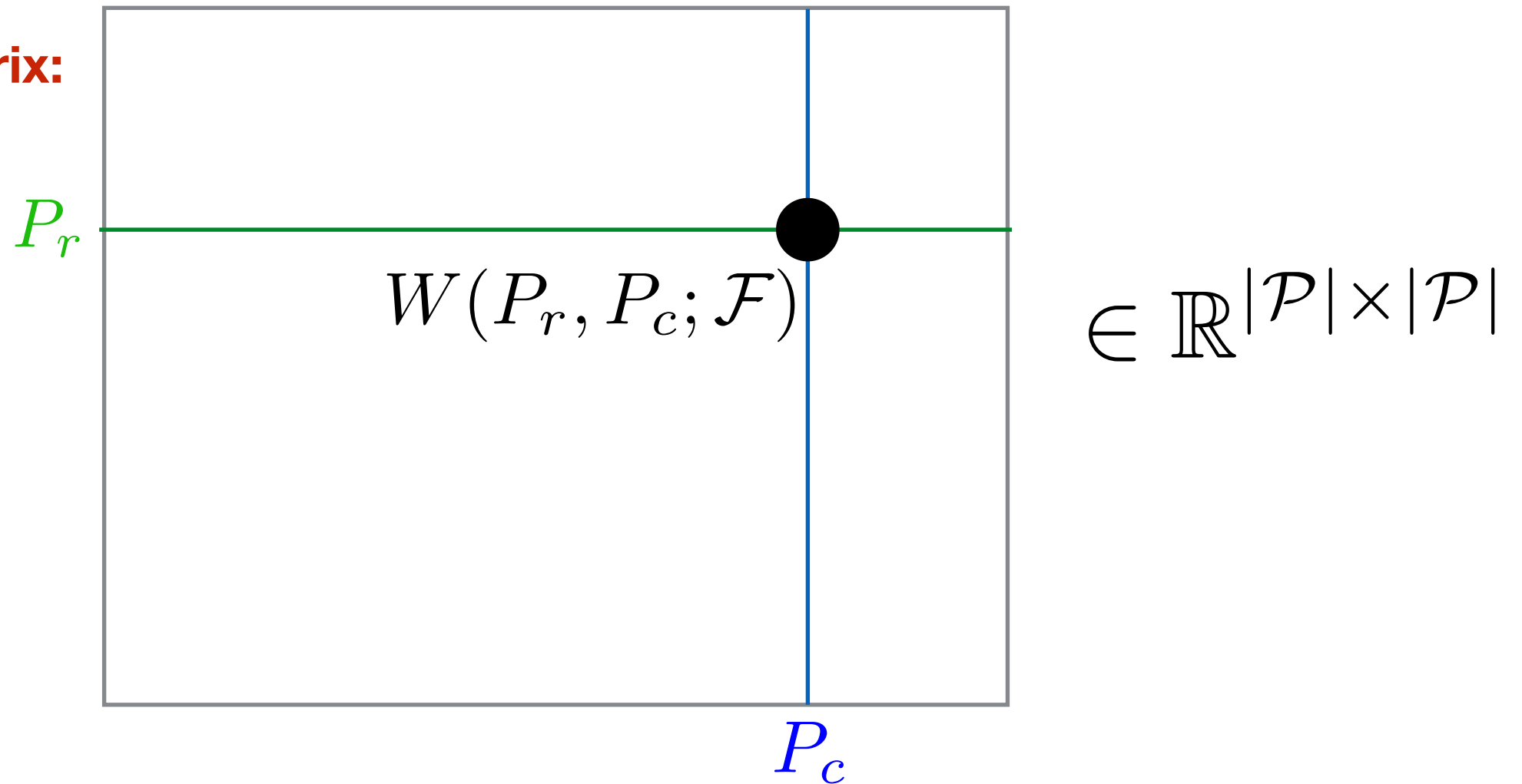


$P_c$

$$\in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$$

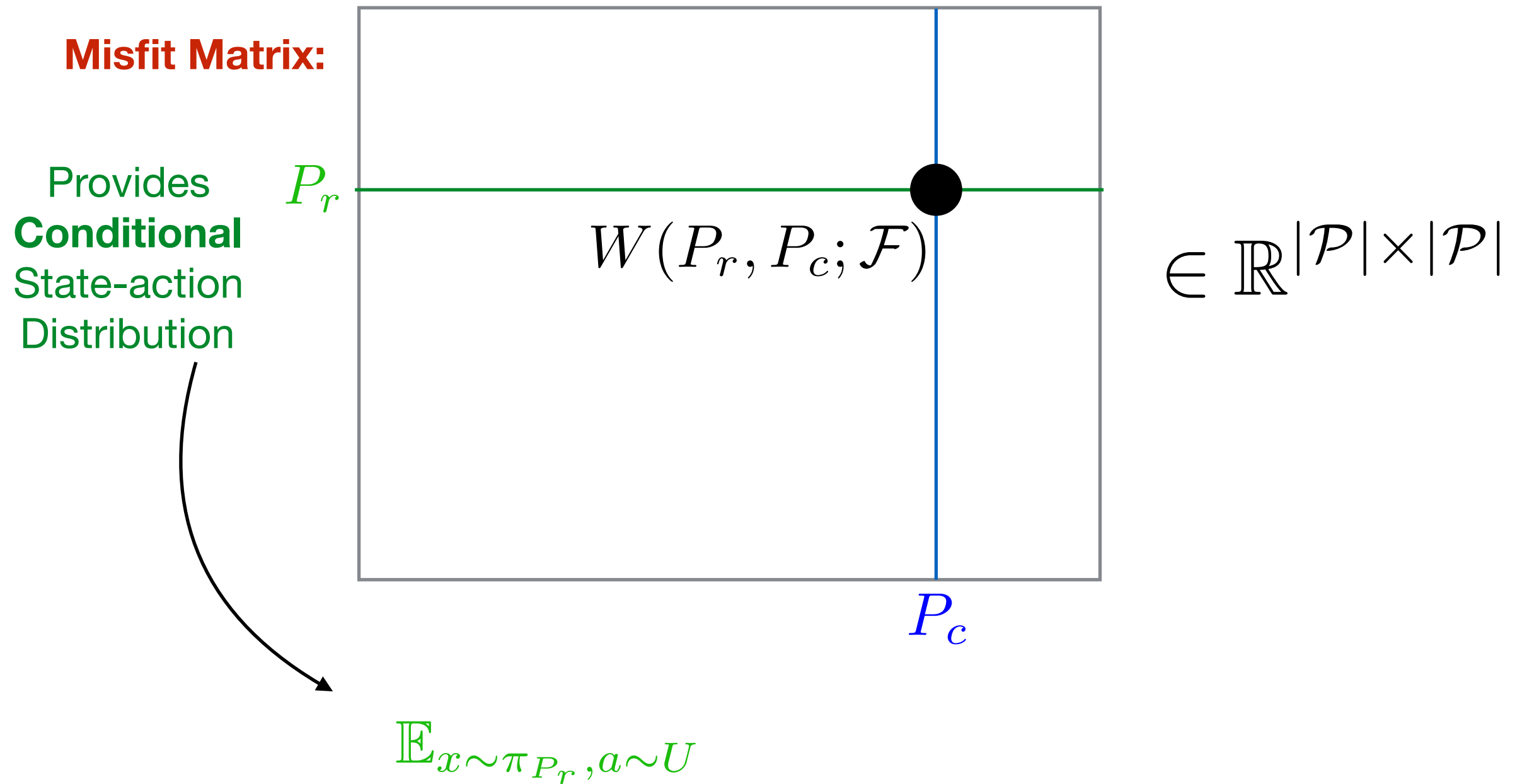
# A Unified Measure—*Witness Rank*

Misfit Matrix:

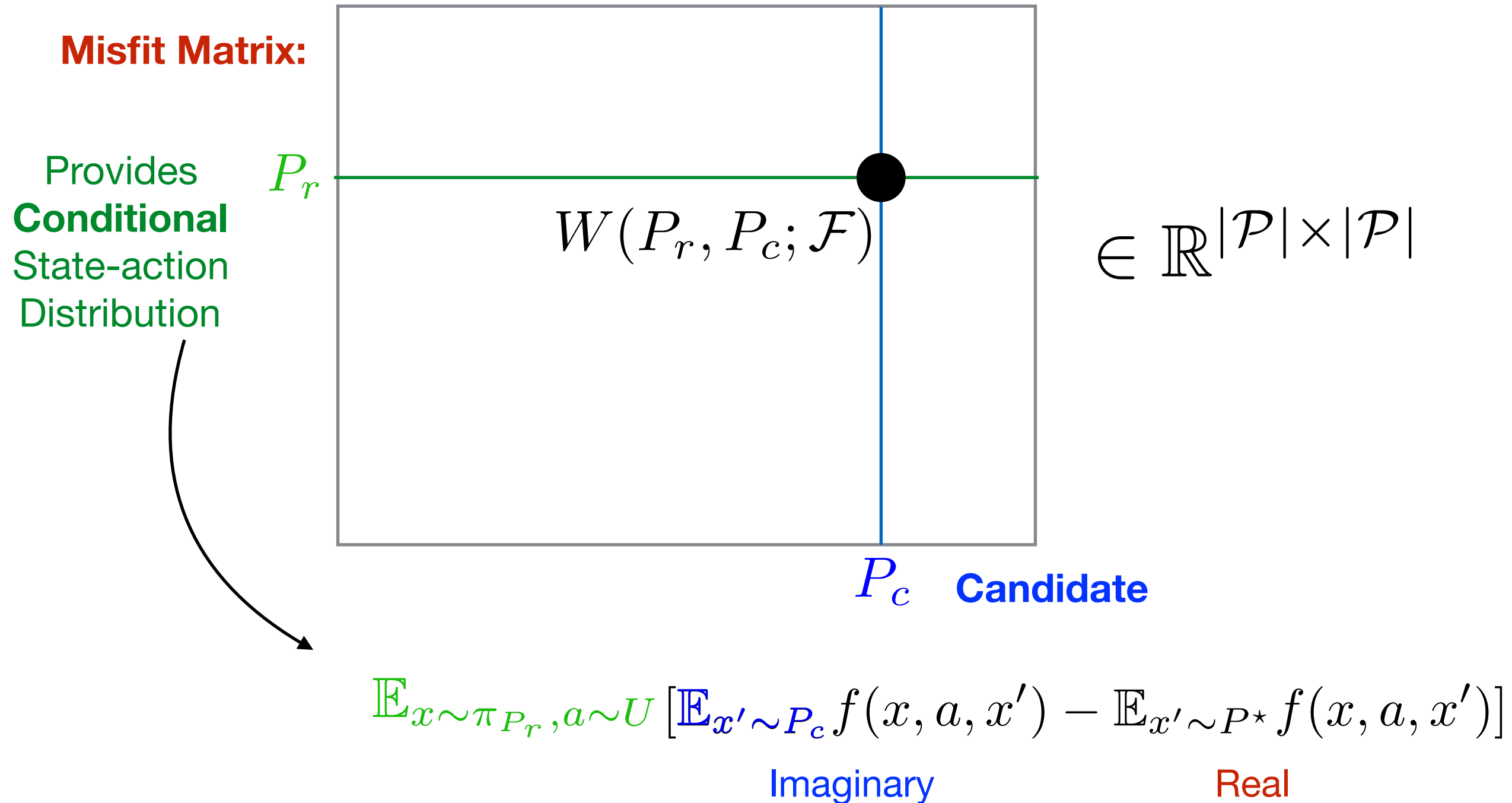




# A Unified Measure—*Witness Rank*



# A Unified Measure—*Witness Rank*

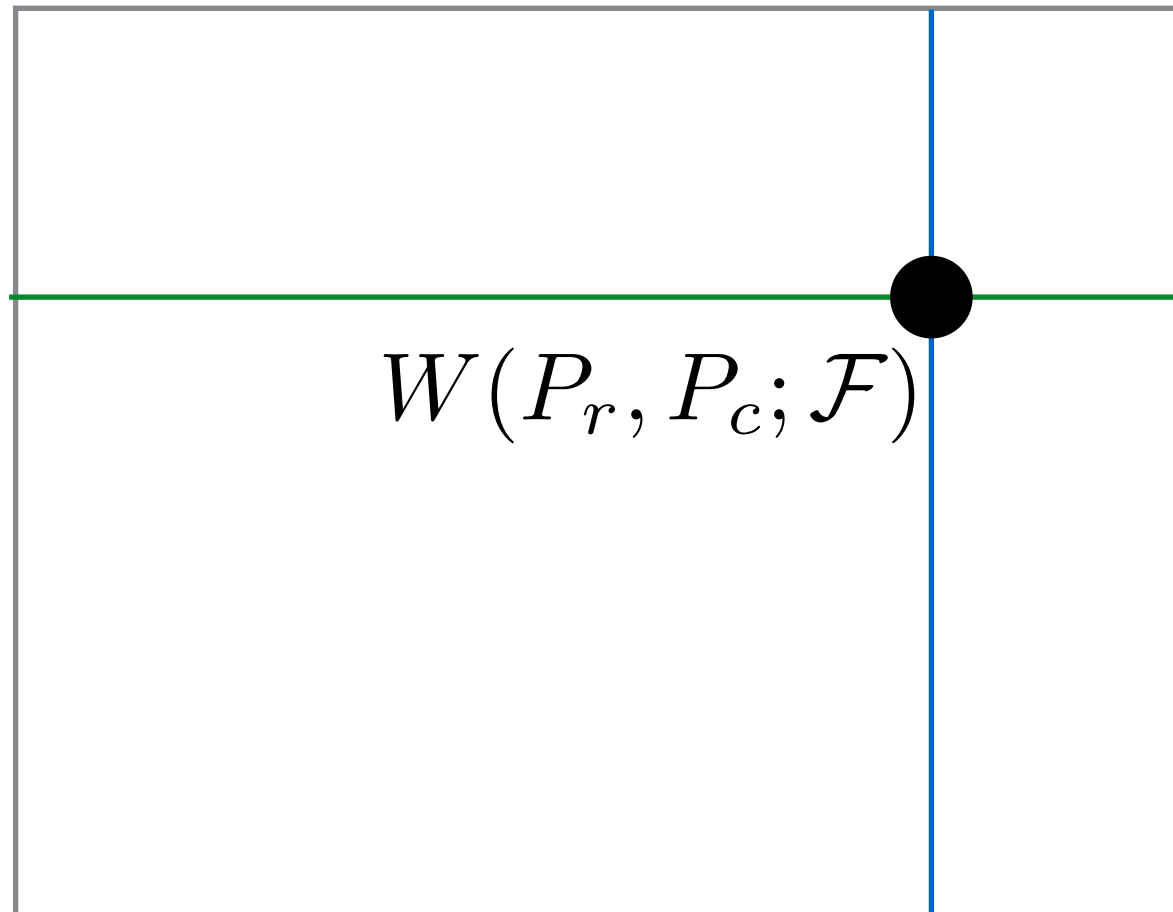


# A Unified Measure—*Witness Rank*

**Misfit Matrix:**

Provides  
**Conditional**  
State-action  
Distribution

$P_r$



$\in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$

$P_c$  **Candidate**

$$W(P_r, P_c; \mathcal{F}) = \max_{f \in \mathcal{F}} \mathbb{E}_{x \sim \pi_{P_r}, a \sim U} [\mathbb{E}_{x' \sim P_c} f(x, a, x') - \mathbb{E}_{x' \sim P^*} f(x, a, x')]$$

Witness functions  
(aka discriminators)

Imaginary

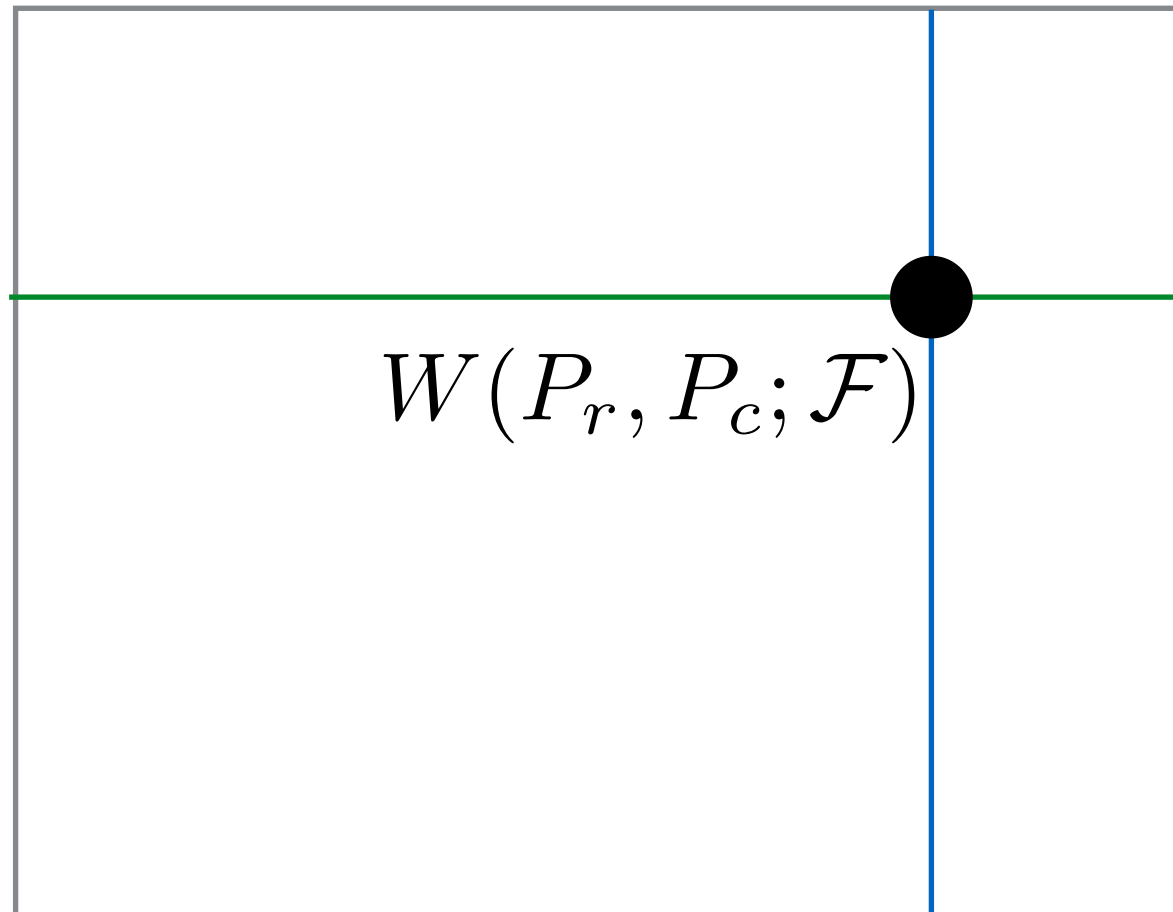
Real

# A Unified Measure—*Witness Rank*

**Misfit Matrix:**

Provides  
**Conditional**  
State-action  
Distribution

$P_r$



$W(P_r, P_c; \mathcal{F})$

$\in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$

$P_c$  **Candidate**

$$W(P_r, P_c; \mathcal{F}) = \max_{f \in \mathcal{F}} \mathbb{E}_{x \sim \pi_{P_r}, a \sim U} [\mathbb{E}_{x' \sim P_c} f(x, a, x') - \mathbb{E}_{x' \sim P^*} f(x, a, x')]$$

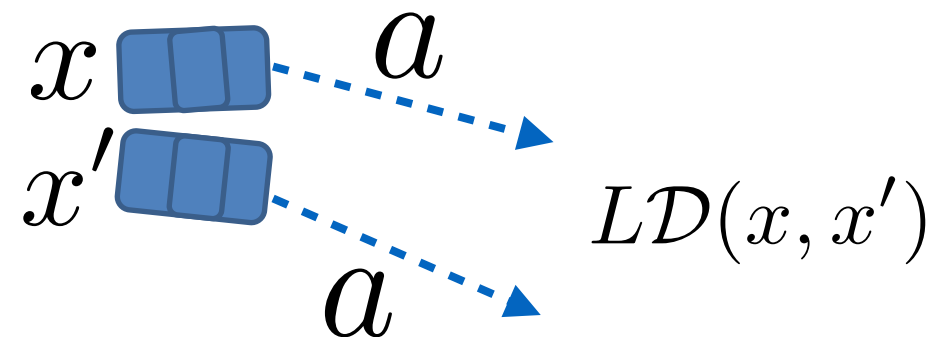
Imaginary

Real

Witness functions  
(aka discriminators)

**Witness Rank**  $\triangleq$  rank of this misfit matrix

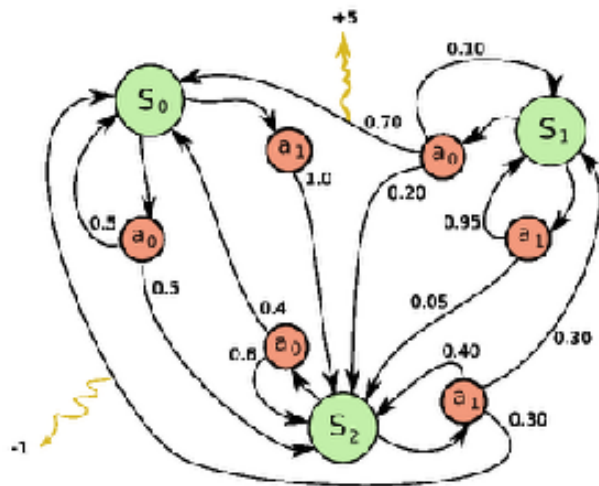
# Capture Complexities of Existing RL problems



Lipschitz Continuous MDPs

[Kearns, Langford, Kakade, 03]

**Rank  $\leq$  Covering number  
of state space**



Small Discrete MDP

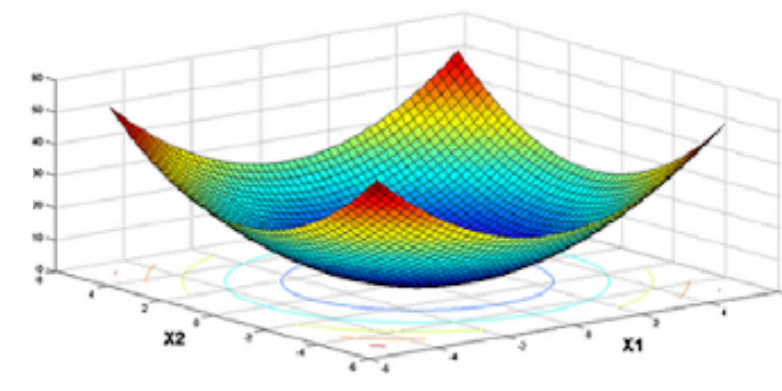
**Rank  $\leq$  # of state**



Factored MDPs

[Guestrin et.al, 03; Osband & Van Roy, 13]

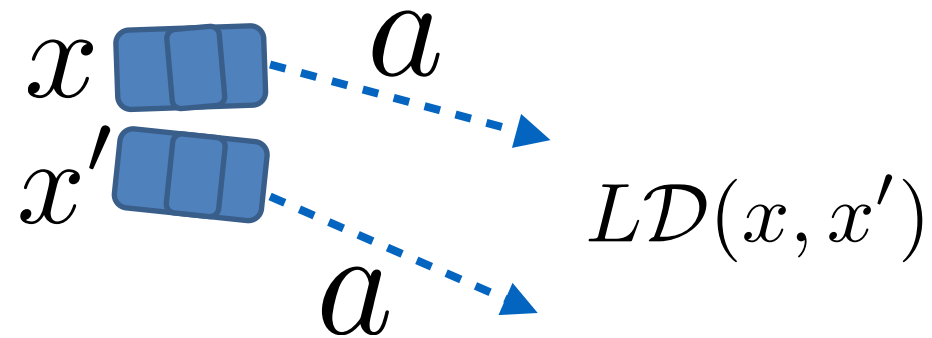
**Rank  $\leq \exp(\text{in-degree})$**



LQR

**Rank  $\leq O(d^2)$**

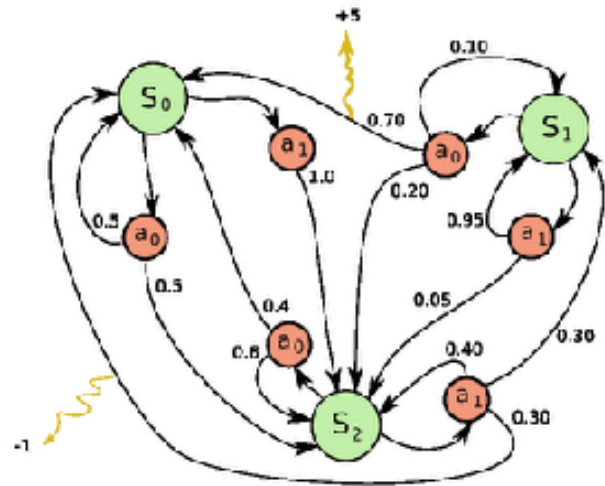
# Capture Complexities of Existing RL problems



Lipschitz Continuous MDPs

[Kearns, Langford, Kakade, 03]

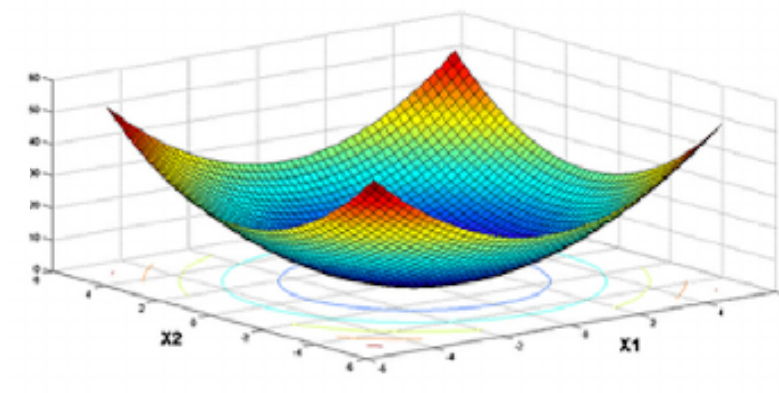
**Rank  $\leq$  Covering number  
of state space**



Small Discrete MDP

**Rank  $\leq$  # of state**

**A Unified  
Algorithm!**



LQR

**Rank  $\leq O(d^2)$**



Factored MDPs

[Guestrin et.al, 03; Osband & Van Roy, 13]

**Rank  $\leq \exp(\text{in-degree})$**

# Sample Complexity

To achieve  $\epsilon$  near-optimal policy (w/ high probability)

# Sample Complexity

To achieve  $\epsilon$  near-optimal policy (w/ high probability)

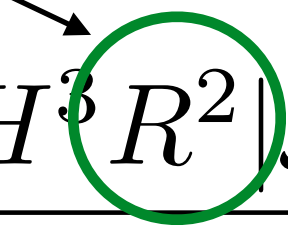
$$\tilde{O}\left(\frac{H^3 R^2 |\mathcal{A}|}{\epsilon^2} \log\left(\frac{|\mathcal{F}| |\mathcal{P}|}{\delta}\right)\right)$$



# Sample Complexity

To achieve  $\epsilon$  near-optimal policy (w/ high probability)

Witness Rank


$$\tilde{O}\left(\frac{H^3 R^2 |\mathcal{A}|}{\epsilon^2} \log\left(\frac{|\mathcal{F}| |\mathcal{P}|}{\delta}\right)\right)$$

# Sample Complexity

To achieve  $\epsilon$  near-optimal policy (w/ high probability)

Witness Rank

$$\tilde{O}\left(\frac{H^3 R^2 |\mathcal{A}|}{\epsilon^2} \log\left(\frac{|\mathcal{F}| |\mathcal{P}|}{\delta}\right)\right)$$

Complexities of  
model class & discriminator class

**Poly Dependency on # of states**

# Take-home Messages

# Take-home Messages

Model-based RL could be exponentially more sample efficient than model-free ones

# Take-home Messages

Model-based RL could be exponentially more sample efficient than model-free ones

Sample efficiency is possible when Witness Rank is small