# Towards Generalization and Efficiency in Reinforcement Learning

## Wen Sun

Carnegie Mellon University

Joint work with Drew Bagnell, Geoff Gordon, Byron Boots, John Langford, Stephane Ross, Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, Arun Venkatraman
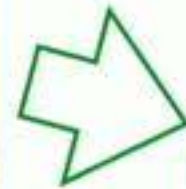
**Carnegie Mellon**
**THE ROBOTICS INSTITUTE**

# Goal:

Design Algorithms that have
**Generalization & Sample Efficiency**
in learning to make decisions
in complex environments

# My Research

## 1. Expert Demonstration



[**Sun**, Venkatraman, Gordon, Boots, Bagnell, 17, ICML]

[**Sun**, Gordon, Boots, Bagnell, 18, NeurIPS]

### All Sequential Decision Making Problems
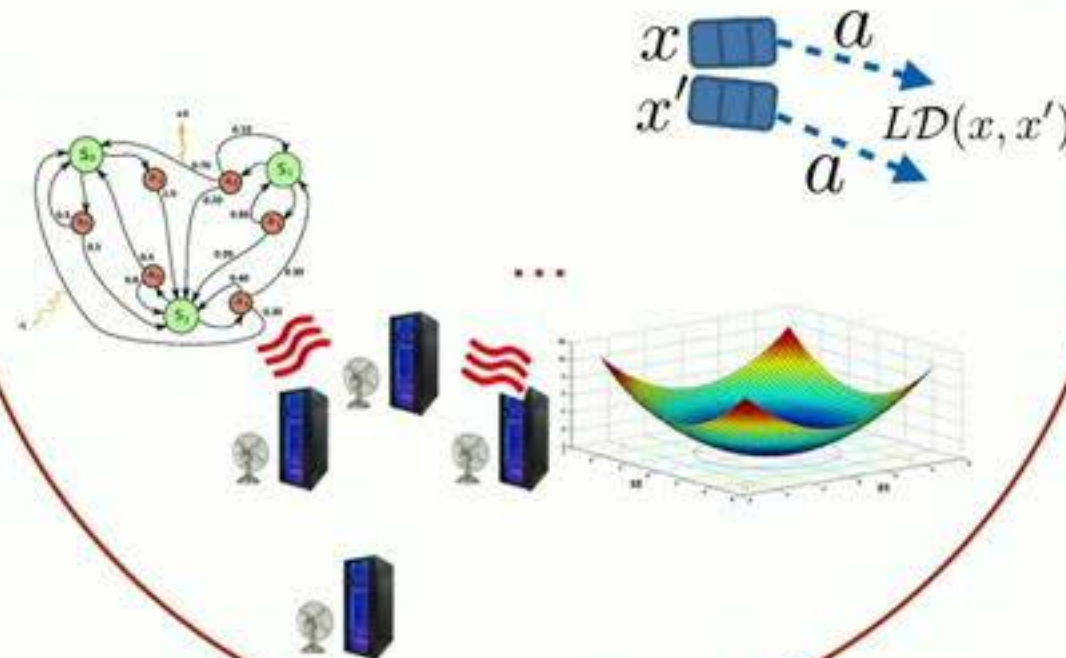
# My Research

## 1. Expert Demonstration



[**Sun**, Venkatraman, Gordon, Boots, Bagnell, 17, ICML]
[**Sun**, Gordon, Boots, Bagnell, 18, NeurIPS]
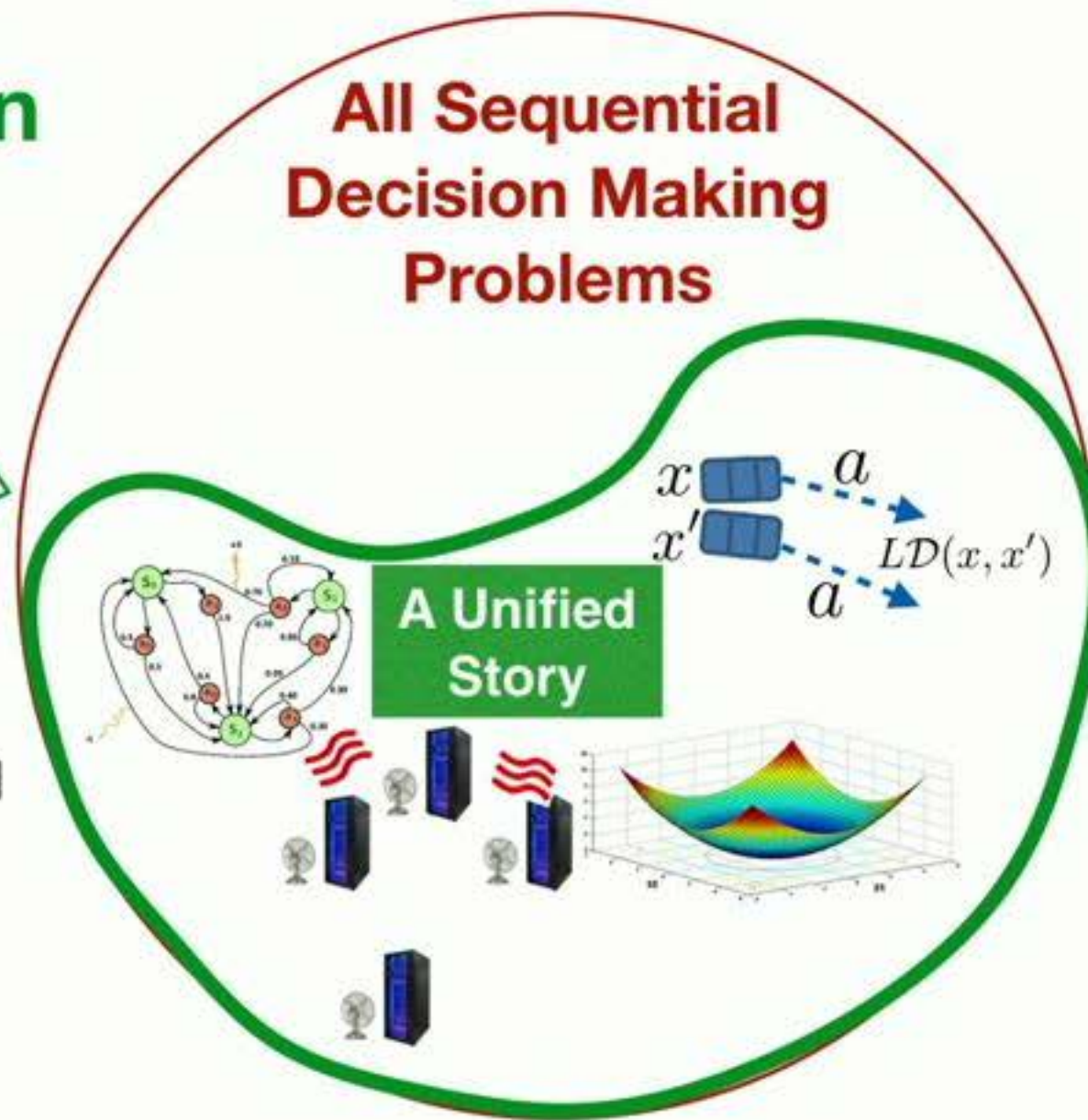
**All Sequential Decision Making Problems**

A Unified Story

$x$ $a$

$x'$ $L\mathcal{D}(x, x')$

$a$

## 2. Exploiting Structures

[**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

3

# Supervised Learning VS Sequential Decision Making

Given i.i.d examples at training:



**Passive:**

MARIO: 1024    COINS    DIFFICULTY TIME
                 05      1           052
FPS: 24                 WorldPause
Attempt: 1 of 1         false
AgentLinear
Selected Actions:

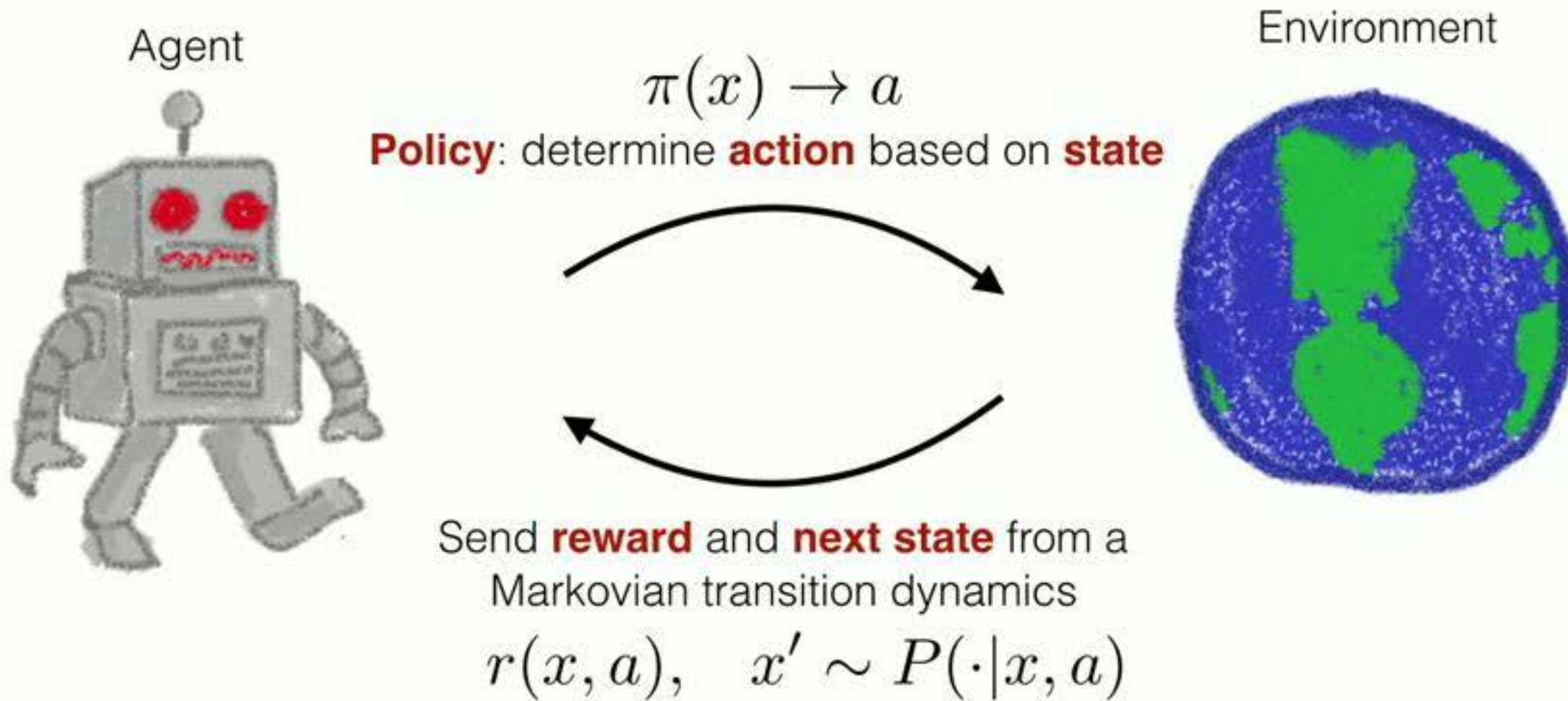RIGHT            SPEED

Active: Decisions ➡ Data Distribution

[Ross&Bagnell, 11, AISTATS]

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \to a$$

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

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \to a$$

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

Send **reward** and **next state** from a
Markovian transition dynamics

$$r(x, a), \quad x' \sim P(\cdot | x, a)$$

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \rightarrow a$$

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

**H Steps**

Send **reward** and **next state** from a
Markovian transition dynamics

$$r(x, a), \quad x' \sim P(\cdot|x, a)$$

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \rightarrow a$$

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

**H Steps**

Send **reward** and **next state** from a
Markovian transition dynamics

$$r(x,a), \quad x' \sim P(\cdot|x,a)$$

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \to a$$

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

**H Steps**

Send **reward** and **next state** from a
Markovian transition dynamics

$$r(x, a), \quad x' \sim P(\cdot | x, a)$$

# Reinforcement Learning

## Markov Decision Process

Agent

Environment

$$\pi(x) \rightarrow a$$

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

**H Steps**

Send **reward** and **next state** from a Markovian transition dynamics
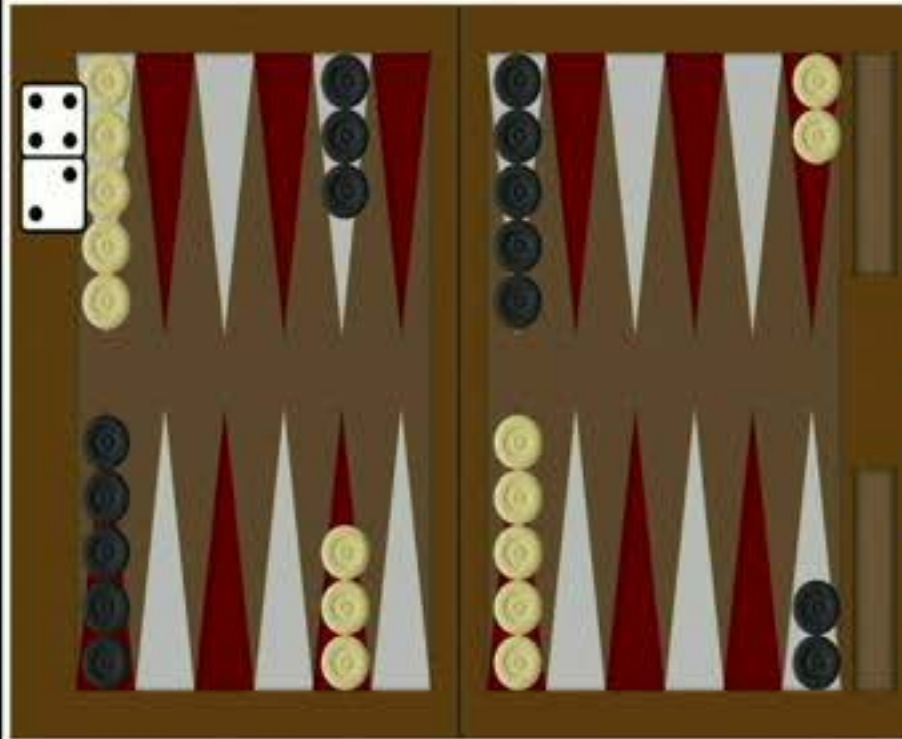
$$r(x, a), \quad x' \sim P(\cdot | x, a)$$

Maximize expected total reward:

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

# Progress of RL in Practice



TD GAMMON [Tesauro 95]

[AlphaZero, Silver et.al, 17]

[OpenAI Five, 18]

# Progress of RL in Practice

*OpenAI Five plays 180 years worth of games against itself every day….***running on 256 GPUs and 128,000 CPU cores**

— *Open AI Five Blog*



[OpenAI Five]

# Progress of RL in Practice



[OpenAI Five]

8

# Inefficient Exploration



Random Trial and error via
massive simulation
(i.e., **128,000** CPUs)

# Inefficient Exploration



Random Trial and error via
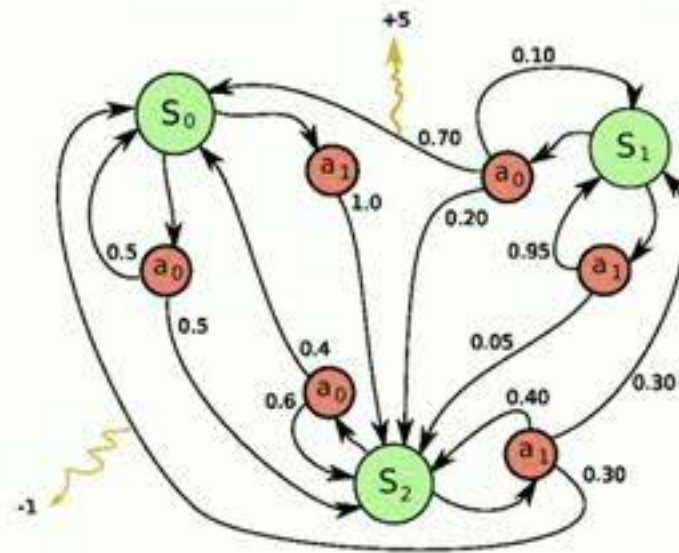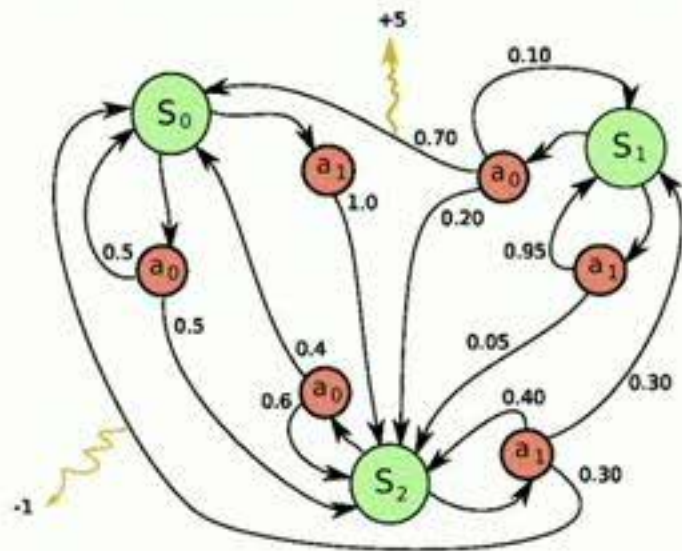massive simulation
(i.e., **128,000** CPUs)

9

# Inefficient Exploration



Random Trial and error via
massive simulation
(i.e., **128,000** CPUs)

$\neq$

**Sample Efficiency**

9

# Progress of RL in Theory



**Sample Efficiency in Small Discrete MDPs**

# Progress of RL in Theory
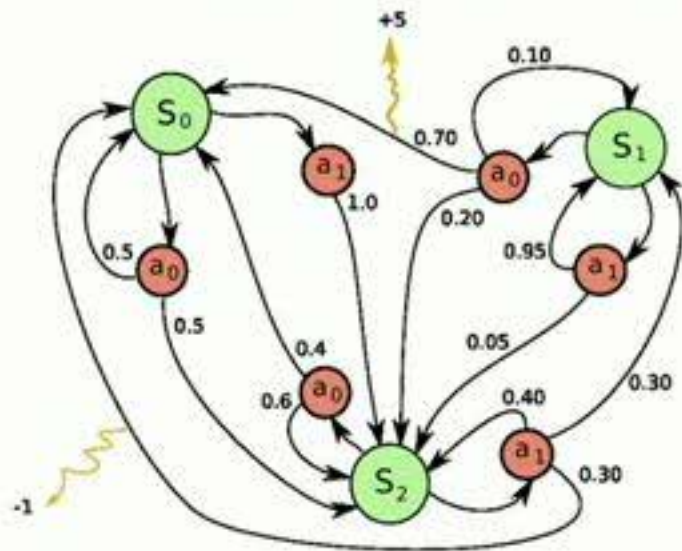
**Sample Efficiency
in Small Discrete MDPs**



**Sample Complexity:**

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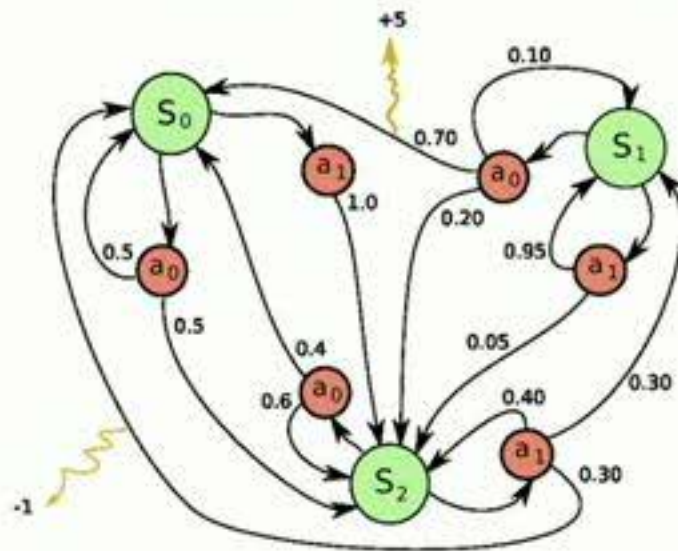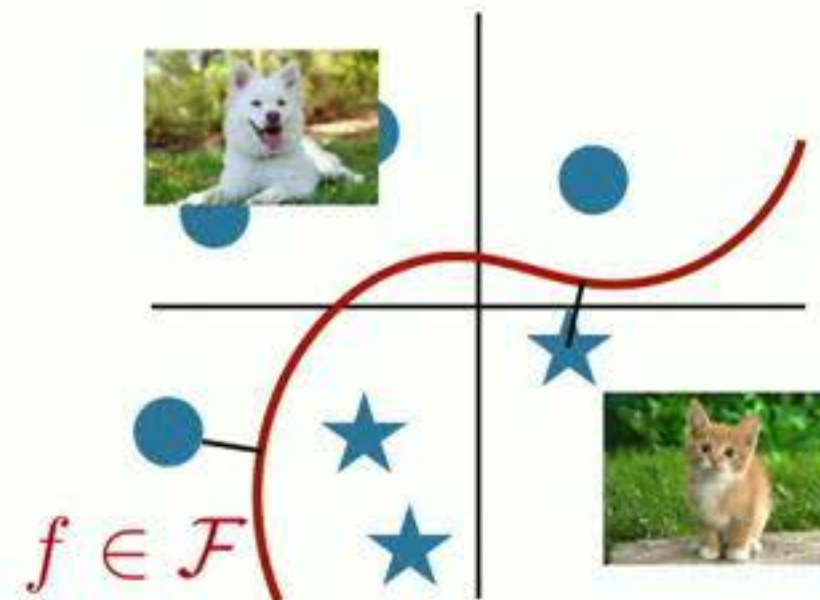need at most

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

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

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

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

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

many interactions

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

$\neq$

## Large-Scale Decision Making Problems

# What We Understand:
## Supervised Learning



[ImageNet]

$f \in \mathcal{F}$

Polynomial Dependency of # of unique images

# What We Understand:
## Supervised Learning
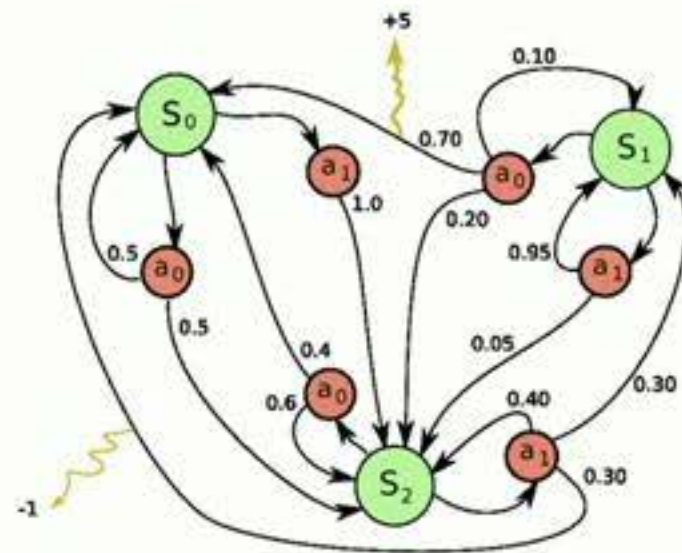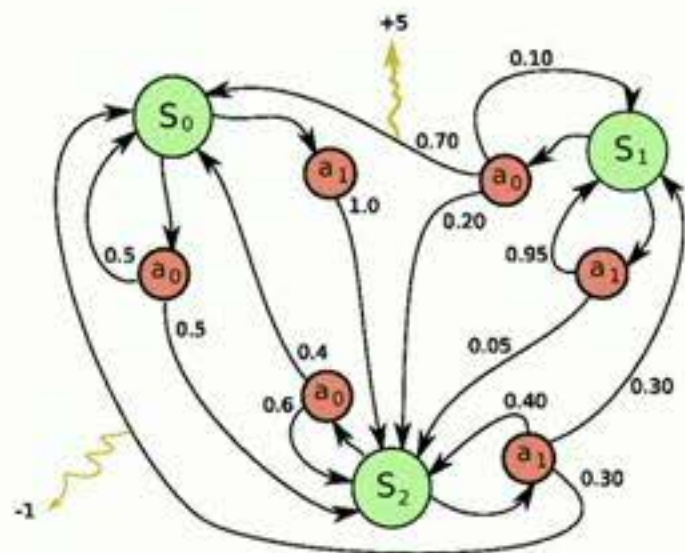


$f \in \mathcal{F}$

Polynomial Dependency of # of unique images

Generalization via Function Approximation

[ImageNet]

11

# What We Want:
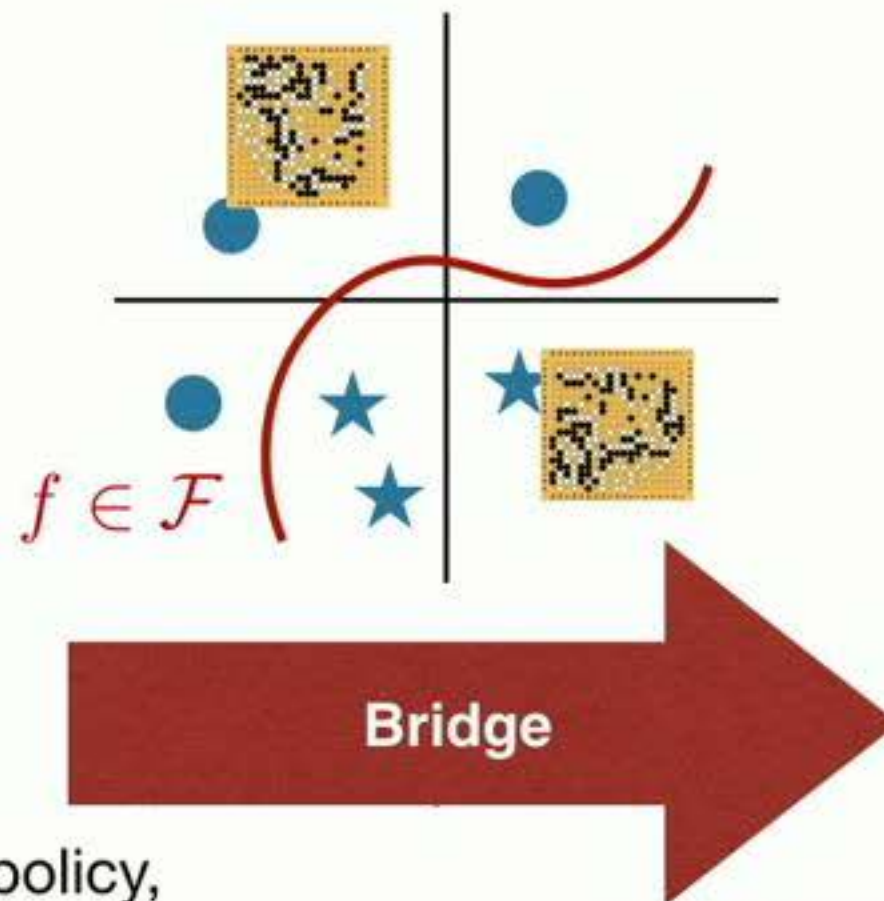## Generalization in Large-Scale MDPs

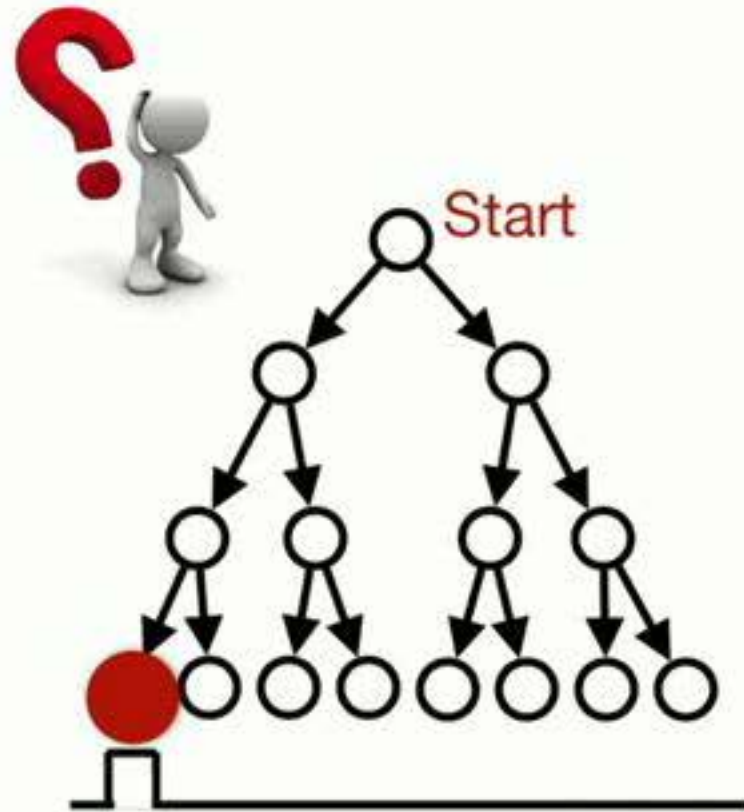**Sample Efficiency**



**Sample Complexity:**
To achieve $\epsilon$ near-optimal policy,
we 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]

$\neq$



12

# What We Want:
## Generalization in Large-Scale MDPs

**Sample Efficiency**



$f \in \mathcal{F}$

**Bridge**

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

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

many interactions

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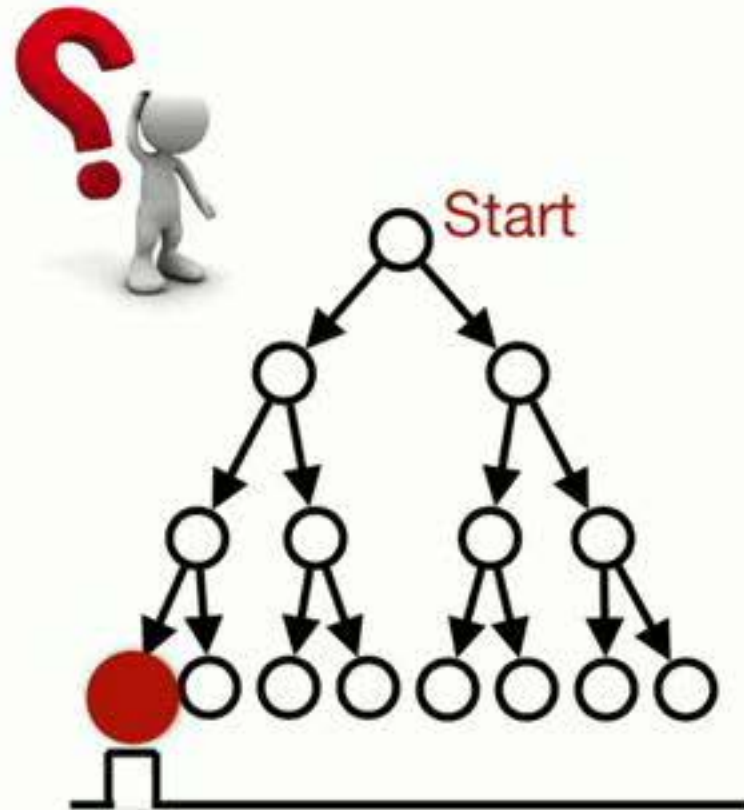
# BUT...



Start

Reward only at one leaf

[e.g., Krishnamurthy et.al 16, Jiang et.al 17]

## Needle in a haystack

**Discrete MDPs**

**H: horizon, S: # of states, A: # of actions**

# BUT...



Start

Reward only at one leaf
[e.g., Krishnamurthy et.al 16, Jiang et.al 17]

**Needle in a haystack**

**Discrete MDPs**

**H: horizon, S: # of states, A: # of actions**

# of Interactions
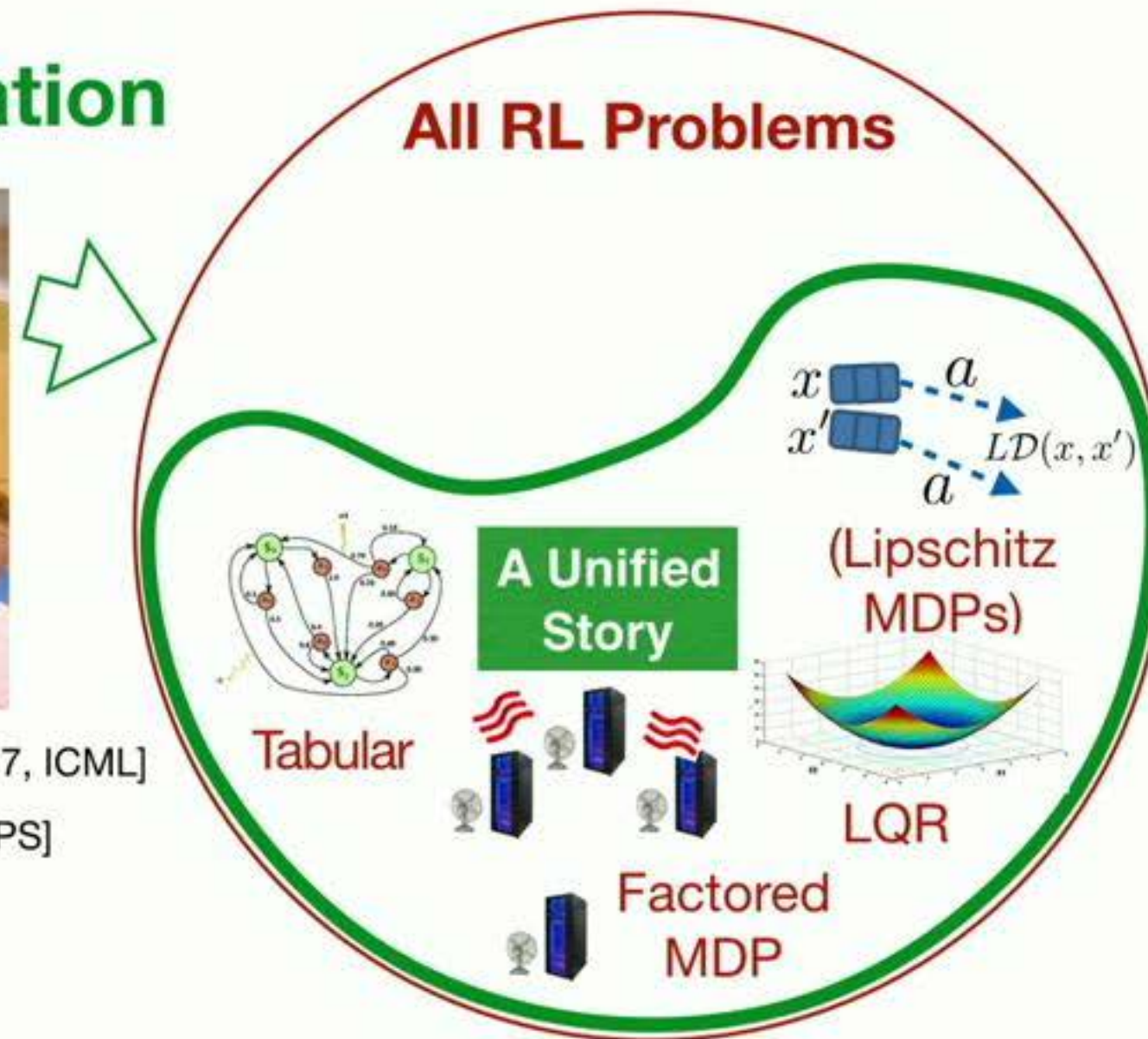with environment $>$ $\Omega(S)$

[e.g.,Dann & Brunskill, 15]

# Generalization & Sample Efficiency via...

## 1. Expert Demonstration

**All RL Problems**

[**Sun**, Venkatraman, Gordon, Boots, Bagnell, 17, ICML]

[**Sun**,, Gordon, Boots, Bagnell, 18, NeurIPS]

A Unified Story

Tabular

Factored MDP

(Lipschitz MDPs)

$LD(x, x')$

LQR

## 2. Exploiting Structures

[**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

# Generalization & Sample Efficiency via...

## 1. Expert Demonstration



**All RL Problems**

[**Sun**, Venkatraman, Gordon, Boots, Bagnell, 17, ICML]

[**Sun**,, Gordon, Boots, Bagnell, 18, NeurIPS]

- **Why** IL (i.e., IL VS RL)

- **How** to reduce RL to **Supervised Learning**

- **Generalize** from **Local** Experts

14

# Imitation Learning

Global Expert $\pi^*$ → Machine Learning → Policy $\pi$

Maps states to actions

- SVM
- Gaussian Process
- Deep Networks

Apprenticeship Learning [Abbeel & Ng 05, Syed & Schapire 08]

Inverse Optimal Control [Ziebart & Bagnell, 10]

Interactive Imitation Learning [Ross& Bangell, 11; Chang et.al., 15]

Generative Adversarial Imitation Learning [Ho & Ermon 16]

# Interactive Imitation Learning w/ Reward

**A global expert is available during training**

# Interactive Imitation Learning w/ Reward

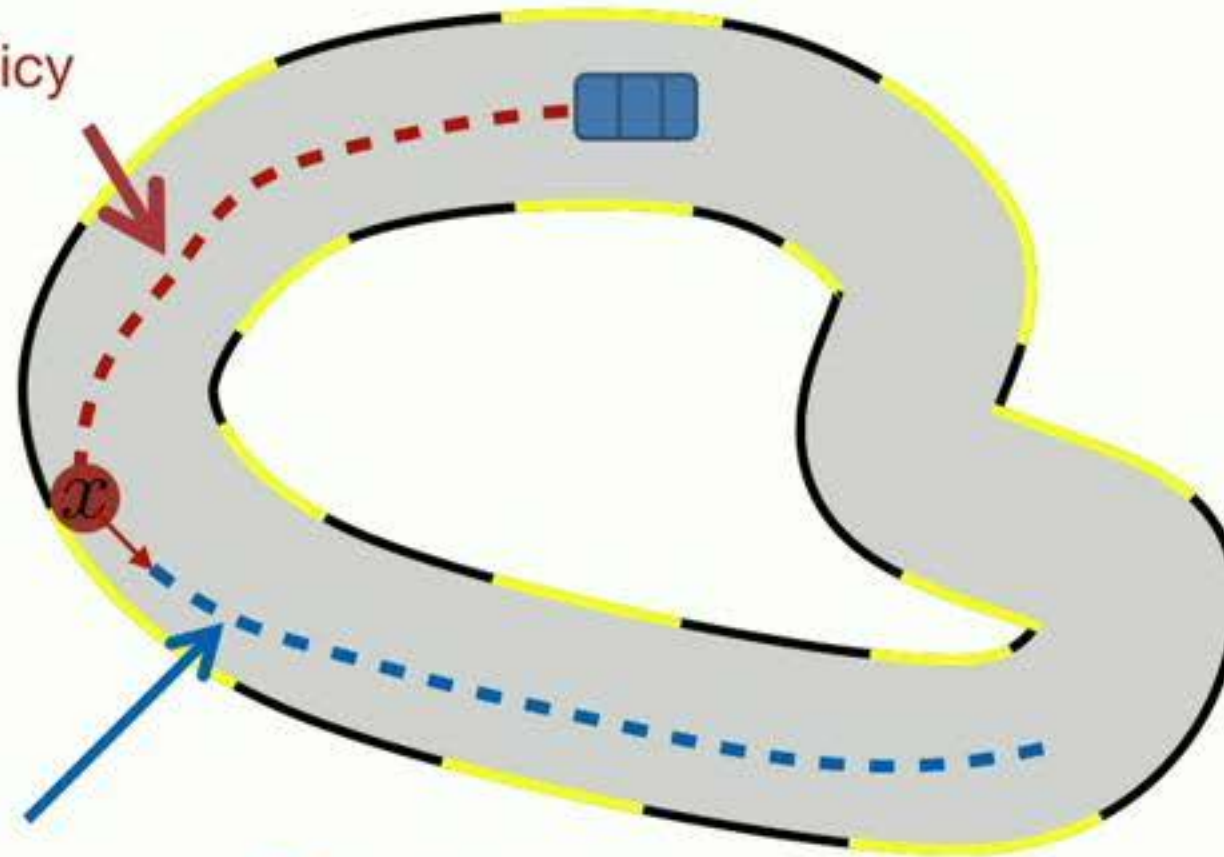**A global expert is available during training**

# Interactive Imitation Learning w/ Reward

**A global expert is available during training**

Execute Learned Policy

# Interactive Imitation Learning w/ Reward

**A global expert is available during training**
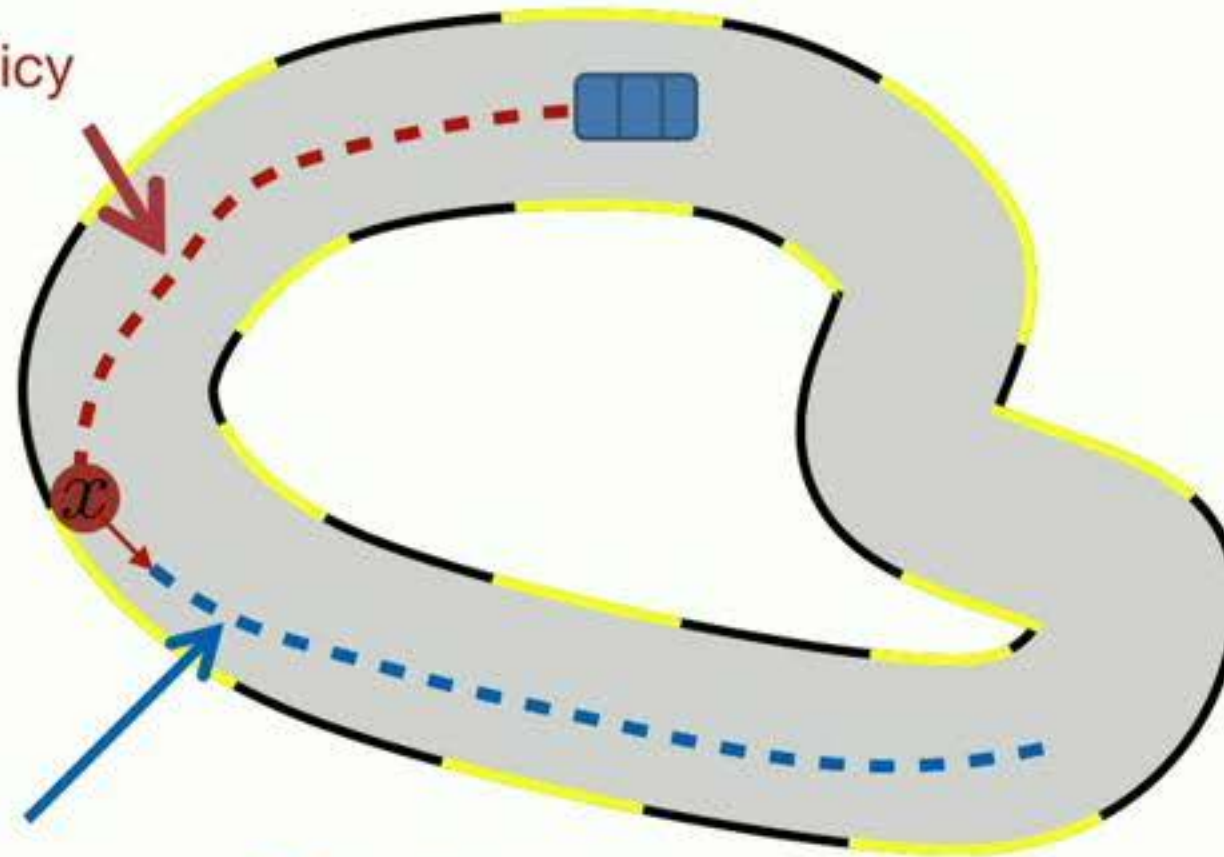
Execute Learned Policy

$x$

Ask a globally optimal Expert
to Take Over

# Interactive Imitation Learning w/ Reward

**A global expert is available during training**



Execute Learned Policy

$x$

Ask a globally optimal Expert
to Take Over

**Record**: Expert trajectory's total cost

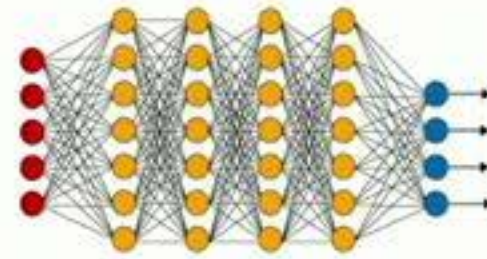How easy to recover from the learner's mistake

# Examples of Interactive Experts

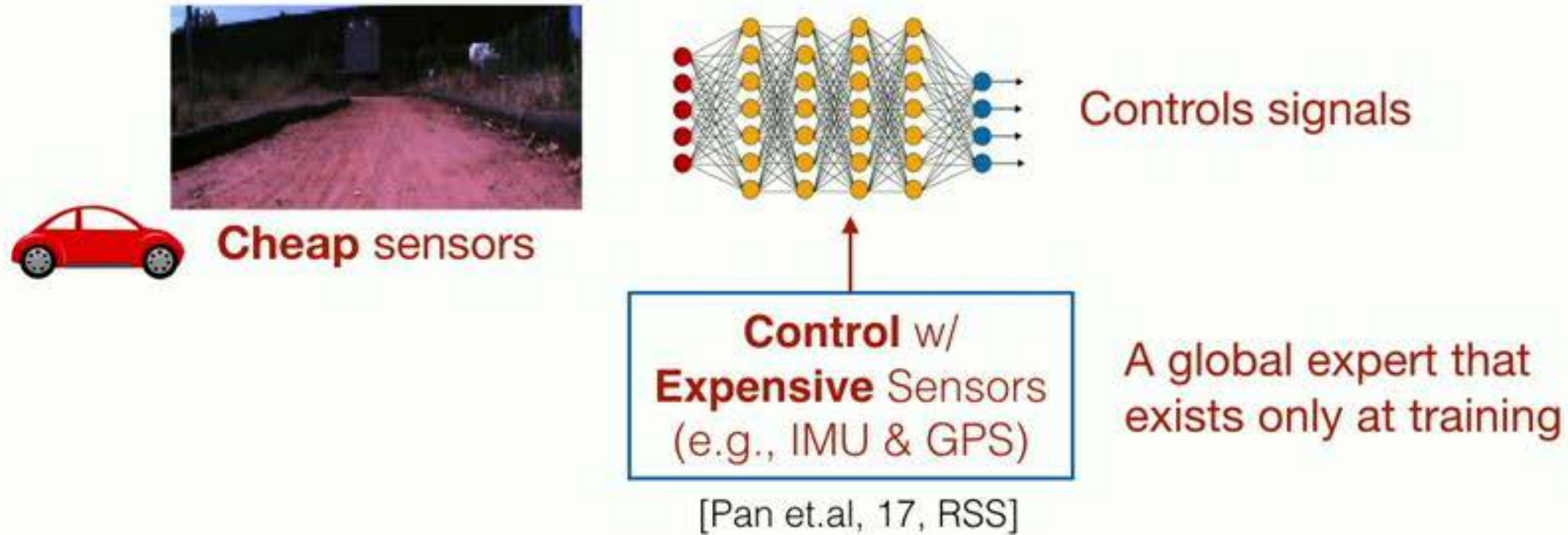# Examples of Interactive Experts

## 1. Planner/Control (e.g., Robotics)
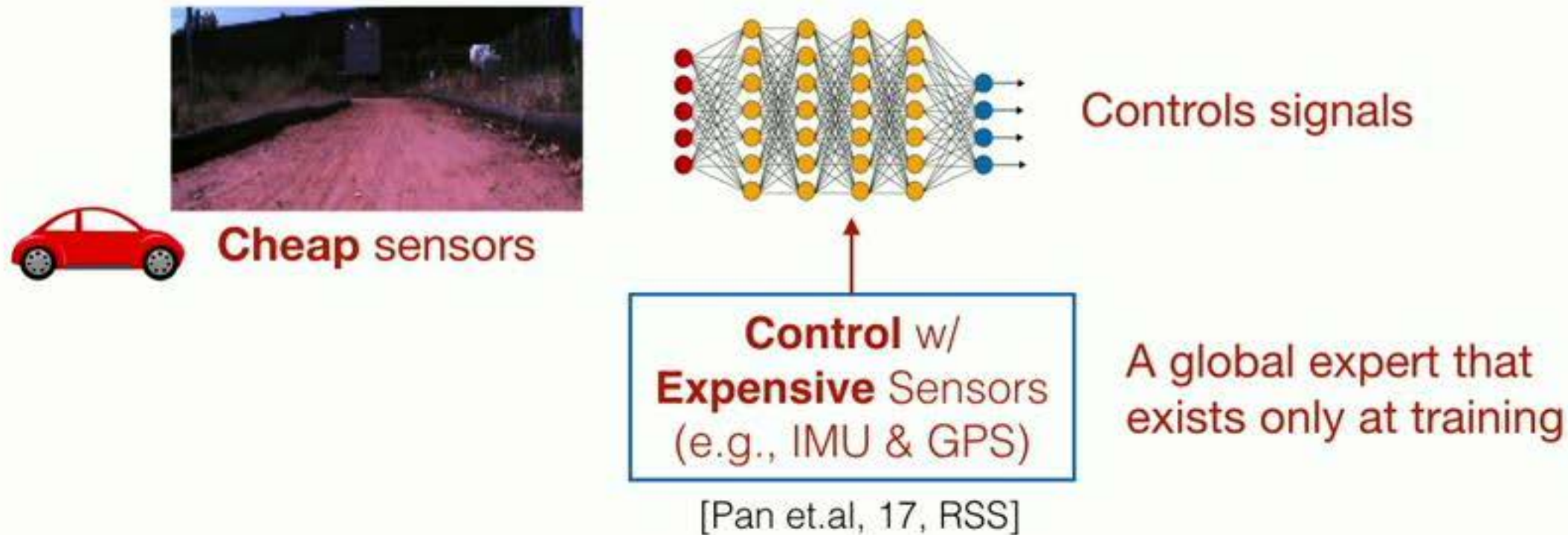


**Cheap** sensors

Controls signals

# Examples of Interactive Experts
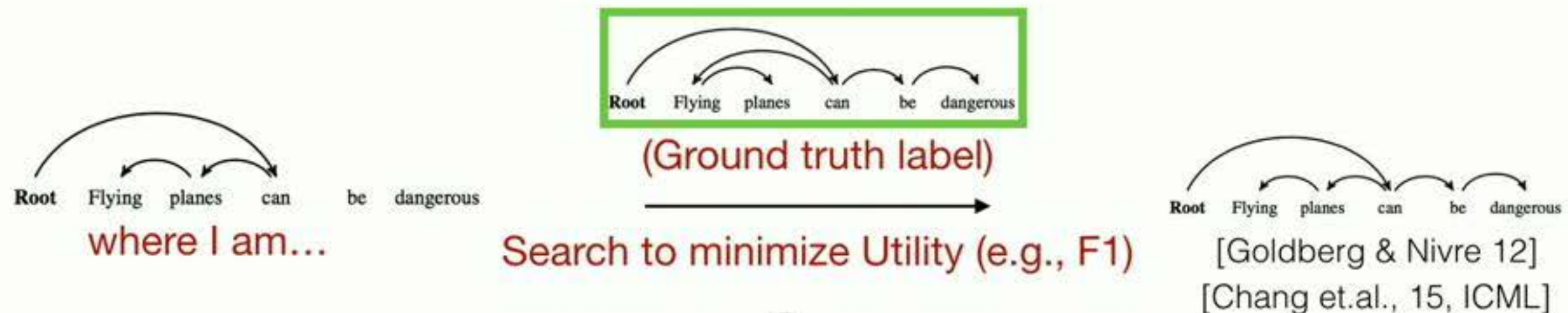
## 1. Planner/Control (e.g., Robotics)



**Cheap** sensors

Controls signals

**Control** w/
**Expensive** Sensors
(e.g., IMU & GPS)

A global expert that
exists only at training

[Pan et.al, 17, RSS]

# Examples of Interactive Experts

## 1. Planner/Control (e.g., Robotics)



**Cheap** sensors

Controls signals

**Control** w/ **Expensive** Sensors (e.g., IMU & GPS)

A global expert that exists only at training

[Pan et.al, 17, RSS]

## 2. Search Algorithms (e.g., NLP)



Root  Flying  planes  can  be  dangerous

(Ground truth label)

Root  Flying  planes  can  be  dangerous

where I am…

Search to minimize Utility (e.g., F1)

[Goldberg & Nivre 12]
[Chang et.al., 15, ICML]

17

Why bother imitating when you have reward signals?

# Why IL: Formalizing Advantages

## 1. Global Optimality

Global Optimal Expert: $\pi^\star$

AggreVaTe (Aggregate with Values) [Ross&Bagnell14]

$$J(\hat{\pi}) \approx J(\pi^\star)$$

# Why IL: Formalizing Advantages

## 1. Global Optimality

Global Optimal Expert: $\pi^{\star}$

AggreVaTe (Aggregate with Values) [Ross&Bagnell14]

$$J(\hat{\pi}) \approx J(\pi^{\star})$$

## 2. Sample Efficiency (i.e., Learns faster)

There exist MDPs, s.t. with global optimal expert, to learn near-optimal solution,
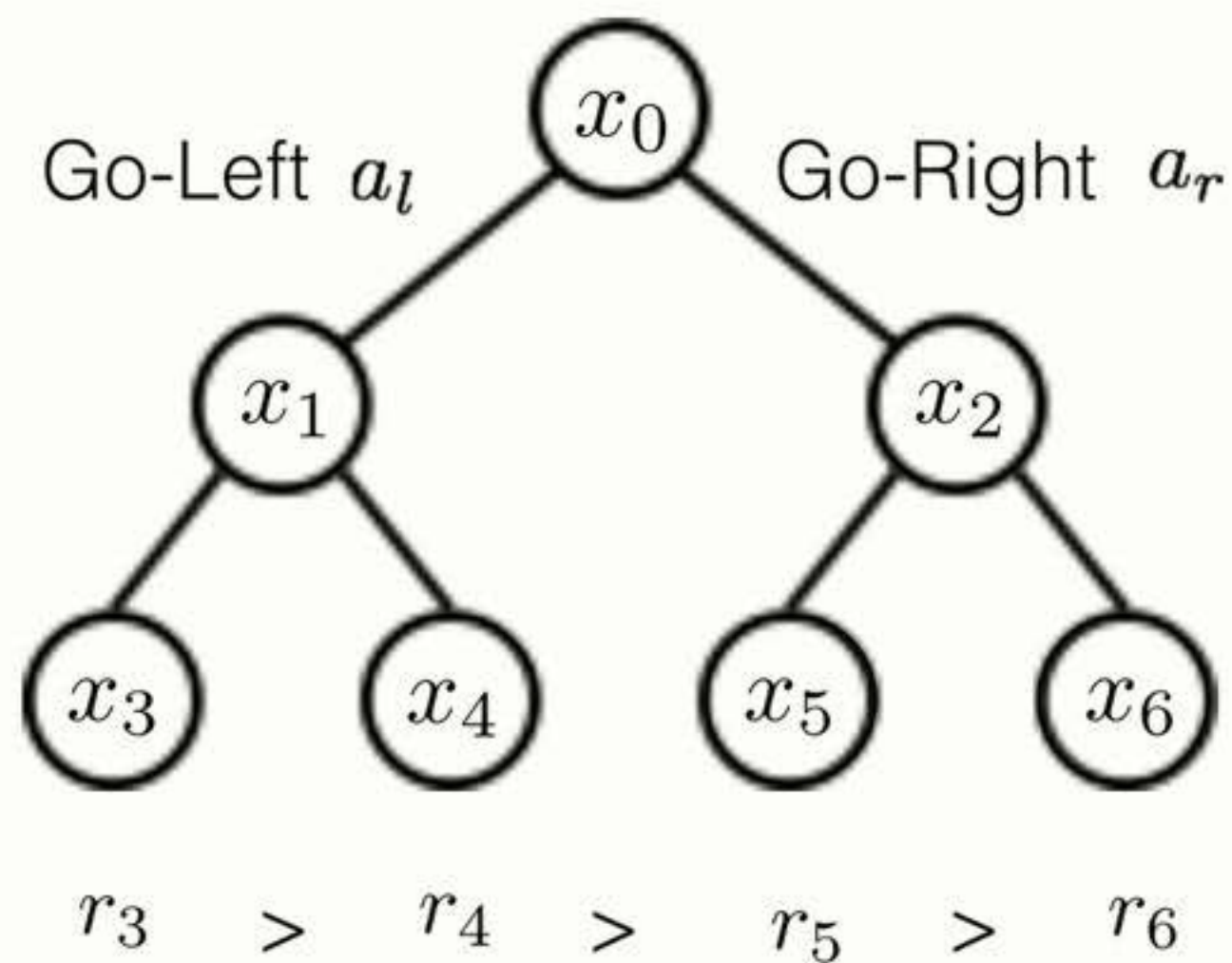
**IL** (e.g., AggreVaTe)          **ANY** RL

$$O(\log(S)) \quad \textbf{vs} \quad \Omega(S)$$

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
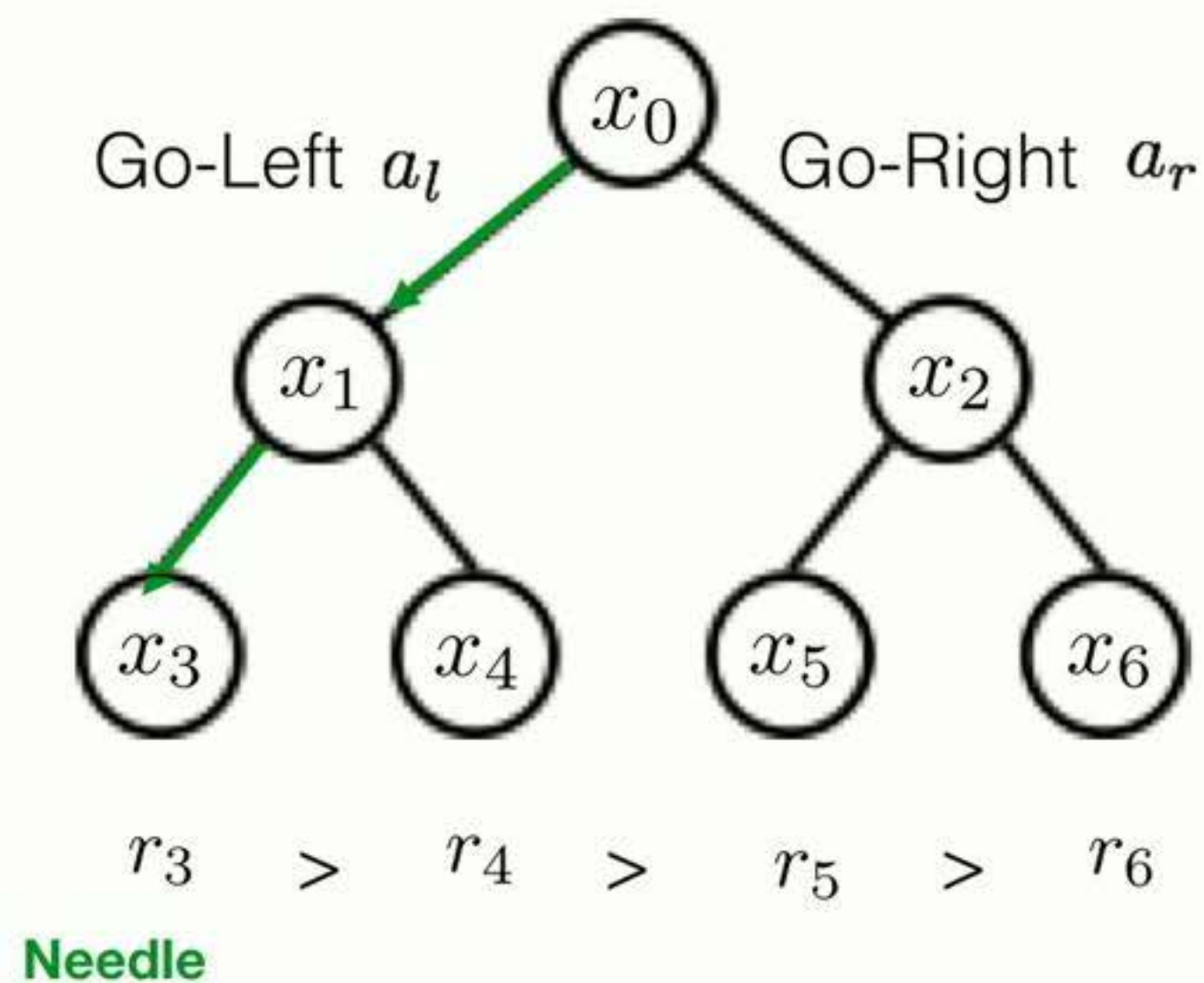**Sun**, Venkatraman, Gordon, Boots, Bagnell, ICML, 17

# Why IL: Formalizing Advantages

## 1. Global Optimality

Global Optimal Expert: $\pi^{\star}$

AggreVaTe (Aggregate with Values) [Ross&Bagnell14]

$$J(\hat{\pi}) \approx J(\pi^{\star})$$

## 2. Sample Efficiency (i.e., Learns faster)

There exist MDPs, s.t. with global optimal expert, to learn near-optimal solution,

**IL** (e.g., AggreVaTe)  **VS**  **ANY** RL

✓ $O(\log(S))$ $\quad$ $\Omega(S)$

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
**Sun**, Venkatraman, Gordon, Boots, Bagnell, ICML, 17

# Deterministic MDP



Global Optimal Expert: An Optimal Planner

Go-Left $a_l$ ... Go-Right $a_r$

$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$

# Deterministic MDP

Global Optimal Expert: An Optimal Planner



$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

**Needle**

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



Go left could lead to $r_3$

$$r_3 \; > \; r_4 \; > \; r_5 \; > \; r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



Go left could lead to $r_3$

$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



Go left could lead to $r_3$

$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



Go left could lead to $r_3$

Go left could lead to $r_3$

$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner



Go left could lead to $r_3$

Go left could lead to $r_3$

$$r_3 \quad > \quad r_4 \quad > \quad r_5 \quad > \quad r_6$$

Halving: Eliminate half of the nodes each round

[**Sun**,et.al., 17, ICML]

# Reduction to Supervised Learning
## Easy Credit Assignment

Global Optimal Expert: An Optimal Planner

Go left could lead to $r_3$

$x_0$

Go left could lead to $r_3$

$x_1$

$x_2$

$x_3$ $x_4$ $x_5$ $x_6$

$$r_3 > r_4 > r_5 > r_6$$

Halving: Eliminate half of the nodes each round

**IL:** $\log(S)$ **vs** **RL:** $\Omega(S)$

21

[**Sun**,et.al., 17, ICML]

# Ex: AggreVaTe [Ross & Bagnell, 14]

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step
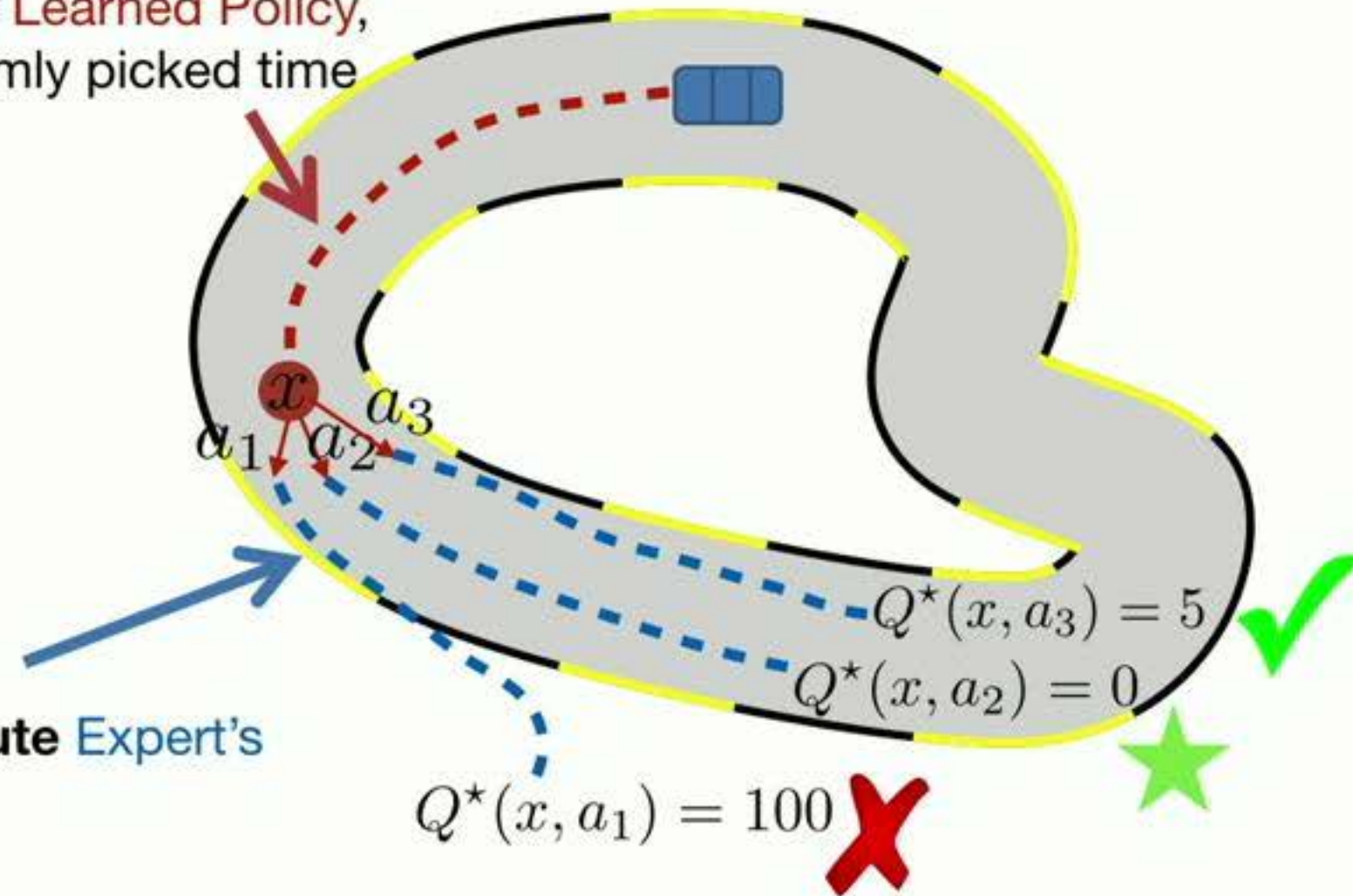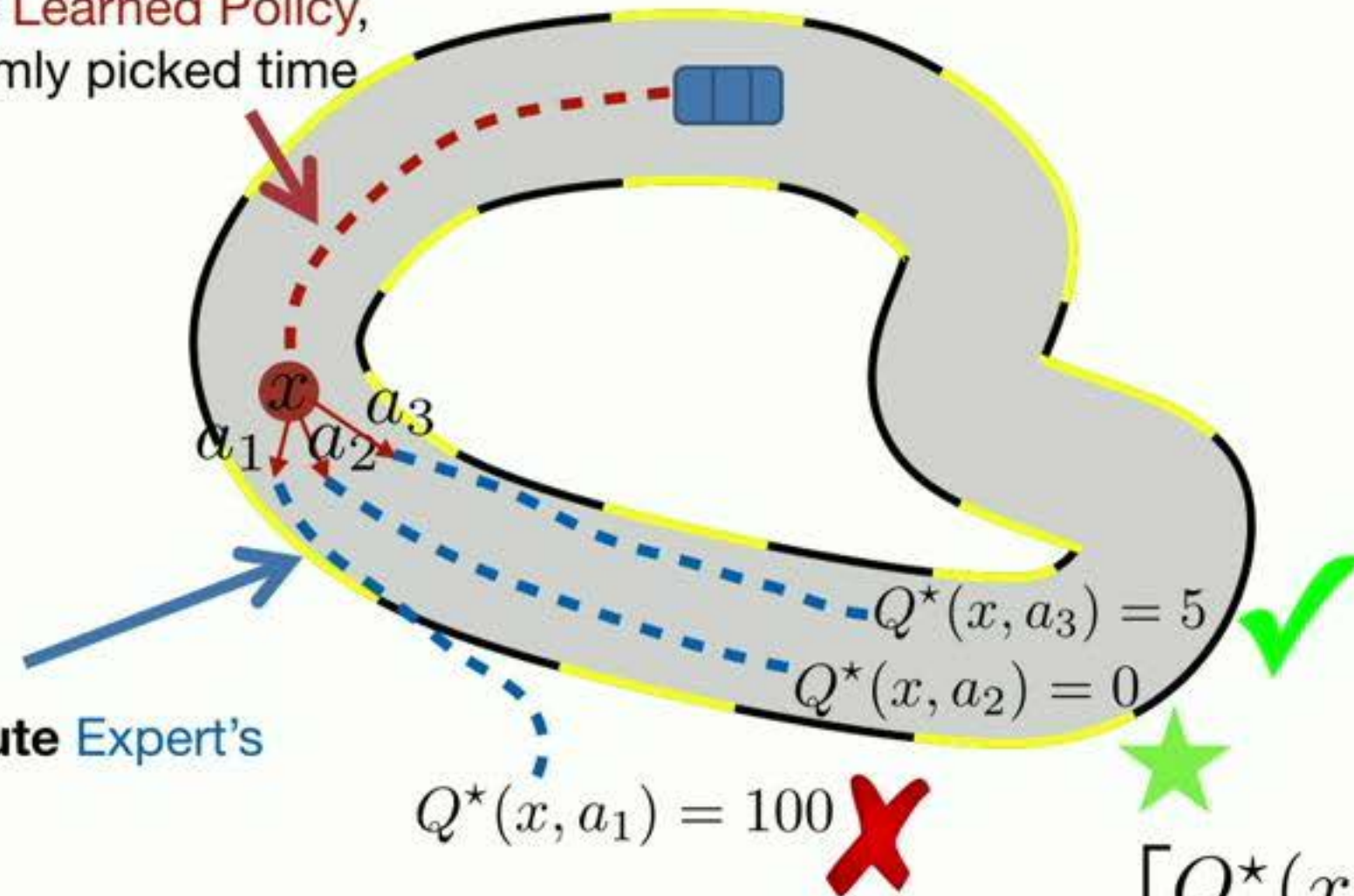
**Rollout: Execute** Expert's Policy

$$x$$

$$a_1$$

$$Q^\star(x, a_1) = 100$$

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step

**Rollout: Execute** Expert's Policy

$x$

$a_1$   $a_2$

$Q^\star(x, a_2) = 0$

$Q^\star(x, a_1) = 100$

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step

**Rollout: Execute** Expert's Policy

$Q^{\star}(x, a_3) = 5$

$Q^{\star}(x, a_2) = 0$

$Q^{\star}(x, a_1) = 100$

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step

$x$

$a_1$ $a_2$ $a_3$

$Q^\star(x, a_3) = 5$ ✔

$Q^\star(x, a_2) = 0$

**Rollout: Execute** Expert's Policy

$Q^\star(x, a_1) = 100$ ✘

# Ex: AggreVaTe [Ross & Bagnell, 14]

**Rollin: Execute** Learned Policy, Stop at a randomly picked time step

**Rollout: Execute** Expert's Policy

$Q^\star(x, a_3) = 5$

$Q^\star(x, a_2) = 0$

$Q^\star(x, a_1) = 100$

$$\left\{ x, \begin{bmatrix} Q^\star(x, a_1) \\ Q^\star(x, a_2) \\ Q^\star(x, a_3) \end{bmatrix} \right\}_N$$

Cost-Sensitive Classification Dataset
(A **Supervised Learning** Dataset)

22

# Differentiable AggreVaTe (AggreVaTeD)

[**Sun**,et.al., 17, ICML]



$$\pi_{\theta_n} \text{ (e.g., neural net)}$$

$$\Rightarrow \left\{ x, \begin{bmatrix} Q^{\star}(x, a_1) \\ Q^{\star}(x, a_2) \\ Q^{\star}(x, a_3) \end{bmatrix} \right\}_N$$

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
[**Sun**, Venketraman, Gordon, Boots, Bagnell, ICML, 17]

# Differentiable AggreVaTe (AggreVaTeD)

$$\pi_{\theta_n} \text{ (e.g., neural net)}$$

$$\left\{ x, \begin{bmatrix} Q^\star(x, a_1) \\ Q^\star(x, a_2) \\ Q^\star(x, a_3) \end{bmatrix} \right\}_N$$

$$\nabla_\theta \ell_n(\theta)|_{\theta_n} \quad \Leftarrow \quad \ell_n(\theta) = \sum_i \sum_a \pi(a|x_i; \theta) Q^\star(x_i, a)$$

Cost-Sensitive loss

23

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
[**Sun**, Venketraman, Gordon, Boots, Bagnell, ICML, 17]

# Differentiable AggreVaTe (AggreVaTeD)

[**Sun**,et.al., 17, ICML]



$\pi_{\theta_n}$ (e.g., neural net)

$$\left\{ x, \begin{bmatrix} Q^\star(x, a_1) \\ Q^\star(x, a_2) \\ Q^\star(x, a_3) \end{bmatrix} \right\}_N$$

$\pi_{\theta_{n+1}}$

$\nabla_\theta \ell_n(\theta)|_{\theta_n}$

e.g., Gradient or natural gradient descent

$$\ell_n(\theta) = \sum_i \sum_a \pi(a|x_i; \theta) Q^\star(x_i, a)$$

Cost-Sensitive loss

23

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
[**Sun**, Venketraman, Gordon, Boots, Bagnell, ICML, 17]

# Differentiable AggreVaTe
## (AggreVaTe**D**)

$\pi_{\theta_n}$ (e.g., neural net)

$\pi$

$a$

$\pi^*$

$\pi^*$

$\pi^*$

$$\left\{ x, \begin{bmatrix} Q^\star(x, a_1) \\ Q^\star(x, a_2) \\ Q^\star(x, a_3) \end{bmatrix} \right\}_N$$

$\pi_{\theta_{n+1}}$

$\nabla_\theta \ell_n(\theta)|_{\theta_n}$

$$\ell_n(\theta) = \sum_i \sum_a \pi(a|x_i; \theta) Q^\star(x_i, a)$$

e.g., Gradient or
natural gradient
descent

Cost-Sensitive loss

**Use rich function
approximators for complex
features**

23

Deeply AggreVaTeD: Differential Imitation Learning for Sequential Prediction
[**Sun**, Venketraman, Gordon, Boots, Bagnell, ICML, 17]

# Dependency Parsing

## Handwritten Algebra Equations & Solutions

[Duyck & Gordon 15]

Input:

$$-5(x-1) = -20$$
$$x - 1 = 4$$
$$x = 5$$

Output:

# Dependency Parsing as Sequential Decision Making

[e.g., Chang, Krishnamurthy, Agarwal, Daume´ III, Langford, 15, ICML]

# Dependency Parsing as Sequential Decision Making

[e.g., Chang, Krishnamurthy, Agarwal, Daume´ III, Langford, 15, ICML]

Encoder (LSTM)



$$-5(x-1) = -20$$
$$x-1 = 4$$
$$x = 5$$

# Dependency Parsing as Sequential Decision Making

[e.g., Chang, Krishnamurthy, Agarwal, Daume´ III, Langford, 15, ICML]

Encoder (LSTM)

# Dependency Parsing as Sequential Decision Making

[e.g., Chang, Krishnamurthy, Agarwal, Daume' III, Langford, 15, ICML]

# What if we *do not* have a Globally Optimal Expert?

### ...we can learn from Local Experts!

# Example: AlphaGo-Zero

[Silver, et.al, 17, Nature]

## Known & Deterministic Transition Dynamics

Fast
Reactive
Policy $\pi$



Slow
Policy $\eta$
(MCTS)

**AlphaZero** leverages transition dynamics to build local experts

# What if we *do not* have any prior knowledge of transition dynamics?

# Dual Policy Iteration
[**Sun** et.al., 18, NeurIPS]
## Imitating a Locally Optimal Control



Current Policy

**New Transitions**

State    Action    Next State

# Dual Policy Iteration
[**Sun** et.al., 18, NeurIPS]

## Imitating a Locally Optimal Control

Current Policy

**New Transitions**

State  Action  Next State

...

$$x' \sim \hat{P}(\cdot|x, a)$$

**Supervised Learning**
Dynamics

# Dual Policy Iteration
## Imitating a Locally Optimal Control

[Sun et.al., 18, NeurIPS]



Current Policy

New Transitions

State   Action   Next State

Local Control

$$x' \sim \hat{P}(\cdot|x,a)$$

Supervised Learning
Dynamics

30

# Dual Policy Iteration
[Sun et.al., 18, NeurIPS]
## Imitating a Locally Optimal Control



Current Policy

**New Transitions**

State   Action   Next State

...

$x' \sim \hat{P}(\cdot|x,a)$

**Imitation**
e.g., AggreVaTe

**Local Control**

**Supervised Learning**
Dynamics

Adapt towards...

30

# Dual Policy Iteration

[**Sun** et.al., 18, NeurIPS]

## Imitating a Locally Optimal Control

Current Policy

New Transitions

State   Action   Next State

Policy Improvement

Convergence

Imitation
e.g., AggreVaTe

Local Control

$x' \sim \hat{P}(\cdot|x,a)$

Adapt towards...

Supervised Learning
Dynamics

30

# Helicopter Funnel

[**Sun** et.al., 18, NeurIPS]

**Instantiation 1:**
Linear Regressors + iLQR + AggreVaTeD w/ Natural Gradient

iLQR: [Li & Todorov, 05]  AggreVaTeD: [Sun, 17, ICML]

# Helicopter Funnel

**Learned Policy from DPI**
(Simulator from Abbeel et.al, 06)

**Instantiation 1:**
Linear Regressors + iLQR + AggreVaTeD w/ Natural Gradient

31

iLQR: [Li & Todorov, 05]   AggreVaTeD: [Sun, 17, ICML]

# Synthetic Discrete MDPs

[**Sun** et.al., 18, NeurIPS]

Randomly Generated
Discrete MDPs

[Archibald et.al., 95]

Conservative Policy Iteration
[Kakade & Langford, 02]

discrete_MDP_1000

Lower equal to better (log-scale)

# of Iterations

**Instantiation 2:**
Maximum Likelihood Estimation + Value Iteration + AggreVaTeD

AggreVaTeD: [Sun, 17, ICML]

32

# Generalization & Sample Efficiency via...



All RL Problems

# Generalization & Sample Efficiency via...

- Why Model-Based RL?
- A Unified Measure

**All RL Problems**

A Unified Story

(Lipschitz MDPs)

Tabular

Factored MDP

LQR

## 2. Exploiting Structures

[**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

33

**Modeling Dynamics**

**Known**

[Sun et.al, ISRR 13]

**Control**

e.g., iterative LQR

[Li & Todorov 03]

# Modeling Dynamics

## Known



[Sun et.al, ISRR 13]

## Learned



[Williams et.al, 17, ICRA]

## Control

e.g., iterative LQR

[Li & Todorov 03]

# Modeling Dynamics

**Known**  **Learned**

[Sun et.al, ISRR 13]  [Williams et.al, 17, ICRA]

**Control**  **Model-Based RL**

e.g., iterative LQR  $\hat{P}(\cdot|x,a) \approx P^{\star}(\cdot|x,a)$

[Li & Todorov 03]  Approximator   Real Transition

# Modeling Dynamics

**Known**



[Sun et.al, ISRR 13]

**Control**
e.g., iterative LQR
[Li & Todorov 03]

**Learned**



[Williams et.al, 17, ICRA]

**Model-Based RL**

$$\hat{P}(\cdot|x, a) \approx P^\star(\cdot|x, a)$$

Approximator    Real Transition

**Ignored**



Directly learn policy
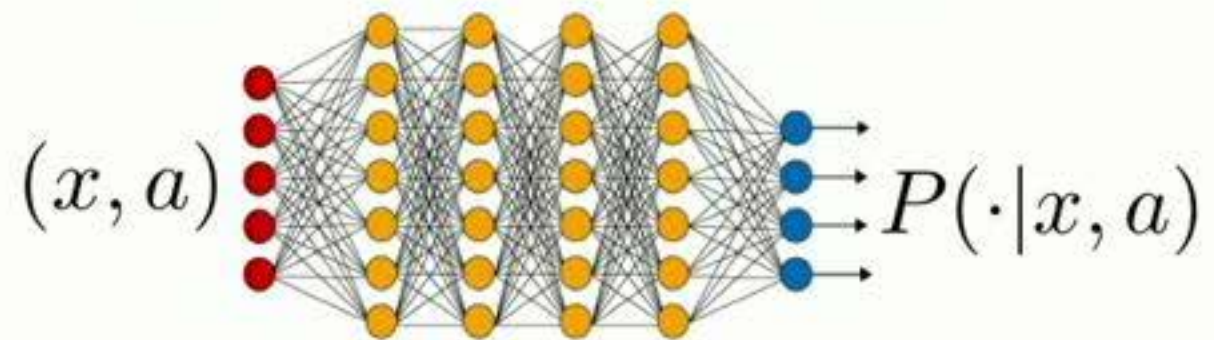**Model-Free RL**
e.g., Q-Learning
[Watkins & Dayan, 92]

**Modeling Dynamics**

**Known**      **Learned**      **Ignored**

[Sun et.al, ISRR 13]      [Williams et.al, 17, ICRA]

**Control**      **Model-Based RL**      Directly learn policy
e.g., iterative LQR      $$\hat{P}(\cdot|x,a) \approx P^\star(\cdot|x,a)$$      **Model-Free RL**
[Li & Todorov 03]      Approximator      Real Transition      e.g., Q-Learning

[Watkins & Dayan, 92]

# Setup of Model-Based RL



$$x' \sim P^{\star}(\cdot | x, a)$$

Real Transition Dynamics

# Setup of Model-Based RL

**Function Approximators**

$$(x, a) \quad P(\cdot | x, a)$$

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

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

**Real Transition Dynamics**
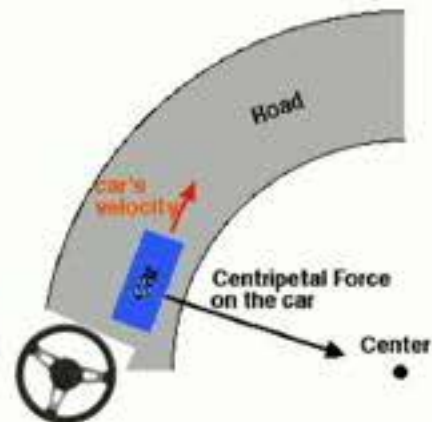
# Setup of Model-Based RL

**Function Approximators**



$$(x, a) \quad P(\cdot | x, a)$$

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

**Realizability:** $P^\star \in \mathcal{P}$

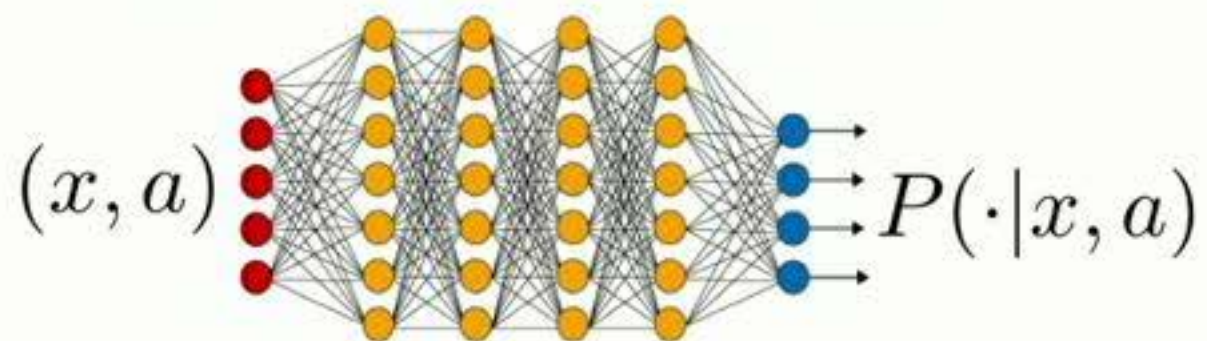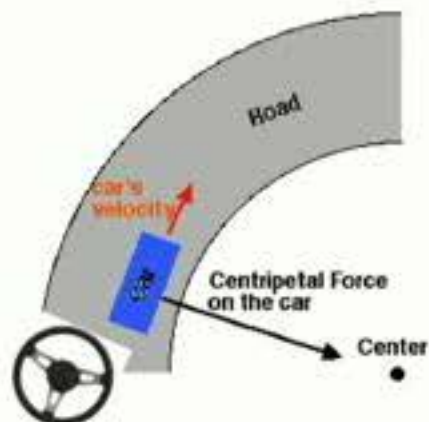$$x' \sim P^\star(\cdot | x, a)$$

Real Transition Dynamics

# Setup of Model-Based RL

**Function Approximators**

$(x, a)$  $P(\cdot | x, a)$

**Optimal Planner (OP)**

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

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

**Realizability:** $P^\star \in \mathcal{P}$

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

Real Transition Dynamics



35

# Setup of Model-Based RL

**Function Approximators**



$$(x, a) \qquad P(\cdot | x, a)$$

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

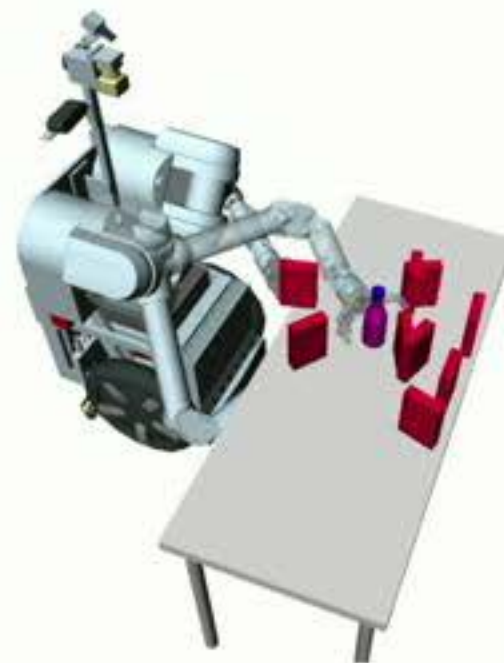**Realizability:** $P^\star \in \mathcal{P}$

$$x' \sim P^\star(\cdot | x, a)$$
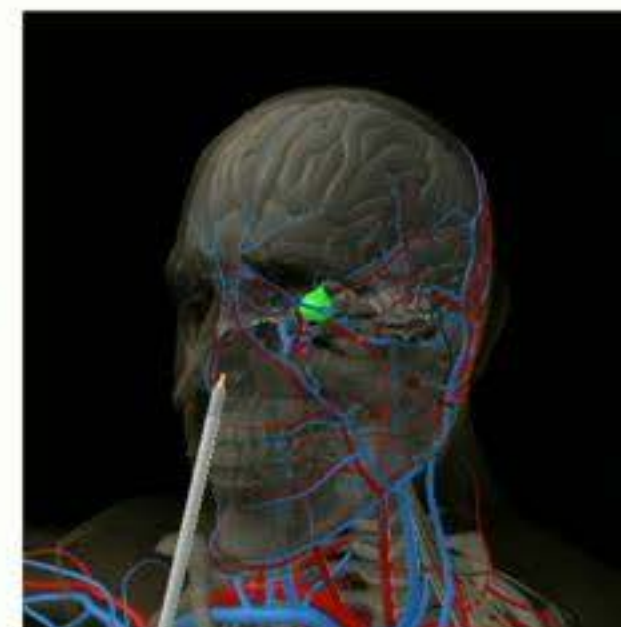
**Real Transition Dynamics**

**Optimal Planner (OP)**

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



[Zucker et.al, IJRR 13]          [Sun et.al, ISRR 13]

e.g., iLQR [Li & Todorov 03]

CHOMP [Ratliff et.al, 09]

SE-LQR [Sun et.al, 16, TASE]

35

# Why Model-Based RL?
## Debate: Model-Based or Model-Free

Iterative Learning Control
(e.g., An & Atkeson & Hollerbach 88, Abbeel 06)

Nonparametric Model-based RL
(e.g.,  Atkeson 98, Deisenroth et.al., 11)

Guided Policy Search
(e.g.,Levine & Abbeel 16)

**Dual Policy Iteration**
[Sun et.al, 18]

…

# Why Model-Based RL?
## Debate: Model-Based or Model-Free

**Model-Based** is often more sample **efficient**
than Model-Free **in practice**…

Iterative Learning Control
(e.g., An & Atkeson & Hollerbach 88, Abbeel 06)

Nonparametric Model-based RL
(e.g., Atkeson 98, Deisenroth et.al., 11)

Guided Policy Search
(e.g.,Levine & Abbeel 16)

**Dual Policy Iteration**
[**Sun** et.al, 18]

…

# In Theory?

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

Model-Based RL:

Polynomial Sample Complexity

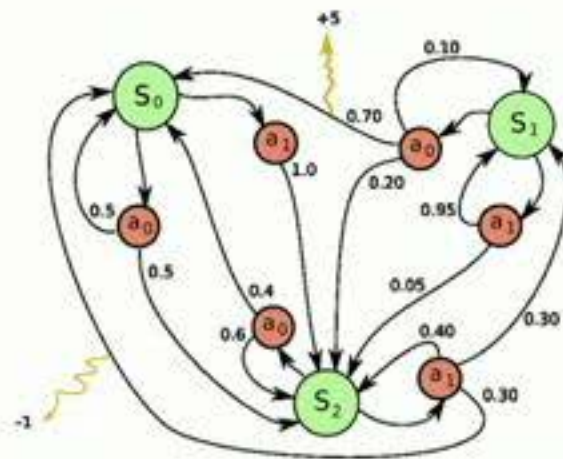**VS**

Any Model-Free RL:

$$\Omega(\exp(H))$$

Model-based Reinforcement Learning in Contextual Decision Processes
**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18.

# We have been exploiting the structures of models, BUT...



$x$

$x'$

$a$

$LD(x, x')$

$a$

**Lipschitz Continuous MDPs**

[Kearn, Langford, Kakade, 03]

**Small Tabular MDP**

[Kearn & Singh, 02]

**Linear Quadratic Regulator (LQR)**

[Dean et.al, 18]

**Factored MDPs**

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

38

# We have been exploiting the structures of models, BUT...



$LD(x, x')$

## Lipschitz Continuous MDPs

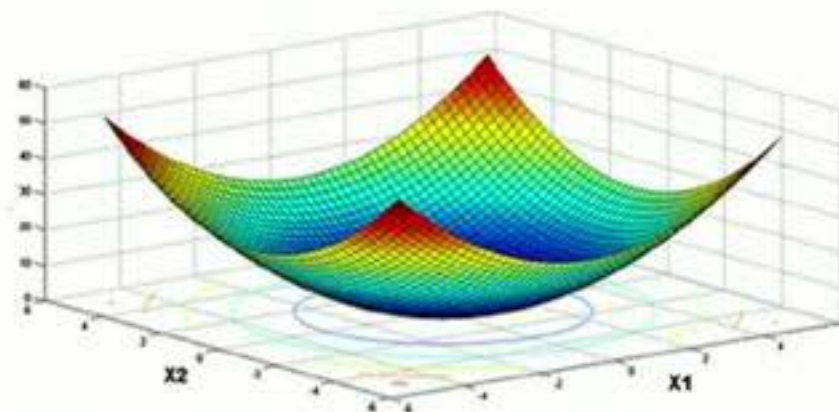[Kearn, Langford, Kakade, 03]

**A Unified Algorithm?**

## Small Tabular MDP

[Kearn & Singh, 02]

## Factored MDPs

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

## Linear Quadratic Regulator (LQR)

[Dean et.al, 18]

# Distinguish Two Distributions:

## Integral Probability Metric (IPM) [Muller et.al, 97]

Distinguish two distributions $P, Q$

# Distinguish Two Distributions:

## Integral Probability Metric (IPM) [Muller et.al, 97]

Distinguish two distributions $P, Q$



**Real** bedroom images
[LSUN dataset]

**Imaginary** samples from
a generative model
[e.g., Wasserstein GAN, 17]

# Distinguish Two Distributions:

## Integral Probability Metric (IPM) [Muller et.al, 97]

Distinguish two distributions $P, Q$



**Real** bedroom images
[LSUN dataset]

**Imaginary** samples from
a generative model
[e.g., Wasserstein GAN,17]

Discriminators $\max_{f \in \mathcal{F}} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$

# Distinguish Two Distributions:
## Integral Probability Metric (IPM) [Muller et.al, 97]

Distinguish two distributions $P, Q$



**Real** bedroom images
[LSUN dataset]

**Imaginary** samples from
a generative model
[e.g., Wasserstein GAN,17]

Discriminators $\max_{f \in \mathcal{F}} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$

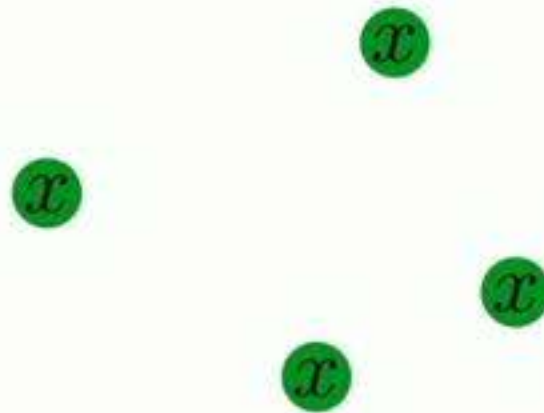$\mathcal{F} \triangleq \{f : \|f\|_\infty \leq 1\} \Rightarrow \|P - Q\|_1$ Total Variation

$\mathcal{F} \triangleq \{f : \|f\|_L \leq 1\} \Rightarrow$ Wasserstein Distance

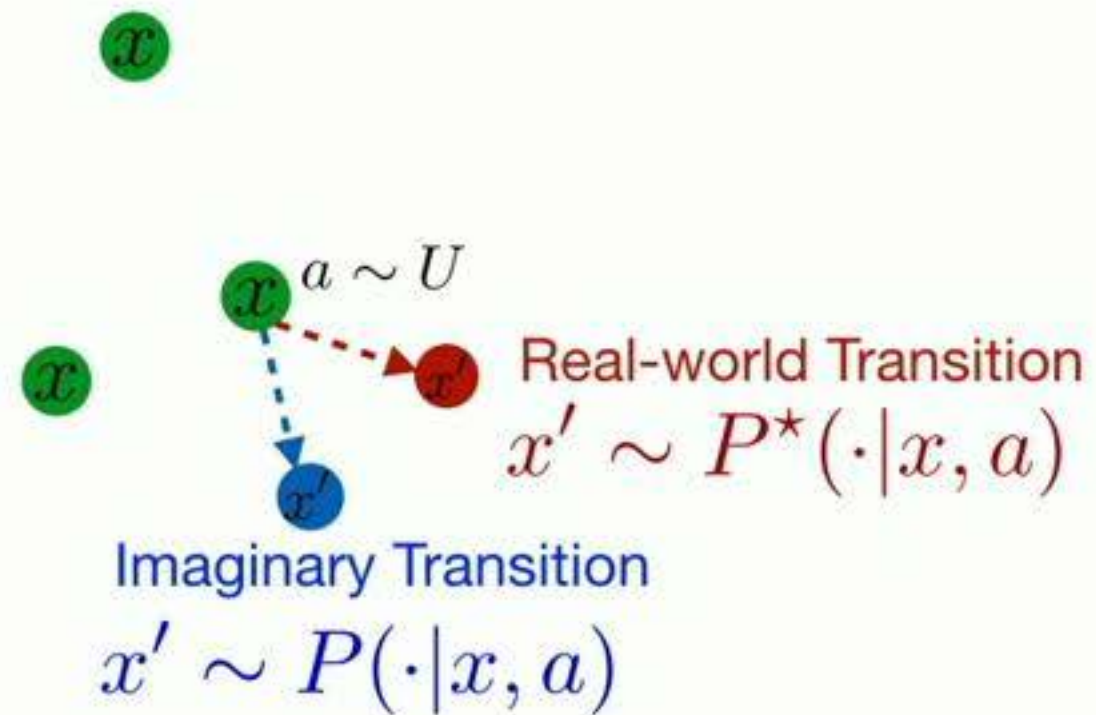# Distinguish a Candidate from the Real

Candiate: $P(\cdot|x,a)$ $\overset{\mathbf{?}}{=}$ Real: $P^{\star}(\cdot|x,a)$

# Distinguish a Candidate from the Real

Candiate: $P(\cdot|x,a)$ $\overset{?}{=\!=}$ Real: $P^{\star}(\cdot|x,a)$

# Distinguish a Candidate from the Real

Candiate: $P(\cdot|x,a)$ **?** Real: $P^{\star}(\cdot|x,a)$



Real-world Transition
$x' \sim P^{\star}(\cdot|x,a)$

Imaginary Transition
$x' \sim P(\cdot|x,a)$

$a \sim U$

# Distinguish a Candidate from the Real

Candiate: $P(\cdot|x,a)$ $\overset{=}{=}\mathbf{?}$ Real: $P^\star(\cdot|x,a)$



Real-world Transition
$$x' \sim P^\star(\cdot|x,a)$$

Imaginary Transition
$$x' \sim P(\cdot|x,a)$$

$$\max_{f \in \mathcal{F}} \mathbb{E}_{x \sim \pi, a \sim U}\left[\mathbb{E}_{x' \sim P} f(x,a,x') - \mathbb{E}_{x' \sim P^\star} f(x,a,x')\right]$$

# Model Rank

**Misfit Matrix:**

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

# Model Rank

**Misfit Matrix:**



$$W(P_r, P_c; \mathcal{F}) \quad \in \mathbb{R}^{|\mathcal{P}| \times |\mathcal{P}|}$$

# Model 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**

$$\mathbb{E}_{x \sim \pi_{P_r}, a \sim U}\left[\mathbb{E}_{x' \sim P_c} f(x, a, x') - \mathbb{E}_{x' \sim P^\star} f(x, a, x')\right]$$

Imaginary        Real-world

# Model 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**

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

Imaginary          Real-world

**Model Rank** is defined as the rank of this misfit matrix

41

# Model Rank

**Misfit Matrix:**

Provides **Conditional** State-action Distribution

$P_r$

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

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

"Similar" ← Low Rank

$P_c$ **Candidate**

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

Imaginary      Real-world

**Model Rank** is defined as the rank of this misfit matrix

# A Unified Framework



$$LD(x, x')$$

Lipschitz Continuous MDPs
**Rank <= Covering number
of state space**
[KLK, 03]

# A Unified Framework



$x$ $a$

$x'$

$LD(x, x')$

Lipschitz Continuous MDPs
**Rank <= Covering number
of state space**

[KLK, 03]

Factored MDPs
**Rank <= exp(in-degree)**

[GKPV, 03; OV, 13, NIPS ]

# A Unified Framework



$x$, $x'$ — $a$, $a$ → $L\mathcal{D}(x, x')$

Lipschitz Continuous MDPs
**Rank <= Covering number of state space**

[KLK, 03]

Factored MDPs
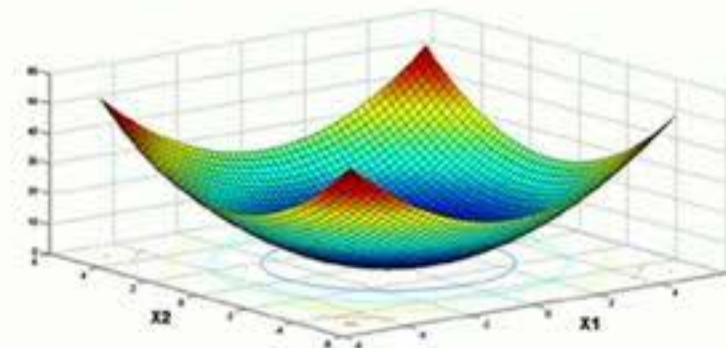**Rank <= exp(in-degree)**

[GKPV, 03; OV, 13, NIPS ]

POMDP
**Rank <= # of hidden states**

[KAL, 16 NIPS]

42

# A Unified Framework



$x$, $x'$ ... $a$, $a$ ... $L\mathcal{D}(x, x')$

**Lipschitz Continuous MDPs**
**Rank <= Covering number of state space**
[KLK, 03]

**Factored MDPs**
**Rank <= exp(in-degree)**
[GKPV, 03; OV, 13, NIPS ]

**POMDP**
**Rank <= # of hidden states**
[KAL, 16 NIPS]

**Linear Quadratic Regulator**
**Rank = O(d^2)**

42

# Sample Complexity

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

Model-based Reinforcement Learning in Contextual Decision Processes
[Sun, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

# Sample Complexity

Model Rank

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

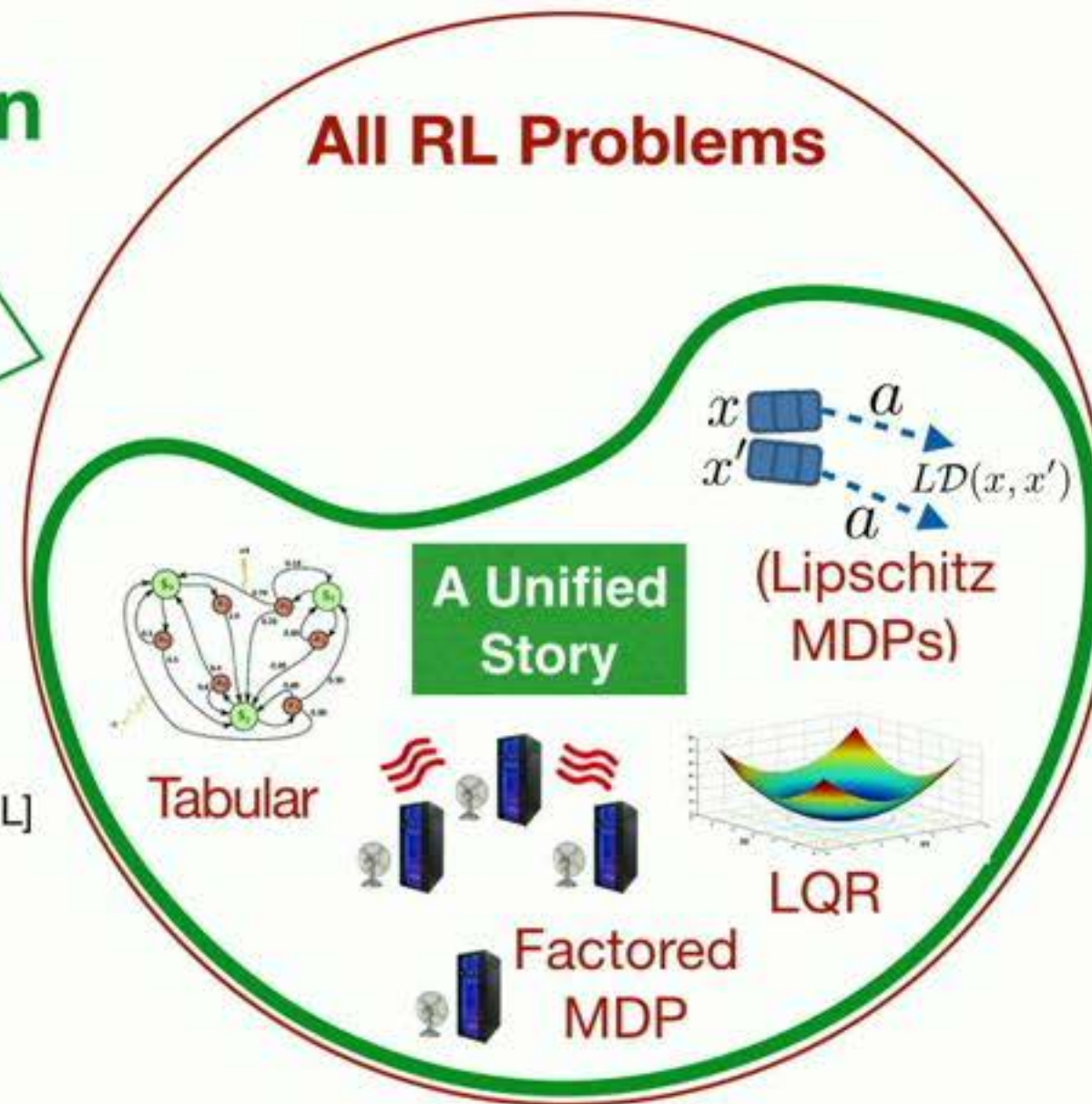Complexity of Discriminators & Models

**Poly Dependency on # of States**

Model-based Reinforcement Learning in Contextual Decision Processes
[Sun, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

# Sample Complexity

Model Rank

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

Complexity of Discriminators & Models

**Poly Dependency on # of States**

**Supervised Learning Type Generalization !**

Model-based Reinforcement Learning in Contextual Decision Processes
[**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

43

# Generalization & Sample Efficiency via...

## 1. Expert Demonstration



All RL Problems

[**Sun**, Venkatraman, Gordon, Boots, Bagnell, 17, ICML]

[**Sun**, Gordon, Boots, Bagnell, 18, NeurIPS]

# Sample Complexity

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

Model-based Reinforcement Learning in Contextual Decision Processes
[Sun, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

# Sample Complexity

Model Rank

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

Complexity of Discriminators & Models

**Poly Dependency on # of States**

Model-based Reinforcement Learning in Contextual Decision Processes
[**Sun**, Jiang, Krishnamurthy, Agarwal, Langford, arXiv, 18]

# Future Work



Medical Treatment

Education

Autonomous Driving

Waseda University's Manga club

Assistance in Disaster Recovery

# 1. Leverage Expert Demonstrations

# 1. Leverage Expert Demonstrations

## Interactive Imitation Learning

# 1. Leverage Expert Demonstrations

# 1. Leverage Expert Demonstrations

No Interaction

No Expert Action

No Reward



**Imitation Learning from Observations**

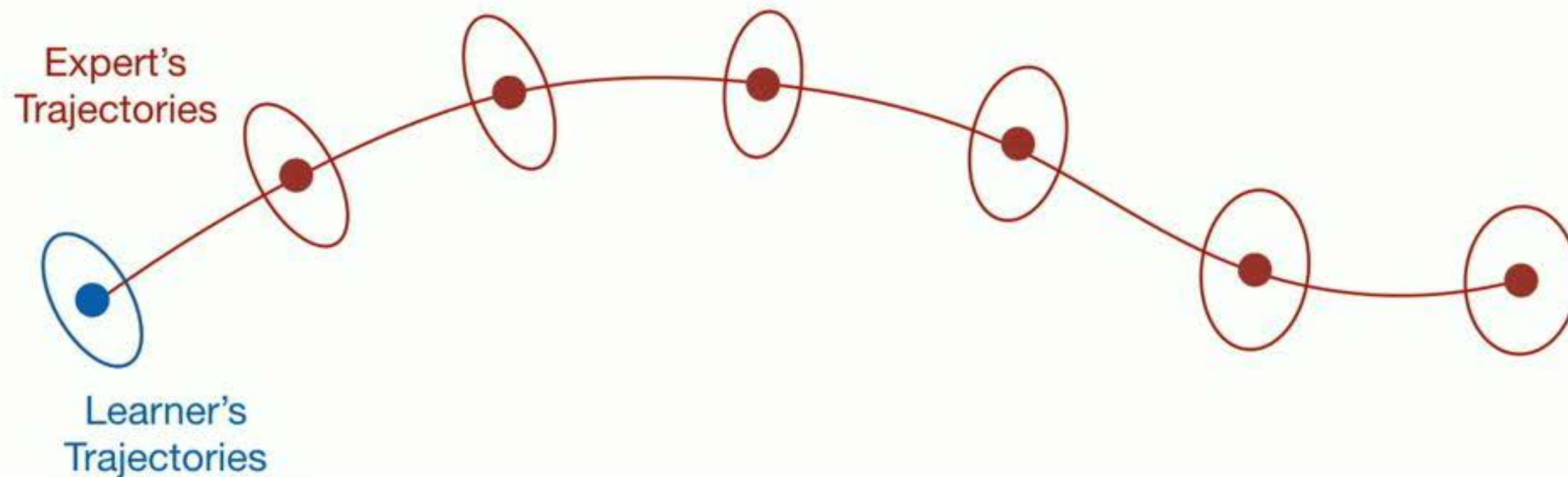# Forward Adversarial Imitation Learning (FAIL):

[Sun et.al, In Submission, 19]

## Learn policies using Integral Probability Metric



Expert's
Trajectories

# 1. Leverage Expert Demonstrations

No Interaction
No Expert Action
No Reward



**Imitation Learning from Observations**

# Forward Adversarial Imitation Learning (FAIL):

[Sun et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

Expert's Trajectories

# Forward Adversarial Imitation Learning (FAIL):

[**Sun** et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

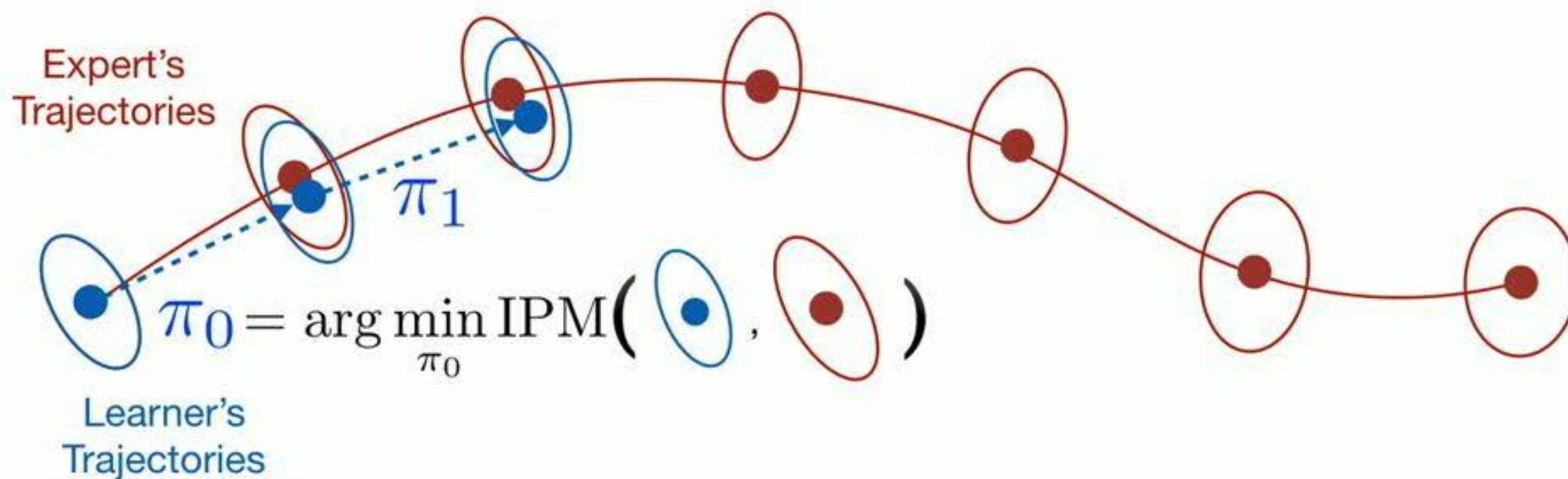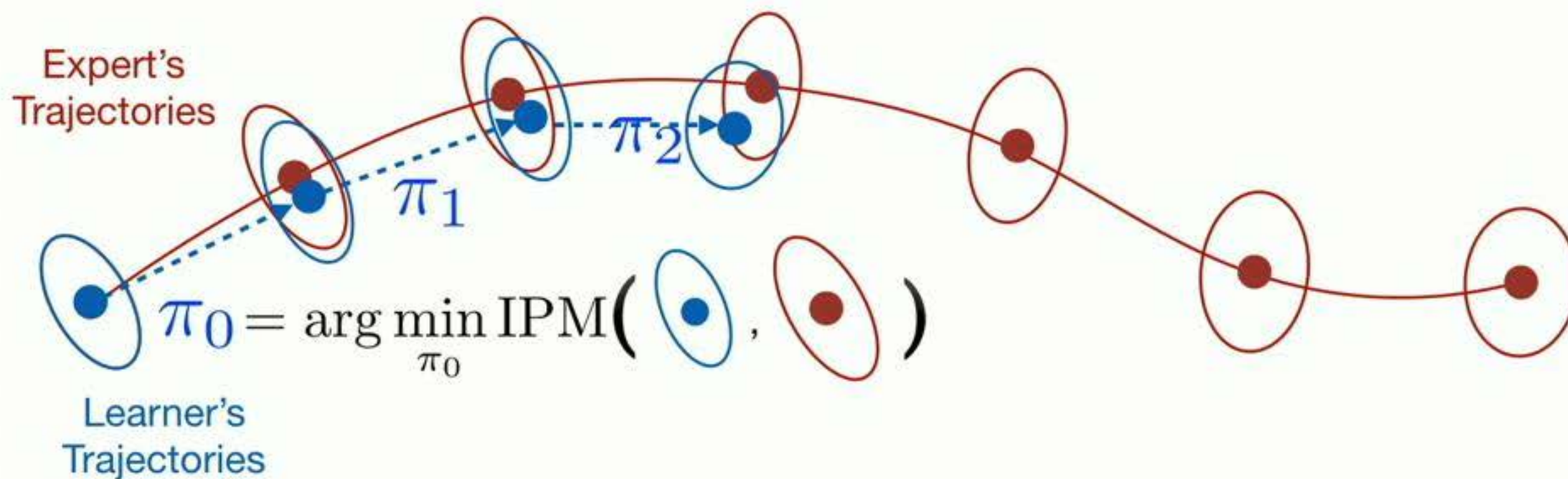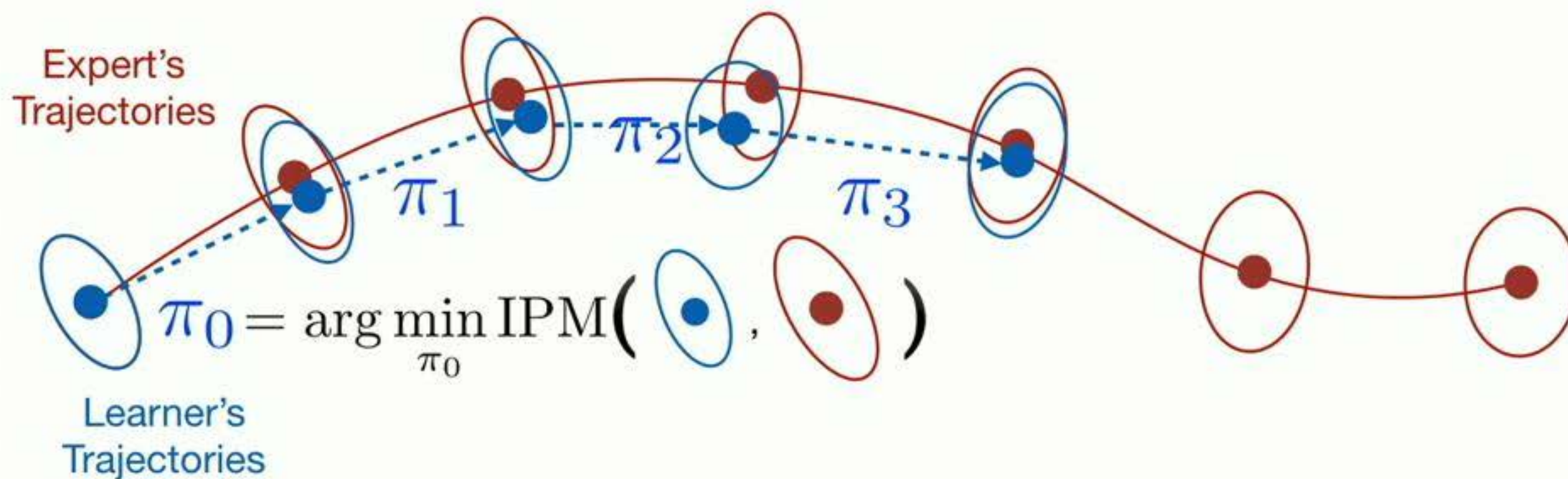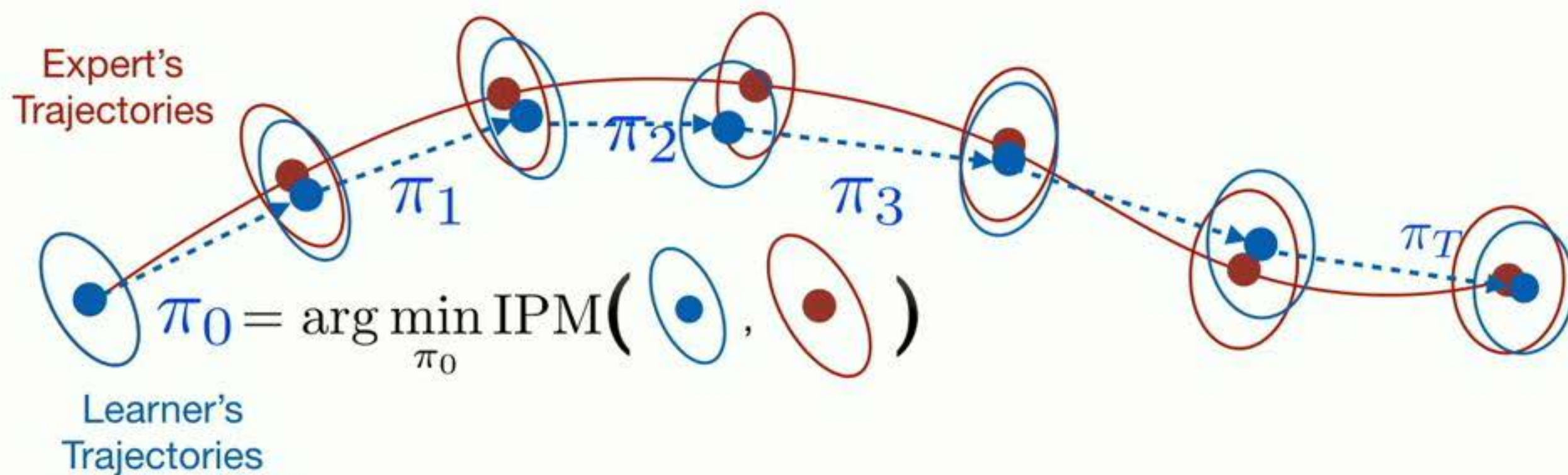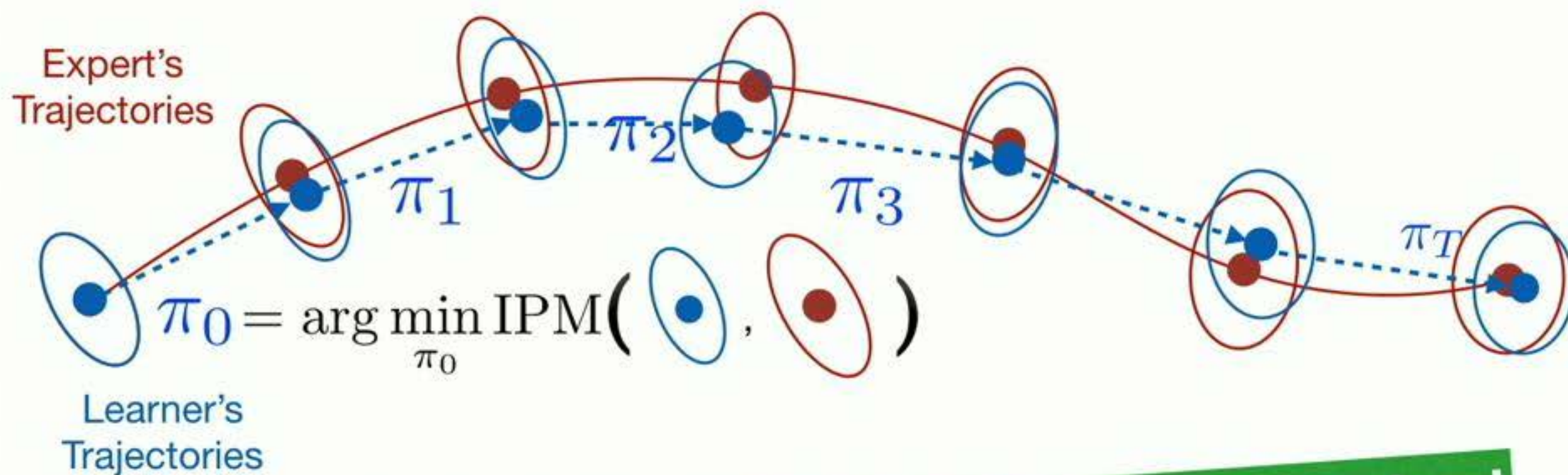# Forward Adversarial Imitation Learning (FAIL):

[Sun et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

# Forward Adversarial Imitation Learning (FAIL):

[**Sun** et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

Expert's Trajectories

Learner's Trajectories

$$\pi_0 = \arg\min_{\pi_0} \text{IPM}\left( \bullet , \bullet \right)$$

# Forward Adversarial Imitation Learning (FAIL):

[Sun et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

Expert's Trajectories

$\pi_1$

Learner's Trajectories

$$\pi_0 = \arg\min_{\pi_0} \text{IPM}\left( \bullet , \bullet \right)$$

# Forward Adversarial Imitation Learning (FAIL):

[**Sun** et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

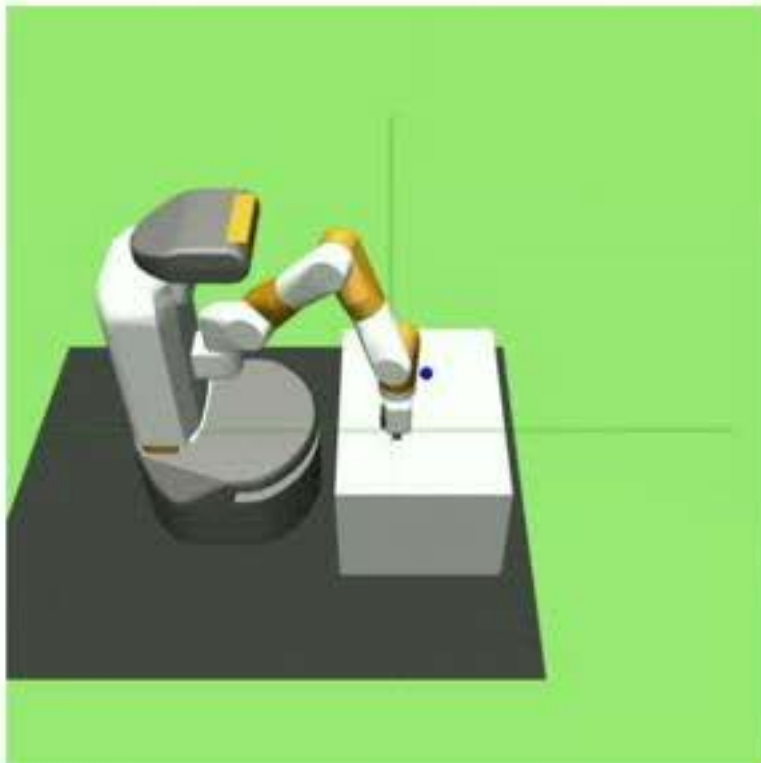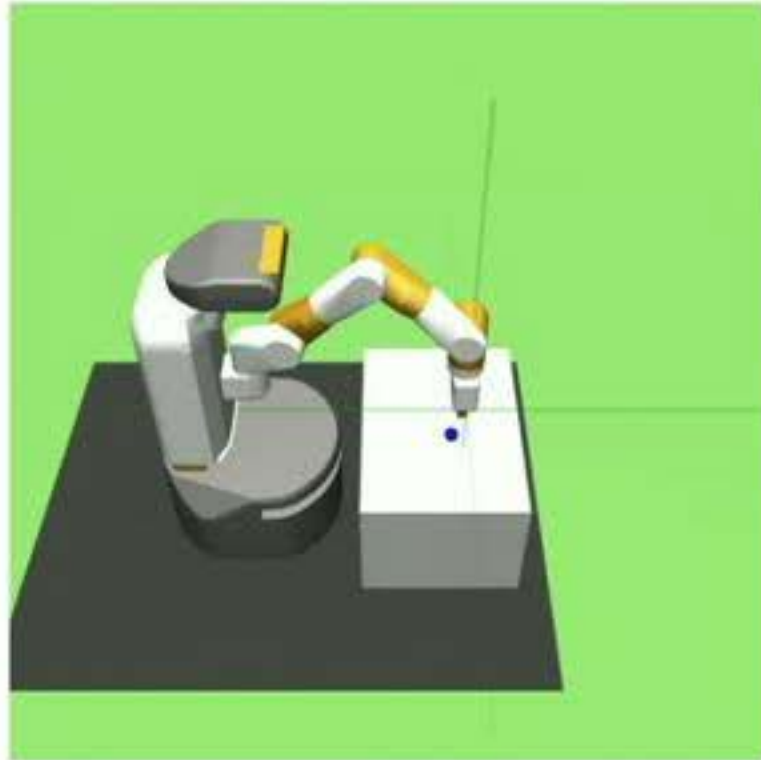# Forward Adversarial Imitation Learning (FAIL):

[Sun et.al, In Submission, 19]

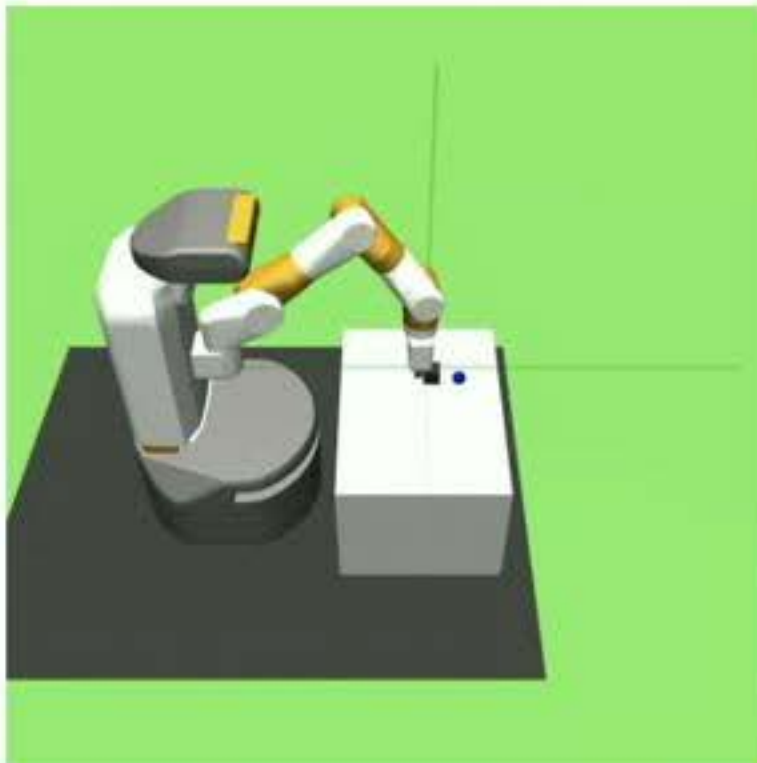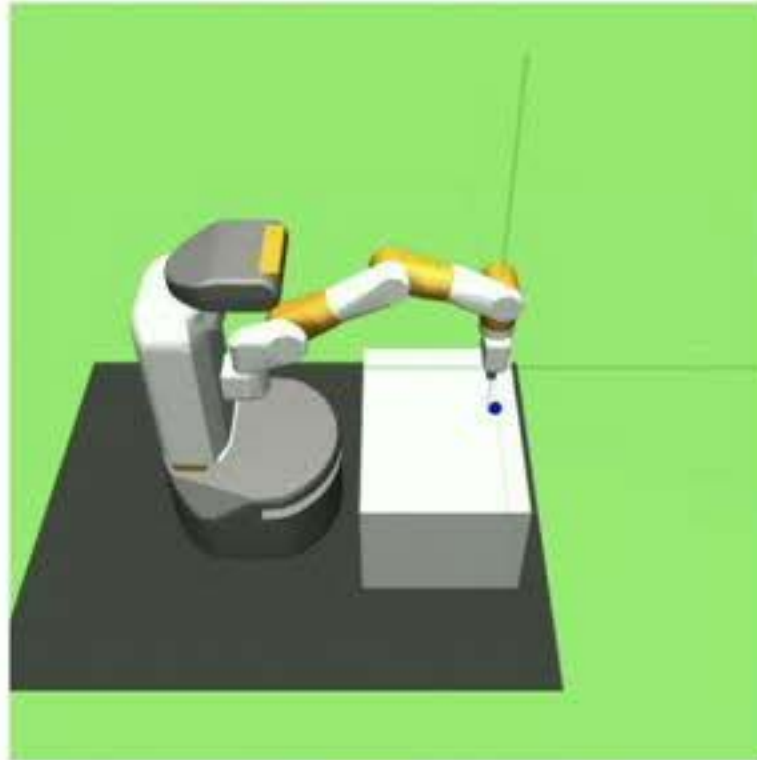## Learn policies using Integral Probability Metric

Expert's Trajectories

$\pi_2$

$\pi_1$

$\pi_3$

$\pi_0 = \arg \min_{\pi_0} \mathrm{IPM}\big( \bullet, \bullet \big)$

Learner's Trajectories

# Forward Adversarial Imitation Learning (FAIL):

[**Sun** et.al, In Submission, 19]

## Learn policies using Integral Probability Metric



Expert's Trajectories

$\pi_2$

$\pi_1$

$\pi_3$

$\pi_T$

$$\pi_0 = \arg \min_{\pi_0} \mathrm{IPM}\left( \bullet , \bullet \right)$$

Learner's Trajectories

# Forward Adversarial Imitation Learning (FAIL):

[**Sun** et.al, In Submission, 19]

## Learn policies using Integral Probability Metric

$$\pi_0 = \arg\min_{\pi_0} \text{IPM}\left( \bullet , \bullet \right)$$

Expert's Trajectories

Learner's Trajectories

$\pi_1$  $\pi_2$  $\pi_3$  $\pi_T$

**Supervised Learning Type Generalization !**

# Promising Simulation Results...



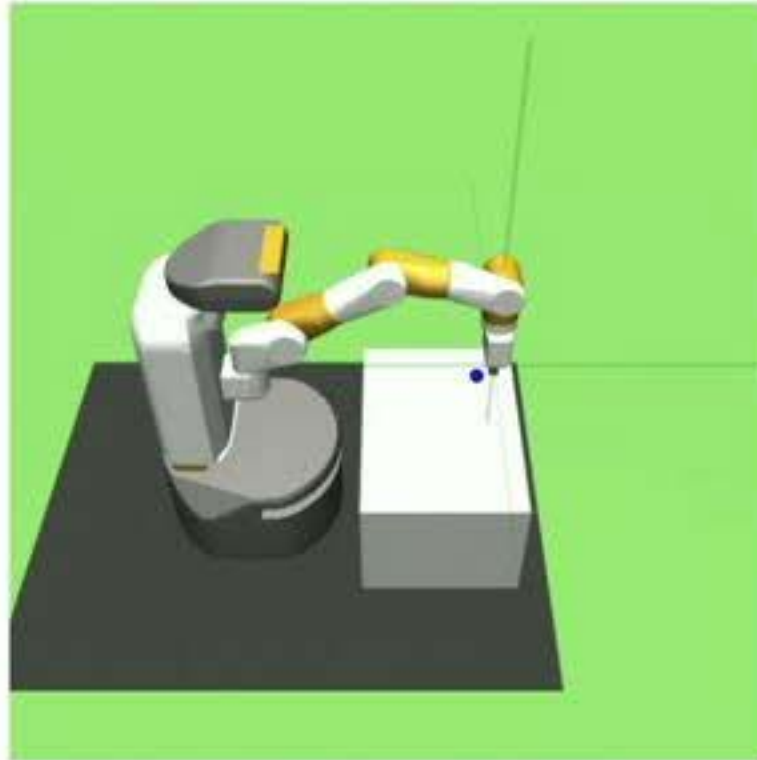[Fetch Robot Simulator from OpenAI Gym]

48

# Promising Simulation Results...



Image from https://www.asme.org/engineering-topics/articles/
robotics/robots-kitchen-at-the-table

[Fetch Robot Simulator from OpenAI Gym]

# Promising Simulation Results...



[Fetch Robot Simulator from OpenAI Gym]



Image from https://www.asme.org/engineering-topics/articles/robotics/robots-kitchen-at-the-table

## Lots of Challenges:

—Learn from videos

—Interaction with experts

# 2. Generalization from Prior Experiences

Medical Treatment

Assistance in Disaster Recovery

Waseda University's Manga club

# 2. Generalization from Prior Experiences

# 2. Generalization from Prior Experiences

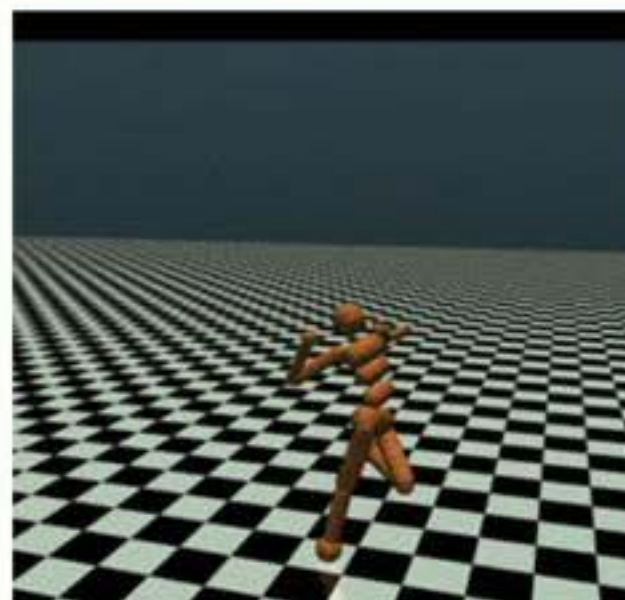# 2. Generalization from Prior Experiences



Right Leg Jump Demo

Videos from Ben Recht's Blog (http://www.argmin.net/2018/03/20/mujocoloco/)

# 2. Generalization from Prior Experiences



Right Leg Jump Demo          Backward Demo

# 2. Generalization from Prior Experiences



Right Leg Jump Demo        Backward Demo        Forward Demo
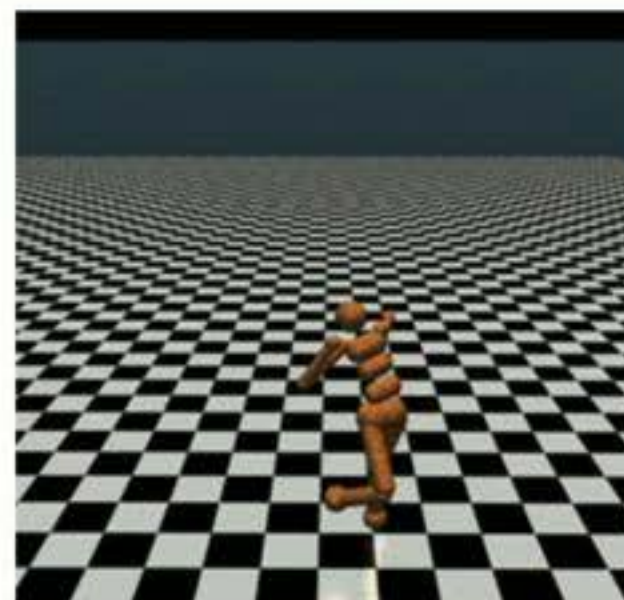
Videos from Ben Recht's Blog (http://www.argmin.net/2018/03/20/mujocoloco/)

# 2. Generalization from Prior Experiences



Right Leg Jump Demo

Backward Demo

Forward Demo

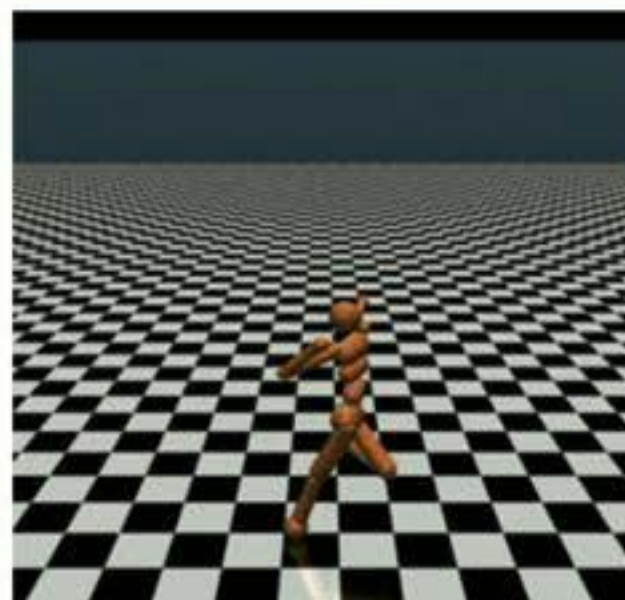...

**New task:**
**Stand up** with little to no training?

Videos from Ben Recht's Blog (http://www.argmin.net/2018/03/20/mujocoloco/)

# 2. Generalization from Prior Experiences



Right Leg Jump Demo

Backward Demo

Forward Demo
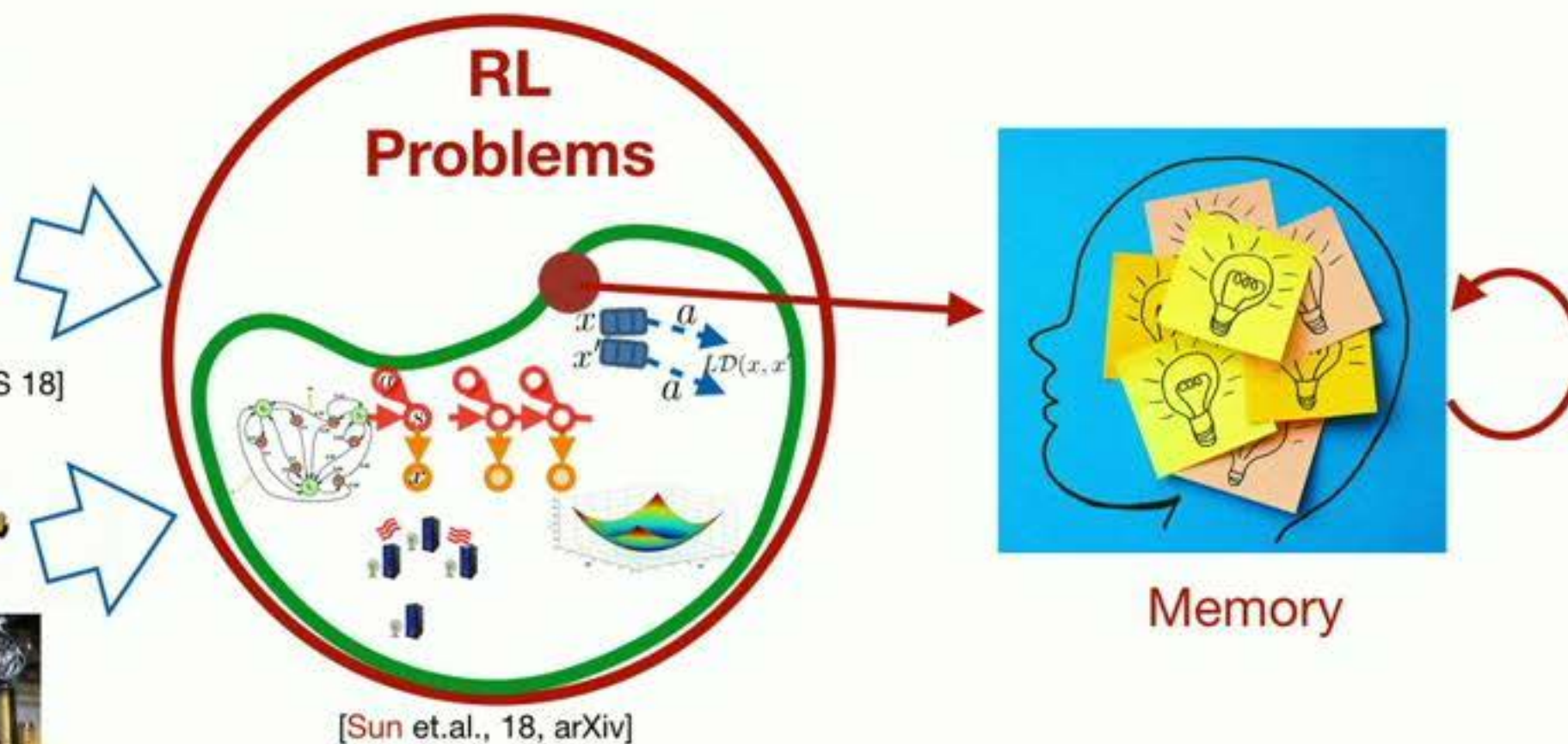
**Offline Learning
From Prior Relevant
Experiences**
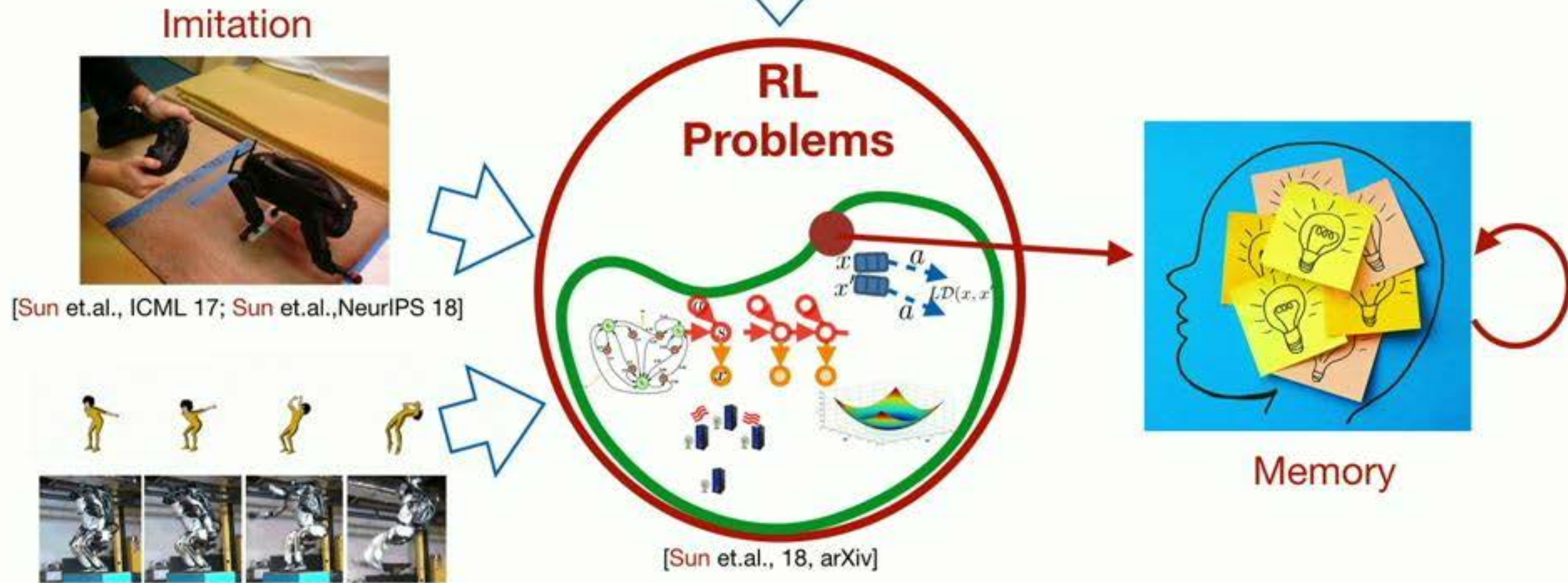
**New task:
Stand up** with little to no training?

Videos from Ben Recht's Blog (http://www.argmin.net/2018/03/20/mujocoloco/)

Imitation

[Sun et.al., ICML 17; Sun et.al.,NeurIPS 18]

RL Problems

[Sun et.al., 18, arXiv]

Memory

**Online Policy Evaluation**

**Reduction** from Policy Evaluation to No-Regret Online Learning
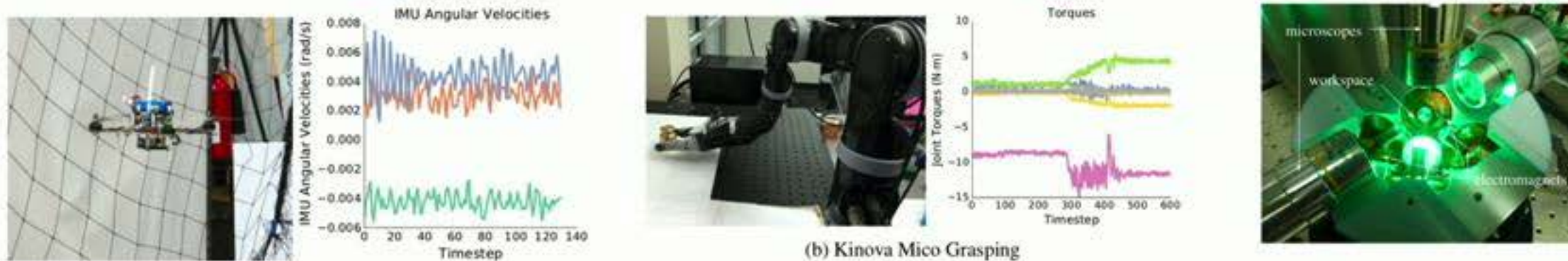[Sun, Bagnell, UAI 15, Best Student Paper]

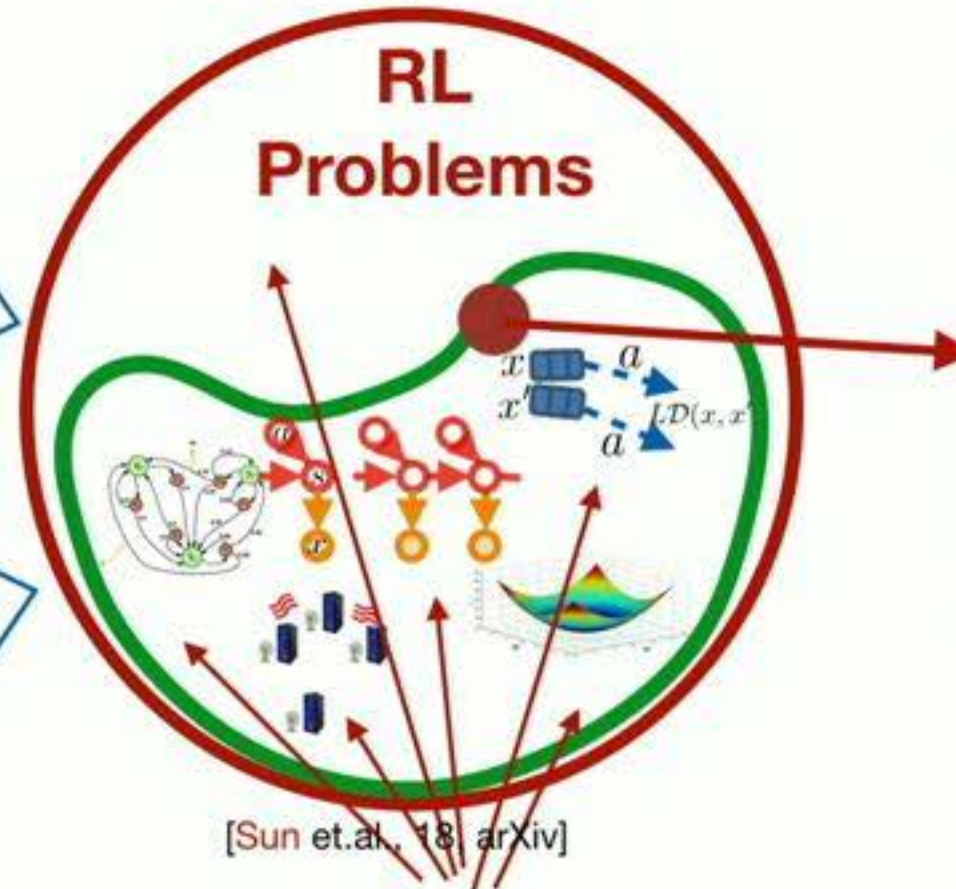Imitation

[Sun et.al., ICML 17; Sun et.al.,NeurIPS 18]

RL
Problems

[Sun et.al., 18, arXiv]

Memory

**Online Policy Evaluation**

**Reduction** from Policy Evaluation to No-Regret Online Learning
[Sun, Bagnell, UAI 15, Best Student Paper]

Imitation

[Sun et.al., ICML 17; Sun et.al.,NeurIPS 18]

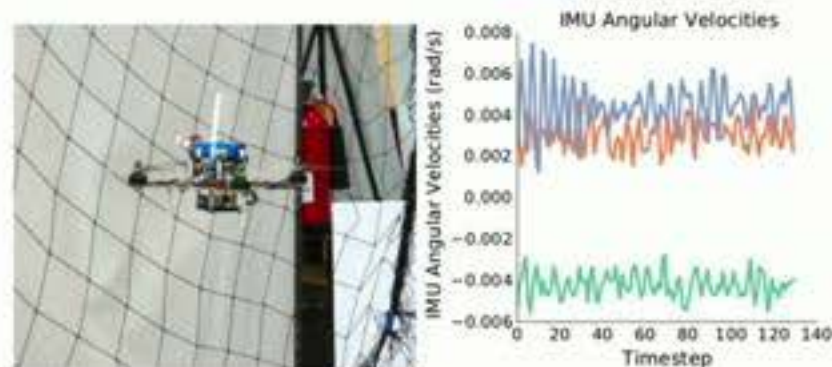**RL Problems**

[Sun et.al., 18, arXiv]

Memory

**System ID**

**Predictive State Inference Machines**
[Sun et.al., ICML 16; Venkatraman & Sun et.al., IJCAI 16, Sun et.al., ICRA 14]

(b) Kinova Mico Grasping

52

# Online Policy Evaluation

**Reduction** from Policy Evaluation to No-Regret Online Learning
[Sun, Bagnell, UAI 15, Best Student Paper]

Imitation

[Sun et.al., ICML 17; Sun et.al.,NeurIPS 18]
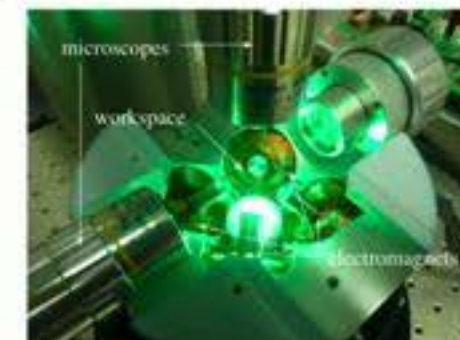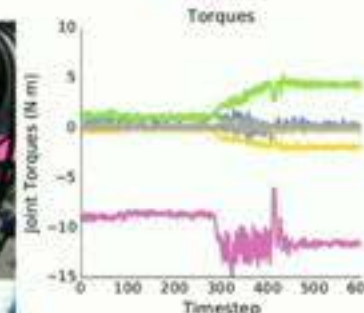
RL Problems

[Sun et.al., 18, arXiv]

Memory

## System ID

**Predictive State Inference Machines**
[Sun et.al., ICML 16; Venkatraman & Sun et.al., IJCAI 16, Sun et.al., ICRA 14]

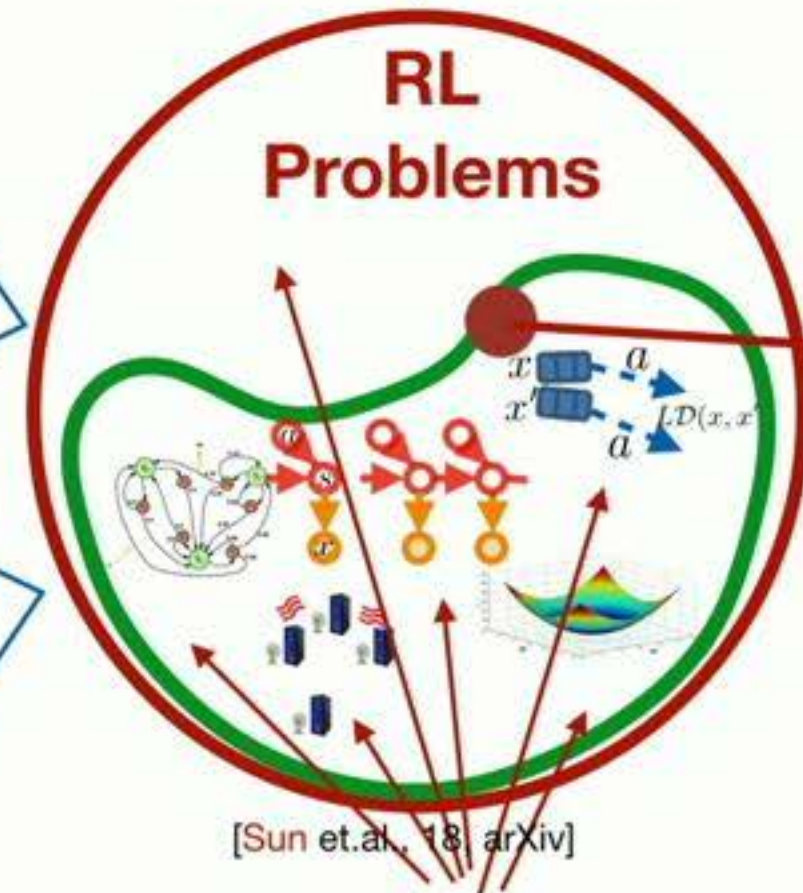(b) Kinova Mico Grasping

52

# Thank You

Online Policy Evaluation

Reduction from Policy Evaluation to No-Regret Online Learning
[Sun, Bagnell, UAI 15, Best Student Paper]
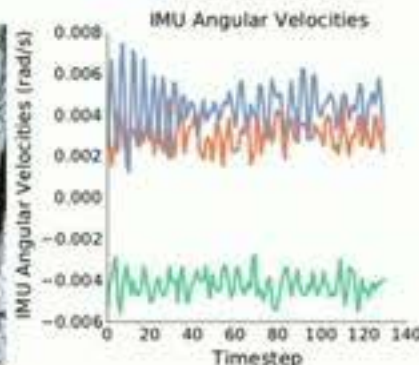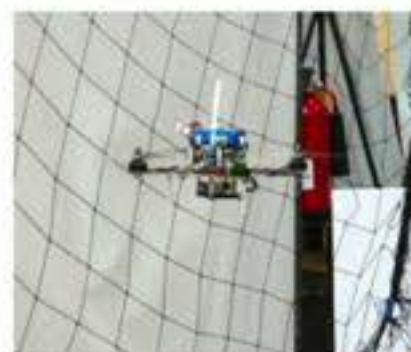
Imitation

[Sun et.al., ICML 17; Sun et.al.,NeurIPS 18]

RL Problems

[Sun et.al., 18, arXiv]

Memory

System ID

Predictive State Inference Machines
[Sun et.al., ICML 16; Venkatraman & Sun et.al., IJCAI 16, Sun et.al., ICRA 14]

(b) Kinova Mico Grasping