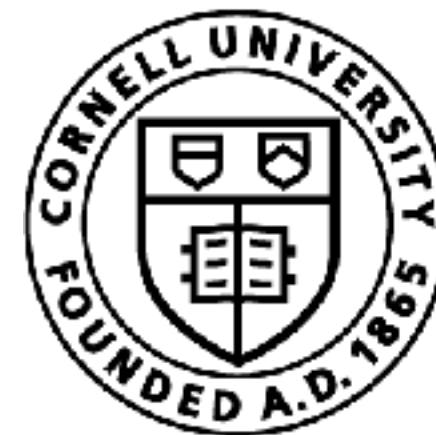


Learning for Robot Decision Making: The Big Picture

Sanjiban Choudhury



Cornell Bowers CIS
Computer Science



How should robots **learn** to make **good** decisions?



WHY ask this question?



Formulate as a Markov Decision Problem (MDP)



Solve MDPs using an all-purpose toolkit
(Imitation/Reinforcement learning, Model based/free)



Deploy learners in real-world
(Safety, distribution shift, value alignment)

Take *any* robot application

HOW can we answer this question?

Solve

*How do you want to represent your policy?
Model-based? Model-free?
Learning: Data? Loss?*

Formulate

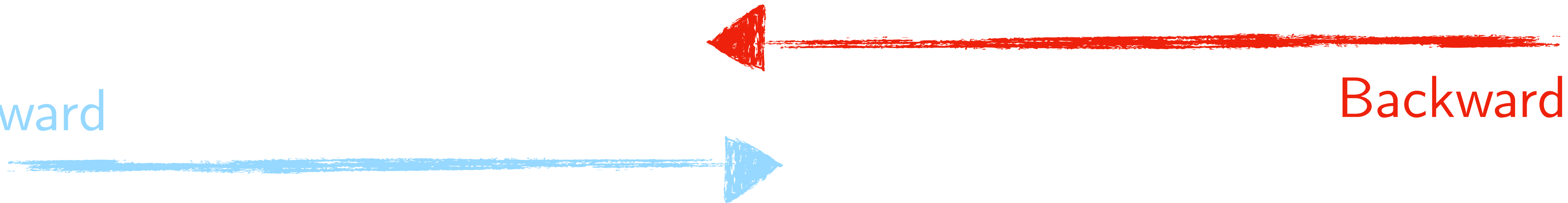
*What is the MDP?
Discrete/Stochastic/Time?
What is known/unknown?*

Application

*What is the robot?
What is the task?
What are the metrics?
What is good enough?*

Forward

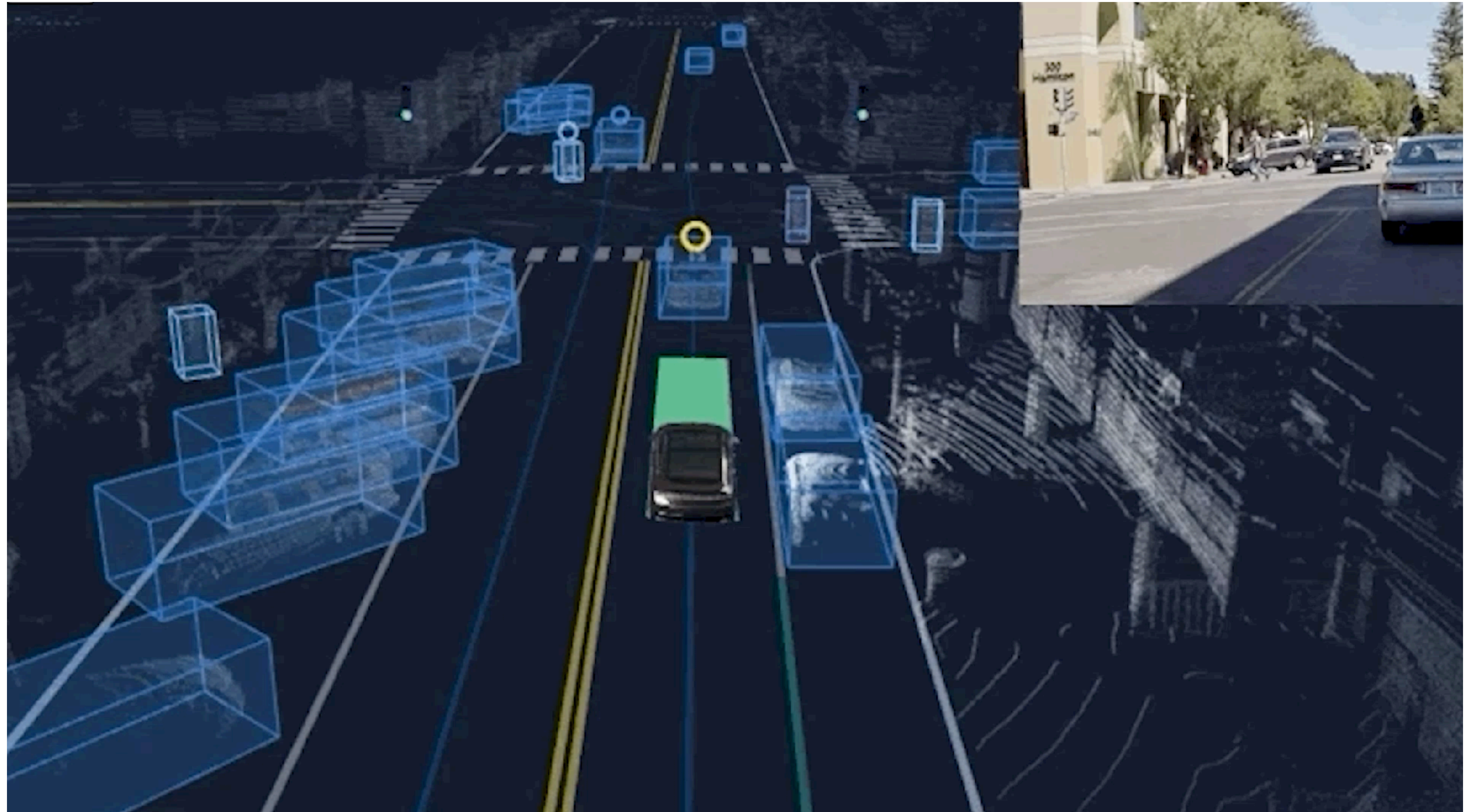
Backward



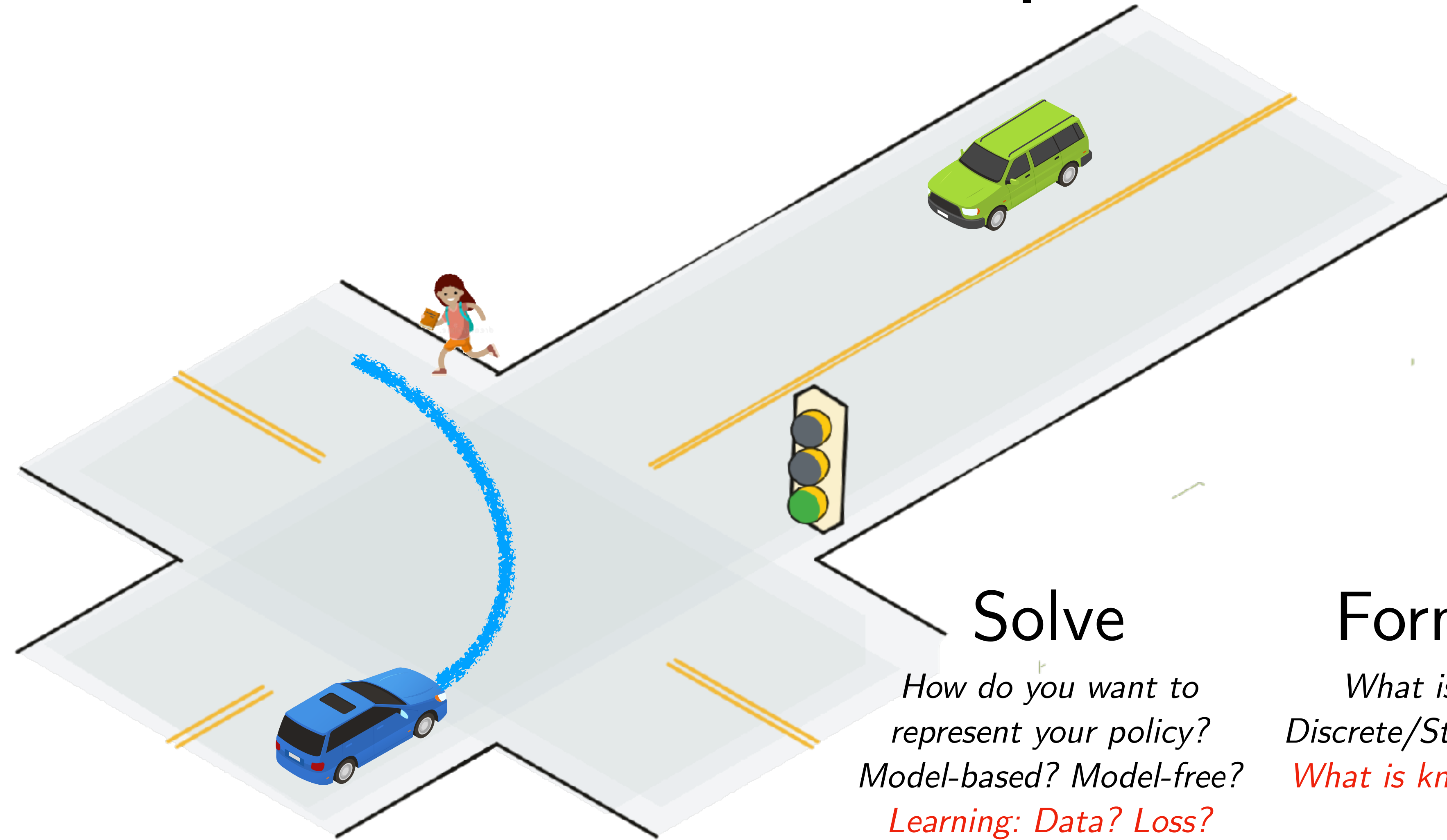
Activity!



Let's solve the Unprotected Left Turn



Let's solve the Unprotected Left Turn



Solve

*How do you want to represent your policy?
Model-based? Model-free?
Learning: Data? Loss?*

Formulate

*What is the MDP?
Discrete/Stochastic/Time?
What is known/unknown?*

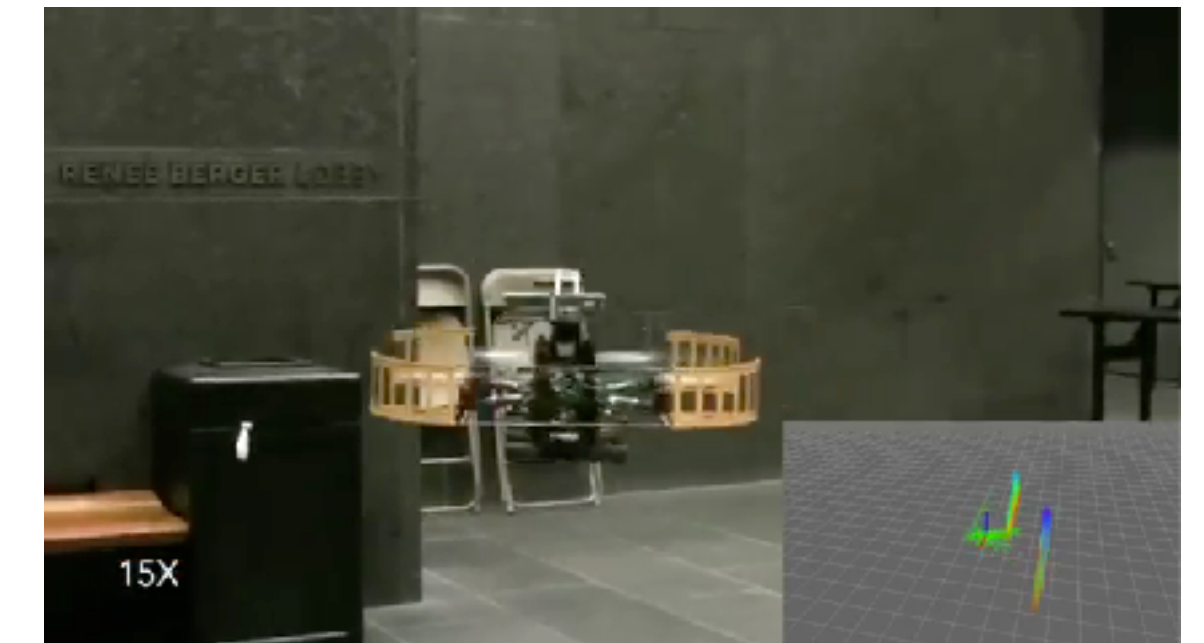
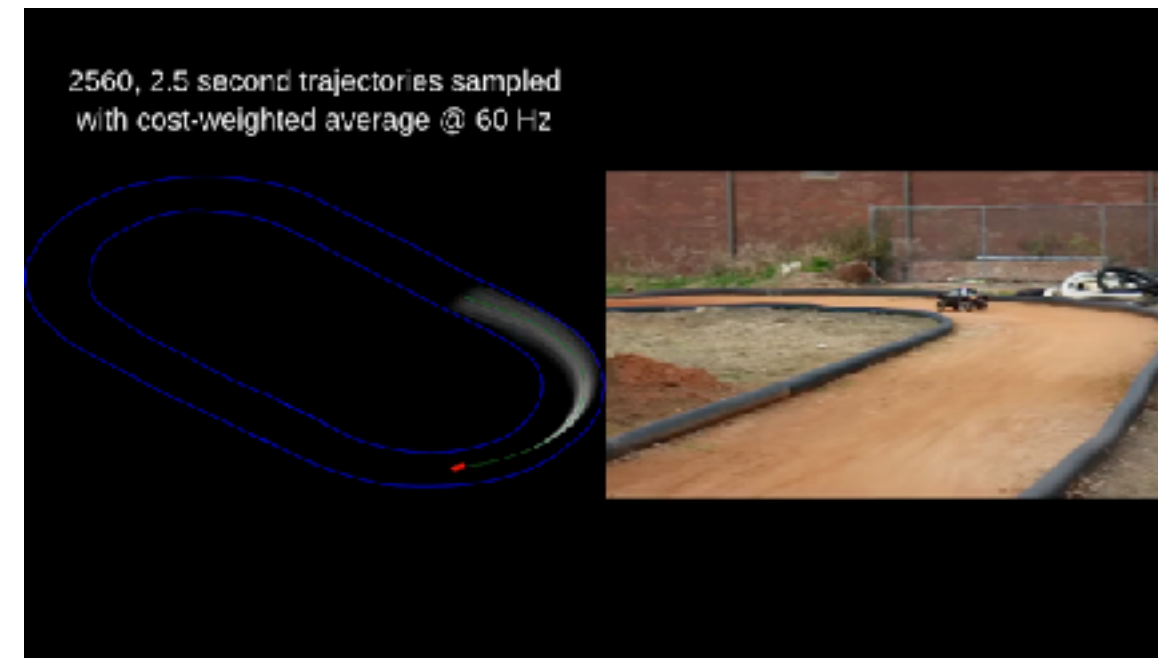
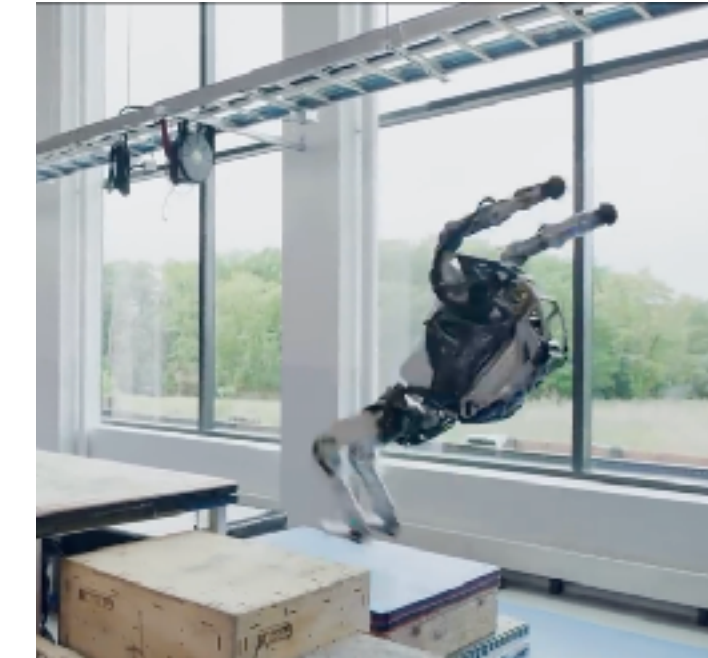
Application

*What is the robot?
What is the task?
What are the metrics?
What is good enough?*

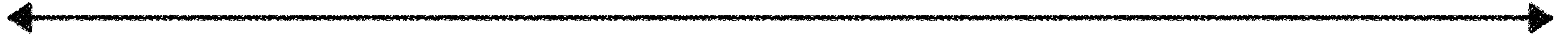




We have worked through many applications in this class ...

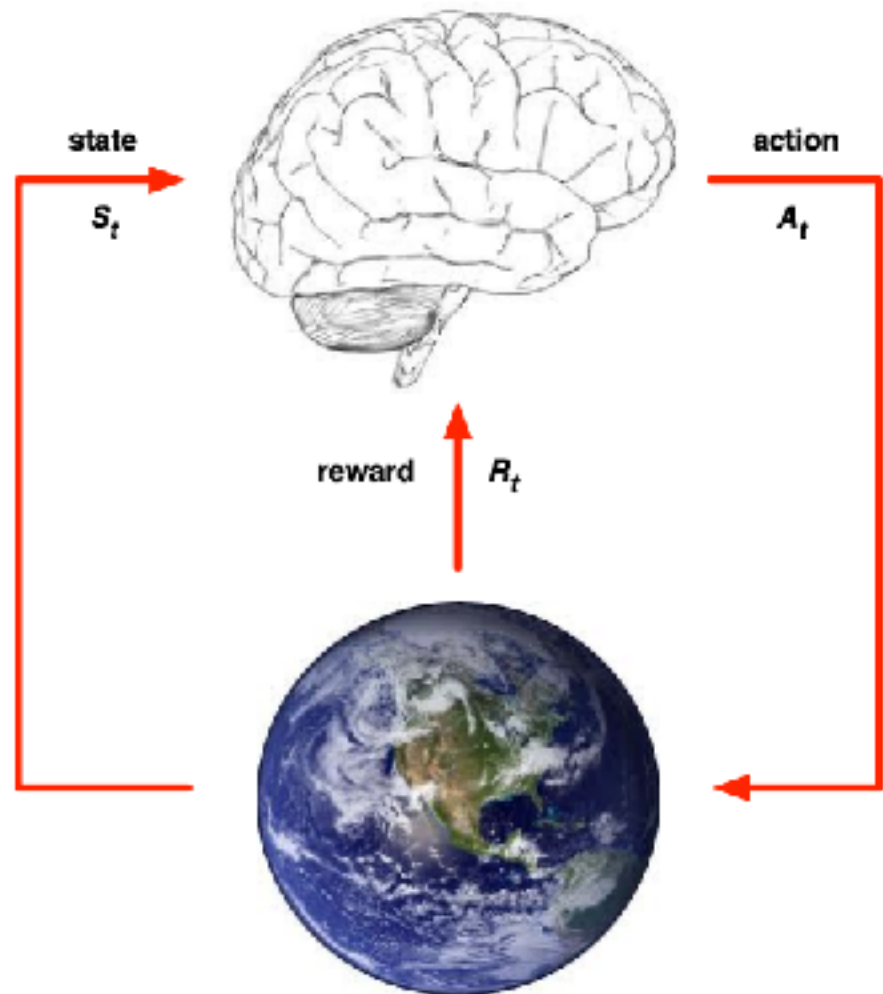


Model-Based OR Model Free?



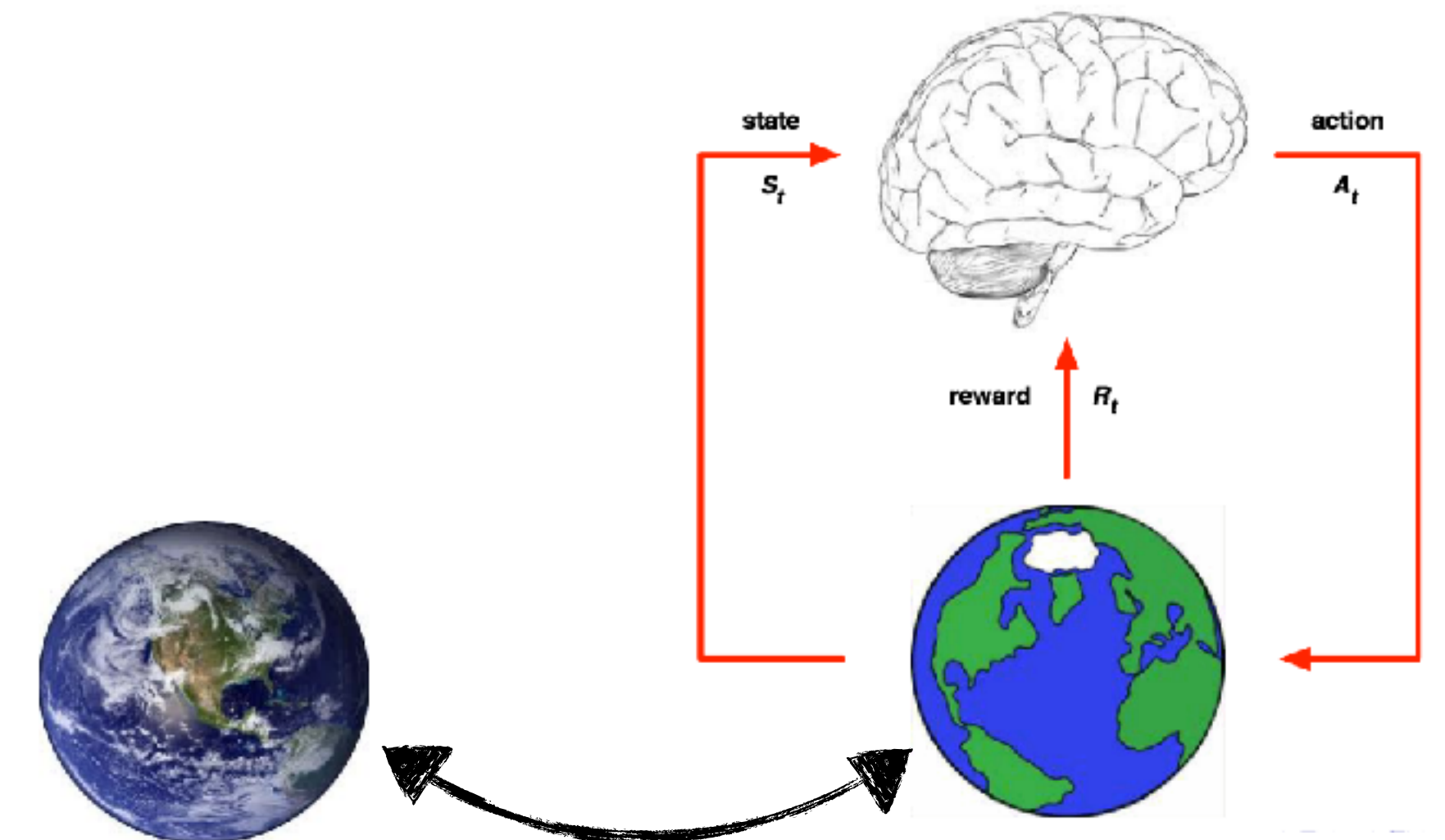
Model Free

Directly learn
 π or $Q(s, a)$



Model Based

Learn a model
 $P(s' | s, a)$, plan with
model to find π



Model-Based OR Model Free?



Model Free

There exists a good
enough reactive policy

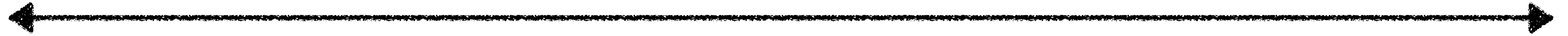
State space is too big to
search exhaustively

Model Based

You need to reason about
many likely options

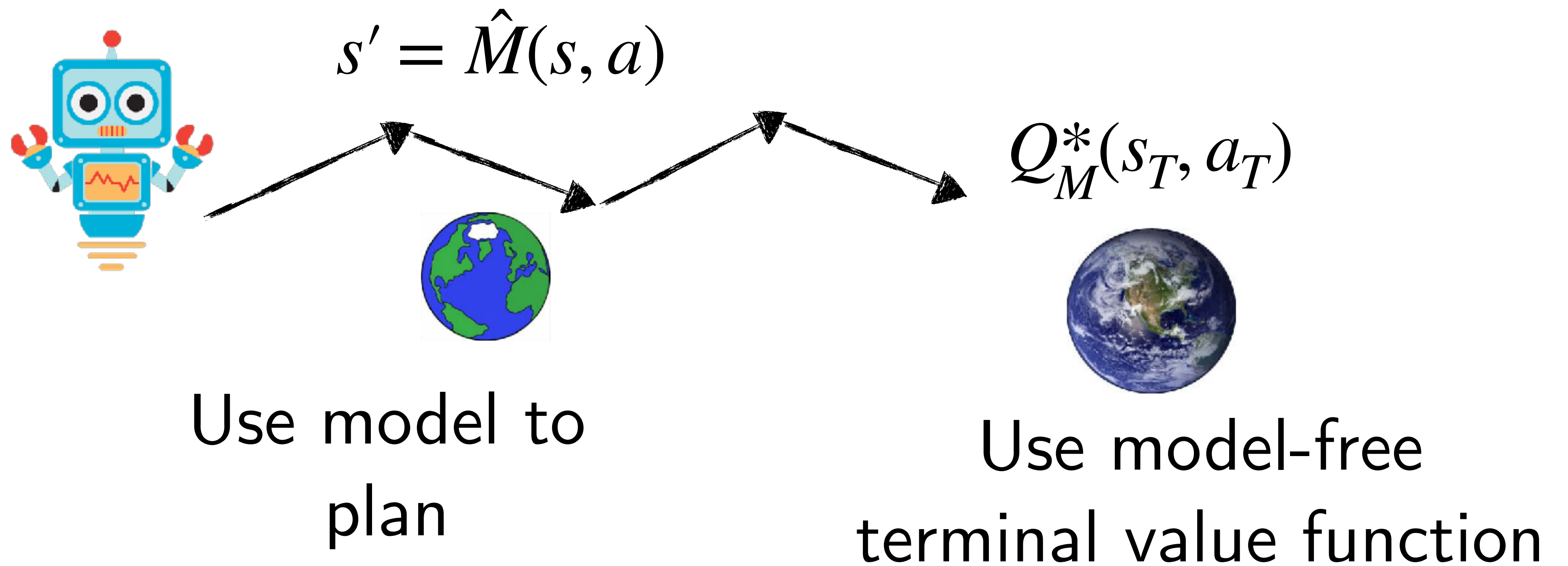
Small state space /
compressible state space

Model-Based ~~OR~~ *AND* Model Free



Model Based

Model Free



HOW can we answer this question?

Solve

*How do you want to represent your policy?
Model-based? Model-free?
Learning: Data? Loss?*

Formulate

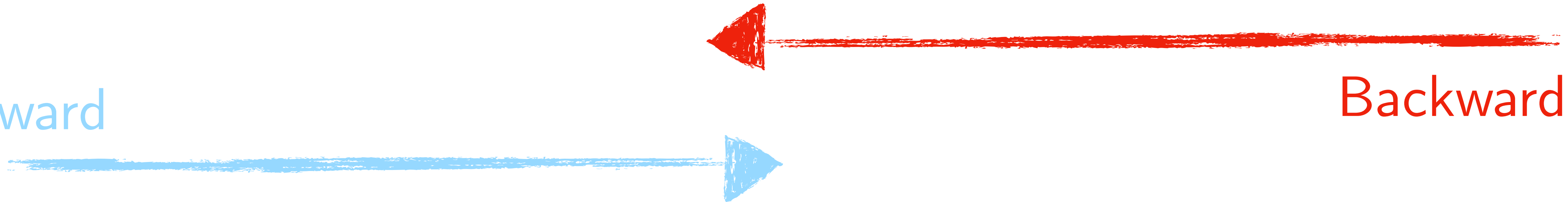
*What is the MDP?
Discrete/Stochastic/Time?
What is known/unknown?*

Application

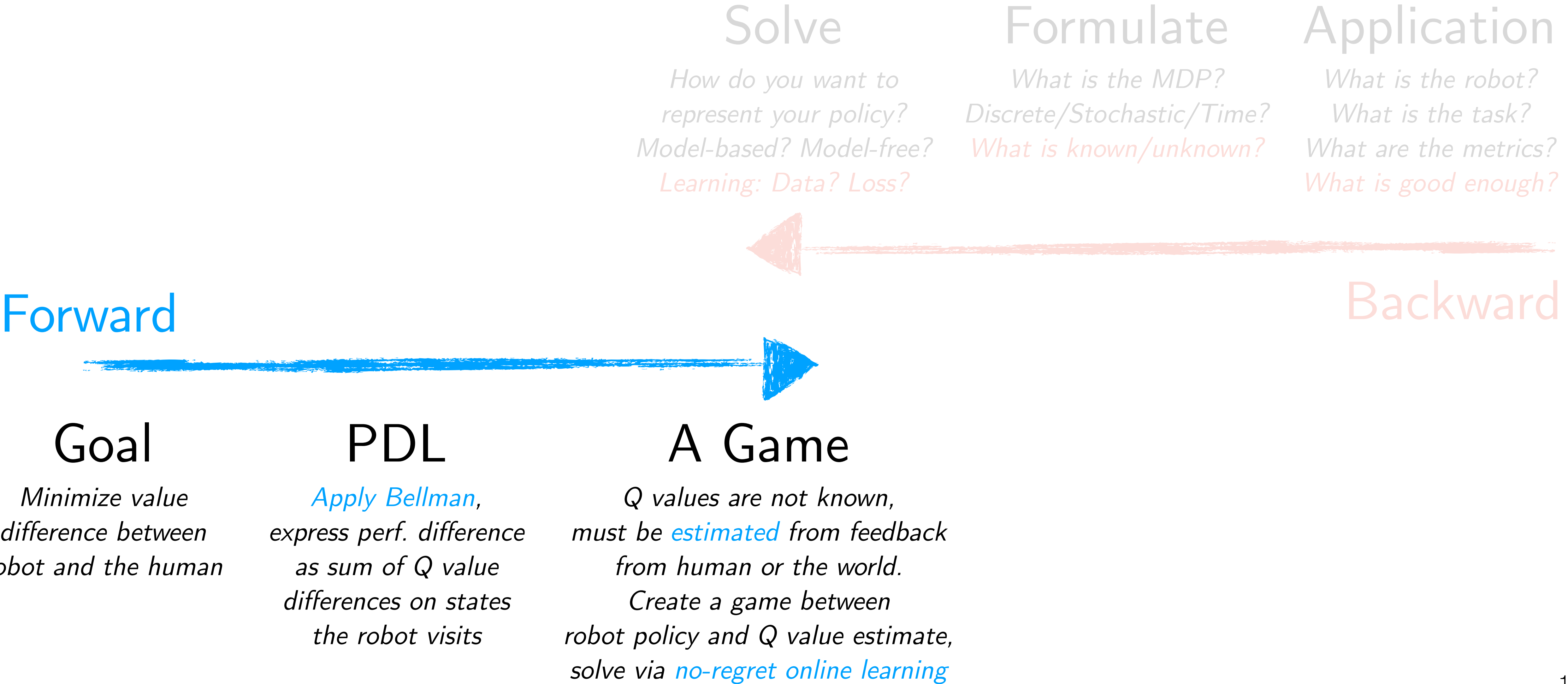
*What is the robot?
What is the task?
What are the metrics?
What is good enough?*

Forward

Backward



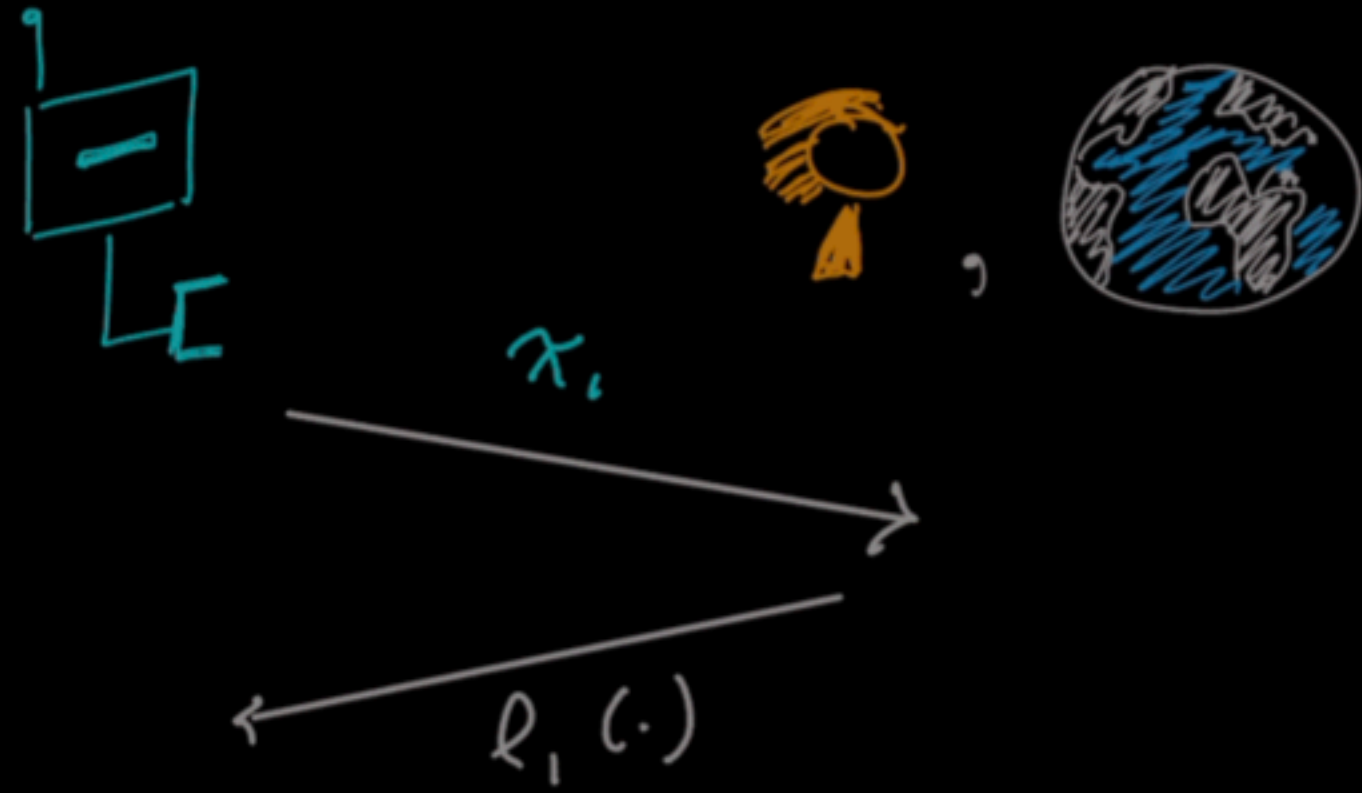
HOW can we answer this question?



5 Levels

of

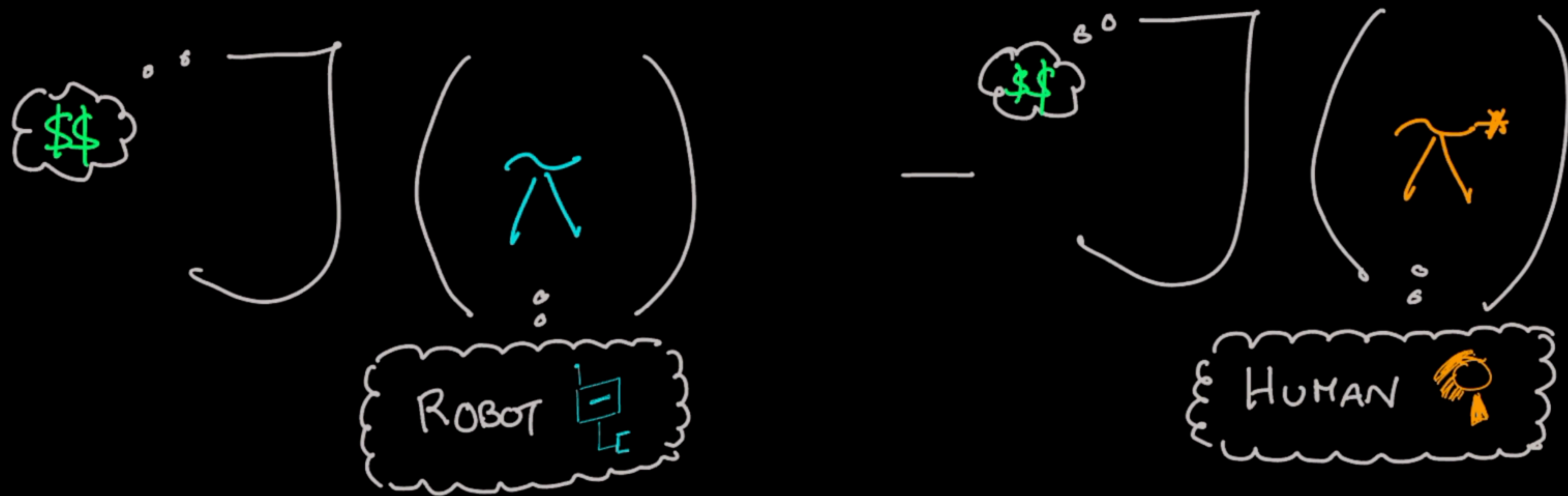
Robot Learning




$$\min_{\pi} \sum_{i=1}^{\infty} Q^*(s, \pi(s)) - Q^*(s, \pi^*(s))$$

min ROBOT max ACTION VALUE

FIND A POLICY $\pi: S \rightarrow A$ SUCH THAT THE
VALUE DIFFERENCE



FOR ANY VALUE $J(\cdot)$ THE HUMAN  CARES ABOUT
\$\$

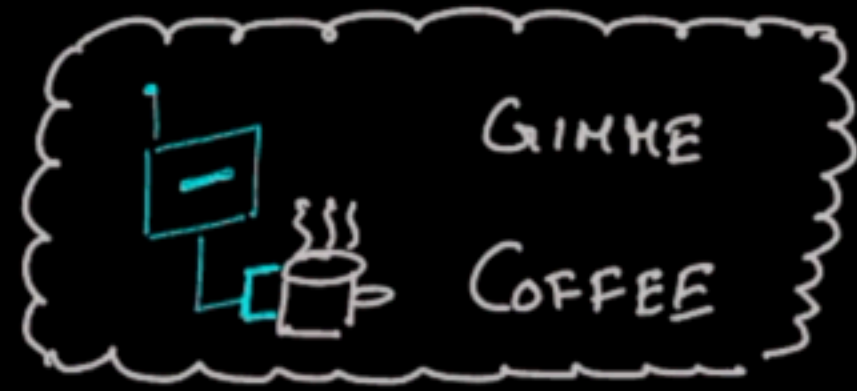
$$V^*(s_t)$$

$$= \min_{a_t}$$

$$C(s_t, a_t) +$$

$$E[V^*(s_{t+1})]$$

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$



HUMAN

$s_t, a_t?$



WORLD

LEARN OPTIMAL VALUE BY INTERACTIONS W/

HUMAN +
WORLD

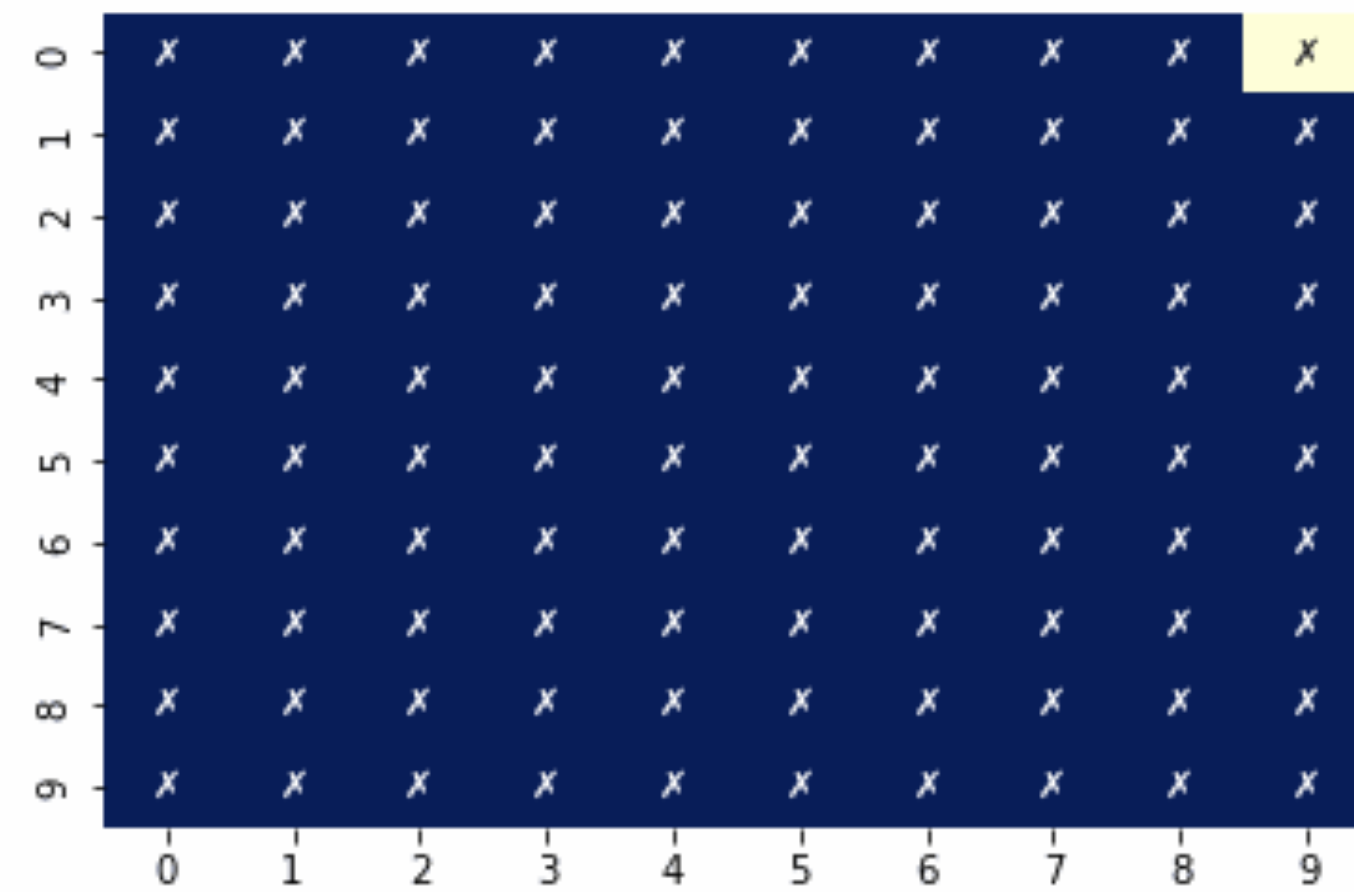
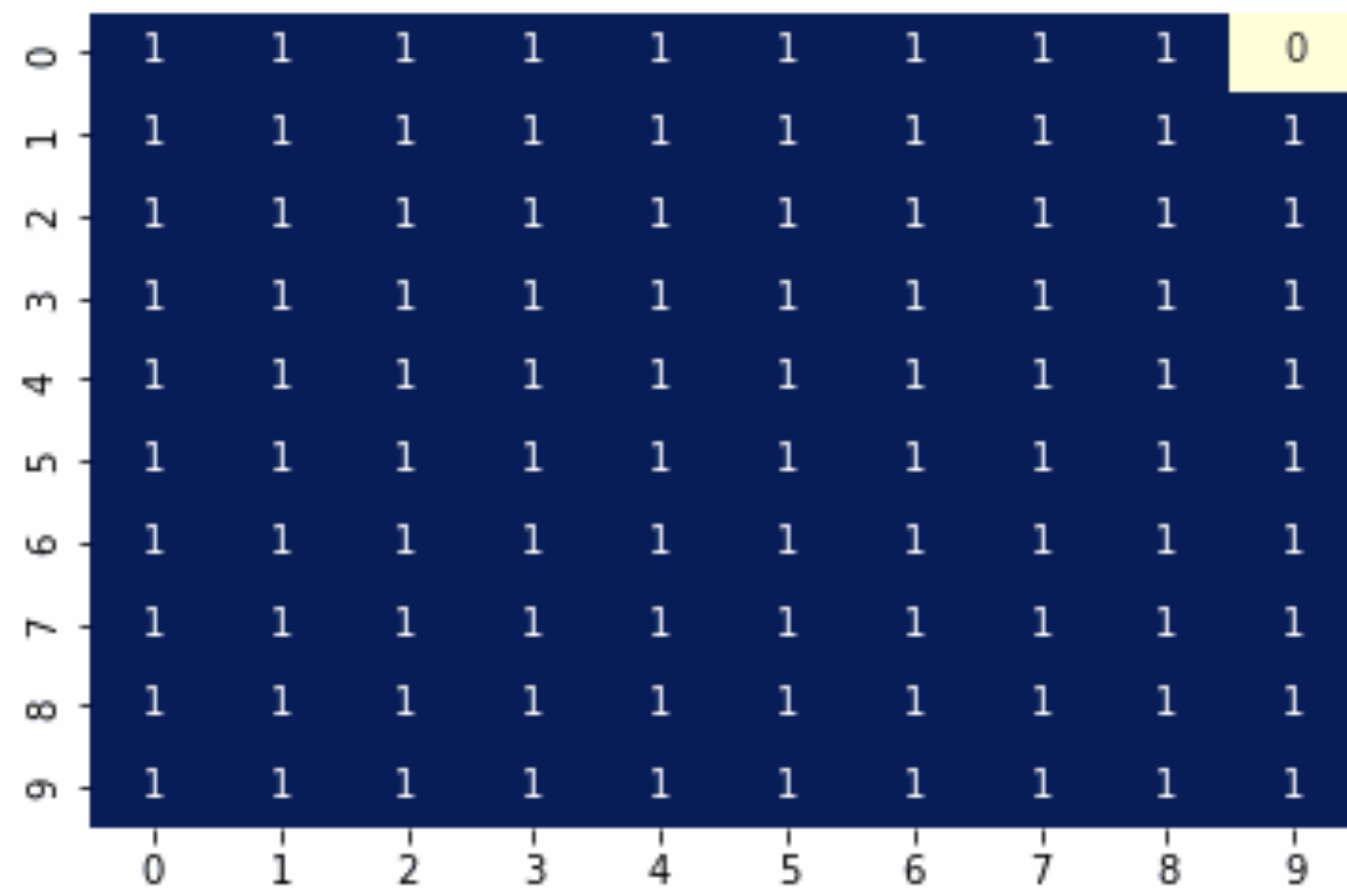
Two Fundamental Approaches

Value Iteration

Iterate over *optimal* value

min() operator in value iteration step

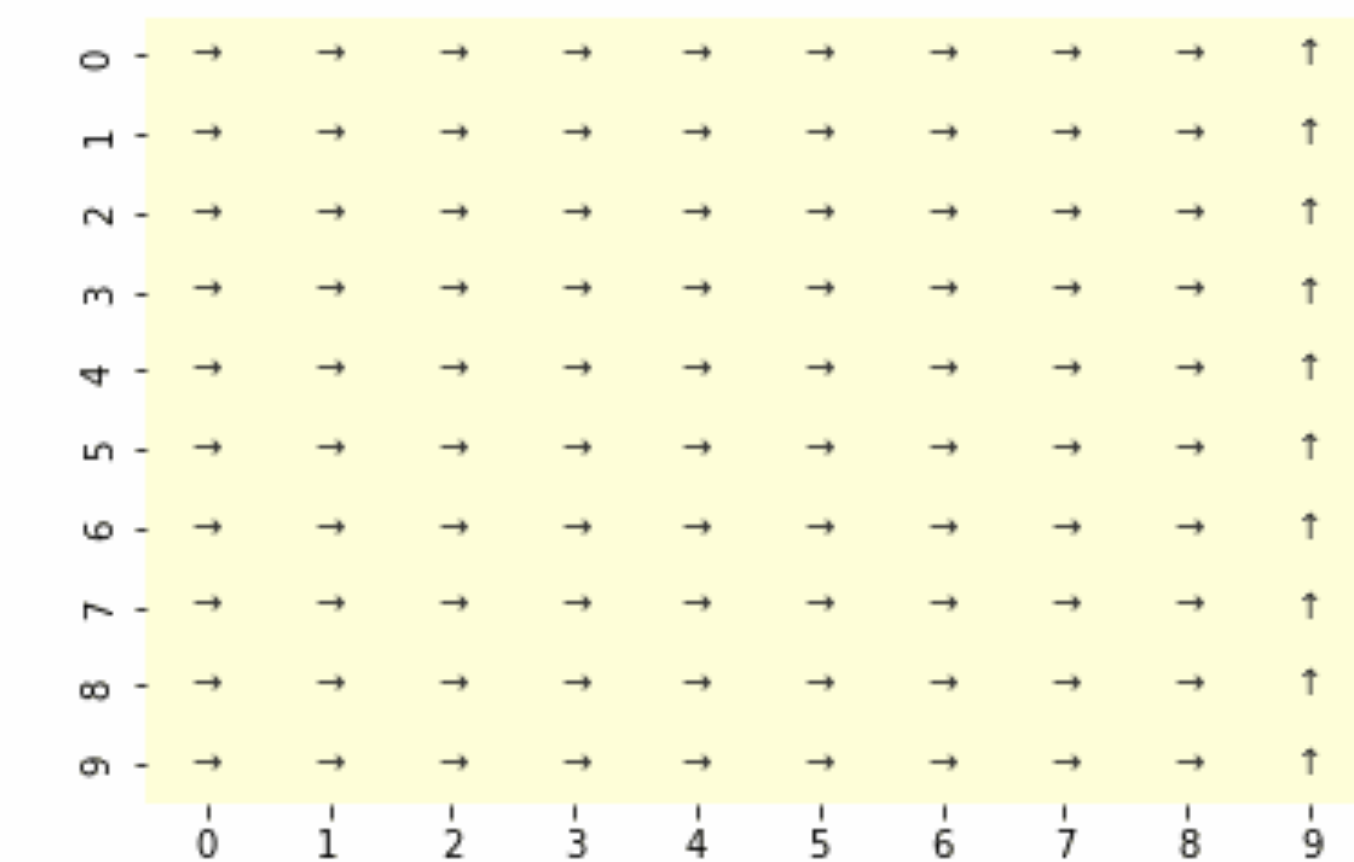
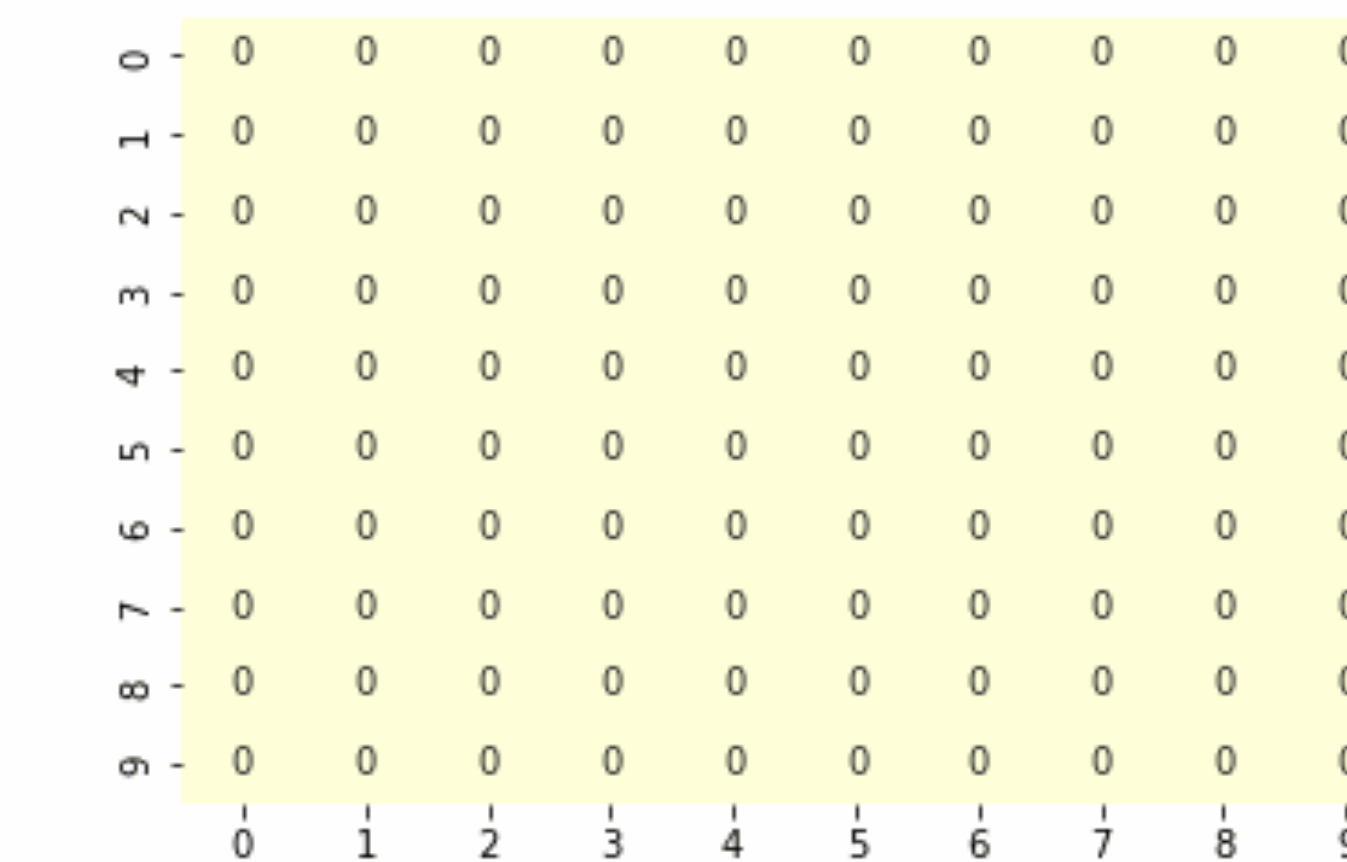
Time: 29



$$V^*(s_t) = \min_a [c(s_t, a) + V^*(s_{t+1})]$$

$$\pi^*(s_t) = \arg \min_a [c(s_t, a) + V^*(s_{t+1})]$$

Iter: 0



Policy Iteration

Evaluate value of current policy,
then improve

min() operator in policy improvement

$$V^\pi(s) = c(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')$$

$$\pi^+(s) = \arg \min_a [c(s, a) + \gamma \mathbb{E}_{s' \sim \mathcal{T}(s, a)} V^\pi(s')]$$



How do we scale these approaches?

For continuous MDP (but linear)?

For non-linear MDP?

Handle constraints?

The LQR Algorithm

Initialize $V_T = Q$

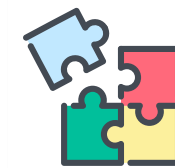
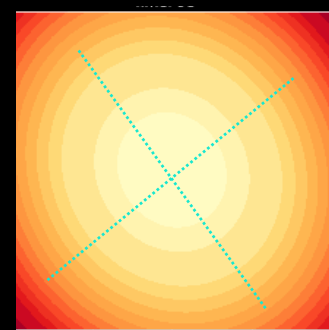
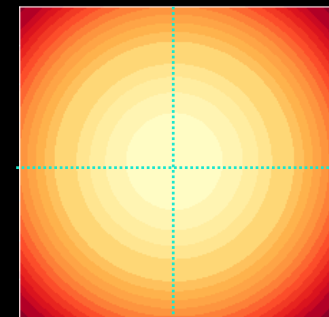
For $t = T \dots 1$

Compute gain matrix

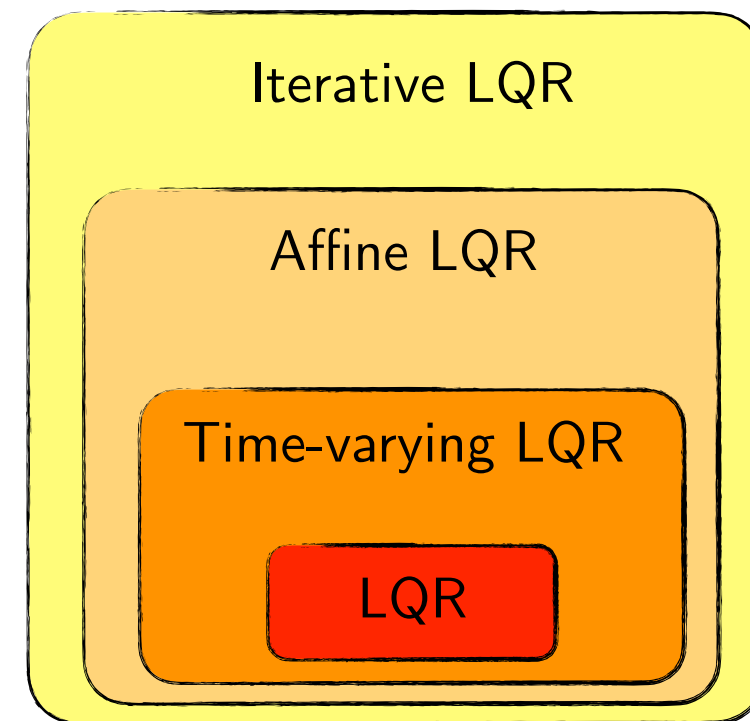
$$K_t = (R + B^T V_{t+1} B)^{-1} B^T V_{t+1} A$$

Update value

$$V_t = Q + K_t^T R K_t + (A + B K_t)^T V_{t+1} (A + B K_t)$$



Strategy: Build up on LQR



$$x_{t+1} = \frac{\partial f}{\partial x} \Big|_{x_t} \delta x_t + \frac{\partial f}{\partial u} \Big|_{u_t} \delta u_t + f(x_t^*, u_t^*)$$

$$x_{t+1} = A_t x_t + B_t u_t + x_t^{off}$$

$$x_{t+1} = A_t x_t + B_t u_t$$

$$x_{t+1} = A x_t + B u_t$$

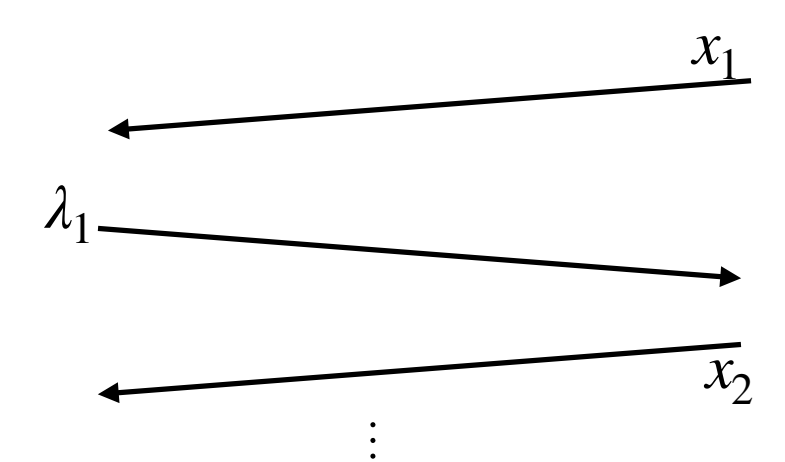
Dual Game: We control lambdas!

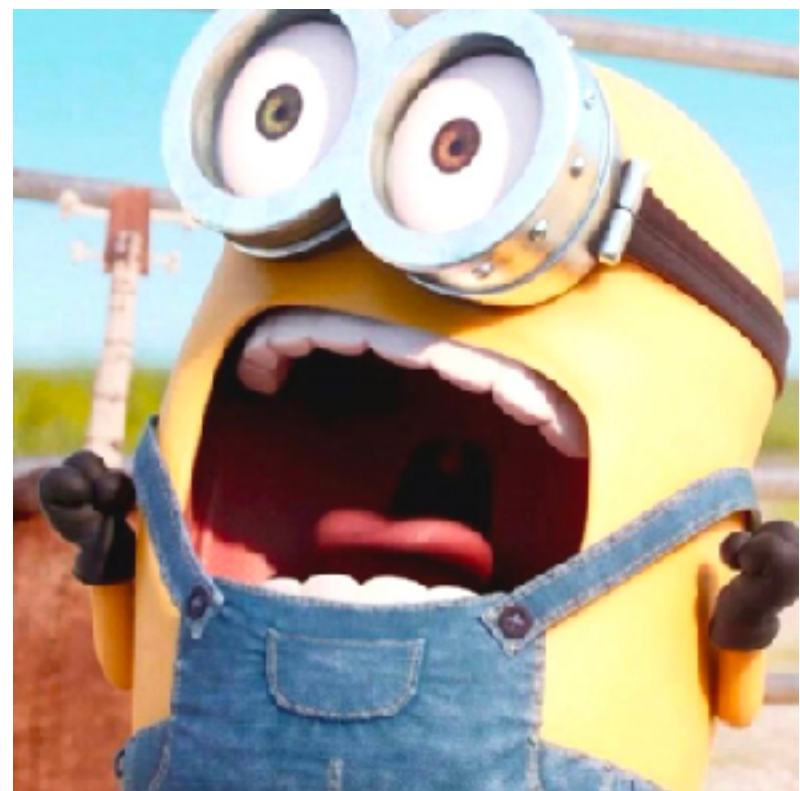
$$\min_x \max_{\lambda} f(x) - \lambda^T g(x)$$

Dual λ



Primal x



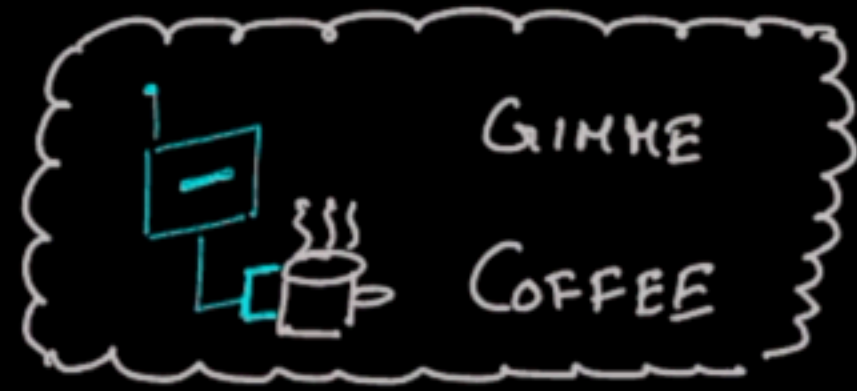


... What if your MDP is really complex?

Large state space, stochastic, continuous actions ...



$$V^*(s_t) = \min_{a_t} [C(s_t, a_t) + \gamma \sum_{s_{t+1} \sim P(\cdot | s_t, a_t)} V^*(s_{t+1})]$$



HUMAN

$C(s_t, a_t)$



s_{t+1}



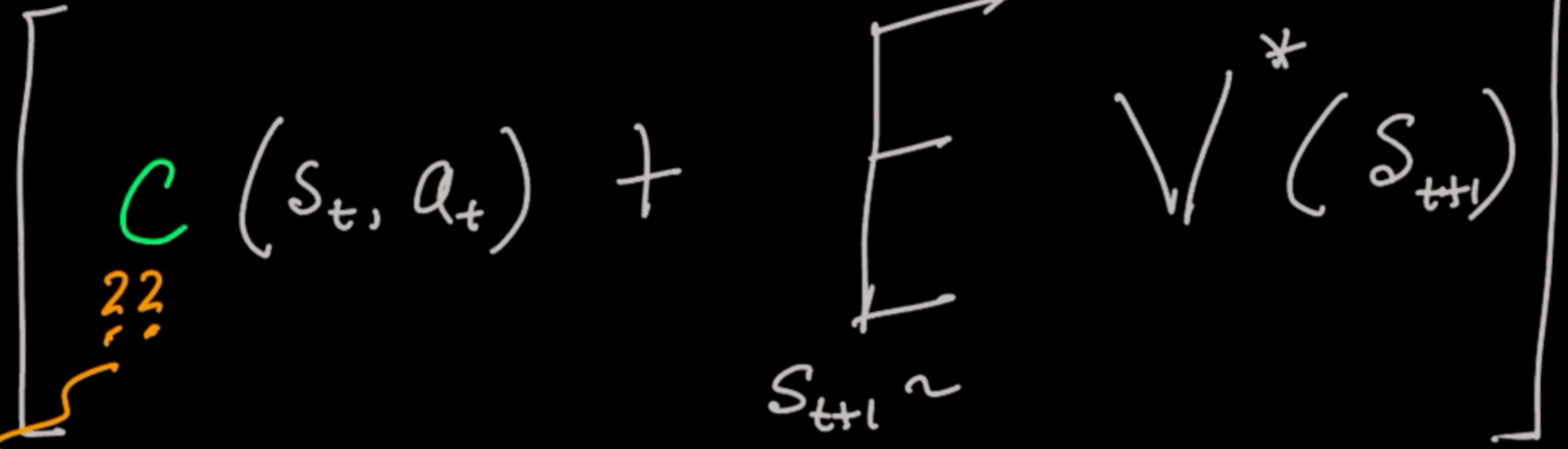
WORLD

$s_t, a_t?$

$s_t, a_t?$

HUMAN +
WORLD

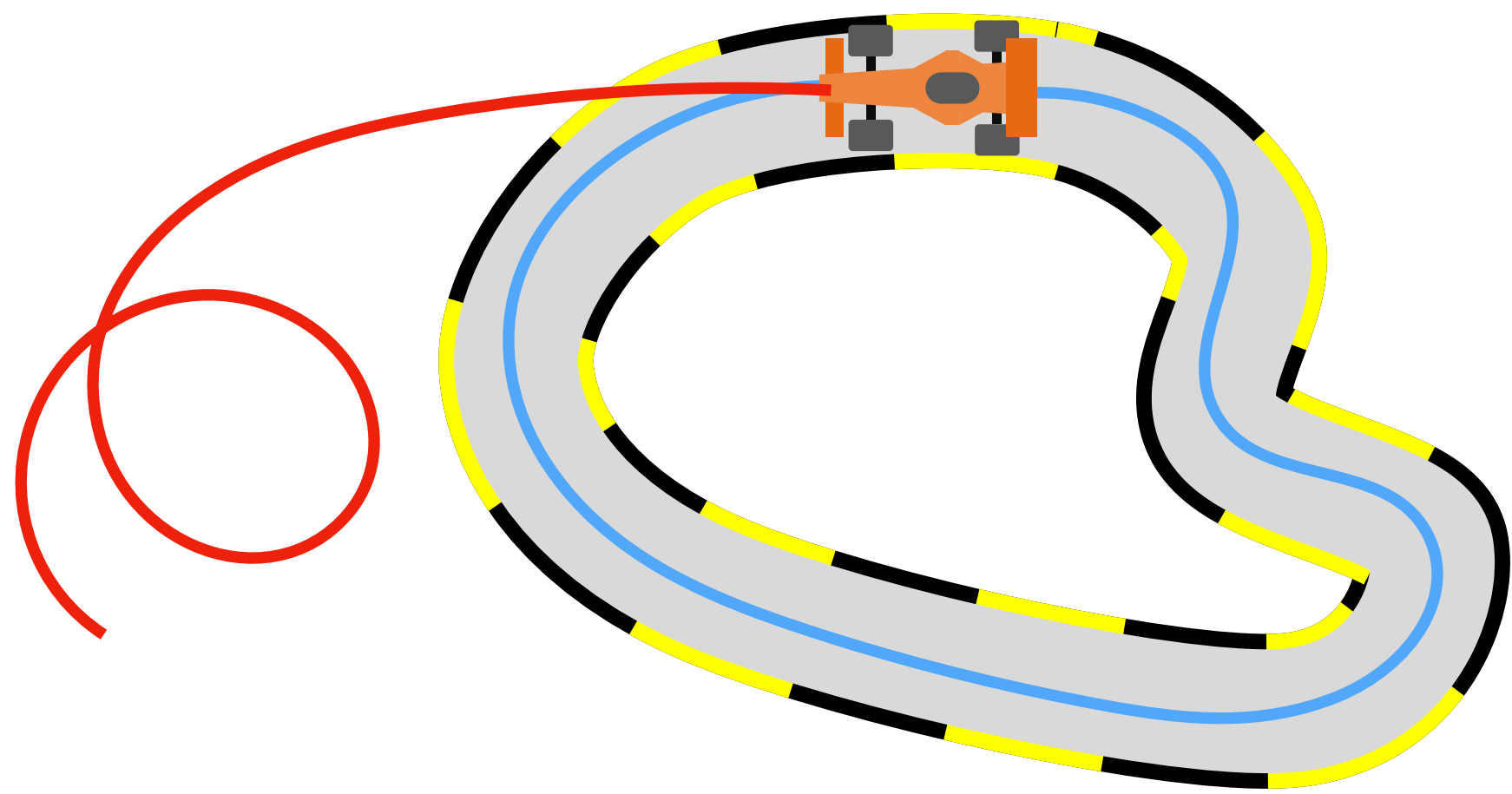
LEARN OPTIMAL VALUE BY INTERACTIONS W/



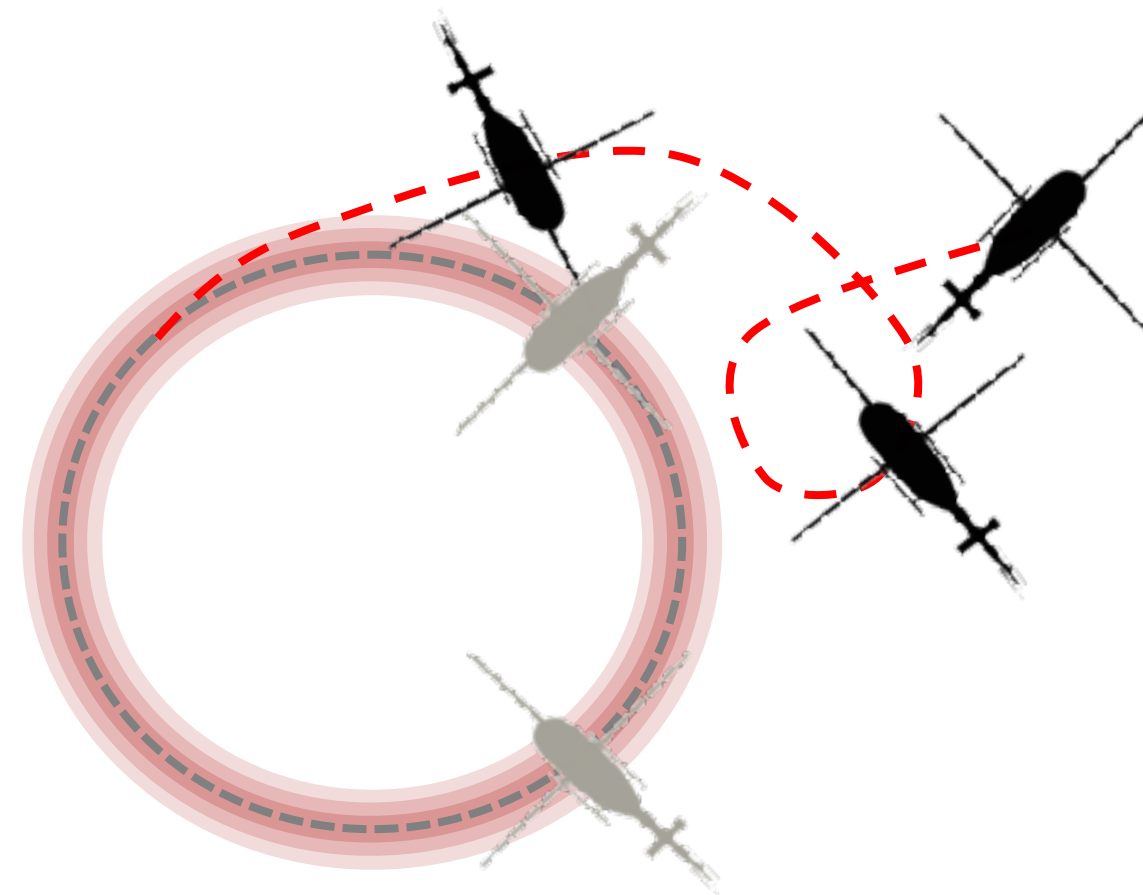


Where all did we see covariate shift?

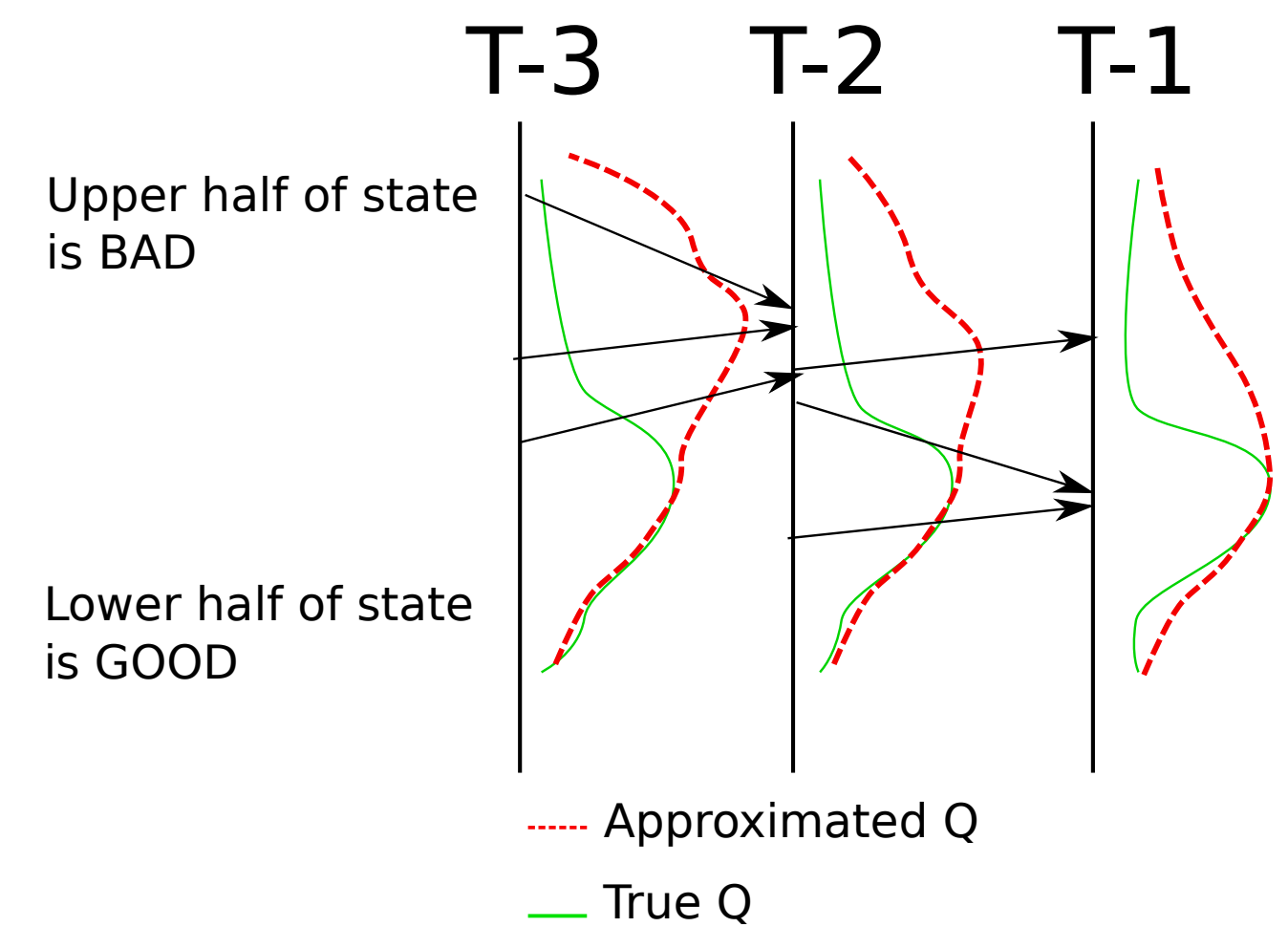
Imitation Learning?



Model Based RL?

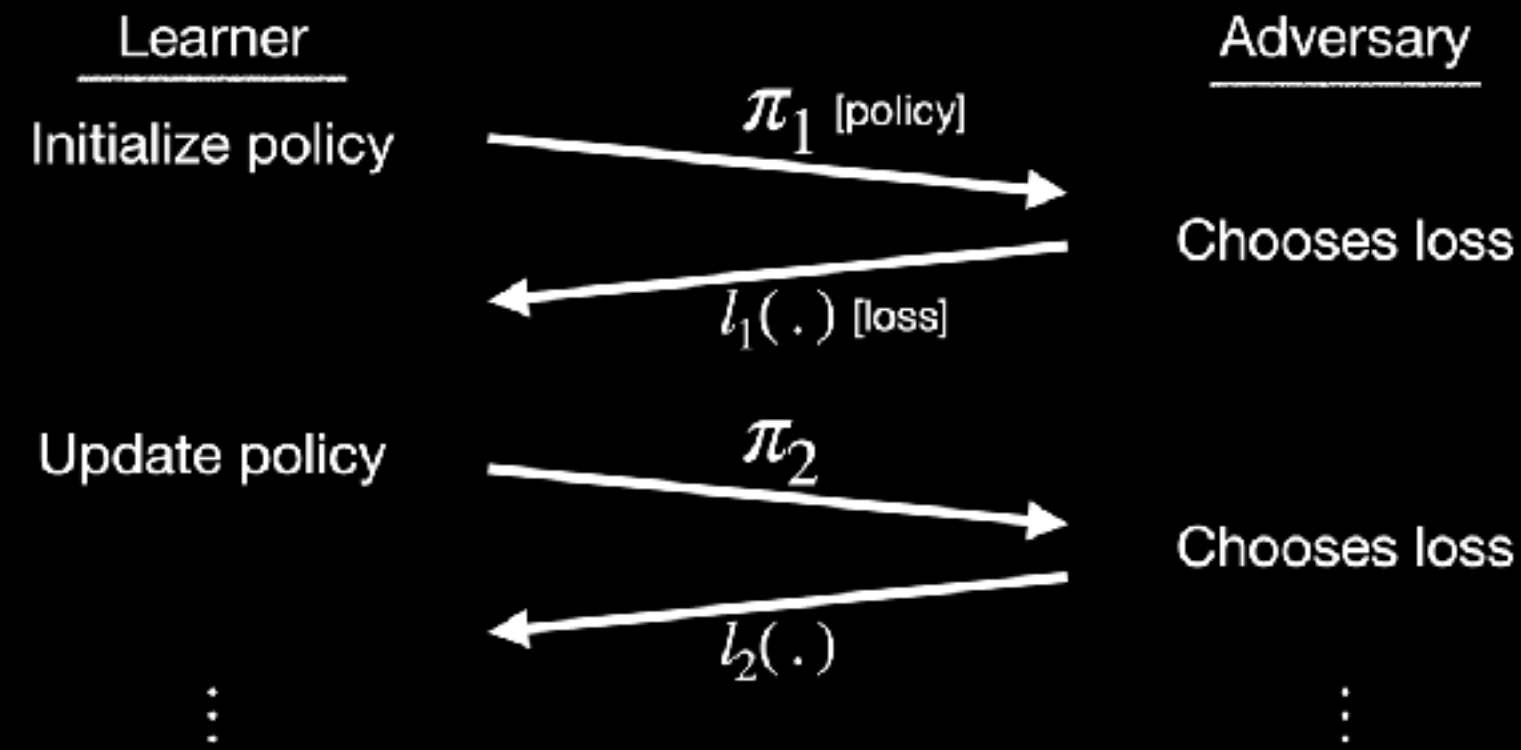


Approximate
Dynamic
Programming?





Interactive Learning



Learning is
a Game!

Follow the leader
is aggressive

Slowly change predictions,
achieve no-regret



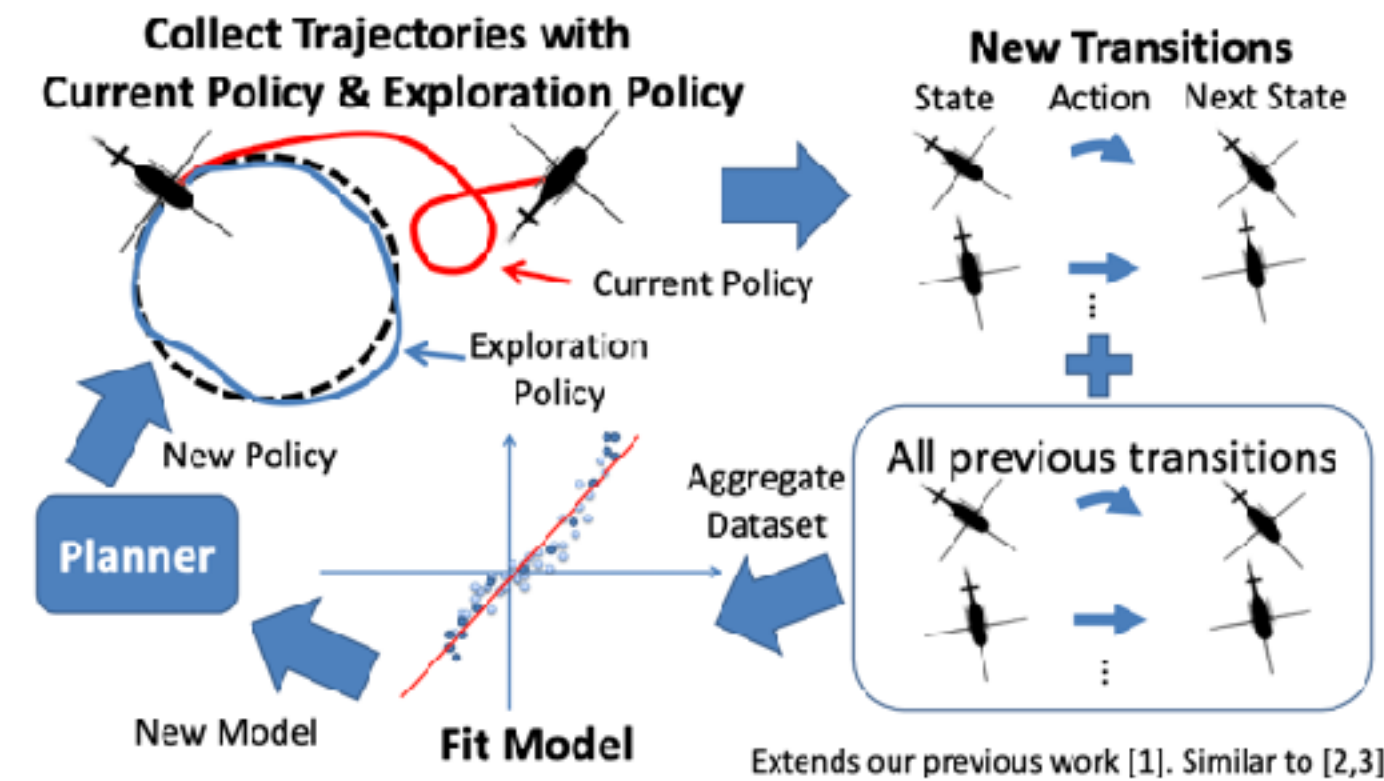
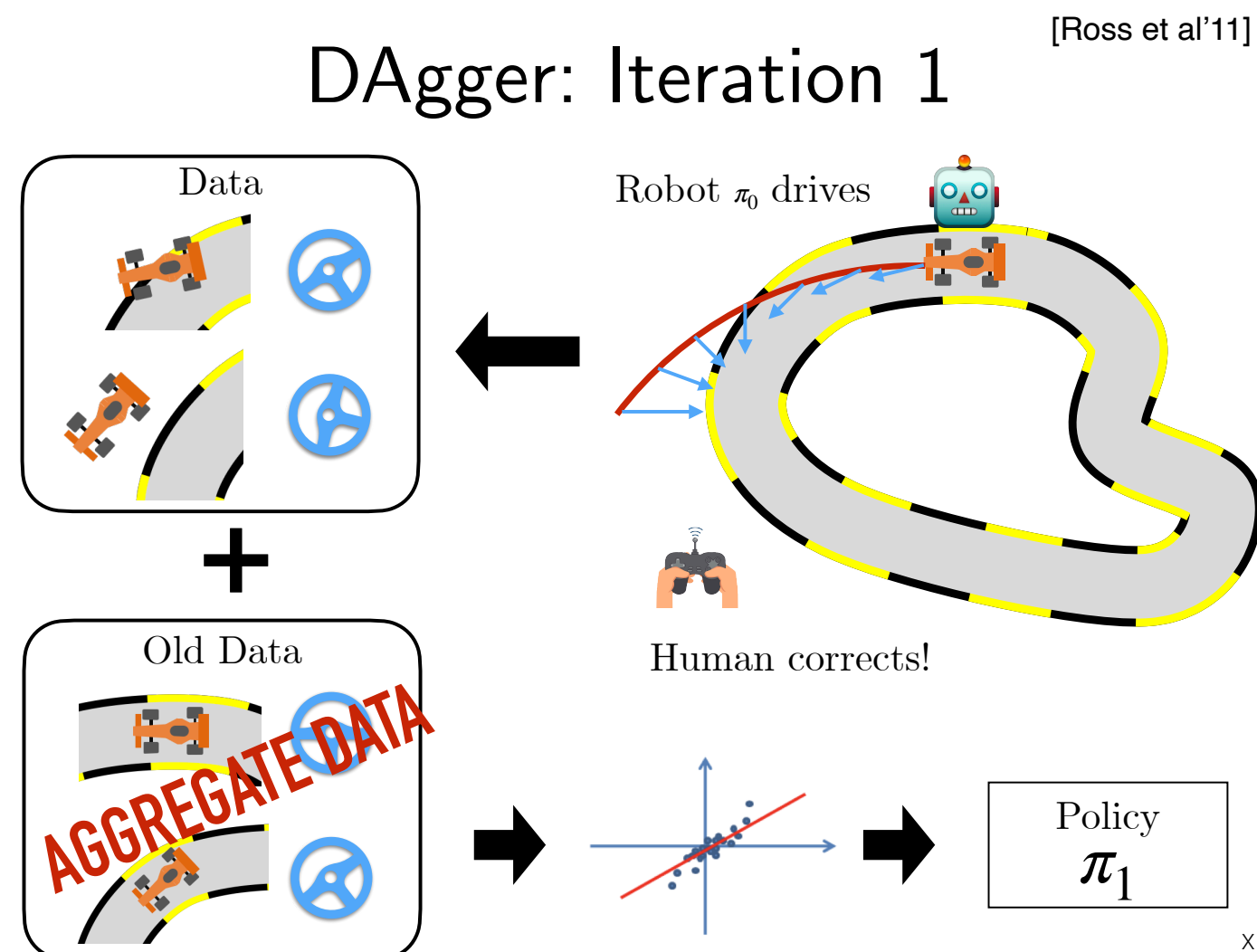


No regret solves all!

DAgger

DAgger for SysID

Conservative policy iteration



Idea 1: Conservative Policy Iteration (CPI)

$$\pi' = (1 - \alpha)\pi + \alpha\pi_{greedy}$$

Mix in old policy and greedy policy

Can prove that performance difference is bounded by

$$V^{\pi'}(s) - V^{\pi}(s) \geq \alpha A_{greedy} - 2\alpha^2 \frac{\gamma}{1 - \gamma}$$

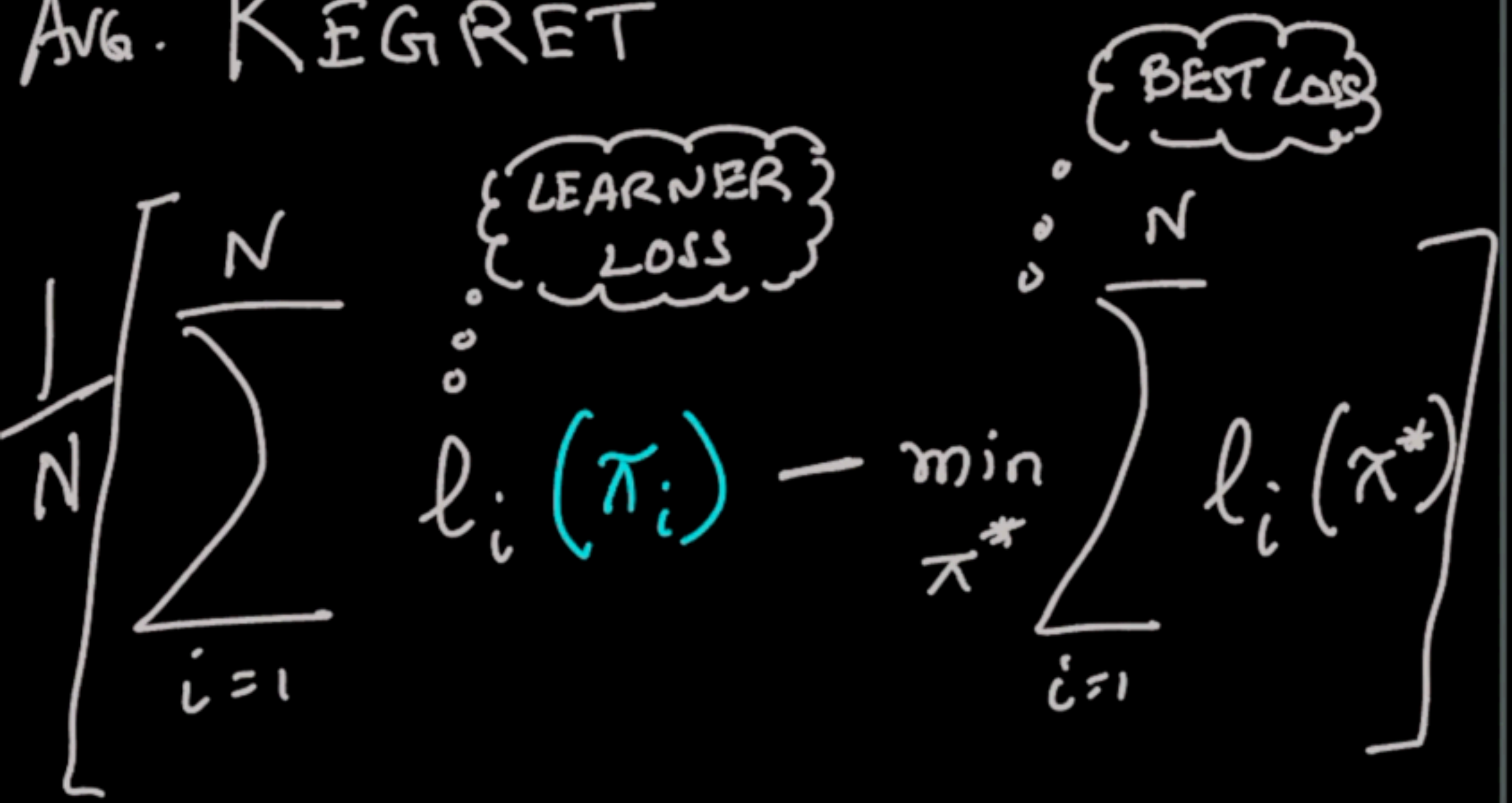
How much greedy policy improves based on estimate

How much distribution shift hurts!

Approximately Optimal Approximate Reinforcement Learning
 Shao-Kai Shiu
 Google DeepMind, University of Cambridge, London WC2N 3AR, UK
 John Langford
 Computer Science Department, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15261

GOAL: MINIMIZE

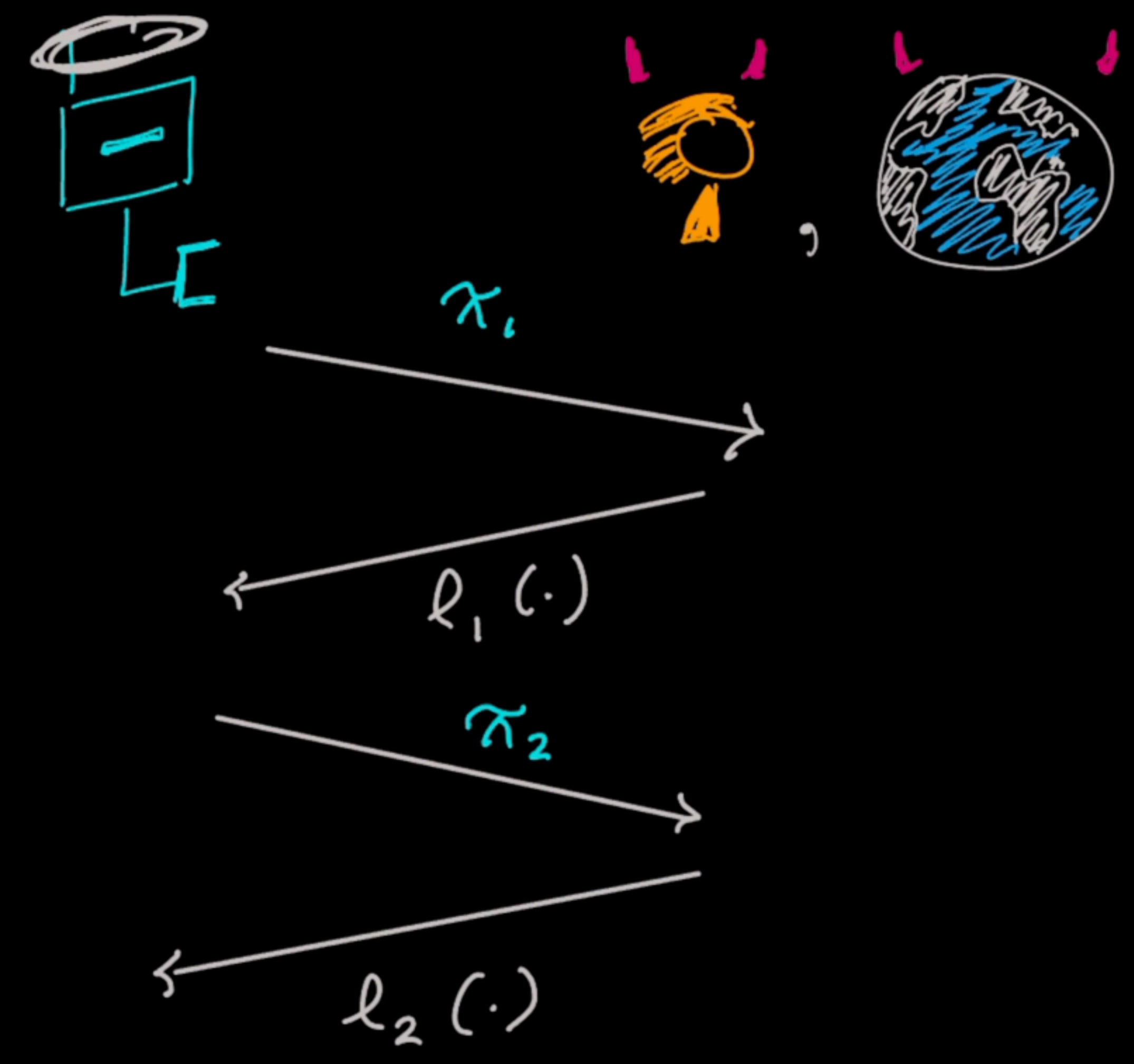
AVG. REGRET



As $N \rightarrow \infty$, Avg REGRET $\rightarrow 0$

"DO AS WELL AS YOU COULD HAVE IF YOU HAD ALL THE DATA UPFRONT"

INTERACTIVE ONLINE LEARNING



Reinforcement Learning: Brass Tacks

We don't know the MDP, all we see are **traces** (s, a, s')

Model Based:

Learn a model. Plan with the model.

Model Free:

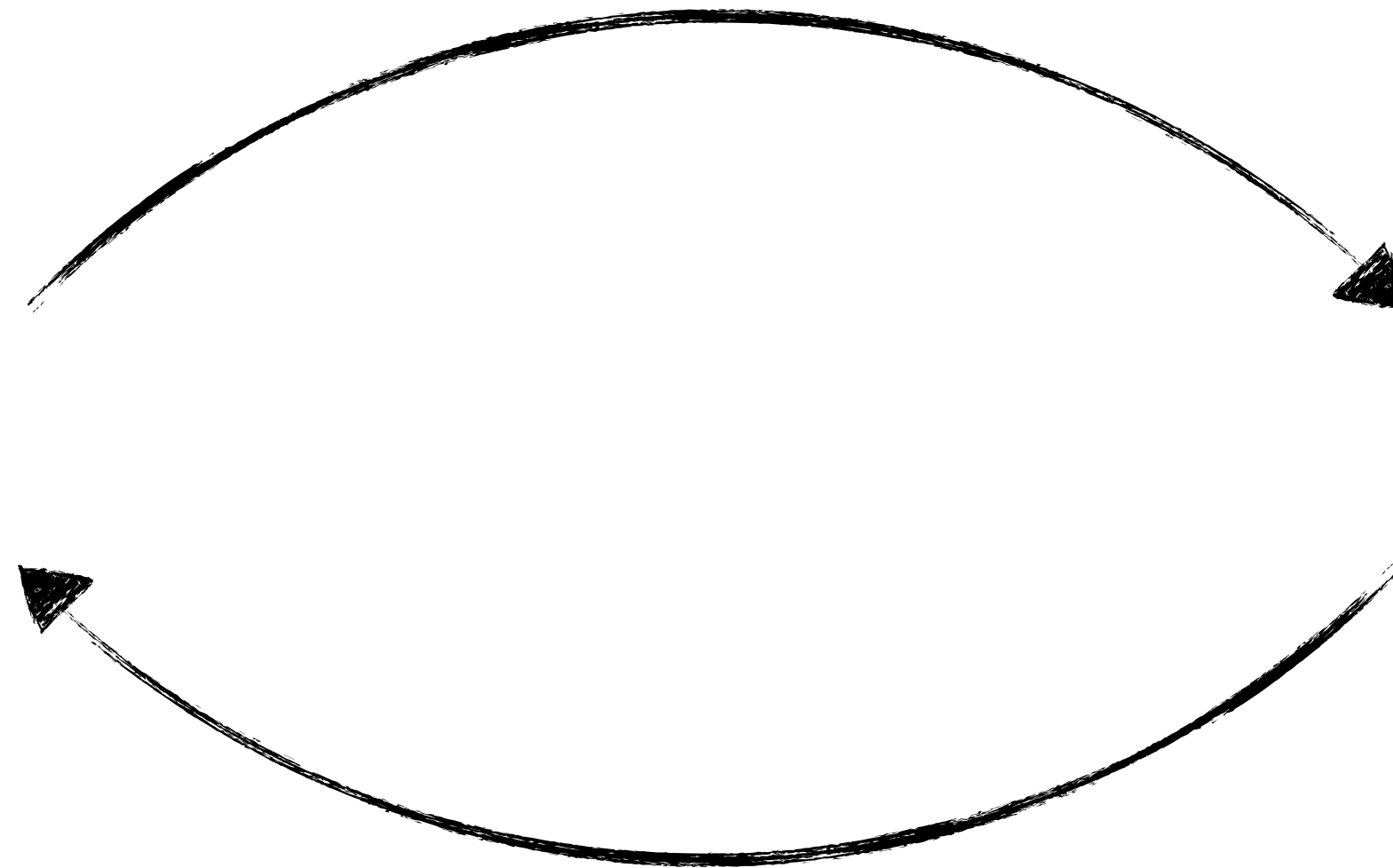
Forget about models. Learn the policy.

Model Free RL: Actor Critic

Actor



Critic



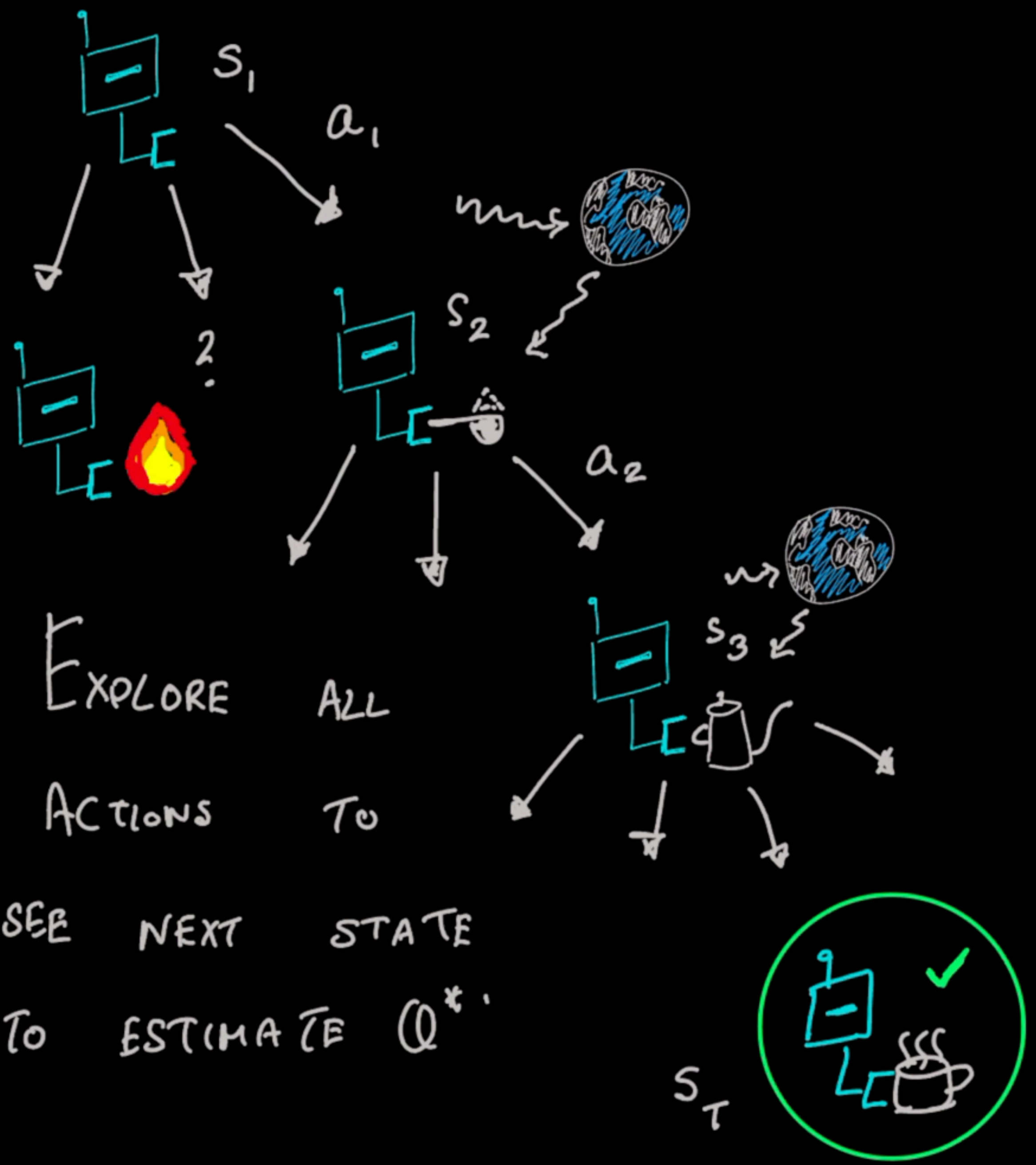
Policy improvement
of π

Estimates value

functions $Q_{\phi}^{\pi} / V_{\phi}^{\pi} / A_{\phi}^{\pi}$

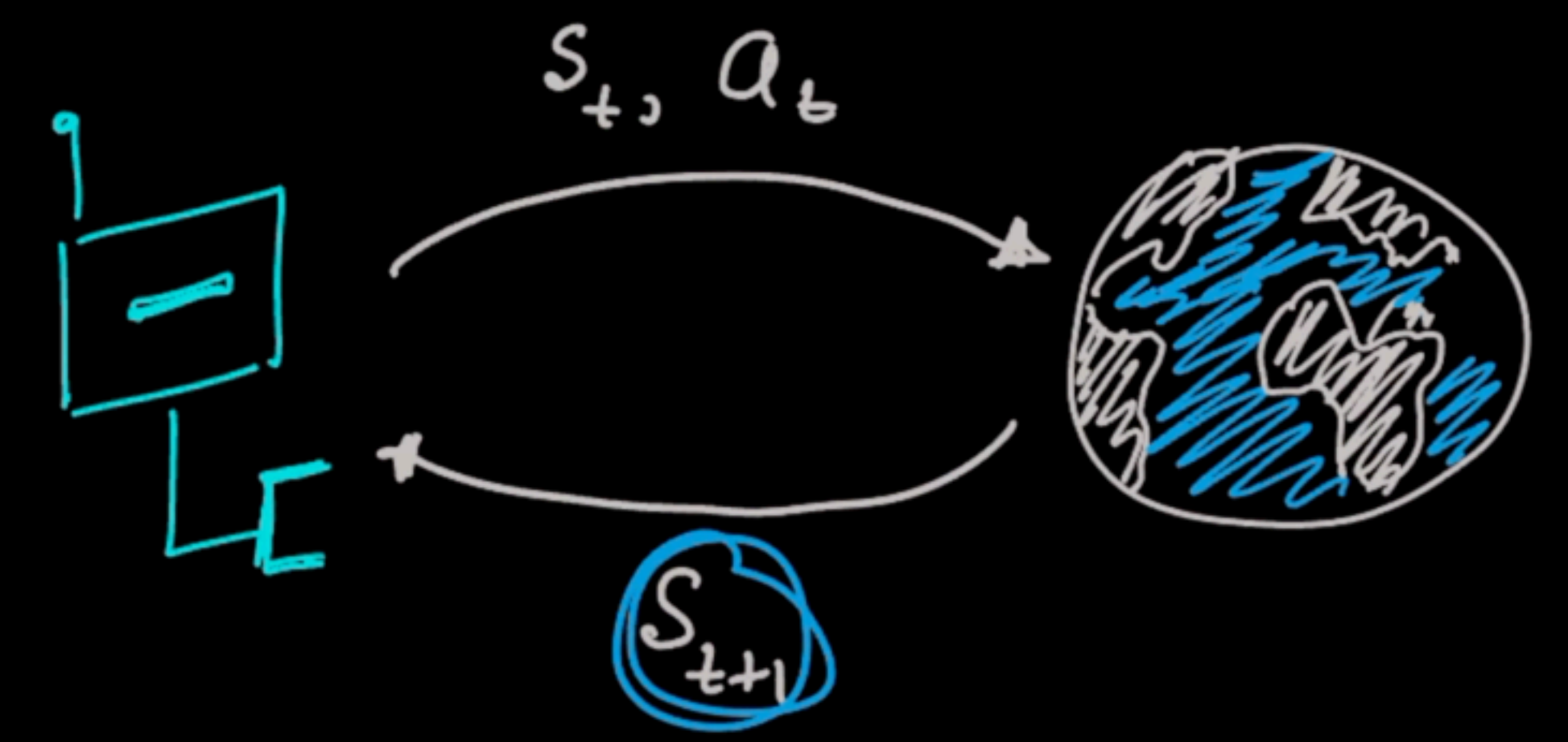
(Natural) Policy Gradient

TD, MC



EXPLORE ALL ACTIONS TO SEE NEXT STATE TO ESTIMATE Q^* .

REINFORCEMENT LEARNING



ESTIMATE $Q^*(s_t, a_t)$ SUCH THAT

$$Q^*(s_t, a_t) = c(s_t, a_t)$$

i.e. BELLMAN CONSISTENT + $E \min_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$

Imitation Learning: Brass Tacks

We don't know the MDP, all we see are **human actions (a^*)**

Learn Cost:

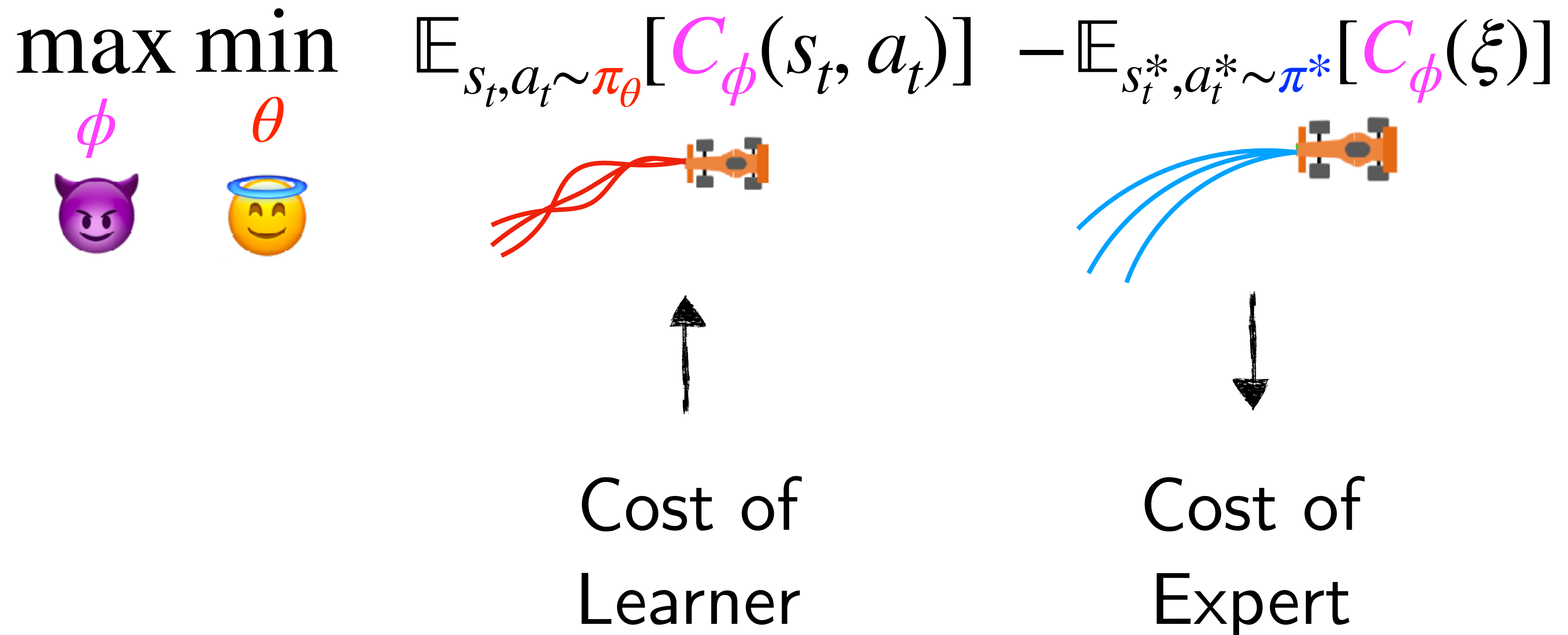
Learn a cost that makes human look cheap, learner look expensive

Learn Values:

Learn Q^* that makes human look cheap, learner look expensive

Inverse Optimal Control (Learn Cost)

Make human look cheap, learner look expensive

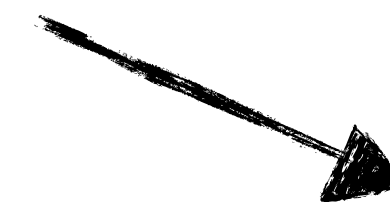


Learn Values

Estimate Q^* from demonstrations, interventions, preferences, ..
and even E-stops!



Demonstrations



Interventions



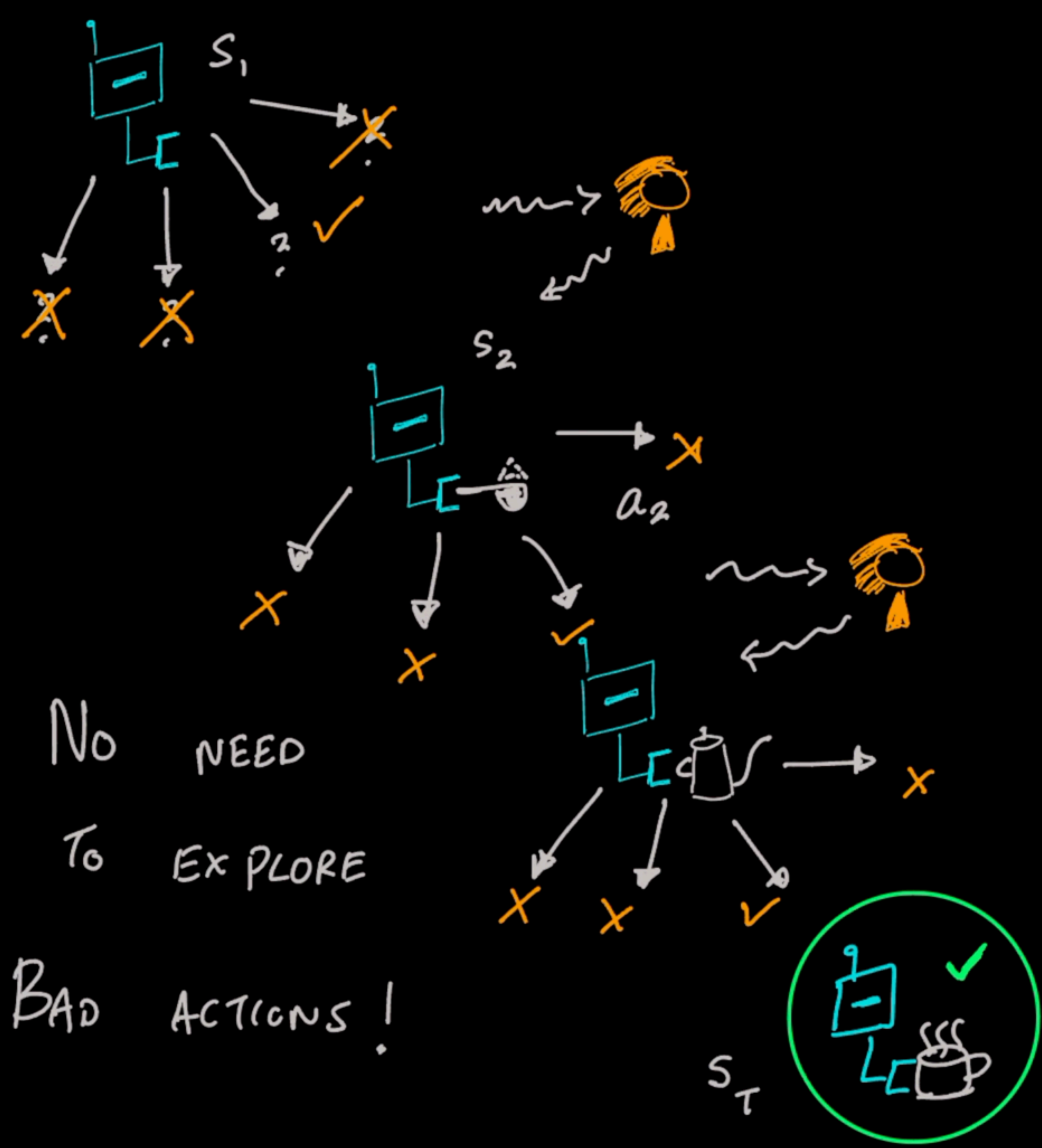
Preferences



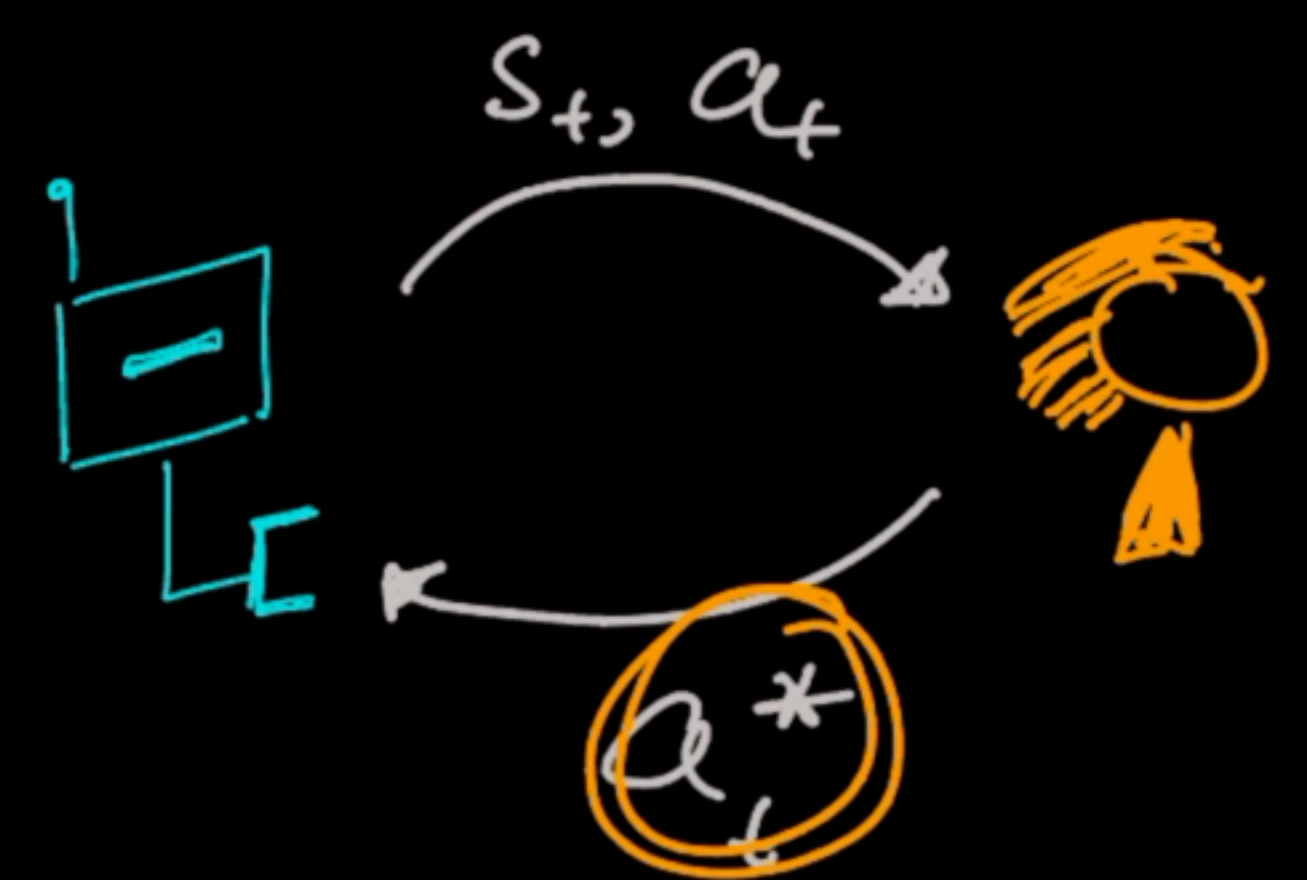
E-stops



$\mathcal{L}(Q_\theta^*)$
Loss



IMITATION LEARNING



ESTIMATE

$Q^*(s_t, a_t)$ SUCH THAT

$$Q^*(s_t, a_t) = \min_{a_t'} Q^*(s_t, a_t')$$

i.e. HUMAN IS OPTIMAL

The Imitation Game

We have an interactive expert.

Apply PDL in forward direction: roll-in learner, roll-out expert

$$\min_{\pi} \max_{Q^*} \sum_{t=1}^T \mathbb{E}_{s_t \sim \pi} [Q^*(s_t, \pi(s_t)) - Q^*(s_t, \pi^*(s_t))]$$



Use no-regret learning to solve the game!



$$O(\epsilon T)$$

The RL Game

We don't have interactive expert.

Apply PDL in reverse direction: roll-in expert, roll-out learner

$$\min_{\pi} \max_{Q^\pi} \sum_{t=1}^T \mathbb{E}_{s \sim \pi^*} Q^\pi(s, \pi(s)) - Q^\pi(s, \pi^*(s))$$

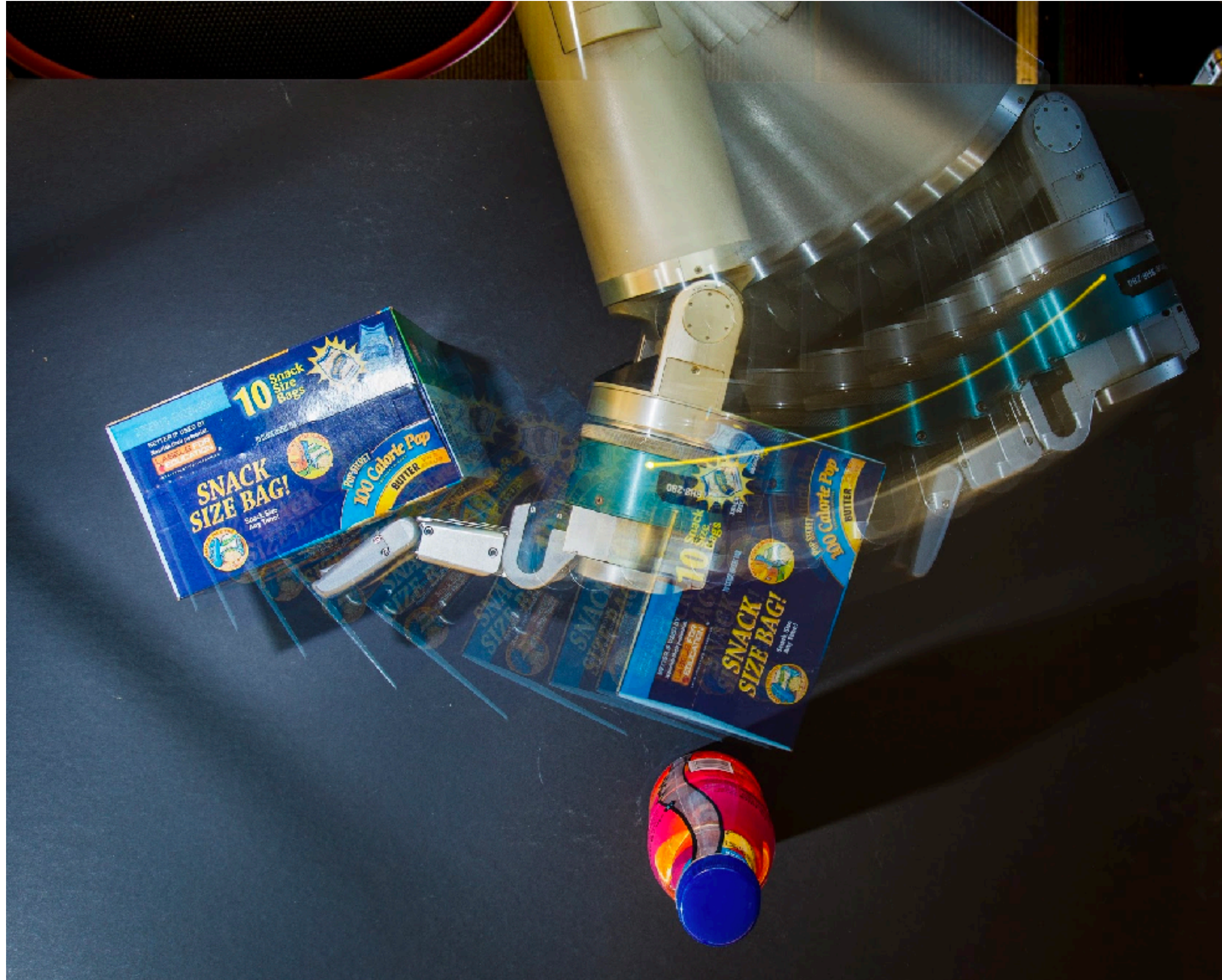
 

Use no-regret learning to solve the game! $O(\epsilon T^2)$

A grand unification of IL / RL Games?



A simple question:
Can learning help us build better planners?

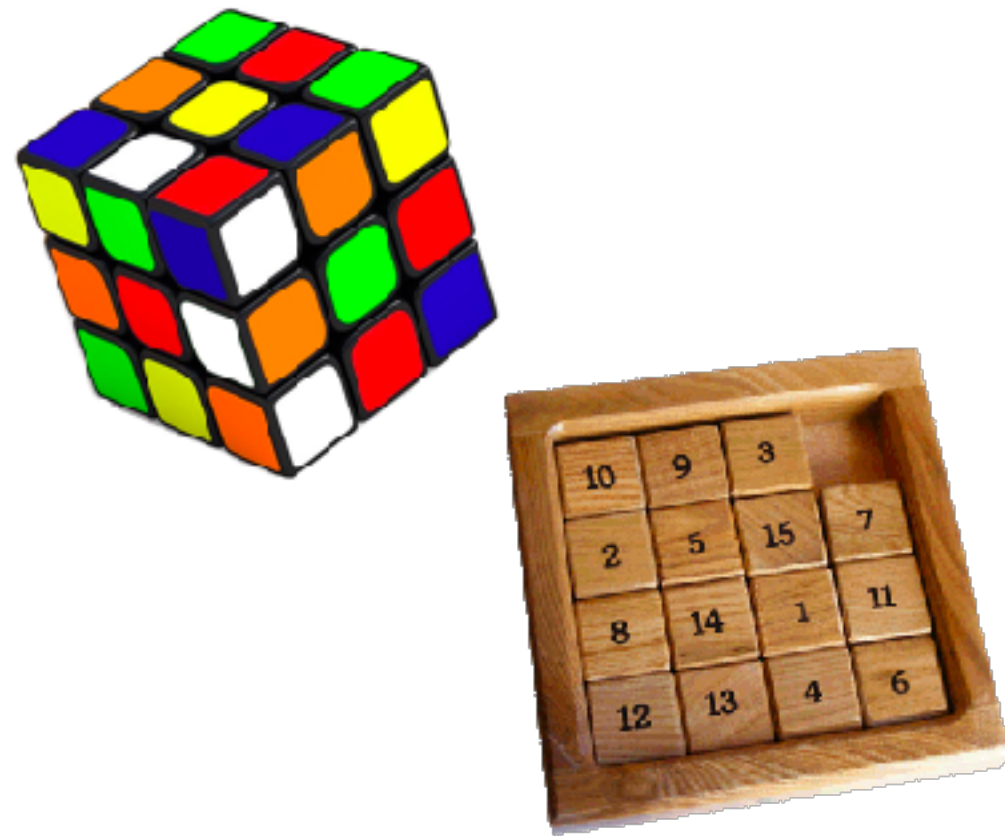


A prospective grad student:
“Is planning just A*?”



Motion Planning: Dealing with expensive collision checking

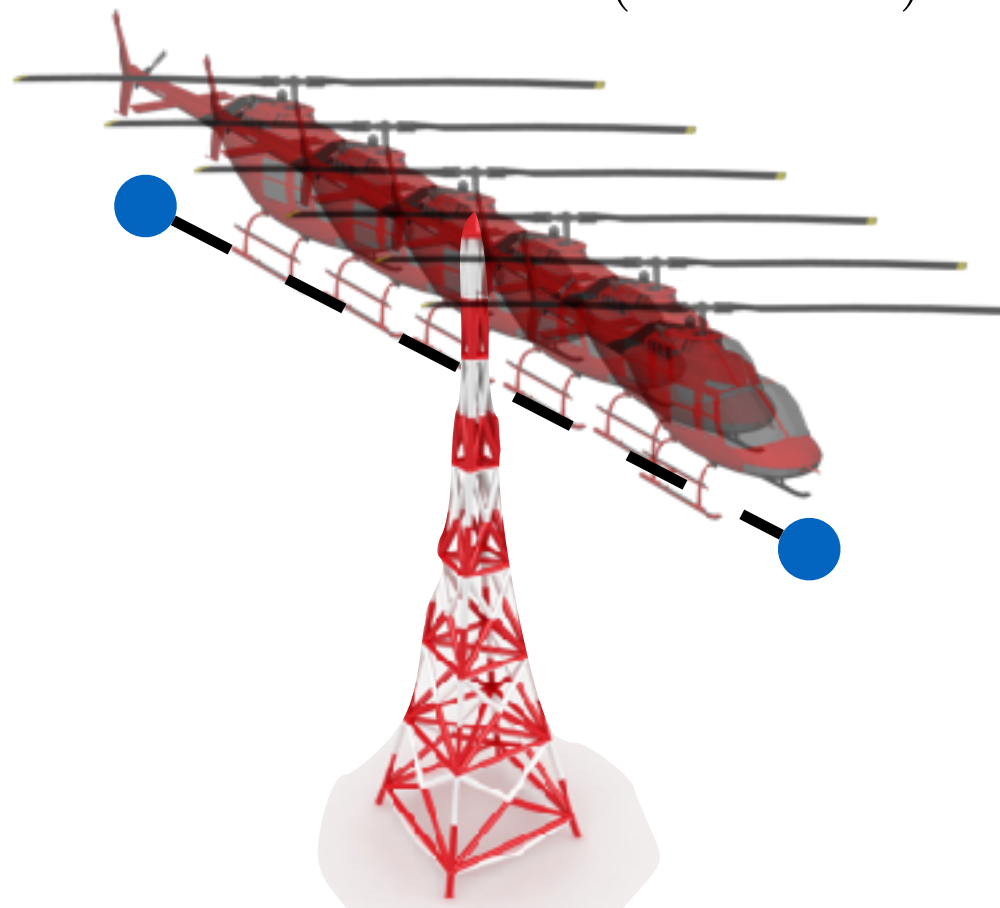
Trivial



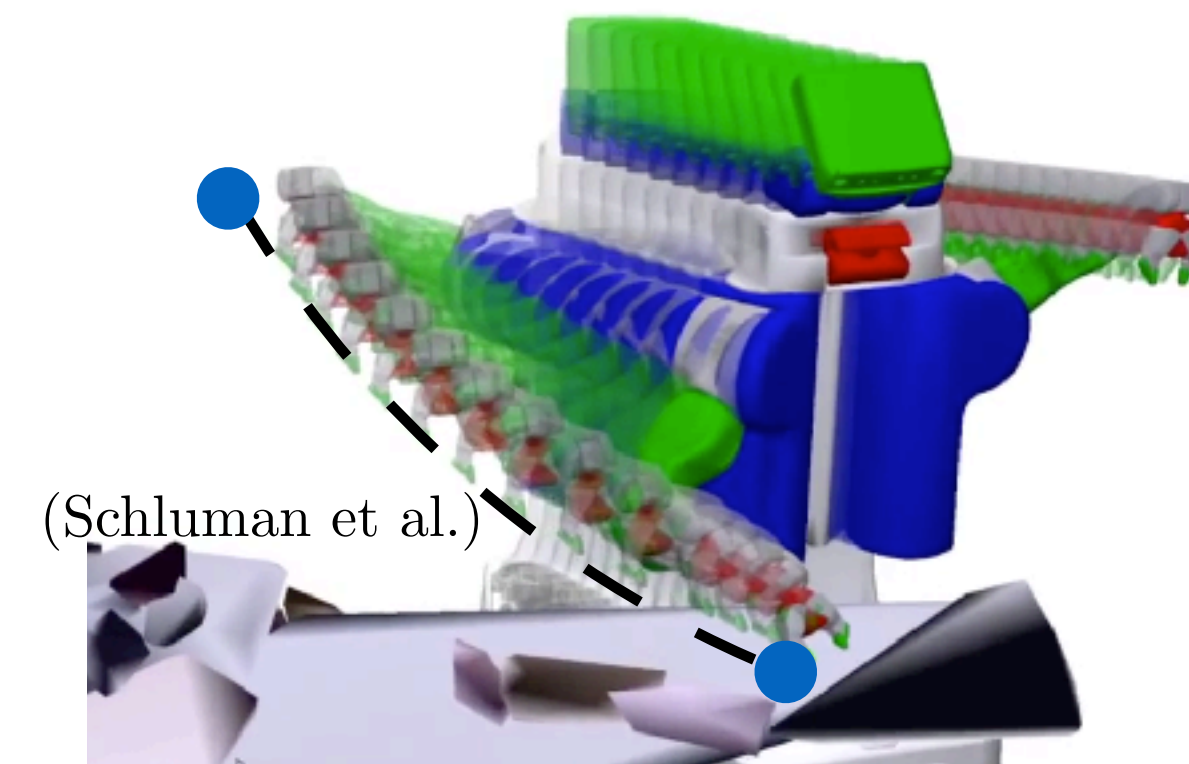
Medium



(Ross et al.)



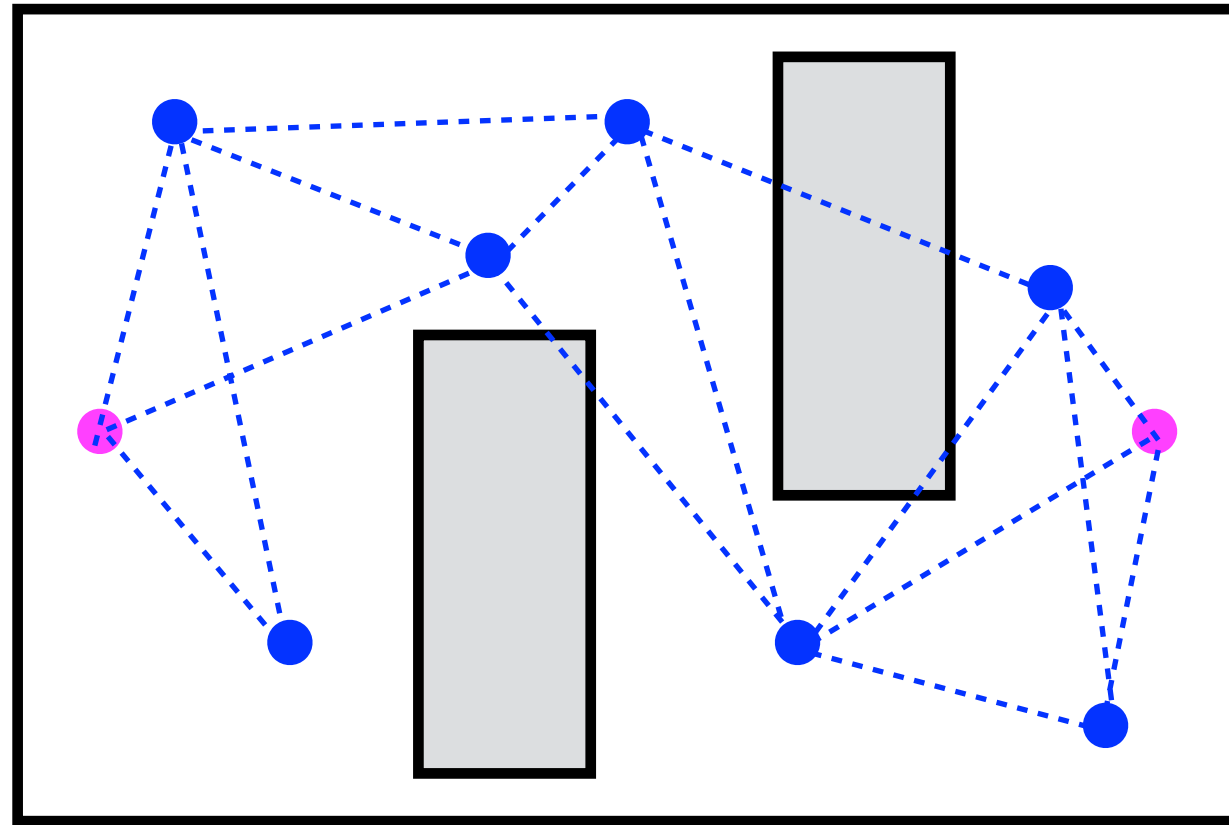
Expensive



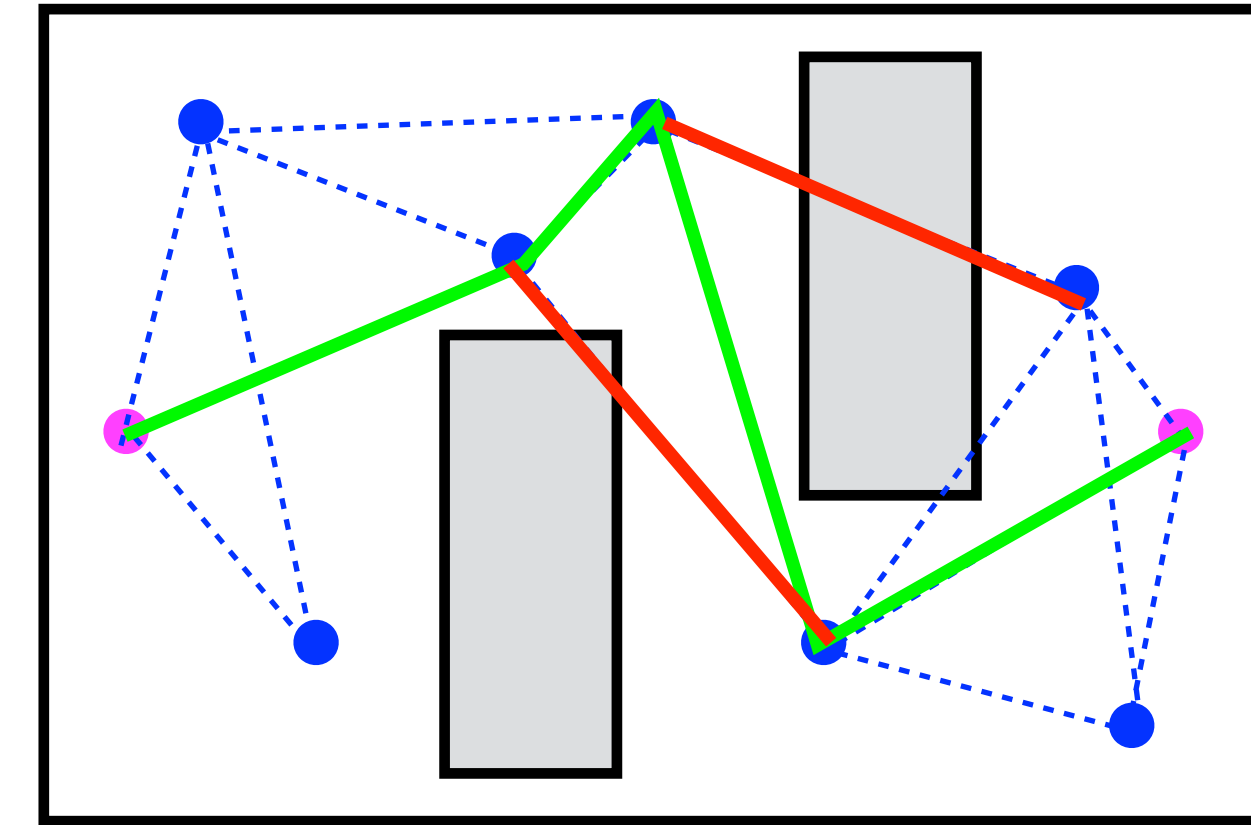
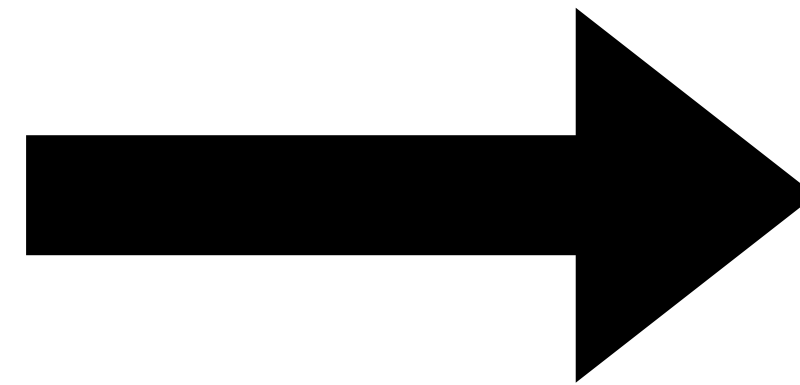
(Schluman et al.)

(LaValle'06, Bialkowski'11,
Hauser'15,)

General framework for motion planning



Create a graph



Search the graph



Interleave

General framework for motion planning

Any planning algorithm

Create graph

Search graph

Interleave



=

e.g. fancy
random
sampler

×

e.g. fancy
heuristic

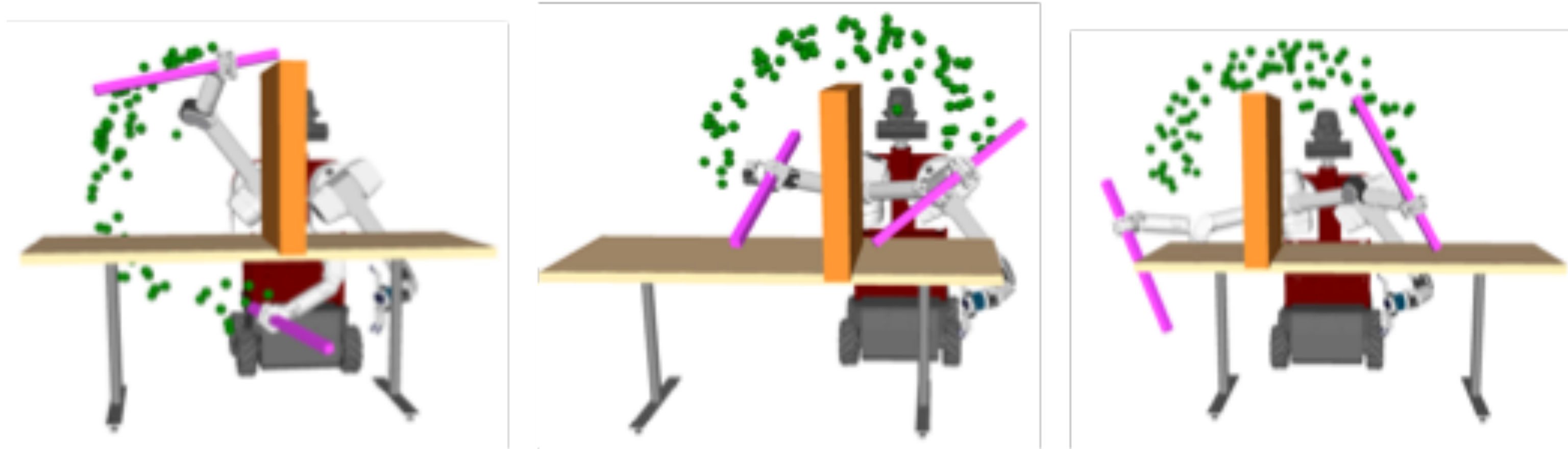
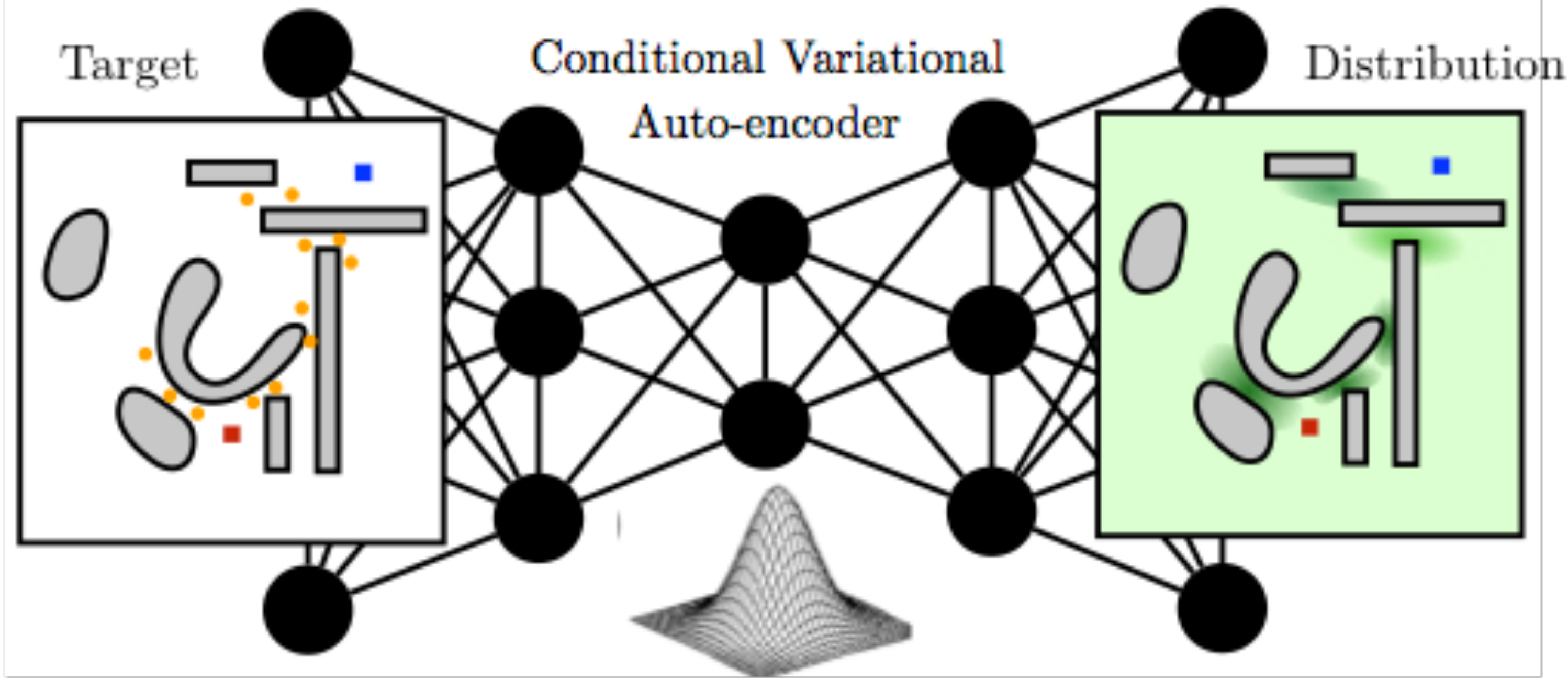
×

e.g. fancy
way of
densifying

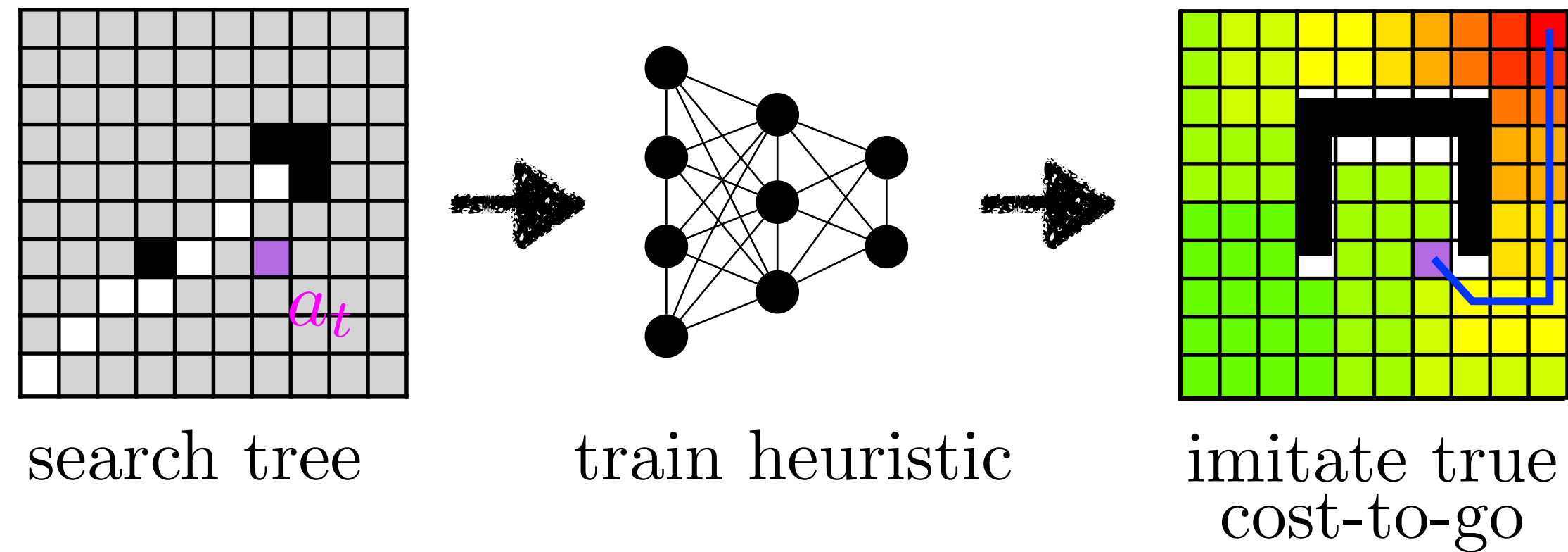
Learn
sampler!

Learn
heuristic!

Learning a Sampler



Learn a Heuristic



				Ours	Baseline Learning			Baseline Handcrafted			
				SAIL	SL	CEM	QL	h_{EUC}	h_{MAN}	A*	MHA*
alternating_gap				0.039	0.432	0.042	1.000	1.000	1.000	1.000	1.000
single_gap				0.158	0.214	0.057	1.000	0.184	0.192	1.000	0.286
shifting_gap				0.104	0.464	1.000	1.000	0.506	0.589	1.000	0.804
forest				0.036	0.043	0.048	0.121	0.041	0.043	1.000	0.075
bugtrap_forest				0.147	0.384	0.182	1.000	0.410	0.337	1.000	0.467
gaps_forest				0.221	1.000	1.000	1.000	1.000	1.000	1.000	1.000
maze				0.103	0.238	0.479	0.399	0.185	0.171	1.000	0.279
multiple_bugtrap				0.479	0.480	1.000	0.835	0.648	0.617	1.000	0.876 ₄₃



Open Challenges