

Reinforcement Learning

Computer Engineering Department

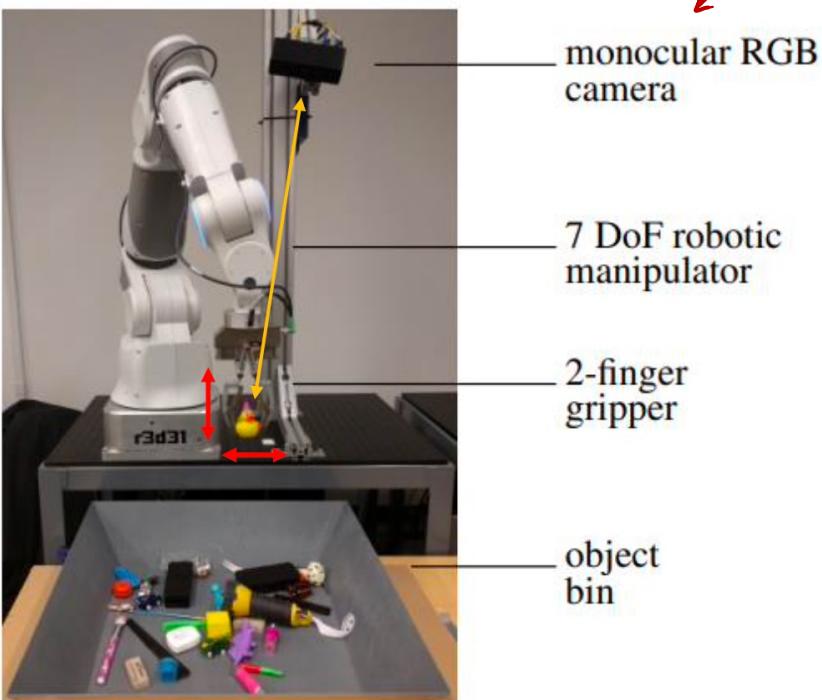
Sharif University of Technology

Mohammad Hossein Rohban, Ph.D.

Spring 2026

Courtesy: Some slides are adopted from CS 285 Berkeley, and CS 234 Stanford, and Pieter Abbeel's compact series on RL.

Motivation



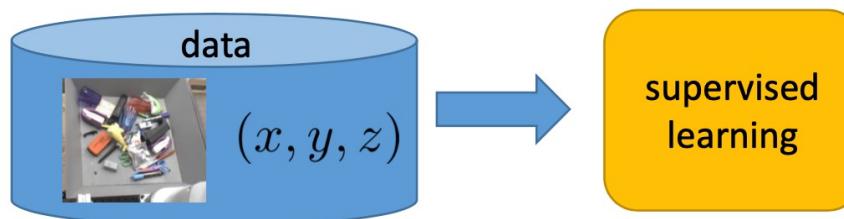
Option 1:

Understand the problem, design a solution



Option 2:

Set it up as a machine learning problem

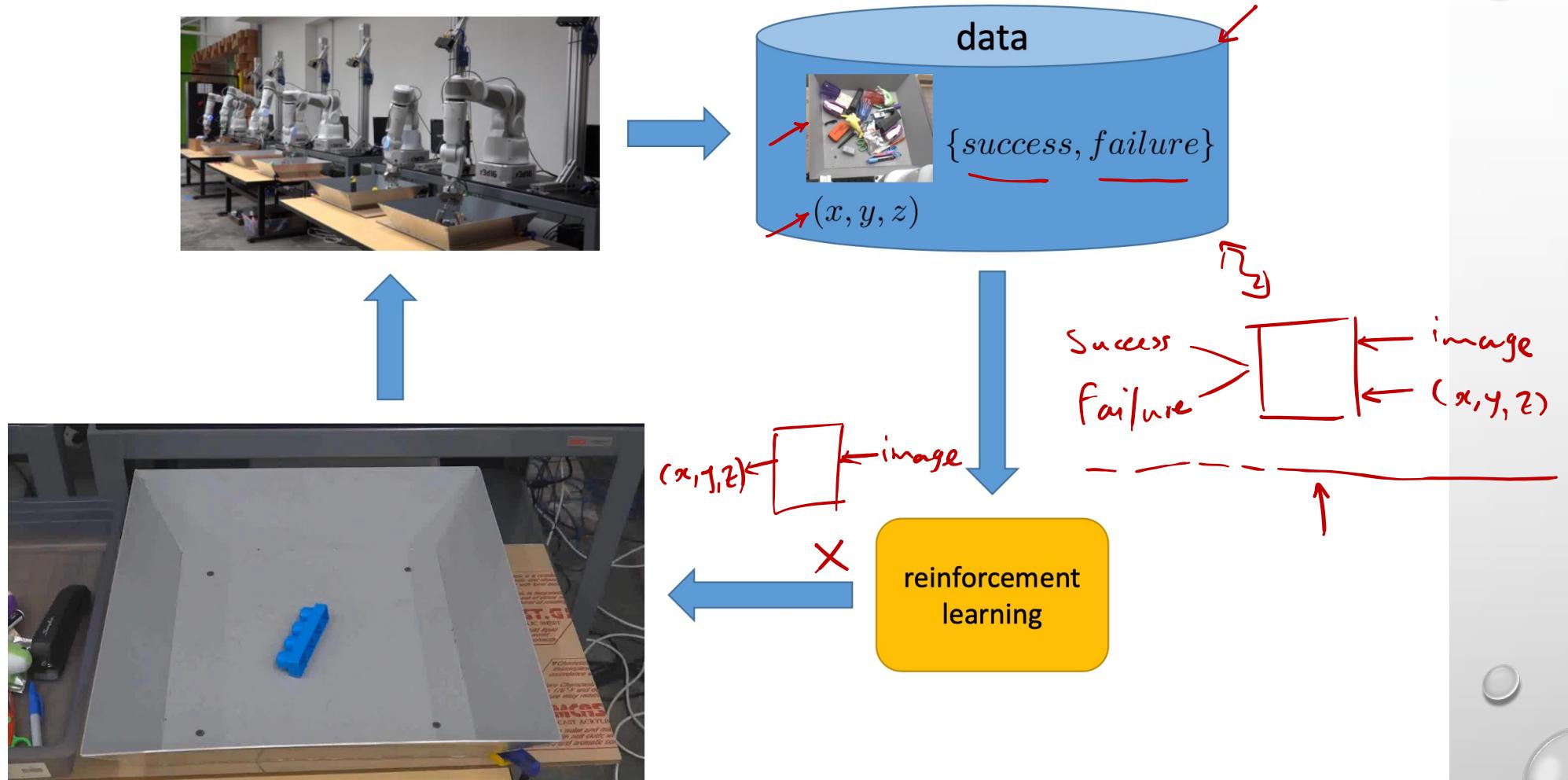


Motivation (cont.)



Courtesy: CS 285 course, Berkeley

Motivation (cont.)



Courtesy: CS 285 course, Berkeley

Motivation (cont.)

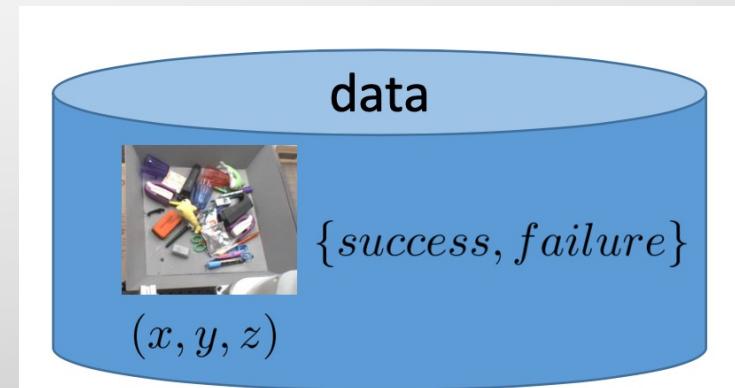
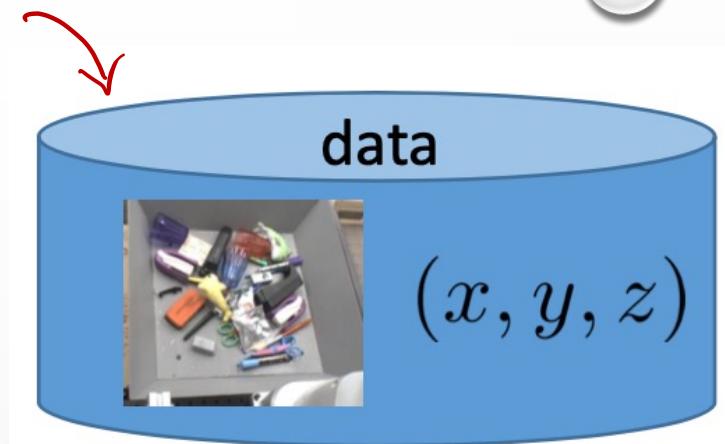
Credit Assignment

- Supervised learning:

- Ground truth is **known** in advance.
- Training data are usually **static** and **iid**.

- Reinforcement learning:

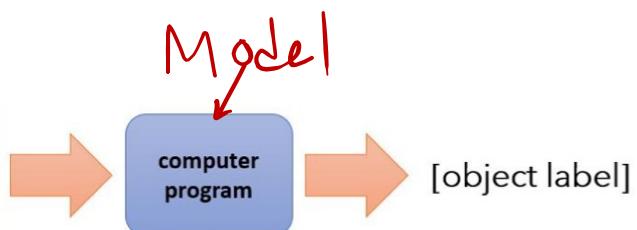
- The best action (**policy**) is usually **unknown a priori**.
- **Sequence of actions** is needed.
- A series of trial and error (**search**) is performed.
 - Usually **delayed reward** shows goodness of the trial.
- Data is **dynamic (exploration)** and **non-iid**.



What is Reinforcement Learning?

$$\max_{\theta} \mathbb{E}_{\pi_\theta} \left\{ \sum_{t=1}^T r_t \right\}$$

supervised learning



input: \underline{x}

output: \underline{y}

data: $\mathcal{D} = \{(x_i, y_i)\}$

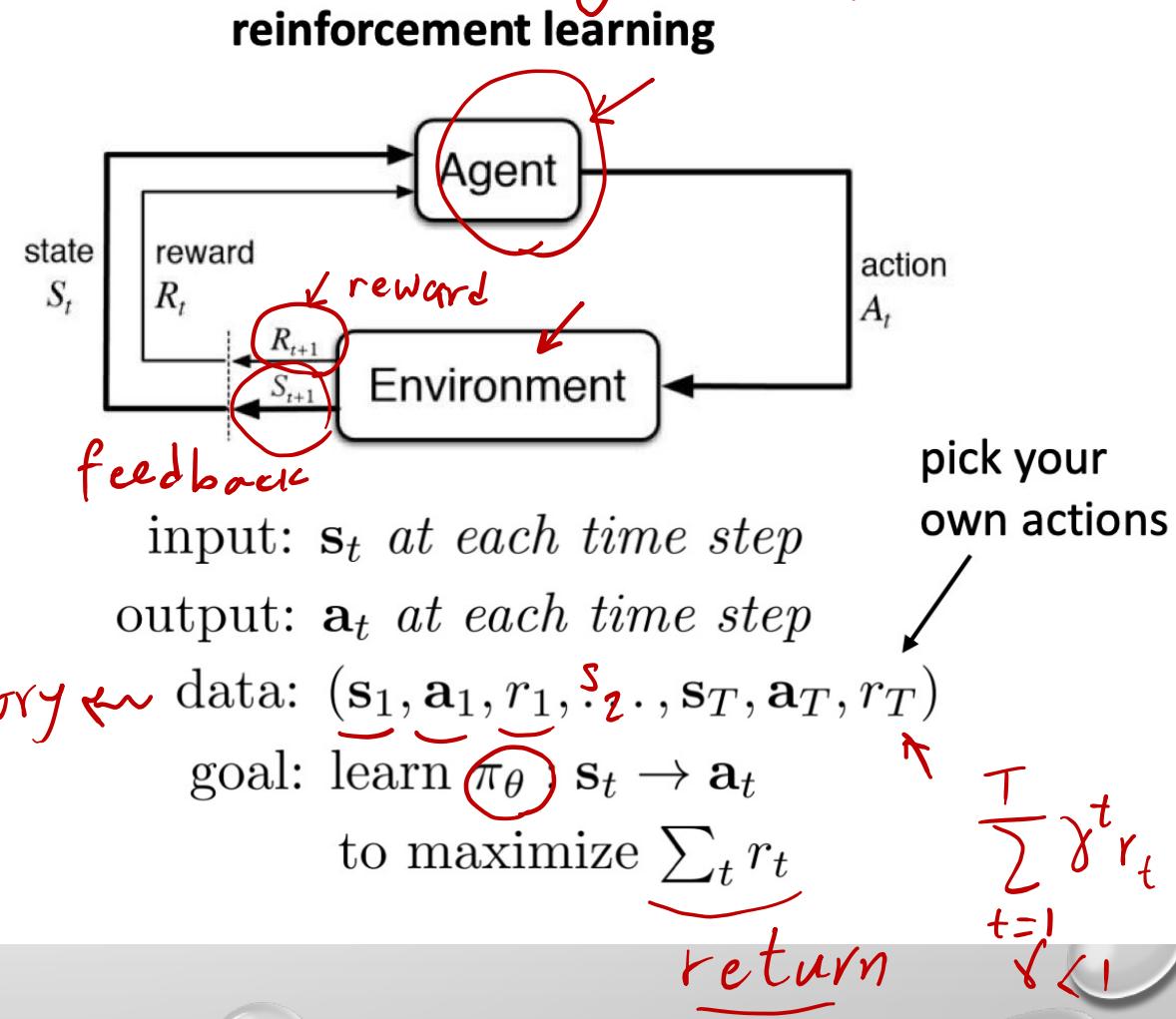
goal: $f_\theta(x_i) \approx y_i$

$$\min_{\theta}$$

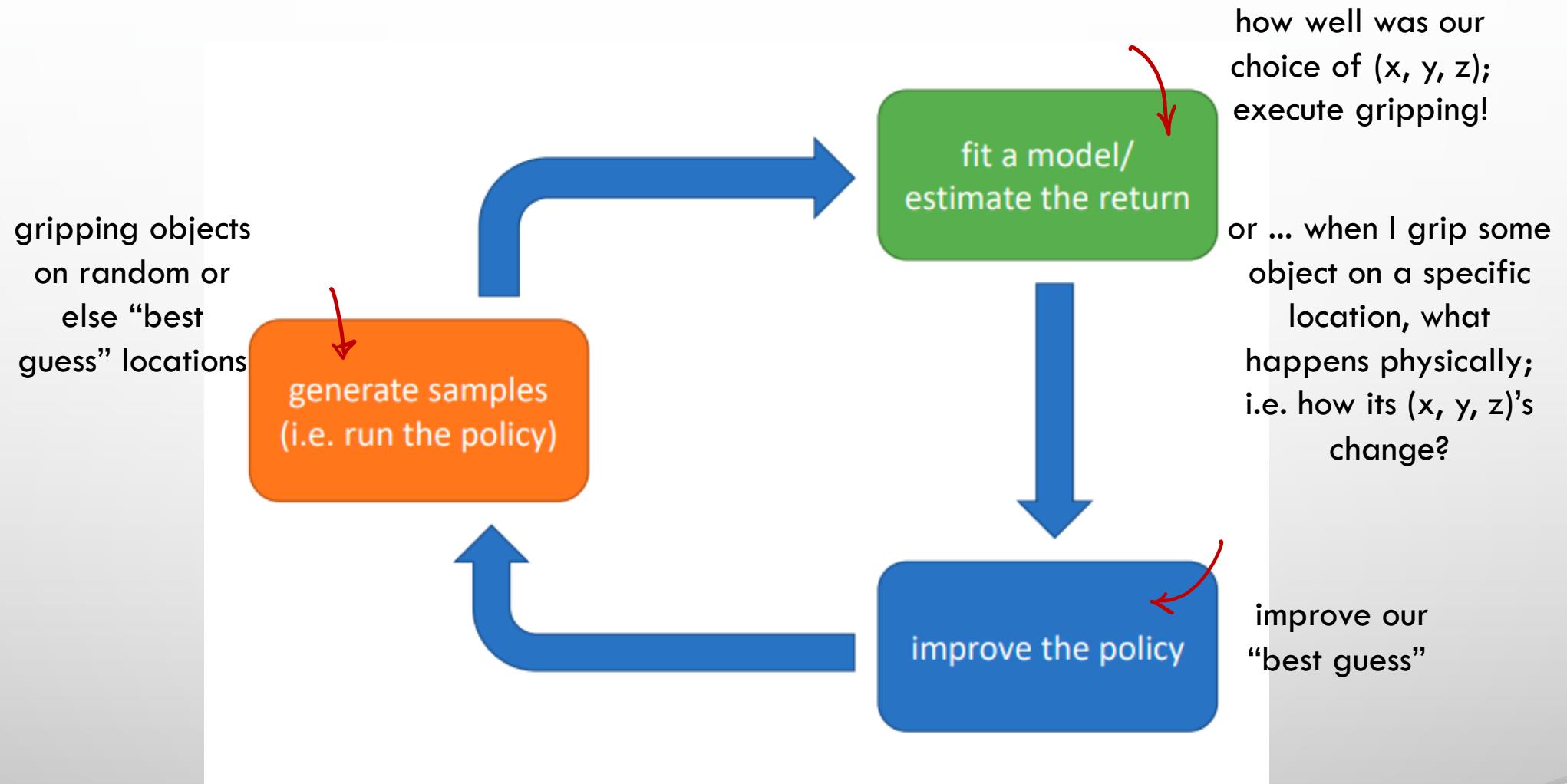
$$\sum_i l(f_\theta(x_i), y_i)$$

someone gives
this to you

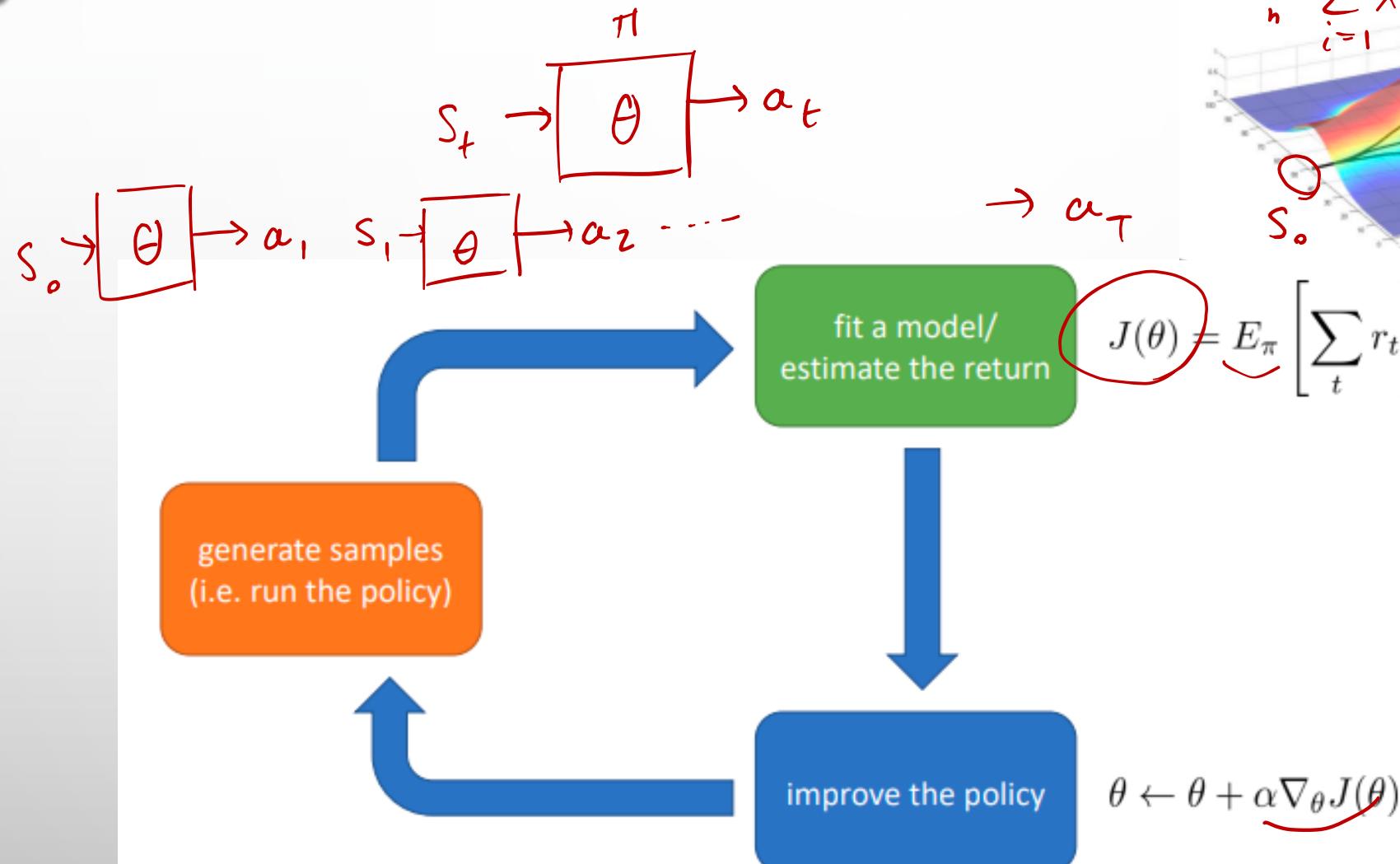
Courtesy: CS 285 course, Berkeley



The Anatomy of Reinforcement Learning



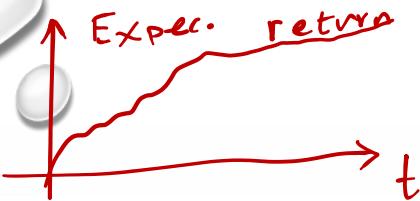
A Simple Example



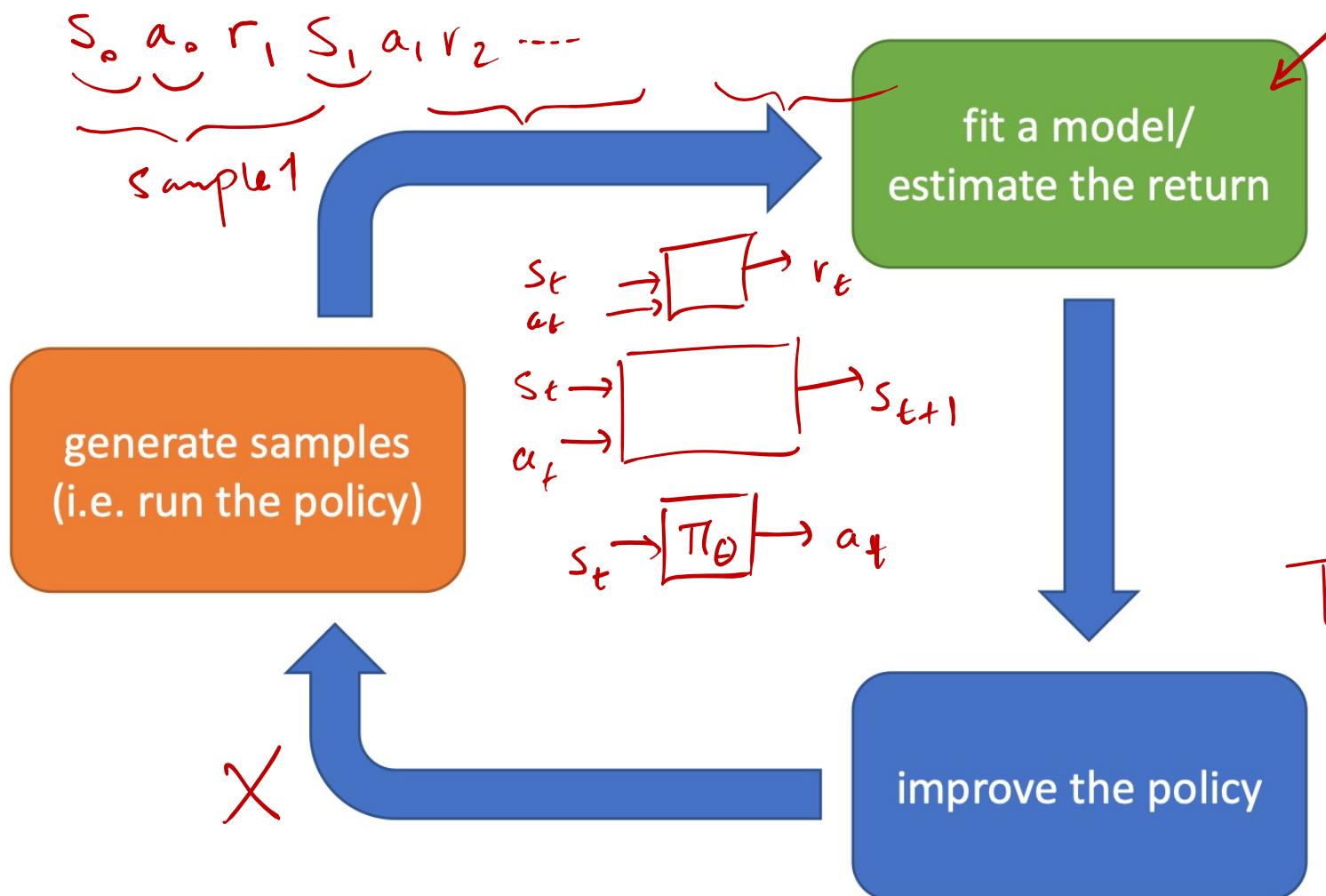
$$\begin{aligned} IE(X) &= \int x \cdot f_X(x) dx \\ &= \frac{1}{n} \sum_{i=1}^n x_i \end{aligned}$$

density

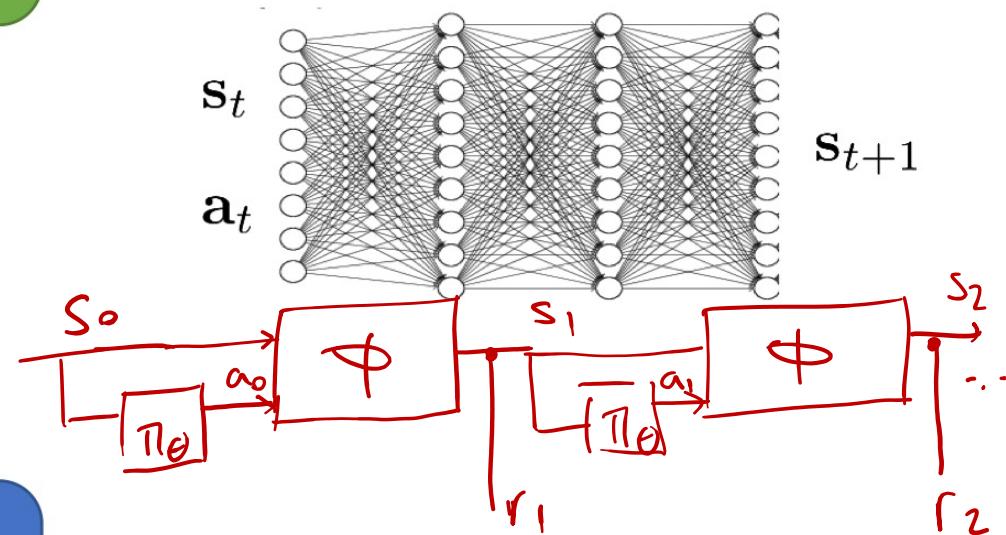
A 3D surface plot illustrating a probability density function. The surface is bell-shaped, representing the distribution of a variable X . Arrows point towards the peak of the distribution, indicating the direction of increasing density.



Another Example



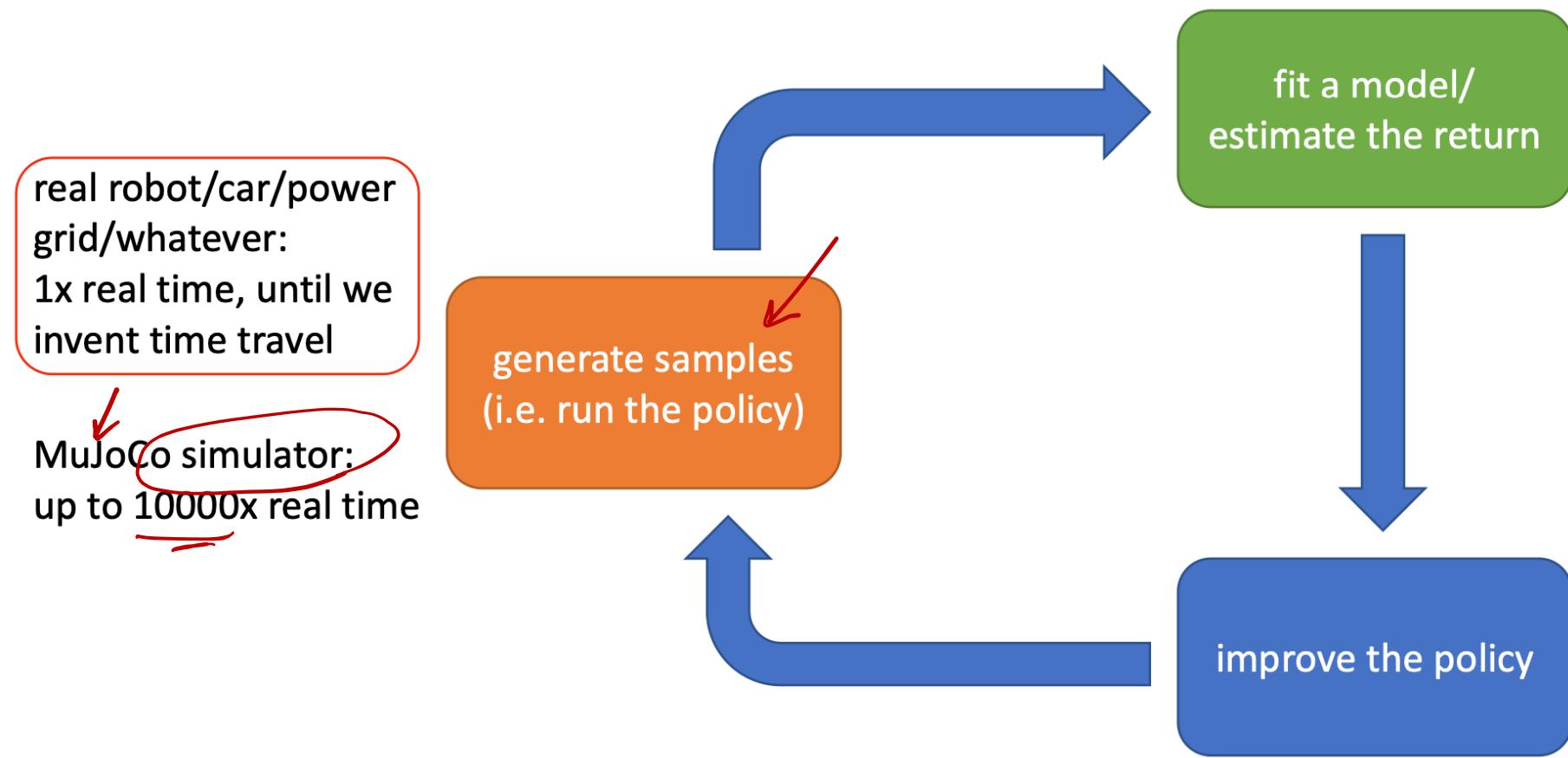
World Model
learn f_ϕ such that $s_{t+1} \approx f_\phi(s_t, a_t)$



backprop through f_ϕ and r to
train $\pi_\theta(s_t) = a_t$

$$\max_{\theta} \sum r_t$$

Which parts are expensive?



$$J(\theta) = E_{\pi} \left[\sum_t r_t \right] \approx \frac{1}{N} \sum_{i=1}^N \sum_t r_t^i$$

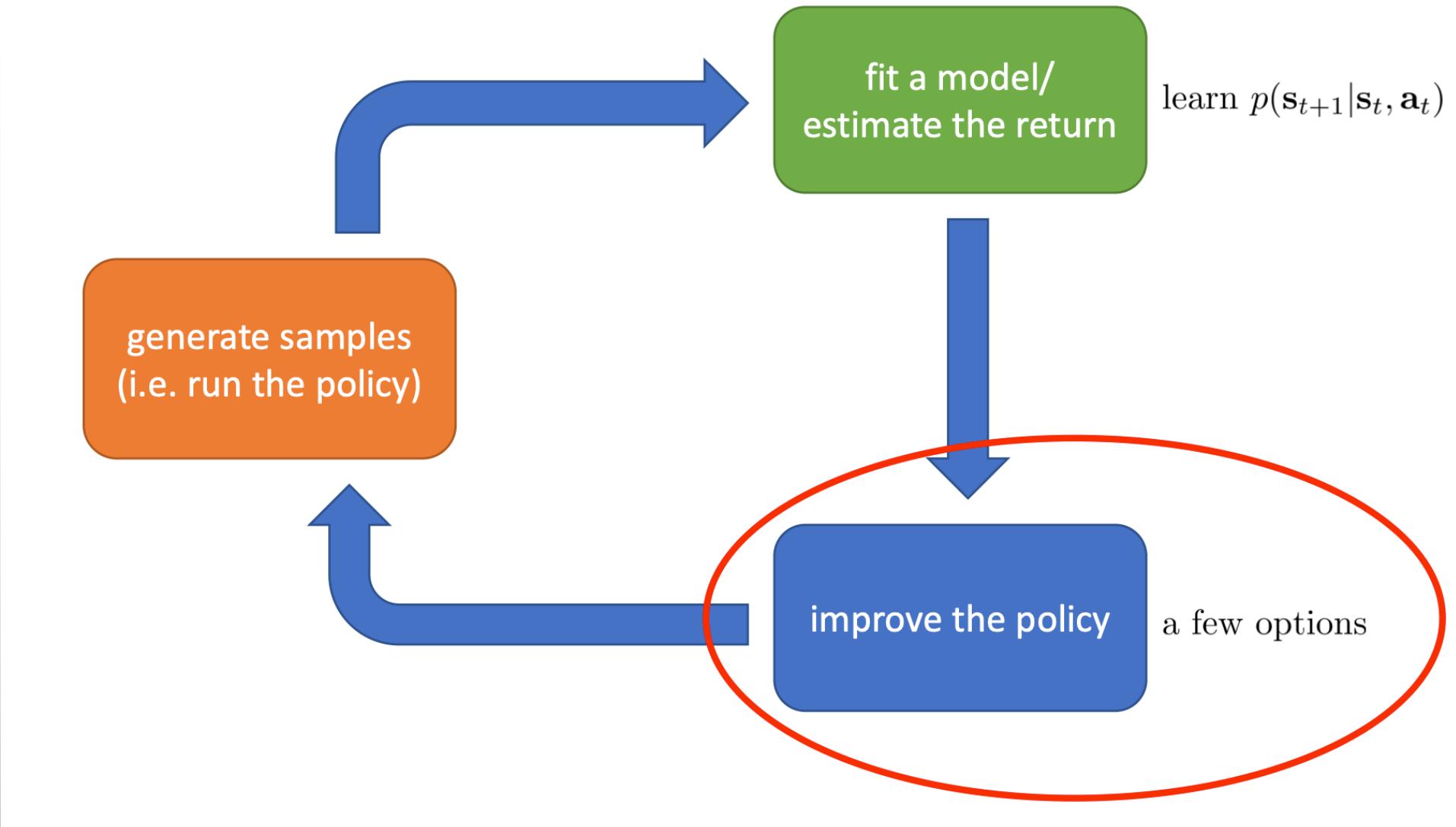
trivial, fast

learn $s_{t+1} \approx f_{\phi}(s_t, a_t)$
expensive

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

backprop through f_{ϕ} and r to
train $\pi_{\theta}(s_t) = a_t$

Model-based RL



Value-based RL

$$\pi_{\text{new}}(s) = \arg \max_a Q^\pi(s, a)$$

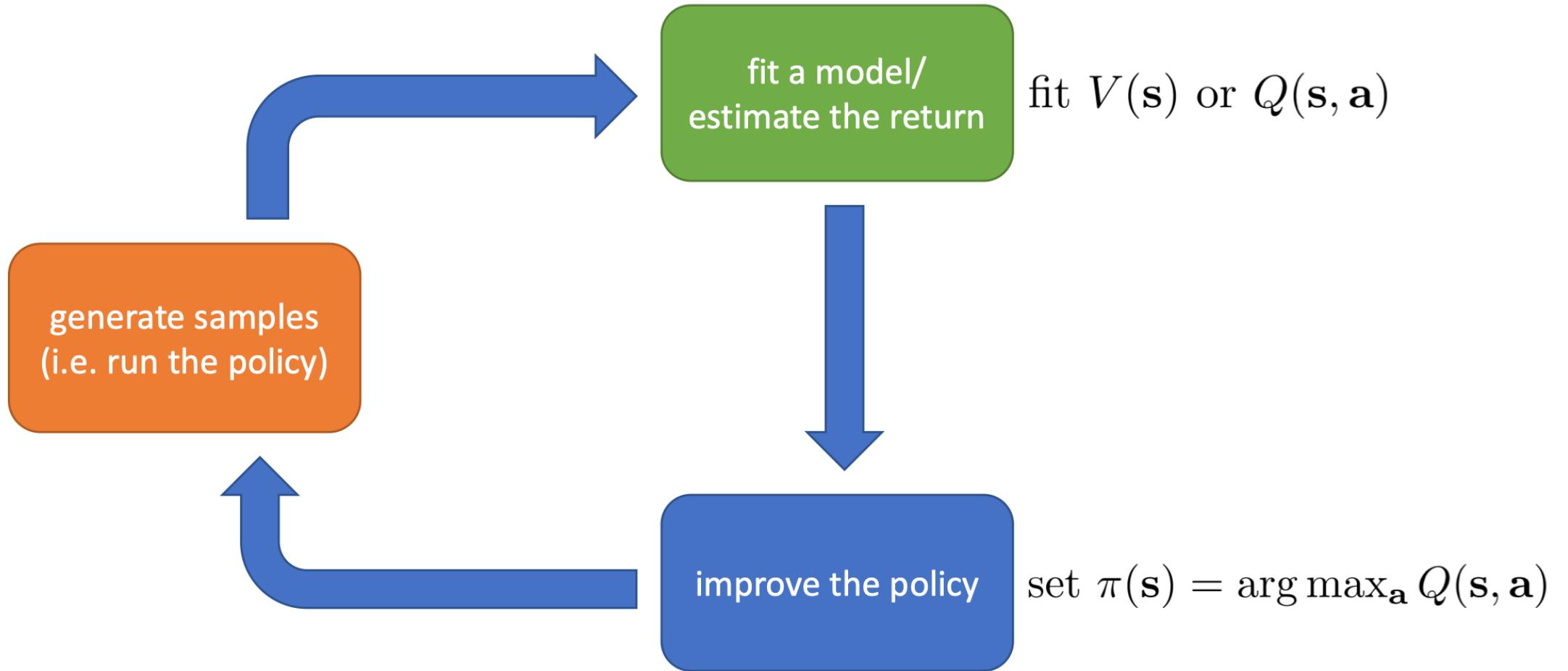
$Q^\pi(s_t, a_t)$ = $\sum_{t'=t}^T E_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t, a_t]$: total reward from taking a_t in s_t

$V^\pi(s_t)$ = $\sum_{t'=t}^T E_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$: total reward from s_t

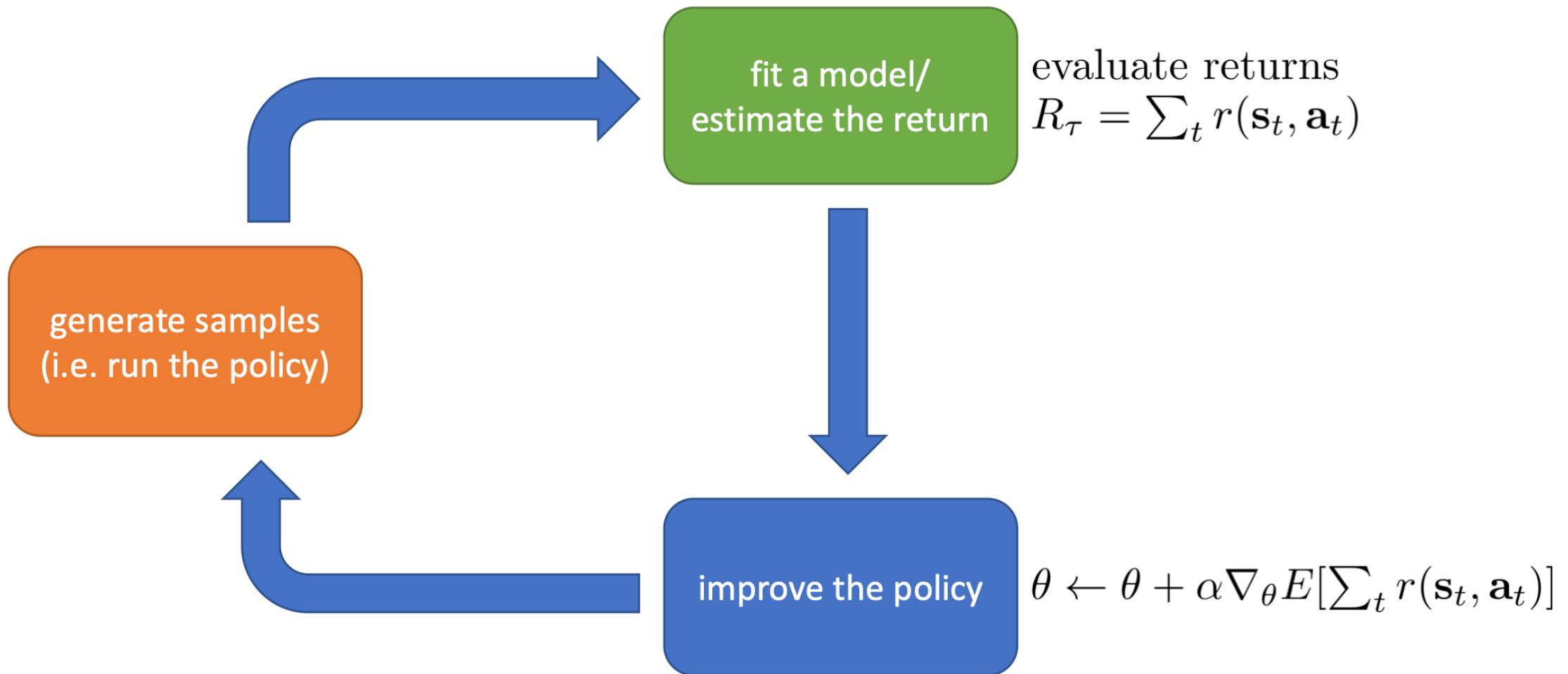
$$V^\pi(s_t) = E_{a_t \sim \pi(a_t | s_t)} [Q^\pi(s_t, a_t)]$$

$E_{s_1 \sim p(s_1)} [V^\pi(s_1)]$ is the RL objective!

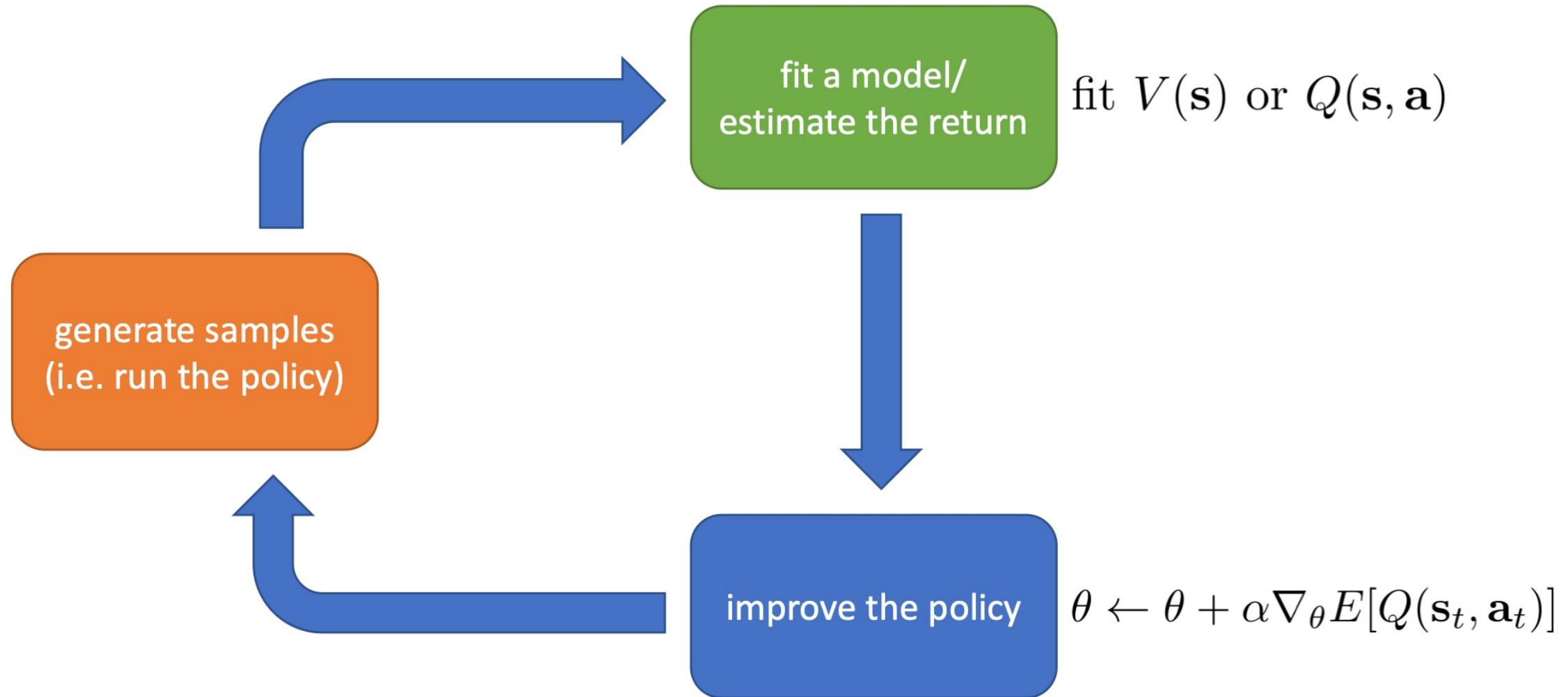
Value-based RL (cont.)



Direct Policy Gradient



Actor-critic: value functions + policy gradients



Where do rewards come from?

- An **expert** gives us the reward
- Learning from **demonstrations**
 - Directly **copying** observed behavior
 - **Inferring rewards** from observed behavior (inverse reinforcement learning)



Motivation (cont.)

	AI Planning	SL	UL	RL	IL
Optimization	X			X	X
Learns from experience		X	X	X	X
Generalization	X	X	X	X	X
Delayed Consequences	X			X	X
Exploration				X	

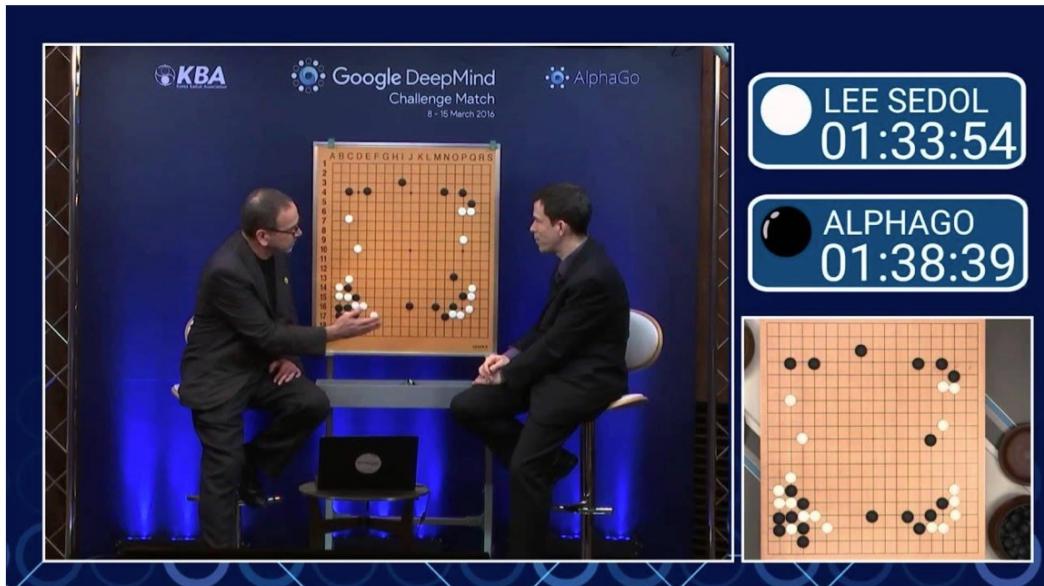
- SL = supervised learning; UL = unsupervised learning; RL = reinforcement learning; IL = imitation learning
- Imitation learning typically assumes input demonstrations of good policies
- IL reduces RL to SL. IL + RL is promising area

Planning vs learning

- Two fundamental problems in sequential decision making
 - Reinforcement learning:
 - The environment is initially **unknown**
 - The agent **interacts** with the environment
 - The agent **improves** its policy
 - Planning:
 - A model of the environment is **known**
 - The agent performs computations with its model (**without any external interaction**)
 - The agent **improves** its policy
 - a.k.a. deliberation, reasoning, introspection, pondering, thought, search

Why should we study deep reinforcement learning?

Impressive because no person had thought of it!



“Move 37” in Lee Sedol AlphaGo match: reinforcement learning “discovers” a move that surprises everyone

Impressive because it looks like something a person might draw!



Data-driven AI vs. RL

Data-Driven AI



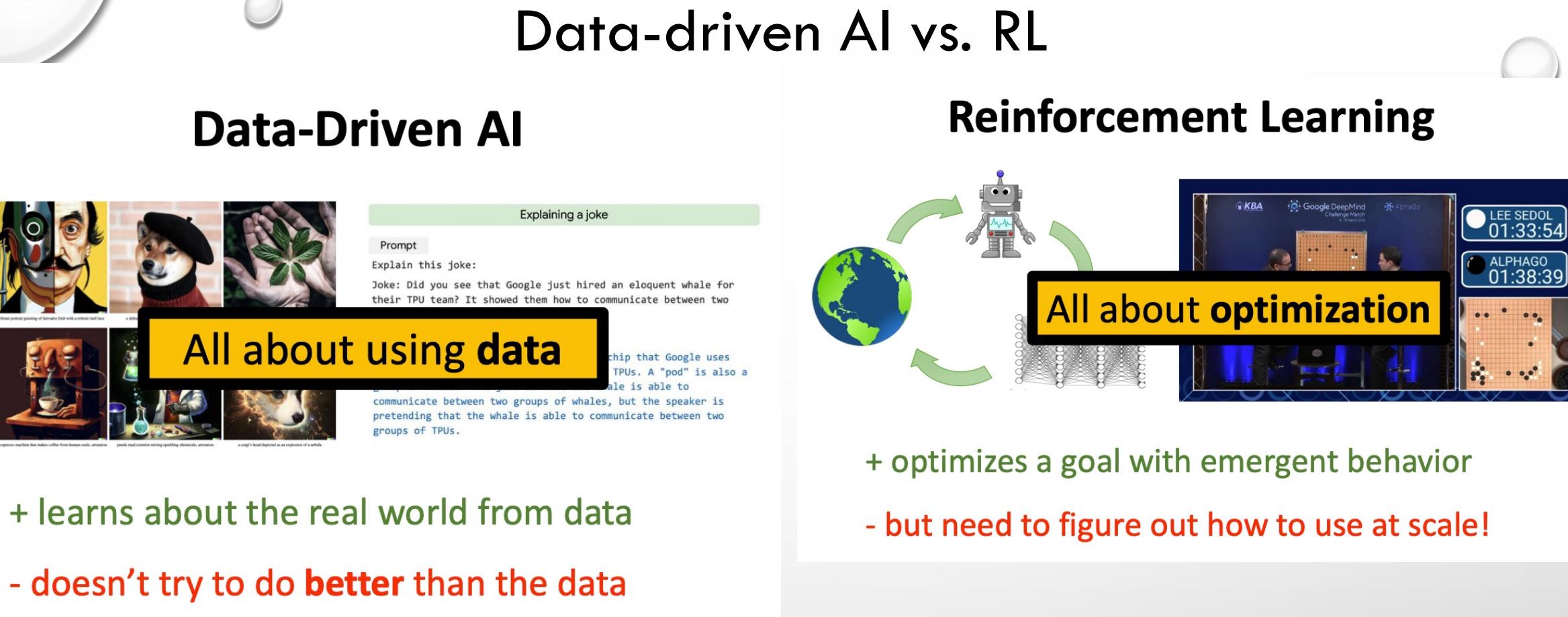
All about using data

Explaining a joke

Prompt

Explain this joke:

Joke: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two



- + learns about the real world from data
- doesn't try to do better than the data

Data without optimization
doesn't allow us to solve new
problems in new ways



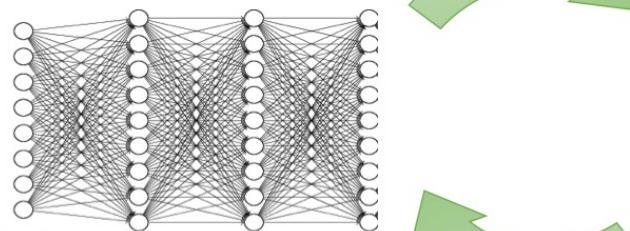
A Bitter Lesson (Richard Sutton)

“We have to learn the bitter lesson that **building in how we think we think** does not work in the long run. The two methods that seem to scale arbitrarily ... are **learning** and **search**

<http://www.incompleteideas.net/Incldeas/BitterLesson.html>

Learning

use **data** to extract **patterns**



allows us to **understand** the world

Search

use **computation** to extract **inferences**

optimization

some optimization process that uses (typically iterative) computation to make rational decisions

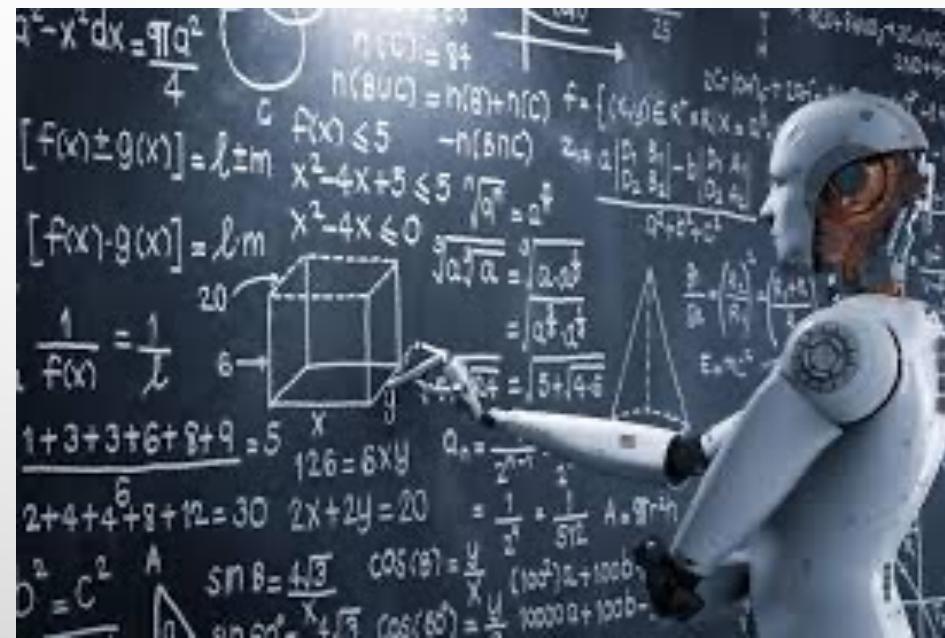
leverages that **understanding** for **emergence**

Data without **optimization** doesn't allow us to solve new problems in new ways

Optimization without **data** is hard to apply to the real world outside of simulators

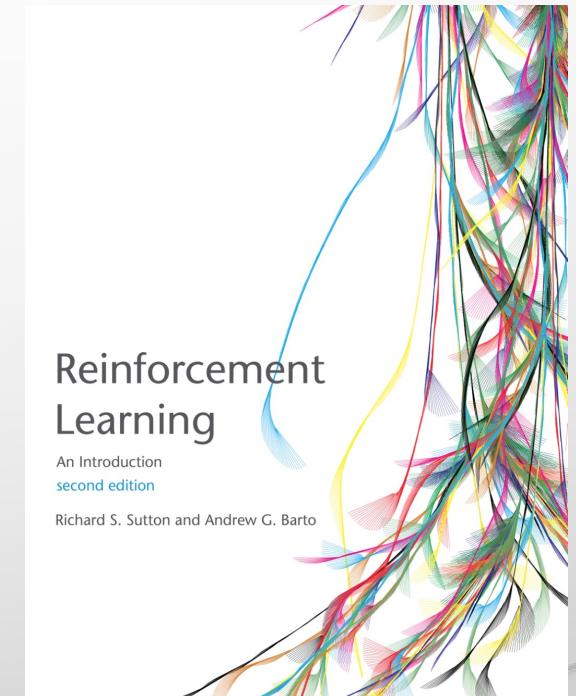
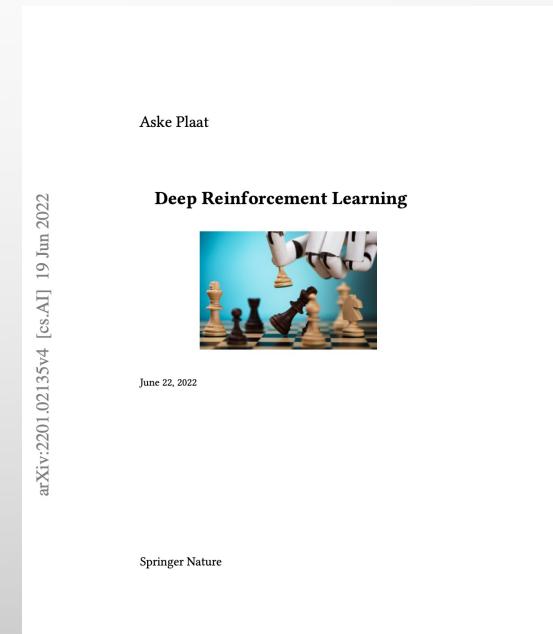
Superintelligence

- The models are trained based on **human annotations** and **preferences**.
 - Can they get **smarter** than humans?



References

- Reinforcement Learning: An Introduction by R. Sutton and A. Barto,
2nd Edition, 2020.
- Deep Reinforcement Learning by A. Plaat, 2022.
- Original papers of some methods.



Teaching Assistants

- Ali Najar (Head TA)
- Pouya Toroghi (Head TA)

Prereqs.

- Stochastic Processes (Prob. And Stats, Markov Processes, Estimation Theory, Information Theory)
- Optimization (Lagrange Multipliers)
- Deep Learning (Concepts and Pytorch)

Motivation (cont.) ChatGPT; Why RL?!

Step 1

Collect demonstration data
and train a supervised policy.

A prompt is sampled from our prompt dataset.

Explain reinforcement learning to a 6 year old.

A labeler demonstrates the desired output behavior.

We give treats and punishments to teach...

This data is used to fine-tune GPT-3.5 with supervised learning.

SFT



Step 2

Collect comparison data and train a reward model.

A prompt and several model outputs are sampled.

Explain reinforcement learning to a 6 year old.

A
In reinforcement learning, the agent is...
B
Explain rewards...
C
In machine learning...
D
We give treats and punishments to teach...

A labeler ranks the outputs from best to worst.

D > C > A > B

This data is used to train our reward model.

RM



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.

Write a story about otters.

The PPO model is initialized from the supervised policy.

PPO


The policy generates an output.

Once upon a time...

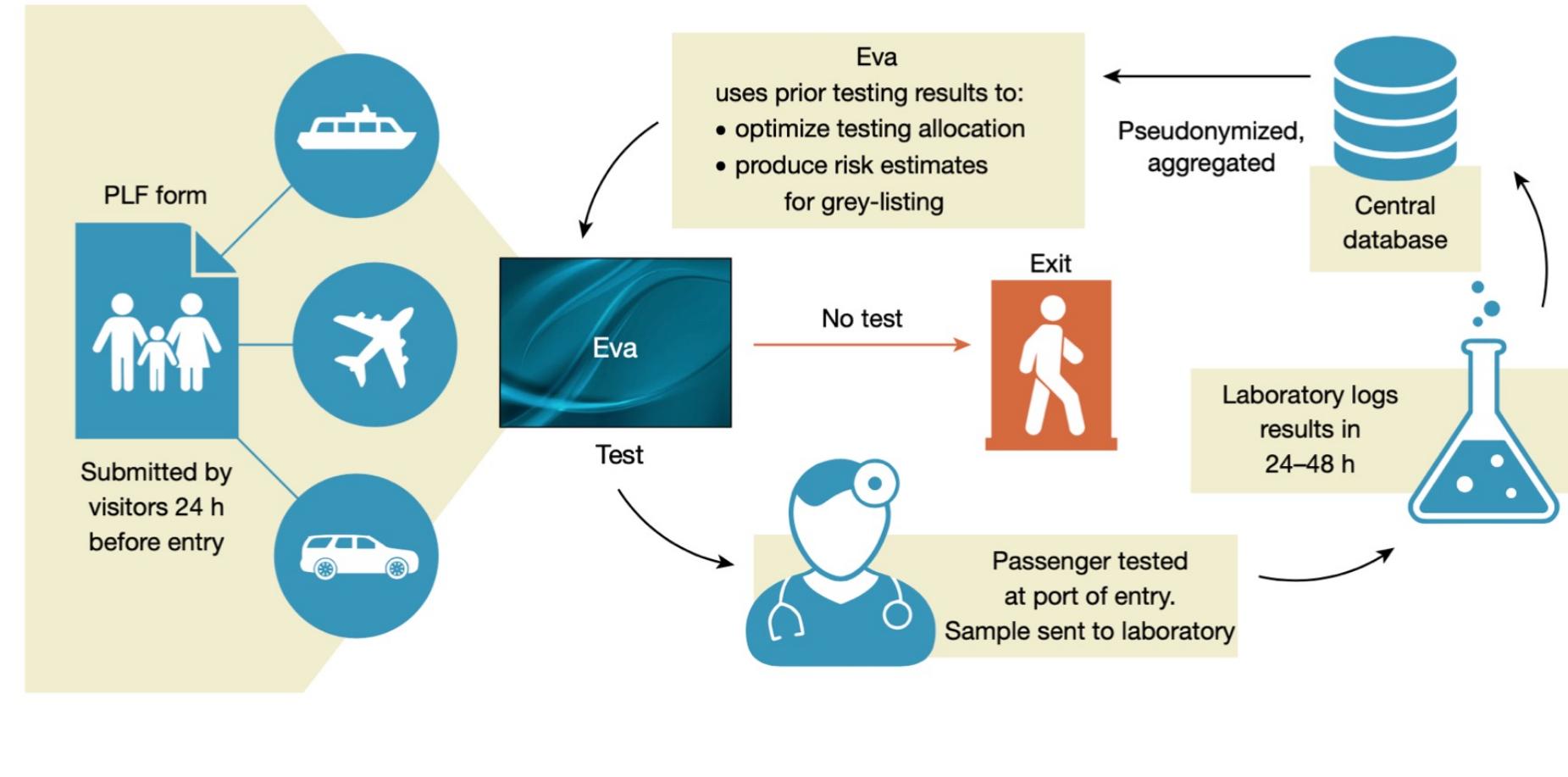
The reward model calculates a reward for the output.

RM


The reward is used to update the policy using PPO.

r_k

Motivation (cont.)

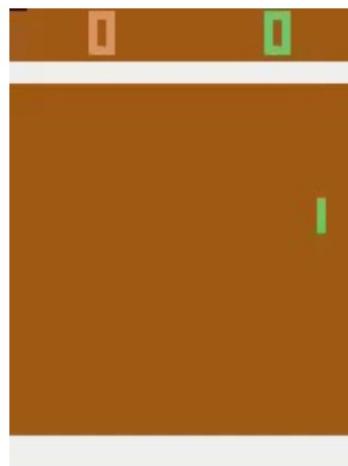


- <https://www.nature.com/articles/s41586-021-04014-z>

History

2013

Atari (DQN)
[Deepmind]



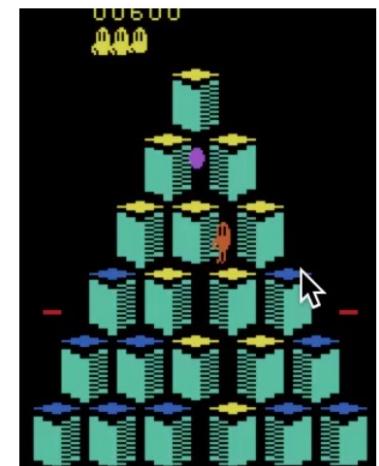
Pong



Enduro



Beamrider



Q*bert

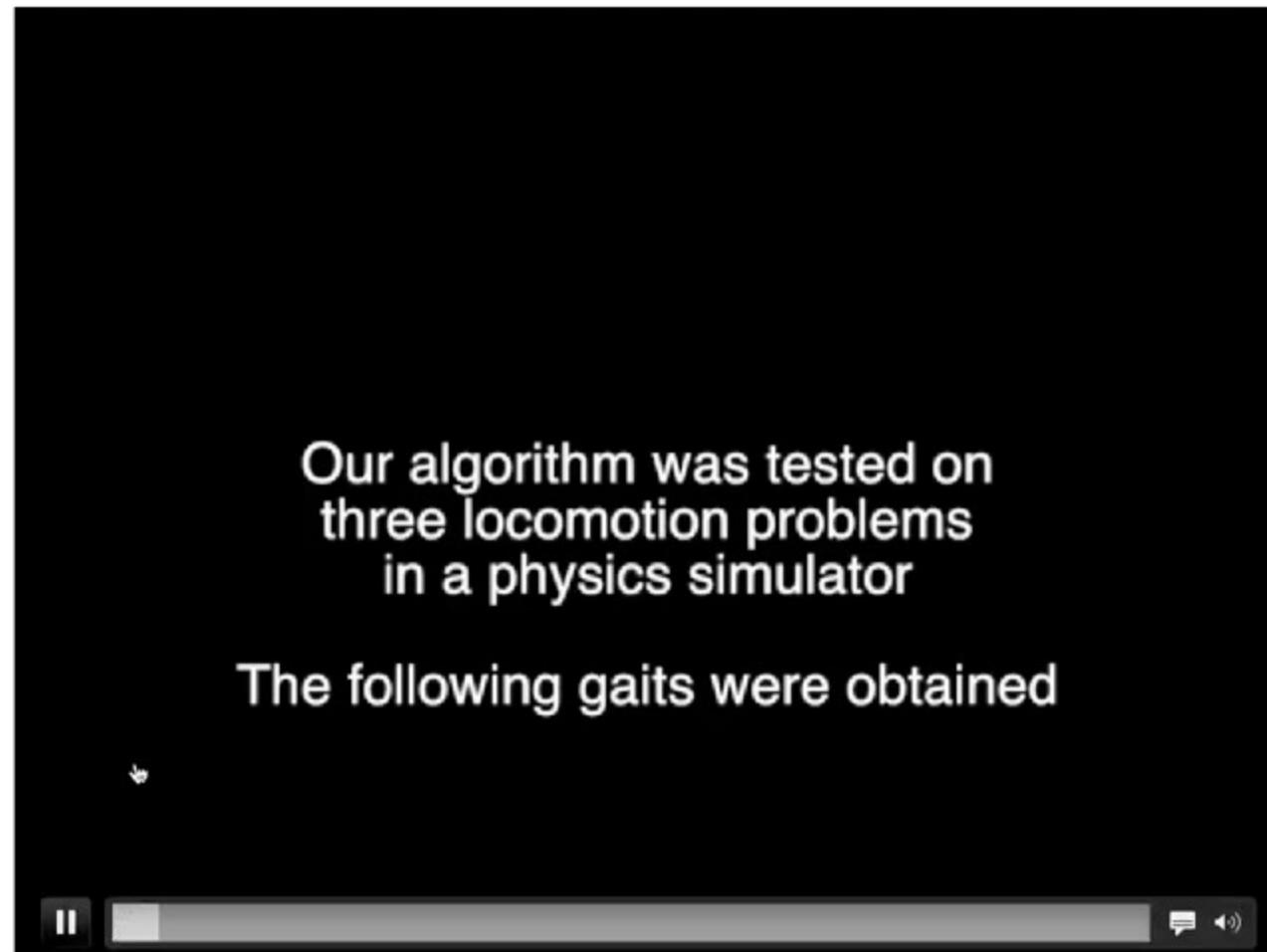
A Few Deep RL Highlights

2013

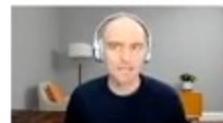
Atari (DQN)
[Deepmind]

2014

2D locomotion (TRPO)
[Berkeley]



Play 0:06 – 0:25



History

2013	Atari (DQN) [Deepmind]
2014	2D locomotion (TRPO) [Berkeley]
2015	AlphaGo [Deepmind]



Tian et al, 2016; Maddison et al, 2014; Clark et al, 2015

A Few Deep RL Highlights

2013

Atari (DQN)
[Deepmind]

2014

2D locomotion (TRPO)
[Berkeley]

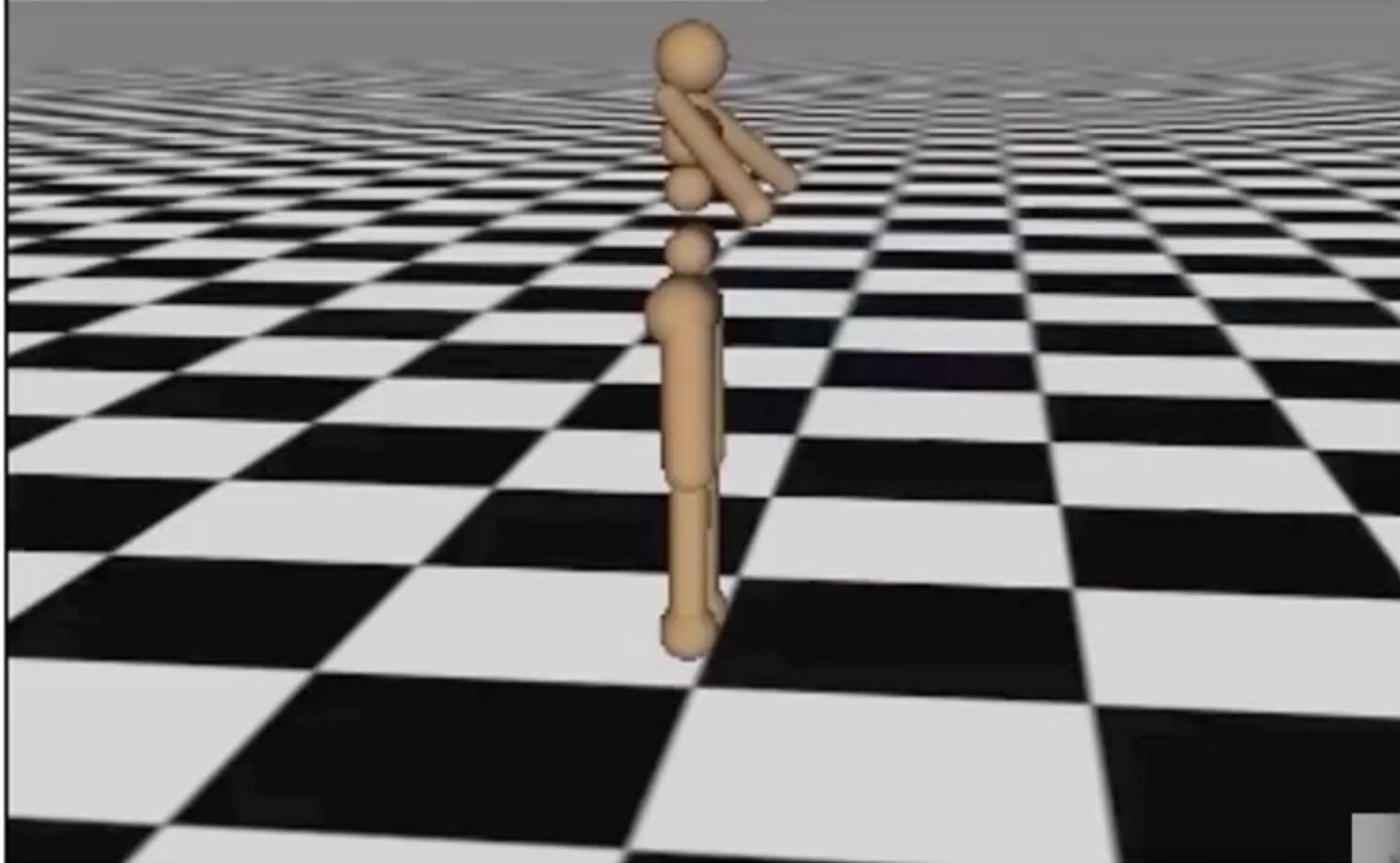
2015

AlphaGo
[Deepmind]

2016

3D locomotion (TRPO+GAE)
[Berkeley]

Iteration 0



[Schulman, Moritz, Levine, Jordan, Abbeel, ICLR 2016]



A Few Deep RL Highlights

2013	Atari (DQN) [Deepmind]
2014	2D locomotion (TRPO) [Berkeley]
2015	AlphaGo [Deepmind]
2016	3D locomotion (TRPO+GAE) [Berkeley]
2016	Real Robot Manipulation (GPS) [Berkeley]

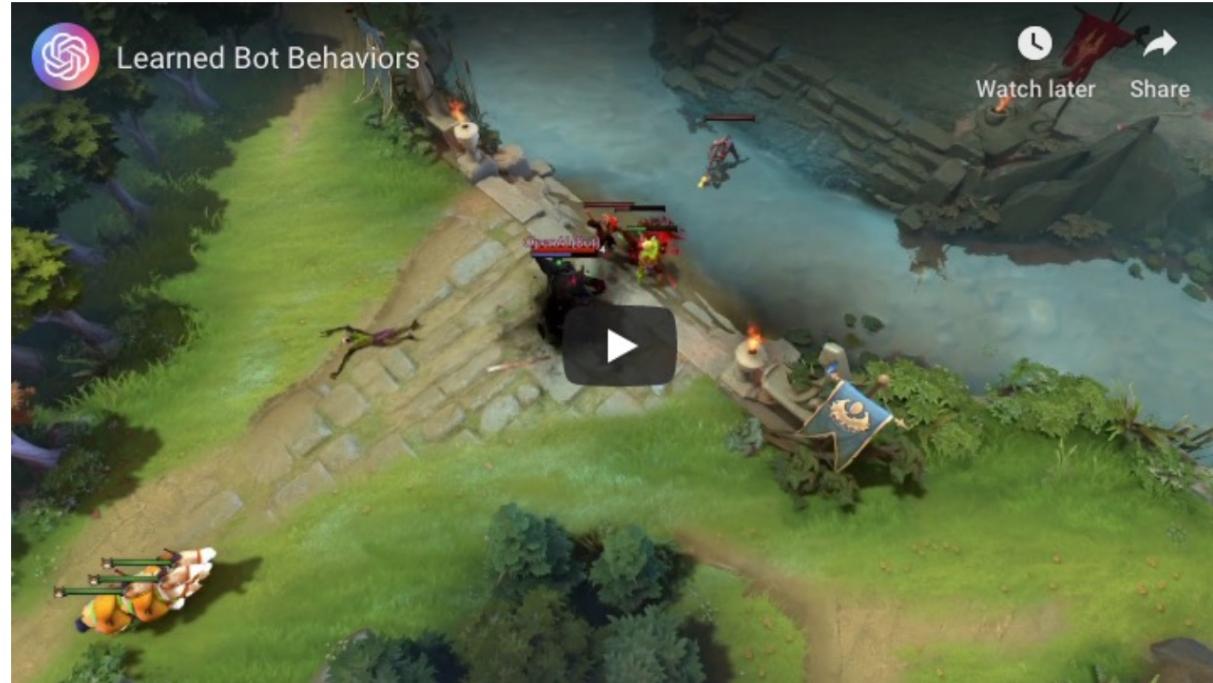


[Levine*, Finn*, Darrell, Abbeel, JMLR 2016]



History

2013	Atari (DQN) [Deepmind]
2014	2D locomotion (TRPO) [Berkeley]
2015	AlphaGo [Deepmind]
2016	3D locomotion (TRPO+GAE) [Berkeley]
2016	Real Robot Manipulation (GPS) [Berkeley, Google]
2017	Dota2 (PPO) [OpenAI]



OpenAI Dota Bot beat best humans 1:1 (Aug 2018)

A Few Deep RL Highlights

2013

Atari (DQN)
[Deepmind]

2014

2D locomotion (TRPO)
[Berkeley]

2015

AlphaGo
[Deepmind]

2016

3D locomotion (TRPO+GAE)
[Berkeley]

2016

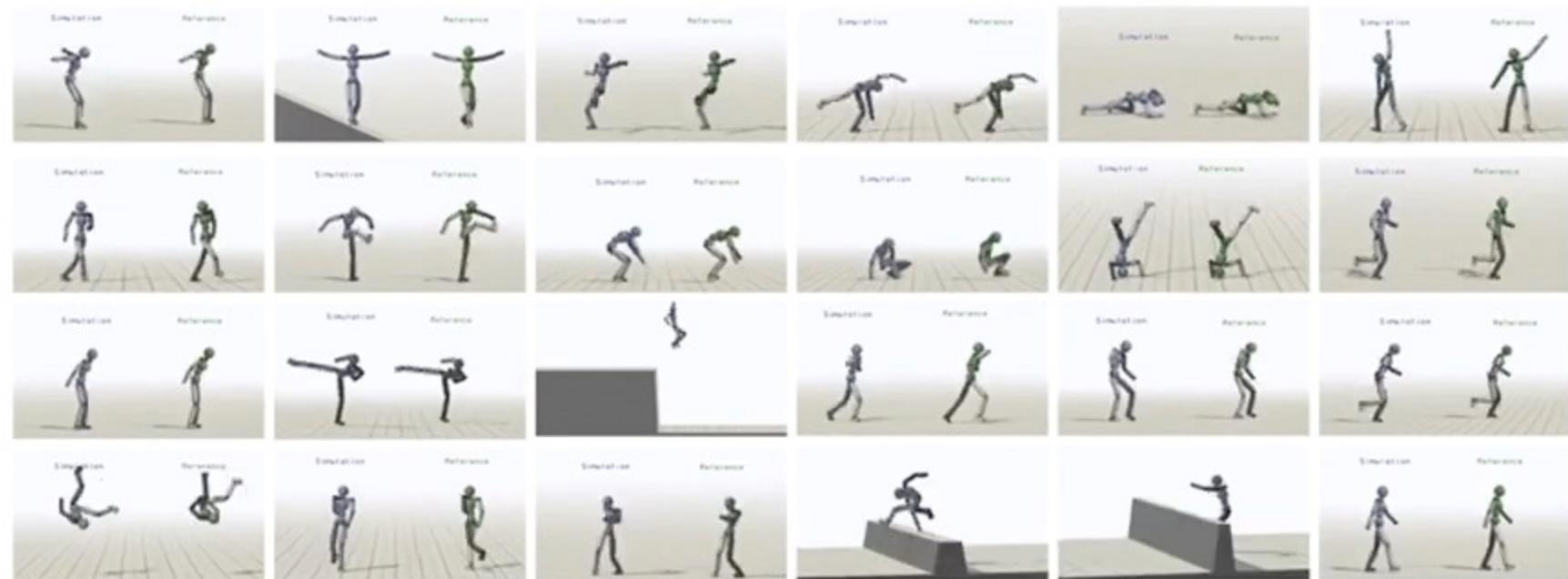
Real Robot Manipulation
(GPS) [Berkeley, Google]

2017

Dota2
(PPO) [OpenAI]

2018

DeepMimic
[Berkeley]



[Peng, Abbeel, Levine, van de Panne, 2018]



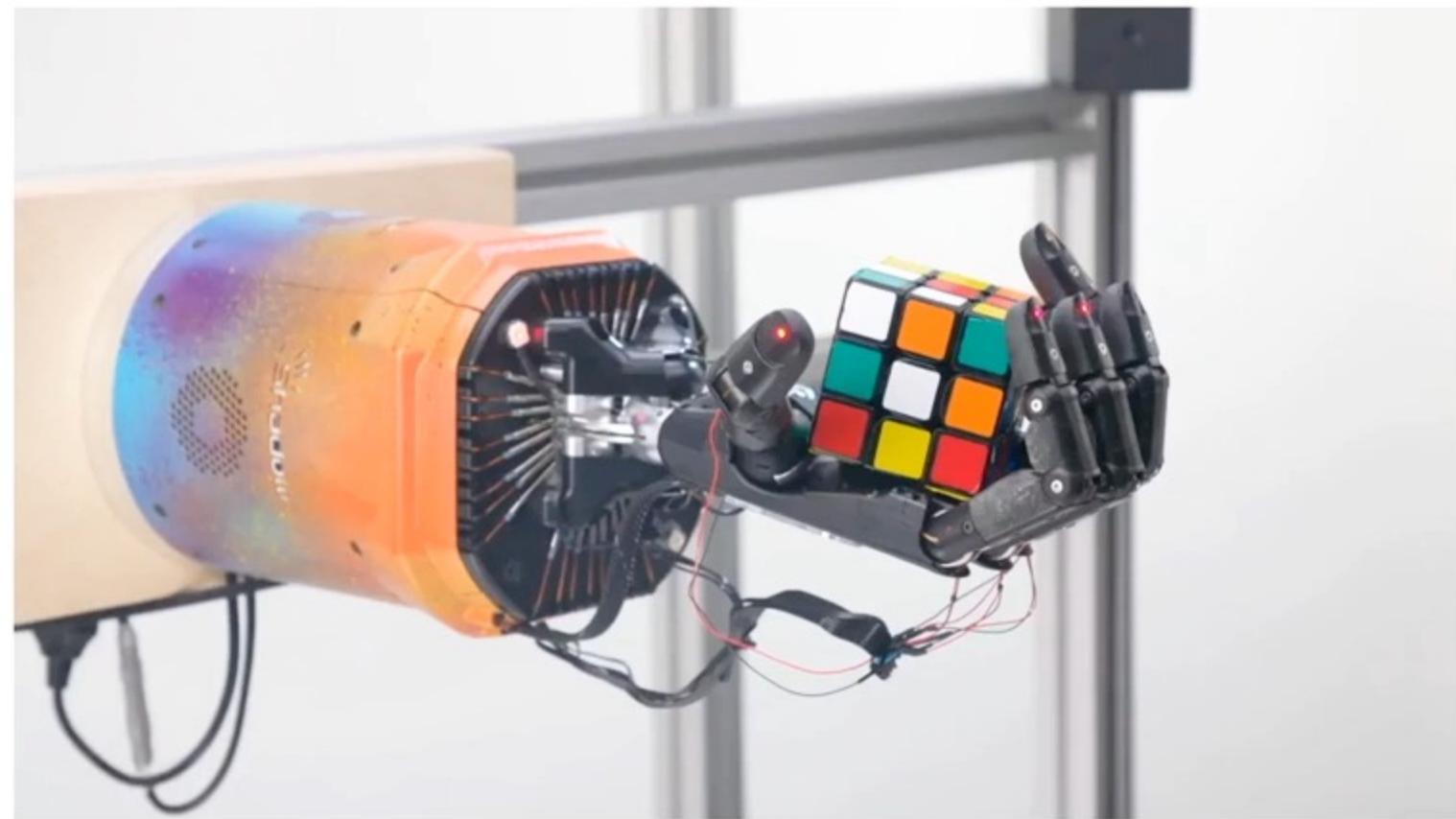
History

2013	Atari (DQN) [Deepmind]
2014	2D locomotion (TRPO) [Berkeley]
2015	AlphaGo [Deepmind]
2016	3D locomotion (TRPO+GAE) [Berkeley]
2016	Real Robot Manipulation (GPS) [Berkeley, Google]
2017	Dota2 (PPO) [OpenAI]
2018	DeepMimic [Berkeley]
2019	AlphaStar [Deepmind]



A Few Deep RL Highlights

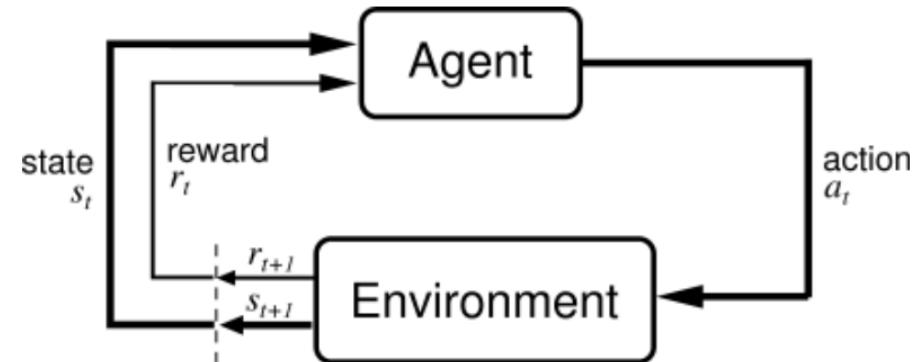
2013	Atari (DQN) [Deepmind]
2014	2D locomotion (TRPO) [Berkeley]
2015	AlphaGo [Deepmind]
2016	3D locomotion (TRPO+GAE) [Berkeley]
2016	Real Robot Manipulation (GPS) [Berkeley, Google]
2017	Dota2 (PPO) [OpenAI]
2018	DeepMimic [Berkeley]
2019	AlphaStar [Deepmind]
2019	Rubik's Cube (PPO+DR) [OpenAI]



Let's Begin: Markov Decision Processes (MDPs)

An MDP is defined by:

- Set of states S
- Set of actions A
- Transition function $P(s' | s, a)$
- Reward function $R(s, a, s')$
- Start state s_0
- Discount factor γ
- Horizon H



The Goal

- The policy is $\pi_\theta: S \rightarrow A$ for infinite horizon or
 $\pi_\theta: S \times \{0, \dots, H\} \rightarrow A$ for finite horizon MDP.

MDP (S, A, T, R, γ, H) ,

goal: $\max_{\pi} \mathbb{E} \left[\sum_{t=0}^H \gamma^t R(S_t, A_t, S_{t+1}) | \pi \right]$

Sometimes the policy could be stochastic: $\pi : S \times A \rightarrow [0,1]$, which is
 $\pi(a|s) = \Pr(A_t = a | S_t = s)$.

Example: Grid World

An MDP is defined by:

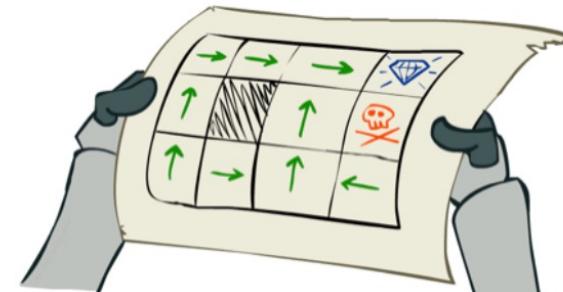
- Set of states S
- Set of actions A
- Transition function $P(s' | s, a)$
- Reward function $R(s, a, s')$
- Start state s_0
- Discount factor γ
- Horizon H



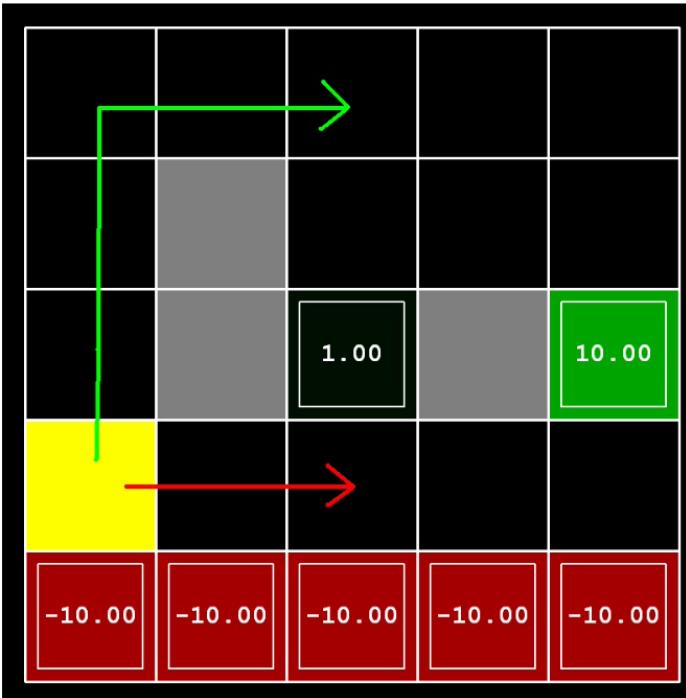
Goal:

$$\max_{\pi} \mathbb{E} \left[\sum_{t=0}^H \gamma^t R(S_t, A_t, S_{t+1}) \middle| \pi \right]$$

π^* :



Exercise



- (a) Prefer the close exit (+1), risking the cliff (-10) (1) $\gamma = 0.1$, noise = 0.5
- (b) Prefer the close exit (+1), but avoiding the cliff (-10) (2) $\gamma = 0.99$, noise = 0
- (c) Prefer the distant exit (+10), risking the cliff (-10) (3) $\gamma = 0.99$, noise = 0.5
- (d) Prefer the distant exit (+10), avoiding the cliff (-10) (4) $\gamma = 0.1$, noise = 0