

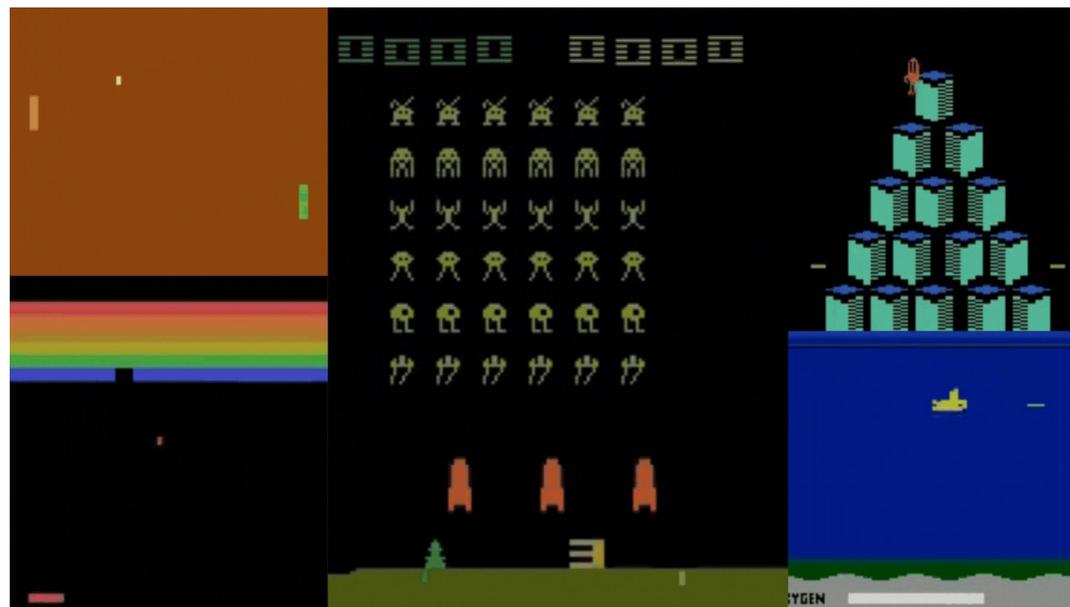


# Deep Reinforcement Learning

## Spring 2026

Abhishek Gupta

TA: Mateo Guaman Castro



Welcome to Deep Reinforcement Learning – Sp 26!

# Lecture Outline

---

Course logistics and scope



What is RL, a formal definition



Why should we care?



Going beyond RL

# Ok so what is CSE 579

---

## Study of applied reinforcement learning algorithms

- Sprinklings of optimal control and trajectory optimization
- Sprinklings of reinforcement learning theory

Not a pure optimal control class, would recommend Karen Leung's classes 😊

Not a pure theory class, would recommend Kevin Jamieson's class – CSE 542 😊

# Course Logistics

---

- Where: Gates G01
- When: Mon/Wed 1-2:20
- Who:
  - Abhishek Gupta (Instructor)
  - Mateo Guaman Castro (TA)
- Office hours:
  - Abhishek: Gates 215, Mon 4-5pm
  - Mateo: ???

# Course Logistics

- Grading: Seminar style
  - 35 % final project
  - 15% seeded paper discussion
  - 45 % for HWs – 15% for each of 3 HWs
  - 5% participation
- Communications through EdStem/e-mail
- Mix of lectures and seeded idea discussion
- Final projects will be presented in a poster session.
  - Intermediate project proposals and milestone check ins.
- Please participate, otherwise it will be boring for all of us!

# Course Logistics - Project

- Final project (35% of grade):
  - Project proposal (1 page) [Due ??]
  - Milestone report (3-4 pages) [Due ?? (subject to change)]
  - Final report (6-8 pages) [Due ?? (subject to change)]
- Project can be investigating any question related to reinforcement learning, imitation learning or sequential decision making
  - New algorithm
  - Performant/stable implementation
  - Empirical investigation
  - New application or domain
  - ...
- Can be done in groups of 1-2 students.

# Course Logistics – Seeded Paper Ideas

---

- We will try out a new format for discussions
- Key idea: we will seed ideas with a "seed paper". Your job is to build from the seed paper and suggest a new paper-level idea, and defend it to the class.
  - **Critiques of prior work:** Tell us why the prior work is not sufficient
  - **Motivation:** Tell us why we should care about your idea
  - **Technical Idea:** Tell us your idea
  - **Experiments:** Tell us how you would validate your idea and what experiments you'd run
  - **Related Work:** Tell us how your idea will position itself in the literature
  - **Downstream Impacts:** Tell us why your idea matters.
- Presentation will be done by groups of 5-6
- Everyone not presenting posts constructive commentaries on EdStem!

# Course Logistics – Homework

---

- 3 HW assignments + 1 extra credit, each Python programming of different algorithms
- HW 1 – Imitation Learning
  - Implement and test out imitation learning algorithms in simulation
- HW 2 – Model-free RL
  - Implement and test out policy gradient and actor critic methods
- HW 3 – Model-based RL
  - Implement and test out model-based RL algorithms
- HW 4 – Diffusion Model Finetuning with RL
  - Implement and test out algorithms for finetuning diffusion models with RL
- Submit through canvas with a small written report.

# Who am I?

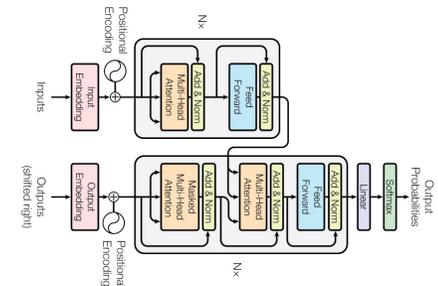
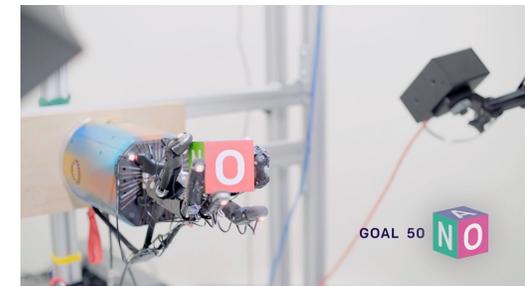
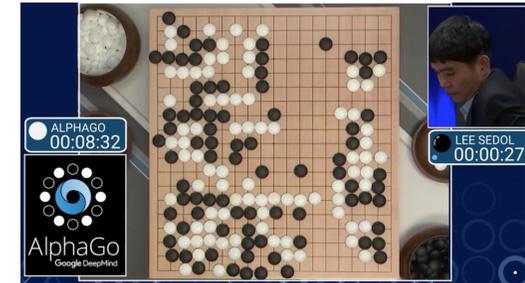
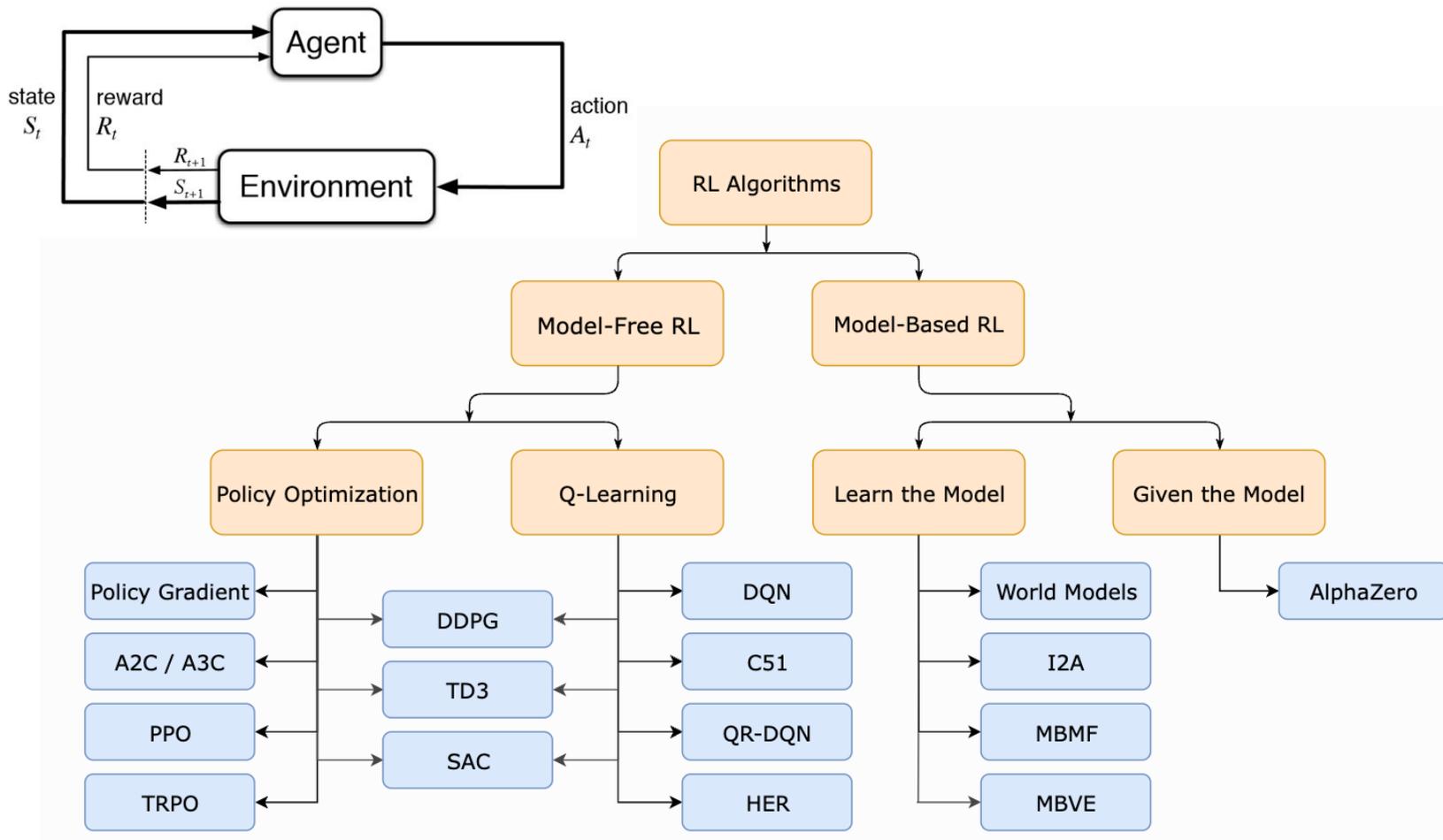


- Assistant professor in CSE
- Grew up in Oregon/India, last 10 years in Berkeley
- Undergrad Berkeley, Ph.D. Berkeley, Postdoc MIT.
- Interests: RL/robotics/optimization and control/robustness and generalization
- Outside of work: Running/soccer/sketching/dog enthusiast

Who are you?

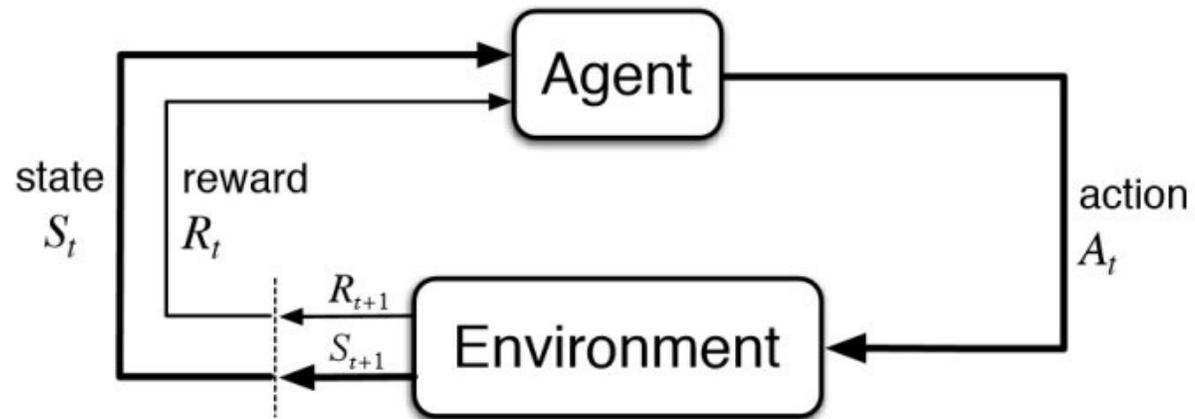
# What is this course about?

- The design and **practice** of reinforcement learning algorithms

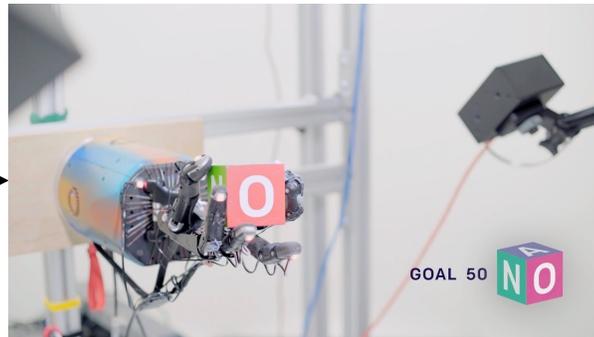
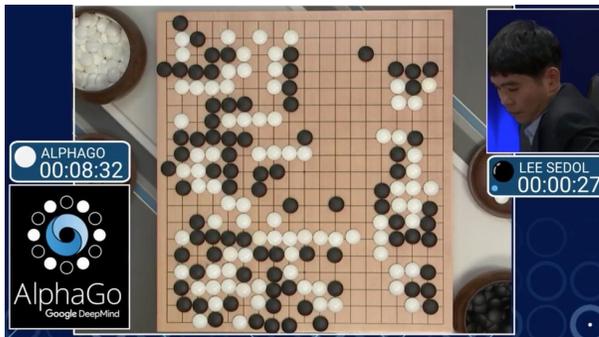


# What is this course about?

- Building RL algorithms that are practical for real applications



- Sample efficient
- Operates from high-dimensional observations
- Continually improving



RL algorithms were not conceived to operate under practical assumptions, needs some extra work

# What is this course about?

## ■ Understanding what powers current RL applications

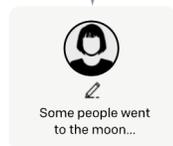
Step 1

**Collect demonstration data, and train a supervised policy.**

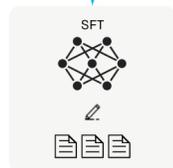
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.



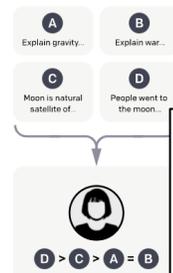
Step 2

**Collect comparison data, and train a reward model.**

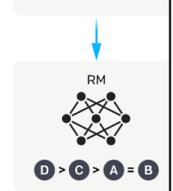
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



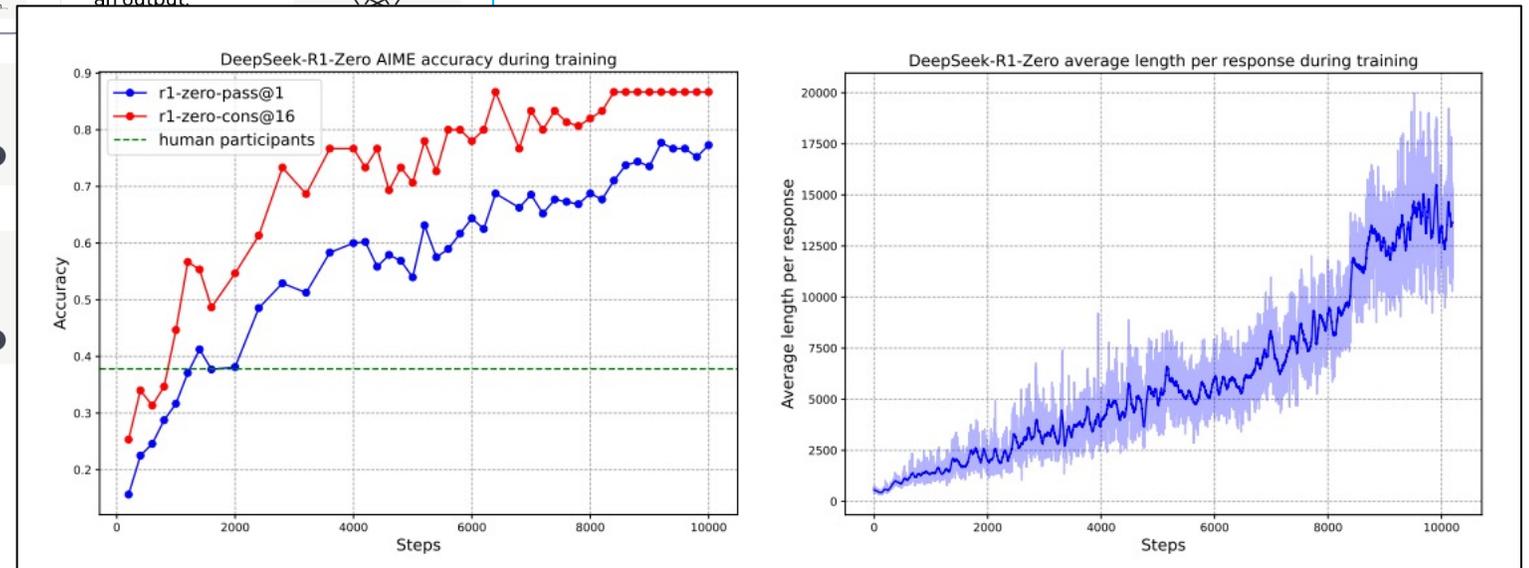
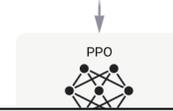
Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.



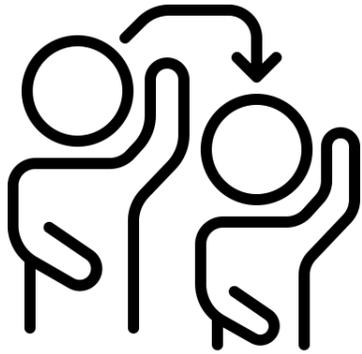
The policy generates an output.



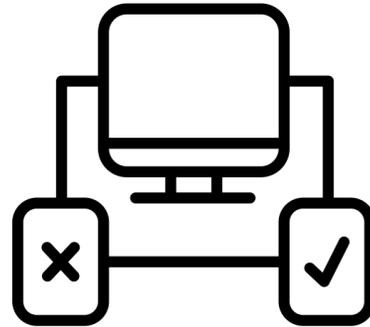
# What is this course about?

- Practice implementing and **tuning** sequential decision making algorithms

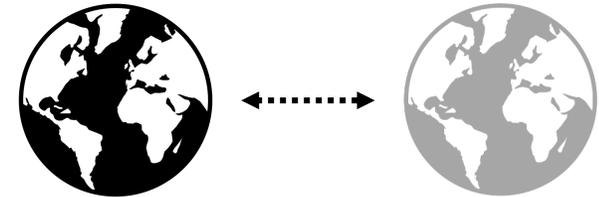
Imitation learning



Model-Free RL



Model-Based RL

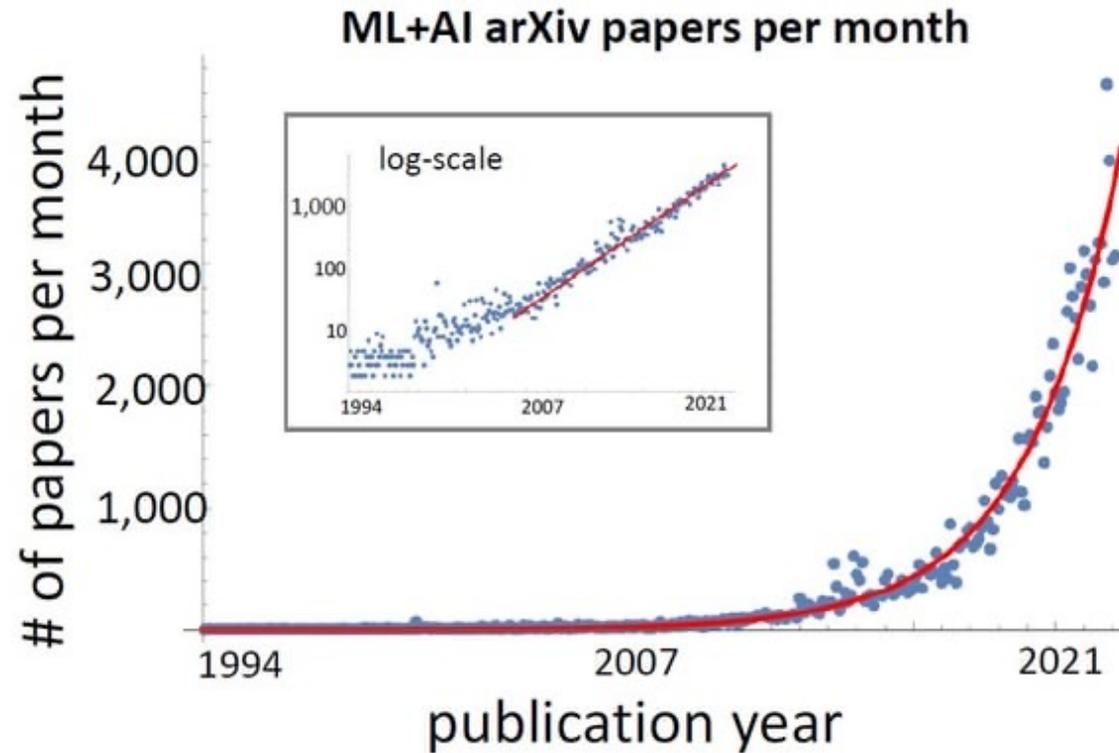


...

- Most RL algorithms require tips and tricks, we will study them

# What is this course about?

- Building the skill of reading and critiquing papers, building research taste



- Understanding the space of “good” ideas, and being able to produce them

# What is this course not about?

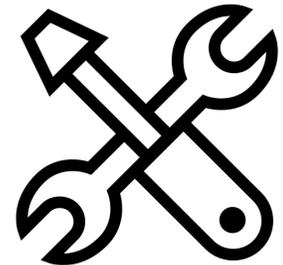
- Not a pure theory course, more an applied-RL course
  - For pure theory classes, recommend CSE 541/CSE 542
  - RL theory book ([https://rltheorybook.github.io/rltheorybook\\_AJKS.pdf](https://rltheorybook.github.io/rltheorybook_AJKS.pdf))

**Lemma 2.10.** Let  $\delta > 0$ . With probability greater than  $1 - \delta$ ,

$$|(P - \hat{P})V^*| \leq \sqrt{\frac{2 \log(2|\mathcal{S}||\mathcal{A}|/\delta)}{N}} \sqrt{\text{Var}_P(V^*)} + \frac{1}{1-\gamma} \frac{2 \log(2|\mathcal{S}||\mathcal{A}|/\delta)}{3N} \mathbf{1}.$$

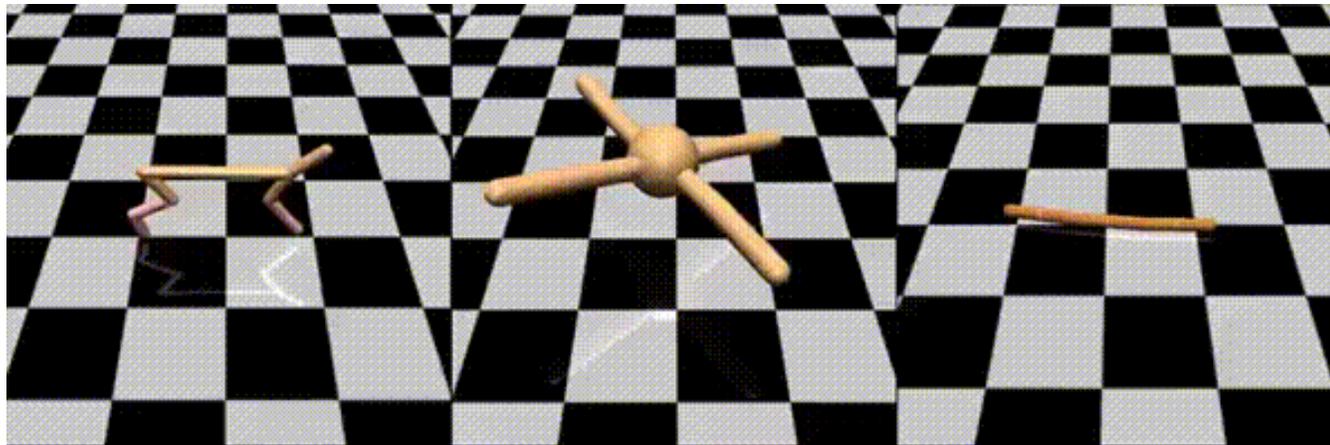
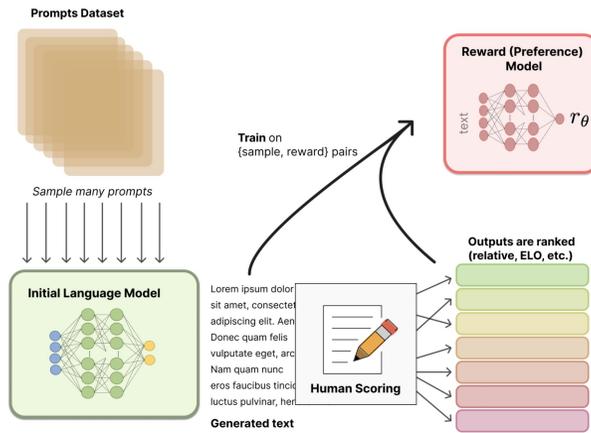
**Theorem 4.3** (FQI guarantee). Fix  $K \in \mathbb{N}^+$ . Fitted Q Iteration guarantees that with probability  $1 - \delta$ ,

$$V^* - V^{\pi^K} \leq \frac{1}{(1-\gamma)^2} \left( \sqrt{\frac{22CV_{\max}^2 \ln(|\mathcal{F}|^2 K/\delta)}{n}} + \sqrt{20C\epsilon_{\text{approx},\nu}} \right) + \frac{\gamma^K V_{\max}}{(1-\gamma)}$$



- Only cover the theory needed to derive algorithms

# What should we be able to do post CSE 579?

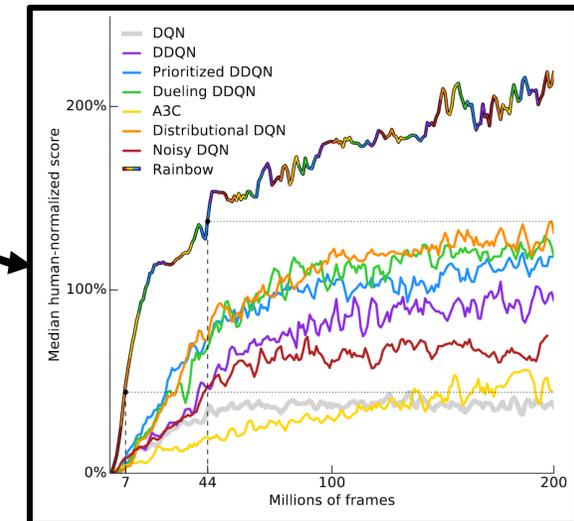


```
def update_critic(self, obs, action, reward, next_obs, not_done, logger, step):
    dist = self.actor(next_obs)
    next_action = dist.rsample()
    log_prob = dist.log_prob(next_action).sum(-1, keepdim=True)
    target_Q1, target_Q2 = self.critic_target(next_obs, next_action)
    target_V = torch.min(target_Q1,
                        target_Q2) - self.alpha.detach() * log_prob
    target_Q = reward + (not_done * self.discount * target_V)
    target_Q = target_Q.detach()

    # get current Q estimates
    current_Q1, current_Q2 = self.critic(obs, action)
    critic_loss = F.mse_loss(current_Q1, target_Q) + F.mse_loss(
        current_Q2, target_Q)
    logger.log('train_critic/loss', critic_loss, step)

    # Optimize the critic
    self.critic_optimizer.zero_grad()
    critic_loss.backward()
    self.critic_optimizer.step()

    self.critic.log(logger, step)
```



# Lecture Outline

---

Course logistics and scope



What is RL, a formal definition



Why should we care?



Going beyond RL

# Ok so let's try and define Reinforcement Learning

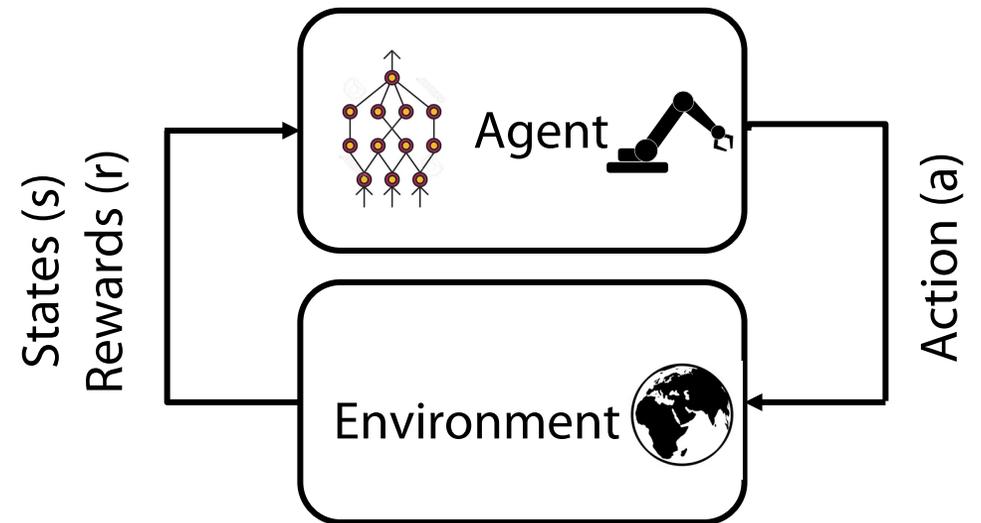
Can learn in arbitrary settings

Go from expert label  $\rightarrow$   
scalar measure of success

No expert corpus needed

Can generalize to new states

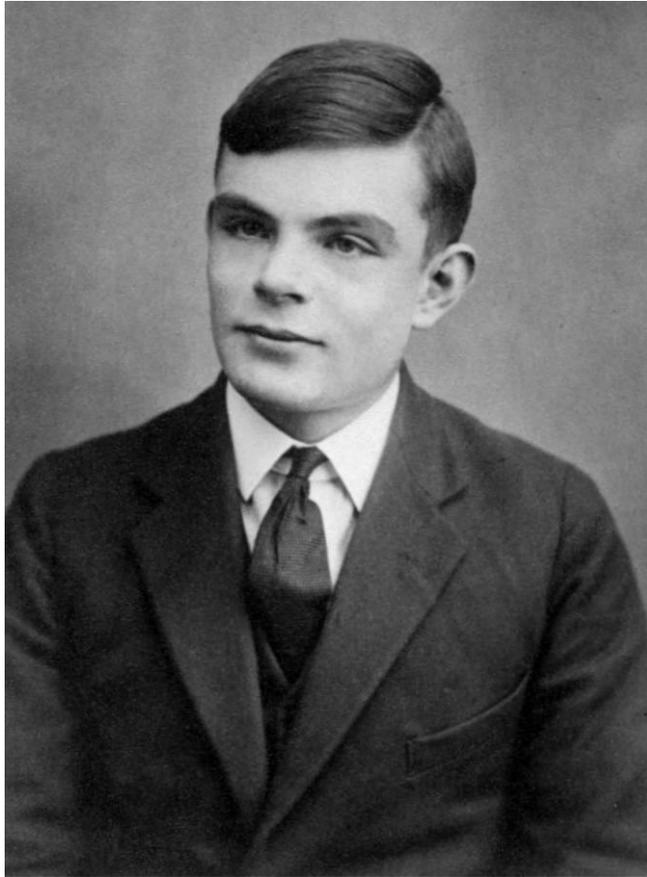
Using **trial and error** in an **environment** to learn a strategy to maximize some notion of **"reward"**



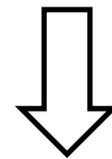
Easy (?) way for agents to continue improving their own behavior on deployment

# Why reinforcement learning?: Philosophical

Hypothesis: By designing algorithms that can improve themselves, we can reach fully intelligent systems



“Instead of trying to produce a programme to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain” – Alan Turing



Rather than try to directly replicate behaviors, try to replicate adaptative learning mechanisms

# Why reinforcement learning?: Practical

A useful tool for building continually improving robots!

Robot learning is amenable to RL for several reasons:

1. Sequential decision making problem (Non IID)
2. Large amounts of expert robot data may be expensive
3. Naturally multi-task and continual
4. Behaviors may be hard to pre-program

Robots that collect their own data to improve!

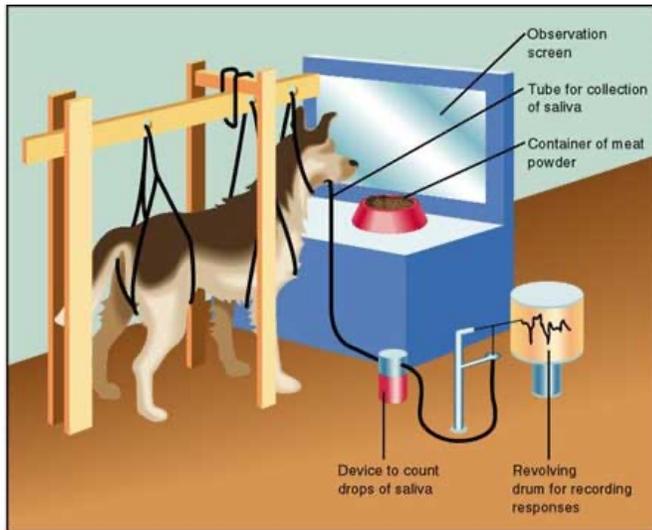
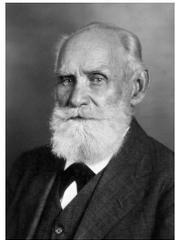


# A Little History on Reinforcement Learning

Two distinct threads converged to give rise to modern RL

Animal Psychology

Optimal Control



$$\min_{x,u} \int_0^x L(t, x(t), u(t)) . dx$$

w.r.t

$$x'(t) = f(x(t), u(t))$$

Ideas from temporal difference learning/dynamic programming united these fields!

# A Little History on Reinforcement Learning

Klopf/Sutton/Barto brought together ideas from psych/neuro and computational TD learning

Harry Klopf



Introduced the ideas of  
“generalized reinforcement”  
– linked together TD  
learning and trial and error  
learning from psychology

**Brain Function and Adaptive Systems -  
A Heterostatic Theory**

A. HARRY KLOPF

Sutton



Barto



AD A101476

AFWAL-TR-81-1070

GOAL SEEKING COMPONENTS FOR ADAPTIVE INTELLIGENCE: AN INITIAL ASSESSMENT



Also had major contributions from Watkins, Shannon, Minsky, Tesauro, Michie, Samuel, etc!

# A Little History on Reinforcement Learning

---

**Some** evidence about RL in the brain

## **Reinforcement learning in the brain**

Yael Niv

Psychology Department & Princeton Neuroscience Institute, Princeton University

Shows the importance of temporal difference reward prediction error in processes in the brain

Dopamine != reward, rather dopamine corresponds strongly to errors in long term reward prediction (aka TD errors) (Montague '96, Schultz '97). Some inconsistencies, e.g. Dealing with aversive events like pain

Likely much more research needed, since decisions can be made in the absence of dopamine → multiple different RL processes in the brain

# A Little History on Modern Reinforcement Learning (my view)

---

## Playing Atari with Deep Reinforcement Learning

---

Volodymyr Mnih   Koray Kavukcuoglu   David Silver   Alex Graves   Ioannis Antonoglou

Daan Wierstra   Martin Riedmiller

DeepMind Technologies



# A Little History on Modern Reinforcement Learning (my view)

---

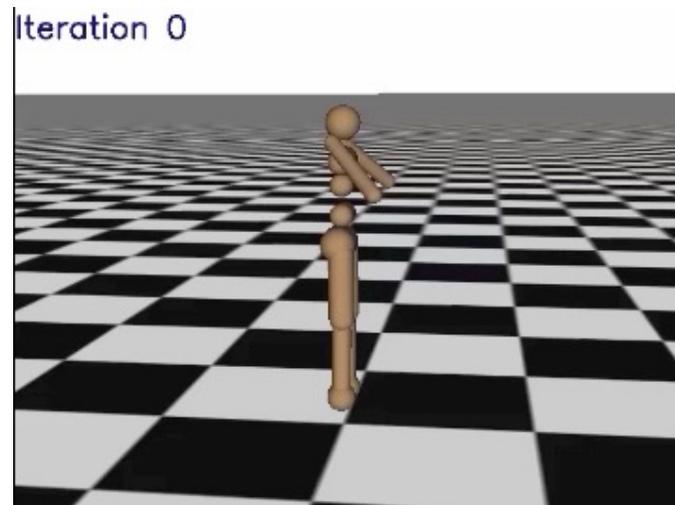
## Trust Region Policy Optimization

---

**John Schulman**  
**Sergey Levine**  
**Philipp Moritz**  
**Michael Jordan**  
**Pieter Abbeel**

University of California, Berkeley, Department of Electrical Engineering and Computer Sciences

JOSCHU@EECS.BERKELEY.EDU  
SLEVINE@EECS.BERKELEY.EDU  
PCMORITZ@EECS.BERKELEY.EDU  
JORDAN@CS.BERKELEY.EDU  
PABBEEL@CS.BERKELEY.EDU



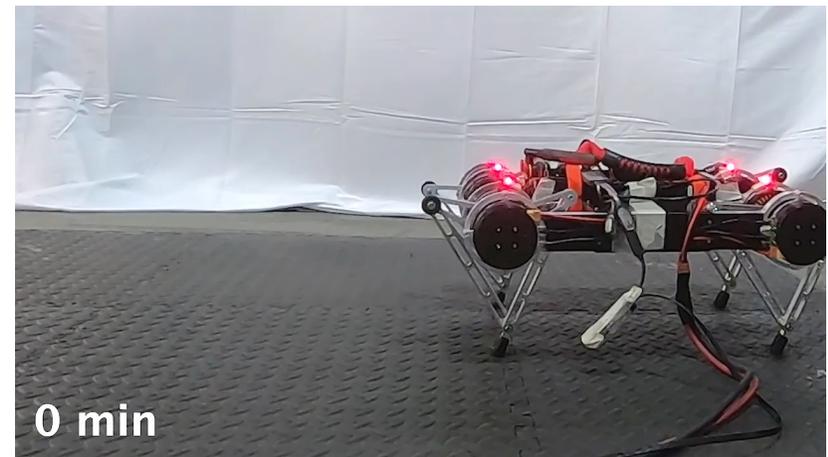
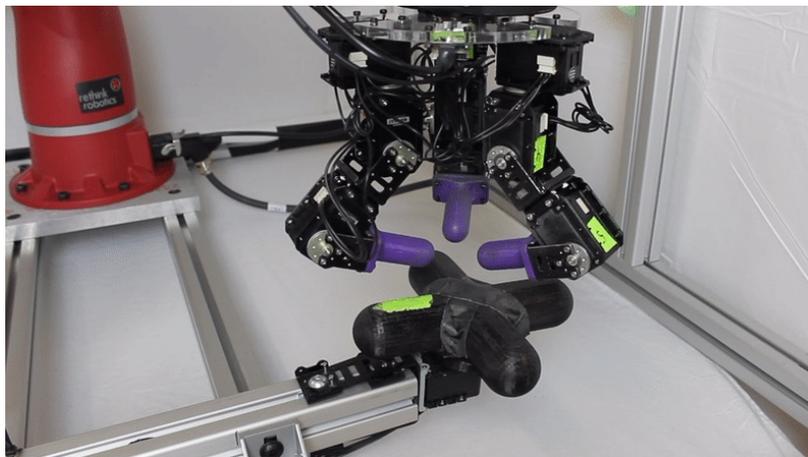
# A Little History on Modern Reinforcement Learning (my view)

---

## **Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor**

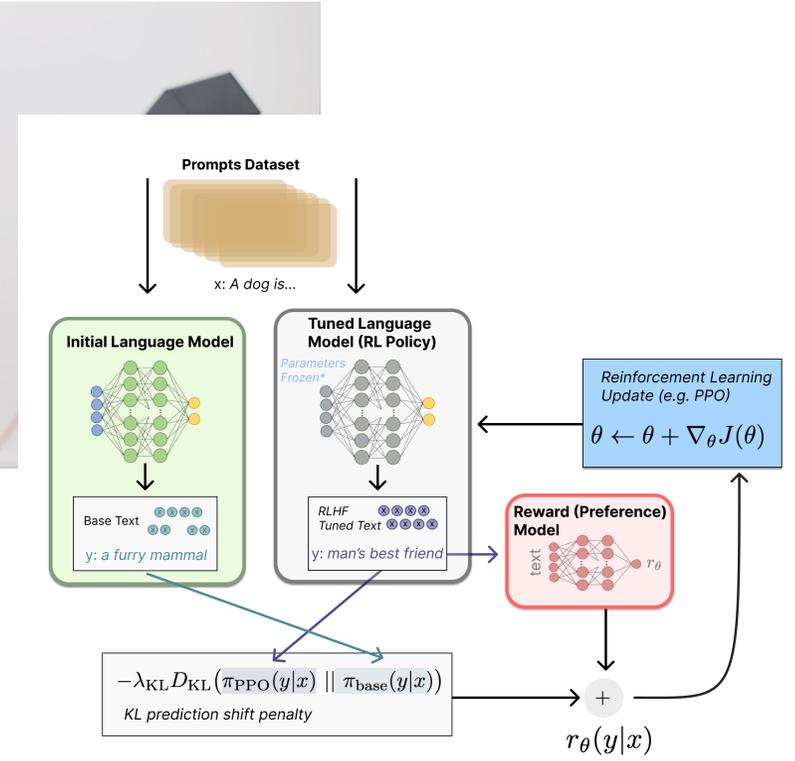
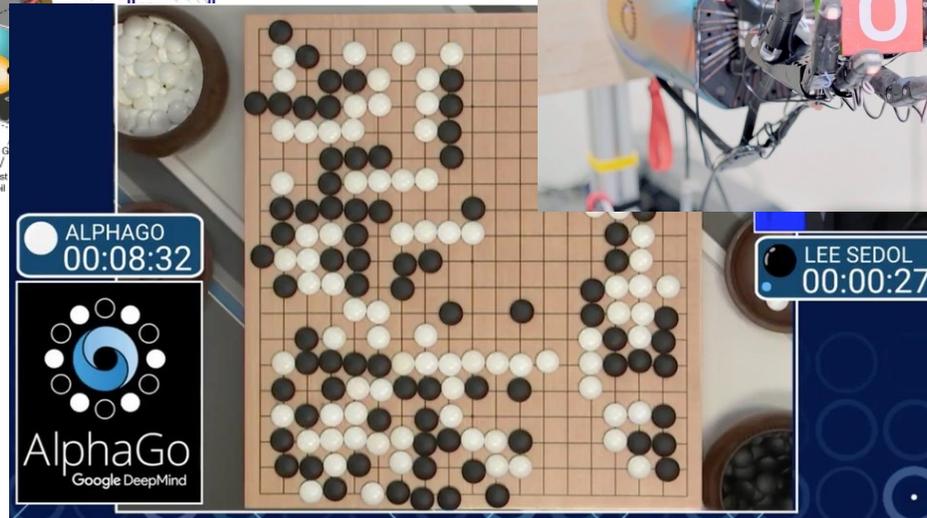
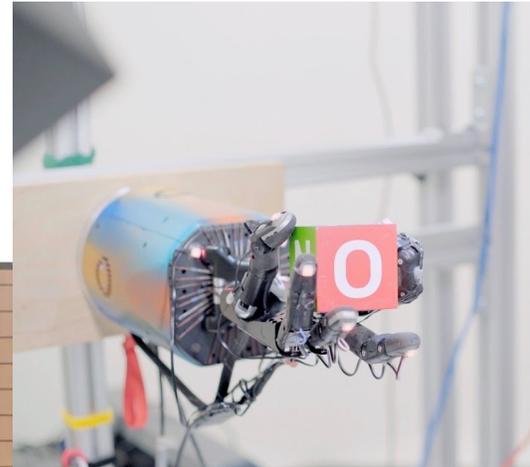
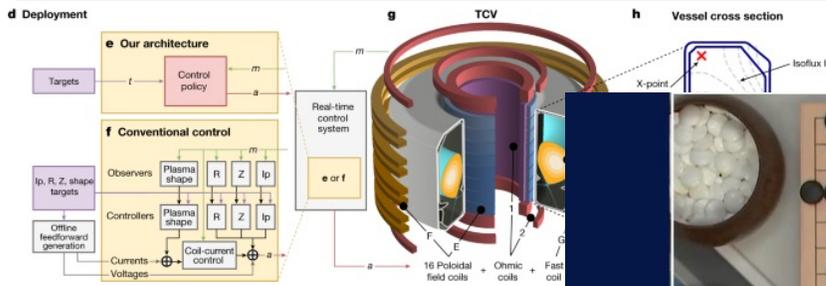
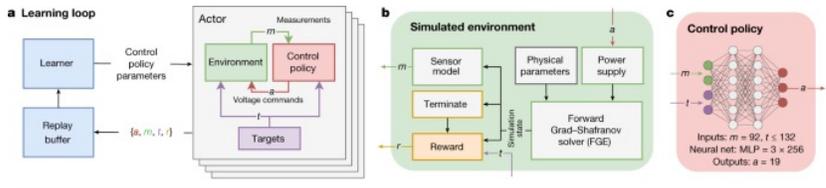
---

**Tuomas Haarnoja<sup>1</sup> Aurick Zhou<sup>1</sup> Pieter Abbeel<sup>1</sup> Sergey Levine<sup>1</sup>**

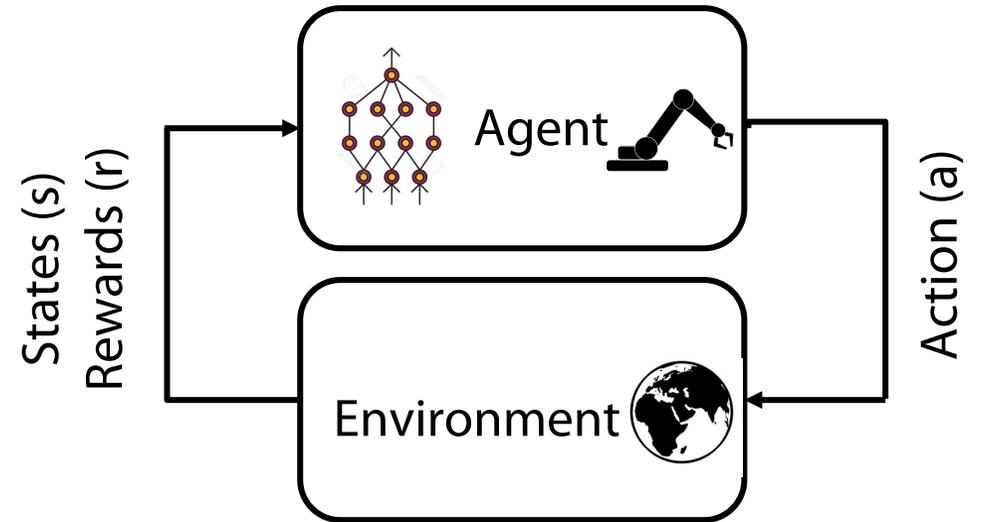


# A Little History on Modern Reinforcement Learning (my view)

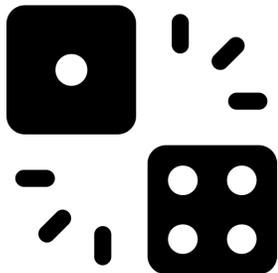
Since then, we have gotten RL now to power a variety of high-impact applications



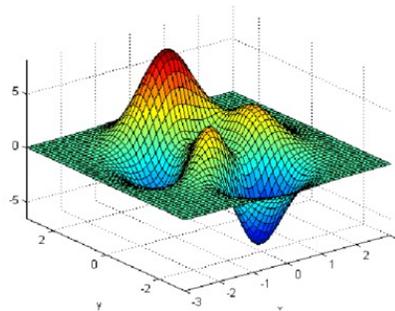
# Let's define a formalism



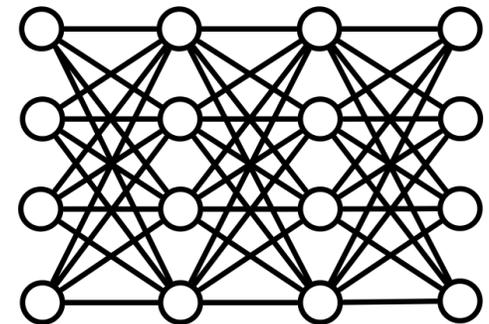
## Probability theory



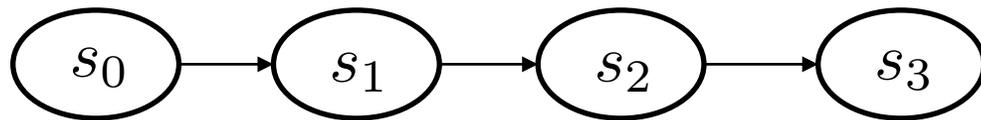
## Optimization



## Machine Learning



# Preliminaries: Markov Chains



Initial state distribution

Transition Kernel (T)



Future is independent of past, conditioned on the present

$$p(s_1, s_2, s_3) = p(s_3|s_2)p(s_2|s_1)p(s_1)$$

Goal of Markov chain: running the Markov chain leads to sampling from stationary distribution  $d^\pi$

Balance equation

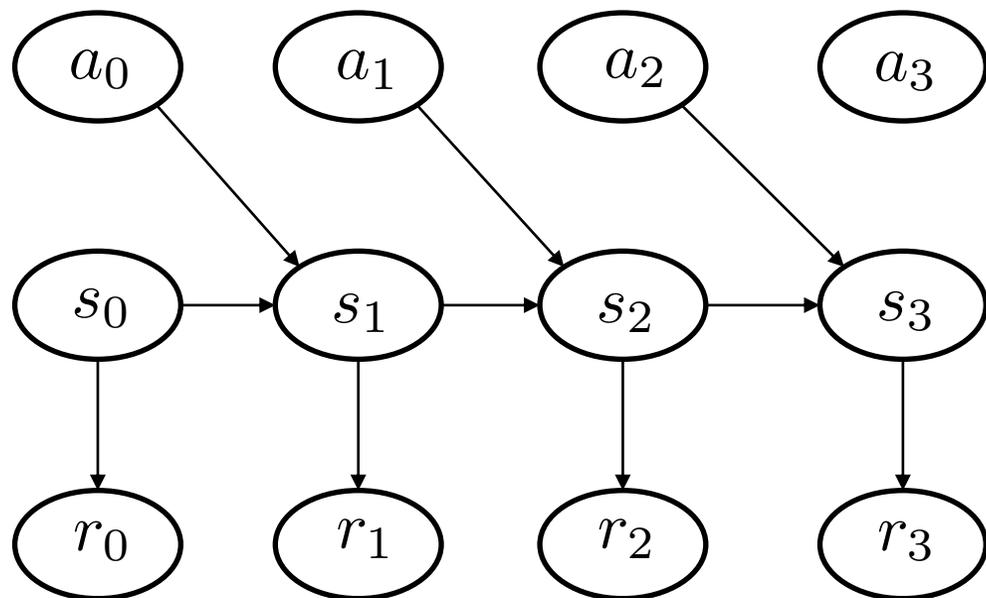
$$T d^\pi = d^\pi$$

Useful in sampling based inference eg MCMC

How can we make this useful for decision making?

# Framework for RL - Markov Decision Process

Augment Markov chain with rewards and actions



States:  $\mathcal{S}$

Initial state dist:  $\rho_0(s)$

Actions:  $\mathcal{A}$

Discount:  $\gamma$

Rewards:  $\mathcal{R}$

Transition Dynamics -  $p(s_{t+1}|s_t, a_t)$

Markov property  $p(s_1, s_2, s_3) = p(s_3|s_2)p(s_2|s_1)p(s_1)$

Trajectory  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T)$

# Mapping MDPs to the Real World

Task: Place kettle in sink



State: Camera Images / Joint Encoders

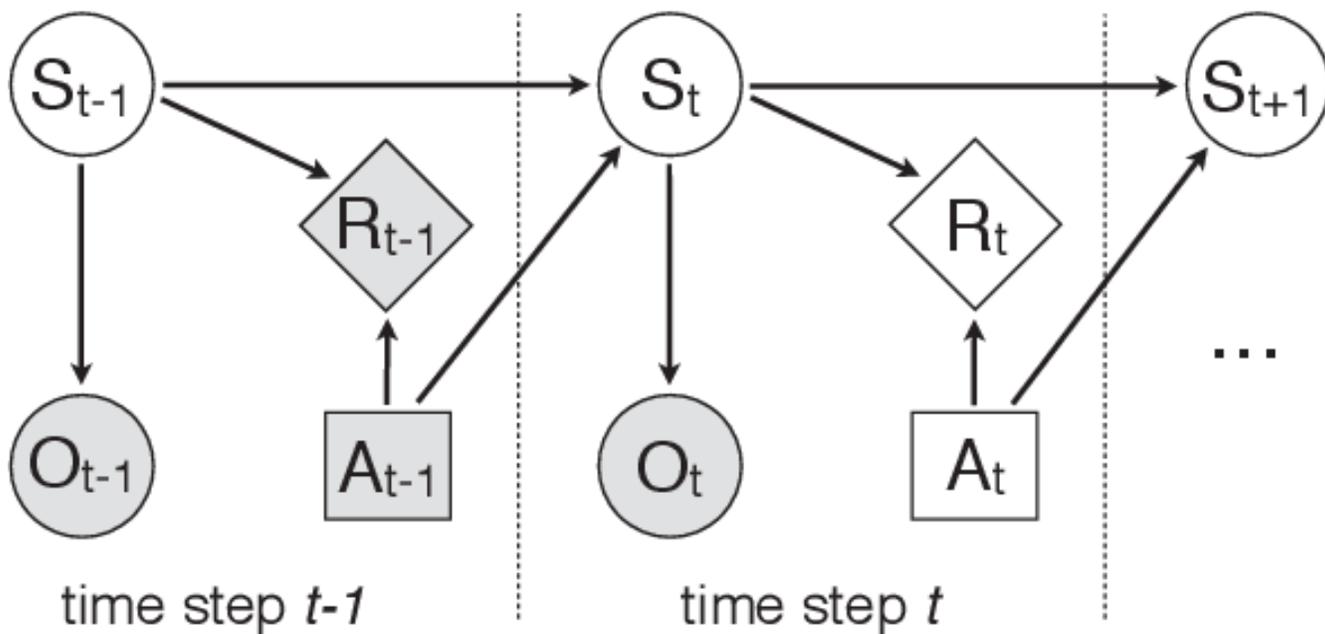
Action: Joint torques/velocities

Reward: Distance from kettle to sink

Transition: World physics

# Aside: Partially Observed MDPs

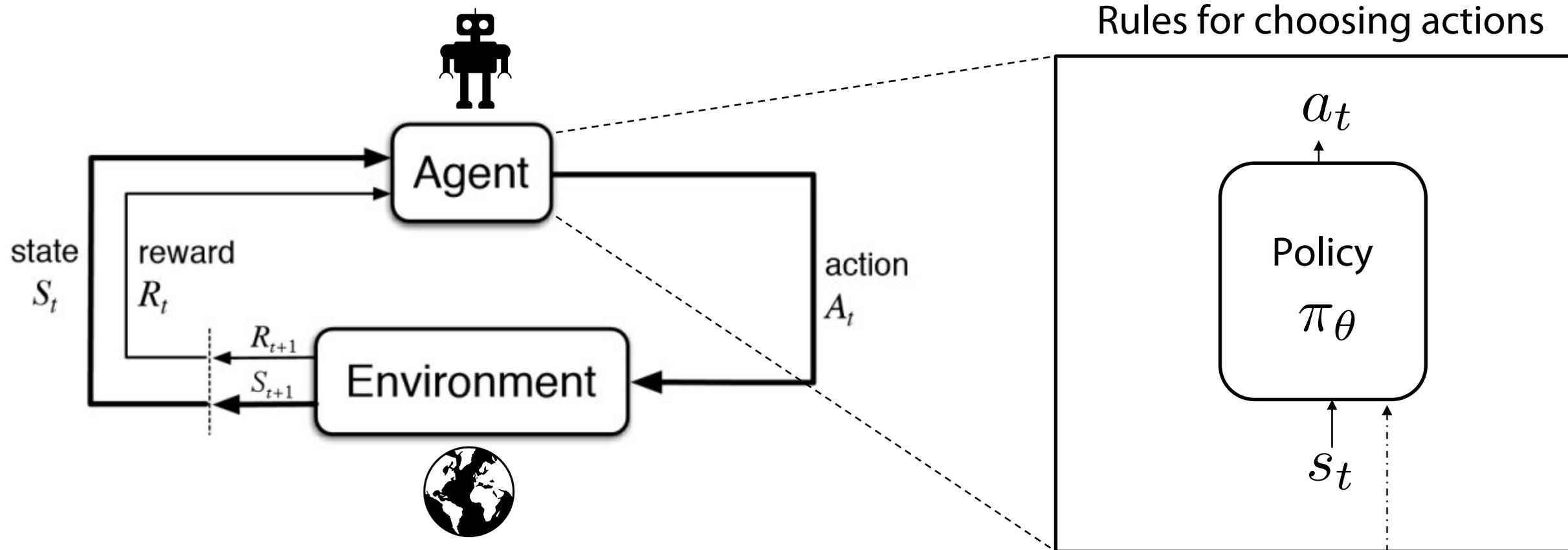
Not every environment is an MDP, in-fact most are POMDPs



POMDPs are hard, we will try to avoid them!



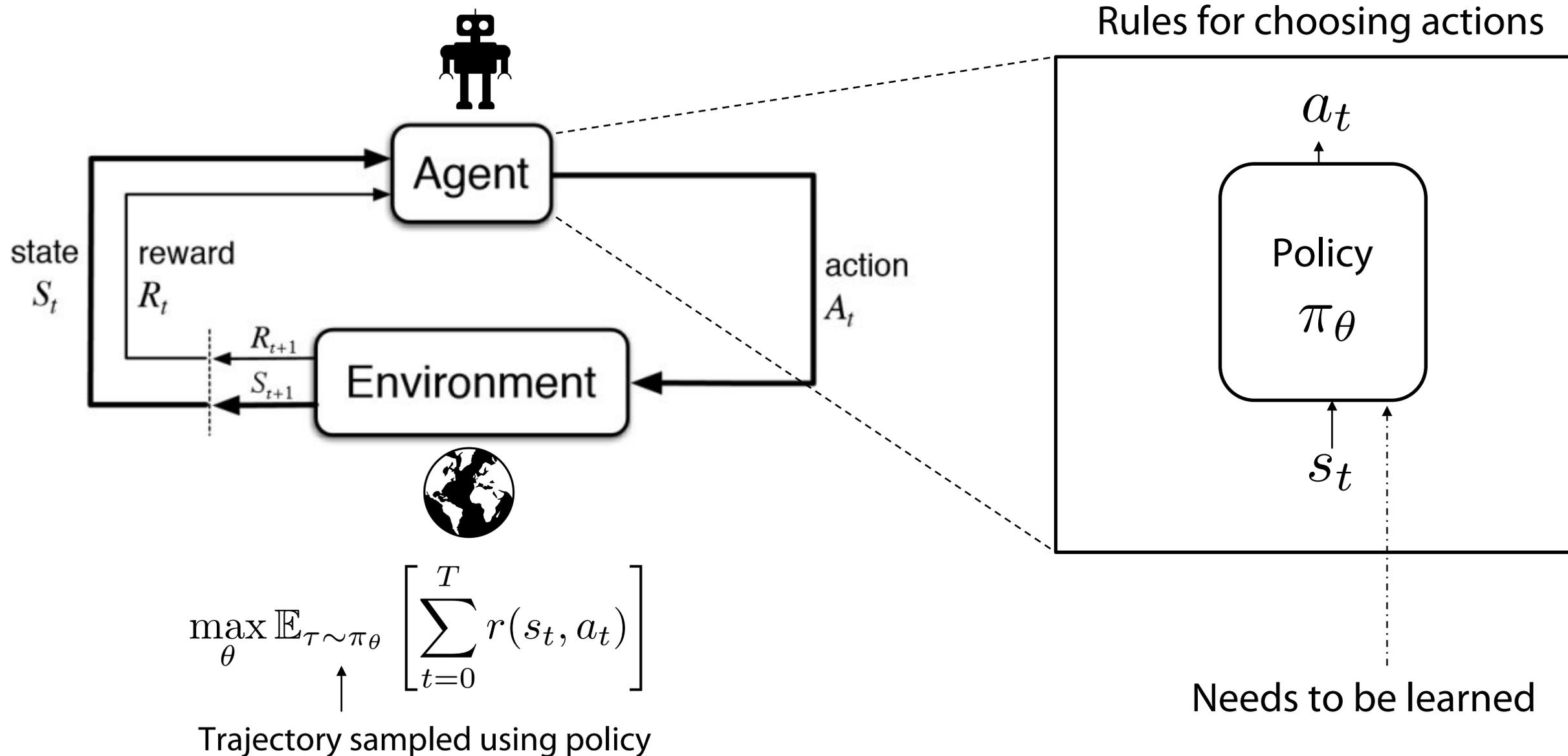
# Reinforcement Learning Formalism



Maximize the sum of expected rewards under policy

Needs to be learned

# Reinforcement Learning Formalism



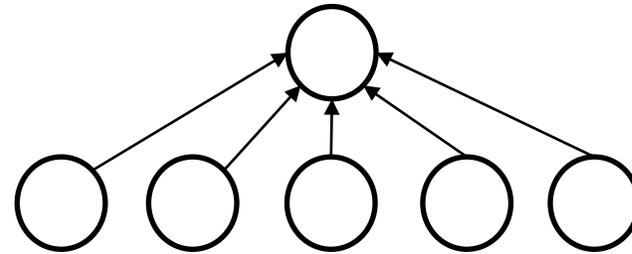
# Main thing to learn - Policies

Policies are **mappings** from states to distributions over actions

Tabular

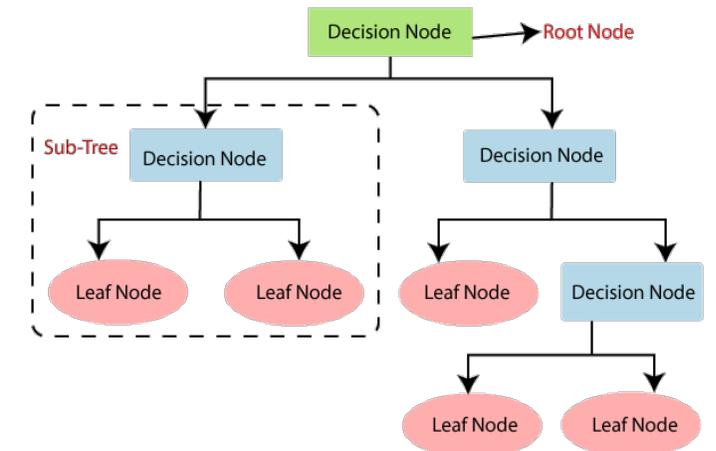
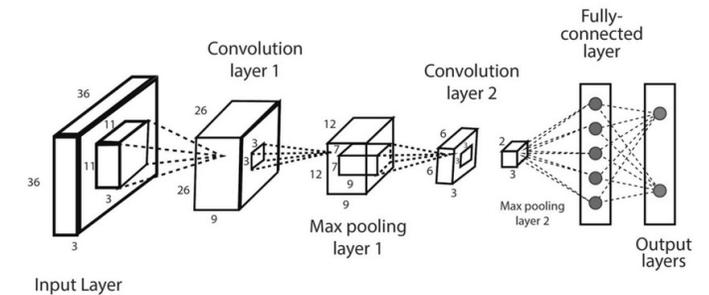
8.67	8.93	9.11	9.30	9.42
8.49		9.09	9.42	9.68
8.33		1.00		10.00
7.13	5.04	3.15	5.68	8.45
-10.00	-10.00	-10.00	-10.00	-10.00

Linear



$$\pi(a|s) = \langle \phi(s, a), w \rangle$$

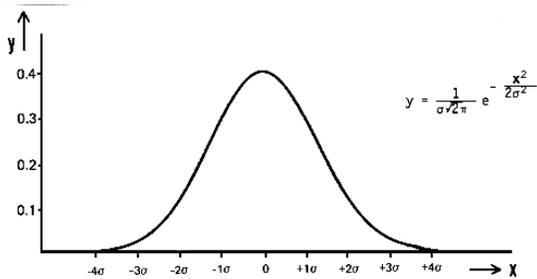
Arbitrary function approx



# Main thing to learn - Policies

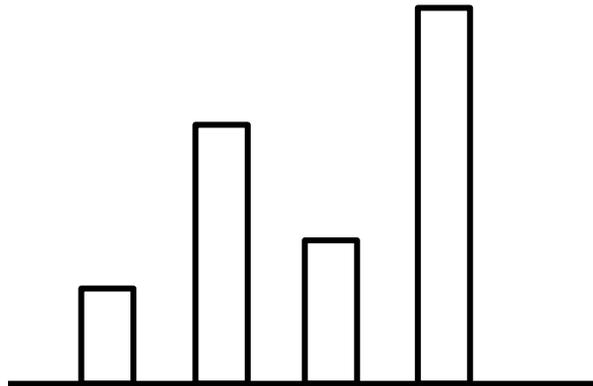
Policies are mappings from states to **distributions** over actions

Gaussian



$$\mu(s), \Sigma(s)$$

Categorical



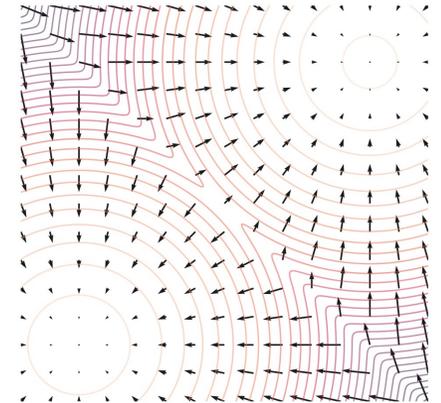
$$p_1(s), p_2(s), \dots, p_k(s)$$

Mixture of Gaussians



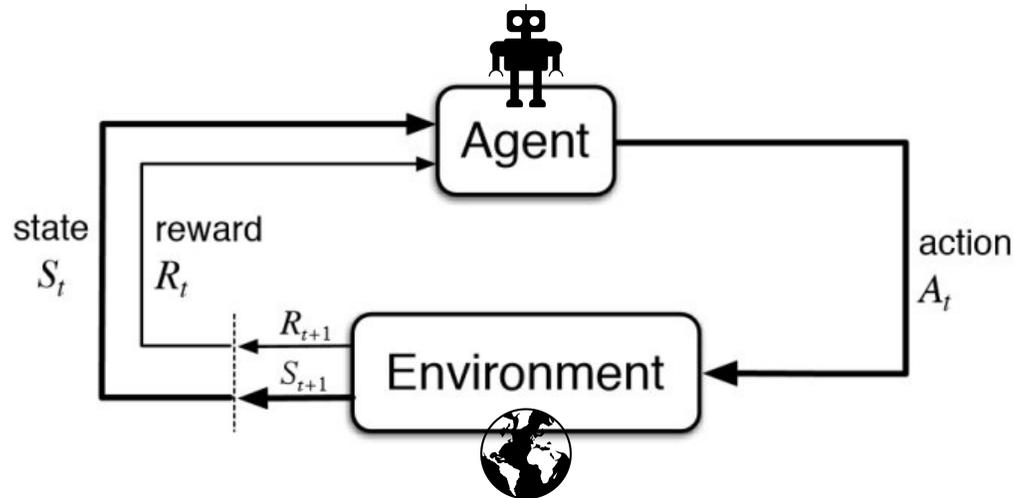
$$\begin{aligned} &\mu_1(s), \Sigma_1(s), w_1 \\ &\mu_2(s), \Sigma_2(s), w_2 \\ &\dots \\ &\mu_N(s), \Sigma_N(s), w_N \end{aligned}$$

Diffusion Models



$$\nabla_a \log \pi(a|s)$$

# Let's revisit policies: stochastic vs deterministic



$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

Lemma 1: Every MDP has atleast one optimal \*deterministic\* policy

**Theorem 1.7.** Let  $\Pi$  be the set of all non-stationary and randomized policies. Define:

$$V^*(s) := \sup_{\pi \in \Pi} V^{\pi}(s)$$
$$Q^*(s, a) := \sup_{\pi \in \Pi} Q^{\pi}(s, a).$$

which is finite since  $V^{\pi}(s)$  and  $Q^{\pi}(s, a)$  are bounded between 0 and  $1/(1 - \gamma)$ .

There exists a stationary and deterministic policy  $\pi$  such that for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}$ ,

$$V^{\pi}(s) = V^*(s)$$
$$Q^{\pi}(s, a) = Q^*(s, a).$$

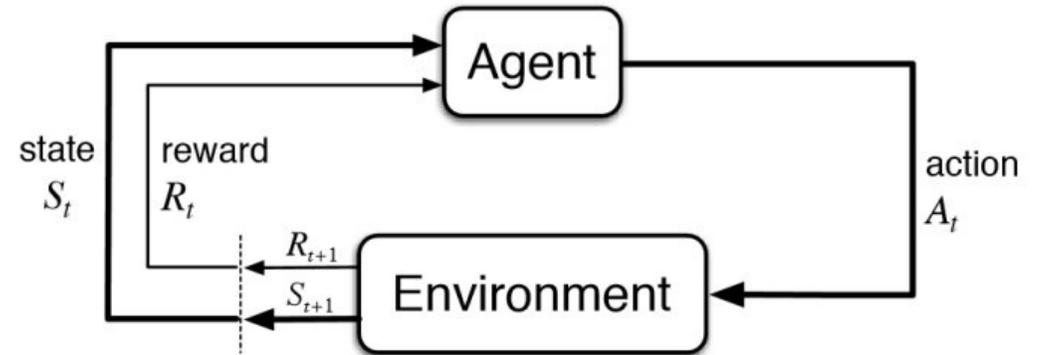
We refer to such a  $\pi$  as an optimal policy.

Intuition: pick the best possible action at every state

Stochastic policies will help in the search/optimization process for finding (close to) deterministic policies

# Let's take a closer look at the objective: horizon

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$



## Finite horizon

$$\mathbb{E}_{\pi_{\theta}^t} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

Time-dependent policy  
(not stationary)

---


$$\mathbb{E}_{\pi_{\theta}^t} \left[ \sum_{t=0}^T \gamma^t r(s_t, a_t) \right]$$

## Infinite horizon discounted

$$\mathbb{E}_{\pi_{\theta}} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

Time-independent (stationary) policy  
→ Need discount to prevent blow up

**Lemma:** there always exists a stationary optimal policy

# Unpacking the Expectation

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

Trajectory View - Ancestral sampling along MDP

Initial state	$\mathbb{E}_{\substack{s_0 \sim \rho_0(s) \\ a_0 \sim \pi_{\theta}(\cdot   s_0) \\ s_1 \sim p(\cdot   s_0, a_0) \\ a_1 \sim \pi_{\theta}(\cdot   s_1) \\ s_2 \sim p(\cdot   s_1, a_1) \\ \dots}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$
Policy	
Dynamics	
Policy	
Dynamics	
...	

Compact representation	$\mathbb{E}_{\substack{s_0 \sim \rho_0(s) \\ a_t \sim \pi_{\theta}(\cdot   s_t) \\ s_{t+1} \sim p(s_{t+1}   s_t, a_t)}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$
------------------------	--

$$\mathbb{E}_{\pi_{\theta}^t} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]$$

Stationary View – sampling from stationary dist

$$d_t^{\pi}(s, a) = \mathbb{P}(s_t = s, a_t = a \mid s_0 \sim \rho_0, \forall i < t, a_i \sim \pi_{\theta}(\cdot | s_i), s_{i+1} \sim p(\cdot | s_i, a_i))$$

(Likelihood of being at state  $s$ , action  $a$  at time step  $t$ )

$$\mu_{\pi}^{\gamma}(s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t d_t^{\pi}(s, a)$$

(Likelihood of being at state  $s$ , action  $a$  across **all** steps)

$\gamma$ subsumed into E	$\mathbb{E}_{(s,a) \sim \mu_{\gamma}^{\pi}(s,a)} \left[ r(s, a) \right]$
--------------------------	--

No sequential sampling

No sum over rewards

# Some notation: Q-functions and V-functions

Estimate of how “good” a policy is – estimate of future returns under a policy  $\pi$

## Q-function

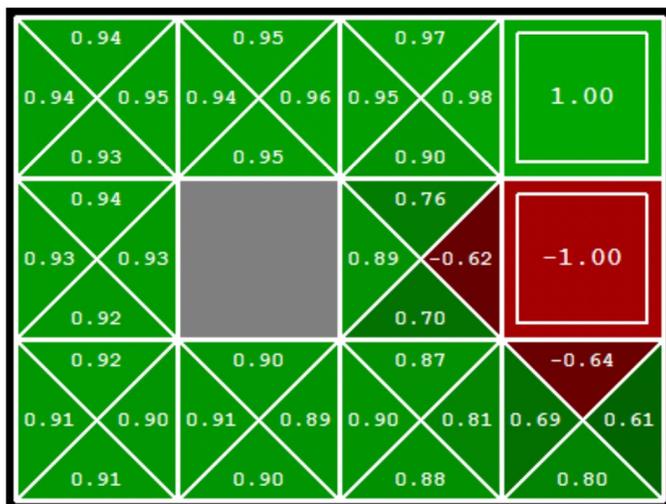
Take one action and then follow policy from  $s$

$$Q^\pi(s, a) = \mathbb{E}_{\pi, p} \left[ \sum_t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]$$

## V-function

Follow policy from  $s$

$$V^\pi(s, a) = \mathbb{E}_{\pi, p} \left[ \sum_t r(s_t, a_t) \mid s_0 = s \right]$$



$$V^\pi(s, a) = \mathbb{E}_{a \sim \pi(\cdot | s)} [Q^\pi(s, a)]$$

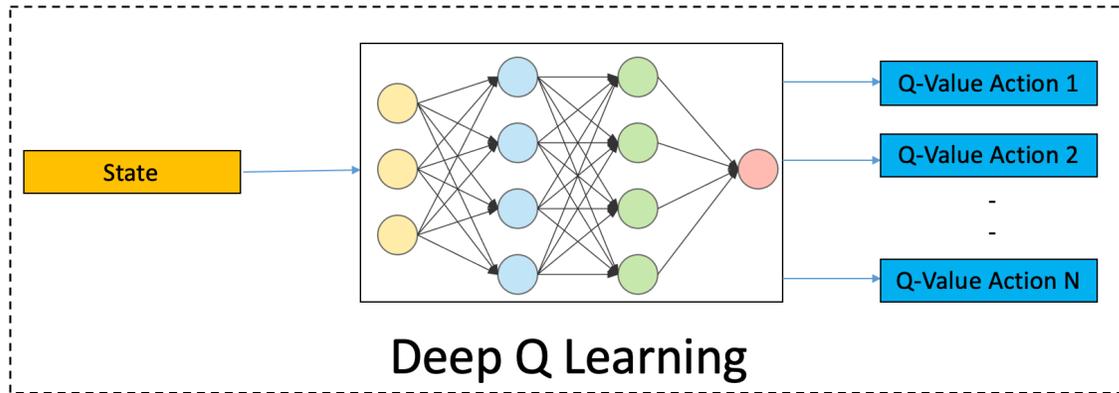
Will be useful soon!

$$J(\pi) = \mathbb{E}_{s \sim \rho_0(s)} [V^\pi(s)]$$

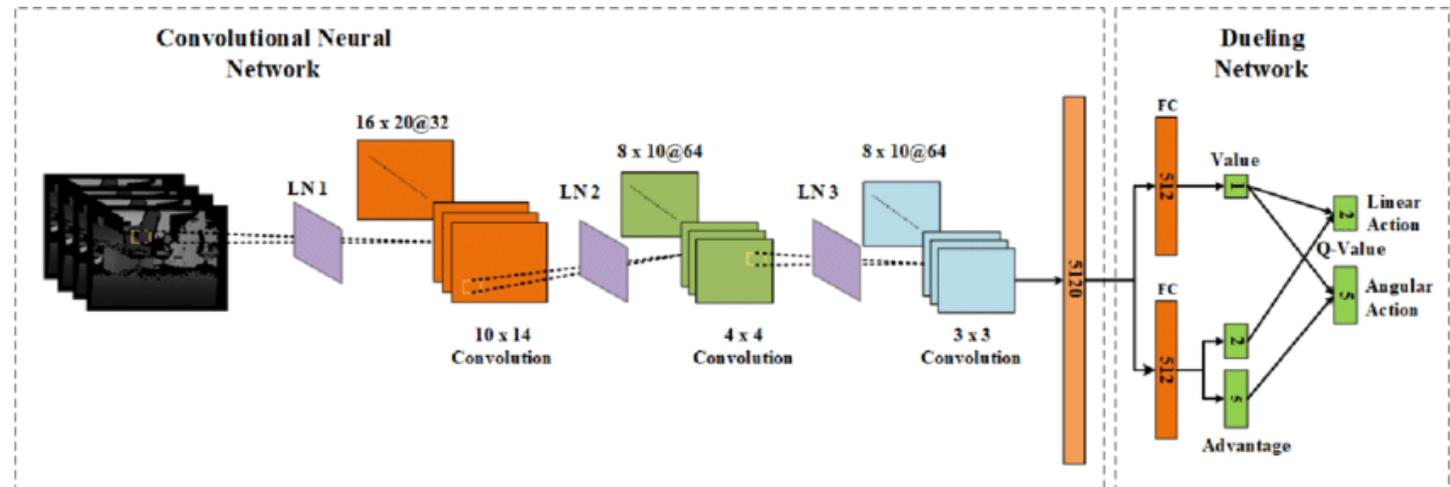
Average value over initial states

# Ok so where does deep learning fit in?

Avoids expensive hand-design for adaptive agents, learn end-to-end: sensors → actions



Policies/Q-values/model are represented as deep neural networks

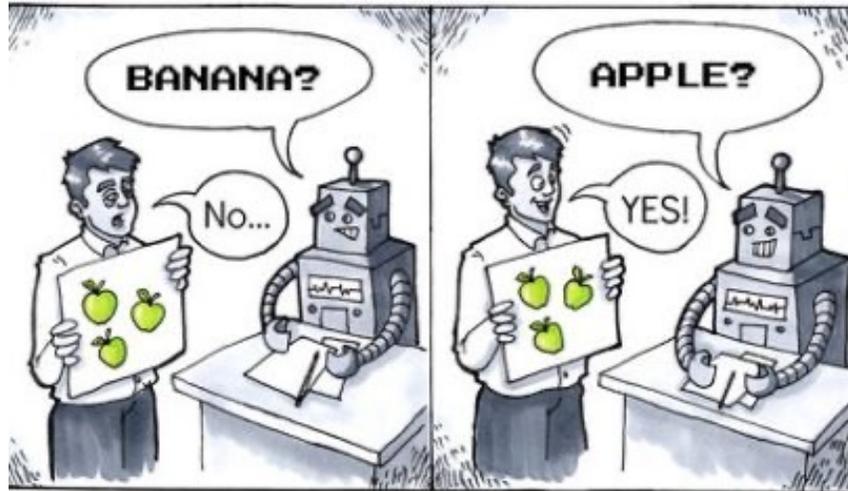


Leads to non-trivial challenges in learning and optimization!



# Ok so is this just supervised learning?

Supervised learning aims to maximize likelihood of observed data under the model



**Supervised Learning**

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log \hat{p}_{\theta}(y|x)]$$

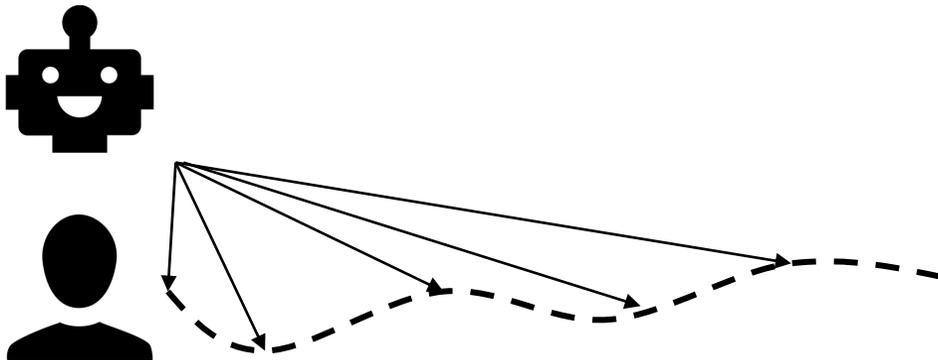
# Why is this not just supervised learning?

## Supervised Learning

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log \hat{p}_{\theta}(y|x)]$$

Sampling from expert

$$D_{\text{KL}}(p^* || p_{\theta}) \quad \text{IID}$$

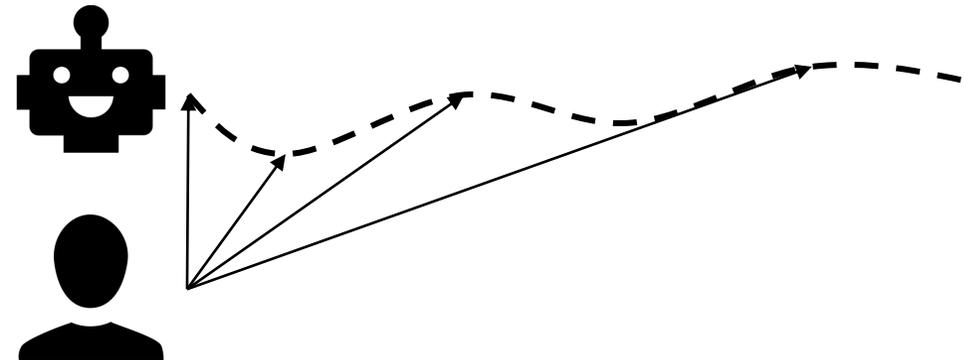


## Reinforcement Learning

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

Sampling from policy

$$D_{\text{KL}}(p_{\theta} || p^*) \quad \text{Non-IID}$$



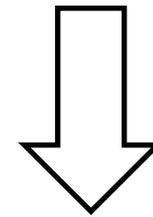
# Why is this not just supervised learning?

## Supervised Learning

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\log \hat{p}_{\theta}(y|x)]$$

The resulting paradigms are different in many ways:

1. Optimization and learning dynamics
2. Balancing exploration and exploitation



## Reinforcement Learning

$$\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r(s_t, a_t) \right]$$

But many overlapping tools! In fact often we try to convert RL into a supervised problem

# Lecture Outline

---

Course logistics and scope



What is RL, a formal definition



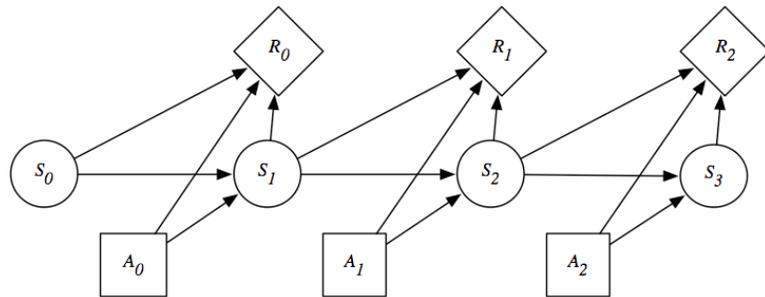
Why should we care?



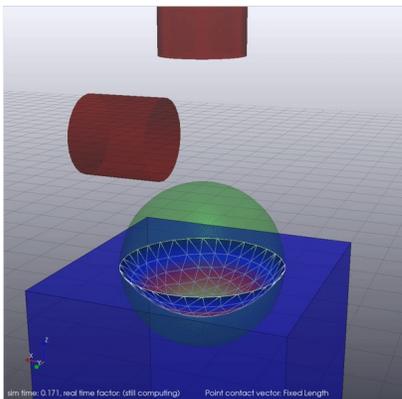
Going beyond RL

# Ok so why should we care about RL?

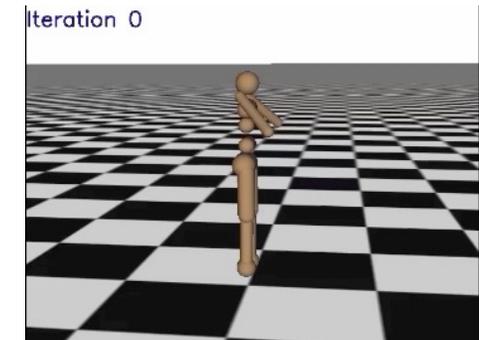
Solves sequential decision making problems



Has black-box assumptions



Enables continual improvement



Reduces burden of human data collection

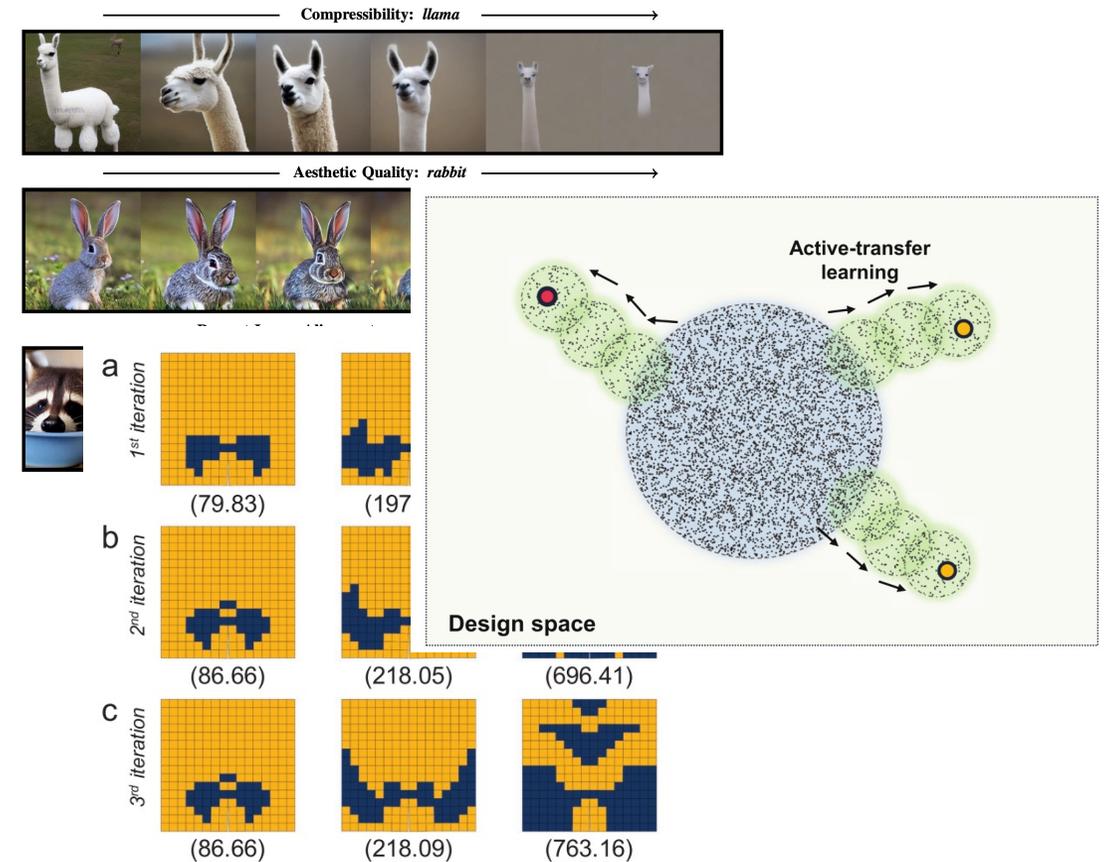


# Reinforcement Learning enables going beyond the data

Generative AI is inherently about **replicating** the data distribution

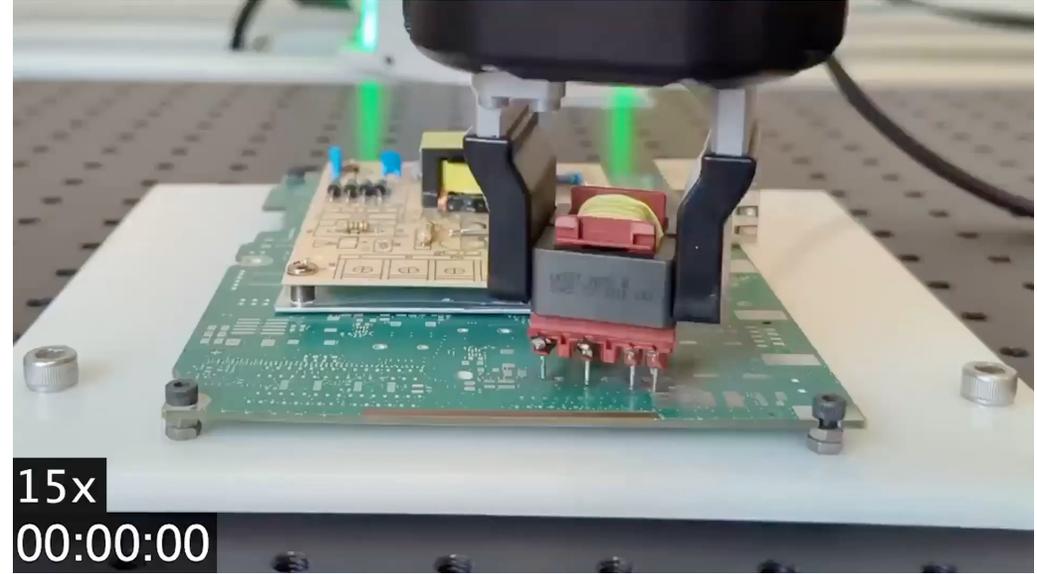
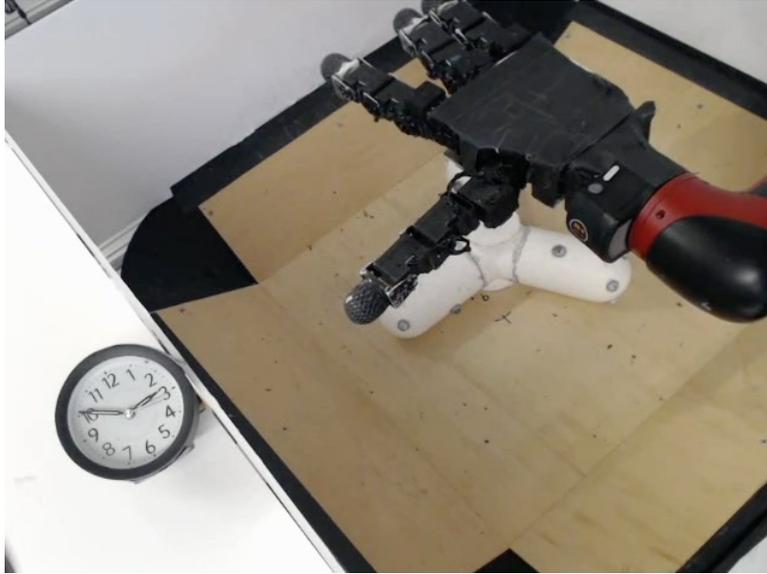


Reinforcement learning can go beyond the data



# Applications of RL: Robotics

RL can enable robotic learning of hard to specify/script behaviors in the presence of contact



# Applications of RL: Large Language Models

Systematically finds and reduces model hallucinations using RLHF

Step 1

**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

Step 2

**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A Explain gravity... B Explain war...  
C Moon is natural satellite of... D People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM  
D > C > A = B

Step 3

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

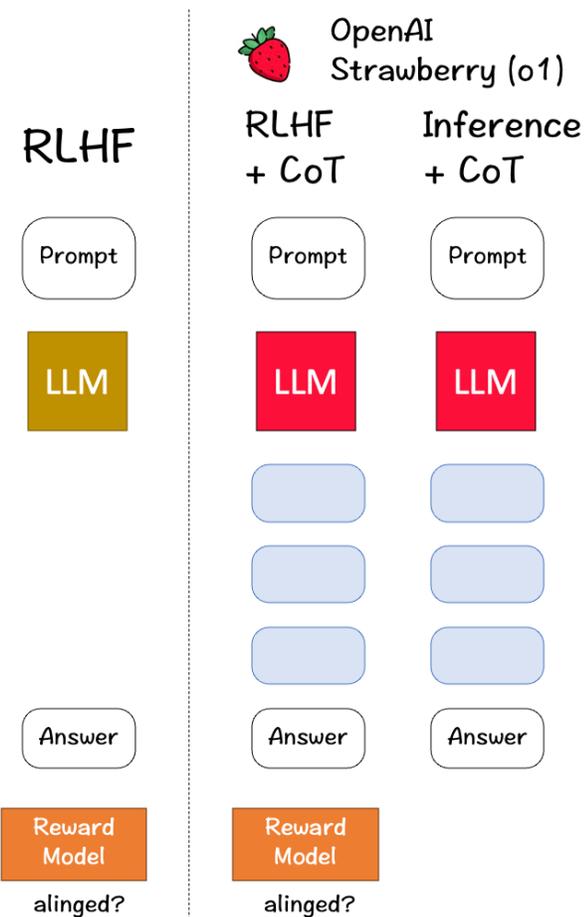
RM

The reward is used to update the policy using PPO.

$r_k$

# Applications of RL: Large Language Models

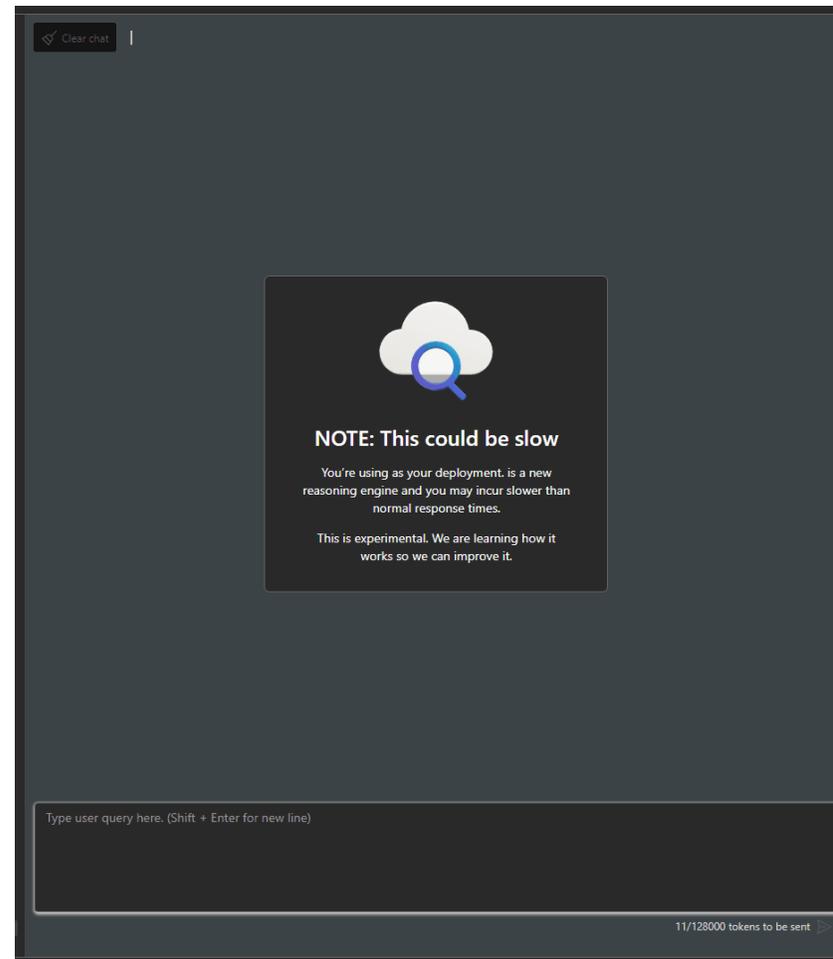
Reasoning-style LLMs are typically RL based



OpenAI Strawberry (o1)

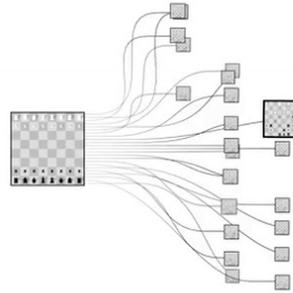
RLHF + CoT      Inference + CoT

Looks a lot like model-based RL



# Applications of RL: Games

Both single and multi-agent RL has proven transformative for game AI

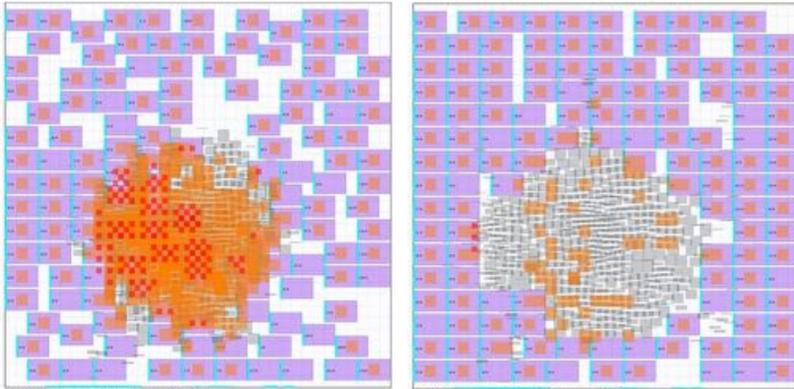


Particularly well suited to RL assumptions

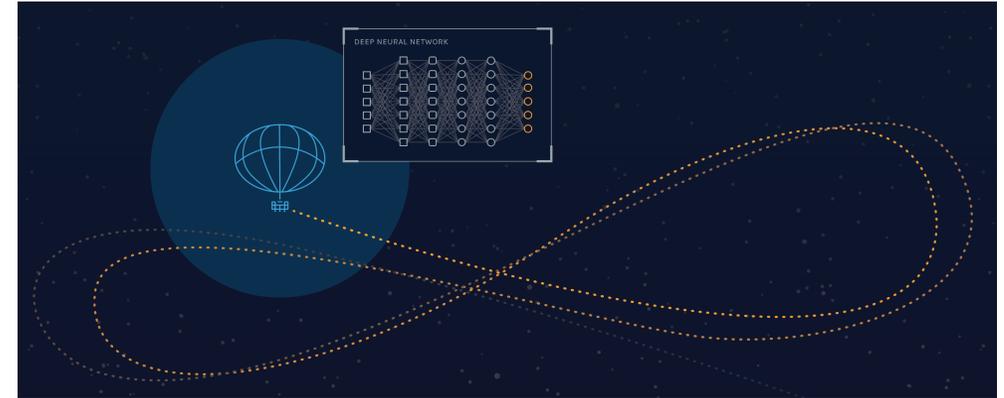
# Applications of RL: Science and Engineering

RL has started to become a useful tool for engineering design

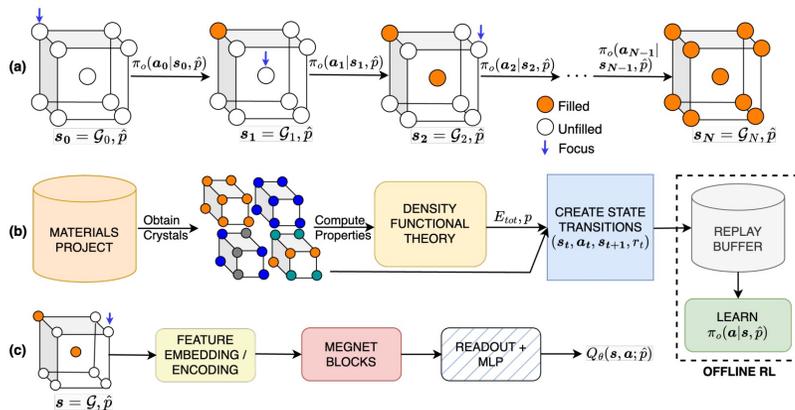
## Chip Design



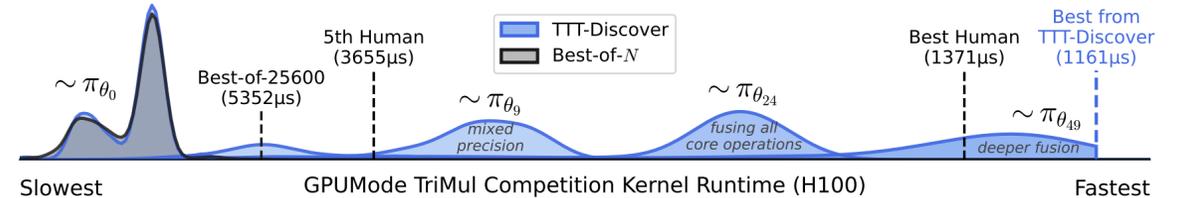
## Weather balloon navigation



## Crystal design

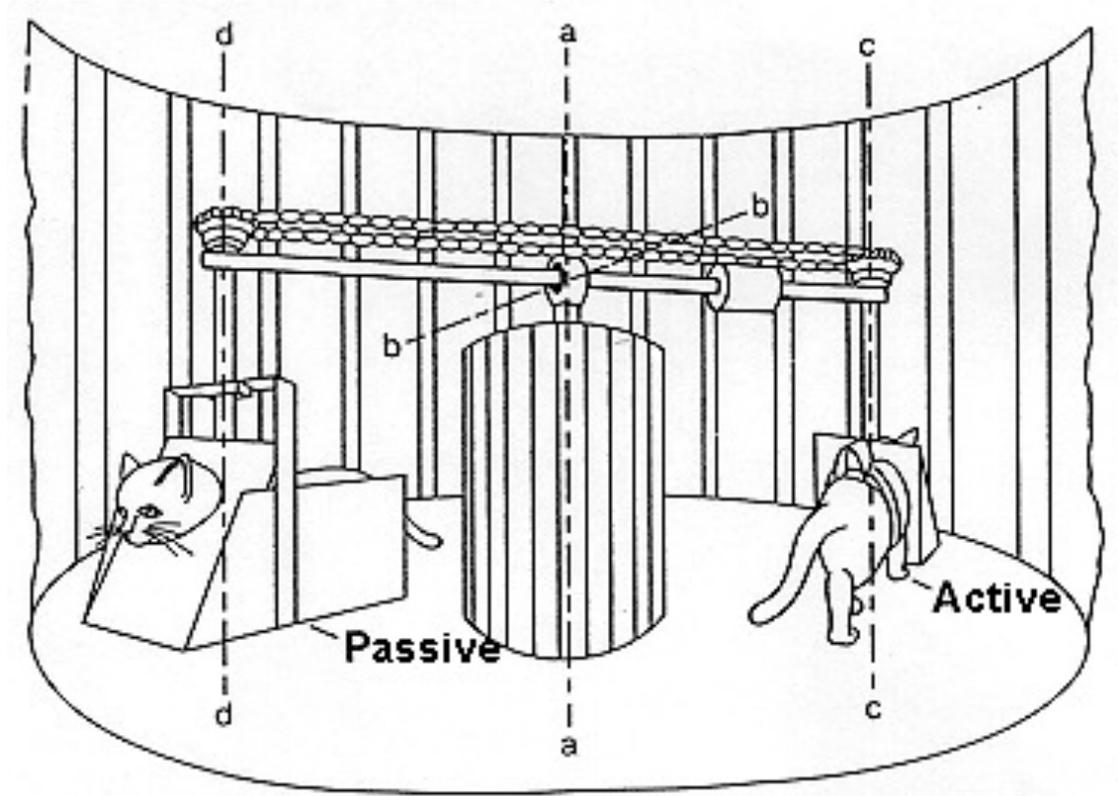
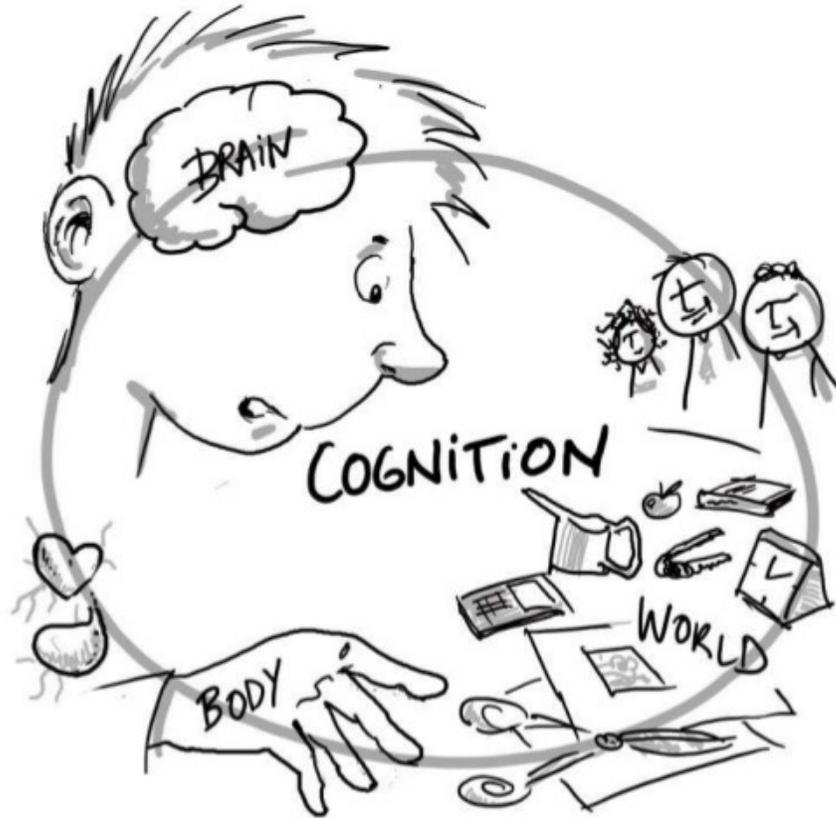


## Kernel Optimization



# Zooming out – why this matters for the study of intelligence?

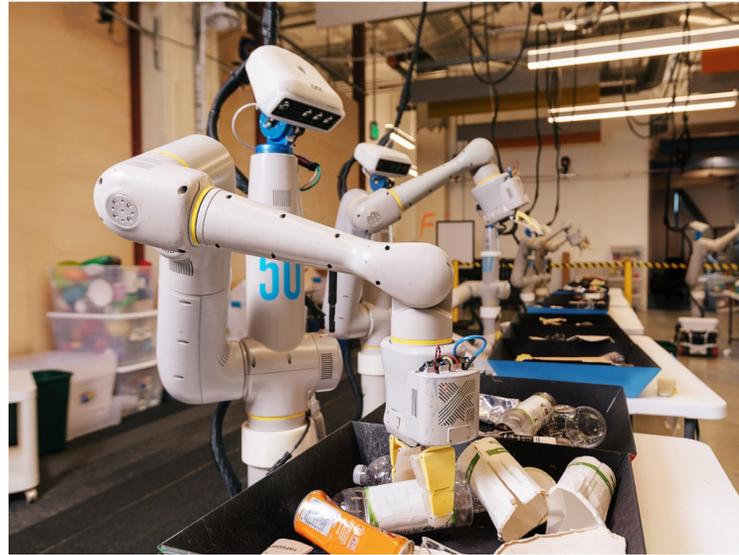
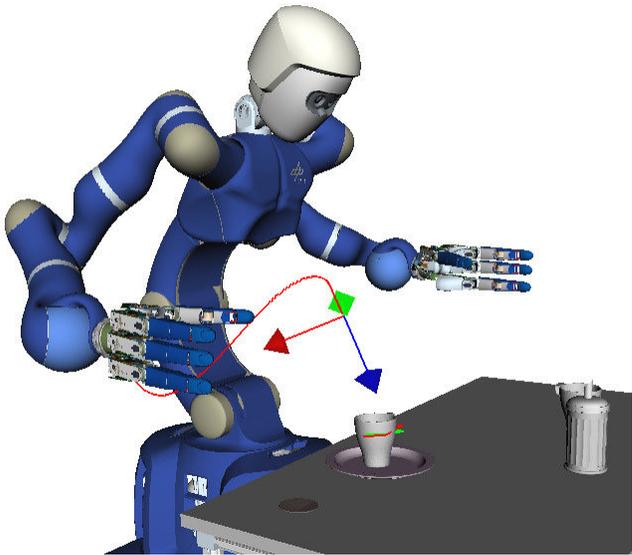
Hypothesis: Intelligence with and without embodiment looks drastically different



Elephants don't play chess!

# Why must we study RL in the real world?

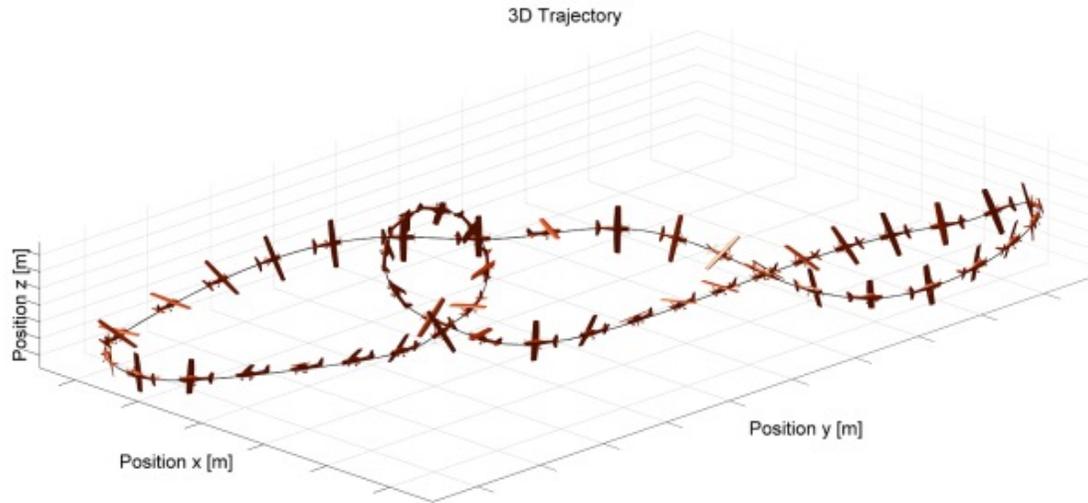
Hypothesis: Agents that learn with embodiment will have emergent complexity in complex, dynamic environments



Increasingly interesting behavior

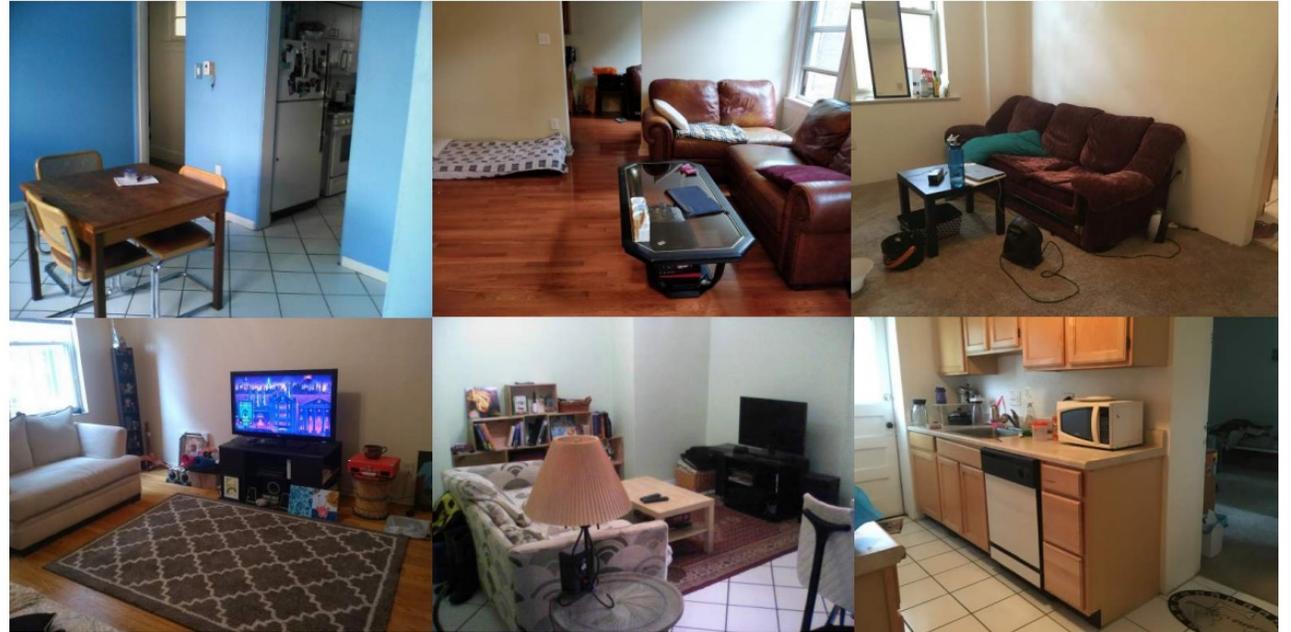
# Where is Reinforcement Learning not useful?

Not the right call for very safety-critical, repetitive applications



# Where is Reinforcement Learning “potentially” useful?

Domains which have high diversity, yet relatively cheap autonomous data collection



But these domains are not as simple as just running RL algorithms!

# Lecture Outline

---

Course logistics and scope



What is RL, a formal definition

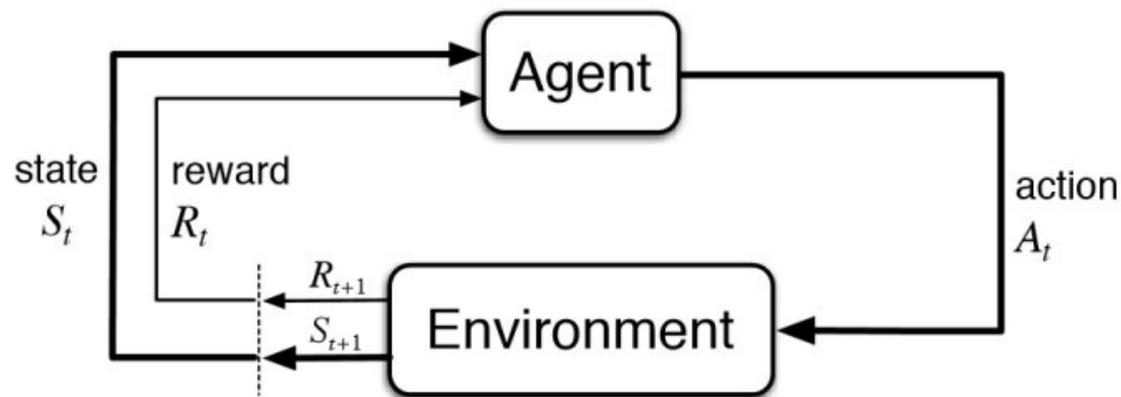


Why should we care?



Going beyond RL

# So is sequential decision making = RL?

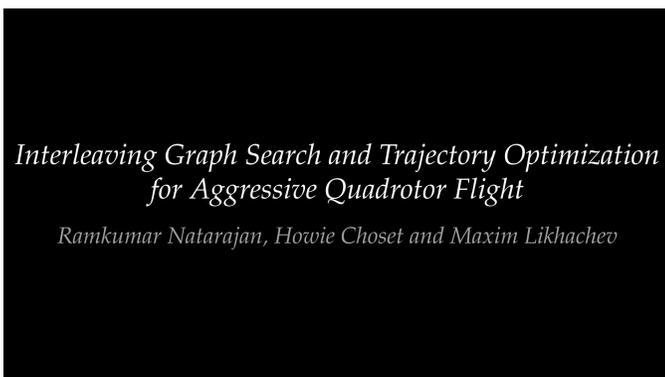


We conflated sequential decision making and RL!

RL is sequential decision making under a particular set of assumptions:

1. Sampling access to the environment
2. Access to reward
3. Goal-directed behavior

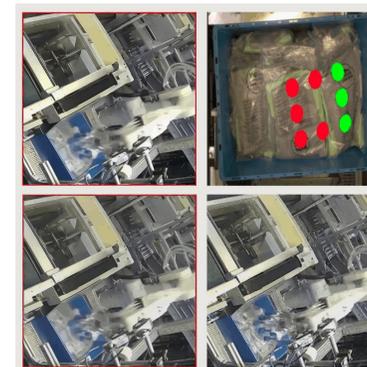
## Trajectory optimization/planning



## Imitation Learning



## Unsupervised Decision Making



# Trajectory Optimization

Sequential decision making with "known" models

*Interleaving Graph Search and Trajectory Optimization  
for Aggressive Quadrotor Flight*

*Ramkumar Natarajan, Howie Choset and Maxim Likhachev*

We combine RRT and local smoothing of contact dynamics to generate complex contact-rich manipulation plans.

May be hard to construct perfect, known models

# Imitation Learning

Sequential decision making provided expert data



Cook Shrimp  
(autonomous)

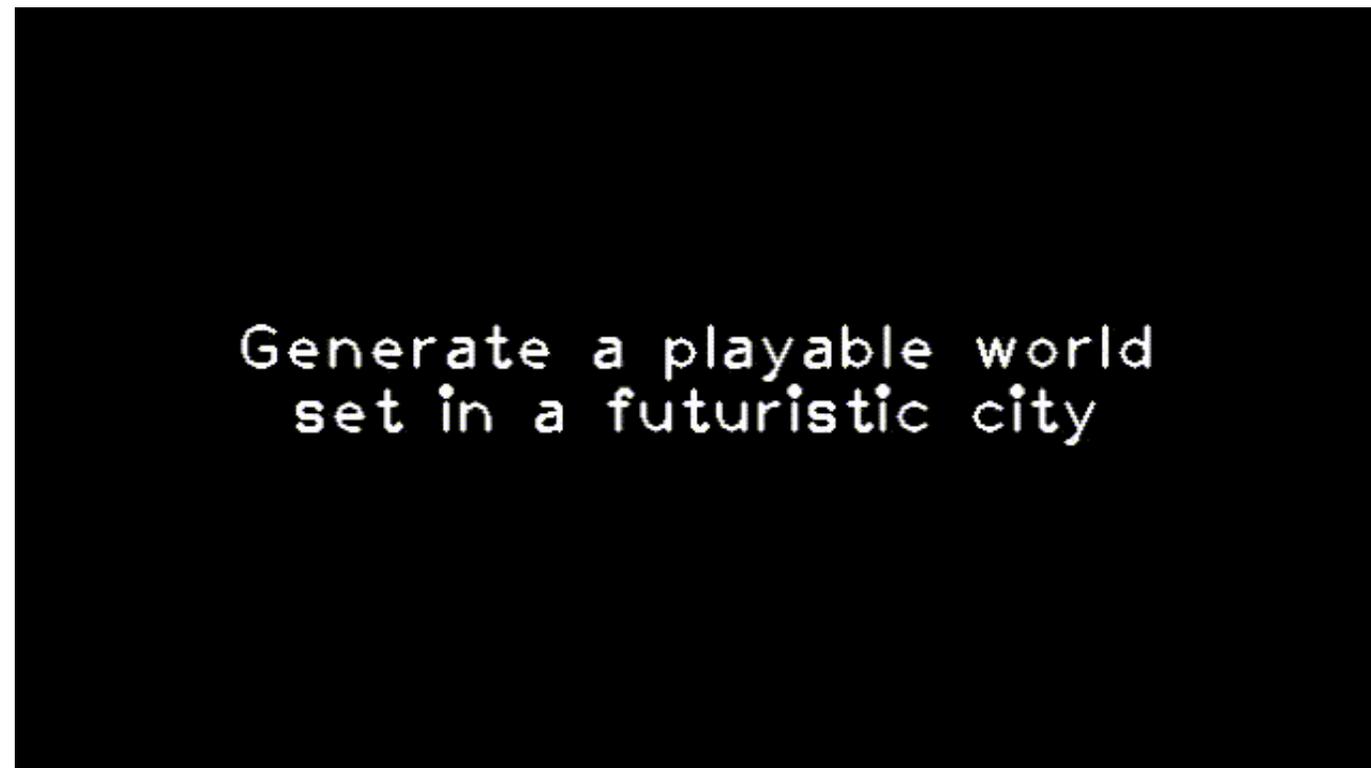
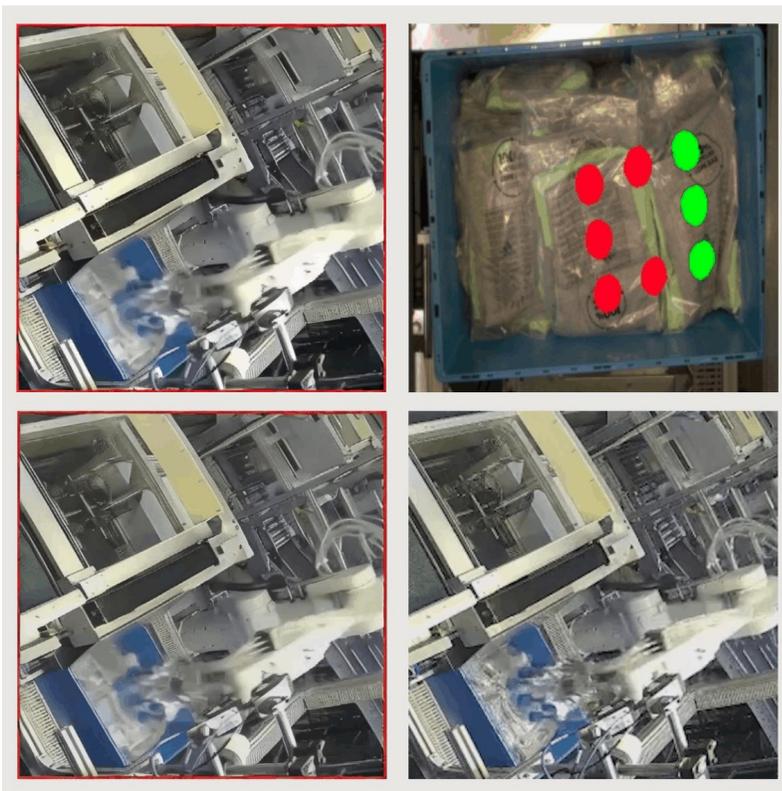


Pick & Place

Often called learning from demonstrations

# Self-Supervised Prediction of the World

Sequential decision making without reward – self-supervised prediction



Often called model-based RL

# How should we think about designing effective RL algorithms?



Easy to specify  
**objectives**

Stable performant  
**optimization** algorithms

Efficient **data** collection

# Class Structure

