# Reinforcement Learning
## CS 59300: RL1

August 26, 2025
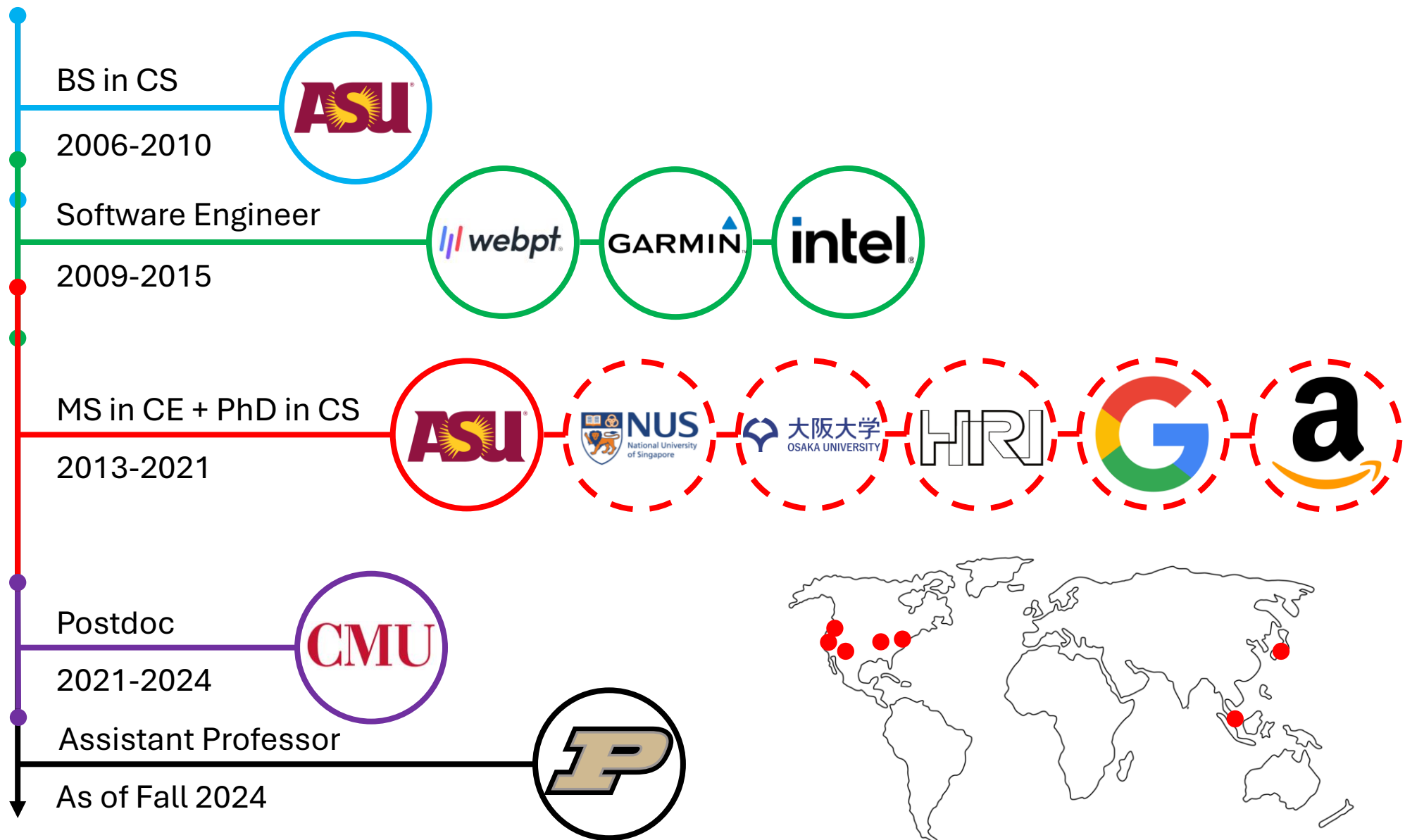
Joseph Campbell

Department of Computer Science

# Today's lecture

**1.** A little about me

**2.** What is reinforcement learning?

**3.** Course logistics

**4.** Intro to reinforcement learning

*Some content inspired by Katerina Fragkiadaki's CMU 10-403, Sergey Levine's Berkeley CS285, and Cathy Wu's MIT 6.7950 courses*

# A little about me

# **C**ollaborative **AI** for **M**achines and **P**eople
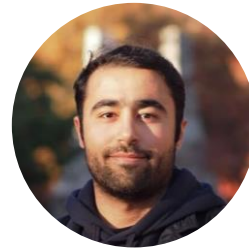
**CAMP Lab**

**PURDUE UNIVERSITY**

## PhD Students

Guven Gergerli | Abel Gurung | Muhan Lin | Shahab Rahimirad | Fiona Xie
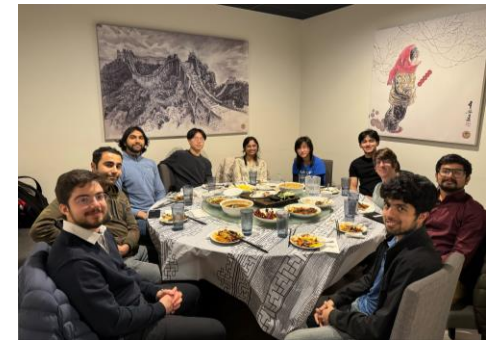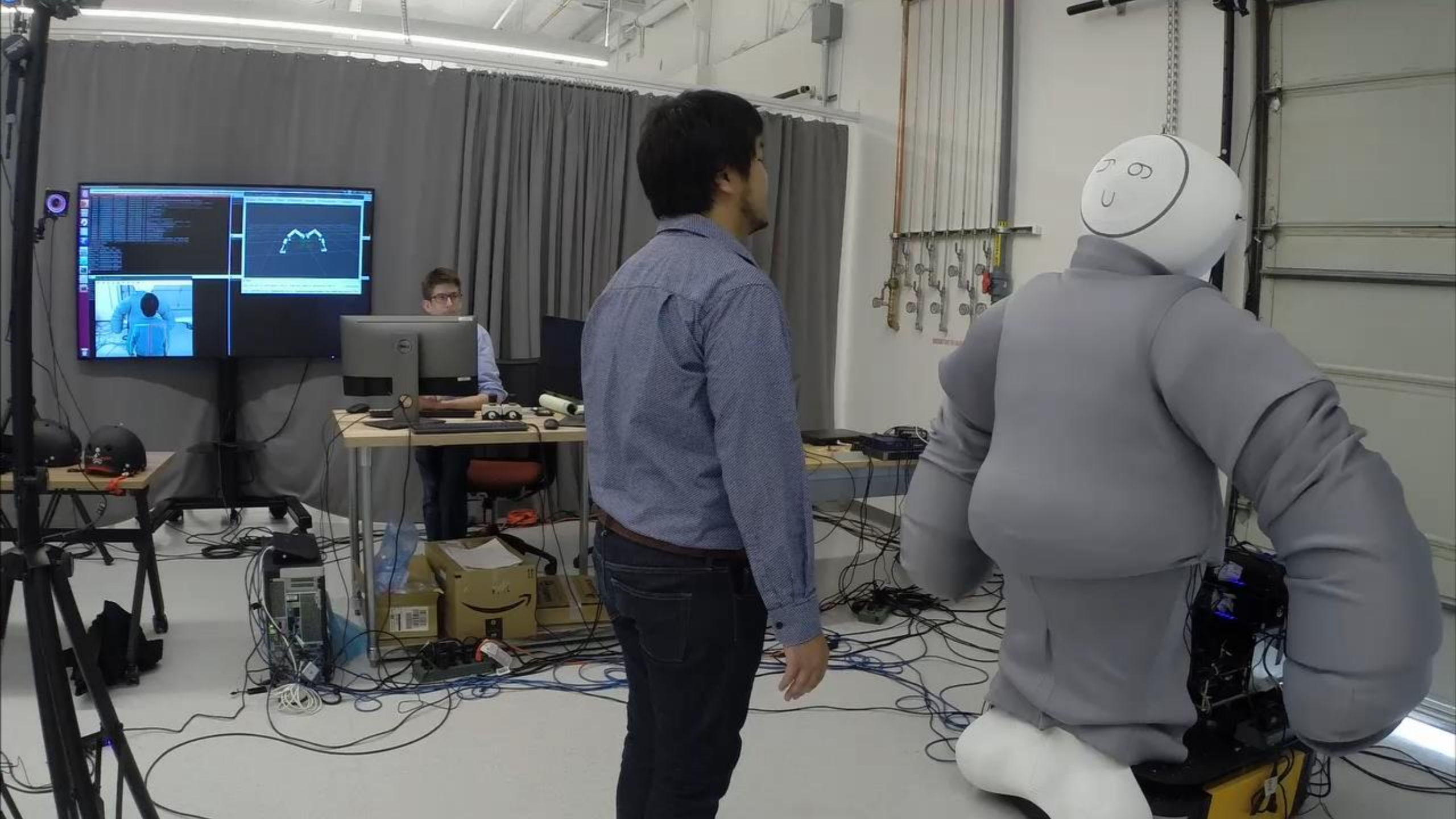
## MS Students

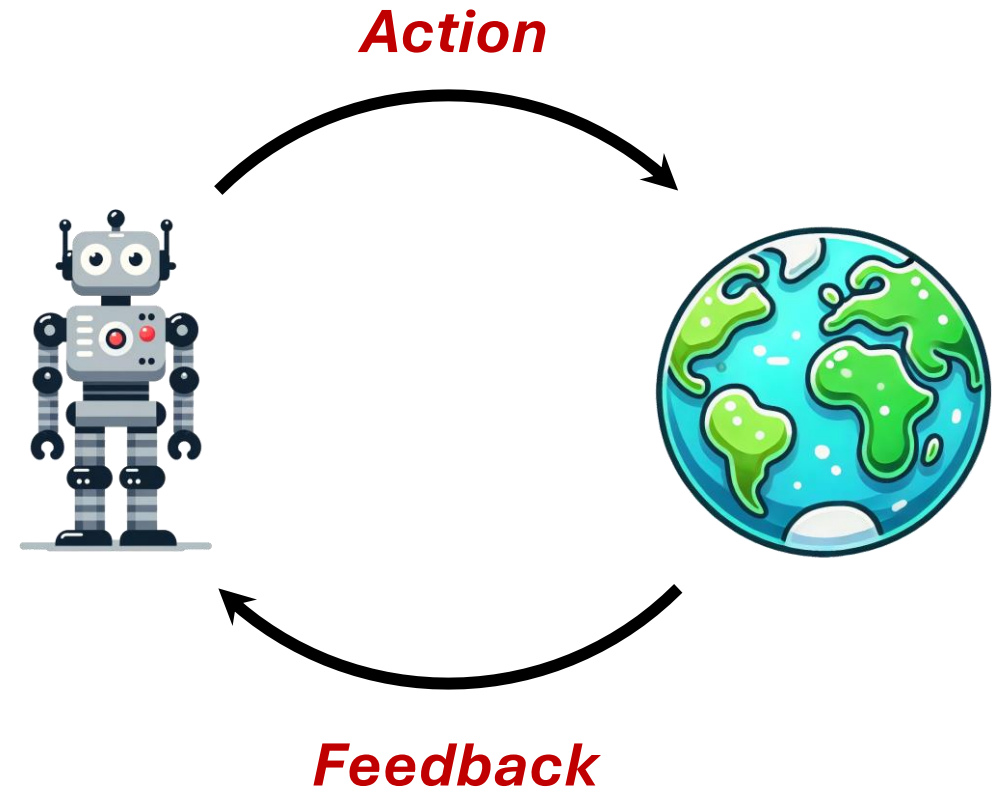Maheep Brar | Zeyun Deng | Luke Luschwitz | Le Mao | Anusha Sarraf

# What is reinforcement learning?

2:11

# Reinforcement learning ≈ trial-and-error!

**Action**

1. Take some action.

2. Get feedback (a reward).

**Feedback**

3. Use feedback to adjust what action we will take next time.

# Instrumental conditioning

**Law of effect**

"...responses that produce a satisfying effect in a particular situation become more likely to occur again in that situation, and responses that produce a discomforting effect become less likely to occur again in that situation."
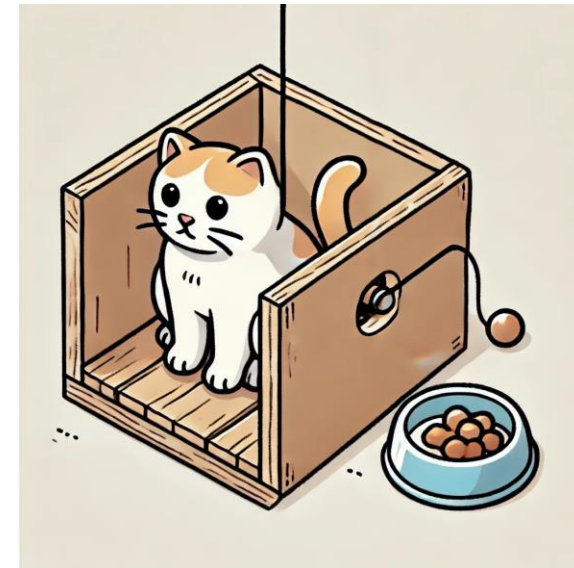
Edward Thorndike
Psychologist @ Columbia

"The cat that is clawing all over the box in her impulsive struggle will probably claw the string or loop or button so as to open the door.

And gradually all the other <span style="color:red">non-successful impulses will be stamped out</span> and the particular impulse leading to the <span style="color:red">successful act will be stamped in</span> by the resulting pleasure, until, after many trials, the cat will, when put in the box, immediately claw the button or loop in a definite way."

Edward Thorndike's famous puzzle box.



Terrible illustration courtesy of ChatGPT.

# Operant conditioning

**Reinforcement schedule**

"...any procedure that delivers reinforcement to an organism according to some well-defined rule."

Record effect of schedule on animal's learning rate.
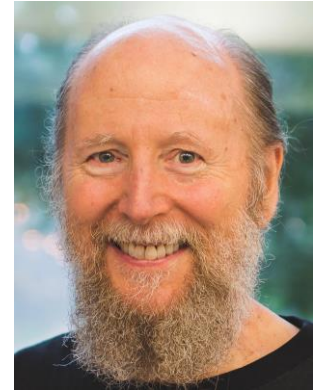


B.F. Skinner
Psychologist @ Harvard

https://youtu.be/yhvaSEJtOV8?si=9ojzWQ9txtAhMkgC

# Goal: take actions that maximize total reward

"Reinforcement learning is learning what to do—how to map situations to actions—so as to maximize a numerical reward signal.

The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them.

In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, all subsequent rewards."



Richard Sutton
University of Alberta



Andrew Barto
UMass Amherst

# Types of feedback: from environment

# Types of feedback: from teacher



EXERCISE TWO:
Remove elastic band and hit forehands
with your elbow close to your body

# Types of feedback: from ourselves

# Frequency of feedback

**Reward shaping**



"…decided to reinforce any response that had the slightest resemblance to a swipe—perhaps, at first, merely the behavior of looking at the ball—and then to select responses which more closely approximated the final form.
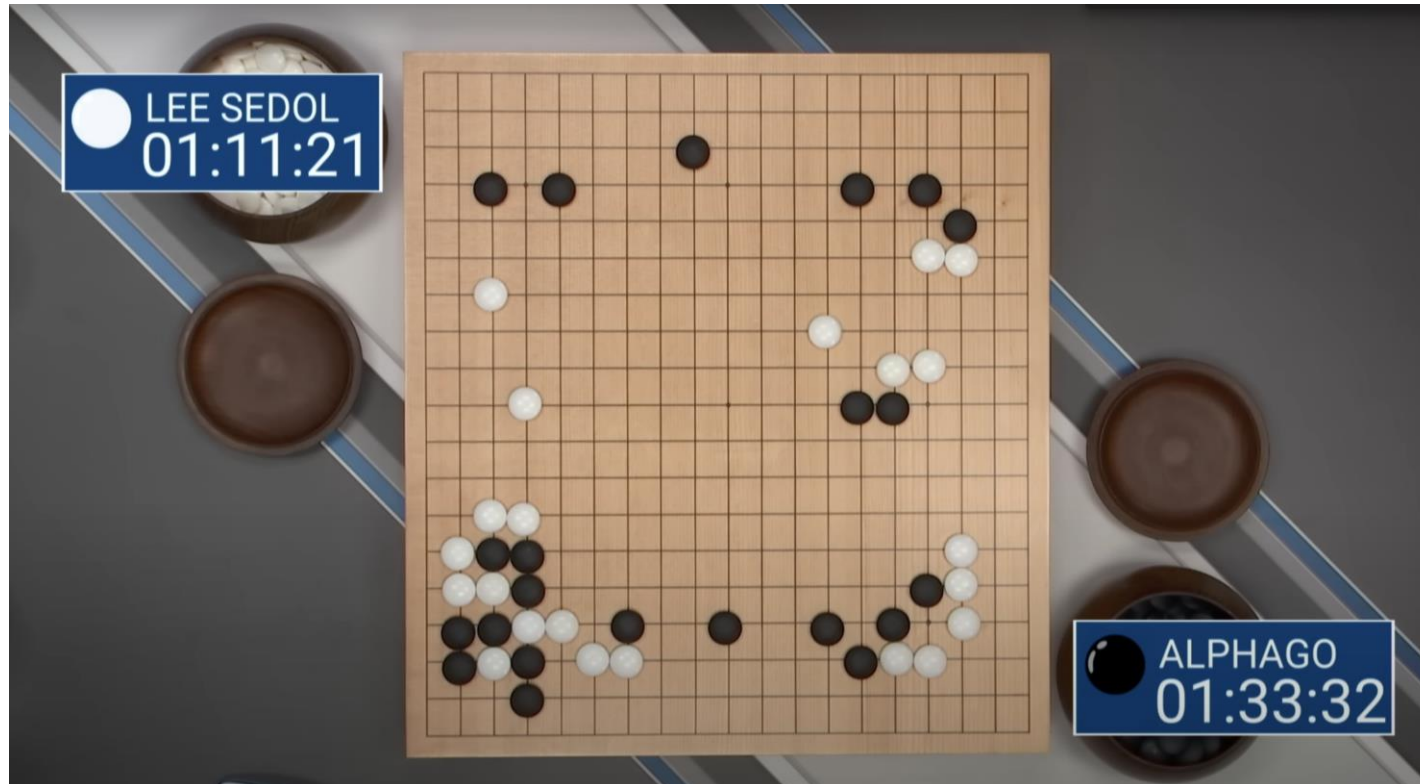
The result amazed us. In a few minutes, the ball was caroming off the walls of the box as if the pigeon had been a champion squash player."

B.F. Skinner
Psychologist @ Harvard

# Why do we care?

We can use this paradigm to build agents that **learn** to **act**!

**2002**

"No simple yet reasonable evaluation function will ever be found for Go." - Martin Müller

**2016**

20

# Robotics



D'Ambrosio et al. Achieving Human Level Competitive Robot Table Tennis. 2024.
https://youtu.be/EqQl-JQxToE?si=sRpkHHcF2UPJF5lC

# Recommendation engines

What artwork is shown to a user in order to maximize watch-rate?

# Image generation

How can we generate human-like images?

Black et al. Training Diffusion Models with Reinforcement Learning. 2024.

# Large language models

How do we make
LLMs generate
human-sounding text?

# In this course, you will learn...

How to build agents that learn to act!

- How reinforcement learning works
- What types of algorithms exist
- When they should be used
- How we can use reinforcement learning to solve real problems

# Course logistics

# What, where, when

Lecture: TTH 10:30 – 11:45 AM
Room: LWSN 1106

Email: joecamp@purdue.edu
Office hours: F 12:00 – 1:00 PM, DSAI 3047

TA email: ggergerl@purdue.edu
TA office hours: W 11:30 – 12:30 PM, DSAI B063

Please include [593RL] in the subject header

Will (try to) post all lectures on Brightspace before class and recordings after class



Guven Gergerli

# Course structure

Intended to provide a **broad overview** of reinforcement learning

- Focus on important concepts and algorithms

- Build theoretical + practical understanding

Mixture of lectures, homework assignments, and a course project

Later in the semester: a couple of guest lectures

# Prerequisites

Introductory-level knowledge of **machine learning** and **linear algebra**

- Do you know what supervised and unsupervised learning is?

- Are you comfortable working with probability distributions?

- Do you understand matrix algebra and other operations (inverse)?

Basic proficiency in **Python**

- *Vast majority* of machine learning libraries are written in Python

- Knowledge of Pytorch will help you greatly

## Schedule (Tentative)

| Week | Tuesday | | Thursday | |
|------|---------|---|----------|---|
| 1 | 8/26 | Introduction to RL | 8/28 | Deep Learning Basics and Behavior Cloning |
| 2 | 9/2 | Multi-Armed Bandits & Markov Decision Proc. | 9/4 | Policy/Value Iteration |
| 3 | 9/9 | Temporal Difference Learning | 9/11 | Deep Q Learning |
| 4 | 9/16 | Deep Q Learning | 9/18 | Policy Gradient |
| 5 | 9/23 | Actor-Critic Methods | 9/25 | Actor-Critic Methods |
| 6 | 9/30 | Model-Based RL | 10/2 | Model-Based RL |
| 7 | 10/7 | Guest Talk | 10/9 | RL Framework Tutorial |
| 8 | 10/14 | No Class Fall Break | 10/16 | Monte Carlo Tree Search |
| 9 | 10/21 | Multi-Agent RL | 10/23 | Multi-Agent RL |
| 10 | 10/28 | Guest Talk | 10/30 | Inverse Reinforcement Learning |
| 11 | 11/4 | Offline Reinforcement Learning | 11/6 | Intelligent Exploration and Curiosity |
| 12 | 11/11 | Intelligent Exploration and Curiosity | 11/13 | Transfer Learning and Lifelong Learning |
| 13 | 11/18 | Good Presentation Tips | 11/20 | Large Language Models |
| 14 | 11/25 | Vision Language Models | 11/27 | No Class Thanksgiving Break |
| 15 | 12/2 | Robotics & Sim-to-Real | 12/4 | Challenges and Open Problems |
| 16 | 12/9 | Poster Presentations | 12/11 | Poster Presentations |

# Grading

Assignments (individual): 40%
- 4x assignments at 10% each

Midterm exam: 15%

Course project (group): 40%
- Midterm update: 10%
- Project report (plus deliverables): 20%
- Poster presentation: 10%

Class participation: 5%

# Discussion is highly encouraged!

Ask questions at any time

If you have a question...odds are others do as well

If I ask questions (which I will), please participate

- This is 5% of your grade!



ONLY Y🞛U CAN ~~PREVENT WILDFIRES~~

*make this class awesome!*

# Course project

Goal: gain experience implementing, applying, evaluating, and improving reinforcement learning algorithms

Form groups of 2-4 students

- Tell me your group by **9/5** (end of week 2)

Project proposal must be approved by me

- Proposal due **9/26** (end of week 5)

# Course project

Midterm-update due by **10/30** (week 10)

- 2-page research-style paper outlining your project and progress
- IEEE paper format. Recommend using Overleaf

Final project due by **12/11** (week 16)

- ≥ 6-page research-style paper outlining your project and findings
- Poster presentation in class on 12/9 and 12/11

Goal is to replicate the lifecycle of a research project

# Example project ideas

**Implementation-based**: implement a deep reinforcement learning algorithm **from scratch** and apply it to a challenging problem, *e.g. robot navigation or manipulation*

- Understand and show me how various hyperparameters and other design choices effect algorithm performance

- If you are interested in robotics or another real-world application, let me know and we can try to find resources to make it work

- Must be a different algorithm/task than in the homeworks!

# Example project ideas

**Analysis-based**: implement multiple algorithms (may use libraries) and apply them to multiple tasks and perform meaningful empirical or theoretical analysis, e.g.

- Try to reproduce a paper's experimental findings

- Compare different variants of curiosity on a benchmark

- Analyze why one algorithm may out-perform another, e.g. when and why does SAC out-perform PPO?

# Example project ideas

**Improvement-based**: identify a weakness with an existing algorithm over a particular benchmark/evaluation metric and introduce a variant of this algorithm which out-performs the original

- This is a classic "research problem" approach

- If you don't make meaningful improvements, that's ok! But show me what you tried and why it did/didn't work

# Example project ideas

I will work with you to develop your project ideas!

- If you have an idea but aren't sure if it is viable, come talk to me at office hours or send an email; we can figure something out

If you are interested in research, we can work towards turning your project into a research publication after the semester

# Collaboration policy and academic integrity

While discussion and exchanging ideas is encouraged, students and groups must perform their own work separately

- For group projects you may not collaborate outside your group

- For homework assignments you may not collaborate with others

All work must be your own unless otherwise stated

- You may use AI tools (e.g. ChatGPT) as an *assistive tool*

- All uses must be explicitly documented in any submitted work

Refer to Purdue's Academic Integrity website

# Intro to reinforcement learning

# The basics

The agent and environment operate at discrete timesteps $t = 0,1,2,\dots$

- The agent observes state $s_t$ at time $t$

- The agent takes action $a_t$

- The agent gets the resulting reward $r_{t+1}$ and the subsequent state $s_{t+1}$

$$a_t$$

$$s_{t+1}, r_{t+1}$$

# The basics

Action $a_t$ is chosen by sampling actions from a probability distribution

$$a_t \sim \pi(a|s)$$

The probability distribution $\pi$ is referred to as a policy.



$a_t$

$s_{t+1}, r_{t+1}$

# Supervised learning vs reinforcement learning

**Supervised/imitation learning**

- Agent learns policy that imitates a set of demonstrations

- Learns quickly
- Performance is bounded by the quality of demonstrations
- Agent doesn't know what to do if it makes a mistake

**Reinforcement learning**

- Agent learns its own policy that maximizes rewards

- Learns slowly
- Can discover novel solutions
- Requires an environment to interact with

## Neurogammon – Supervised

## TD-Gammon – Reinforcement

- Developed by Gerald Tesauro in 1989 at IBM Research

- Trained using supervised learning given expert demonstrations

- Intermediate-level play

- Developed by Gerald Tesauro in 1992 at IBM Research

- Trained against itself using self-play beginning with random actions

- Top-level play

Neurogammon – Supervised



TD-Gammon – Reinforcement

- Deve...                                                    ...sauro
  in 19...
- Train...                                                    ...ing
  learn...
  dem...

> **Why did TD-Gammon play so much better?**
>
> **Because it discovered novel strategies!**

- Intermediate-level play
- Top-level play

49:25

# What is deep reinforcement learning?

How can an agent learn over a complex state space, e.g. images?

Policy $\pi(a|s)$ is represented using a deep neural network.

# Limitations of reinforcement learning

The agent needs to interact with an environment

Requires either a simulator or access to the real world

- If simulator:
  - Is it high-enough fidelity to use the policy in the real-world after training?

- If real world:
  - Takes a long time to collect samples + wear and tear (if embodied)
  - May not be safe (can we train autonomous vehicles on city roads?)

# Brute force solution: more agents!



Ahn et al. AutoRT: Embodied Foundation Models for Large Scale Orchestration of Robotic Agents. 2024.

# Slightly more elegant: fancy simulator!



*NVIDIA*

# Limitations of reinforcement learning

The agent's policy is <span style="color:red">non-stationary</span>

- Actions depend on agent's policy (bad policy = bad trial-and-error)

What if the agent needs to start from a known "good" state?

- Environment may need to be reset somehow
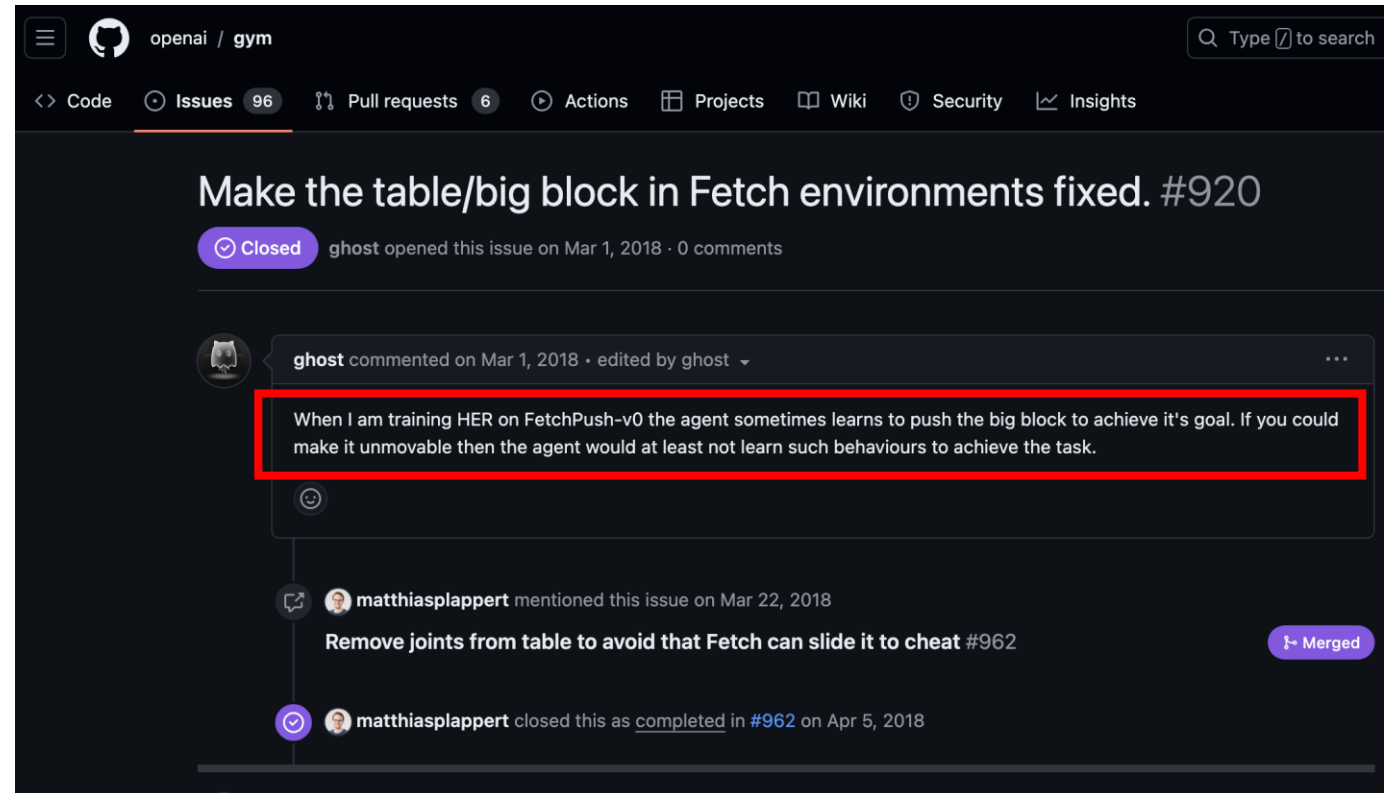
# Limitations of reinforcement learning

Reward functions may be difficult to specify

- How do you specify social/cultural norms in a reward function?
  - Required for chat bots, social navigation, autonomous driving, …

- How frequently does an agent get a reward?
  - Reward may be difficult to obtain through random exploration
  - Reward may be the result of a long series of actions
  - Hint: remember "reward scheduling" and "reward shaping"?

# Reward hacking

Rewards are hard for people to design!

Agents may find "loopholes" to exploit reward process



Skalse et al. Defining and Characterizing Reward Hacking. 2022.

https://openai.com/index/faulty-reward-functions/

https://youtu.be/tlOIHko8ySg?si=bMZ5lpjTqFaoo44j

# Take-aways

- In reinforcement learning an agent discovers its own actions which maximize its received reward

- Differs from supervised learning which imitates demonstrations

- Reinforcement learning has many challenging problems