

# Reinforcement Learning



# Motivation

Create Myanmar Reinforcement Learning Research Community

# Great Expectations

## WILL YOU BE

an expert at RL after the class - NO

able to train world champion chess AI - NO

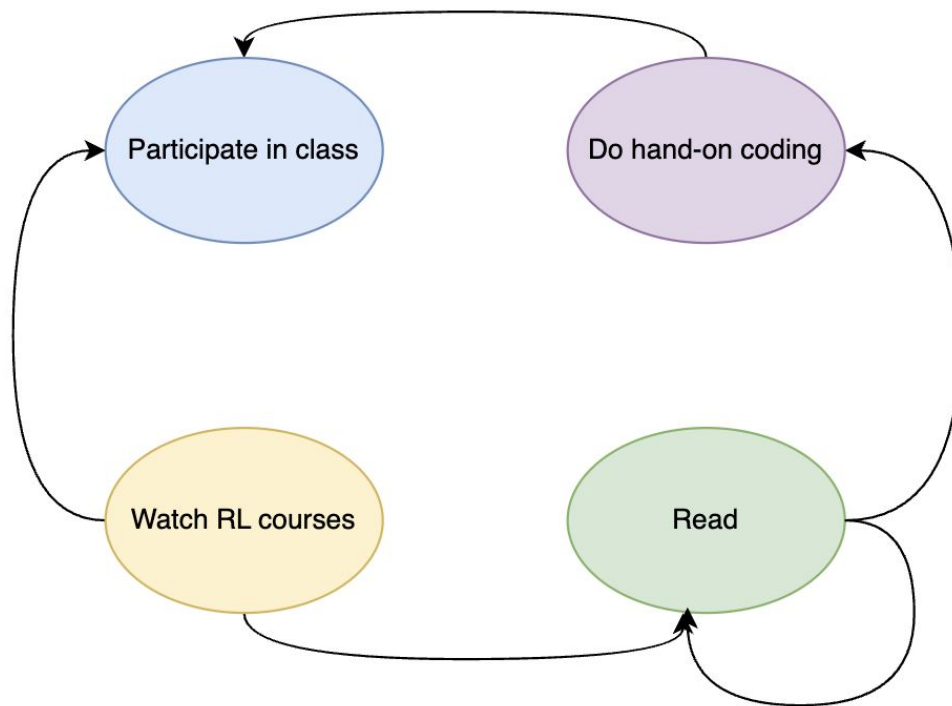
able to implement RL algorithms on your own - MAYBE

able to understand what RL is - YES

# Student Profile

	Student 1	Student 2	Student 3
Machine Learning	Well-versed, have done deep learning problems	Beginner, have done toy problems	Very Basic
Reinforcement Learning (RL)	Know definitions and fundamental concept	Know the terms, but No Idea Yet	No Idea Yet
Motivation	Pursue RL	Understand RL	Get to know RL

# How to Make the Most out of this class



# Types of Machine Learning

Supervised Learning

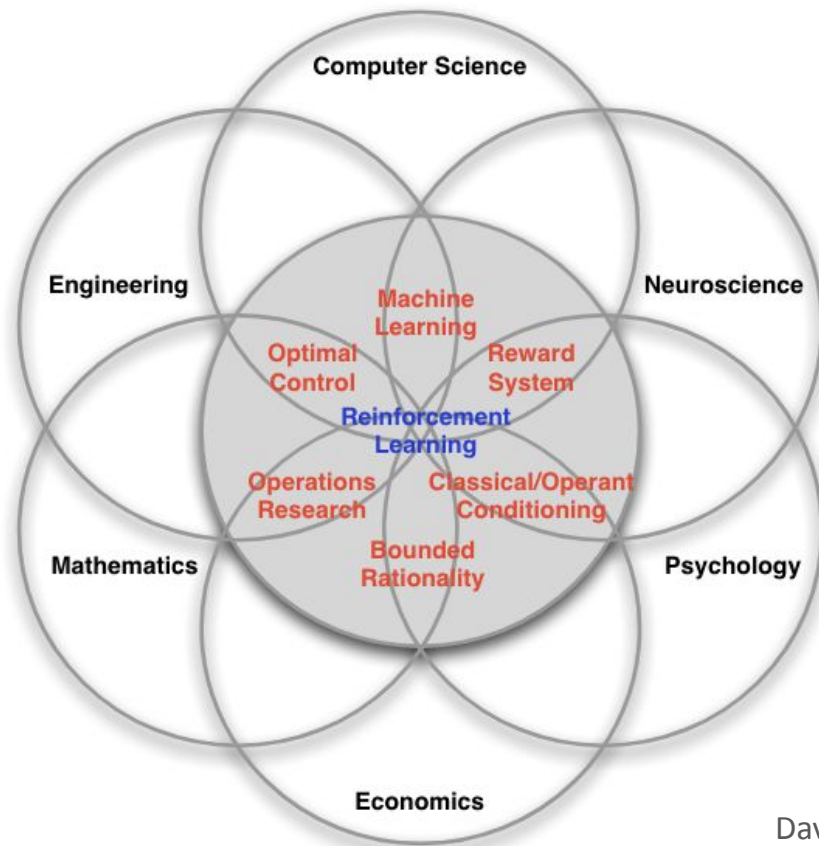
Unsupervised Learning

Reinforcement Learning

# What is Reinforcement Learning (RL)?

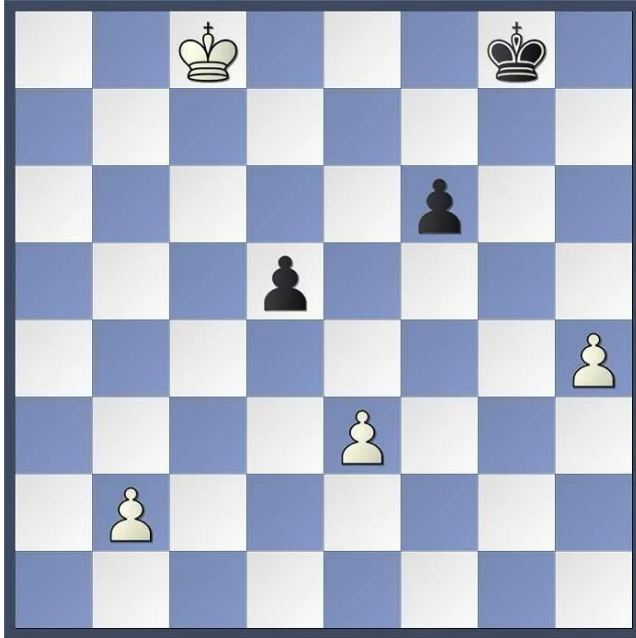


# Many Faces of Reinforcement Learning



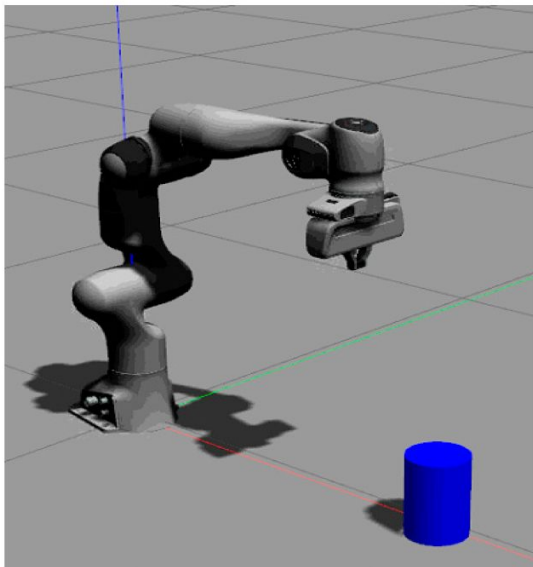


# Board games

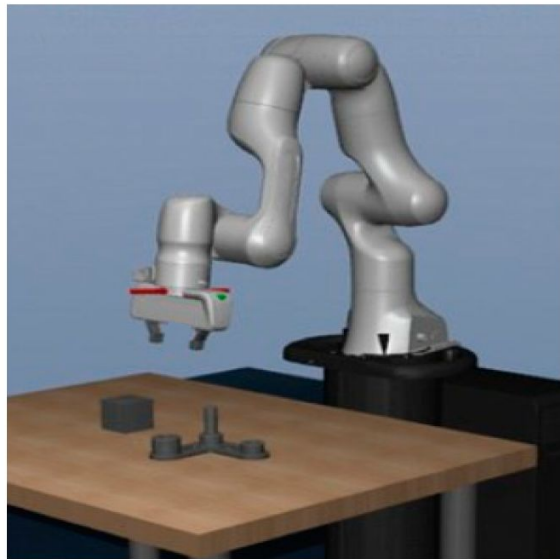


- Chess
- Go
- Connect four
- Basically, all board games

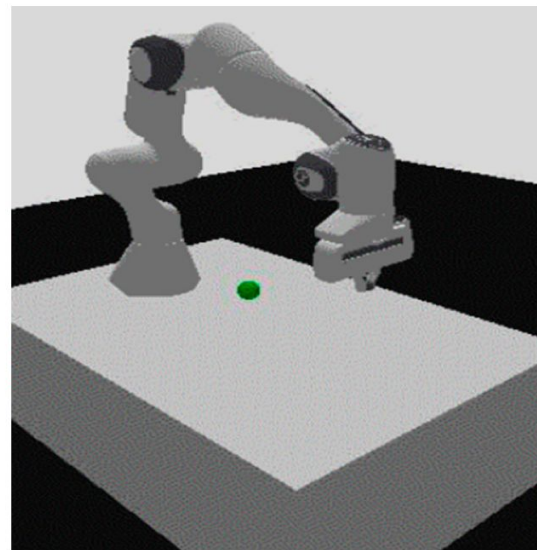
# Robotics



(a)



(b)



(c)

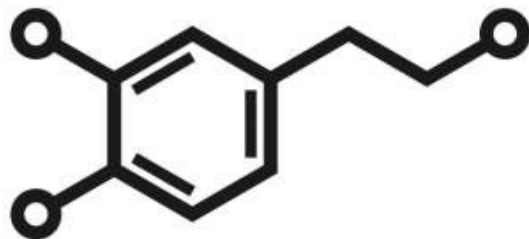
# Other Examples

Recommender systems

Stock trading

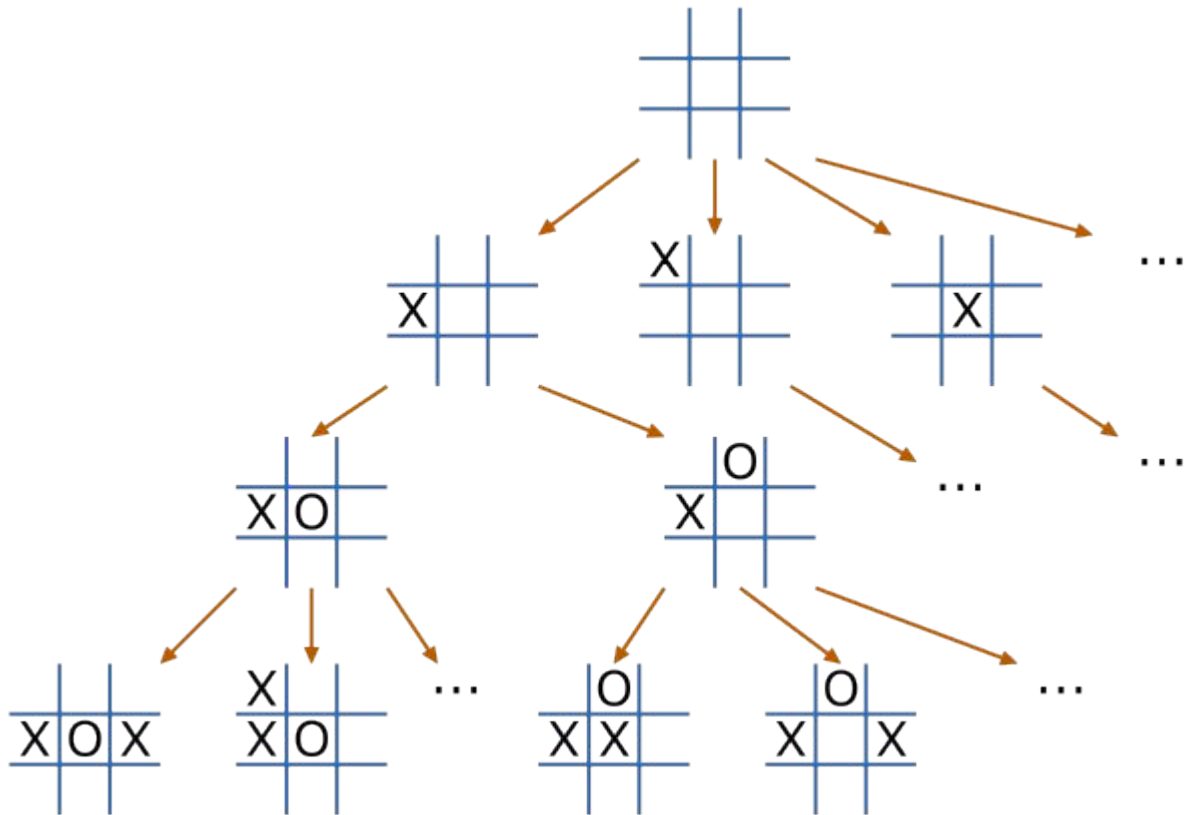
Language model alignment

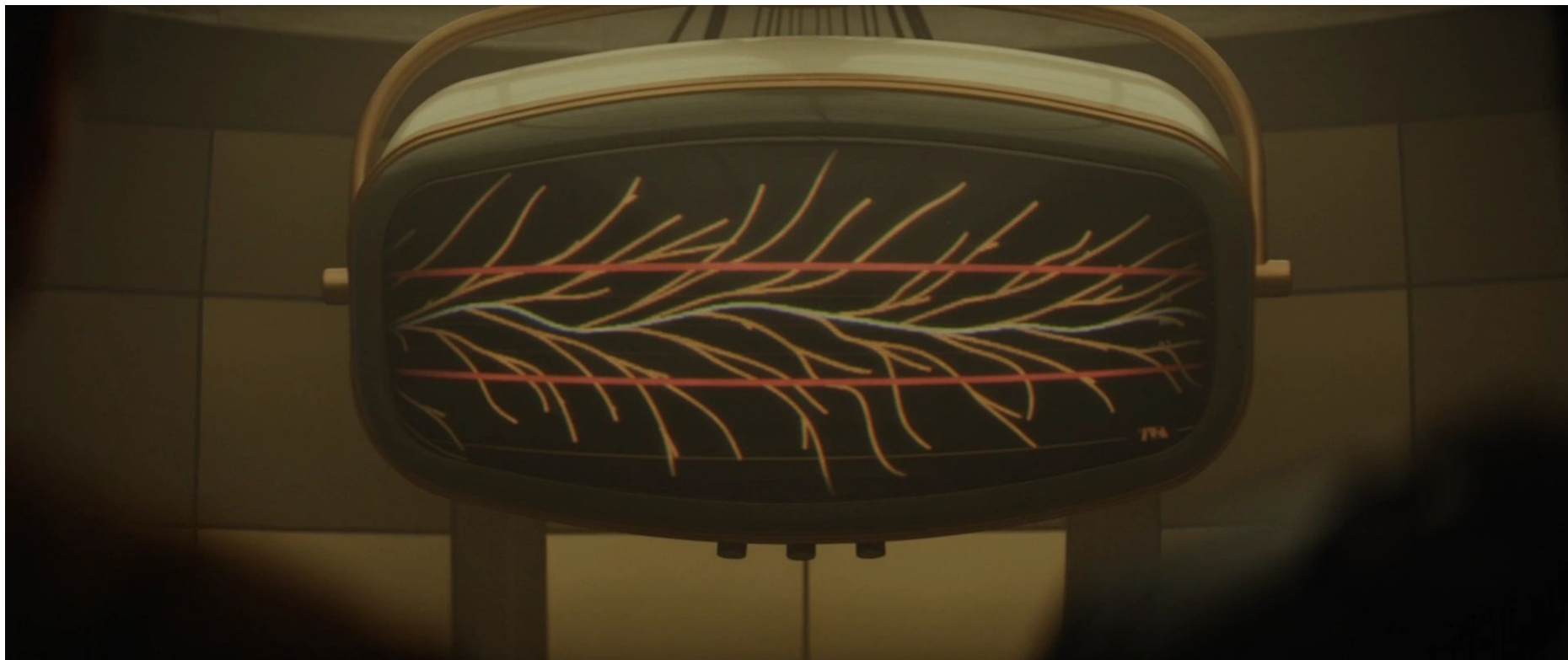
Reward



DOPAMINE

Do we always know the reward?

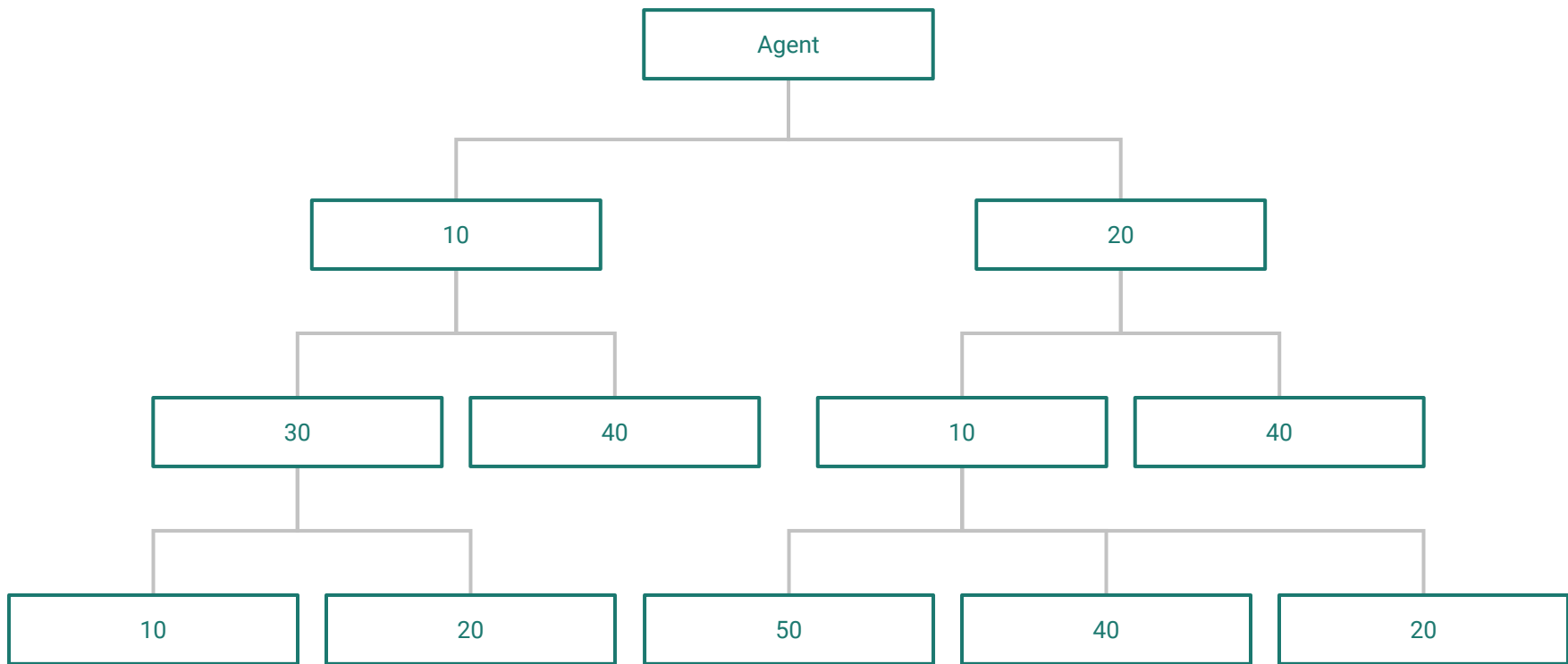




If we know it, it's called supervised learning. Or it becomes a search problem.



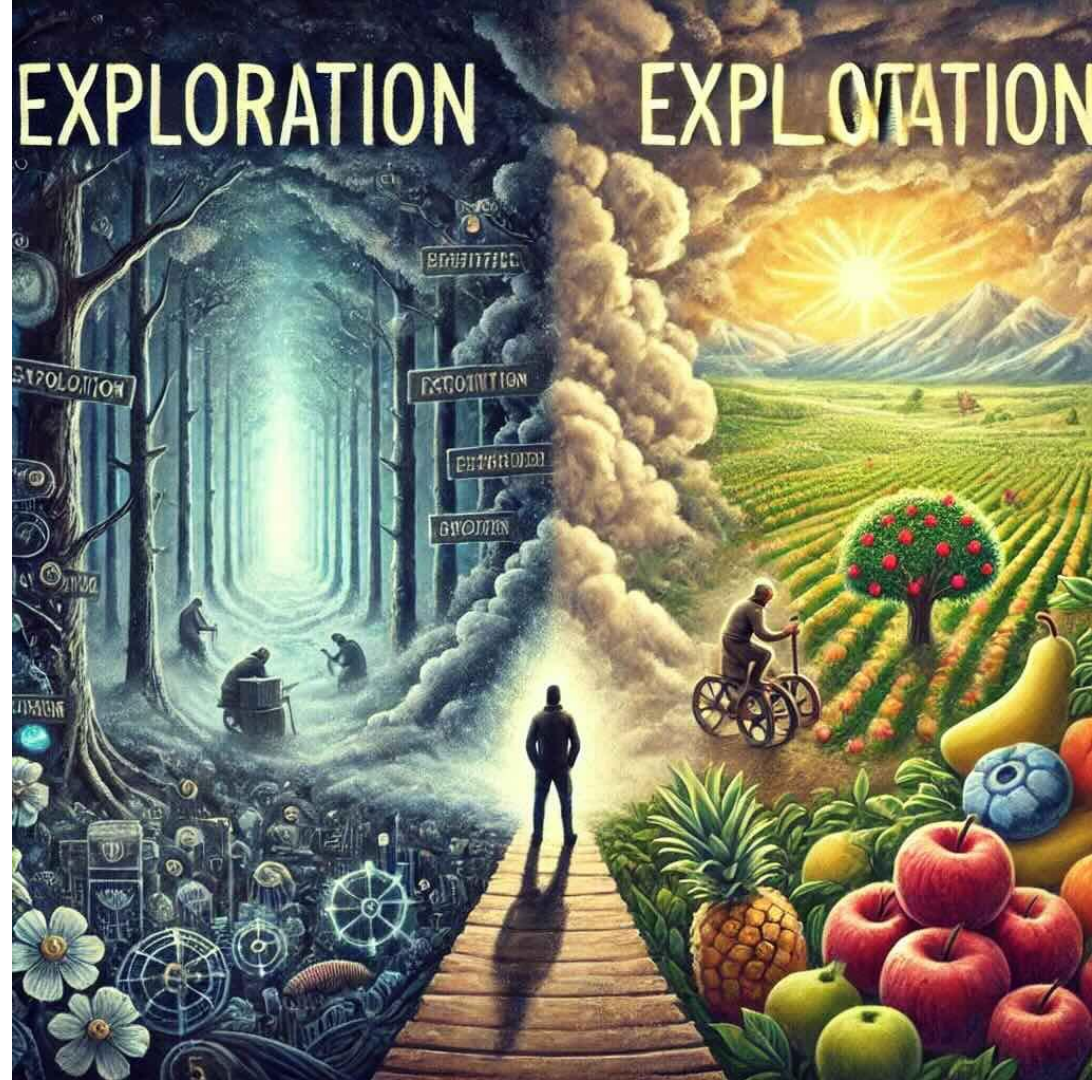
# Chess Engines vs RL



# Exploration

$$V_S$$

# Exploitation



# Life Lesson

- Explore
- But how?
- How do we know the future?
- How can we learn to explore?

# Credit Assignment Problem

- Why do we win?
- Why do we lose?
- What went well?
- What went wrong?

MATHEMATICS

# Expectation

$$\mathbf{E}[X] = \sum_{i=1}^{\infty} x_i p_i$$

[https://en.wikipedia.org/wiki/Expected\\_value](https://en.wikipedia.org/wiki/Expected_value)

# Expected value of a Fair Die

$$p(1) = \frac{1}{6}$$

$$p(2) = \frac{1}{6}$$

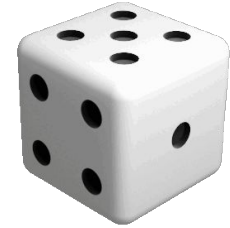
$$p(3) = \frac{1}{6}$$

$$p(4) = \frac{1}{6}$$

$$p(5) = \frac{1}{6}$$

$$p(6) = \frac{1}{6}$$

$$E(X) = \left(\frac{1}{6} * 1\right) + \left(\frac{1}{6} * 2\right) + \left(\frac{1}{6} * 3\right) + \left(\frac{1}{6} * 4\right) + \left(\frac{1}{6} * 5\right) + \left(\frac{1}{6} * 6\right)$$





# Expected value of unfair die

$$p(1) = 0.3$$

$$p(2) = 0.25$$

$$p(3) = 0.1$$

$$p(4) = 0.075$$

$$p(5) = 0.075$$

$$p(6) = ?$$

$$E(X) = ?$$

# Expected value of mystery die

Probabilities of rolls are not given?

How would you find the expected value?



# Law of Large Numbers

In **probability theory**, the law of large numbers (LLN) is a **mathematical theorem** that states that the **average** of the results obtained from a large number of independent random samples converges to the true value, if it exists

$$Q_t(a) = \frac{R_1 + R_2 + \cdots + R_{N_t(a)}}{N_t(a)}$$

Assume  $Q_t(a)$  is expected value of a dice

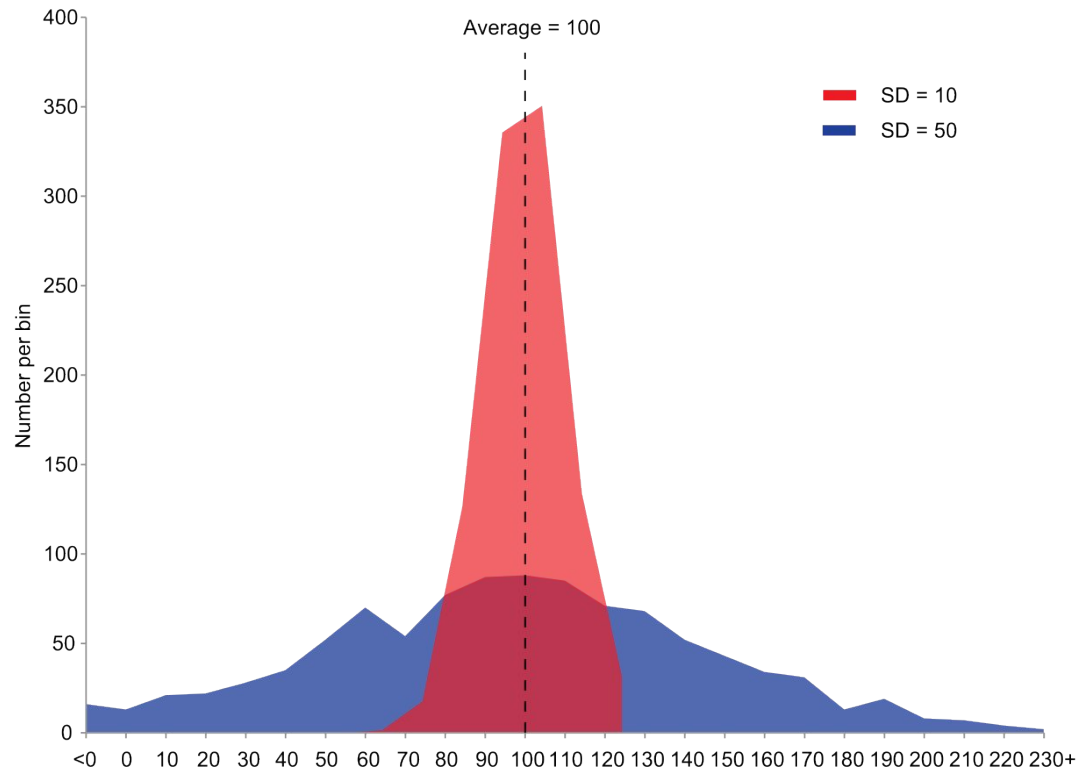
$R_1, R_2 \dots$  are values of dice rolls.

$N_t$  is number of dice rolls.

$$\begin{aligned}
Q_{n+1} &= \frac{1}{n} \sum_{i=1}^n R_i \\
&= \frac{1}{n} \left( R_n + \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + (n-1) \frac{1}{n-1} \sum_{i=1}^{n-1} R_i \right) \\
&= \frac{1}{n} \left( R_n + (n-1) Q_n \right) \\
&= \frac{1}{n} \left( R_n + n Q_n - Q_n \right) \\
&= Q_n + \frac{1}{n} \left[ R_n - Q_n \right],
\end{aligned}$$

$$\textit{NewEstimate} \leftarrow \textit{OldEstimate} + \textit{StepSize} \left[ \textit{Target} - \textit{OldEstimate} \right].$$

# Variance



$$\text{Var}(X) = \sum_{i=1}^n p_i \cdot (x_i - \mu)^2,$$

where  $\mu$  is the expected value. That is,

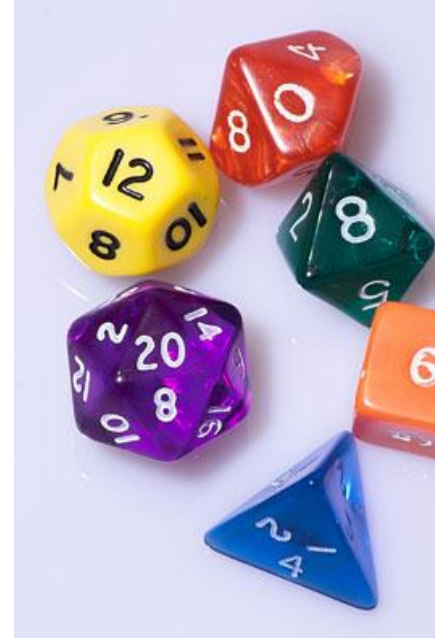
$$\mu = \sum_{i=1}^n p_i x_i.$$



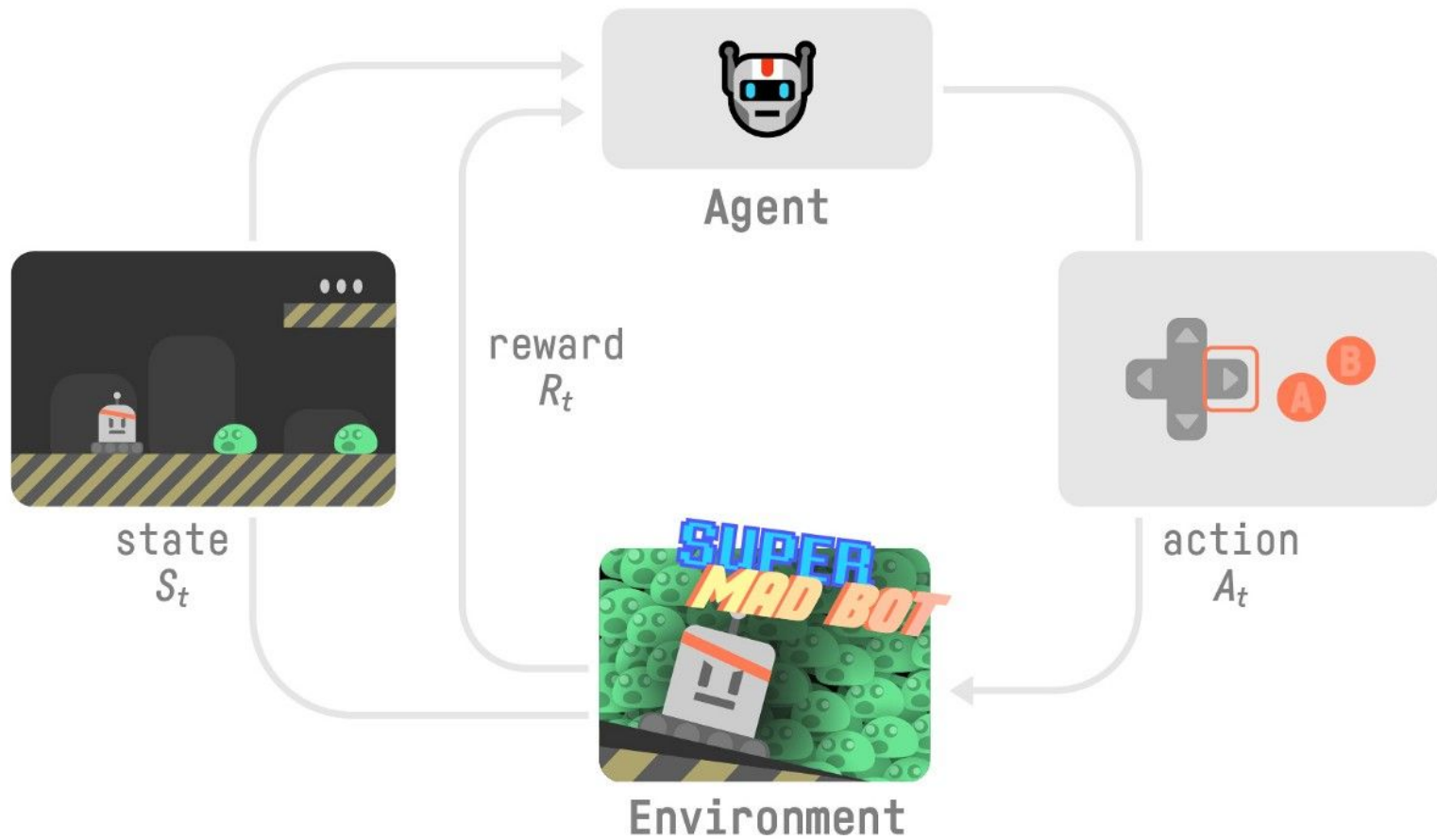
# You are given six dice.

Assume you cannot see the die faces. And the probability is not uniform.

Find the strategy to get maximum sum in 1000 rolls or fewer.



# BASIC TERMINOLOGIES



# Terminologies

**State (S)** - The current situation or configuration of the environment.

**Action (A)** - A decision or move that the agent can take from a given state.

**Policy ( $\pi$ )** - An agent's behavior, specifying which action to take in each state. It can be deterministic (always choosing the same action for a given state) or stochastic (random).

**Reward (R)** - A scalar feedback received by the agent after taking an action in a state and transitioning into a new state.

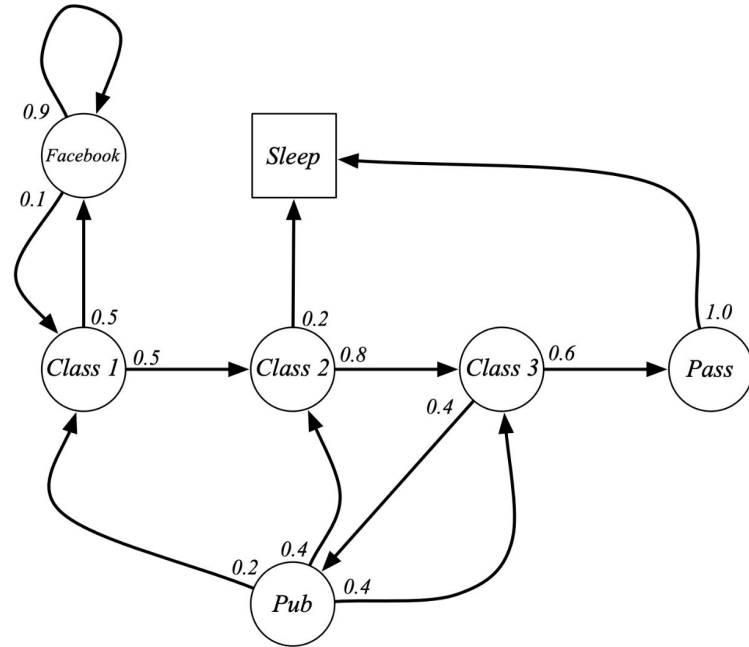
**Transition Probability (P)** - the probability of moving from current state  $s$  to new state  $s'$ .

# Reward hypothesis

*That all of what we mean by goals and purposes can be well thought of as maximization of the expected value of the cumulative sum of a received scalar signal (reward)*

# Markov Decision Process

The future state depends only on the current state, not on the sequence of events that preceded it.

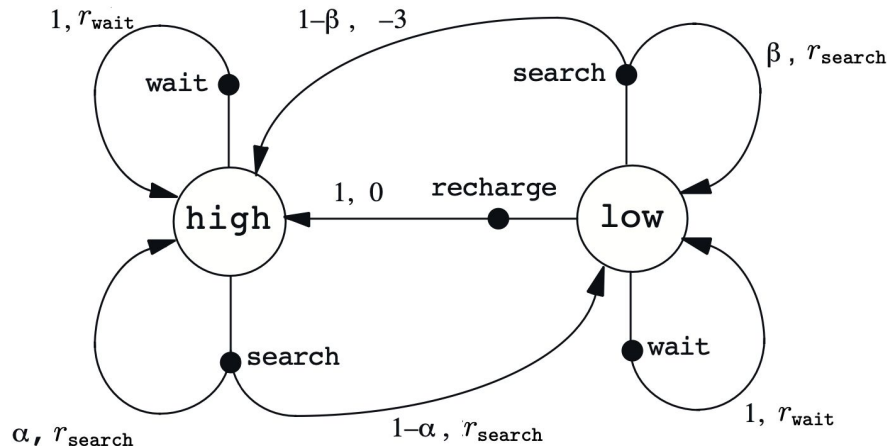


David Silver

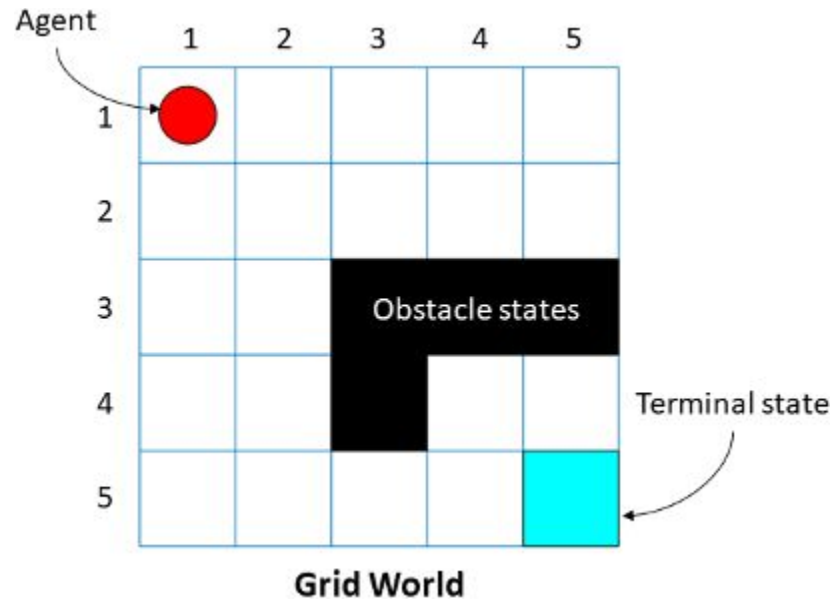
<https://www.davidsilver.uk/teaching/>

# Recycling Robot

$s$	$s'$	$a$	$p(s' s, a)$	$r(s, a, s')$
high	high	search	$\alpha$	$r_{\text{search}}$
high	low	search	$1 - \alpha$	$r_{\text{search}}$
low	high	search	$1 - \beta$	$-3$
low	low	search	$\beta$	$r_{\text{search}}$
high	high	wait	1	$r_{\text{wait}}$
high	low	wait	0	$r_{\text{wait}}$
low	high	wait	0	$r_{\text{wait}}$
low	low	wait	1	$r_{\text{wait}}$
low	high	recharge	1	0
low	low	recharge	0	0.



# Grid World





# Policy

Policy defines the learning agent's way of behaving at a given time

Denoted as  $\pi$

$$\pi(a|s) = p(a \mid s)$$

Probability of taking action  $a$  given state  $s$ .

## Returns

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \cdots + R_T,$$

Value function of state  $s$

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi}[G_t \mid S_t = s]$$