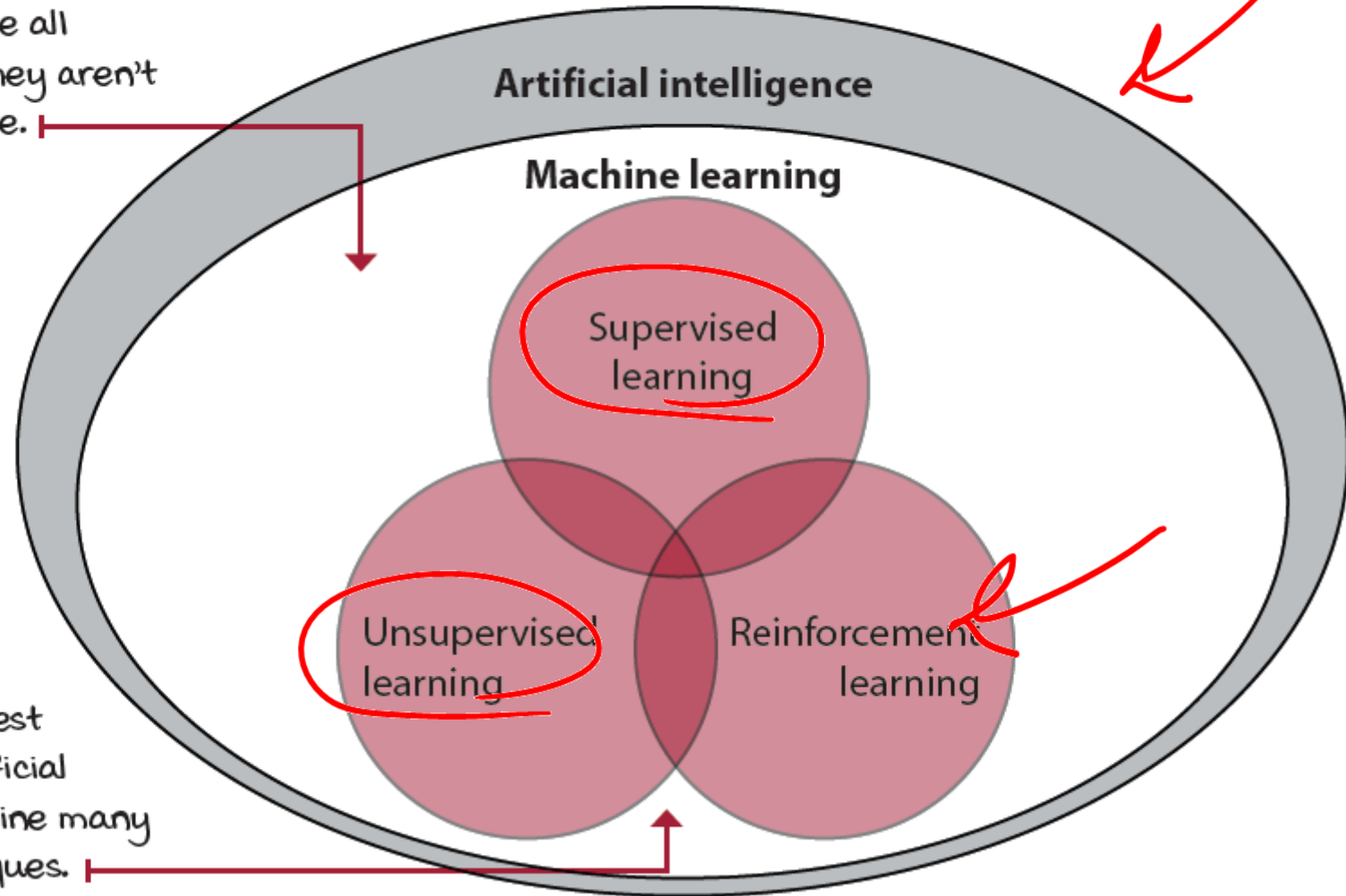# Reinforcement Learning Framework

# Main branches of machine learning

(i) These types of machine learning tasks are all important, and they aren't mutually exclusive.

**Artificial intelligence**

**Machine learning**

Supervised learning

Unsupervised learning

Reinforcement learning

(2) In fact, the best examples of artificial intelligence combine many different techniques.

# MAIN BRANCHES OF MACHINE LEARNING

$$y = f(x)$$

Supervised learning (SL) is the task of learning from labeled data. In SL, a human decides which data to collect and how to label it. The goal in SL is to generalize.

Unsupervised learning (UL) is the task of learning from unlabeled data. Even though data no longer needs labeling, the methods used by the computer to gather data still need to be designed by a human. The goal in UL is to compress.

$$f(x)$$

Reinforcement learning (RL) is the task of learning through trial and error. In this type of task, no human labels data, and no human collects or explicitly designs the collection of data. The goal in RL is to act.
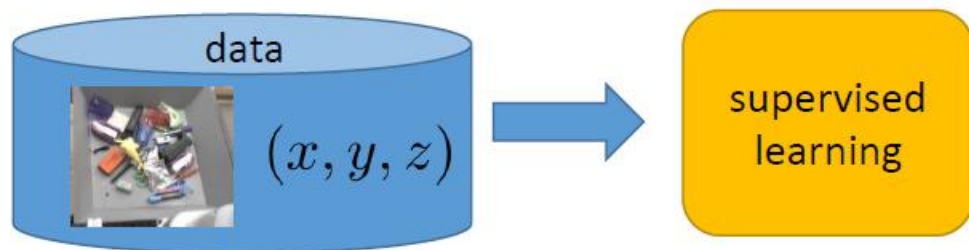
$$y = f(x)$$

$$z$$

## Standard (supervised) machine learning:

given $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$

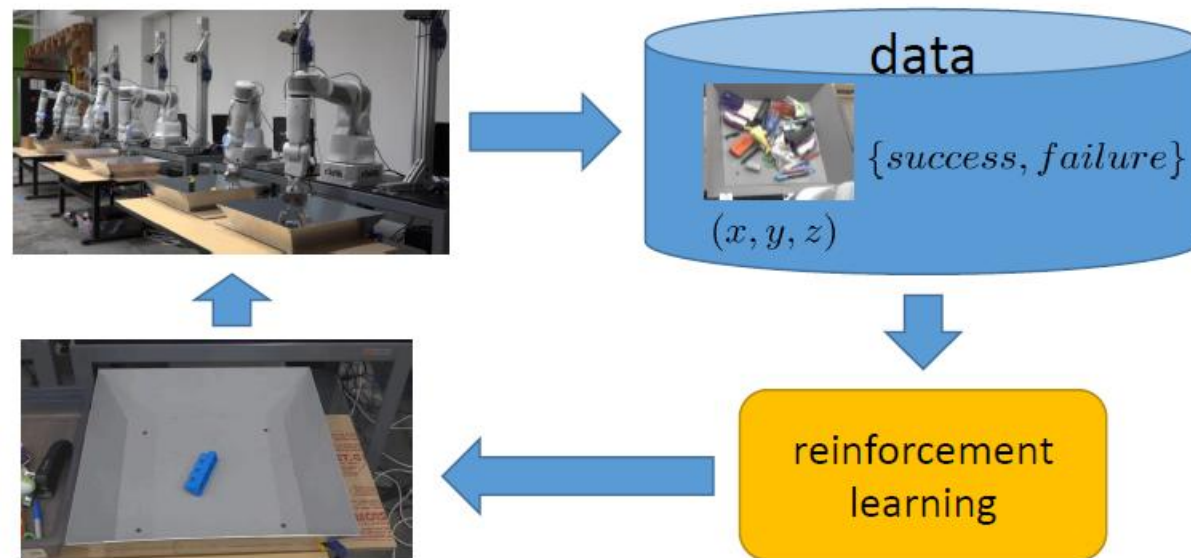learn to predict $y$ from $\mathbf{x}$ $\qquad f(\mathbf{x}) \approx y$

## Usually assumes:

- i.i.d. data
- known ground truth outputs in training



## Reinforcement learning:

- Data is **not** i.i.d.: previous outputs influence future inputs!
- Ground truth answer is not known, only know if we succeeded or failed
  - more generally, we know the reward

*"Almost all young people working on Artificial Intelligence look around and say - What's popular? Statistical learning. So, I'll do that. That's exactly the way to kill yourself scientifically!"*

*Marvin Minsky during his course called Society of Mind at MIT in 2011*

Nathaniel Rochester

Oliver Selfridge

Ray Solomonoff

Marvin Minsky

John McCarthy
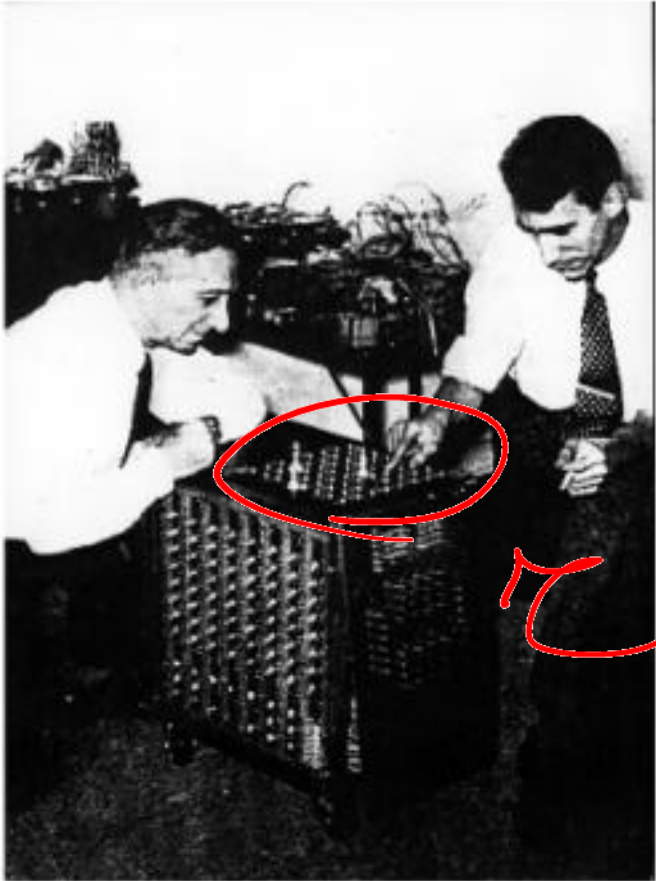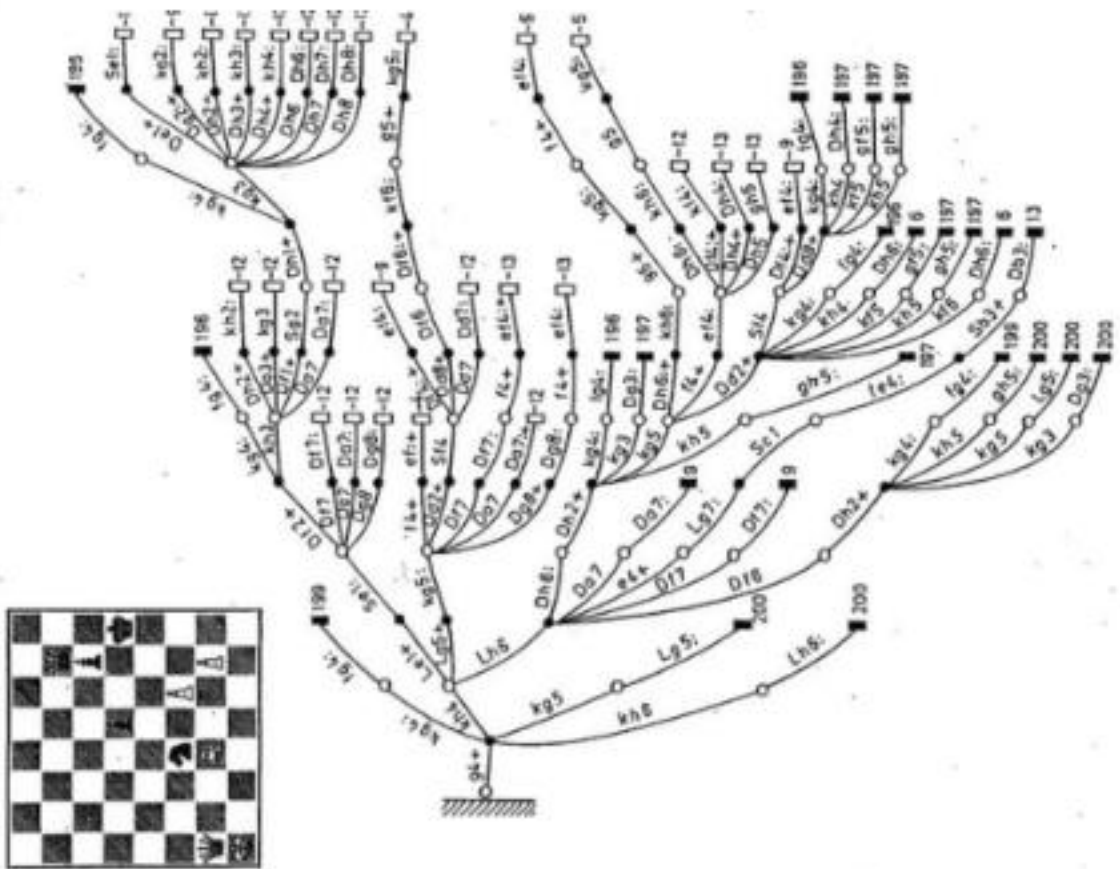
Claude Shannon

Trenchard More

Dartmouth Summer Research Project on Artificial Intelligence, 1956

Computer Science

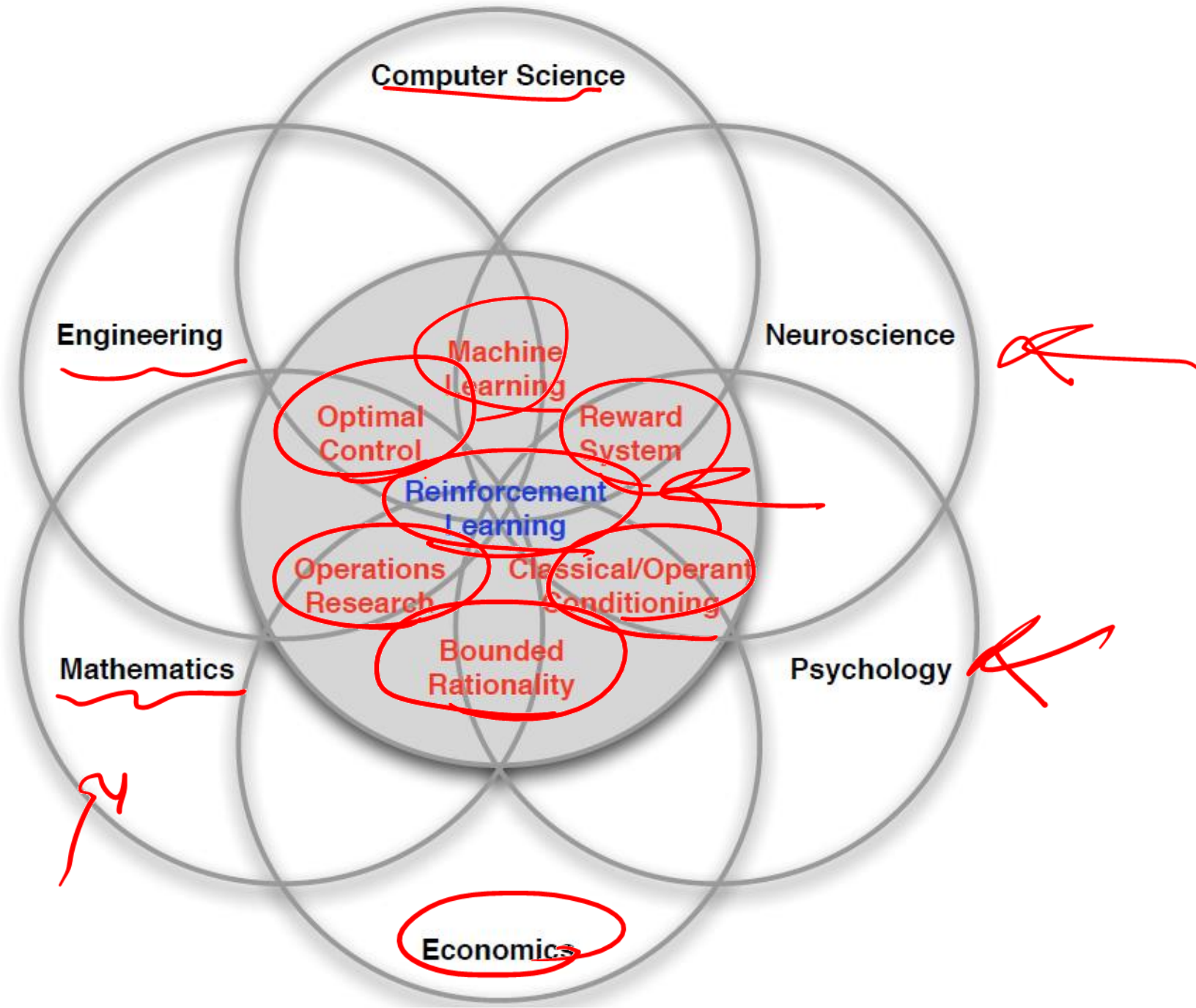Neuroscience

Engineering

Machine Learning

Optimal Control

Reward System

Reinforcement Learning

Operations Research

Classical/Operant Conditioning

Bounded Rationality

Mathematics

Psychology

Economics

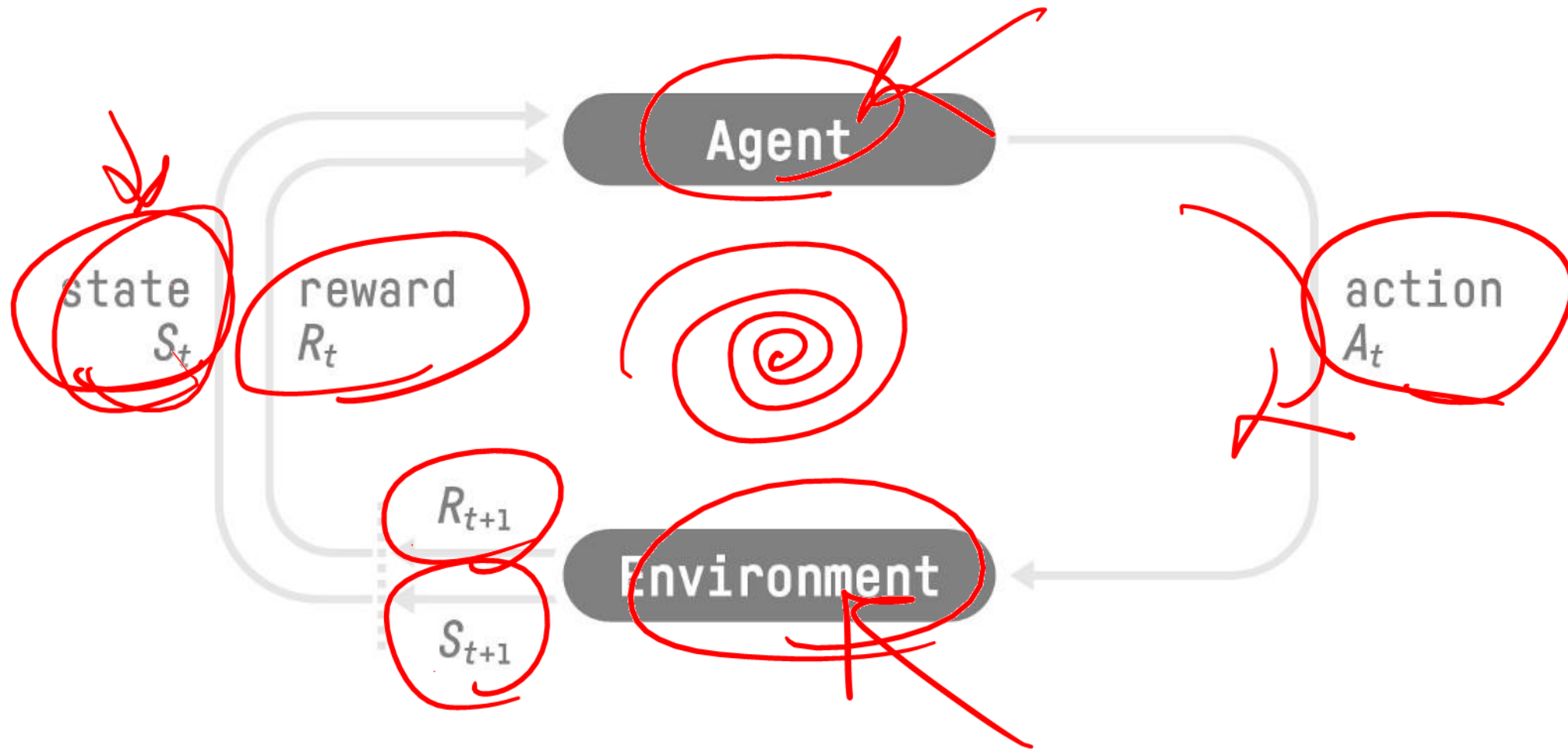Reinforcement learning can be viewed as a microcosm of the whole AI problem

Richard S. Sutton

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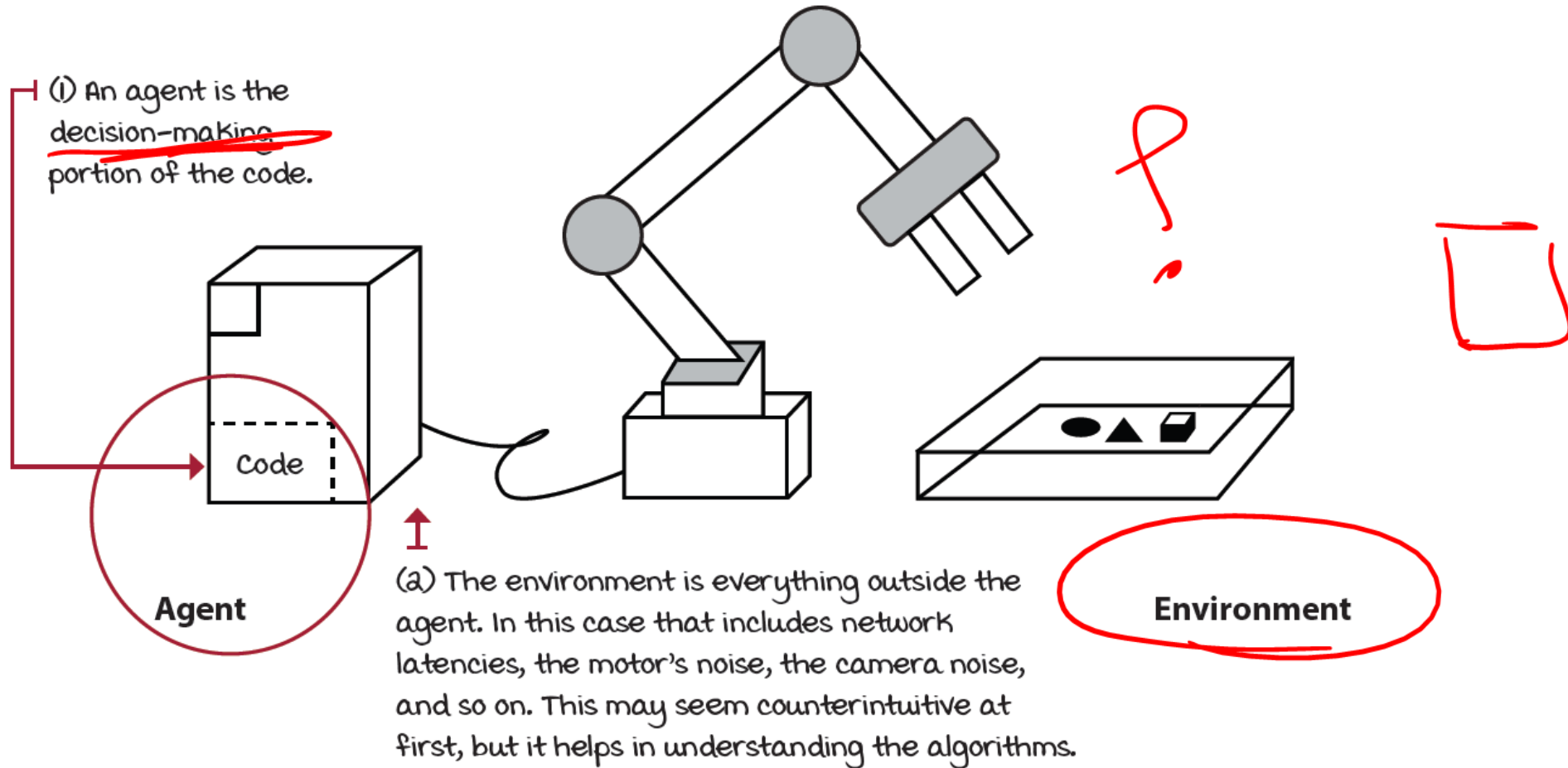
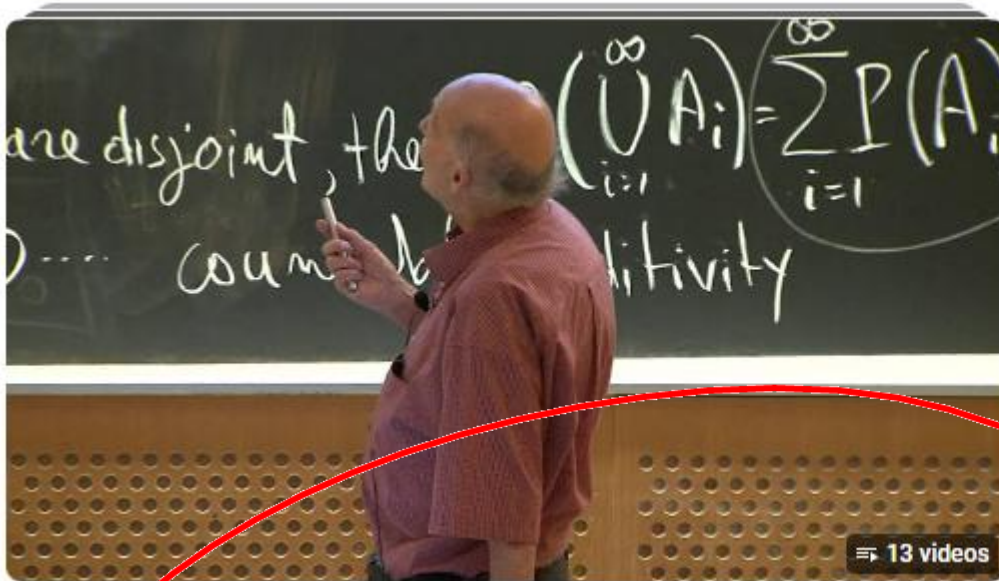| | AI Planning | SL | UL | RL | IL |
|---|---|---|---|---|---|
| Optimization | X | | | X | X |
| Learns from experience | | X | X | X | X |
| Generalization | X | X | X | X | X |
| Delayed Consequences | X | | | X | X |
| Exploration | | | | X | |

Agent

state
$S_t$

reward
$R_t$

$R_{t+1}$

$S_{t+1}$

Environment

action
$A_t$

# How should we define the boundary between agent and environment?

# ENVIRONMENT AND AGENT

(1) An agent is the decision-making portion of the code.

Code

**Agent**

(2) The environment is everything outside the agent. In this case that includes network latencies, the motor's noise, the camera noise, and so on. This may seem counterintuitive at first, but it helps in understanding the algorithms.

**Environment**

MIT 6.868J The Society of Mind, Fall 2011

MIT OpenCourseWare · Playlist

1. Introduction to 'The Society of Mind' · 2:05:54
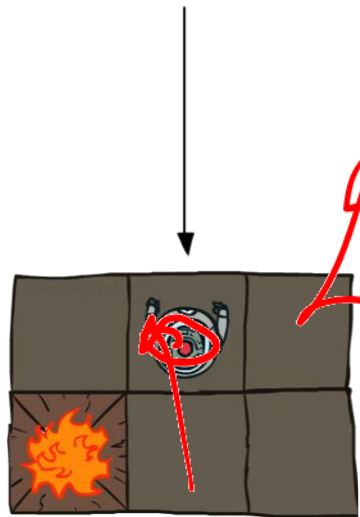2. Falling In Love · 1:45:55

View full playlist

13 videos

https://www.youtube.com/playlist?list=PLUl4u3cNGP61E-vNcDV0w5xpslBYNJDkU

THE
SOCIETY
OF
MIND

"270 brilliantly original essays... on how the mind works."
—Isaac Asimov, Information Week

MARVIN MINSKY

COFOUNDER OF THE ARTIFICIAL INTELLIGENCE LABORATORY, MIT

1997

**Three Dogmas of Reinforcement Learning**

August 6, 2024 · Arash Alikhani

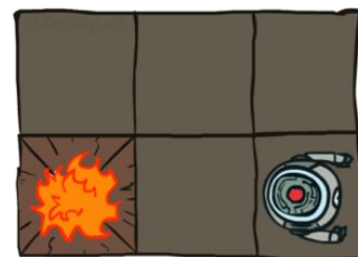https://rljclub.github.io/posts/three-dogmas-of-reinforcement-learning
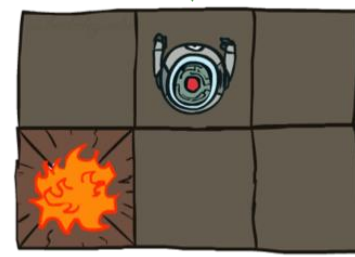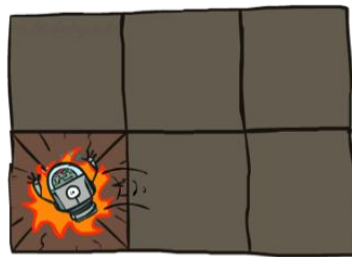
# ENVIRONMENT AND AGENT

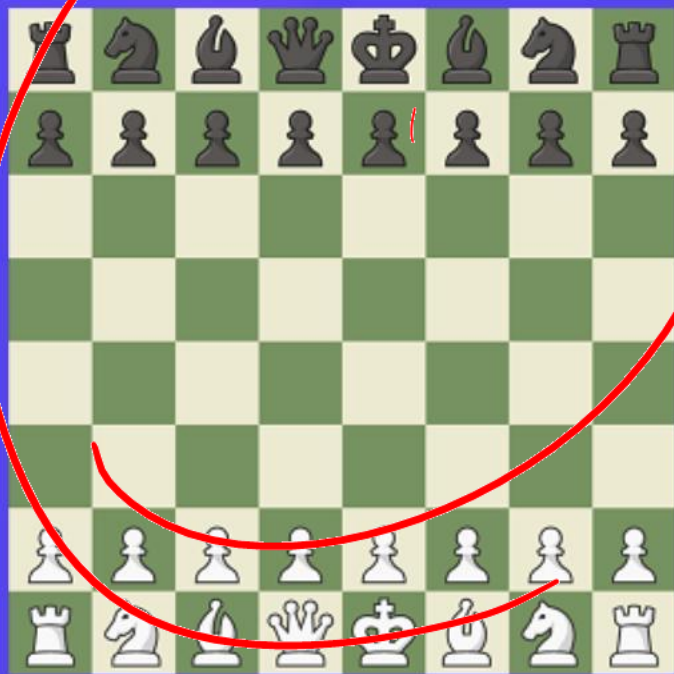Deterministic Grid World

Stochastic Grid World

# Observation Space

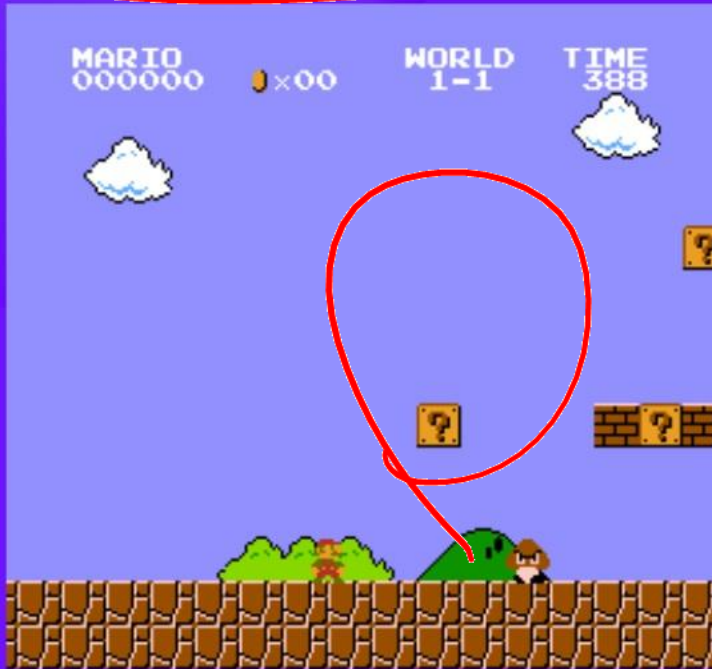*State*: **complete description** of the state of the world (no hidden information).

*Observation*: **partial description** of the state of the world.

# Action Space

*Discrete*: finite number of possible actions

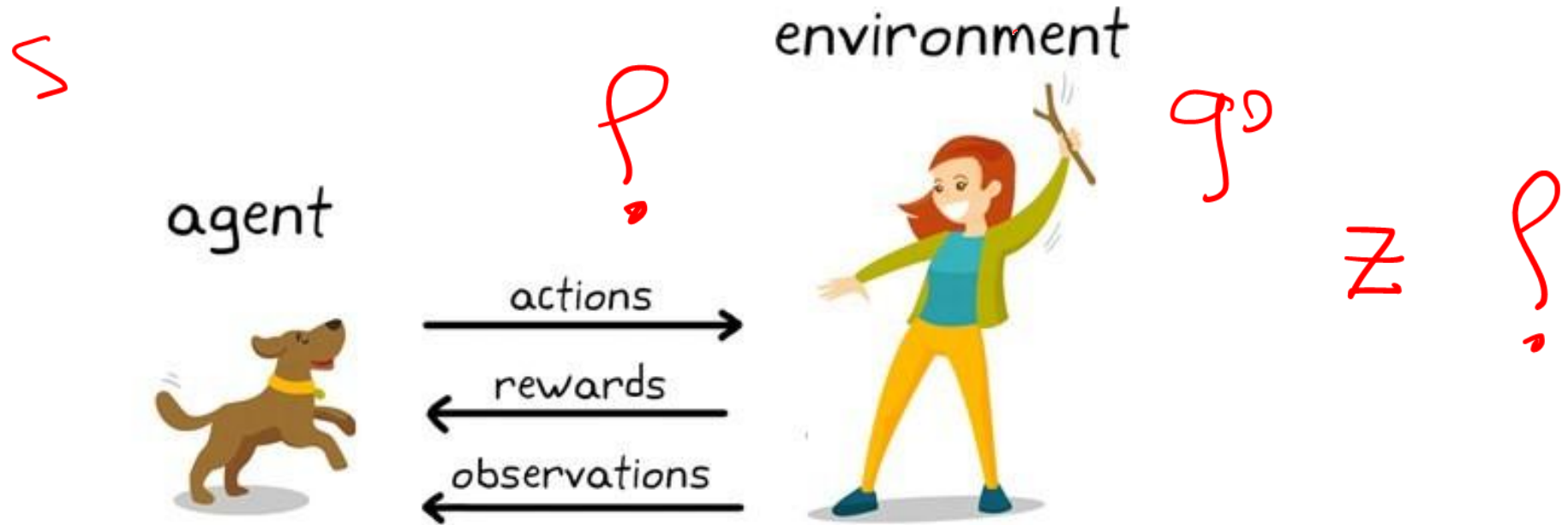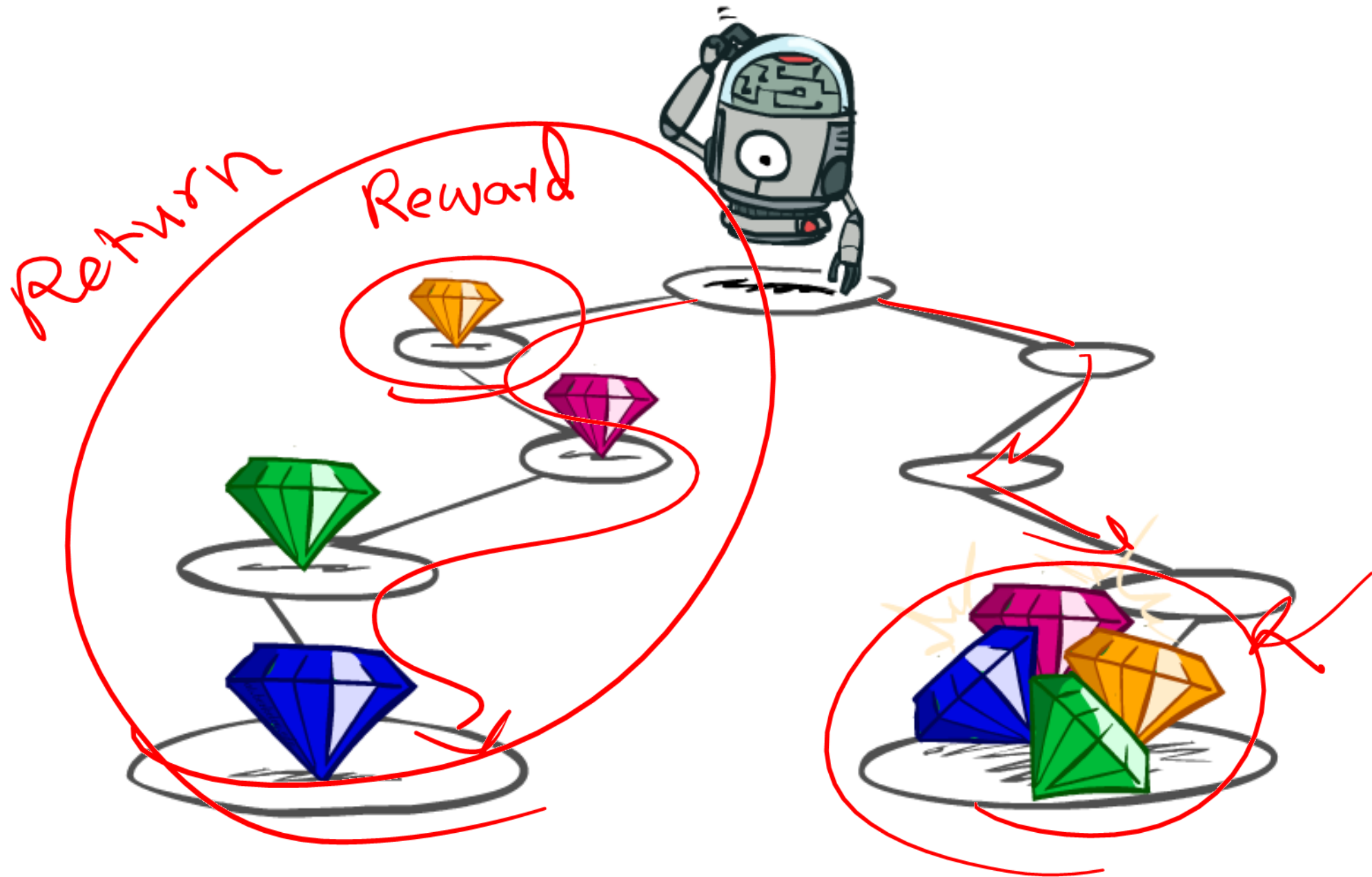*Continuous*: infinite number of possible actions

# REWARD HYPOTHESIS

## The Reward Hypothesis

"...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)"

-- Sutton (2004)

environment

agent

actions →

← rewards

← observations

$$R(\tau) = \sum_{k=0}^{\infty} r_{t+k+1}$$

trajetry

$$R(\tau) = r_{t+1} + r_{t+2} + r_{t+3} + r_{t+4} + \ldots$$

Return: cumulative reward

Trajectory (read Tau)
Sequence of states and actions

Worth Now $1$

Worth Next Step $\gamma$

Worth In Two Steps $\gamma^2$

$$R(\boxed{\tau}) = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \ldots$$
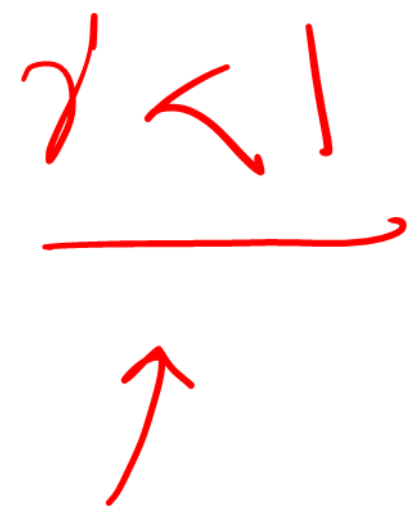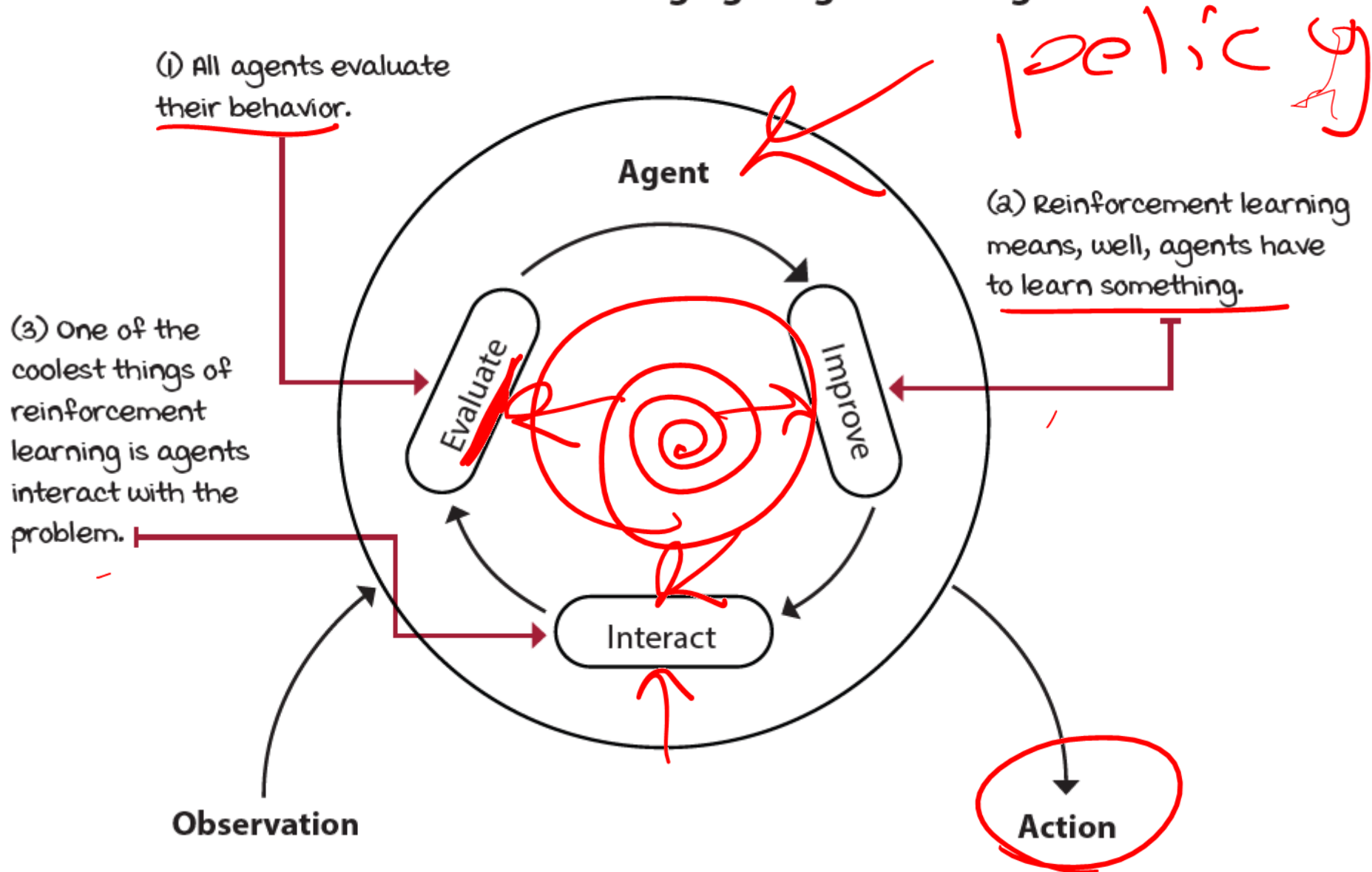
Return: cumulative reward

Gamma: discount rate

Trajectory (read Tau)
Sequence of states and actions

$$R(\tau) = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

$$\gamma < 1$$

# The three internal steps that every reinforcement learning agent goes through



(1) All agents evaluate their behavior.

(2) Reinforcement learning means, well, agents have to learn something.

(3) One of the coolest things of reinforcement learning is agents interact with the problem.

policy

Agent

Evaluate

Improve

Interact

Observation

Action

# The Policy π: the agent's brain

Policy π: is the **brain of our Agent,** it's the function that tell us what **action to take given the state we are.**
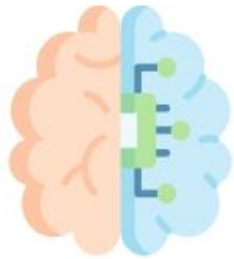→ So it **defines the agent behavior at a given time.**



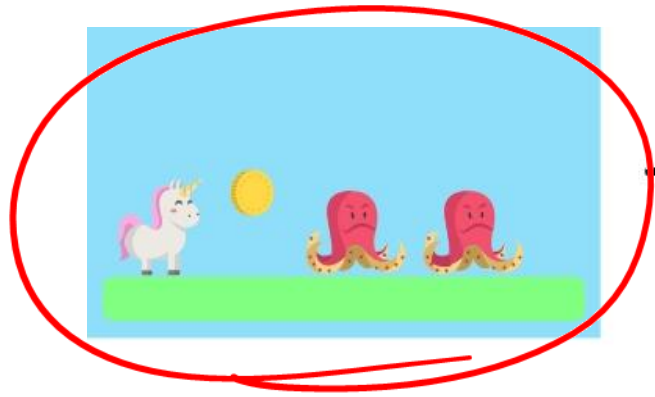State → π(State) → Action

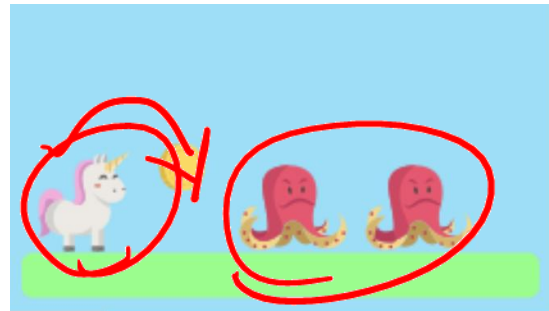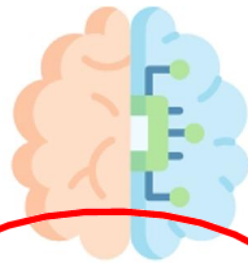$$a = \pi(s)$$



State $s_0$   →   $\pi(s_0)$   →   $a_0$ = Right

$$\pi(a|s) = P[A|s]$$

Probability Distribution over the
set of actions given the state

State $s_0$ → $\pi(A|s_0)$ → [Left: 0.1, Right: 0.7, Jump: 0.2]

# Process the environment goes through as a consequence of agent's actions

(5) Finally, the reaction is passed back to the agent.

**Environment**

Observation, reward

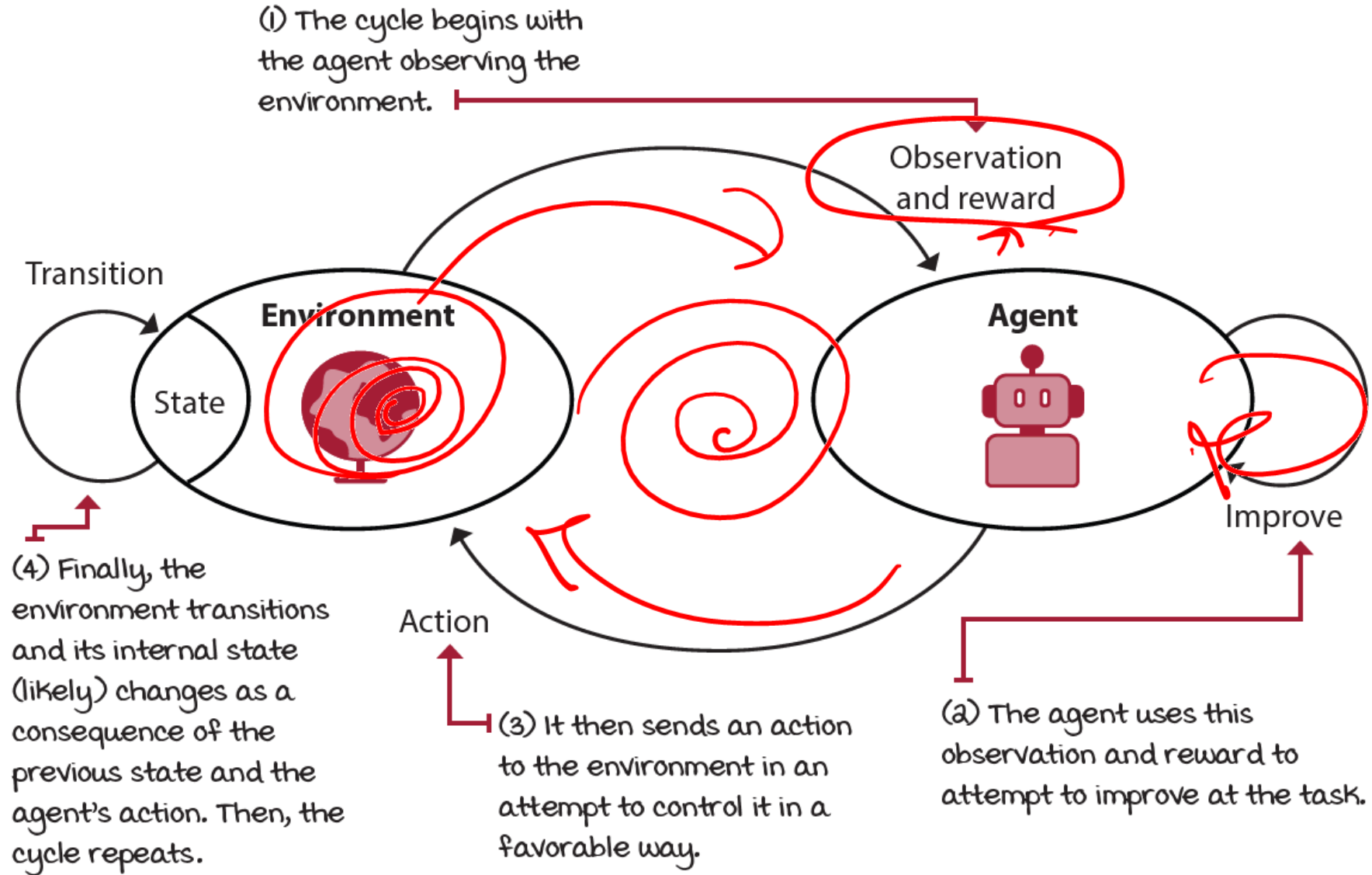(4) The new state and reward are passed through a filter: some problems don't let the true state of the environment be perceived by the agent!

Transition

(1) Environment receives the action selected by the agent.

Action

Next state

Reward

State

(3) . . . the environment will transition to a new internal state.

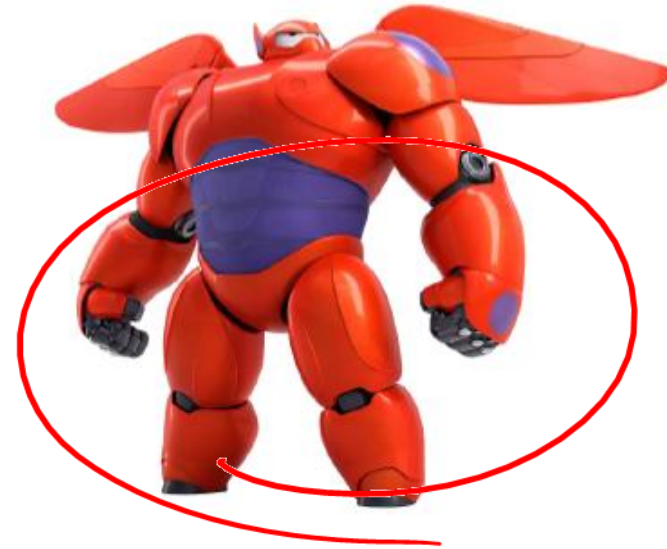(2) Depending on the current environment state, and the agent's chosen action . . .

# The reinforcement learning cycle

(1) The cycle begins with the agent observing the environment.

Observation and reward

Transition

**Environment**

State

**Agent**

Improve

Action

(4) Finally, the environment transitions and its internal state (likely) changes as a consequence of the previous state and the agent's action. Then, the cycle repeats.

(3) It then sends an action to the environment in an attempt to control it in a favorable way.

(2) The agent uses this observation and reward to attempt to improve at the task.
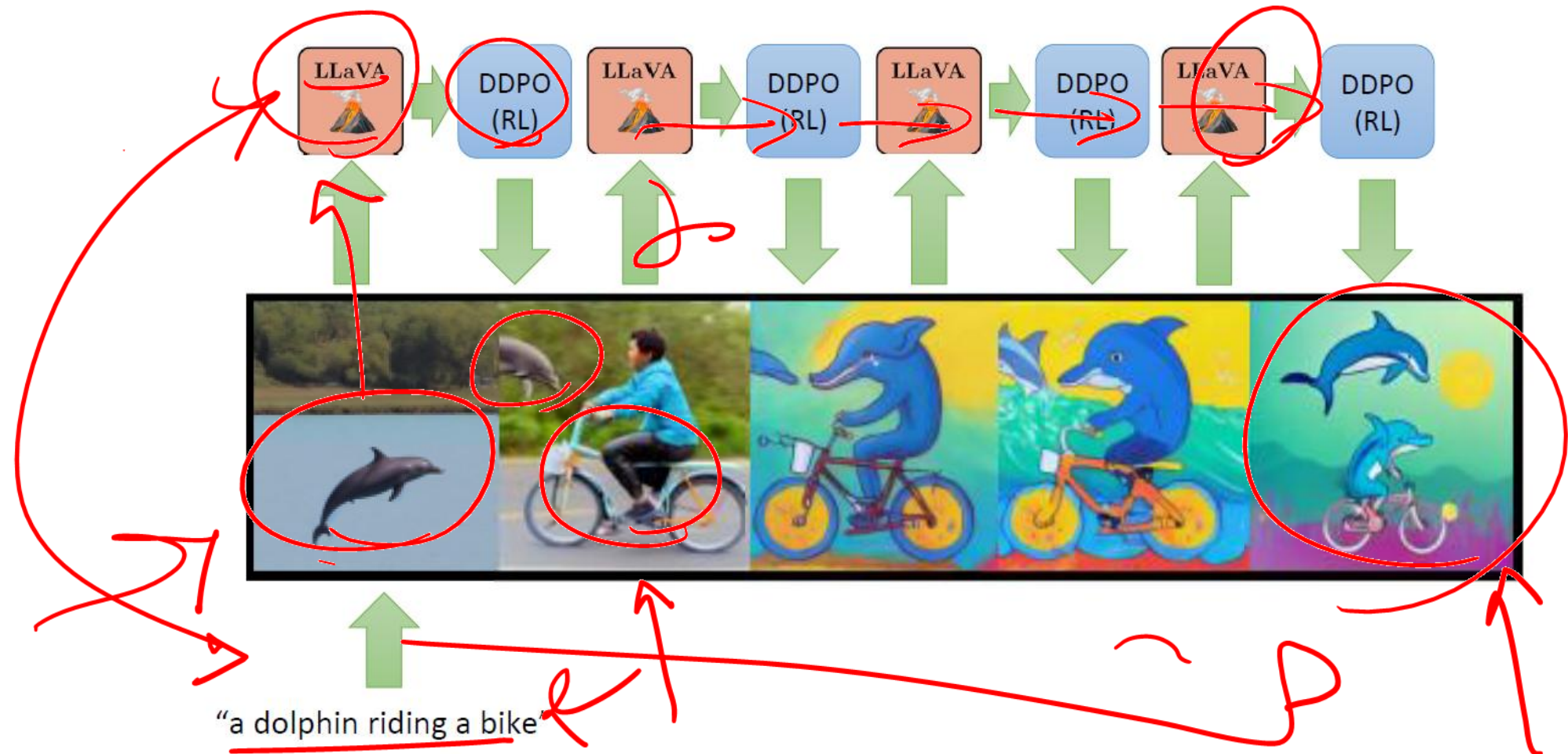
Actions: muscle contractions
Observations: sight, smell
Rewards: food

Actions: motor current or torque
Observations: camera images
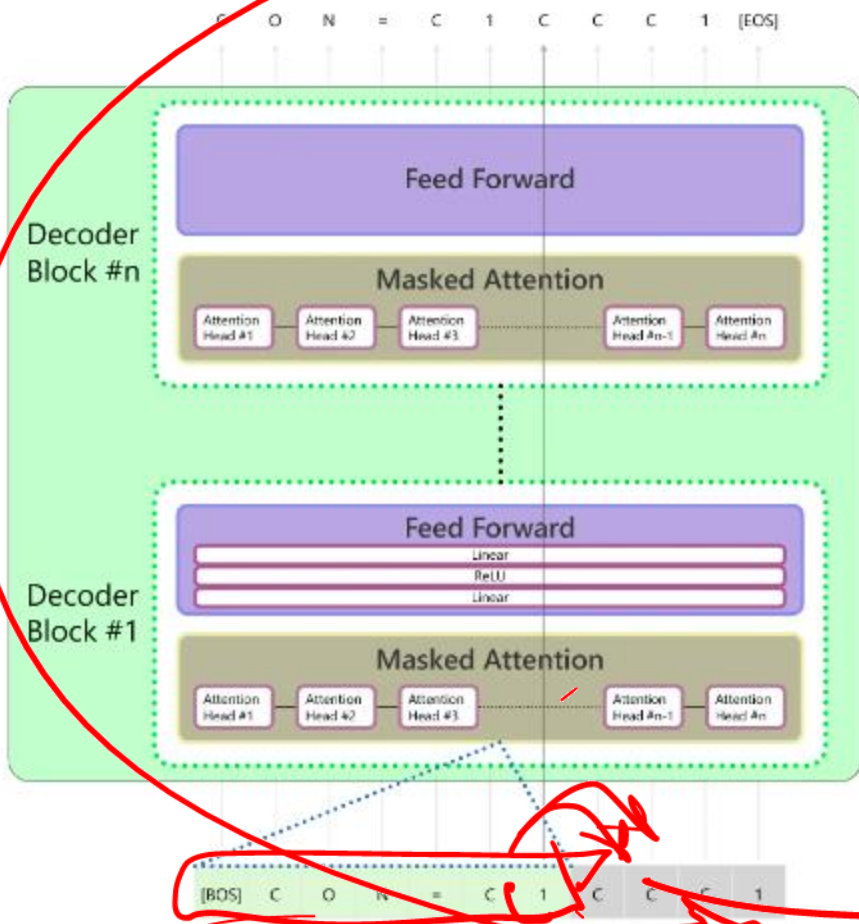Rewards: task success measure (e.g., running speed)

Inventory Management

Actions: what to purchase
Observations: inventory levels
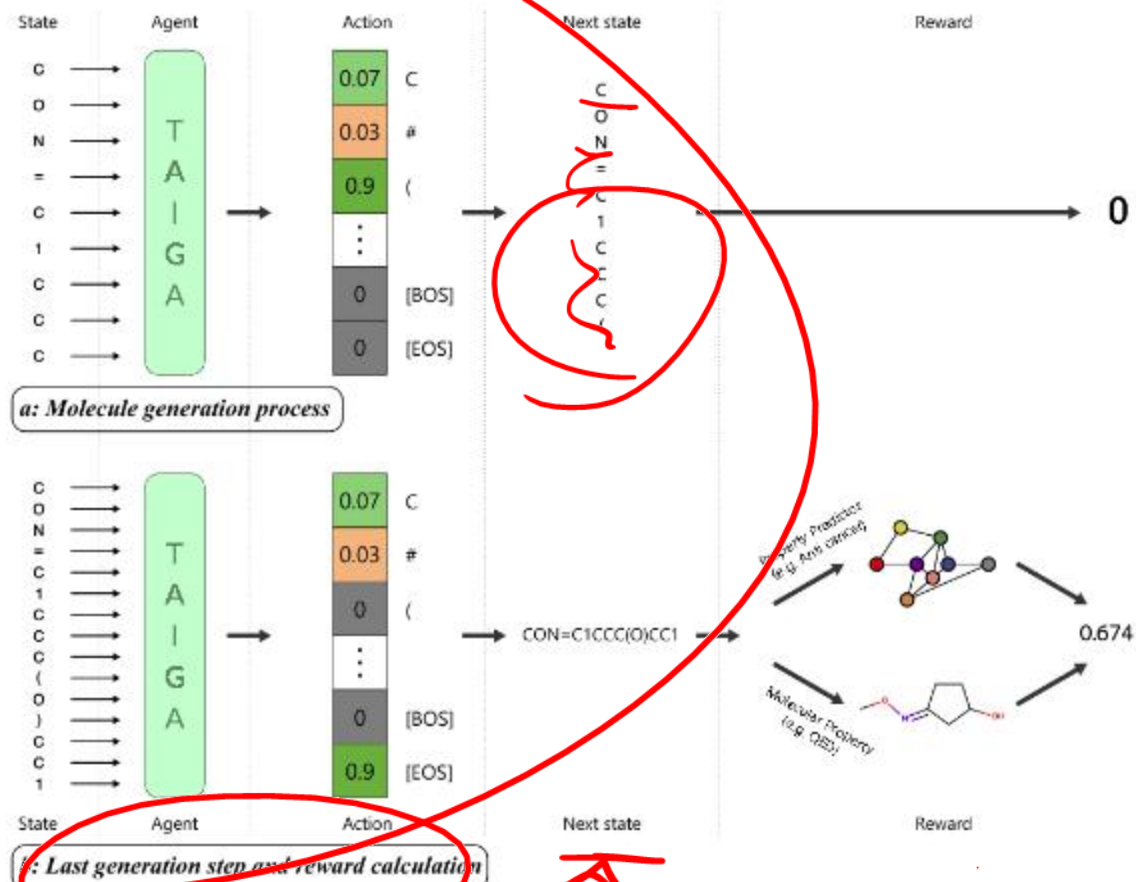Rewards: profit
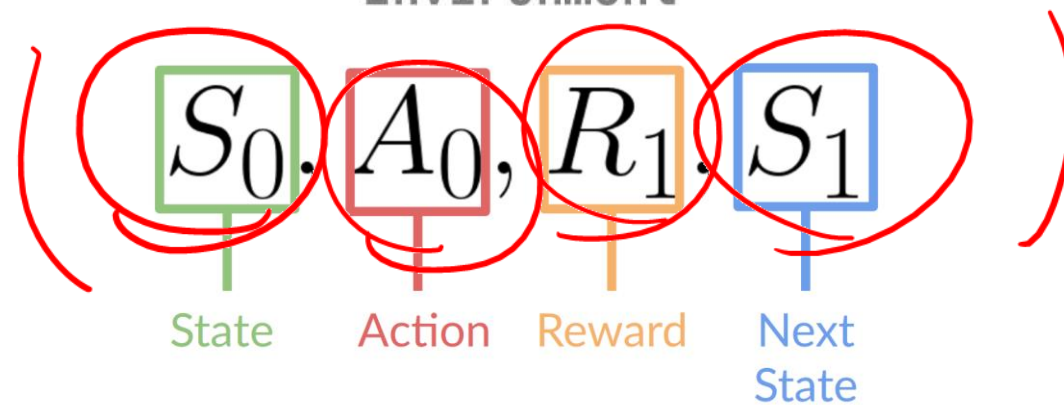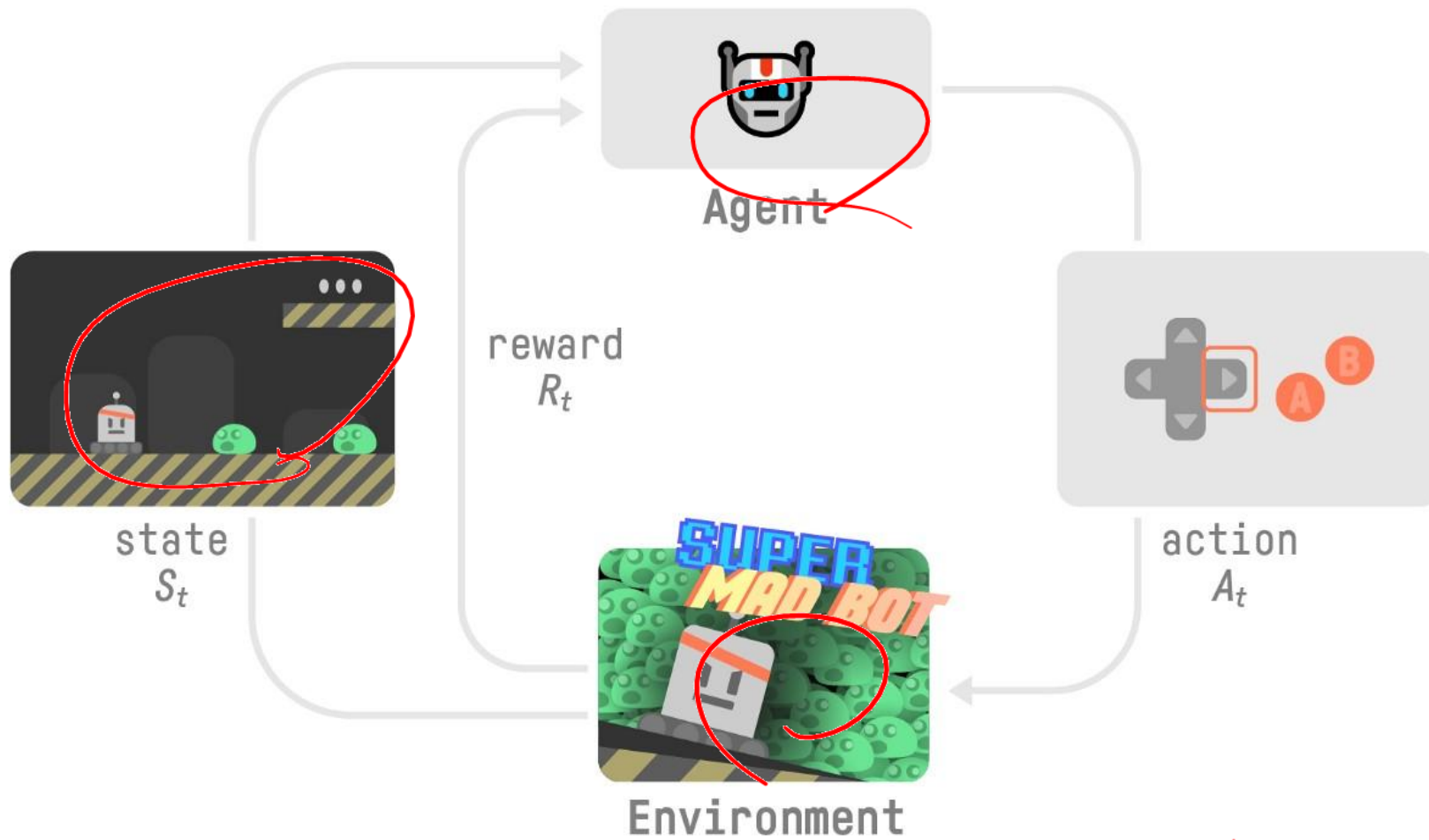
# Reinforcement learning with image generation



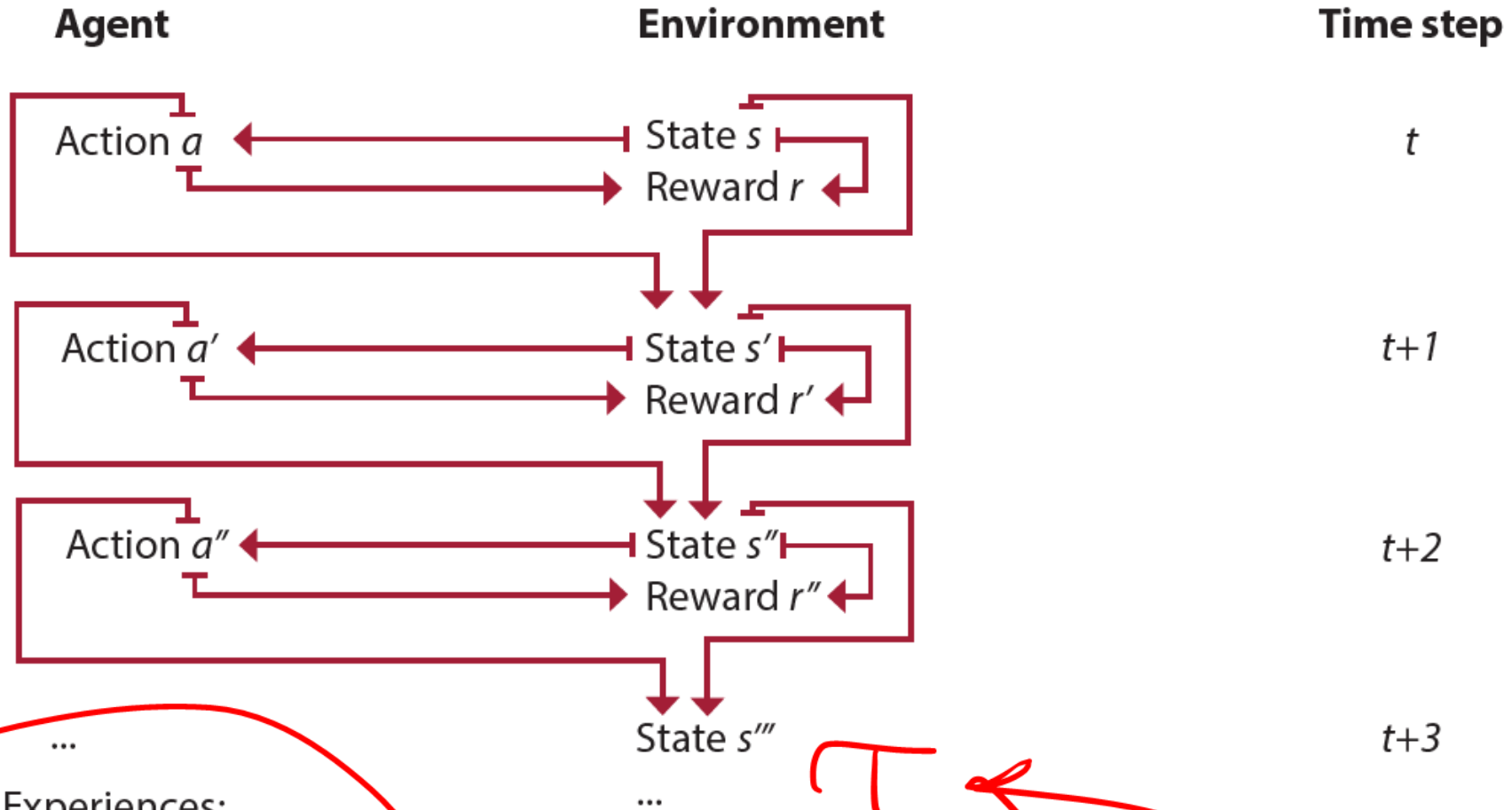"a dolphin riding a bike"

Kevin Black*, Michael Janner*, Yilun Du, Ilya Kostrikov, Sergey Levine. **Training Diffusion Models with Reinforcement Learning**. 2023.

C O N = C 1 C C C 1 [EOS]

Taiga

Decoder Block #n

Feed Forward

Masked Attention

Attention Head #1 → Attention Head #2 → Attention Head #3 ......... Attention Head #n-1 → Attention Head #n

Decoder Block #1

Feed Forward

Linear
ReLU
Linear

Masked Attention

Attention Head #1 → Attention Head #2 → Attention Head #3 ......... Attention Head #n-1 → Attention Head #n

[BOS] C O N = C 1 C C C 1

Stage 1

Stage 2

State    Agent    Action    Next state    Reward

C
O
N
=
C
1
C
C
C

T A I G A

0.07    C
0.03    #
0.9     (
⋮
0       [BOS]
0       [EOS]

CON=C1CCC(

0

a: Molecule generation process

C
O
N
=
C
1
C
C
C
(
O
)
C
C
1

T A I G A

0.07    C
0.03    #
0       (
⋮
0       [BOS]
0.9     [EOS]

CON=C1CCC(O)CC1

0.674

State    Agent    Action    Next state    Reward

b: Last generation step and reward calculation

Agent

state
$S_t$

reward
$R_t$

action
$A_t$

Environment

$$( S_0, A_0, R_1, S_1 )$$

State     Action     Reward     Next State

# Experience tuples

| Agent | Environment | Time step |
|---|---|---|
| Action $a$ | State $s$ | $t$ |
| | Reward $r$ | |
| Action $a'$ | State $s'$ | $t+1$ |
| | Reward $r'$ | |
| Action $a''$ | State $s''$ | $t+2$ |
| | Reward $r''$ | |
| ... | State $s'''$ | $t+3$ |
| | ... | |

Experiences:
$t,$  $(s,\ a,\ r',\ s')$
$t+1,$ $(s',\ a',\ r'',\ s'')$
$t+2,$ $(s'',a'',r''',s''')$
...

$T$

Trajectory

$t(100$

**supervised learning**



input: $\mathbf{x}$

output: $\mathbf{y}$

data: $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

goal: $f_\theta(\mathbf{x}_i) \approx \mathbf{y}_i$

someone gives this to you

**reinforcement learning**



input: $\mathbf{s}_t$ *at each time step*
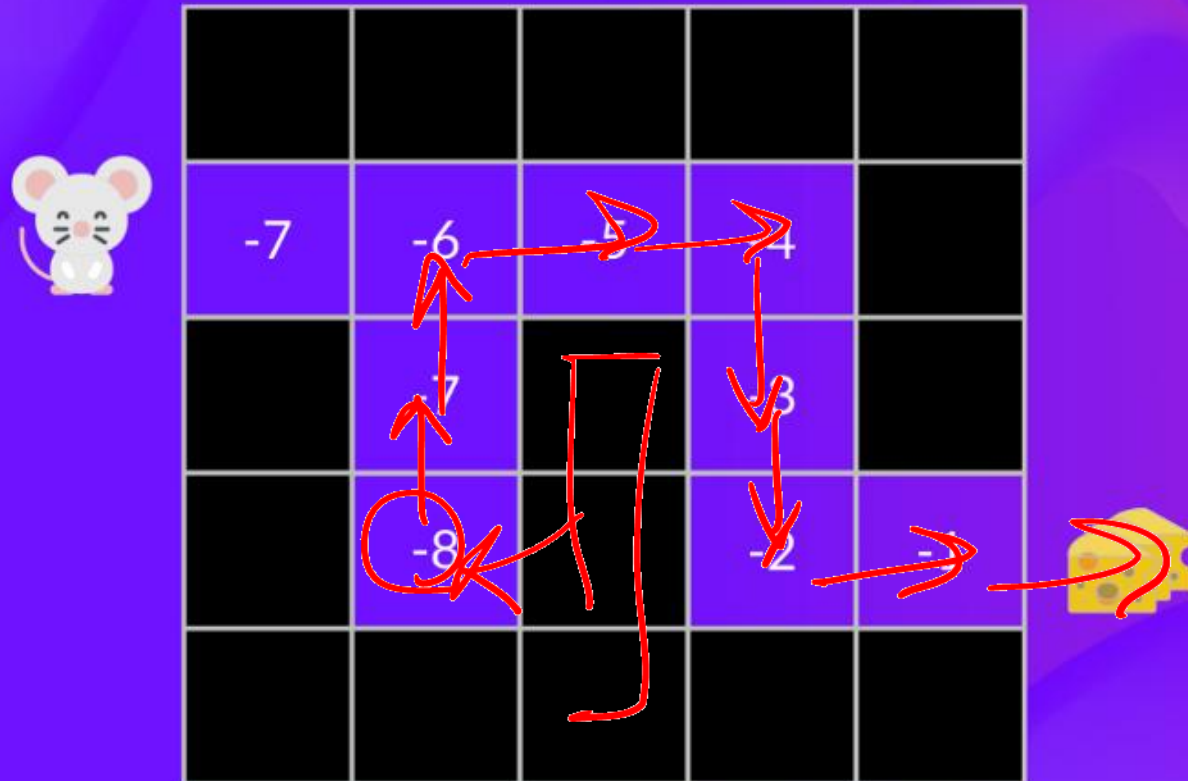
output: $\mathbf{a}_t$ *at each time step*

data: $(\mathbf{s}_1, \mathbf{a}_1, r_1, \ldots, \mathbf{s}_T, \mathbf{a}_T, r_T)$

goal: learn $\pi_\theta : \mathbf{s}_t \to \mathbf{a}_t$
to maximize $\sum_t r_t$

pick your own actions

# The State Value Function

*State Value Function:* calculate the **value of a state.**

*Return*

$$V_\pi(s) = \mathbf{E}_\pi[G_t | S_t = s]$$

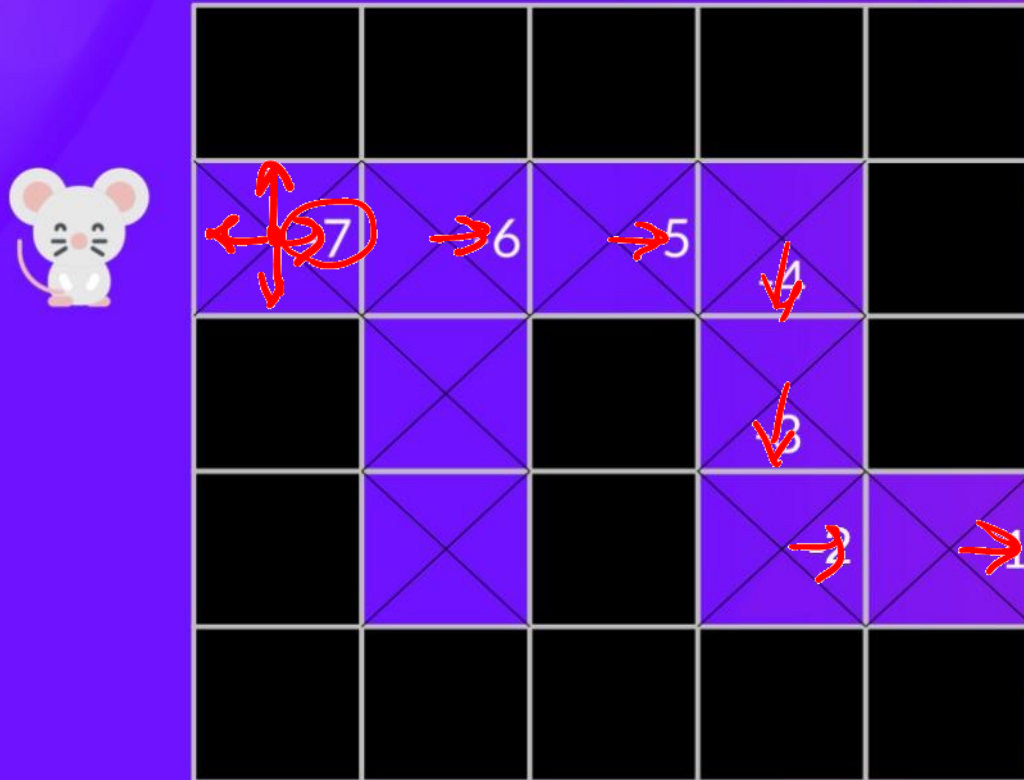Value of state s

Expected return

If the agent starts at state s

And uses the policy to choose its actions for all time steps

For each state,
the state-value function outputs
the expected return
if the agent starts in that state
and then follows the policy forever after.

# The Action Value Function

*Action Value Function:* calculate the **value of state-action pair.**



*We didn't fill all the state-actions pair for the example of Action-value function

$$Q_\pi(s, a) = \mathbf{E}_\pi[G_t | S_t = s, A_t = a]$$
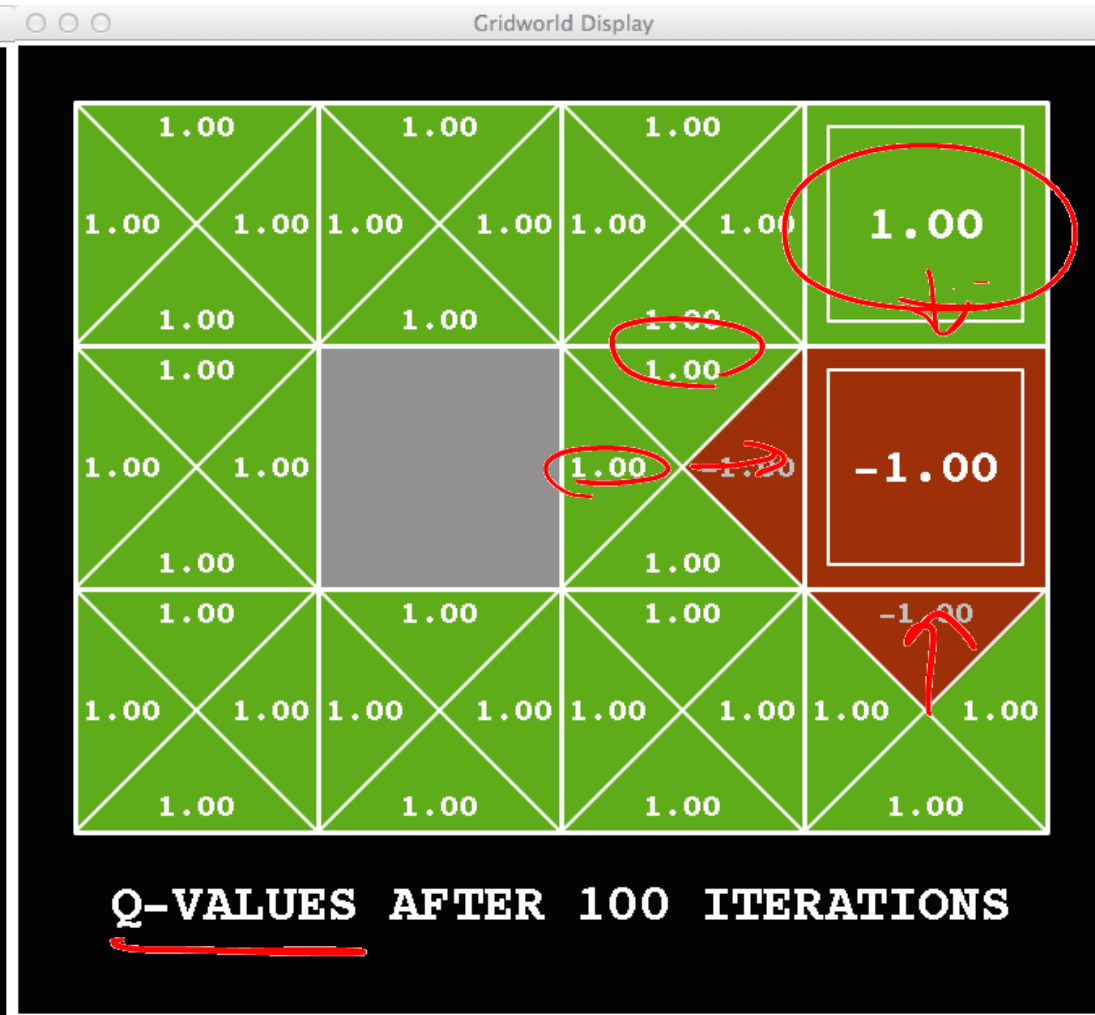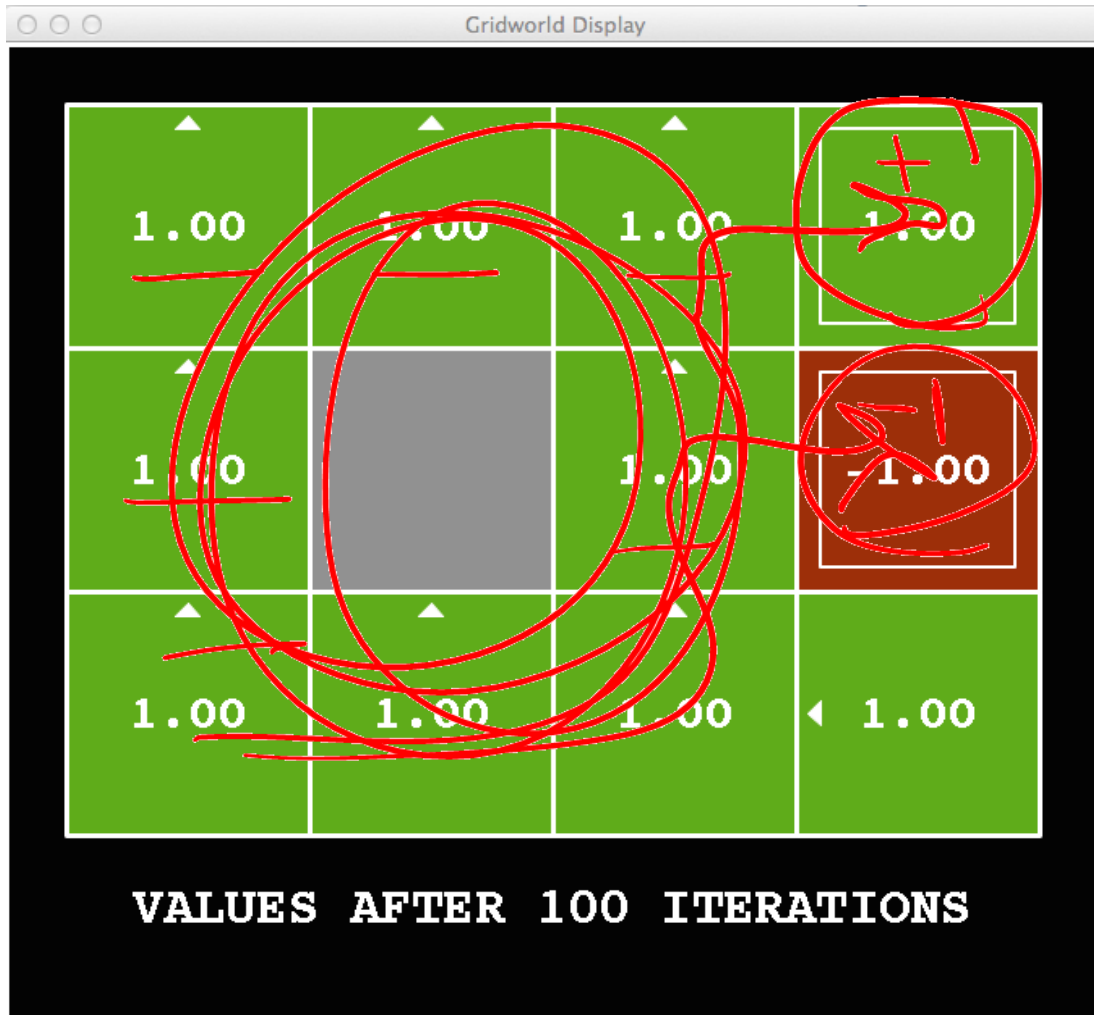
Value of state-action pair s,a

Expected return

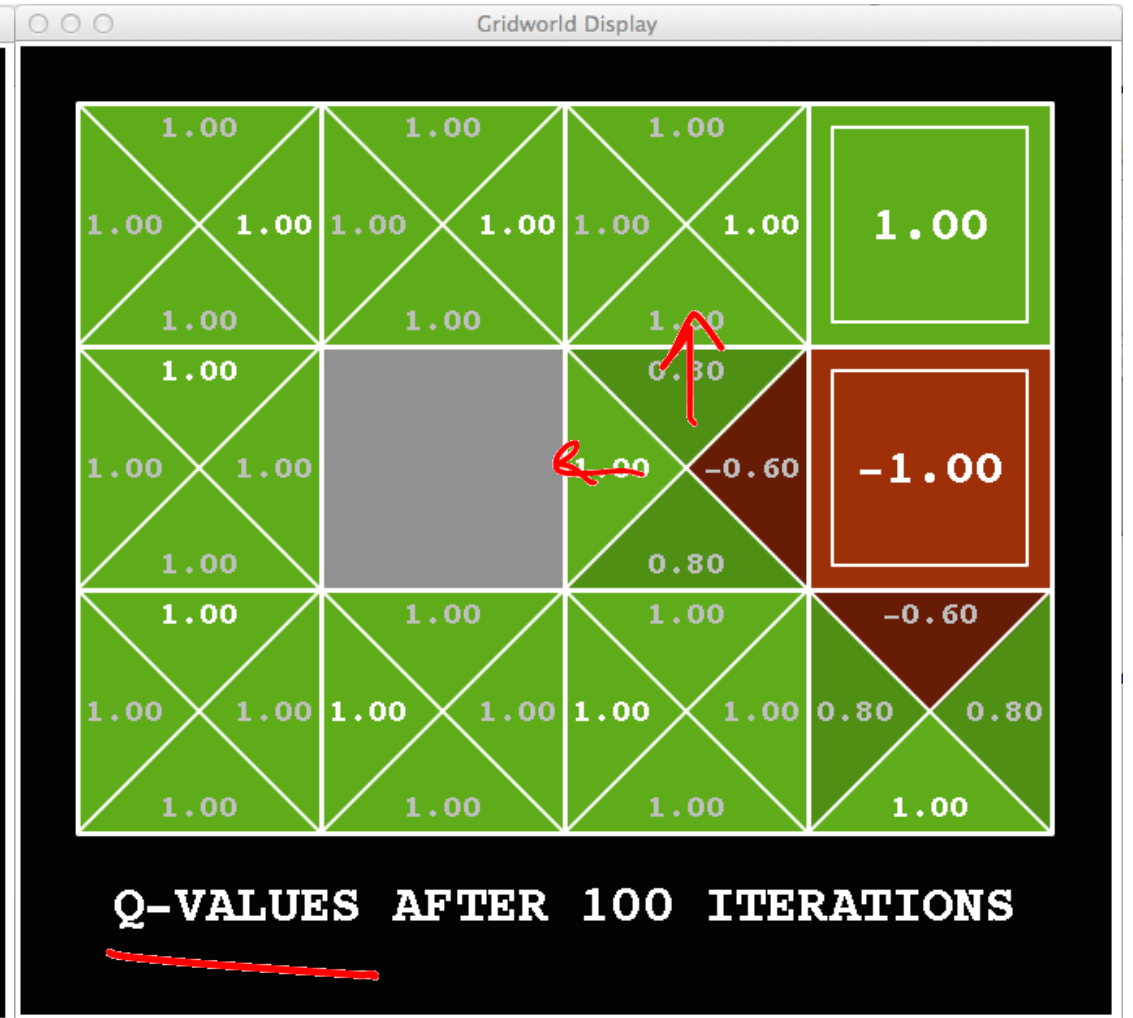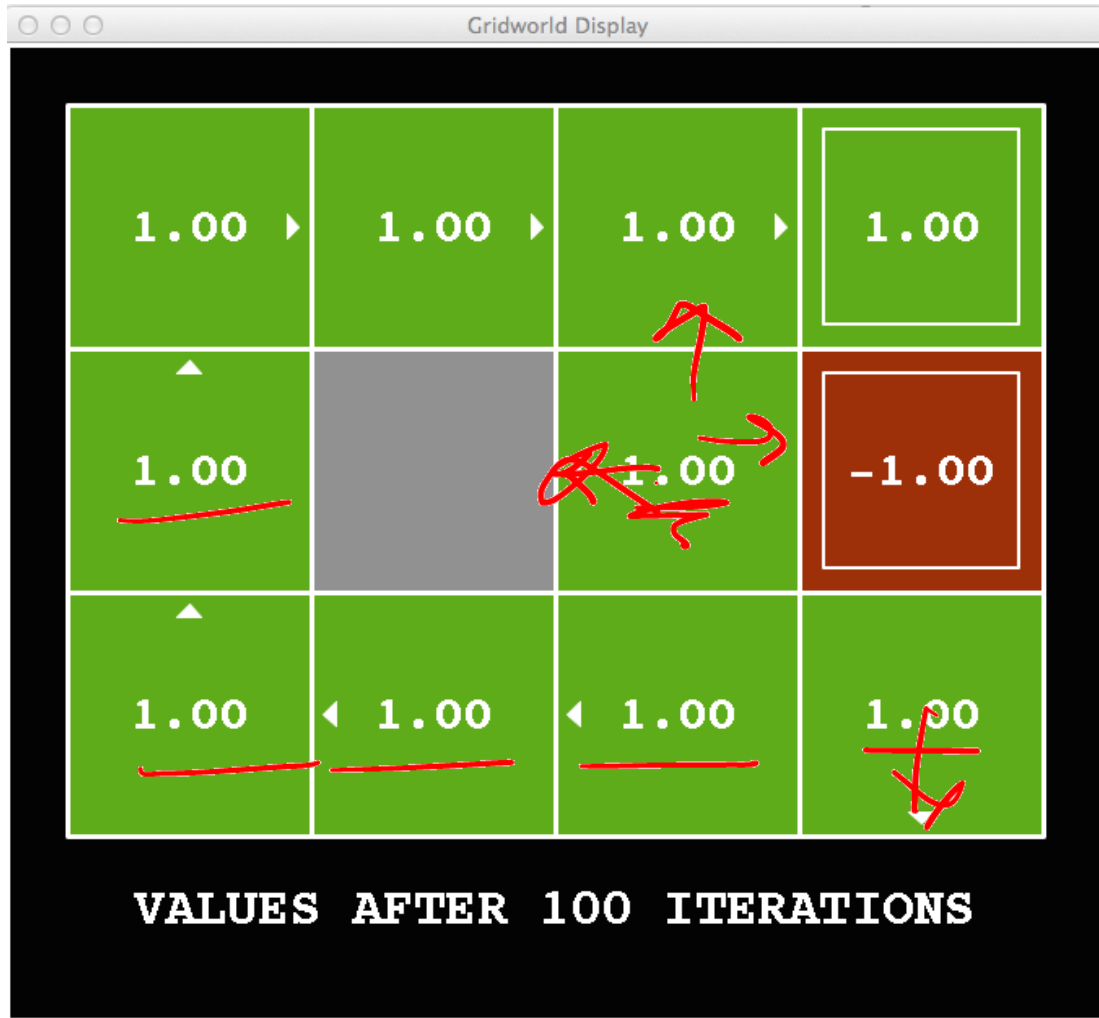If the agent starts at state s

and chooses action a

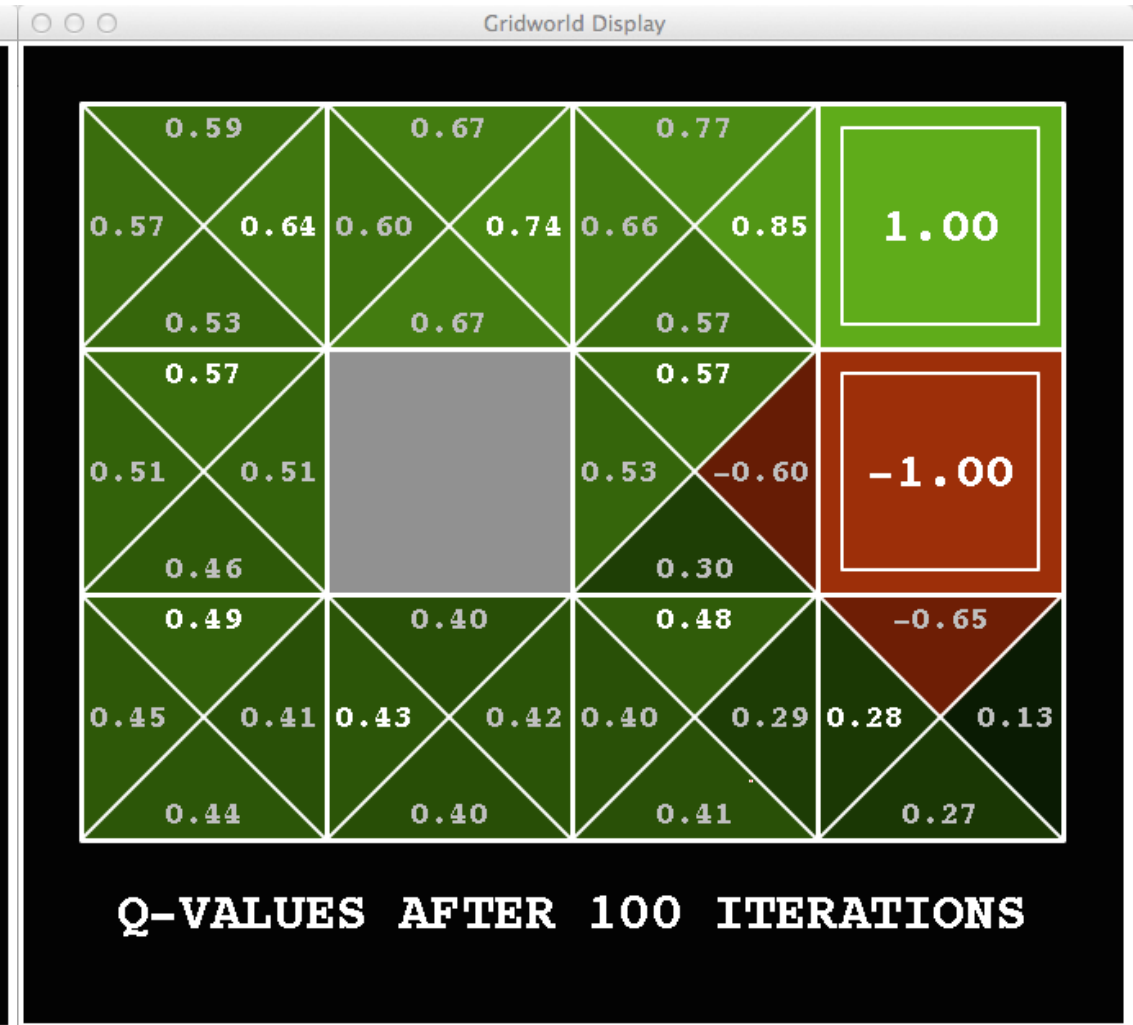And then uses the policy to choose its actions for all time steps

For each state and action,
the action-value function outputs
the expected return
if the agent starts in that state
and takes the action
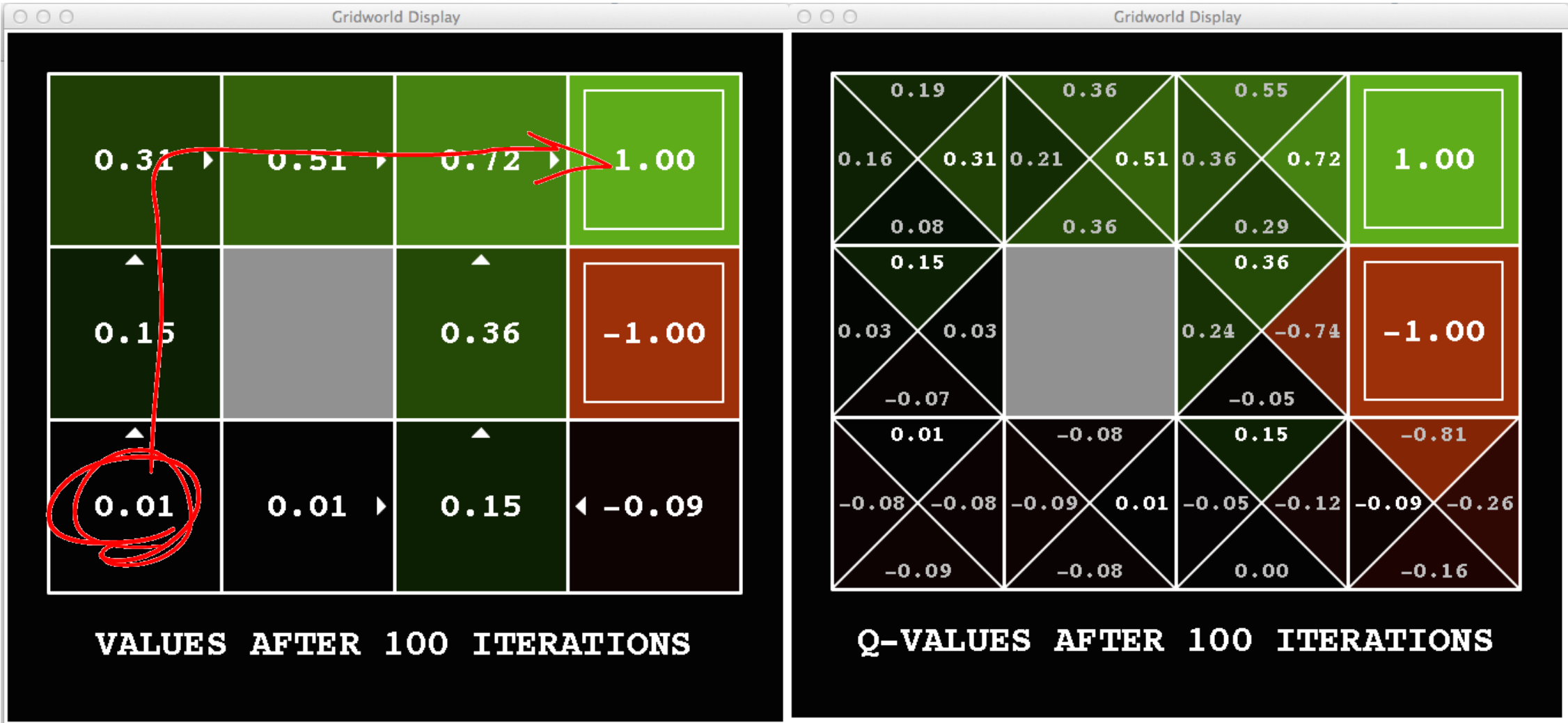and then follows the
policy forever after.

**Gridworld Display** — VALUES AFTER 100 ITERATIONS

**Gridworld Display** — Q-VALUES AFTER 100 ITERATIONS

Noise = 0
Discount = 1
Living reward = 0

VALUES AFTER 100 ITERATIONS

Q-VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 1
Living reward = 0

VALUES AFTER 100 ITERATIONS

| | | | |
|---|---|---|---|
| 0.64 ▶ | 0.74 ▶ | 0.85 ▶ | 1.00 |
| 0.57 ▲ | | 0.57 | -1.00 |
| 0.49 ▲ | ◀ 0.43 | 0.48 ▲ | ◀ 0.28 |

Q-VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = 0

VALUES AFTER 100 ITERATIONS

Q-VALUES AFTER 100 ITERATIONS

Noise = 0.2
Discount = 0.9
Living reward = -0.1

# WHAT WE HAVE LEARNED SO FAR?

- what is reinforcement learning and its actual place & significance

- reinforcement learning framework & basic concepts
  - agent
  - environment
  - state/observation
  - action
  - reward
  - policy
  - model
  - experience/trajectory/horizon
  - discount factor
  - state value function
  - action value function

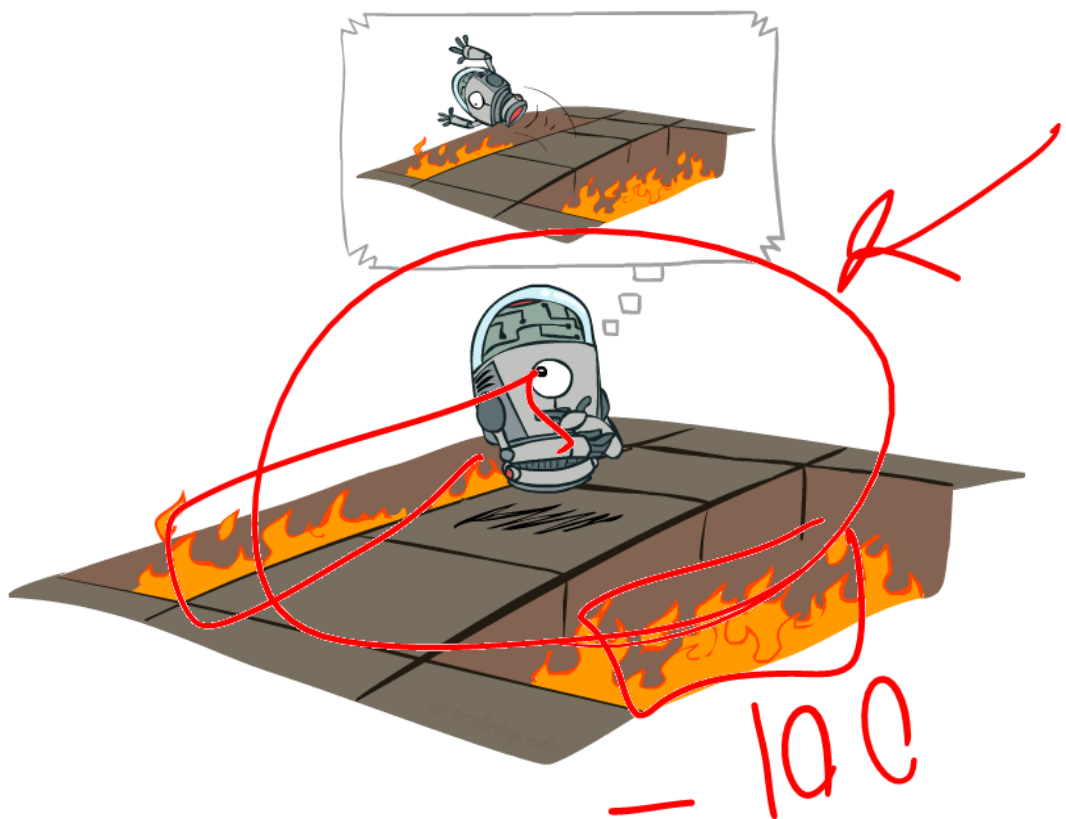# Challenges of Reinforcement Learning

# Type of tasks

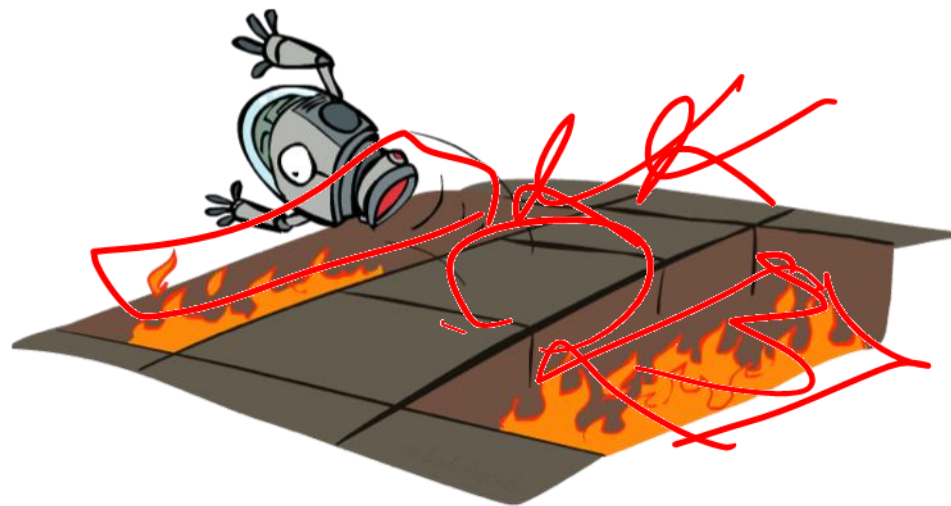*Episodic*: **starting point** and an **ending point** (a terminal state)

*Continuing*: task that **continue forever** (no terminal state)
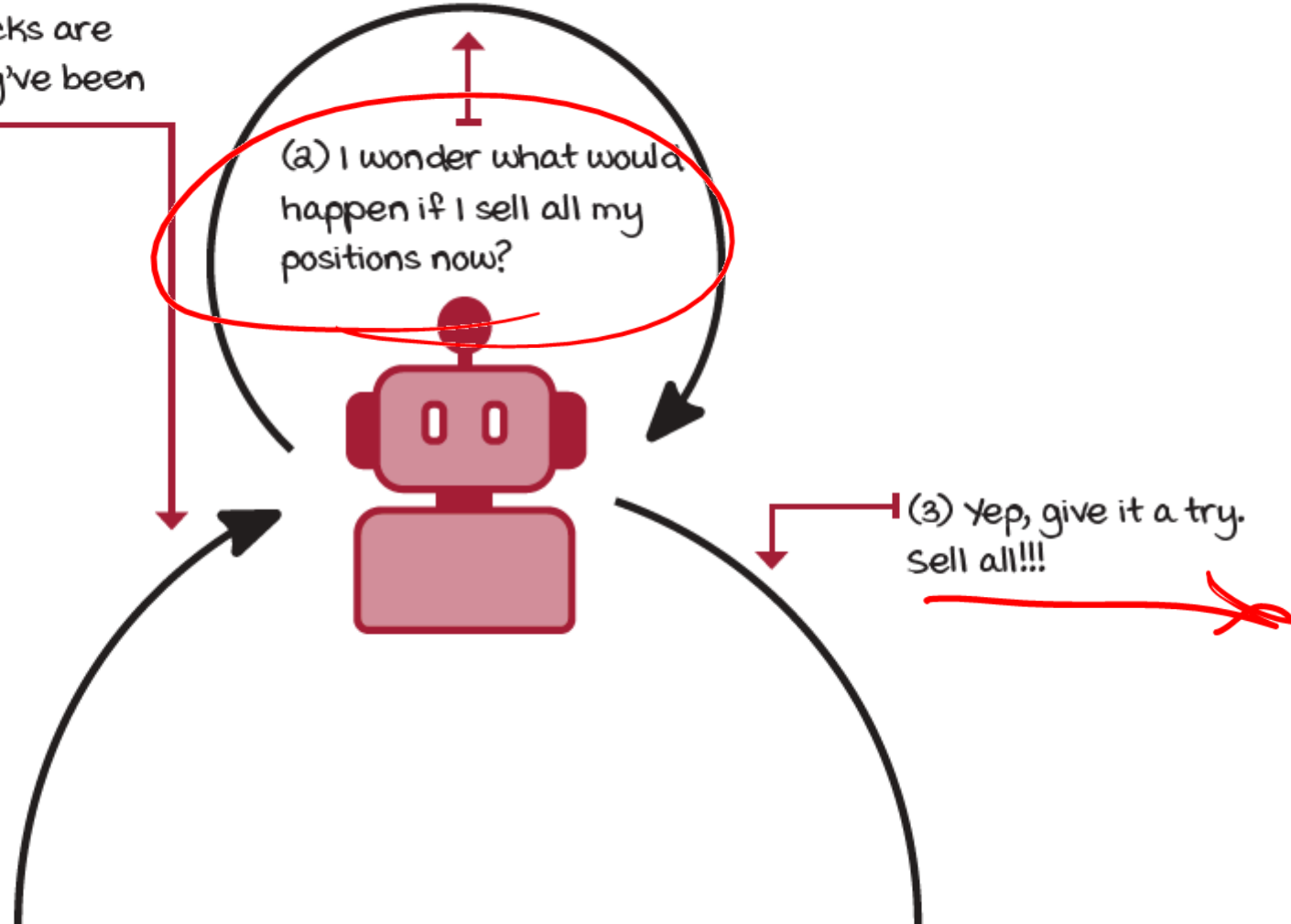
Offline Solution

Online Learning

Deep reinforcement learning agents will explore! Can you afford mistakes?
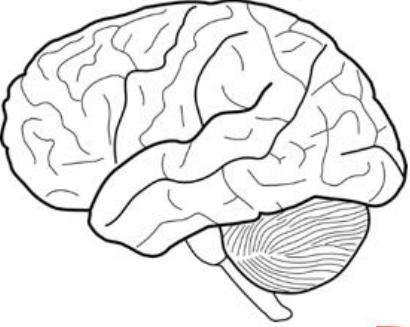
On-policy

Off-policy

(a) online reinforcement learning

(b) off-policy reinforcement learning

rollout data $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$
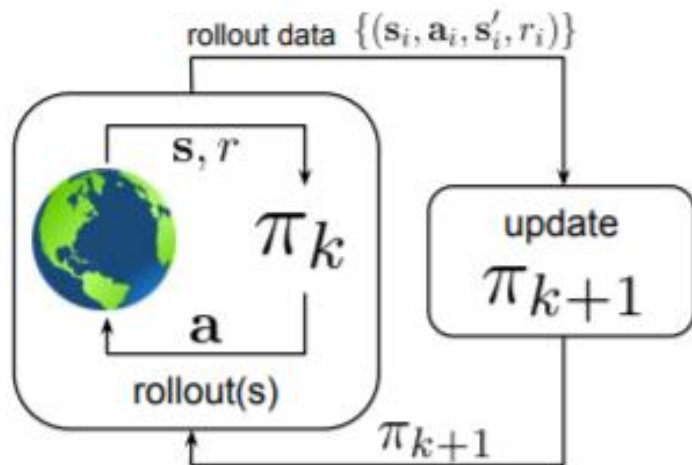
$\pi_k$

update $\pi_{k+1}$

$\mathbf{a}$

rollout(s)

$\pi_{k+1}$

rollout data $\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$

$\pi_k$

buffer $\mathcal{D}$

update $\pi_{k+1}$

$\mathbf{a}$

rollout(s)

$\pi_{k+1}$

(c) offline reinforcement learning

$\{(\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}'_i, r_i)\}$

$\mathbf{s}, r$

$\pi_\beta$

$\mathbf{a}$

rollout(s)

data collected **once** with **any** policy

buffer $\mathcal{D}$

learn $\pi$

training phase

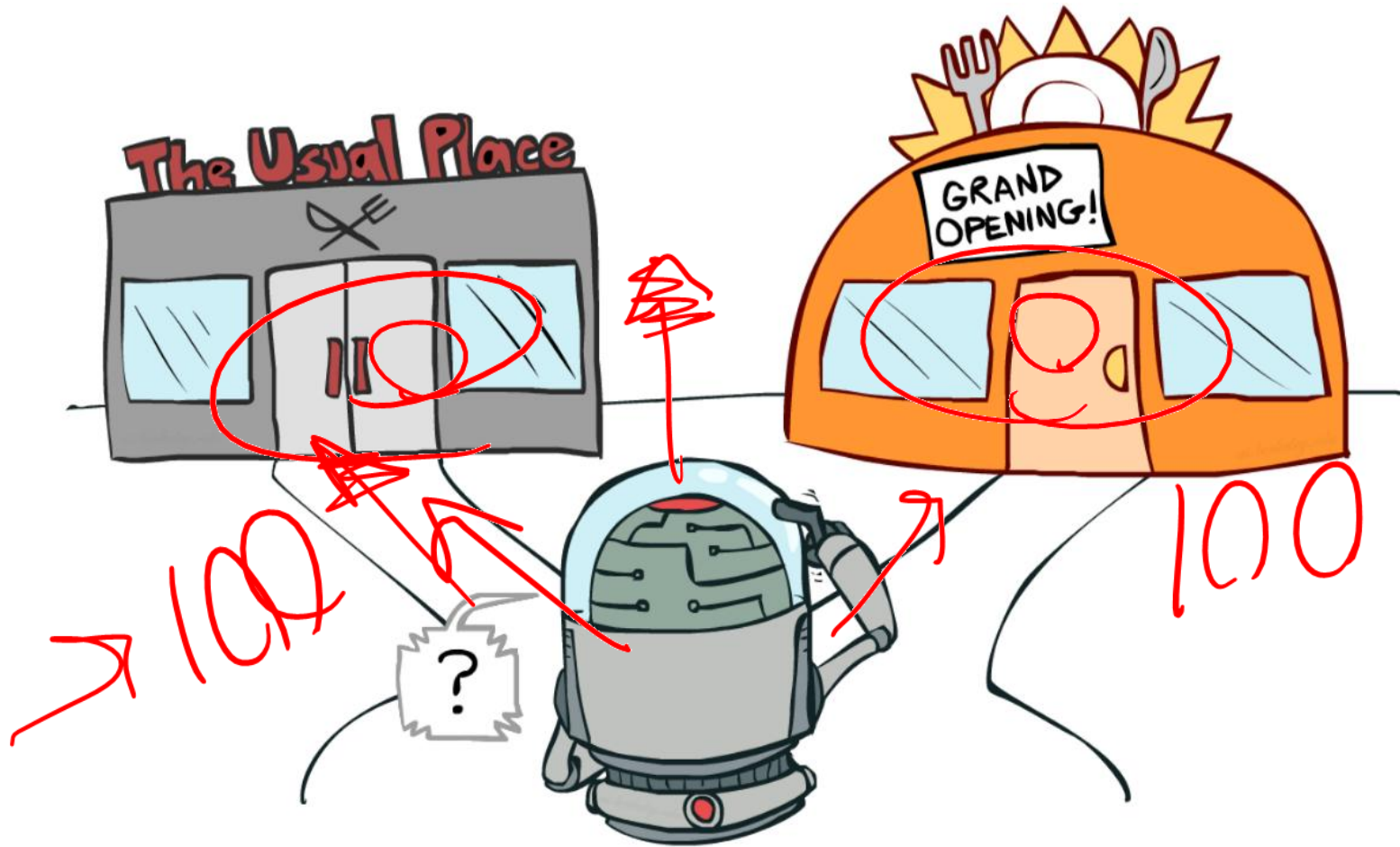$\mathbf{s}, r$

$\pi$

$\mathbf{a}$

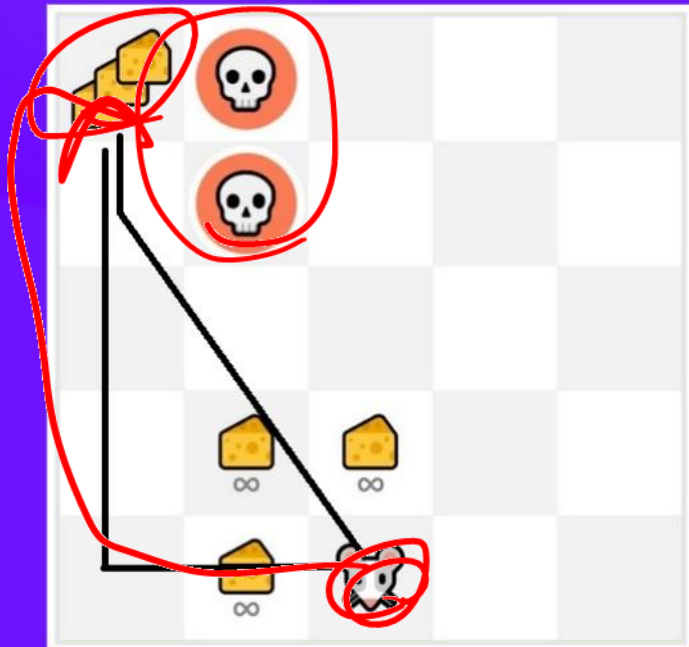deployment

Model-free

Model-based

$$p(s', r \mid s, a)$$

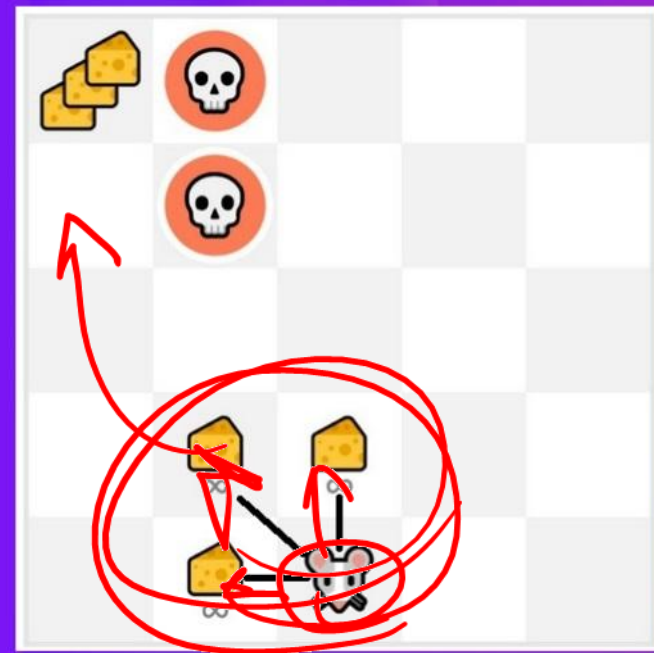# EXPLORATION VS. EXPLOITATION DILEMMA
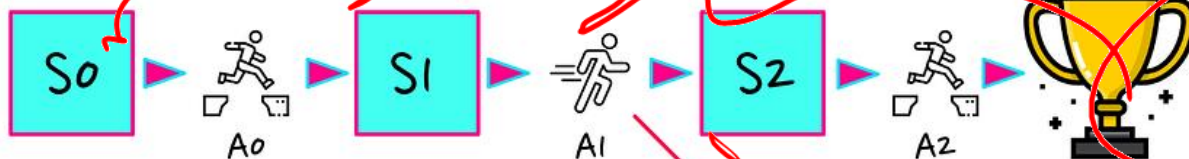
# Exploration/ Exploitation tradeoff

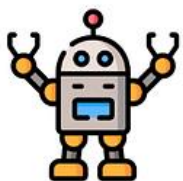*Exploration*: trying **random actions** in order to find **more information** about the environment.

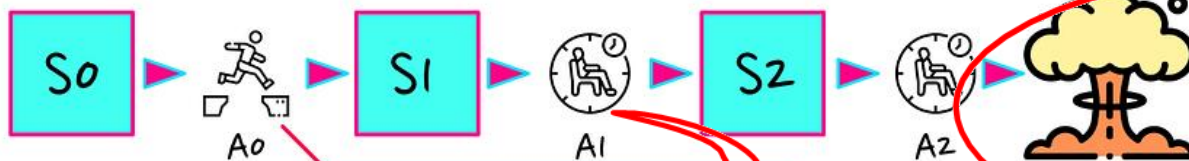*Exploitation*: using known information to **maximize the reward.**

# CREDIT ASSIGNMENT PROBLEM

# REWARD ENGINEERING PROBLEM



reward

IRL

what is the reward?

- Fly a helicopter      → **Reward**: air time, inverse distance, ...
- Manage an investment portfolio      → **Reward**: gains, gains minus risk, ...
- Control a power station      → **Reward**: efficiency, ...
- Make a robot walk      → **Reward**: distance, speed, ...
- Play video or board games      → **Reward**: win, maximise score, ...

If the goal is to learn via interaction, these are all reinforcement learning problems (Irrespective of which solution you use)

# GENERALIZATION PROBLEM



|  | Singleton Environments | IID Generalisation Environments | OOD Generlisation Environments |
|---|---|---|---|
| **Graphical Models** | MDP | CMDP | CMDP |
| **Train and Test Distribution** | $\bullet = \bullet$ <br> Train = Test | $p_{train}(c) = p_{test}(c)$ <br> Train Distribution = Test Distribution | $p_{train}(c) \neq p_{test}(c)$ <br> Train Distribution ≠ Test Distribution |
| **Example Benchmarks** | Atari <br> MuJoCo | OpenAI Procgen <br> Nethack Learning Environment <br> MiniHack | Distracting Control Suite <br> CausalWorld <br> CARLA |

# SAMPLE EFFICIENCY PROBLEM

Value-based

Policy-based

continuous action

# Two approaches to find optimal policy π*:

**Policy-Based methods**: train the agent to learn which **action to take,** given a state.

**Value-Based methods**: train the agent to learn which state **is more valuable** and take the action that **leads to it.**

# Two approaches to find optimal policy $\pi^*$:

*Policy-Based methods*:

- Train **directly the policy.**
- Our policy **is a Neural Network.**
- **No value function.**



State → π(State) → Action

# Two approaches to find optimal policy π*:

## *Value-Based methods*:

- Don't train the policy.
- **Our policy is a function defined by hand.**
- Instead **train a value-function that is a Neural Network.**



State → Q(State) → Values of State Action pair → π(State) → Action

$$\pi(s) = \arg\max_a Q_\pi(s, a)$$

Value based

Q learning

SARSA

Model-free

Actor
Critic

Policy-based

Policy Gradient

Dyna

Model-based

# WHAT WE HAVE LEARNED SO FAR?

- episodic vs continuing reinforcement learning
- offline vs online learning
- safe reinforcement learning
- on-policy vs off-policy vs offline reinforcement learning
- model-free vs model-base reinforcement learning
- exploration vs. exploitation dilemma
- credit assignment problem
- reward engineering problem
- generalization problem
- sample efficiency problem
- value-base vs policy-base vs actor-critic methods

MP, MRP, MDP

agent state $S_t^a$

observation $O_t$

action $A_t$

reward $R_t$

observation $O_t$

action $A_t$

reward $R_t$

environment state $S_t^e$

An information state (a.k.a. Markov state) contains all useful information from the history.

## Definition

A state $S_t$ is Markov if and only if

$$\mathbb{P}[S_{t+1} \mid S_t] = \mathbb{P}[S_{t+1} \mid S_1, ..., S_t]$$

Room 1 | Room 2 | Room 3

Goal

Outside

A Markov chain can be defined as a tuple of $(\mathcal{S}, \mathcal{P})$

- $\mathcal{S}$ is a finite set of states called the state space.

$$\mathcal{P} = \begin{array}{c|cccccc} & \text{Room 1} & \text{Room 2} & \text{Room 3} & \text{Outside} & \text{Found item} & \text{End} \\ \hline \text{Room 1} & 0.2 & 0.8 & 0 & 0 & 0 & 0 \\ \text{Room 2} & 0.2 & 0 & 0.4 & 0.4 & 0 & 0 \\ \text{Room 3} & 0 & 0.2 & 0 & 0 & 0.8 & 0 \\ \text{Outside} & 0 & 0.2 & 0 & 0.8 & 0 & 0 \\ \text{Found item} & 0 & 0 & 0 & 0 & 0 & 1.0 \\ \text{End} & 0 & 0 & 0 & 0 & 0 & 1.0 \end{array}$$

- Episode 1: (Room 1, Room 2, Room 3, Found item, End)
- Episode 2: (Room 3, Found item, End)
- Episode 3: (Room 2, Outside, Room 2, Room 3, Found item, End)
- Episode 4: (Outside, Outside, Outside, ...)

We can define the Markov reward process as a tuple $(\mathcal{S}, \mathcal{P}, \mathcal{R})$

- $\mathcal{S}$ is a finite set of states called the state space.
- $\mathcal{P}$ is the dynamics function (or transition model) of the environment, where $P(s'|s) = P\left[S_{t+1} = s' \mid S_t = s\right]$ specify the probability of environment transition into successor state $s'$ when in current state $s$.
- $\mathcal{R}$ is a reward function of the environment. $R(s) = \mathbb{E}\left[R_t \mid S_t = s\right]$ is the reward signal provided by the environment when the agent is in state $s$.

$p = 0.2$

$p = 0.2$

End

$r = 0$

$p = 1.0$

$p = 0.8$    Room 1    $p = 0.4$    Room 2    Room 3    $p = 0.8$    Found item

$r = -1$     $r = -1$     $r = -1$     $r = +10$

$p = 0.2$

$p = 0.4$

$p = 0.2$

Outside

$r = +1$

$p = 0.8$

- Episode 1: (Room 1, Room 2, Room 3, Found item, End)
  **Total rewards** $= -1 - 1 - 1 + 10 + 0 = 7.0$
- Episode 2: (Room 3, Found item, End)
  **Total rewards** $= -1 + 10 = 9.0$
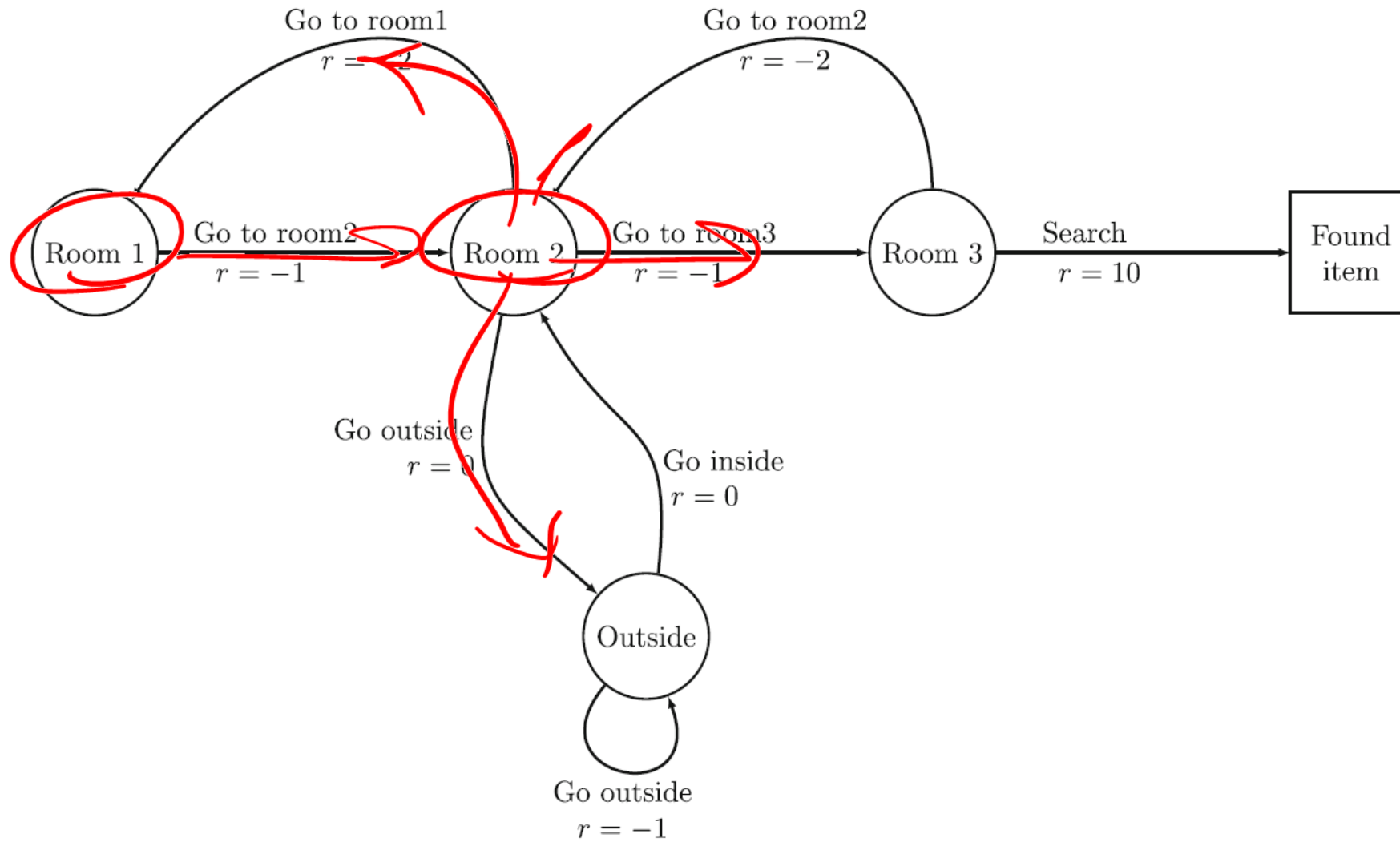- Episode 3: (Room 2, Outside, Room 2, Room 3, Found item, End)
  **Total rewards** $= -1 + 1 - 1 - 1 + 10 + 0 = 8.0$
- Episode 4: (Outside, Outside, Outside … )
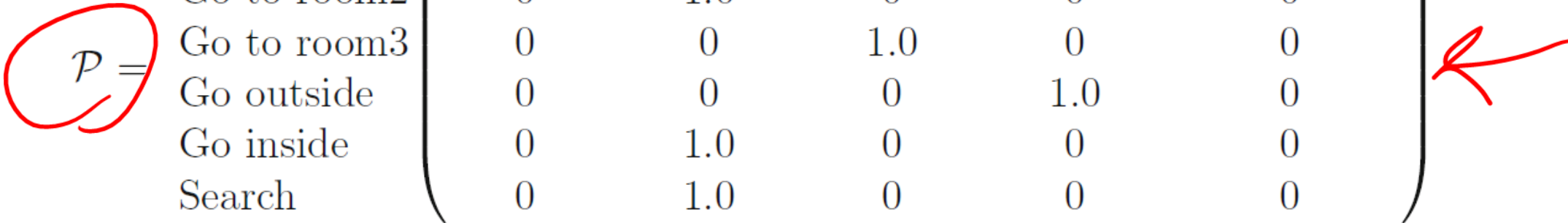  **Total rewards** $= 1 + 1 + \cdots = \infty$

We can define the MDP as a tuple $(S, A, P, R)$:

- $S$ is a finite set of states called the state space.
- $A$ is a finite set of actions called the action space.
- $P$ is the dynamics function (or transition model) of the environment, where $P(s'|s,a) = P\left[ S_{t+1} = s' \mid S_t = s, A_t = a \right]$ specify the probability of environment transition into successor state $s'$ when in current state $s$ and take action $a$.
- $R$ is a reward function of the environment; $R(s,a) = \mathbb{E}\left[ R_t \mid S_t = s, A_t = a \right]$ is the reward signal provided by the environment when the agent is in state $s$ and taking action $a$.

Go to room1
$r = -1$

Go to room2
$r = -2$

Room 1

Go to room2
$r = -1$

Room 2

Go to room3
$r = -1$

Room 3

Search
$r = 10$

Found
item

Go outside
$r = 0$

Go inside
$r = 0$

Outside

Go outside
$r = -1$

- $\mathcal{S} = \{$Room 1, Room 2, Room 3, Outside, Found item$\}$
- $\mathcal{A} = \{$Go to room1, Go to room2, Go to room3, Go outside, Go inside, Search$\}$
- $\mathcal{R} = \{-1, -2, +1, 0, +10\}$

$$\mathcal{P} = \begin{array}{c} \\ \text{Go to room1} \\ \text{Go to room2} \\ \text{Go to room3} \\ \text{Go outside} \\ \text{Go inside} \\ \text{Search} \end{array} \begin{array}{ccccc} \text{Room 1} & \text{Room 2} & \text{Room 3} & \text{Outside} & \text{Found item} \\ \left( \begin{array}{ccccc} 1.0 & 0 & 0 & 0 & 0 \\ 0 & 1.0 & 0 & 0 & 0 \\ 0 & 0 & 1.0 & 0 & 0 \\ 0 & 0 & 0 & 1.0 & 0 \\ 0 & 1.0 & 0 & 0 & 0 \\ 0 & 1.0 & 0 & 0 & 0 \end{array} \right) \end{array}$$

$$\mathcal{P} = \begin{array}{c} \\ \text{Go to room1} \\ \text{Go to room2} \\ \text{Go to room3} \\ \text{Go outside} \\ \text{Go inside} \\ \text{Search} \end{array} \begin{array}{ccccc} \text{Room 1} & \text{Room 2} & \text{Room 3} & \text{Outside} & \text{Found item} \\ \left( \begin{array}{ccccc} 0.6 & 0 & 0 & 0.4 & 0 \\ 0 & 1.0 & 0 & 0 & 0 \\ 0 & 0 & 0.2 & 0.8 & 0 \\ 0 & 0 & 0 & 1.0 & 0 \\ 0 & 1.0 & 0 & 0 & 0 \\ 0 & 1.0 & 0 & 0.0 & 0 \end{array} \right) \end{array}$$
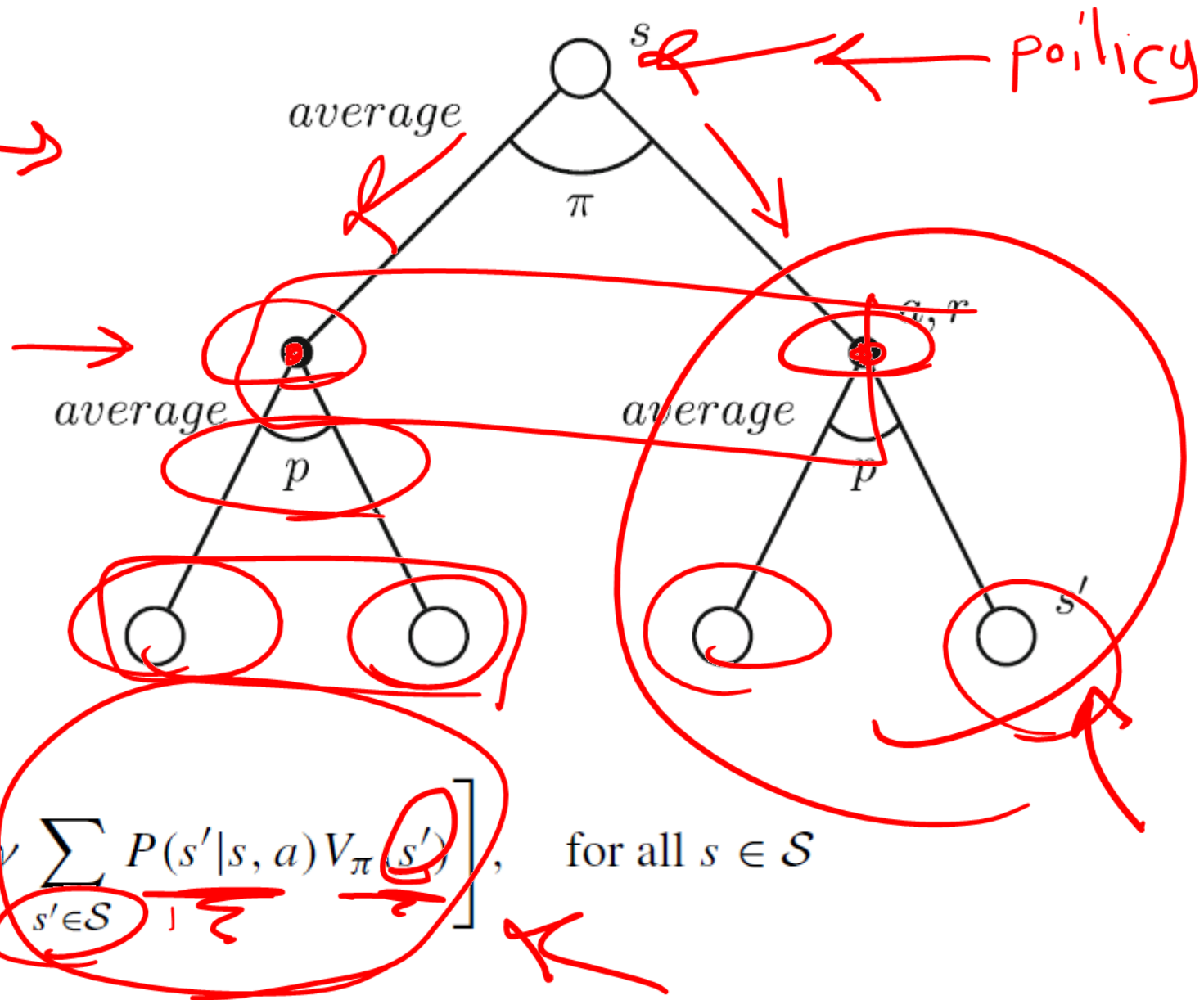
$$V_\pi(s) = \mathbb{E}_\pi\left[G_t \mid S_t = s\right], \quad \text{for all } s \in \mathcal{S}$$

$$Q_\pi(s, a) = \mathbb{E}_\pi\left[G_t \mid S_t = s, A_t = a\right], \quad \text{for all } s \in \mathcal{S}, a \in A$$
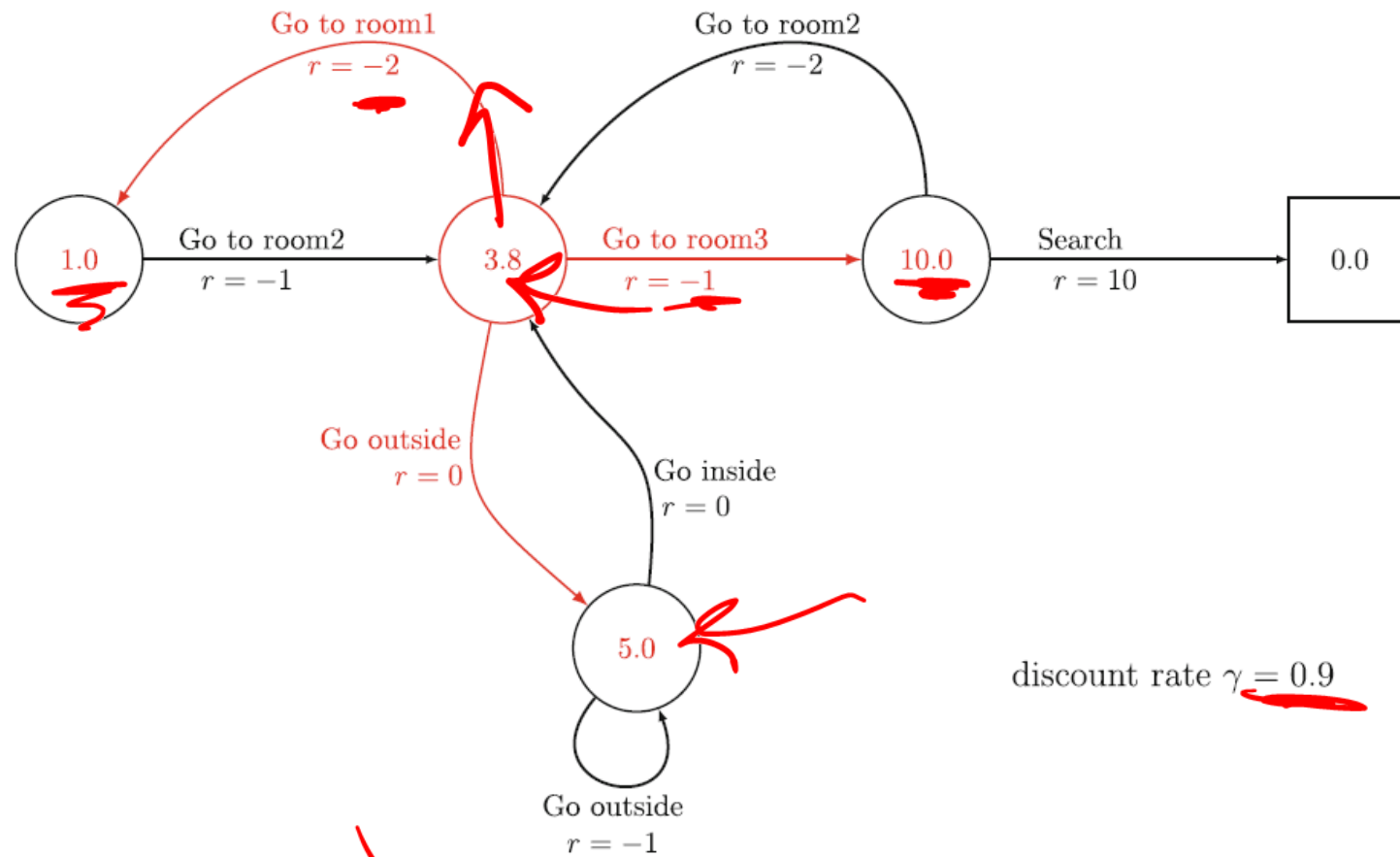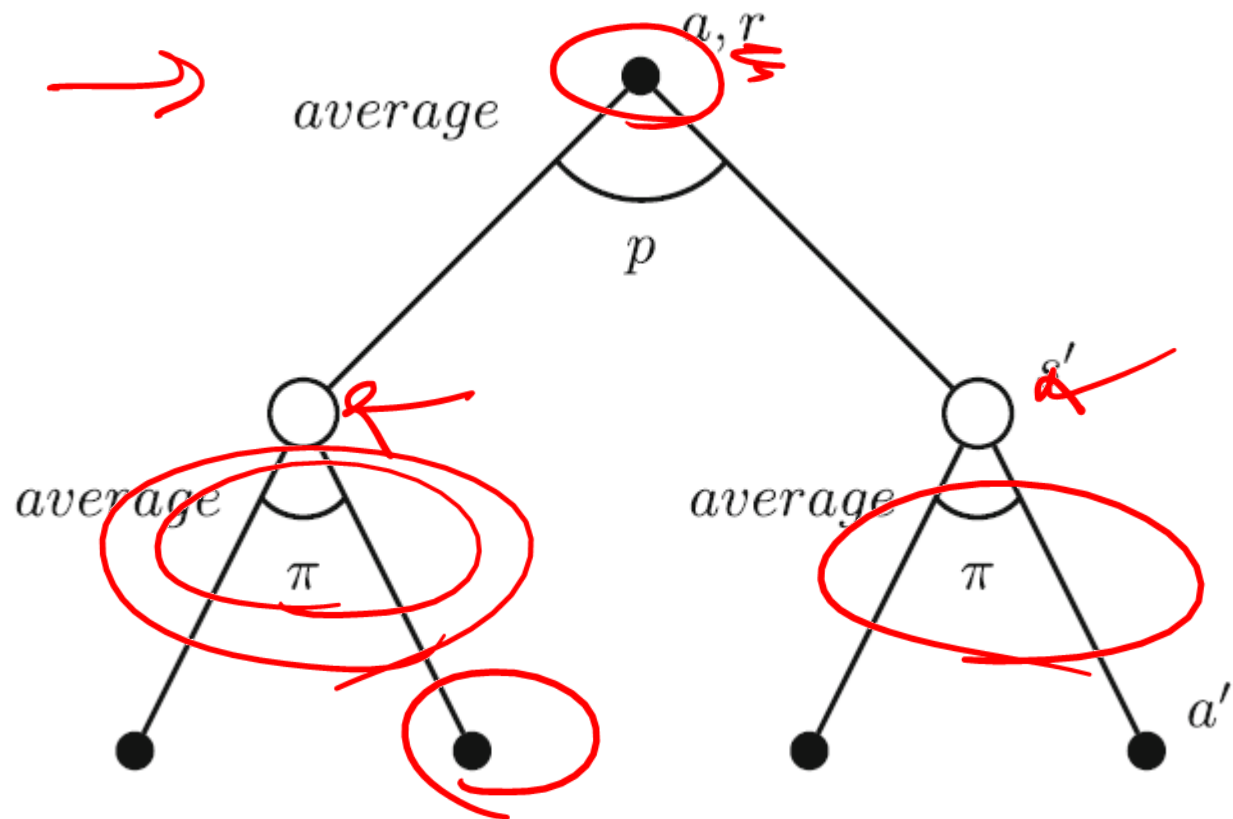
Backup
Diagram

average

$\pi$

$s$

poilicy

env

average

$p$

average

$p$

$s, r$

$s'$

$V_\pi(s) = \mathbb{E}_\pi\left[ G_t \mid S_t = s \right]$

$= \sum_{a \in A} \pi(a|s) \left[ R(s, a) + \nu \sum_{s' \in \mathcal{S}} P(s'|s, a) V_\pi(s') \right], \quad \text{for all } s \in \mathcal{S}$
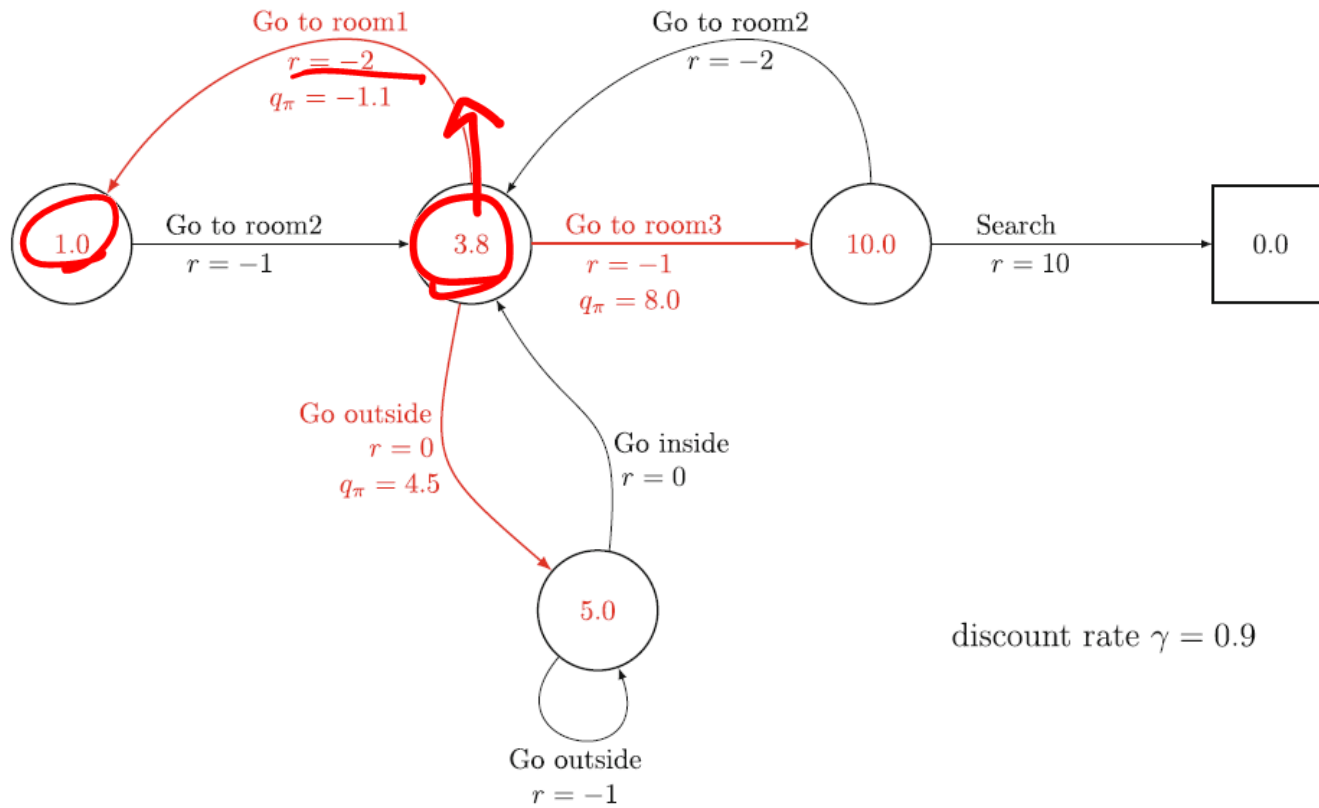
Go to room1
$r = -2$

Go to room2
$r = -2$

1.0

Go to room2
$r = -1$

3.8

Go to room3
$r = -1$

10.0

Search
$r = 10$

0.0

Go outside
$r = 0$

Go inside
$r = 0$

5.0

Go outside
$r = -1$

discount rate $\gamma = 0.9$

$$V_\pi(Room\ 2) = 0.33 * (-2 + 0.9 * 1.0) + 0.33 * (-1 + 0.9 * 10.0)$$
$$+ 0.33 * (0 + 0.9 * 5.0)$$
$$= 0.33 * -1.1 + 0.33 * 8 + 0.33 * 4.5$$
$$= 3.76$$

env $\longrightarrow$



$$Q_\pi(s,a) = \mathbb{E}_\pi\left[ G_t \,\middle|\, S_t = s, A_t = a \right]$$

$$= R(s,a) + \gamma \sum_{s' \in \mathcal{S}} P(s'|s,a) \sum_{a' \in A} \pi(a'|s') Q_\pi(s',a'), \quad \text{for all } s \in \mathcal{S}, a \in A$$

$$Q_\pi(Room\ 2,\ Go\ to\ room1) = -2 + 0.9 * 1.0$$
$$= -1.1$$

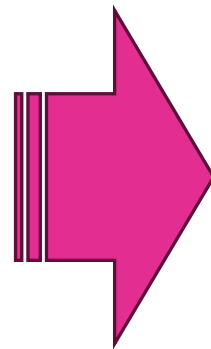$$Q_\pi(Room\ 2,\ Go\ to\ room3) = -1 + 0.9 * 10.0$$
$$= 8.0$$

$$Q_\pi(Room\ 2,\ Go\ outside) = 0 + 0.9 * 5.0$$
$$= 4.5$$

$$V_\pi(Room\ 2) = 0.33 * -1.1 + 0.33 * 8 + 0.33 * 4.5$$
$$= 3.76$$

# WATCH THE FOLLOWING VIDEO



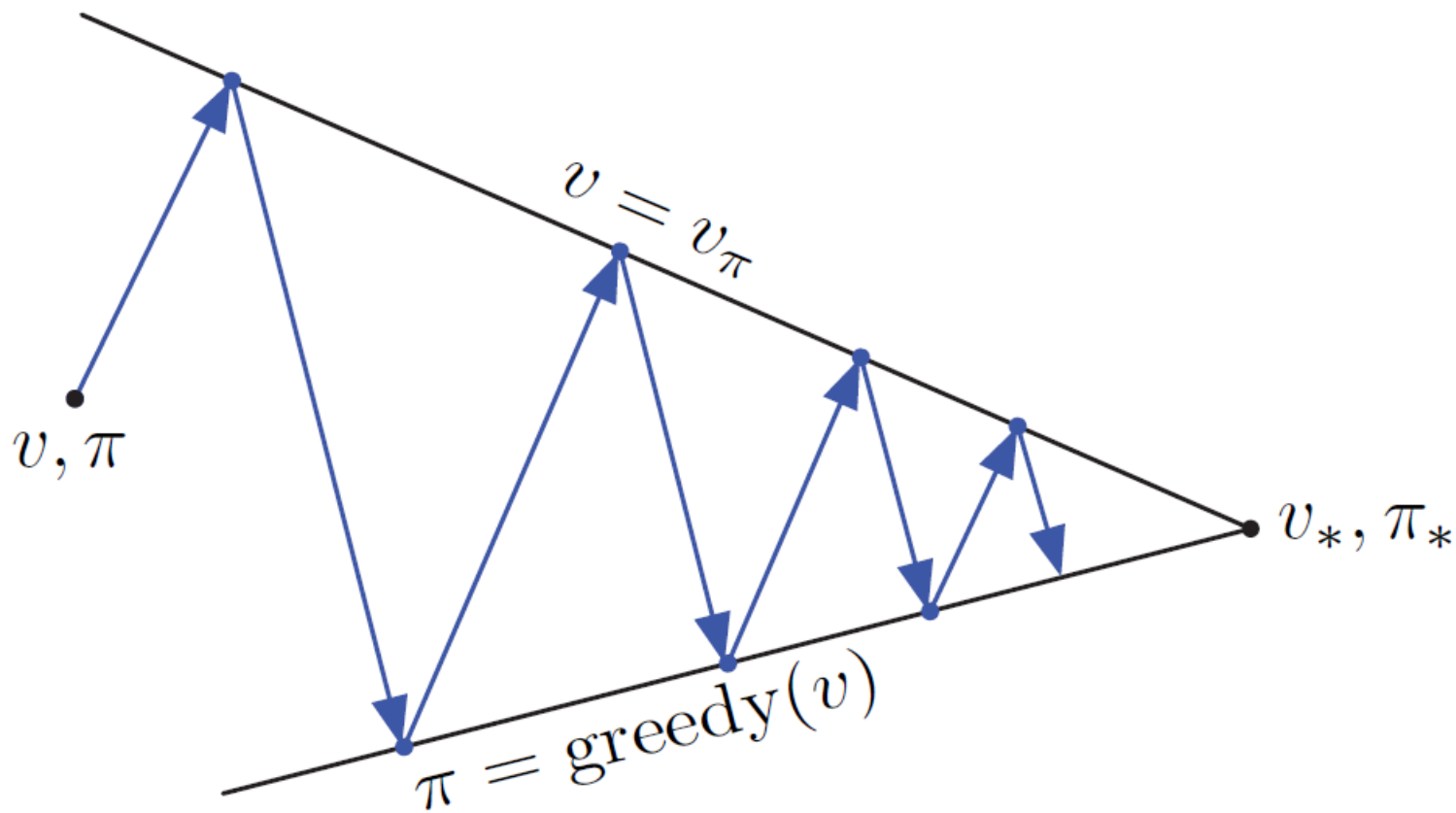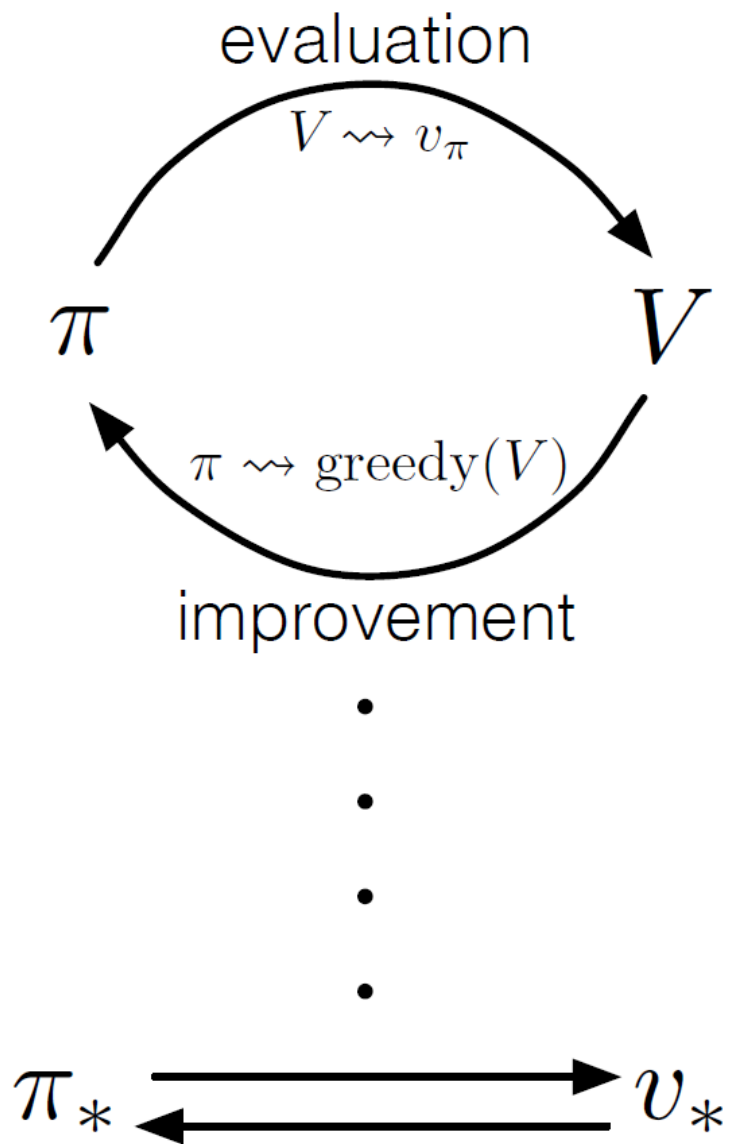https://www.youtube.com/watch?v=NFo9v_yKQXA

# How to solve
# full RL problem?

# When we have:

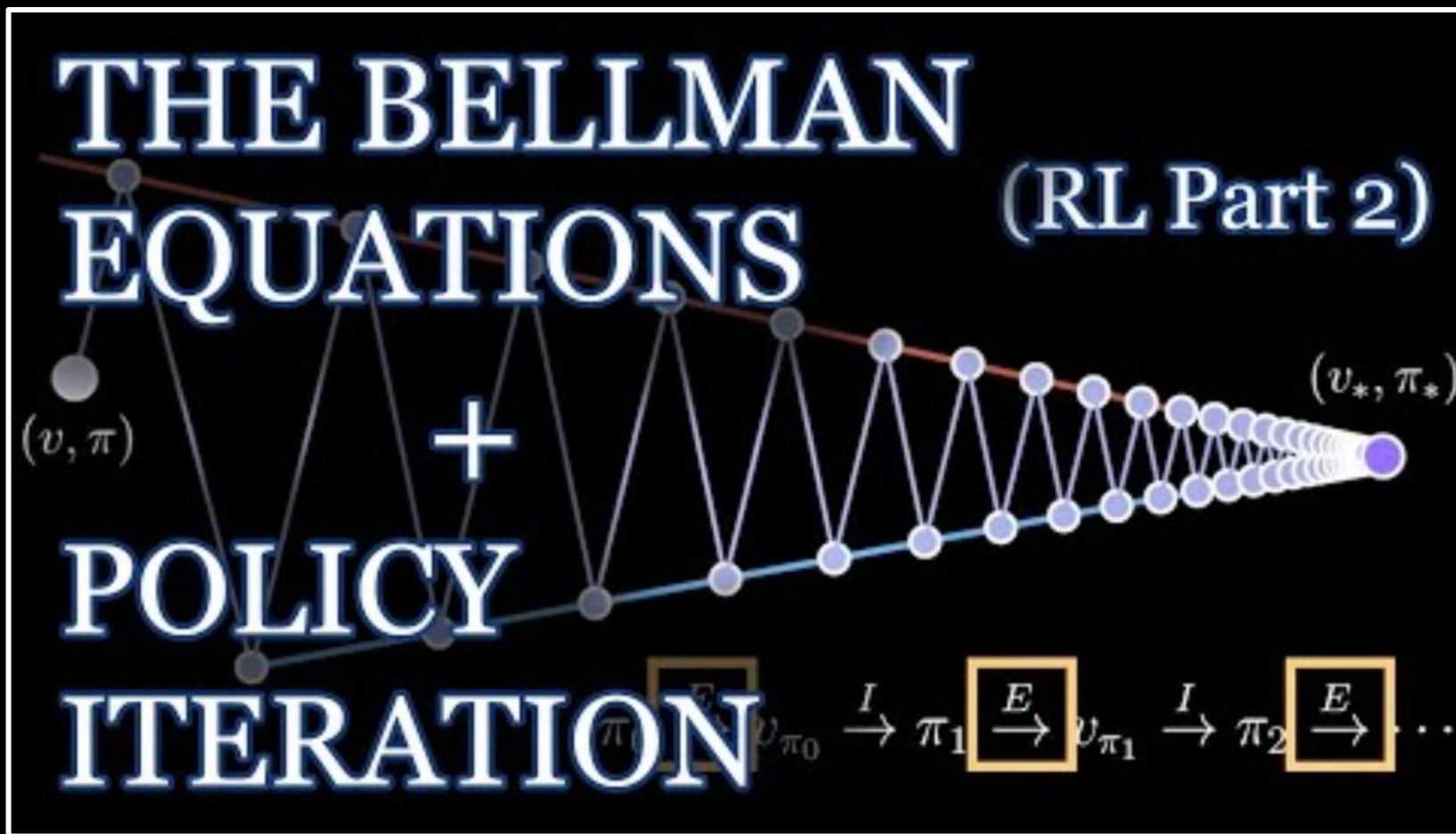$$P(s', r|s, a) = \mathbb{P}[S_{t+1} = s', R_{t+1} = r|S_t = s, A_t = a]$$

# OPTIMAL VALUE AND POLICY

$$Q_*(s, a) = \max_\pi Q_\pi(s, a), \quad \text{for all } s \in \mathcal{S}, a \in \mathcal{A}$$

$$\pi_*(a|s) = \begin{cases} 1, & \text{if} \quad a = \underset{a \in A}{\arg\max} \, Q_*(s, a) \\ 0, & \text{otherwise} \end{cases}$$

evaluation

$V \rightsquigarrow v_\pi$

$\pi$

$V$

$\pi \rightsquigarrow \text{greedy}(V)$

improvement

$\pi_*$     $v_*$

$v, \pi$

$v = v_\pi$

$v_*, \pi_*$

$\pi = \text{greedy}(v)$

# WATCH THE FOLLOWING VIDEO



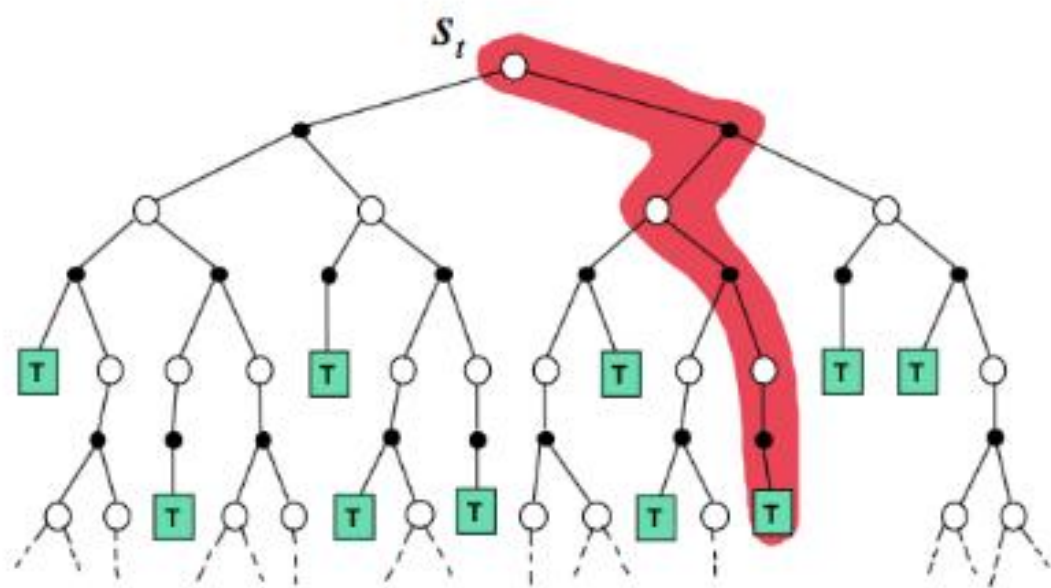https://www.youtube.com/watch?v=_j6pvGEchWU

# When we don't have:

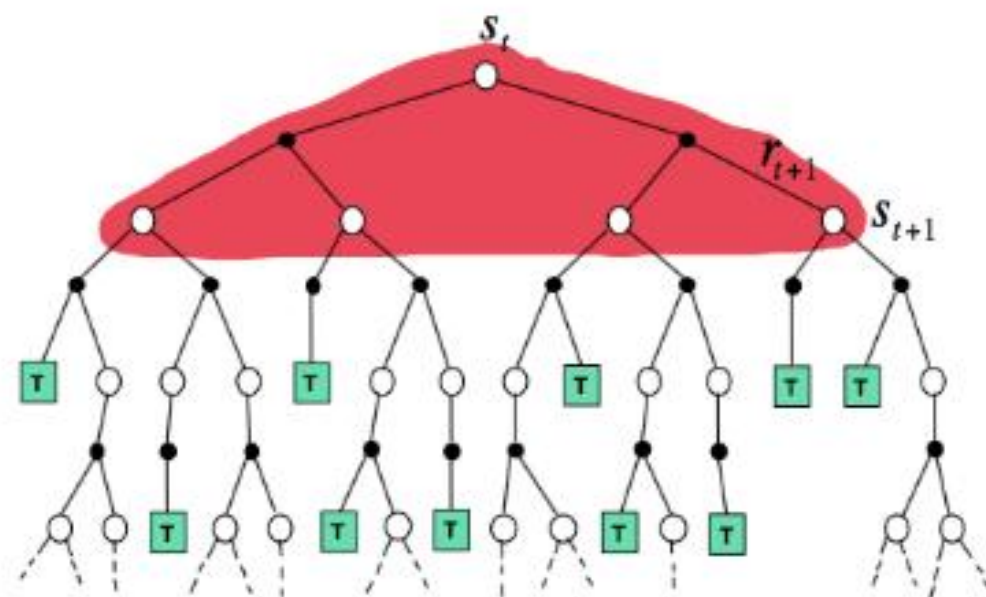$$\cancel{P(s', r|s, a) = \mathbb{P}[S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a]}$$

## Monte-Carlo

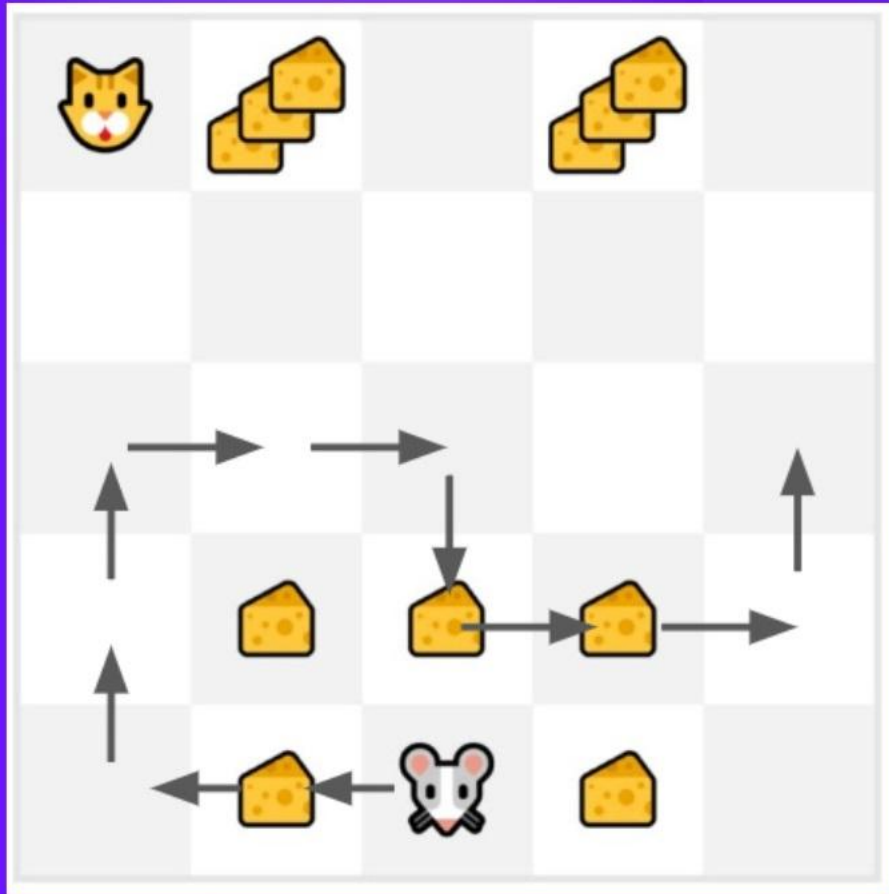$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$

## Dynamic Programming

$$V(S_t) \leftarrow \mathbb{E}_\pi \left[ R_{t+1} + \gamma V(S_{t+1}) \right]$$

# Monte Carlo Approach:

- **Calculate the return Gt.**

Gt = Rt+1 + Rt+2 + Rt+3...

Gt = 1 + 0 + 0 + 0+ 0 + 0 + 1 + 1+ 0 + 0

Gt= 3

- **We can now update V(S0).**
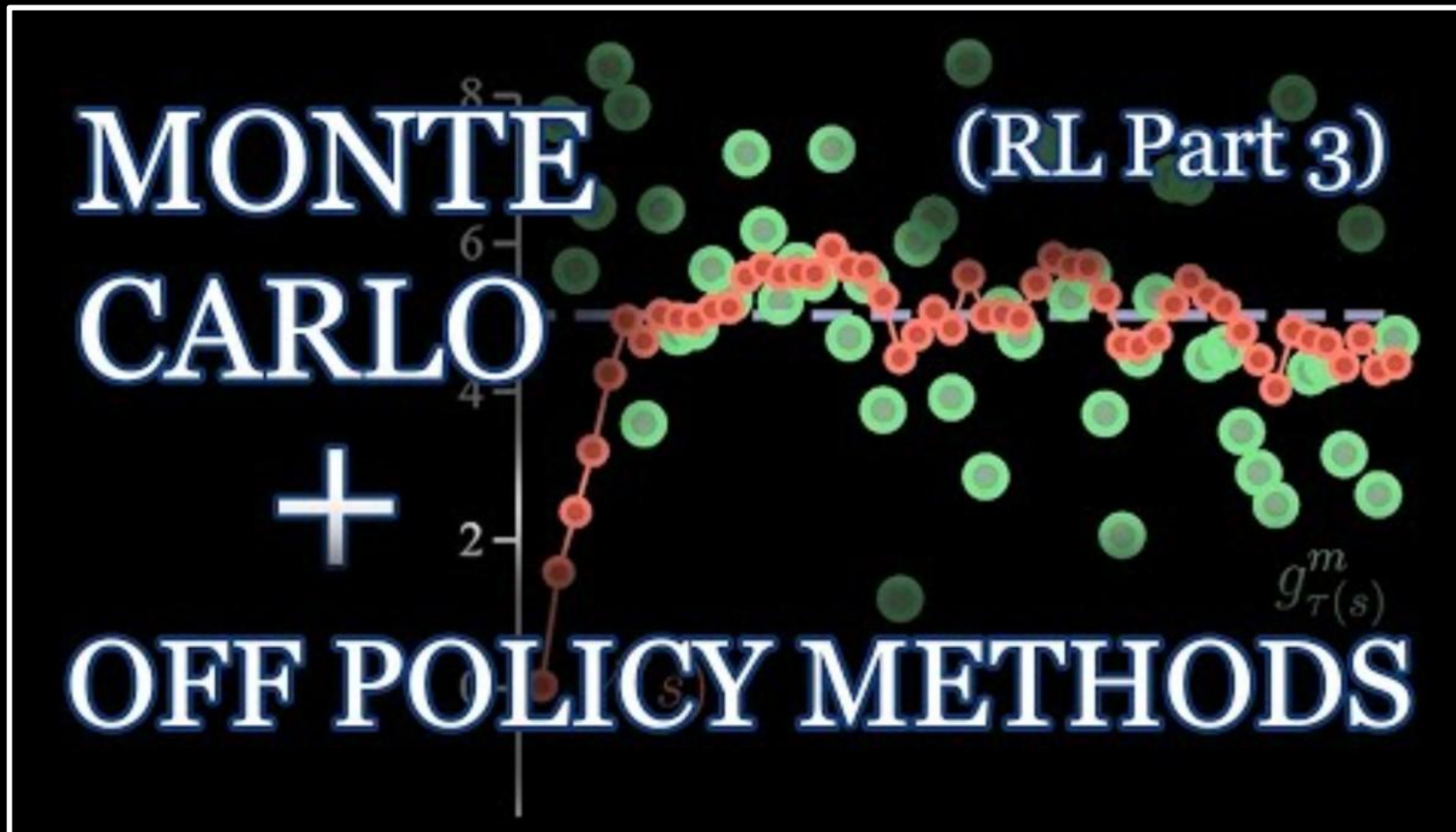
$$V(S_t) \leftarrow V(S_t) + \alpha[G_t - V(S_t)]$$

New V(S0) = V(S0) + lr * [Gt-V(S0)]

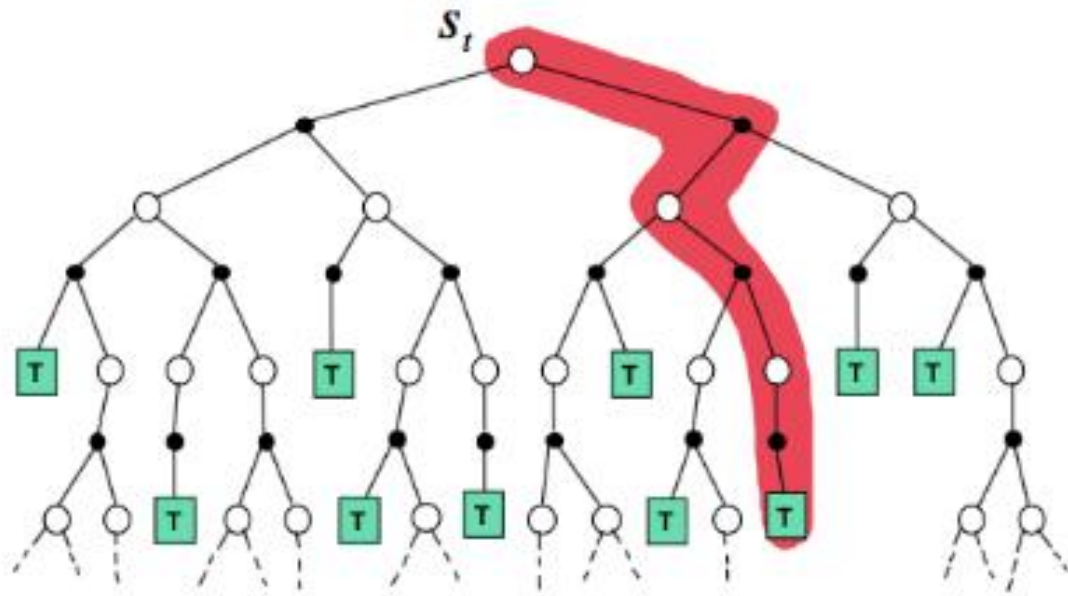New V(S0) = 0 + 0.1 * [3 –0]

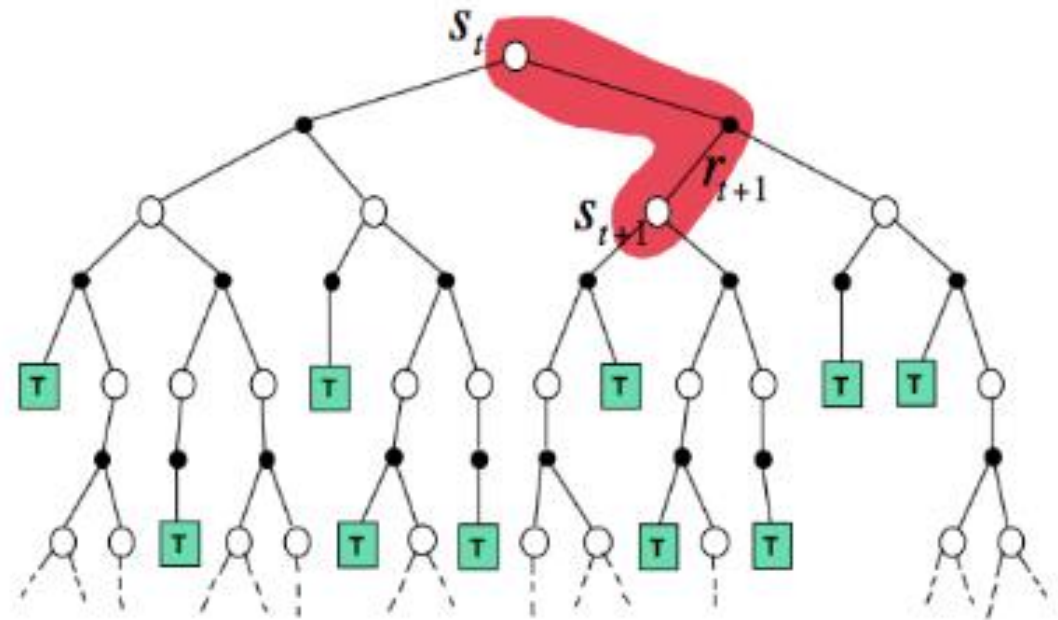New V(S0) = 0.3

# WATCH THE FOLLOWING VIDEO

## Monte-Carlo

$$V(S_t) \leftarrow V(S_t) + \alpha \left( G_t - V(S_t) \right)$$

## Temporal-Difference

$$V(S_t) \leftarrow V(S_t) + \alpha \left( R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \right)$$

https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_td.html

# TD Learning Approach:

*Temporal Difference Learning:* learning at **each time step.**

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

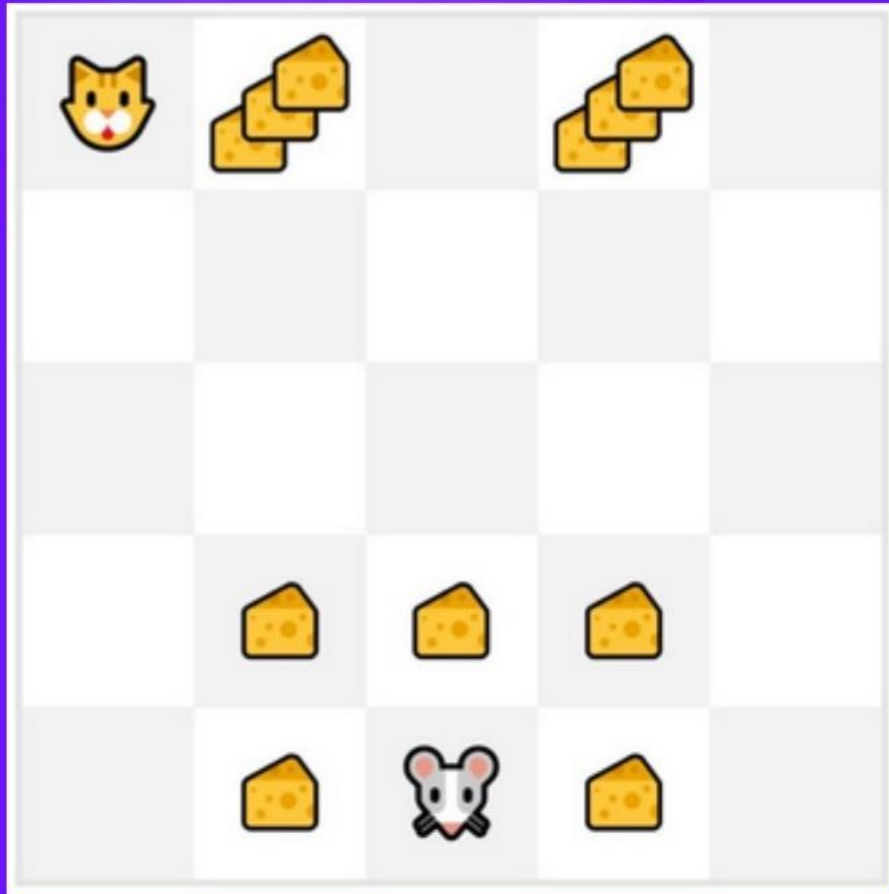New value of state t

Former estimation of value of state t

Learning Rate

Reward

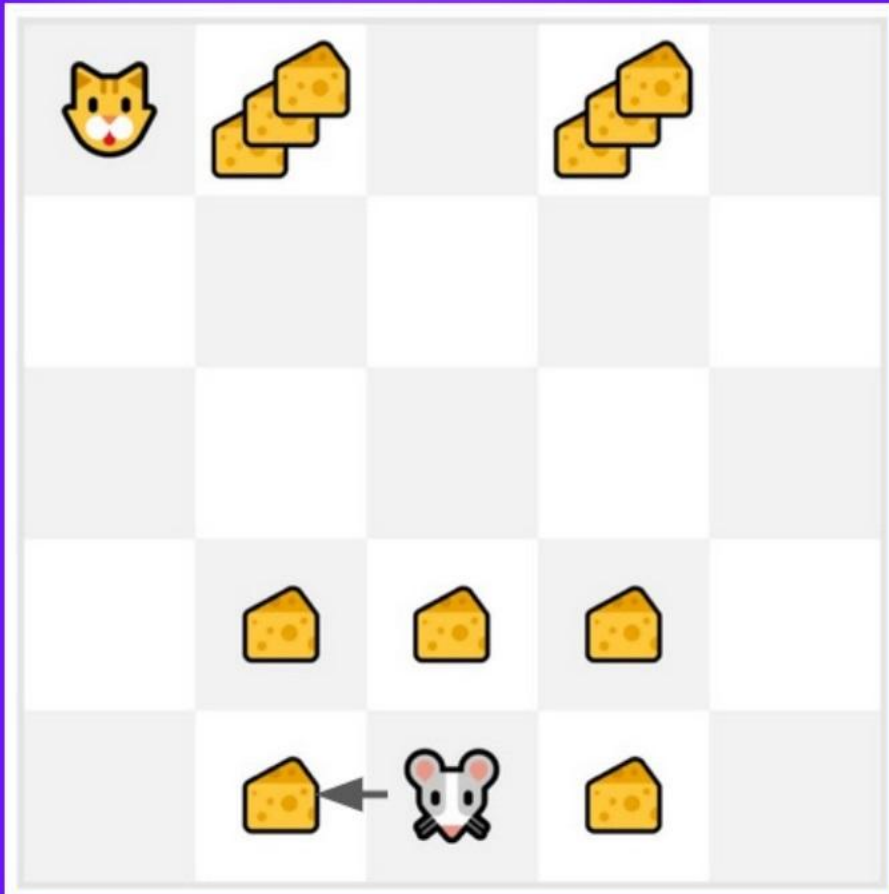Discounted value of next state

TD Target

# TD Approach:



**At the end of one step** (State, Action, Reward, Next State):
- We have Rt+1 and St+1
- We update V(St):
  - **We estimate Gt** by adding Rt+1 and the discounted value of next state.
  **TD target** : Rt+1 + gamma * V(St+1)

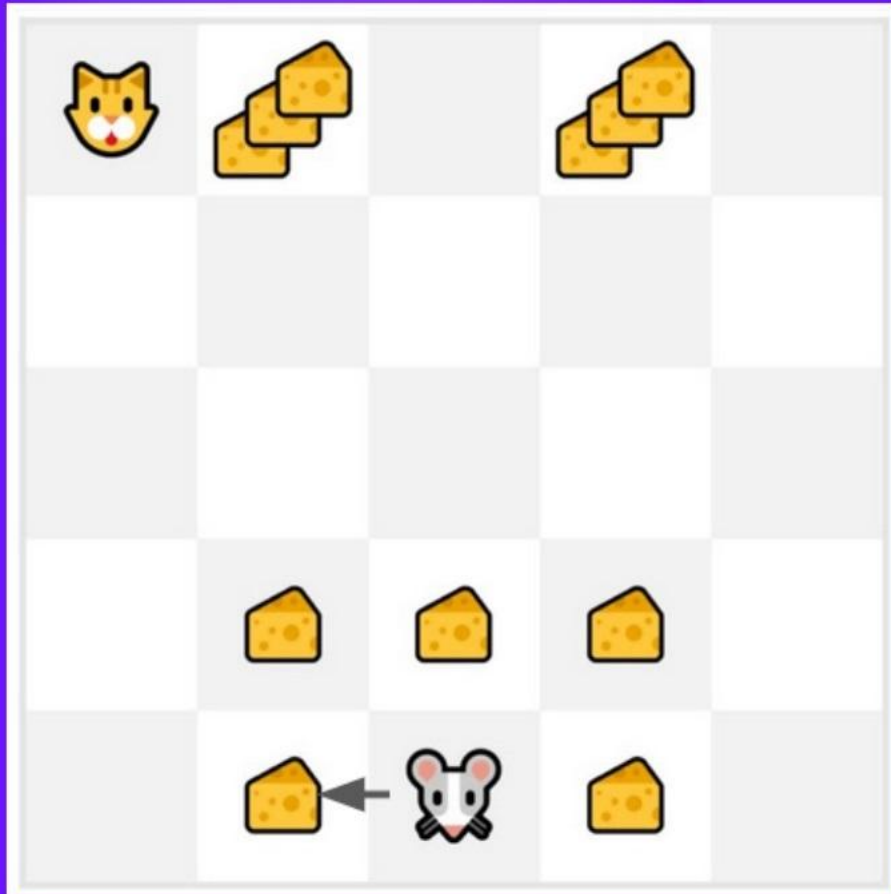$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Now we **continue to interact with this environment** with our updated value function. By running more and more steps, **the agent will learn to play better and better.**

# TD Approach:



- We just started to train our Value function **so it returns 0 value for each state.**

- Learning rate (lr) is 0.1 and our discount rate is 1 **(no discount)**

- Our mouse, **explore the environment** and take a random action: going left.

- It gets **a +1 reward (cheese).**

# TD Approach:



- **We can now update V(S0):**

$$V(S_t) \leftarrow V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

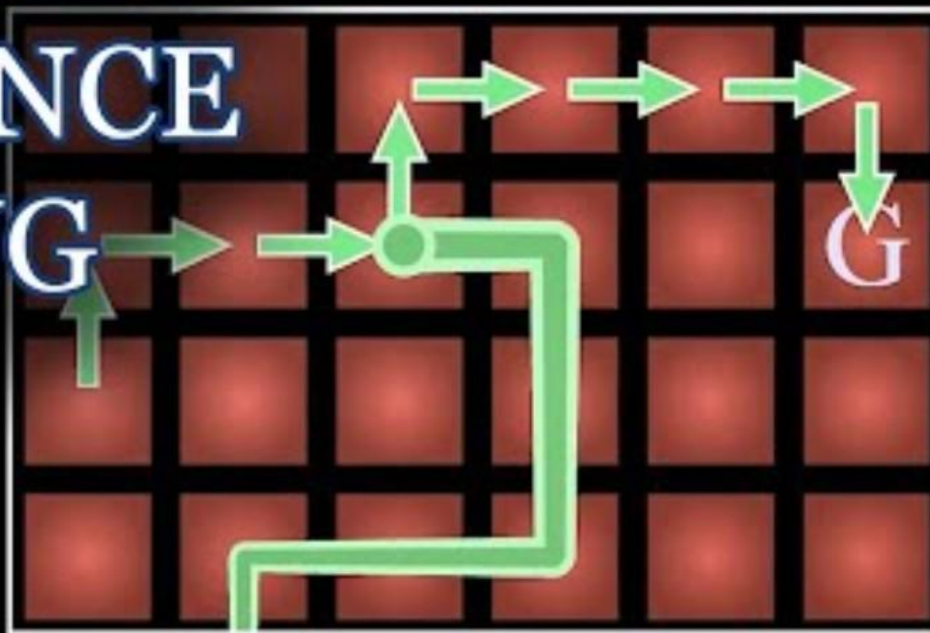New V(S0) = 0 + 0.1 * [1 + 1 * 0–0]
The new V(S0) = 0.1

So we just updated our **value function for State 0.**

Now **we continue to interact with this environment with our updated value function.**

# WATCH THE FOLLOWING VIDEO

Temporal-difference learning

Dynamic programming

width of update

depth (length) of update

Monte Carlo

Exhaustive search