# Reinforcement Learning

Agent

State  Reward  Action
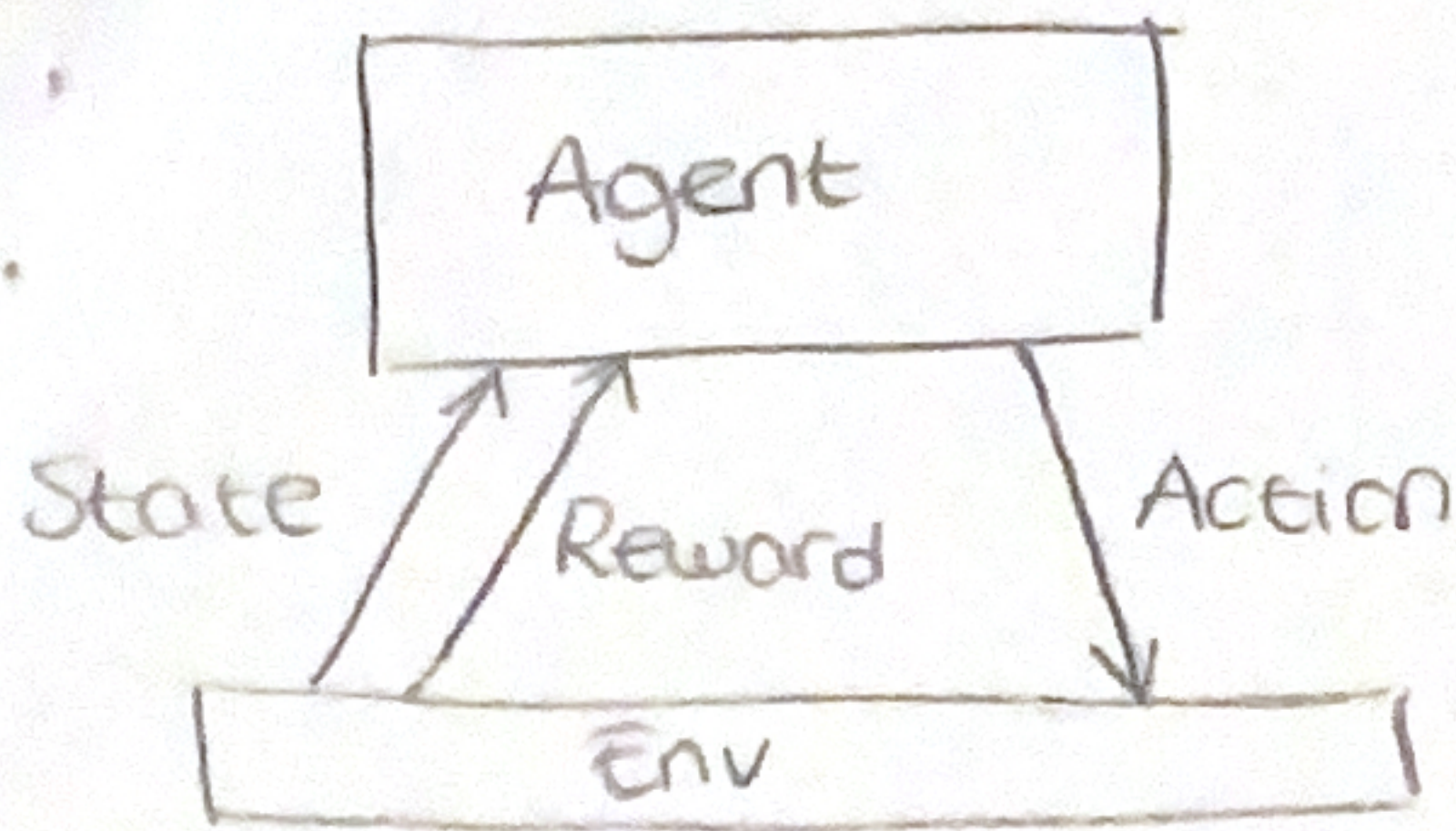
Env

$$s_0 \xrightarrow[r_0]{a_0} s_1 \xrightarrow[r_1]{a_1} \ldots$$

Agent chooses action $a_i$ in state $s_i$ and gets reward $i$. But goal is to maximize:

$$\underbrace{r_0 + \gamma r_1 + \gamma^2 r_2 + \ldots}_{} \quad 0 \le \gamma < 1$$

Reward on the long term
$\gamma$: discount factor
(a parameter on how much we care about immediate & future rewards)

- Learn a control policy
$$\pi: \underset{\substack{\downarrow \\ \text{set of states}}}{S} \to \underset{\substack{\downarrow \\ \text{set of actions}}}{A}$$

- Take $a$ from $A$ given current state $s$ from $S$

- **Problems**
  - Delayed Reward & Temporal Credit Assignment
  Determine which of the actions in its sequence are to be credited for eventual rewards.
  - Exploration vs Exploitation
  Trade-off in choosing exploration of unknown states & actions (more info) or exploitation of states & actions that are known to yield reward.

## The Learning Task
### Markov Decision Process

$S \to$ set of states
$A \to$ set of actions

At each $t_i$, in $s_{t_i}$ perform $a_{t_i}$, get $r(s_{t_i}, a_{t_i})$; produce $s_{t+1} = \delta(s_t, a_t)$

These only depend on current state.

- $\pi: S \to A$
  $\downarrow$
  learn this

- $V^\pi(s_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \ldots$
  $\downarrow$
  $= \sum_{i=0}^{\infty} \gamma^i r_{t+i}$
  cumulative reward  $0 \le \gamma < 1$

  - When $\gamma = 0 \to$ only immediate reward is considered
  - When $\gamma \to 1 \to$ Future rewards are more important

Optimal policy is:

$$\pi \equiv \arg\max_\pi V^\pi(s)$$
$\downarrow$
pick the policy with the most value