

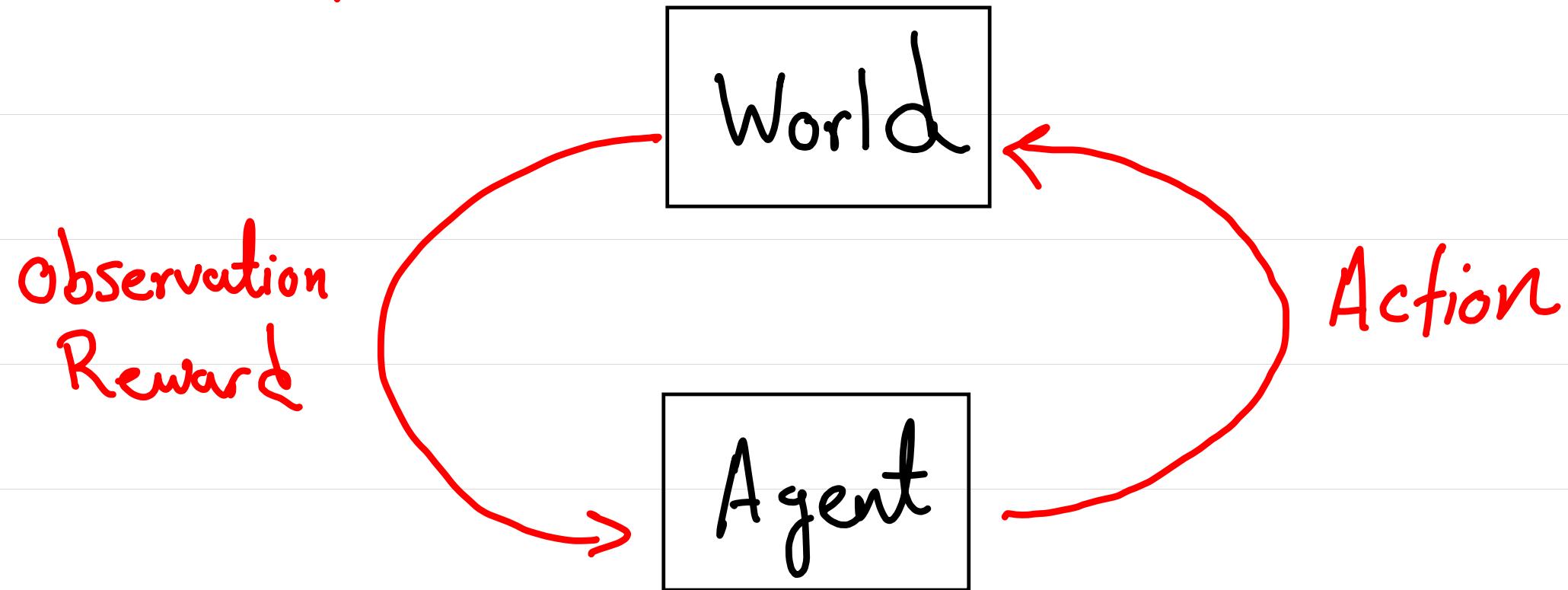
Week 11 - Part 2

■ Outcome of this part

- Sequential Decision process
- Observation, history, and state
- Markov Decision Process (MDP)
 - What is Markov about MDP? Why Markov assumption is common?
- Dynamics Model & Reward Model
 - Transition Graph
- Return / Utility
- Policies
- Finding the "best" Policy (i.e., Solving the MDP)

} ← Next part

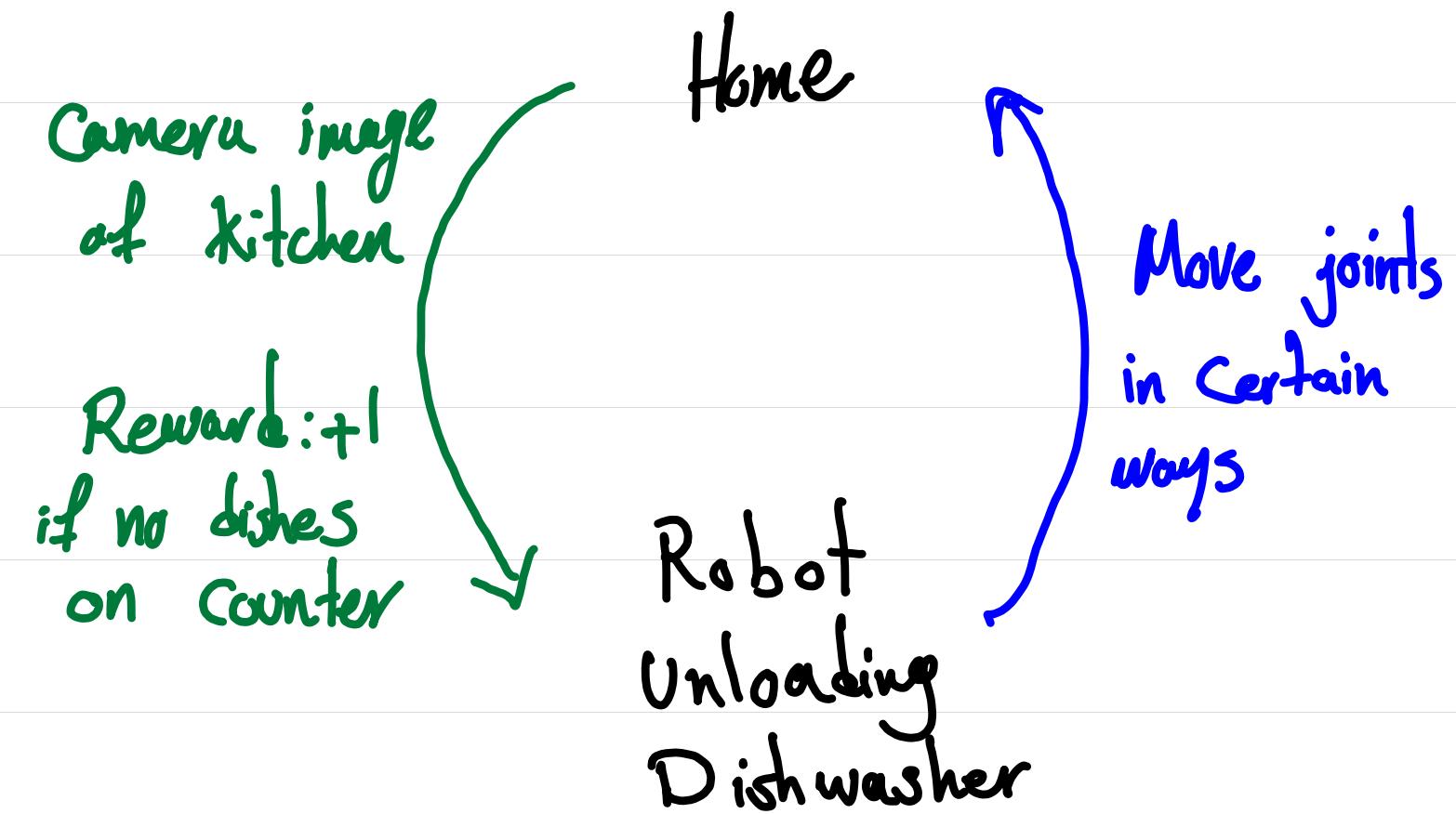
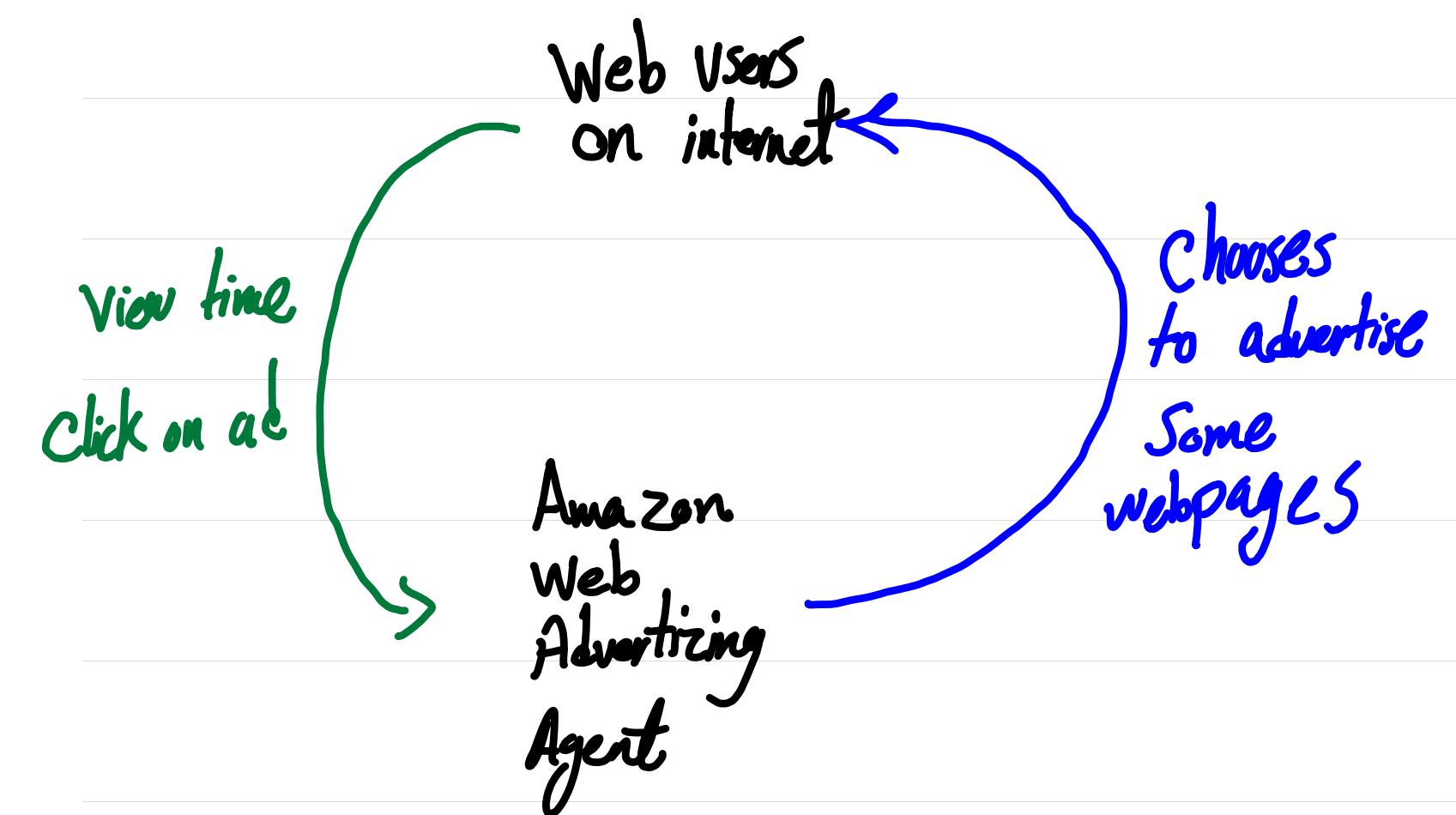
Sequential Decision Making



■ Goal: Select actions to maximize total expected future reward

↳ May require balancing immediate & long-term rewards

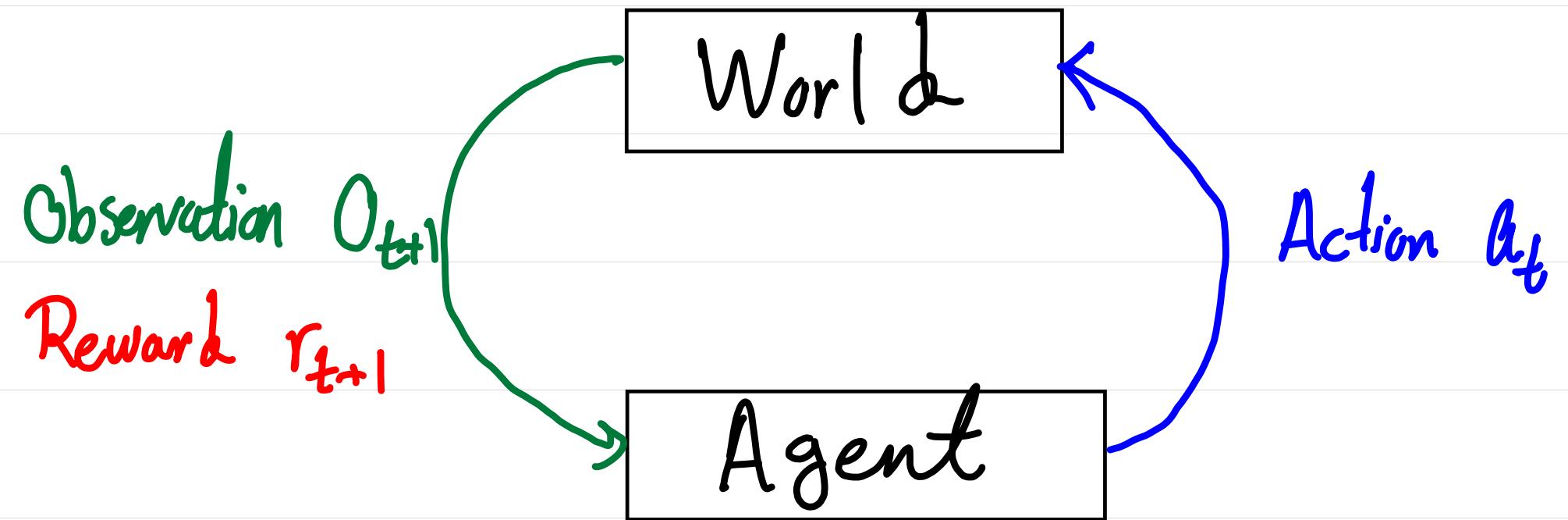
Some Examples



The robot can break all dishes by throwing them on the floor.

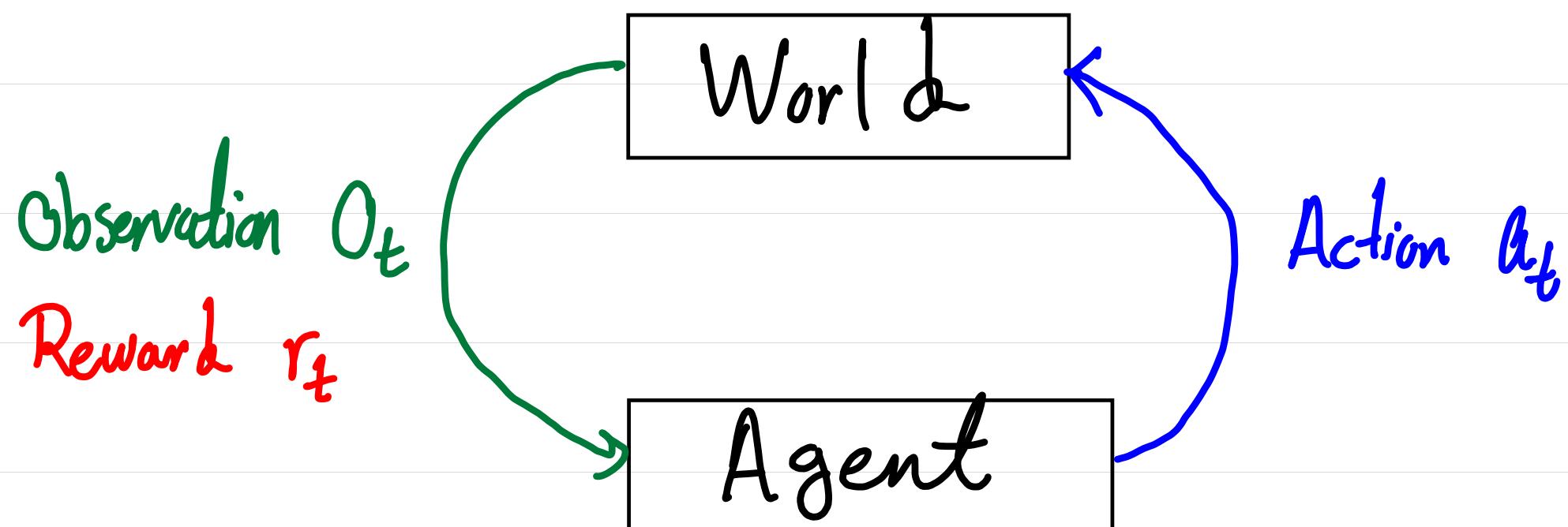
This issue is known as
Reward Hacking

Sequential Decision Process: Agent & the World (Discrete Time)



- Each time step t :
 - Agent take an action a_t
 - Given action a_t , the world updates and emits observation o_{t+1} and reward r_{t+1}
 - the agent receives o_{t+1} and r_{t+1}

History: Sequence of Past Observations, Actions, & Rewards



- History: $h_t = (a_0, O_1, r_1, a_1, O_2, r_2, \dots, a_{t-1}, O_t, r_t)$
- Agent chooses action based on history
- State: information assumed to determine what happens next
 - ↳ State is a Function of history: $S_t = \Psi(h_t)$
 - the simplest State function is $S_t = h_t$
 - $S_t = O_t$

Observation V.s. History V.s State

■ Can anyone tell me what is the difference between observation, history, and state?

43

■ Consider the Atari game. How would you design/choose observation, history, and state?

Observation is the last image snapshot of the game

State would be defined as the last few Observation

- This allows us to find the speed & direction of the ball

Markov Assumption

- Often, to make problems tractable, and because it is not a terrible assumption in reality, we will make **Markov Assumption**.

$$P(S_{t+1} | S_t, a_t) = P(S_{t+1} | S_t, a_t, S_{t-1}, a_{t-1}, S_{t-2}, a_{t-2}, \dots, S_1, a_1)$$

- Future (S_{t+1}) is independent of past ($S_{t-1}, a_{t-1}, \dots, S_1, a_1$) given the present (S_t, a_t)

with markov assumption

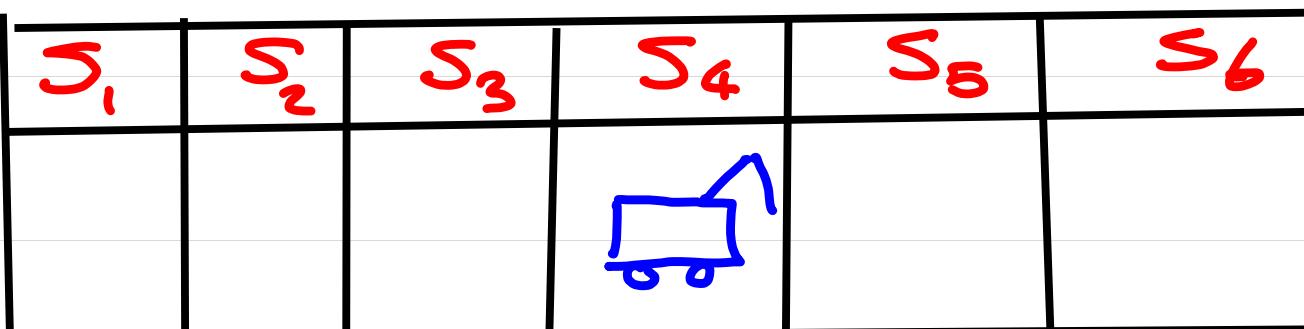
$$P(S_{t+1}, S_t, \dots, S_0 | a_t, \dots, a_0) = \frac{P(S_{t+1} | S_t, S_{t-1}, \dots, S_1, a_t, \dots, a_1) \times P(S_t | S_{t-1}, S_{t-2}, \dots, a_t, \dots, a_1) \times \dots \times P(S_1 | S_0)}{P(S_0)}$$

Why is Markov Assumption Popular?

- It is simple and often can be satisfied if includes some previous observations as part of state
- There are many problems that we can satisfy Markov assumption by simply setting $S_t = O_t$.
- Markov assumption significantly simplifies working with the problem.
- Just like all ML problems we studied so far, there is this trade-off between the expressive power of state representation and how long it takes to train your model.

MDP Example: Mars Rover

locations on Mars →



- Mars Rover is the robot that you designed to collect samples on Mars.
- Your robot has to choose from going to right or left, to move between locations.
- Mars is strangely small. Only 6 locations on Mars.

State Set: $S = \{s_1, \dots, s_6\}$

Action Set: $A = \{\text{"Try Left"}, \text{"Try Right"}, \text{"exit"}\}$

Rewards: reward 1 and 5 when you get to s_1 and s_6 , respectively.
reward 0, otherwise.

s_1 and s_6 are the terminal locations, you obtain 1 and 5 samples from Soil in locations s_1 and s_6 , respectively, and then the robot will shut down

Dynamics Model & Reward Model

- Transition / dynamics model: predicts next state
- In other words, dynamics model specifies the distribution of outcomes when the agent makes a decision

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

In this course, we assume stationary transition / dynamics model.
i.e., $P[S_{t+1} = s' | S_t = s, A_t = a] = P[S_{K+1} = s' | S_K = s, A_K = a]$
for all t and K .

- Reward model: predicts immediate reward: $R(s, a, s')$
 - Sometimes just $R(s)$, $R(s, a)$, or $R(s')$, depending on the model
 - Note that rewards can be zero, negative, or positive.

Example: Mars Rover Dynamics Model (Reward Model)

State →

reward →

s_1	s_2	s_3	s_4	s_5	s_6

States & Actions: $S = \{s_1, \dots, s_6\}$, $A = \{a_L, a_R, a_e\}$

Rewards Model:

$$R(s, a, s') = \begin{cases} 5 & \text{if } s = s_5, a \in \{a_L, a_R\}, s' = s_6 \\ 1 & \text{if } s = s_2, a \in \{a_L, a_R\}, s' = s_1 \\ 0 & \text{Otherwise} \end{cases}$$

$$R(s_6, a_e, s_6) = R(s_1, a_e, s_1) = 0$$

It's Mars we are talking about. It's different from our Earth. Our robot is not completely predictable there. When the robot decides to go toward a direction, with 0.2 probability, it may go the opposite direction.

Example: Mars Rover Dynamics Model (Dynamics Model)

State →

reward →

s_1	s_2	s_3	s_4	s_5	s_6

Transition Model:

$$T(s_i, a_L, s_{i+1}) = T(s_i, a_R, s_{i+1}) = 0.8 \quad \forall i \in \{2, \dots, 5\}$$

not completely predictable

$$T(s_i, a_L, s_{i+1}) = T(s_i, a_R, s_{i-1}) = 0.2 \quad //$$

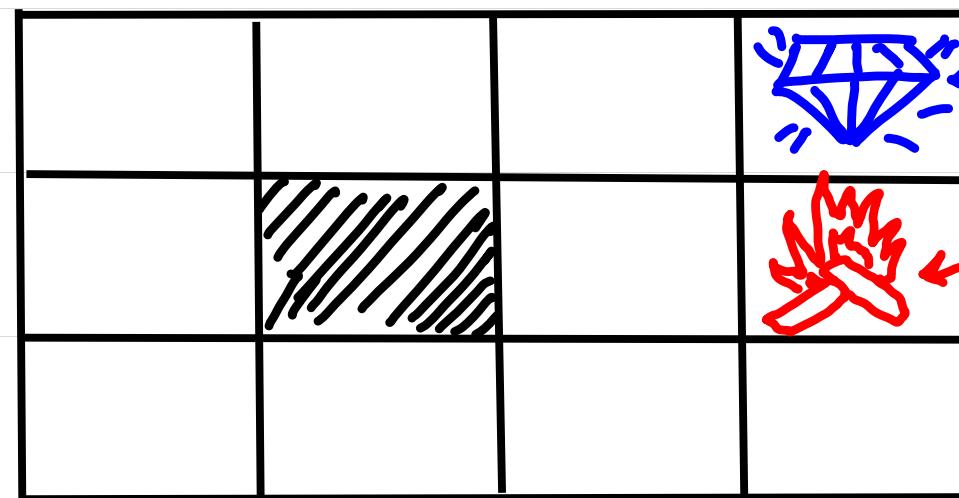
there. When the robot decides to go toward a direction, with 0.2 probability, it may go the opposite direction.

$$T(s_1, a_e, s_1) = T(s_1, a_e, s_2) = 1$$

It's Mars we are talking about. It's different from our Earth. Our robot is

Example: Mars Rover In Grid World

■ Consider the following Grid Model of Mars



Mars Rover gets a reward of (+1)
if it gets to this location

Mars Rover gets a reward of (-1)
if it gets to this location

$$R(s, a, s') = -0.03 \text{ for all non-terminal states } s'$$

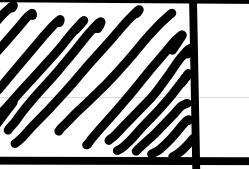
■ Noisy Movements on Mars: actions do not always go as planned

- with prob. 0.8, the action takes the agent to the desired direction
- With prob. 0.1, the actions takes the agent to a direction perpendicular to the desired direction.

- If there's a wall in the direction would have been taken to, it stays put.

Homework: Mars Rover In Grid World

- formulate this MDP.

$S_{2,0}$	$S_{2,1}$	$S_{2,2}$	$S_{2,3}$
$S_{1,0}$		$S_{1,2}$	$2_{1,3}$
$S_{0,0}$	$S_{0,1}$	$S_{0,2}$	$S_{0,3}$

Example: The undergrad Student's life at UoT

- Erfan wants to thrive in his undergrad
- He has three states: happy, tired, burnt-out
 - burnt-out is the terminal state.
- He has two actions: "Work hard", "Party hard"
 - Working hard gets double reward: reward 2 for working hard



s	a	s'	$T(s, a, s')$	$R(s, a)$
	Work hard		0.5	2
	Work hard		0.5	2
	party hard		1	1
	Work hard		1	2
	Party hard		0.5	1
	party hard		0.5	1
	end		1	0



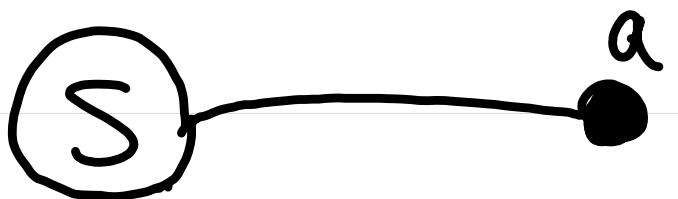
We can summarize
the dynamics
and reward model
with a table

Transition Graph

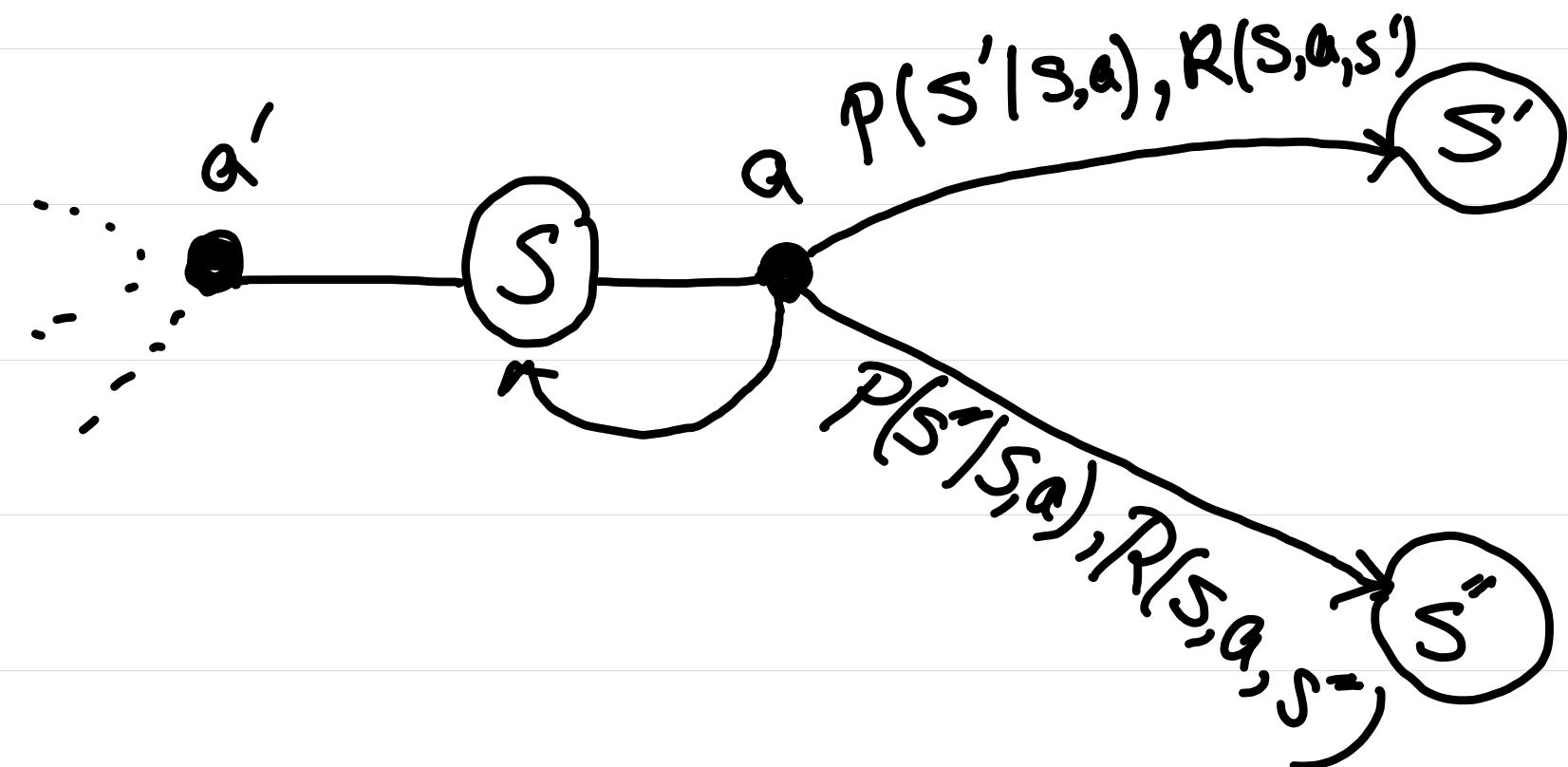
- Transition Graph is a useful way of summarizing the dynamics of a finite MDP
- There are two types of nodes
 - ↳ State nodes: There is a state node for each possible state
 - ▲ Typically a large circle labeled by the name of the state, e.g., 
 - ↳ Action nodes (also called q-state node): There is an action node for each state-action pair.
 - ▲ Typically a small solid circle labeled by the name of action or (state, action) pair and connected by a line to the state node, e.g.,  or 

Transition Graph

- Starting in state "s" and taking action "a" move you along the line from state nodes "s" to action node (s, a) .



- Then the environment responds with a transition to the next state's node.



Activity: Draw the transition Graph for Erfan the Undergrad

