# Individual Project

## Causal Discovery in Explainable Reinforcement Learning

Siran Shen

Imperial College London

# Introduction

- *Reinforcement learning (RL)* is widely applied in society.

- To trust *RL agents*, humans need to understand how they make decisions.

- It is crucial for RL agents to become *explainable*.

# Problem Statement and Goal

- This project is about explain RL models with *causality*.

- Project aim: automatise the process of inferring action influence models (AIM) at training time.

- Goal: to generate a *directed Acyclic graph(DAG)* with its corresponding action matrix (AM) which is needed for the AIM generation automatically.
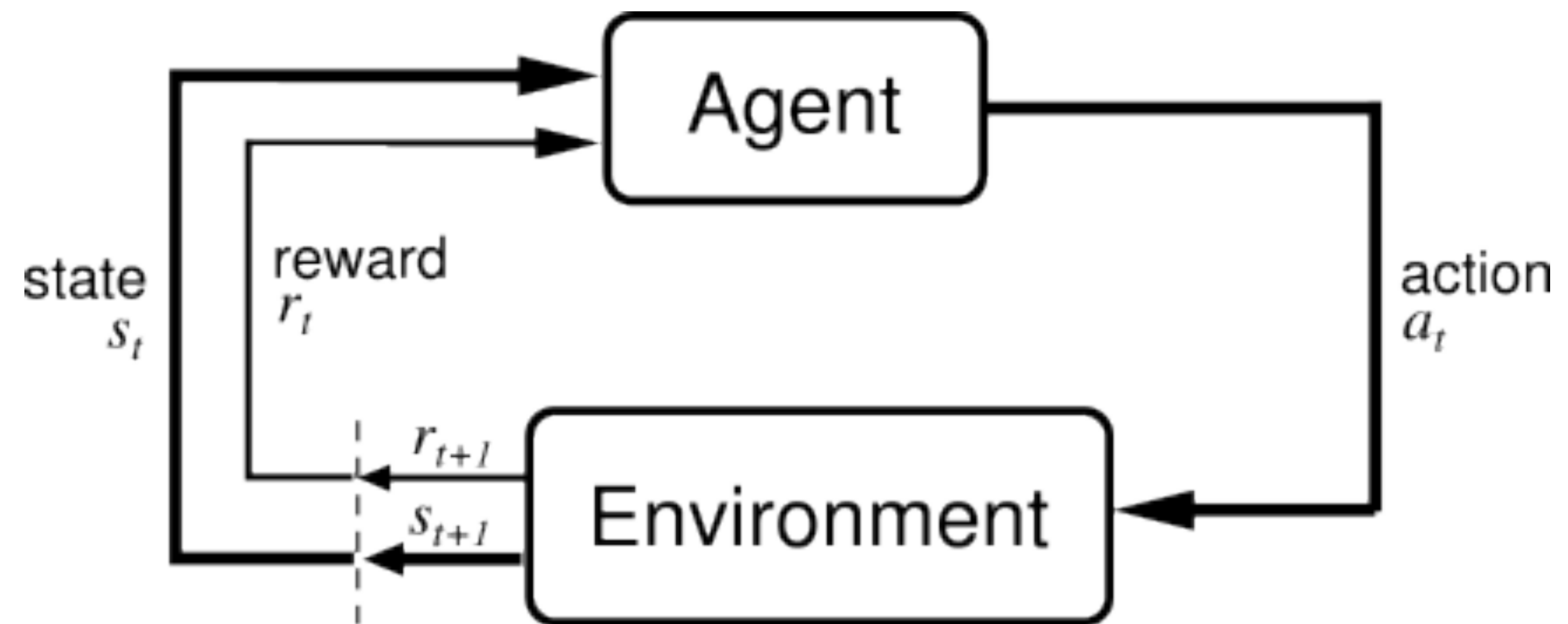
# Proposed Solution

- Main idea: a new method built upon the existing causal discovery algorithm - Greedy Equivalence Search (GES).

- The proposed algorithm can be spilt into two parts:

  1. *Complete partial DAG (CPDAG)* and Action Matrix Generator

  2. Causal DAG Maker

- Three RL environments:

  1. Taxi problem

  2. Cart pole

  3. Lunar landing

# Background
## Reinforcement Learning

- Reinforcement learning (RL): a process where agents try to find a best solution from environment.

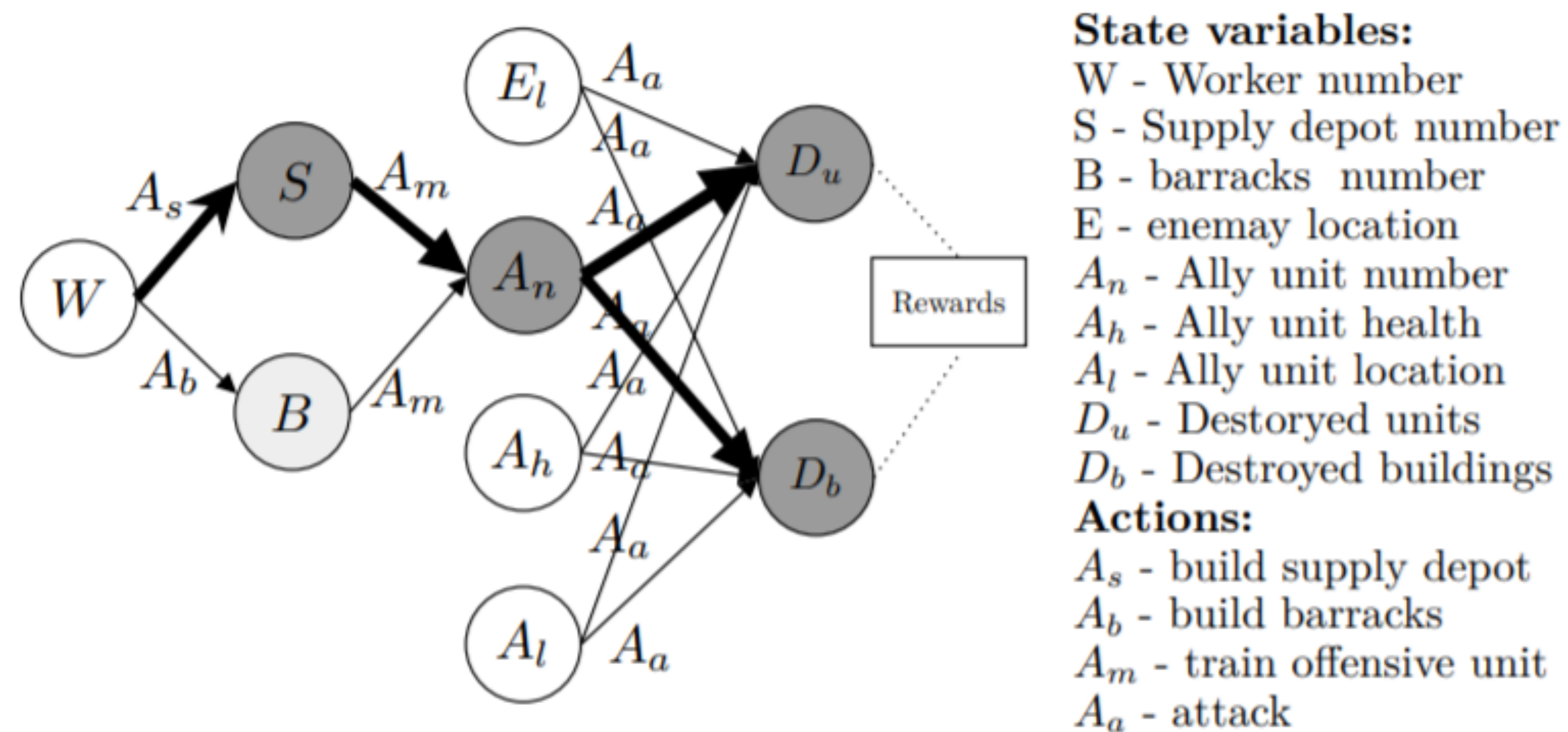- In RL, agent interacts with the environment continually: illustrated on the right.



Agent-Environment Interactions in RL

# Background
## Notions

- Action Influence Model (AIM):



StarCraft Action Influence Graph

# Background

## Notions

- Directed Acyclic Graph (DAG): a graph without cycles and only having directed edges.

- Action Matrix (AM): each non-zero element represents the index of the action performed on that edge.

.

$$\begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

A Causal DAG in matrix form of StarCraft Action Influence Graph

$$\begin{bmatrix} 0 & 91 & 42 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 477 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 477 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 13 & 13 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 13 & 13 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 13 & 13 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 13 & 13 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

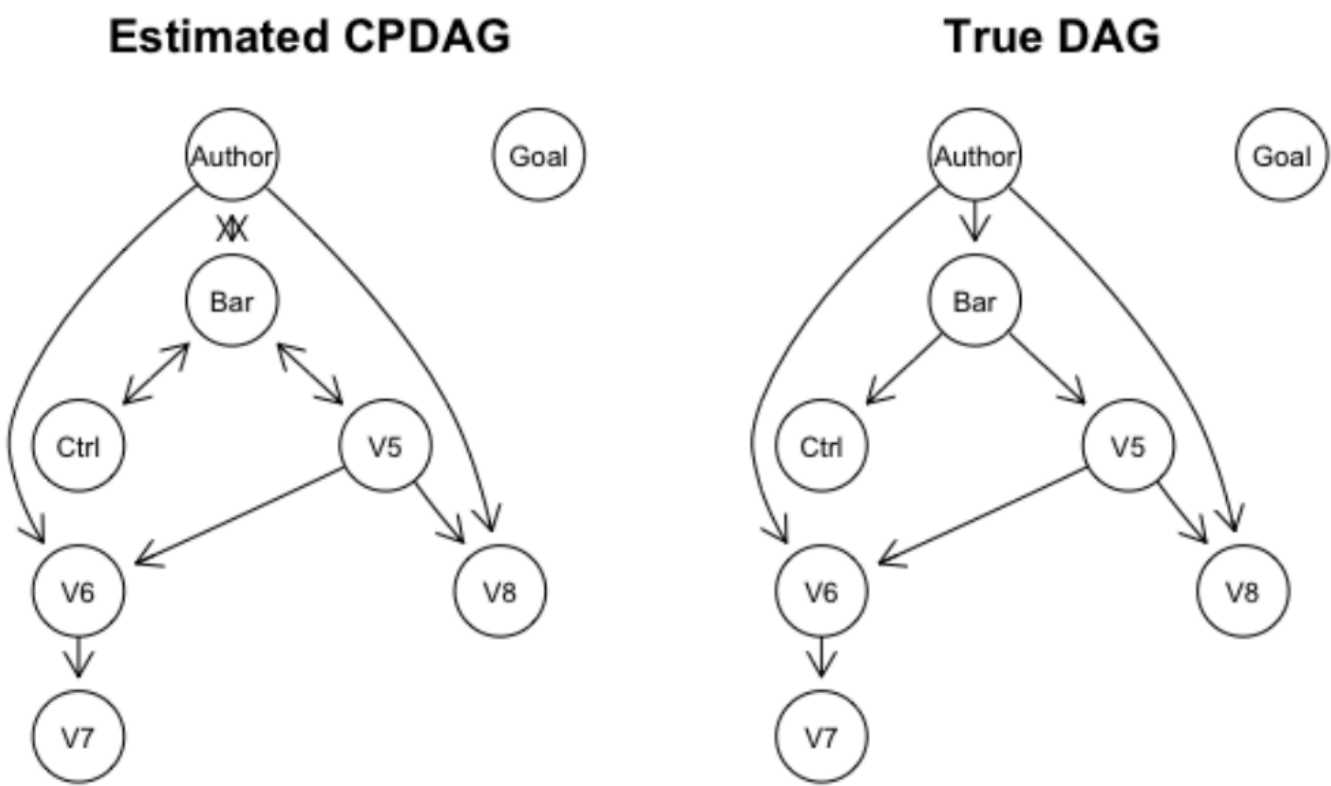The associated action matrix

# Background

## Greedy Equivalence Search(GES) Algorithm

- Input: a dataset along with an empty graph

- Output: a CPDAG with a score.

- Three stage: forward, backward and turning.



gmG Data



Estimated and True DAG from gmG Data

# Background
## Learning Explanations from Casual Graphic Model



Causal graphic DAG along with its action matrix provided

during the training phase

AIM learned

Explanations of "why A" and "why not A" questions produced (A is some action in the RL environment)

# Methodology

**CPDAG and Action Matrix Generator**

- Input: a dataset $D$

- Output: a CPDAG with a score

- Based on GES algorithm

- Group dataset by RL agent's actions first

```
┌─────────────────────────────────┐
│       input dataset D           │
└─────────────────────────────────┘
              │
       group D by actions
              ▼
┌─────────────────────────────────┐
│   datasets grouped by actions   │
└─────────────────────────────────┘
              │
     For each step, for each edge
              ▼
┌─────────────────────────────────┐
│ Select the action whose dataset │
│   gives the highest score       │
└─────────────────────────────────┘
```

# Methodology

## Casual DAG Maker
## - first stage

- Input: a set of CPDAG-Score pairs *S*

- Output: A graph *G* where there are multiple edges between two nodes

- Add the set of CPDAGs to G one by one.

# Methodology

**Casual DAG Maker
- second stage**

- Input: the graph *G* gotten from previous stage

- Output: A DAG with its AM

- Threshold parameter *H:* 0.1

- Get the final causal DAG by merging all the edges between two variables into one unique directed edge.

# Environments and Implementation

## Taxi Problem - Environment

- 5 × 5 grid world.

- Taxi cannot move across obstacles.

- 4 special locations: R, G, Y, B.

- 6 Actions: move north/ east/ west, pick up/ drop off passenger.

- Goal: drive the passenger to destination successfully

# Environments and Implementation

**Taxi Problem - Implementation**

- Train the agent by Q-learning algorithm.

- Run the game for 100 episodes with the trained agent.

- Get a data set with 4 columns: taxi location, passenger location, destination location, action.

- Final DAG produced: $3 \times 3$ matrix in form of:

$$
\begin{bmatrix}
taxi\_position & passenger\_position & destination\_location \\
passenger\_position & \cdots & \\
destination\_location & \cdots &
\end{bmatrix}
$$

# Environments and Implementation

**CartPole Problem - Environment**

- A pole attached to a cart moving along a frictionless track.

- 2 actions:

  1. push cart to the left,

  2. push cart to the right.

- Task: to keep the pole balanced for as long as possible.

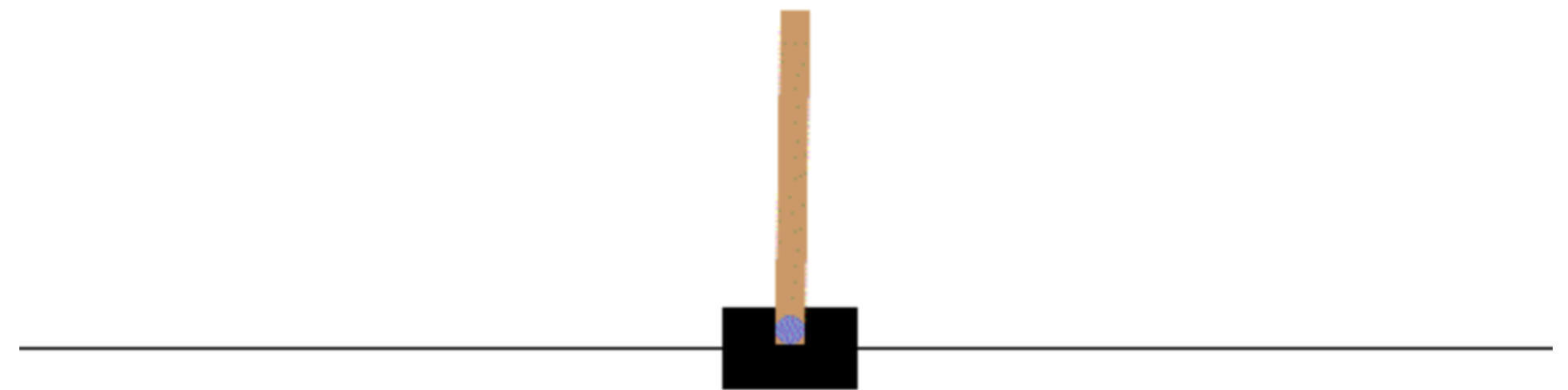# Environments and Implementation

## CartPole Problem - Implementation

- Train the agent by REINFORCE algorithm.

- Run the game for 100 episodes with the trained agent.

- Get a dataset with 5 columns: cart velocity, cart position, pole angle, pole angular velocity, action.

- The final DAG produced: $4 \times 4$ matrix in form of:

$$\begin{bmatrix} cart\_velocity & cart\_position & pole\_angle & pole\_angular\_velocity \\ cart\_position & \cdots & & \\ pole\_angle & \cdots & & \\ pole\_angular\_velocity & \cdots & & \end{bmatrix}$$

# Environments and Implementation
## LunarLand Problem - Environment

- Task: Land the LunarLander safely in the landing pad.

- 4 Actions:
  1. do nothing
  2. fire left orientation engine
  3. fire main engine
  4. fire right orientation engine

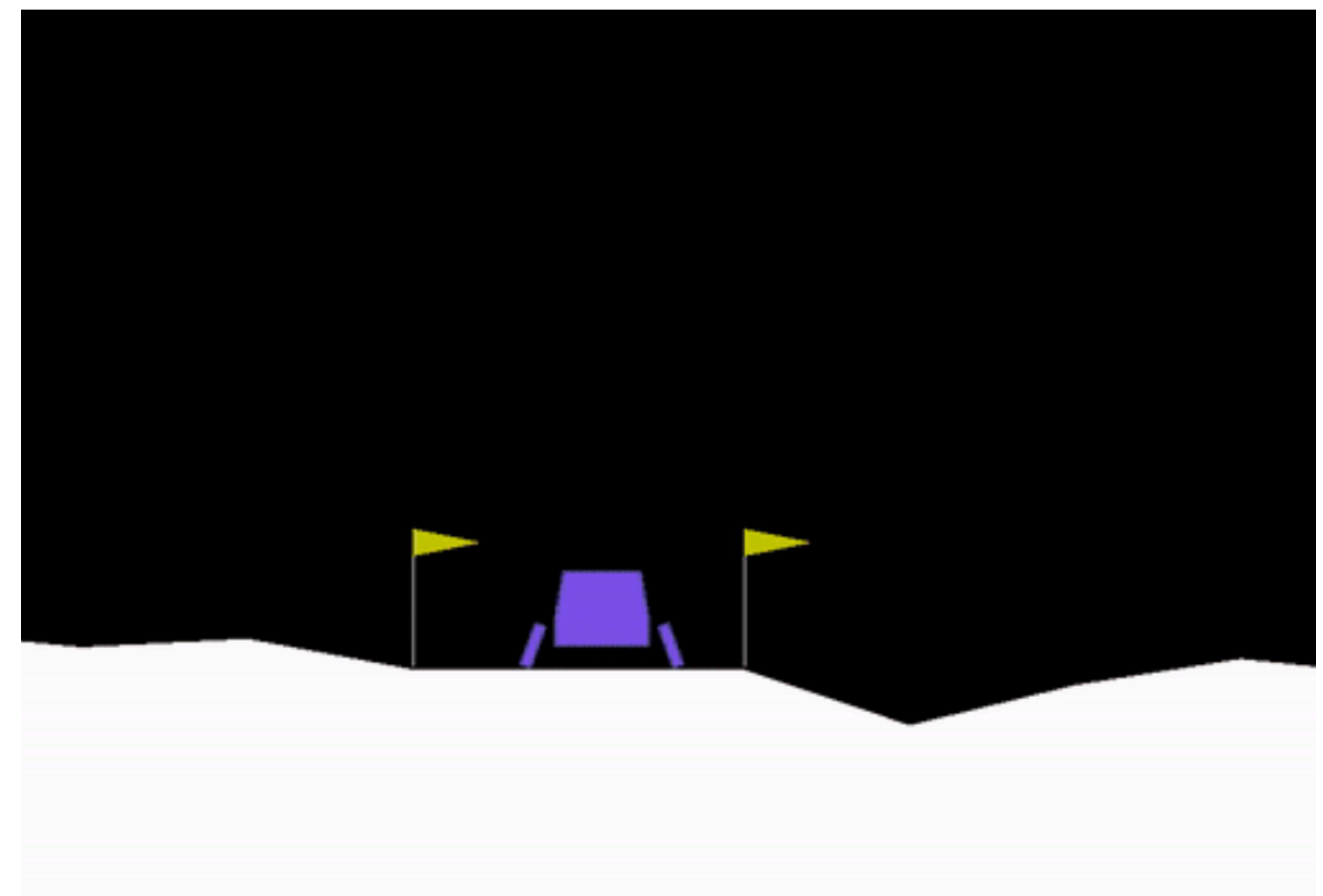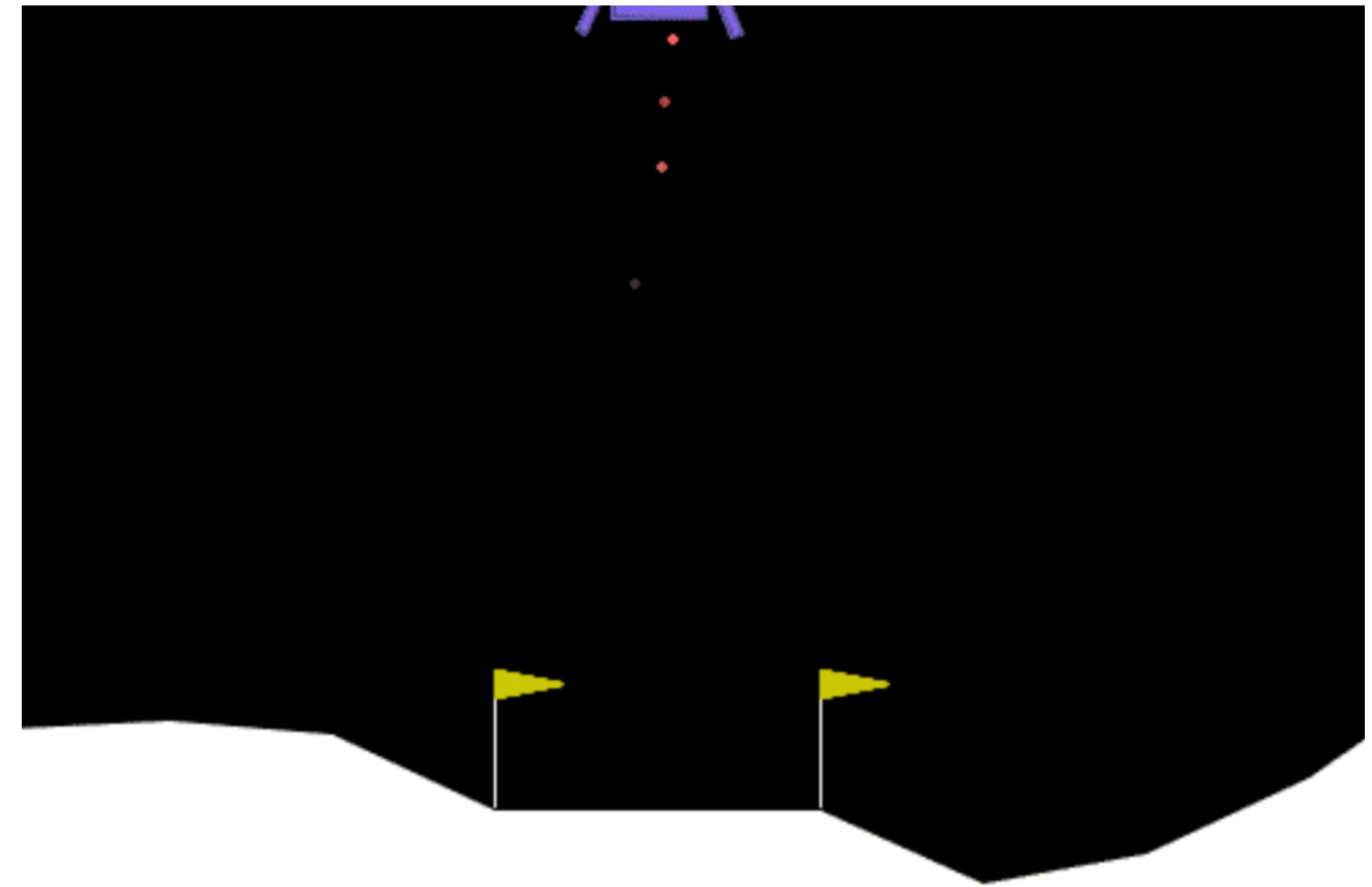# Environments and Implementation

**LunarLand Problem - Implementation**

- Train the agent by Deep Q-Network algorithm.

- Run the game for 100 episodes with the trained agent.

- Get a dataset with 7 columns: lander xCoor, lander yCoor, lander xV, lander yV, lander angle, lander angularV, total reward, action.

- The final DAG produced: $6 \times 6$ matrix in form of:

$$
\begin{bmatrix}
xCoor & yCoor & xVelocity & yVelocity & angle & angularVelocity \\
yCoor & \cdots & & & & \\
xVelocity & \cdots & & & & \\
yVelocity & \cdots & & & & \\
angle & \cdots & & & & \\
angularVelocity & \cdots & & & &
\end{bmatrix}
$$

# Experimental Results
## Taxi Problem -  Causal DAG with AM

# Experimental Results

**Taxi Problem - Explanations**

- "Why?" Explanation: Before taxi pick up the passenger, taxi does "move" action because the goal is to get passenger location.

- "Why not?" Explanation: Before taxi pick up the passenger, taxi does "move" action, not action "drop off" because it is more desirable to do action "move" to change taxi position, as the goal is to get passenger location.

# Experimental Results
## CartPole Problem - Causal DAG with AM

# Experimental Results
## CartPole Problem - Explanations

- "Why?" explanation: The force pushes cart to the right because the goal is to change cart position that depends on pole angular velocity/ pole angle.

- "Why not?" explanation: The force pushes cart to the left, but not right because it is more desirable to do action "push cart to the left", to have more "pole angular velocity"/ "pole angle", as the goal is to have "cart position".

# Experimental Results
## LunarLand Problem - Causal DAG with AM

# Experimental Results

**LunarLand Problem - Explanations**

- why perform "do nothing" action:

  1. Because the goal is to change the coordinate of the lander in x.

  2. Because the goal is to change the coordinate of the lander in x, that depends on the coordinates of the lander in y.

  3. Because the goal is to change the coordinate of the lander in x, which is influenced by its linear velocities in x, the coordinate of the lander in y, that depends on its linear velocities in y.

# Experimental Results

**LunarLand Problem - Explanations**

- why perform "fire main engine" action:

  1. Because the goal is to increase the coordinate of the lander in x-axis.

  2. Because the goal is to increase the coordinate of the lander in x-axis that depends on the coordinate of the lander in y-axis.

  3. Because the goal is to increase the coordinate of the lander in x-axis that depends on its linear velocities in x-axis.

# Experimental Results
## LunarLand Problem - Explanations

- why perform "do nothing" action but not "fire main engine" action:

  1. Because it is more desirable to do action do nothing to have more the coordinate of the lander in y-axis, in order to have more its linear velocity in x-axis, as the goal is to have the coordinate of the lander in x-axis.

  2. Because it is more desirable to do action do nothing to have more its linear velocity in x-axis, as the goal is to have the coordinate of the lander in x-axis.

# Experimental Results

## LunarLand Problem - Explanations

• why perform "fire main engine" action but not "do nothing" action:

1. Because it is more desirable to do action fire main engine to have less its linear velocity in y-axis in order to have more its angle, as the goal is to have the coordinate of the lander in x-axis.

2. Because it is more desirable to do action fire main engine, to have more the coordinate of the lander in y-axis, as the goal is to have the coordinate of the lander in x-axis.

3. Because it is more desirable to do action fire main engine, to have more its linear velocity in x-axis, to have less its linear velocity in y-axis, as the goal is to have the coordinate of the lander in x-axis.

# Evaluation
## Achievements

- Causal DAG along with its action matrix can be generated successfully

- meaningful explanations can be produced as well.

# Evaluation

**Future Work**

- DAG generated from the environment with fewer variables could be more reliable.

- Our designed algorithm might treat correlation between two variables as causal relation.

- a computational evaluation is still needed.

Any questions?