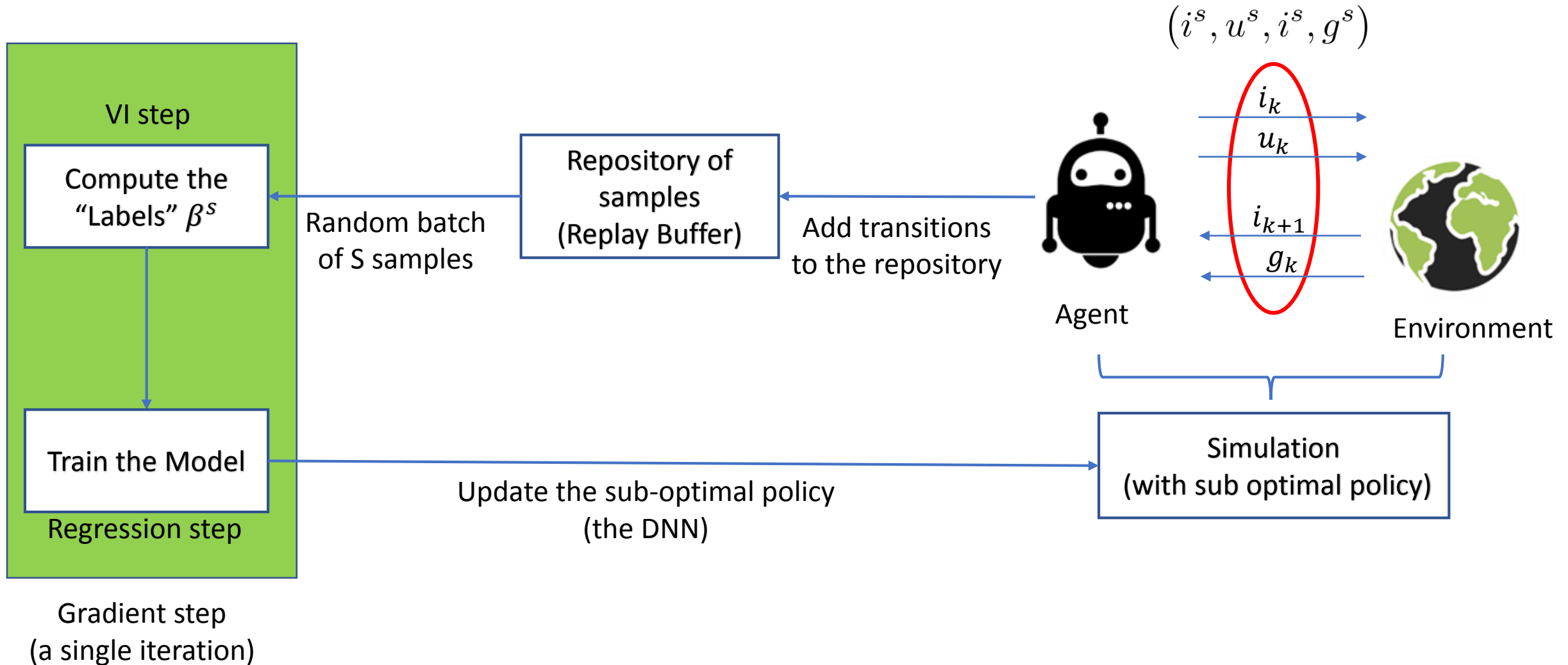


# Deep Q-Learning Algorithm

- First let's recall the following diagram:



# Deep Q Networks Algorithm (DQN)

---

**Algorithm 1** DQN Algorithm (Minh et al, 2015)

---

**Input:** Initial DNN parameters  $\theta^0$  and architecture  $Q_{\theta^0}(i, u)$ . Replay Buffer  $\mathcal{B}$ .

- 1: **for**  $p = 1, \dots, P$ : **do** (updates on the target network  $\theta'$ )
- 2:     Save the network parameters  $\theta' \leftarrow \theta$ .
- 3:     **for**  $k = 0, \dots, K$  **do** (obtaining new samples)
- 4:         Collect  $M$  sample transitions  $\{(i^m, u^m, i^{m+1}, g^m)\}_{m=1}^M$  using the sub-optimal policy:

$$\tilde{\mu}^{(t+1)}(i) = (1 - \epsilon^{(t+1)}) \arg \min_{u \in U(i)} \left\{ \tilde{Q}_{\theta^{(t+1)}}(i, u) \right\} + \epsilon^{(t+1)} A(U(i))$$

and add those sample transitions to the Replay Buffer  $\mathcal{B}$ .

- 5:     **for**  $t = 1, \dots, T$ : **do** (gradient step)
- 6:         Randomly sample a batch of size  $S$  from the Replay Buffer  $\mathcal{B}$ .
- 7:         Perform one gradient step:

$$\theta^{(t+1)} = \theta^{(t)} - \gamma^{(t)} \left\{ \sum_{i=1}^S \nabla_{\theta} \tilde{Q}_{\theta^{(t)}}(i^s, u^s) (\tilde{Q}_{\theta^{(t)}}(i^s, u^s) - g(i^s, u^s) - \alpha \min_{u \in U(i^s)} \{\tilde{Q}_{\theta'}(i^{s+1}, u)\}) \right\}$$

- 8:     **end for**
- 9:     **end for**
- 10: **end for**

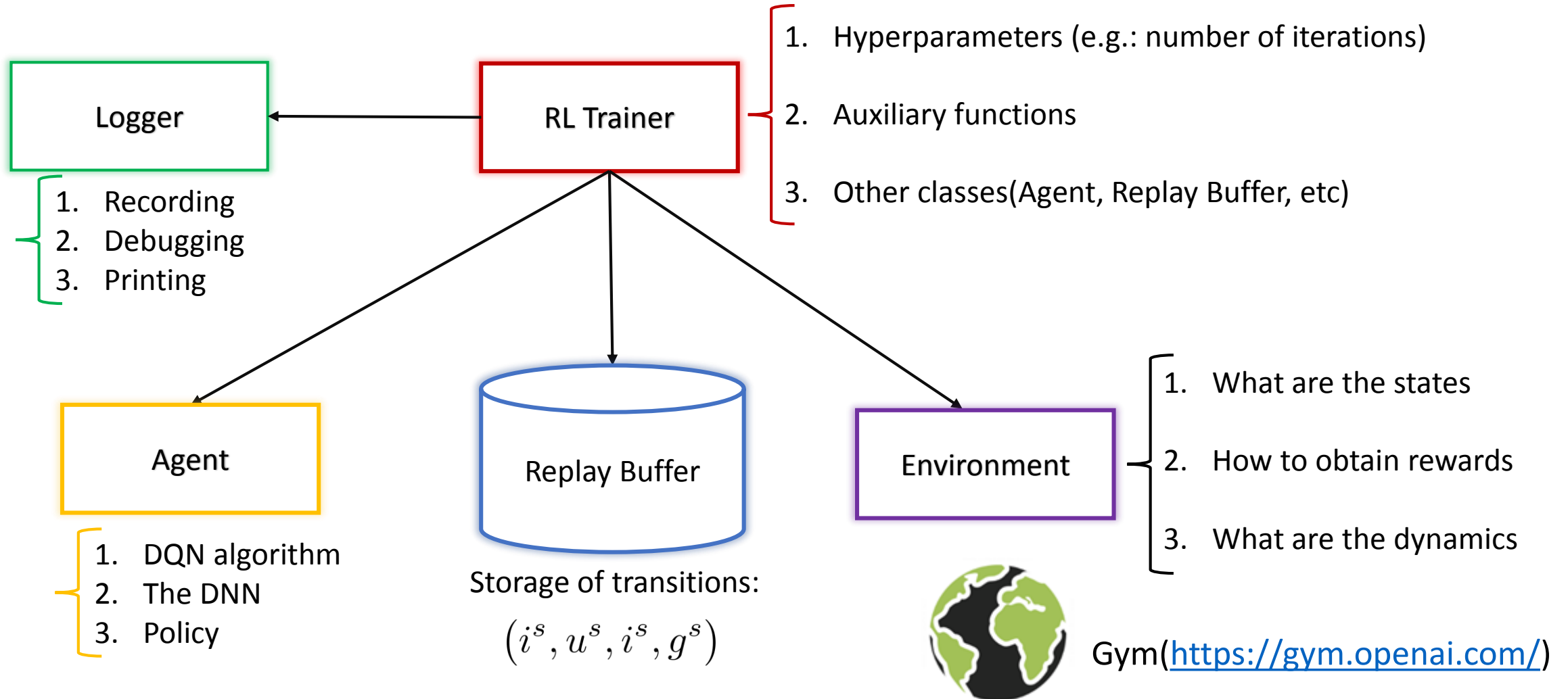
**Output:** The last DNN configuration  $\bar{\theta}$ . A suboptimal policy:

$$\tilde{\mu}(i) = \arg \min_{u \in U(i)} \left\{ \tilde{Q}_{\bar{\theta}}(i, u) \right\}$$

---

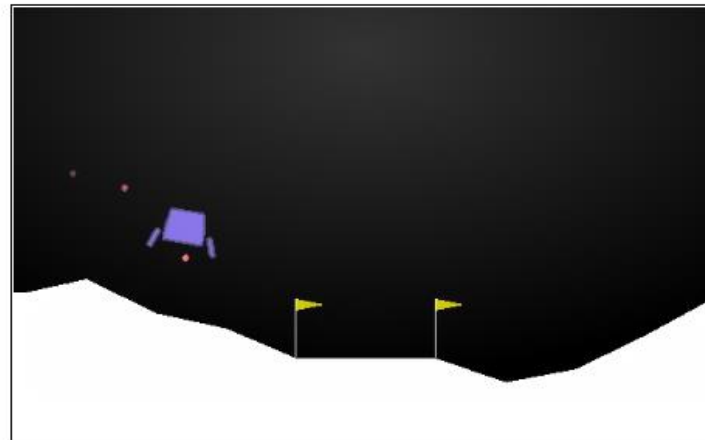
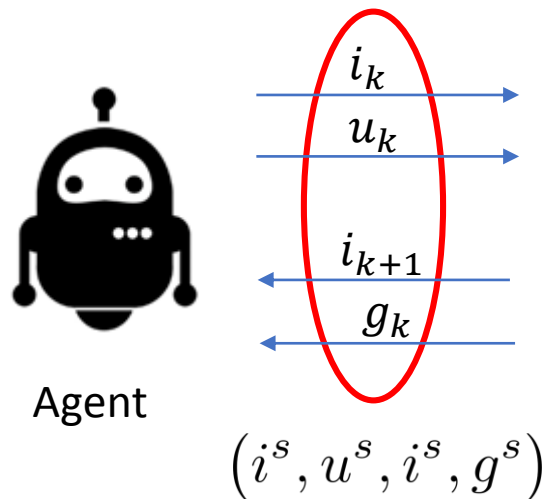
# Implementing DQN Algorithm

- We can divide the algorithm into classes:

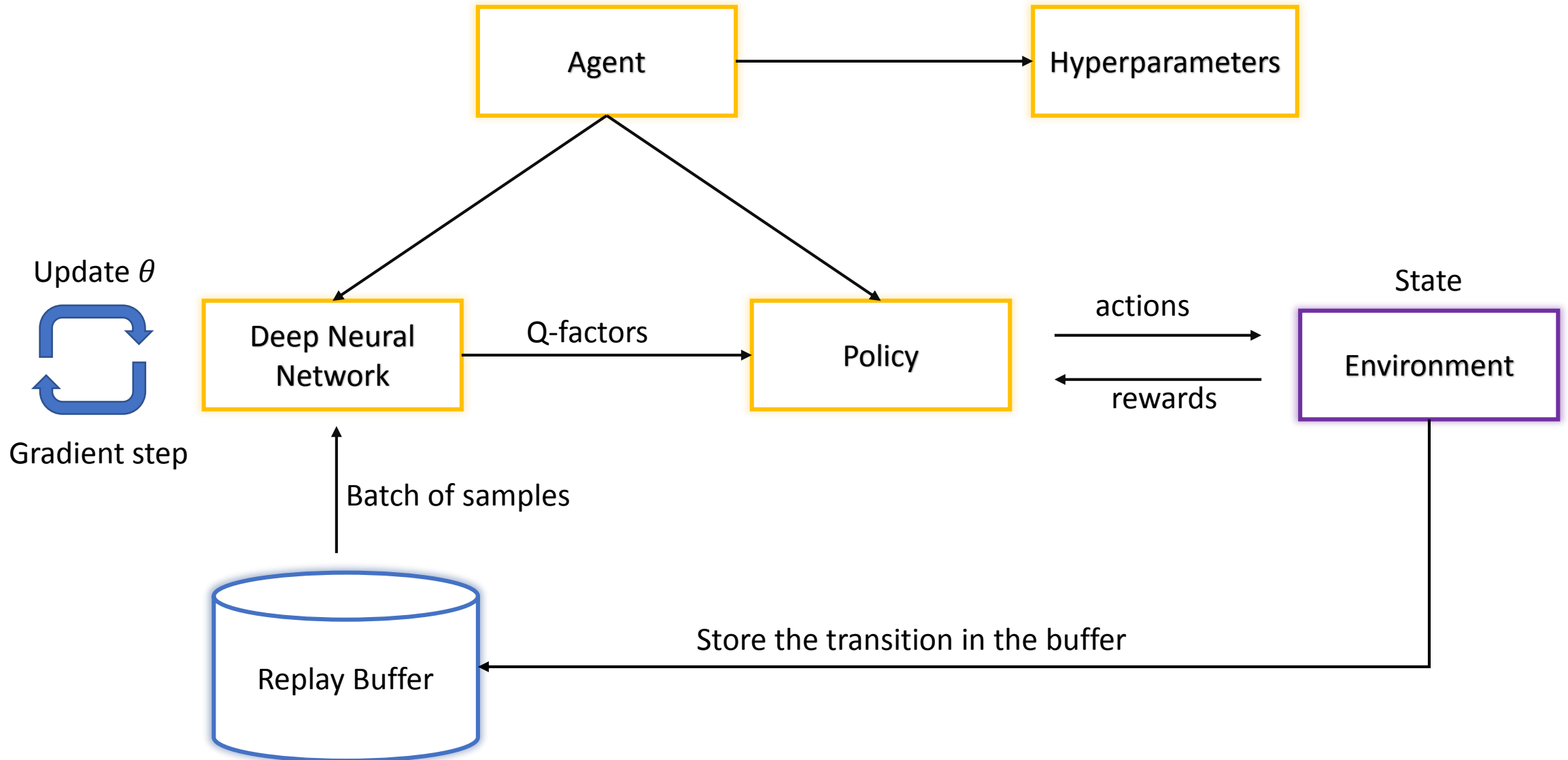


# Environment

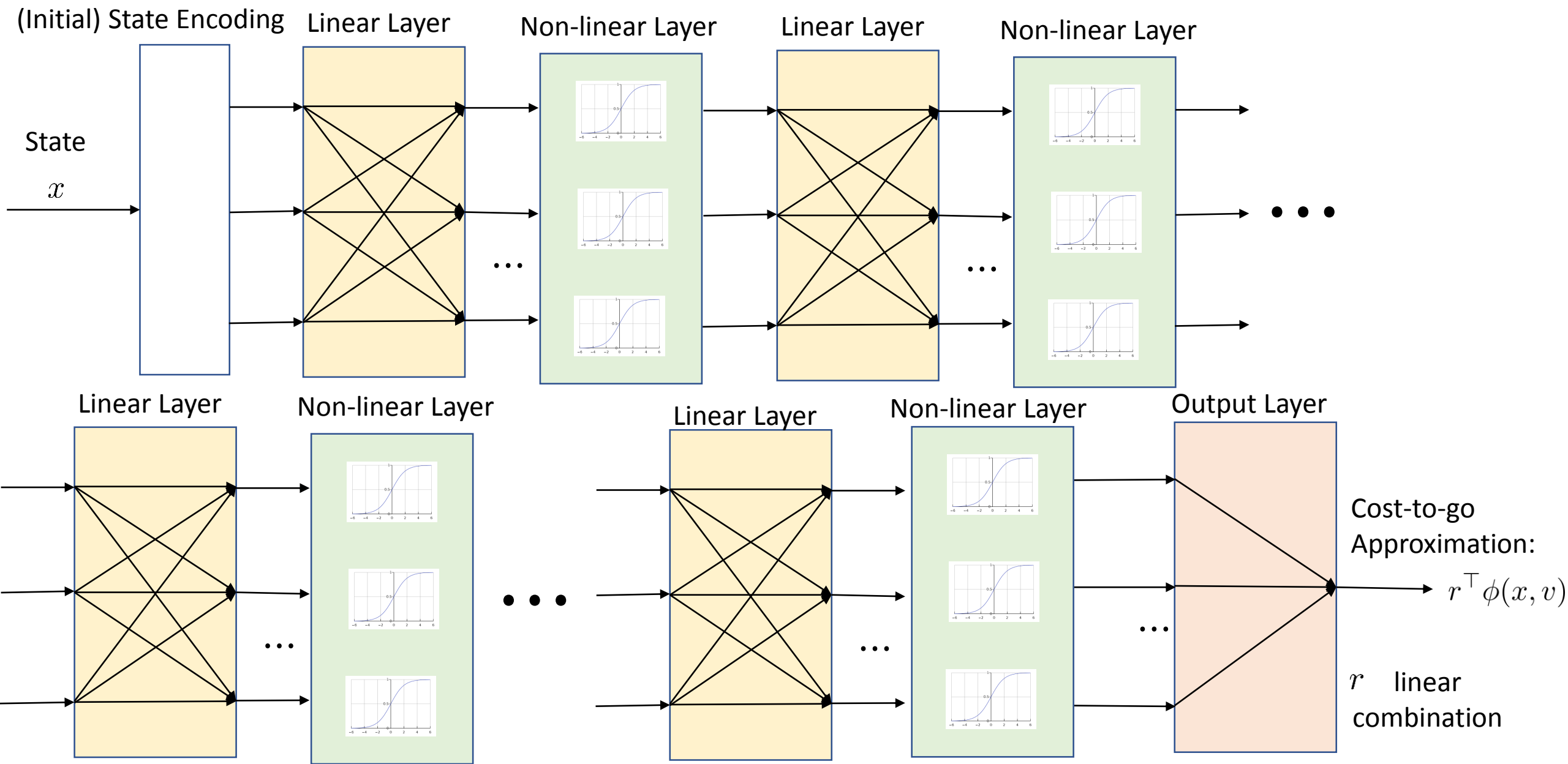
- We will try to design an AI to play the “Lunar Lander” game:
  - The goal is to land a space shuttle on the upright position with zero vertical speed between the flags at coordinate (0,0).
  - Reward for moving from the top of the screen to landing pad and zero speed is about 100..140 points.
  - Landing outside the landing pad loses reward.
  - Episode finishes if the lander crashes or comes to rest, receiving additional -100 or +100 points.
  - Each leg ground contact is +10.
  - Firing main engine is -0.3 points each frame.
  - Max score is 200 points. Landing outside landing pad is possible.
  - Four discrete actions: do nothing, fire left orientation engine, fire main engine, fire right orientation engine.



# Agent configuration

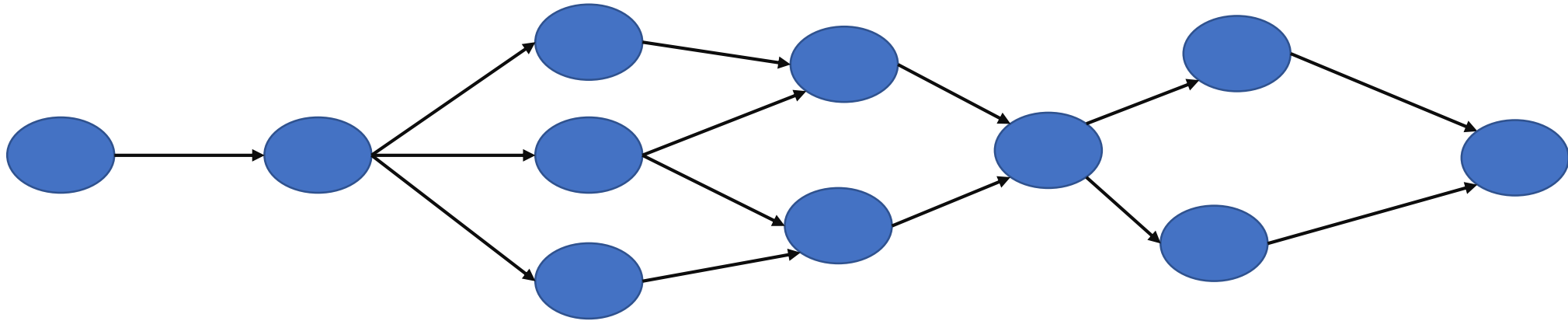


# Implementing a DNN: Tensorflow



# Tensorflow: a general overview

- Tensorflow is a platform that allows us flexible construction of **computational graphs**
  - A neural network is a computational graph
- The key idea is that Tensorflow allow us to first build the graph (with no data):



- And then, on runtime, it let's us “flow” the data through the graph.
- It also allow us to differentiate the operations on each node: **Backpropagation**

# DQN Implementation

Let's jump to the code!

Code is available on Bcourses.

Zoom meeting link with the lecture:

[https://berkeley.zoom.us/rec/share/1O1Yd5b-9EZOXM\\_0zGHEQrQCLqv7aaa80XAX8vFfxRmWoG3gdGTUEvswzMyoig-q?startTime=1583877350000](https://berkeley.zoom.us/rec/share/1O1Yd5b-9EZOXM_0zGHEQrQCLqv7aaa80XAX8vFfxRmWoG3gdGTUEvswzMyoig-q?startTime=1583877350000)