

Introduction to Big Data Science

12-4 Period

Introduction to Reinforcement Learning

Contents

- ◆ Introduction to Reinforcement Learning
- ◆ Q-Learning and Value-Based Deep RL
- ◆ Policy-Based Deep RL
- ◆ Model-Based Deep RL
- ◆ Example of Q-Learning

Reinforcement Learning

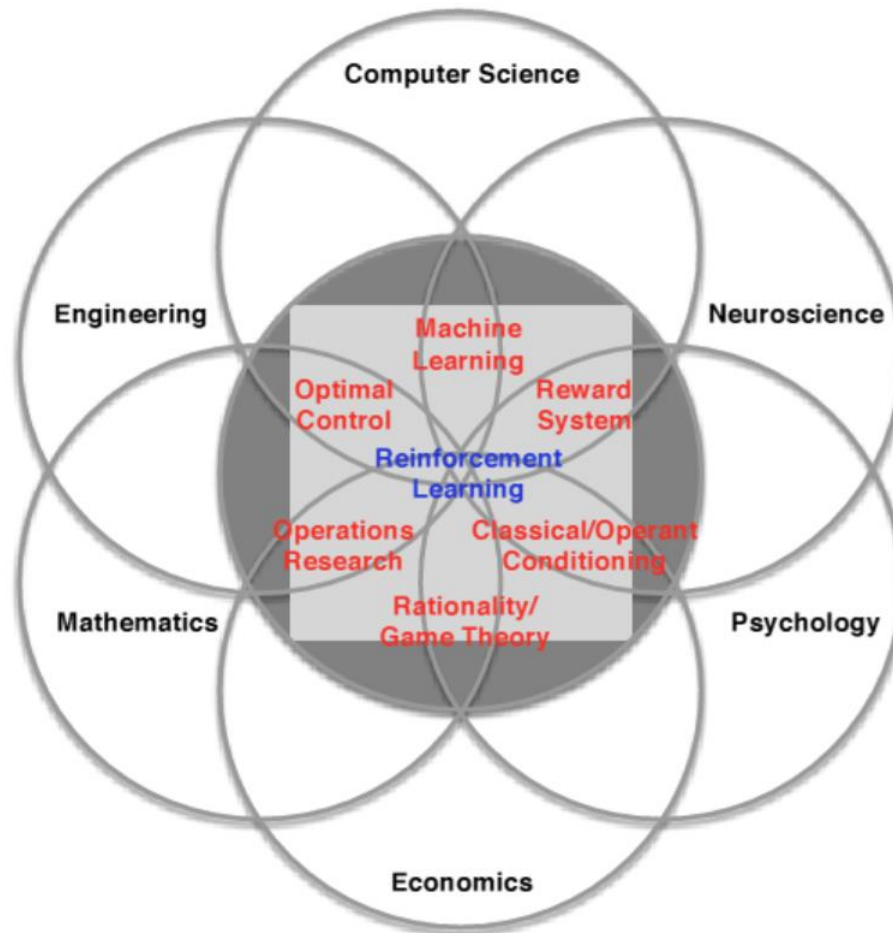
RL is a general-purpose framework for decision-making

- ▶ RL is for an **agent** with the capacity to **act**
- ▶ Each **action** influences the agent's future **state**
- ▶ Success is measured by a scalar **reward** signal
- ▶ Goal: **select actions to maximise future reward**

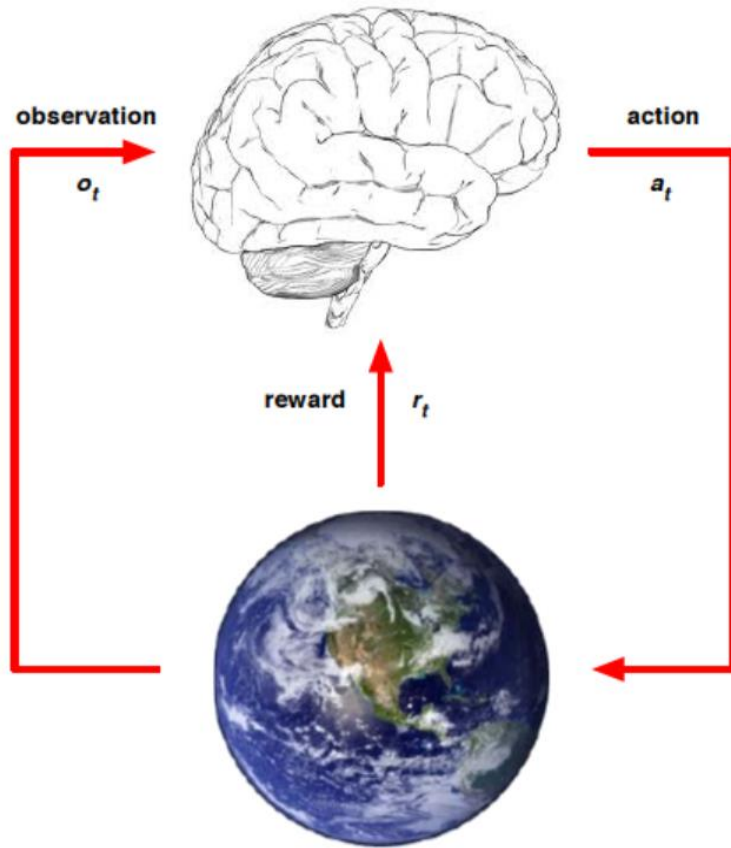
Deep Reinforcement Learning

- ◆ We seek a single agent which can solve any human-level task.
- ◆ RL defines the objective.
- ◆ DL gives the mechanism.
- ◆ RL + DL = general intelligence

Many Faces of RL



Agent and Environment



- ▶ At each step t the agent:
 - ▶ Executes action a_t
 - ▶ Receives observation o_t
 - ▶ Receives scalar reward r_t
- ▶ The environment:
 - ▶ Receives action a_t
 - ▶ Emits observation o_{t+1}
 - ▶ Emits scalar reward r_{t+1}

State



- Experience is a sequence of observations, actions, rewards

$$o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t$$

- The **state** is a summary of experience

$$s_t = f(o_1, r_1, a_1, \dots, a_{t-1}, o_t, r_t)$$

- In a fully observed environment

$$s_t = f(o_t)$$

Major Components of an RL Agent

- ◆ An RL agent may include one or more of these components:
 - Policy: agent's behavior function
 - Value function: how good is each state and/or action
 - Model: Agent's representation of the environment

Policy

- ▶ A **policy** is the agent's behaviour
- ▶ It is a map from state to action:
 - ▶ Deterministic policy: $a = \pi(s)$
 - ▶ Stochastic policy: $\pi(a|s) = \mathbb{P}[a|s]$

Value Function

- ▶ A **value function** is a prediction of future reward
 - ▶ “How much reward will I get from action a in state s ?”
- ▶ **Q -value function** gives expected total reward
 - ▶ from state s and action a
 - ▶ under policy π
 - ▶ with discount factor γ

$$Q^\pi(s, a) = \mathbb{E} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots \mid s, a]$$

- ▶ Value functions decompose into a Bellman equation

$$Q^\pi(s, a) = \mathbb{E}_{s', a'} [r + \gamma Q^\pi(s', a') \mid s, a]$$

Optimal Value Function

- ▶ An optimal value function is the maximum achievable value

$$Q^*(s, a) = \max_{\pi} Q^{\pi}(s, a) = Q^{\pi^*}(s, a)$$

- ▶ Once we have Q^* we can act optimally,

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$$

- ▶ Optimal value maximises over all decisions. Informally:

$$\begin{aligned} Q^*(s, a) &= r_{t+1} + \gamma \max_{a_{t+1}} r_{t+2} + \gamma^2 \max_{a_{t+2}} r_{t+3} + \dots \\ &= r_{t+1} + \gamma \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1}) \end{aligned}$$

- ▶ Formally, optimal values decompose into a Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q^*(s', a') \mid s, a \right]$$

Q Learning

- ◆ A model-free reinforcement learning technique. Specifically, Q-learning can be used to find an optimal action-selection policy for any given (finite) Markov decision process (MDP).
- ◆ It works by learning an action-value function that ultimately gives the expected utility of taking a given action in a given state and following the optimal policy thereafter.
- ◆ A policy is a rule that the agent follows in selecting actions, given the state it is in. When such an action-value function is learned, the optimal policy can be constructed by simply selecting the action with the highest value in each state.
- ◆ One of the strengths of Q-learning is that it is able to compare the expected utility of the available actions without requiring a model of the environment.
- ◆ Q-learning can handle problems with stochastic transitions and rewards, without requiring any adaptations. It has been proven that for any finite MDP, Q-learning eventually finds an optimal policy, in the sense that the expected value of the total reward return over all successive steps, starting from the current state, is the maximum achievable.

Q Learning

The algorithm therefore has a function that calculates the Quality of a state-action combination:

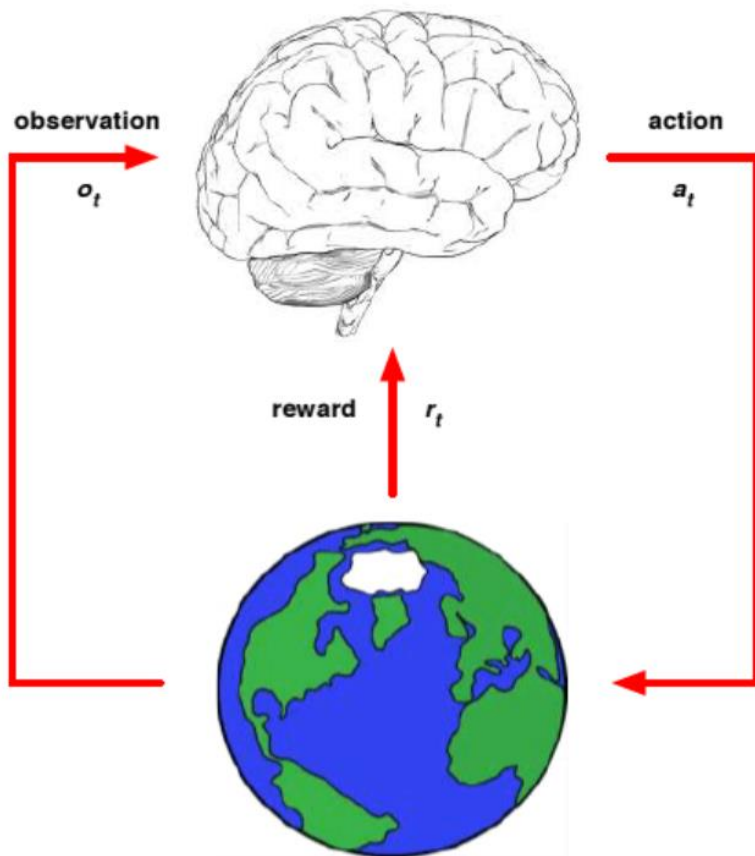
$$Q : S \times A \rightarrow \mathbb{R} .$$

Before learning has started, Q returns an (arbitrary) fixed value, chosen by the designer. Then, at each time t the agent selects an action a_t and observes a reward r_t and a new state s_{t+1} that may depend on both the previous state s_t and the selected action, Q is updated. The core of the algorithm is a simple value iteration update, using the weighted average of the old value and the new information:

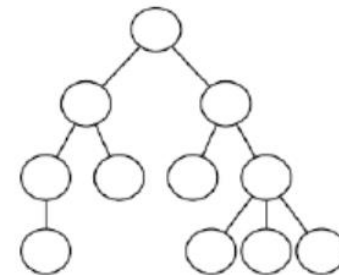
$$Q(s_t, a_t) \leftarrow (1 - \alpha) \cdot \underbrace{Q(s_t, a_t)}_{\text{old value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \overbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} \right)}^{\text{learned value}}$$

where r_t is the reward observed for the current state s_t , and α is the learning rate ($0 < \alpha \leq 1$).

Model



- ▶ **Model** is learnt from experience
- ▶ Acts as proxy for environment
- ▶ Planner interacts with model
- ▶ e.g. using lookahead search



Approaches to Reinforcement Learning

Value-based RL

- ▶ Estimate the **optimal value function** $Q^*(s, a)$
- ▶ This is the maximum value achievable under any policy

Policy-based RL

- ▶ Search directly for the **optimal policy** π^*
- ▶ This is the policy achieving maximum future reward

Model-based RL

- ▶ Build a model of the environment
- ▶ Plan (e.g. by lookahead) using model

Deep Reinforcement Learning

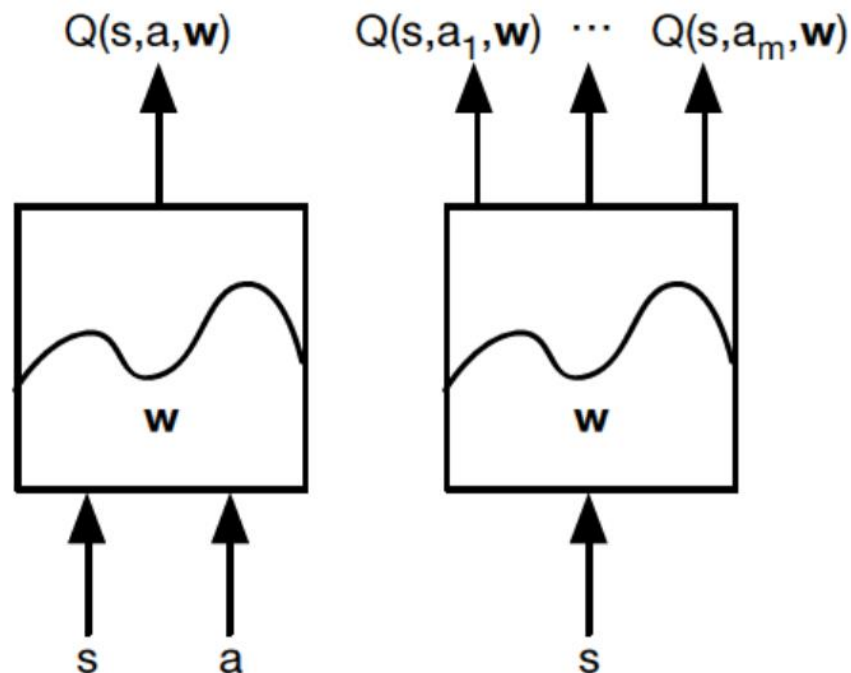
- ◆ Use deep neural networks to represent
 - Value function
 - Policy
 - Model
- ◆ Optimize loss function by stochastic gradient descent

- ◆ Value-based Deep RL

Q-Network

Represent value function by **Q-network** with weights **w**

$$Q(s, a, \mathbf{w}) \approx Q^*(s, a)$$



Q-Learning

- ▶ Optimal Q-values should obey Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left[r + \gamma \max_{a'} Q(s', a')^* \mid s, a \right]$$

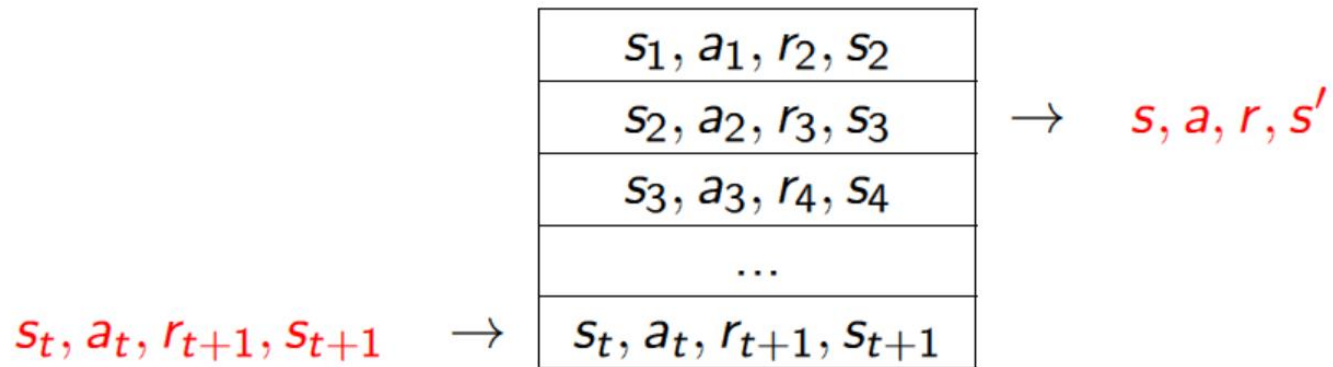
- ▶ Treat right-hand side $r + \gamma \max_{a'} Q(s', a', \mathbf{w})$ as a target
- ▶ Minimise MSE loss by stochastic gradient descent

$$l = \left(r + \gamma \max_a Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

- ▶ Converges to Q^* using table lookup representation
- ▶ But **diverges** using neural networks due to:
 - ▶ Correlations between samples
 - ▶ Non-stationary targets

Deep Q-Network(DQN): Experience Replay

To remove correlations, build data-set from agent's own experience

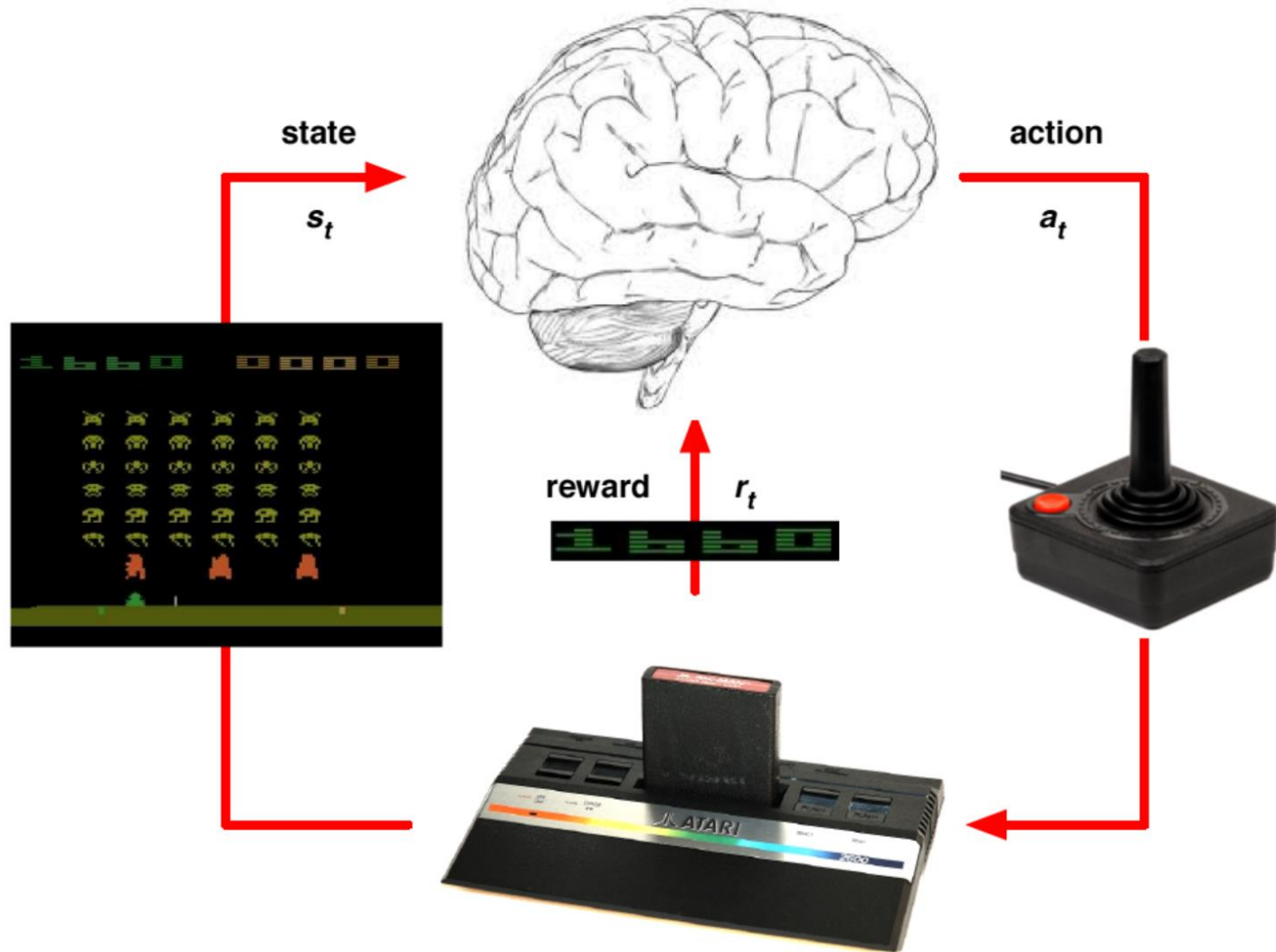


Sample experiences from data-set and apply update

$$l = \left(r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

To deal with non-stationarity, target parameters \mathbf{w}^- are held fixed

Deep Reinforcement Learning in Atari



Improvements since Nature DQN

- ▶ **Double DQN:** Remove upward bias caused by $\max_a Q(s, a, \mathbf{w})$
 - ▶ Current Q-network \mathbf{w} is used to **select** actions
 - ▶ Older Q-network \mathbf{w}^- is used to **evaluate** actions

$$l = \left(r + \gamma Q(s', \underset{a'}{\operatorname{argmax}} Q(s', a', \mathbf{w}), \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

- ▶ **Prioritised replay:** Weight experience according to surprise
 - ▶ Store experience in priority queue according to DQN error

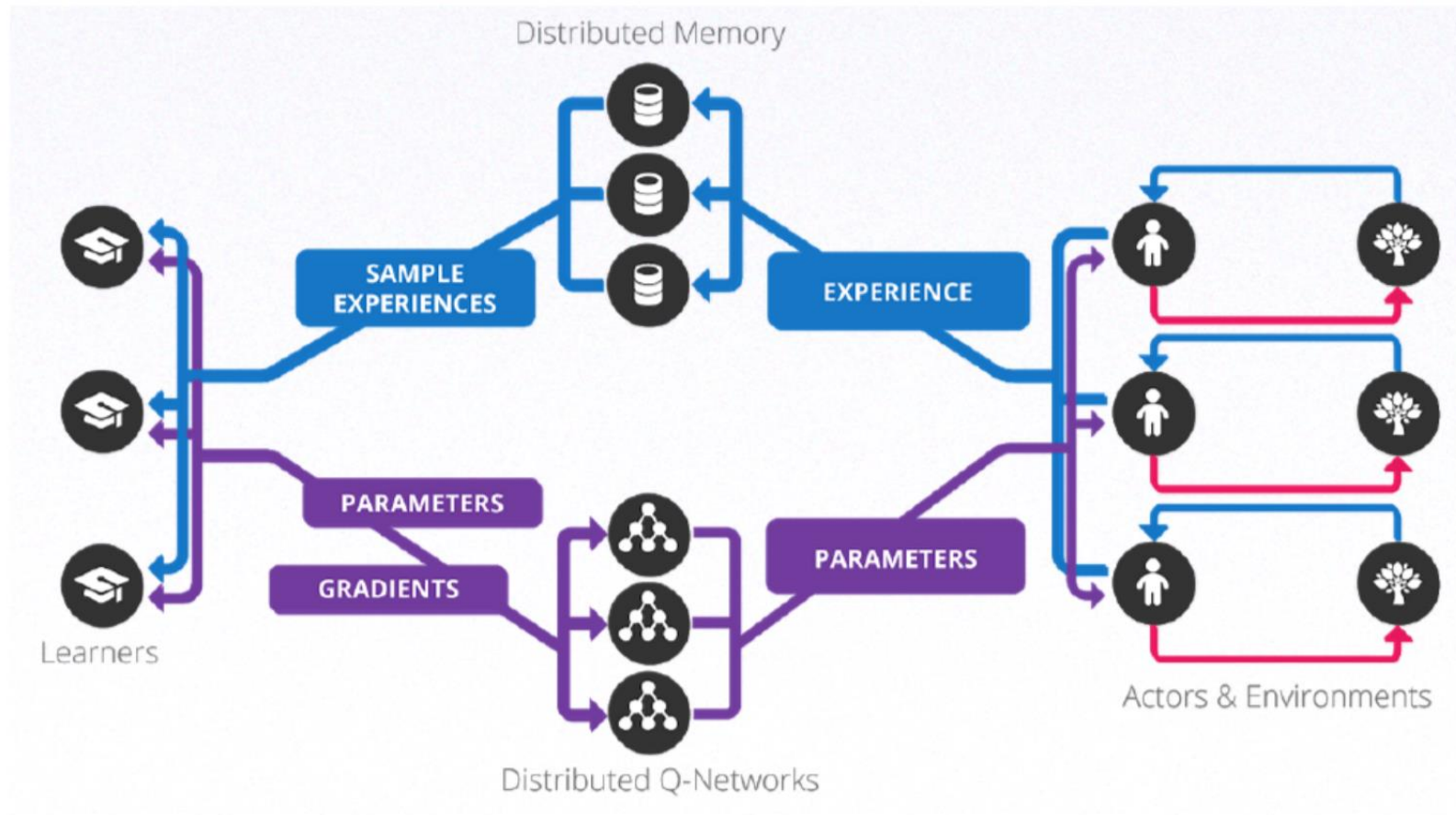
$$\left| r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right|$$

- ▶ **Duelling network:** Split Q-network into two channels
 - ▶ Action-independent **value function** $V(s, v)$
 - ▶ Action-dependent **advantage function** $A(s, a, \mathbf{w})$

$$Q(s, a) = V(s, v) + A(s, a, \mathbf{w})$$

Combined algorithm: 3x mean Atari score vs Nature DQN

Gorila(General Reinforcement Learning Architecture)



- ▶ 10x faster than Nature DQN on 38 out of 49 Atari games
- ▶ Applied to recommender systems within Google

- ◆ Policy-based Deep RL

Deep Policy Networks

- ▶ Represent policy by deep network with weights \mathbf{u}

$$a = \pi(a|s, \mathbf{u}) \text{ or } a = \pi(s, \mathbf{u})$$

- ▶ Define objective function as total discounted reward

$$L(\mathbf{u}) = \mathbb{E} [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots \mid \pi(\cdot, \mathbf{u})]$$

- ▶ Optimise objective end-to-end by SGD
- ▶ i.e. Adjust policy parameters \mathbf{u} to achieve more reward

Policy Gradients

How to make high-value actions more likely:

- ▶ The gradient of a stochastic policy $\pi(a|s, \mathbf{u})$ is given by

$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E} \left[\frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q^\pi(s, a) \right]$$

- ▶ The gradient of a deterministic policy $a = \pi(s)$ is given by

$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E} \left[\frac{\partial Q^\pi(s, a)}{\partial a} \frac{\partial a}{\partial \mathbf{u}} \right]$$

- ▶ if a is continuous and Q is differentiable

Actor-Critic Algorithm

- ▶ Estimate value function $Q(s, a, \mathbf{w}) \approx Q^\pi(s, a)$
- ▶ Update policy parameters \mathbf{u} by stochastic gradient ascent

$$\frac{\partial l}{\partial \mathbf{u}} = \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q(s, a, \mathbf{w})$$

or

$$\frac{\partial l}{\partial \mathbf{u}} = \frac{\partial Q(s, a, \mathbf{w})}{\partial a} \frac{\partial a}{\partial \mathbf{u}}$$

Asynchronous Advantage Actor-Critic (A3C)

- ▶ Estimate state-value function

$$V(s, \mathbf{v}) \approx \mathbb{E}[r_{t+1} + \gamma r_{t+2} + \dots | s]$$

- ▶ Q-value estimated by an n -step sample

$$q_t = r_{t+1} + \gamma r_{t+2} \dots + \gamma^{n-1} r_{t+n} + \gamma^n V(s_{t+n}, \mathbf{v})$$

- ▶ Actor is updated towards target

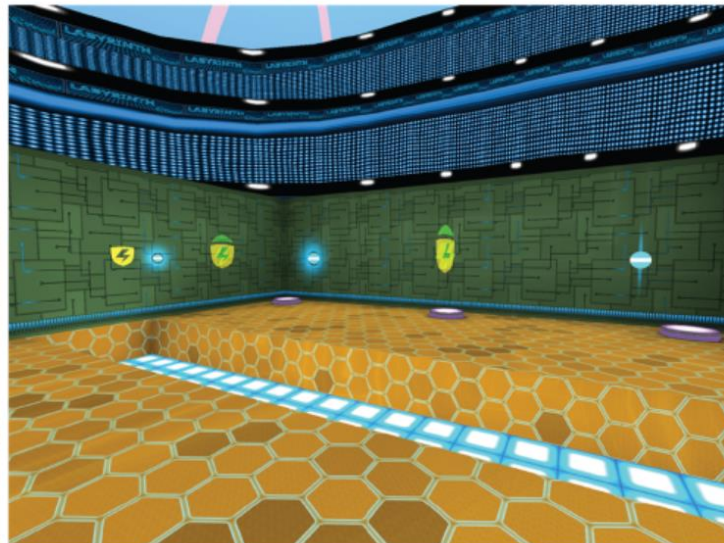
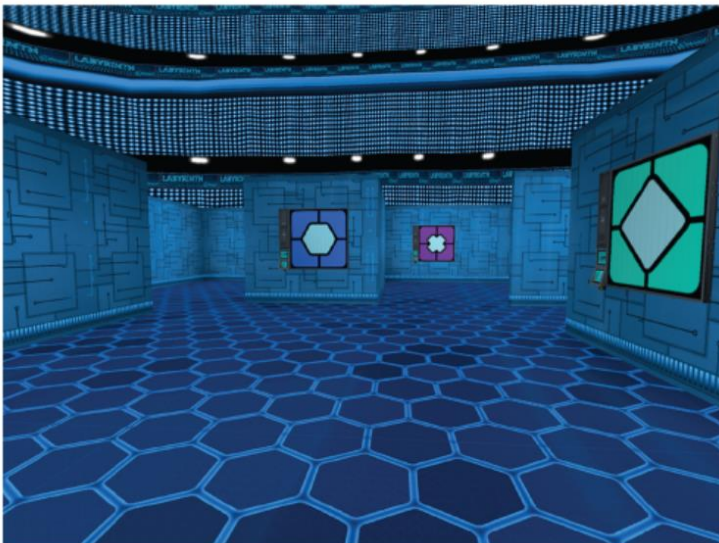
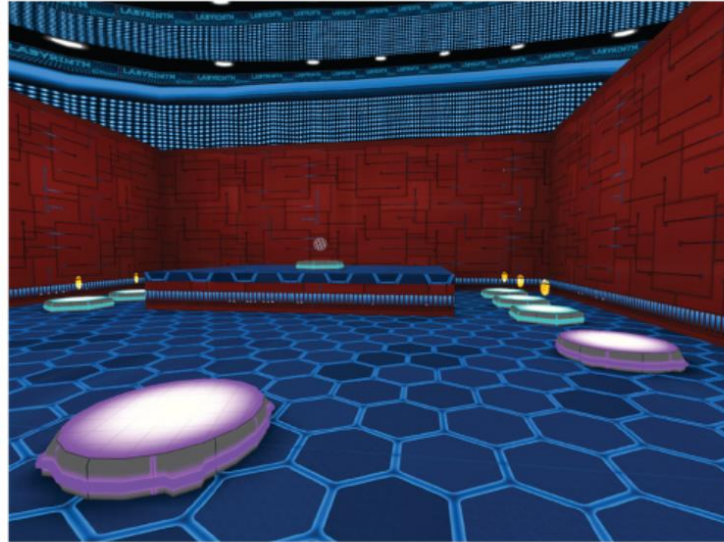
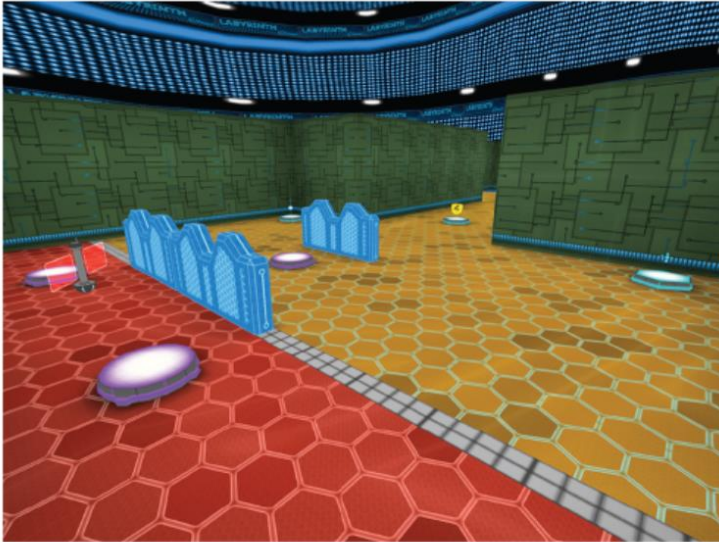
$$\frac{\partial l_u}{\partial \mathbf{u}} = \frac{\partial \log \pi(a_t | s_t, \mathbf{u})}{\partial \mathbf{u}} (q_t - V(s_t, \mathbf{v}))$$

- ▶ Critic is updated to minimise MSE w.r.t. target

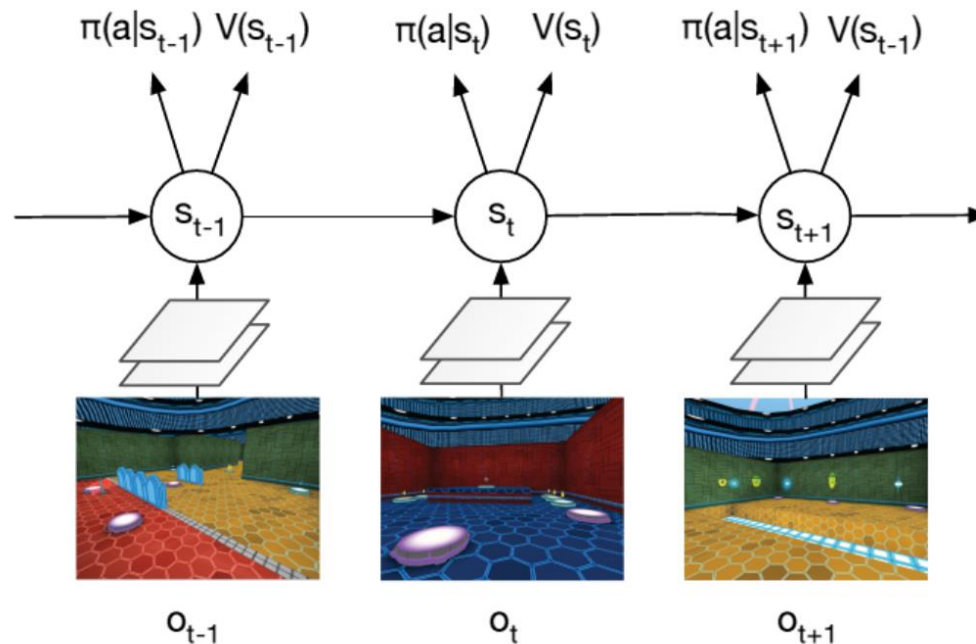
$$l_v = (q_t - V(s_t, \mathbf{v}))^2$$

- ▶ 4x mean Atari score vs Nature DQN

Deep DL in Labyrinth



A3C in Labyrinth



- ▶ End-to-end learning of softmax policy $\pi(a|s_t)$ from pixels
- ▶ Observations o_t are raw pixels from current frame
- ▶ State $s_t = f(o_1, \dots, o_t)$ is a recurrent neural network (LSTM)
- ▶ Outputs both value $V(s)$ and softmax over actions $\pi(a|s)$
- ▶ Task is to collect apples (+1 reward) and escape (+10 reward)

- ◆ Model-based Deep RL

Learning Models of the Environment

- ▶ Demo: generative model of Atari
- ▶ Challenging to plan due to compounding errors
 - ▶ Errors in the transition model compound over the trajectory
 - ▶ Planning trajectories differ from executed trajectories
 - ▶ At end of long, unusual trajectory, rewards are totally wrong

Deep RL in Go

What if we have a perfect model? e.g. game rules are known

AlphaGo paper:

www.nature.com/articles/nature16961

AlphaGo resources:

deepmind.com/alphago/



Example code for Q-Learning

```
def choose_action(state, q_table):
    # This is how to choose an action
    state_actions = q_table.iloc[state, :]
    if (np.random.uniform() > EPSILON) or (state_actions.all() == 0):
        action_name = np.random.choice(ACTIONS)
    else:
        # act greedy
        action_name = state_actions.argmax()
    return action_name

def get_env_feedback(S, A):
    # This is how agent will interact with the environment
    if A == 'right':
        # move right
        if S == N_STATES - 2:
            # terminate
            S_ = 'terminal'
            R = 1
        else:
            S_ = S + 1
            R = 0
    else:
        # move left
        R = 0
        if S == 0:
            S_ = S
            # reach the wall
        else:
            S_ = S - 1
    return S_, R
```

Example Code

```
def rl():
    # main part of RL loop
    q_table = build_q_table(N_STATES, ACTIONS)
    for episode in range(MAX_EPISODES):
        step_counter = 0
        S = 0
        is_terminated = False
        update_env(S, episode, step_counter)
        while not is_terminated:

            A = choose_action(S, q_table)
            S_, R = get_env_feedback(S, A)  # take action & get next state & reward
            q_predict = q_table.ix[S, A]
            if S_ != 'terminal':
                q_target = R + GAMMA * q_table.iloc[S_, :].max()
            else:
                q_target = R  # next state is terminal
                is_terminated = True  # terminate this episode

            q_table.ix[S, A] += ALPHA * (q_target - q_predict)  # update
            S = S_  # move to next state

        update_env(S, episode, step_counter+1)
        step_counter += 1
    return q_table
```

Thanks!

