

# Introduction to Reinforcement Learning

## Lecture 2: Frontiers of RL

Joshua Evans

Bath Reinforcement Learning Laboratory

Department of Computer Science



UNIVERSITY OF  
**BATH**

Dynamic  
Programming

Monte Carlo  
Methods

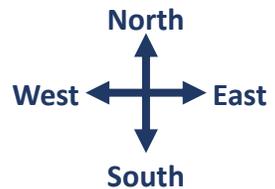
Temporal Difference  
Methods

Learn  $V(s)$  or  $Q(s, a)$ , then derive a policy.

# Tabular Methods

- So far, our value functions could all be represented as tables.
  - For  $V(s)$ , one entry per state.
  - For  $Q(s, a)$ , one entry per state-action pair.

13	14	15	16 
9	10	11	12
5	6	7	8
1 	2	3	4



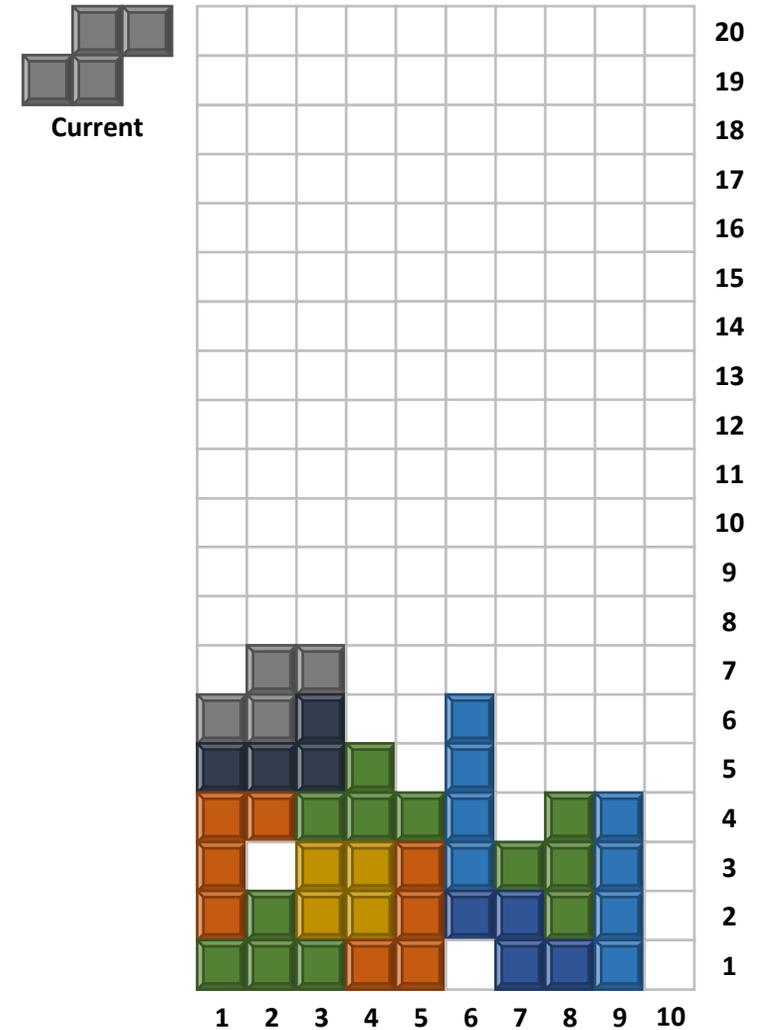
State	Action	$Q(s, a)$
1	N	$Q(0, N)$
1	S	$Q(0, S)$
1	E	$Q(0, E)$
1	W	$Q(0, W)$
...	...	...
16	W	$Q(24, W)$

16 states  $\times$  4 actions  $\rightarrow$  64 state-action pairs

- We need to store a value in memory for each state-action pair.
- We need to experience a state-action pair to learn about it.
- Feasible here, but what about larger problems?

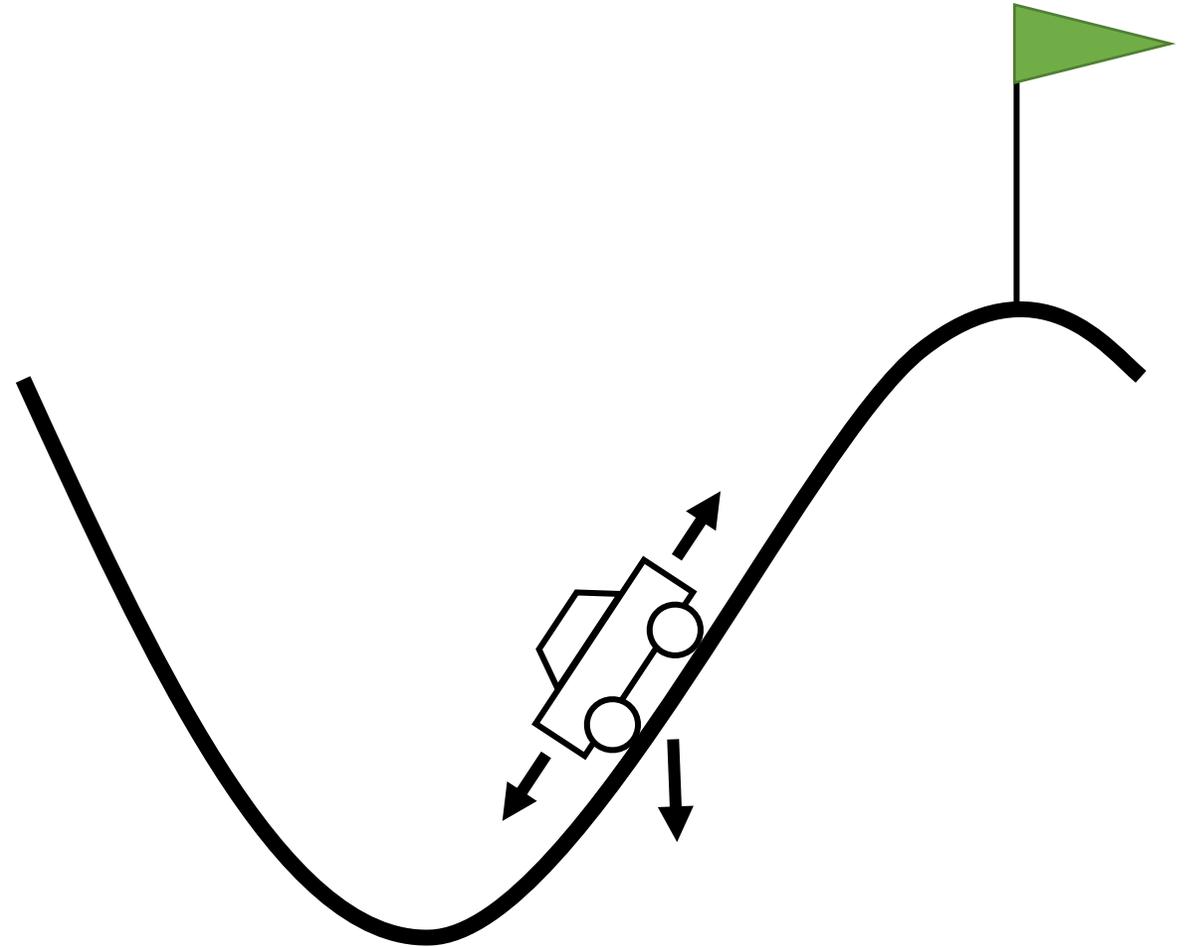
# Larger State-Spaces

- Tetris
  - $10 \times 20$  grid, 7 current tetriminos.
  - Upper bound of  $7 \times 2^{200}$  states.
- Not feasible to learn to play Tetris using tabular methods.
  - Can't fit the Q-table in memory.
  - Cannot reasonably experience all state-action pairs.



# Larger State-Spaces

- Mountain Car Problem
  - Continuous State-Space
  - Position ( $-1.2 < x < 0.5$ )
  - Velocity ( $-0.07 < v < 0.07$ )
- Actions:
  - Full Throttle Forward
  - Full Throttle Backwards
  - Zero Throttle
- Infinite number of states!



## Mountain Car Problem

The underpowered car must reach the top of the hill as quickly as possible.

# Generalisation

- We can only reasonably expect to experience a small subset of the entire state-action-space.
- Can we **generalise** what we know about states we have visited to states we haven't visited?
- Can we do this in a way which avoids storing information about every state separately?

Yes, using **function approximation!**

# Function Approximation

- Function Approximated Methods
  - Represent  $v_{\pi}(s)$  as a real-valued function with parameter vector  $\theta$ .
  - We only store and update  $\theta$ .

$$\hat{v}(s, \theta) \approx v_{\pi}(s)$$

- Function could be a linear function, polynomial, **neural network...**
- Learn a function with parameters  $\theta$  which maps states to values.
  - How can we go about doing this?

# Approximating Value Functions

## State-Values

$$s \mapsto v$$

- Dynamic Programming

$$s \mapsto E_{\pi}\{R_{t+1} + \gamma \hat{v}(S_{t+1}, \boldsymbol{\theta}_t) | S_t = s\}$$

- Monte-Carlo

$$s \mapsto G_t$$

- Temporal-Difference

$$s \mapsto R_{t+1} + \gamma \hat{v}(S_{t+1}, \boldsymbol{\theta}_t)$$

## Action-Values

$$s, a \mapsto q$$

- Dynamic Programming

$$s, a \mapsto E_{\pi}\{\dots | S_t = s, A_t = a\}$$

- Monte-Carlo

$$s, a \mapsto G_t$$

- Temporal-Difference

$$s, a \mapsto R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \boldsymbol{\theta}_t)$$

We can use these as training examples in a supervised learning context!

# Approximating Value Functions

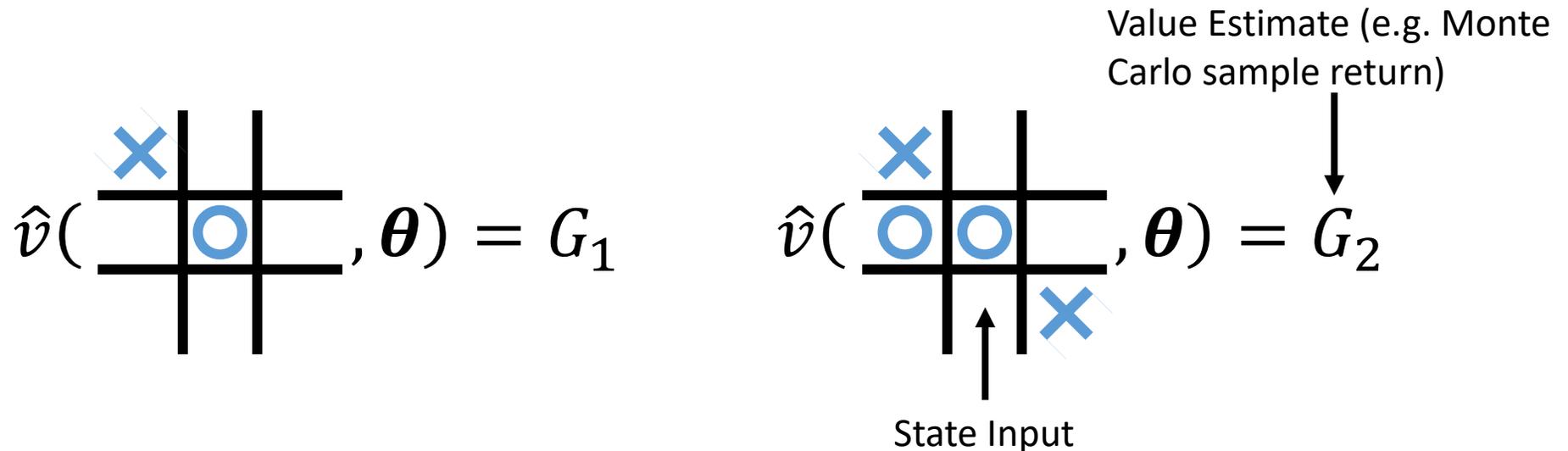
- Our goal is to train a function to correctly map inputs to outputs.
  - We do this all the time in supervised learning!

$$f(\text{Image of a Cat}, \theta) = 0 \quad f(\text{Image of a Dog}, \theta) = 1$$

Change  $\theta$  so that  $f$  maps images to classes correctly.

# Approximating Value Functions

- Our goal is to train a function to correctly map inputs to outputs.
  - We do this all the time in supervised learning!
  - Here, our inputs are states, and our outputs are values.

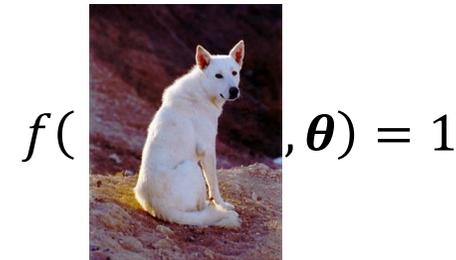
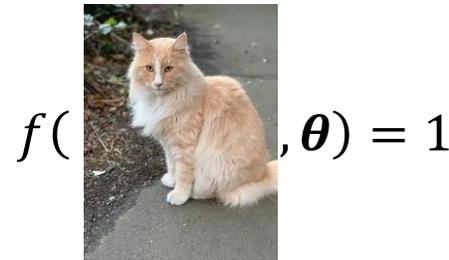
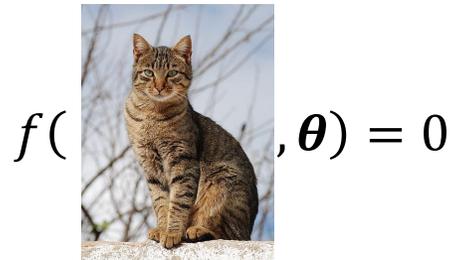


Change  $\theta$  so that  $\hat{v}$  maps states to values correctly.

# Approximating Value Functions

- We can measure how well our function is mapping inputs to outputs using objective functions.

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m [y_i - \hat{y}_i]^2$$



# Approximating Value Functions

- We can measure how well our function is mapping inputs to outputs using objective functions.

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m [y_i - \hat{y}_i]^2$$

                  ↑                  ↑  
                  True          Predicted  
                  Output      Output

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m [(v_{\pi}(S_t) - \hat{v}(S_t, \theta_t))]^2$$

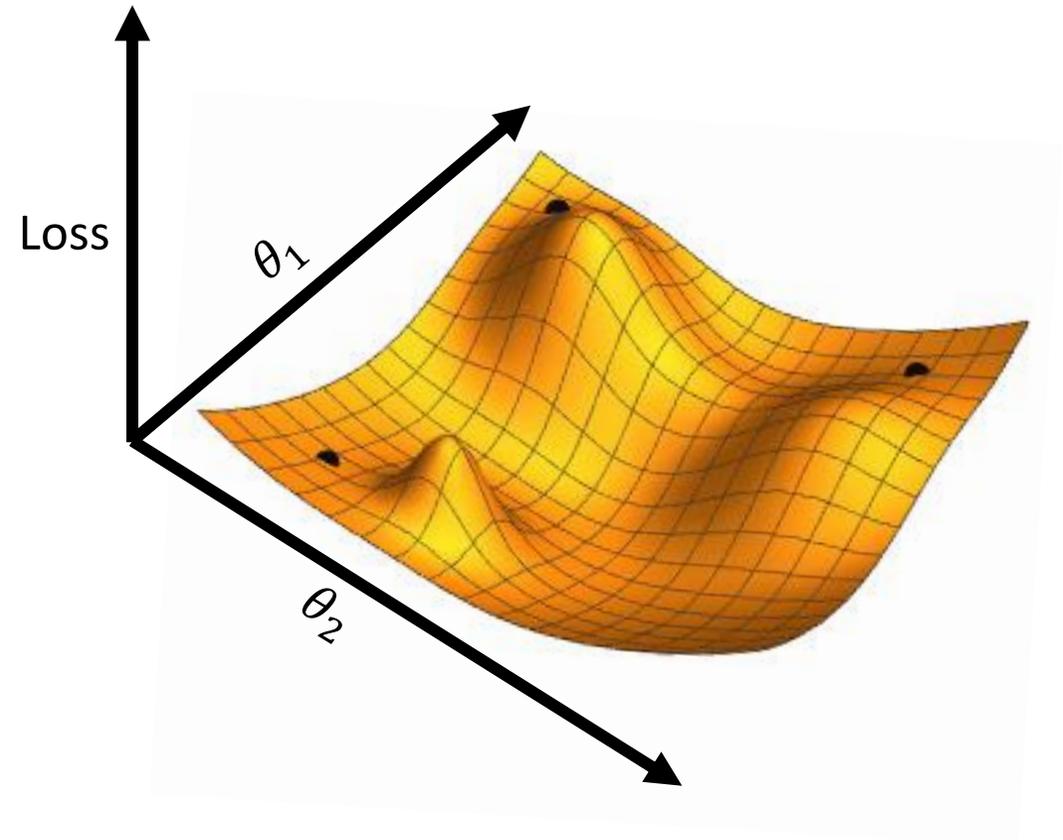
                  ↑                  ↑  
                  New Value      Existing Value  
                  Estimate      Estimate



# Approximating Value Functions

- How do we know how to change  $\theta$  in order to lower our loss?
- Calculate the gradient of our loss function with respect to each of the parameters in  $\theta$ .
- Take a step in the direction which lowers the loss.

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta_t} \text{MSE}$$



# Function Approximation

- Function Approximated Methods
  - Represent  $v_{\pi}(s)$  as a real-valued function with parameter vector  $\theta$ .
  - We only store and update  $\theta$ .

$$\hat{v}(s, \theta) \approx v_{\pi}(s)$$

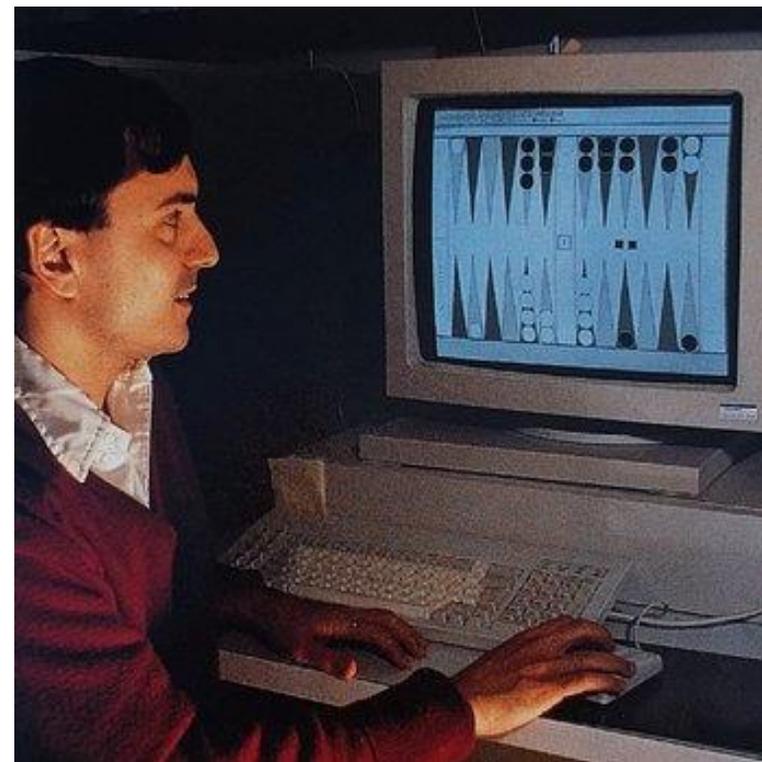
Deep RL



- Function could be a linear function, polynomial, **neural network...**
- Learn a function with parameters  $\theta$  which maps states to values.
  - How can we go about doing this?

# A Brief History of Deep RL

- 1992 – Gerald Tesauro's TD-Gammon



# A Brief History of Deep RL

- 2009 – Riedmiller et al. Robocup



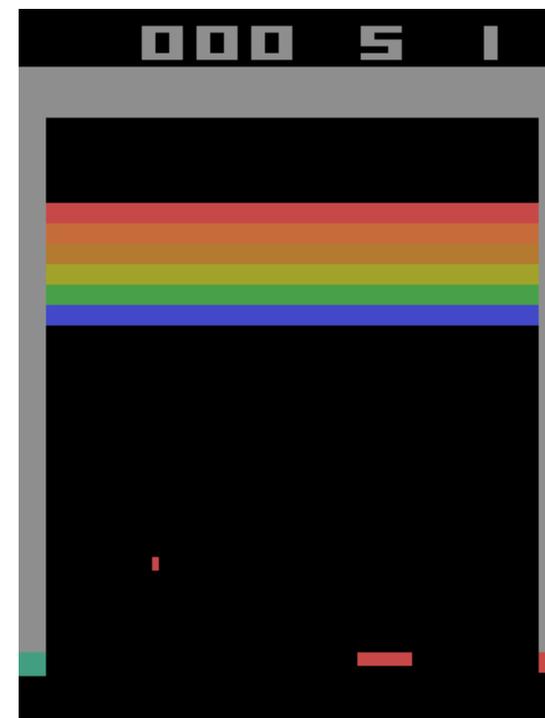
# A Brief History of Deep RL

- 2009 – Riedmiller et al. Robocup



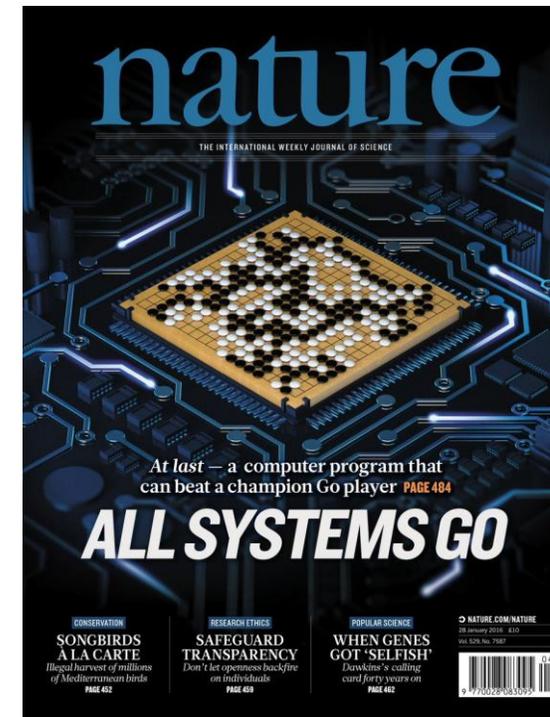
# A Brief History of Deep RL

- 2013-2015 – DeepMind train agents to play Atari games.



# A Brief History of Deep RL

- 2016 – DeepMind train agents to beat world champions at Go.



# More Recent Results

- 2017-2019 – OpenAI Five (DotA 2 Agents)

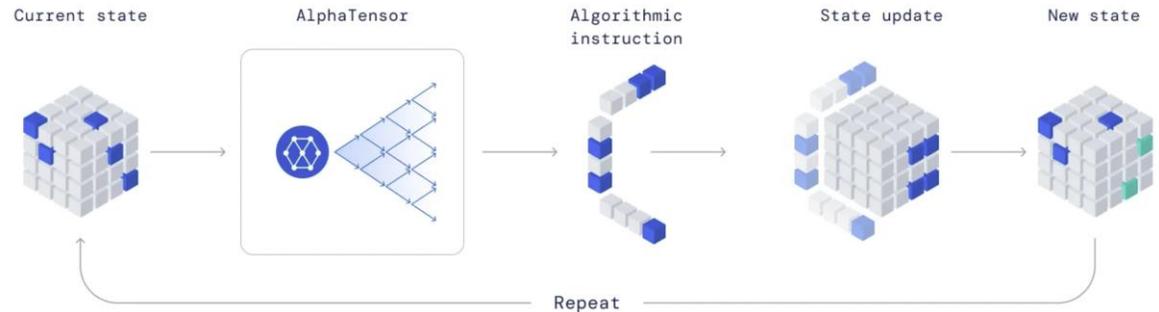




# Cutting-Edge Results

## AlphaTensor (2022)

Discovery of novel, efficient, and provably-correct algorithms.



## Autonomous Driving (Present)

Using deep reinforcement learning to help guide autonomous vehicles.



# The Obligatory “RLHF” Slide

Step 1

Collect demonstration data and train a supervised policy.

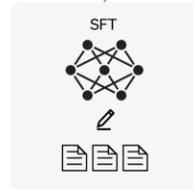
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3.5 with supervised learning.



Step 2

Collect comparison data and train a reward model.

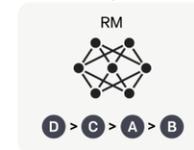
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

A new prompt is sampled from the dataset.



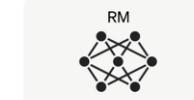
The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.



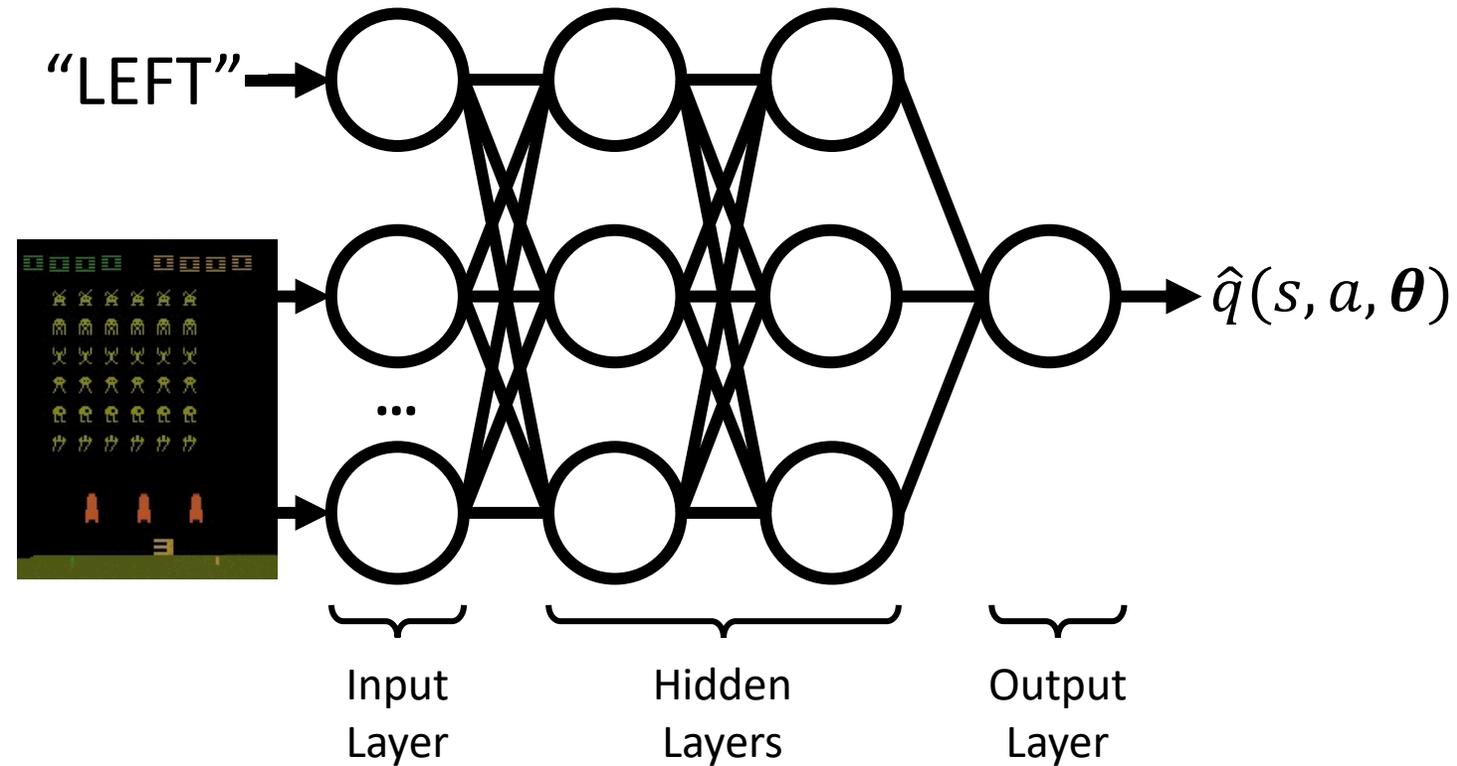
The reward is used to update the policy using PPO.



## Reinforcement Learning from Human Feedback (Present)

Use human preferences to train a reward model. Then, use that reward model to fine-tune a GPT model with reinforcement learning methods.

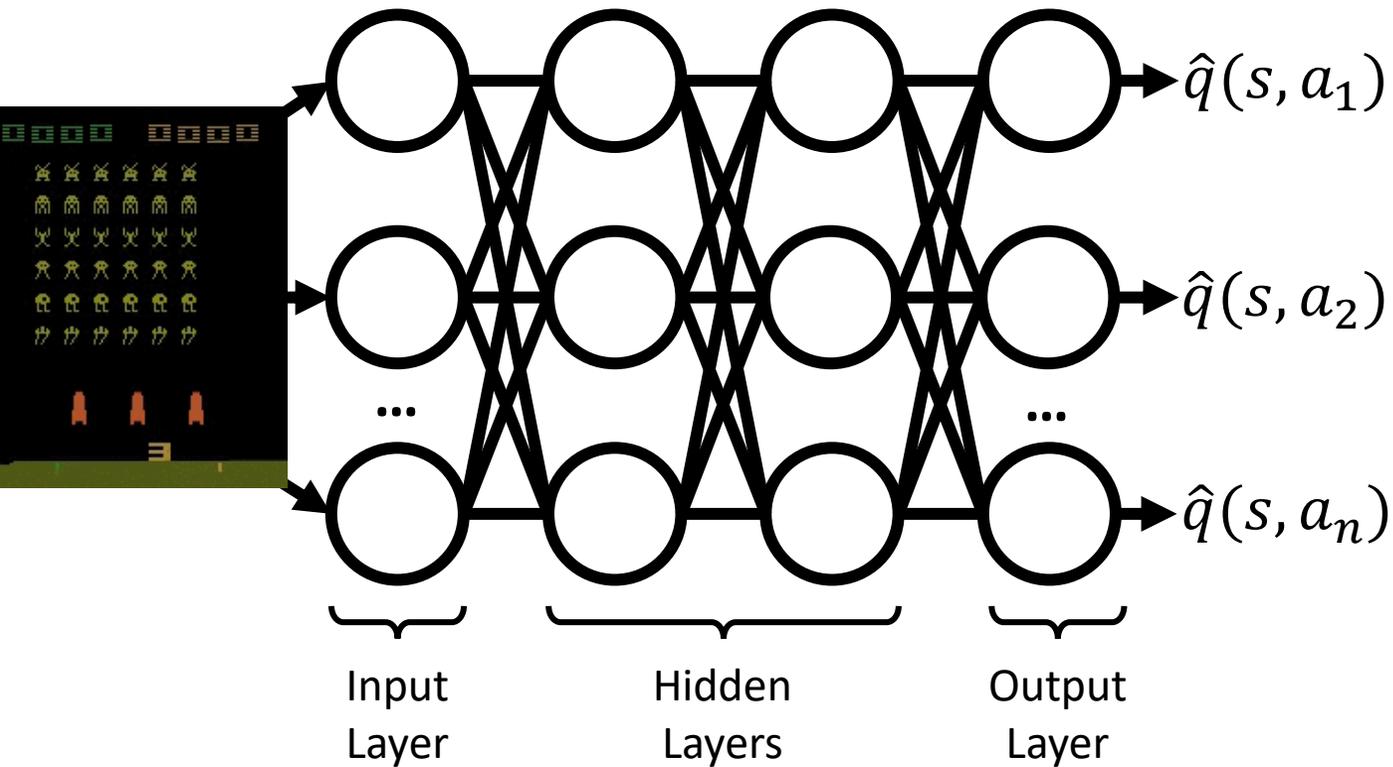
# Deep Q-Networks (DQN)



What should our neural network look like?

- We could have state-action pairs as input, and a single action-value as output.

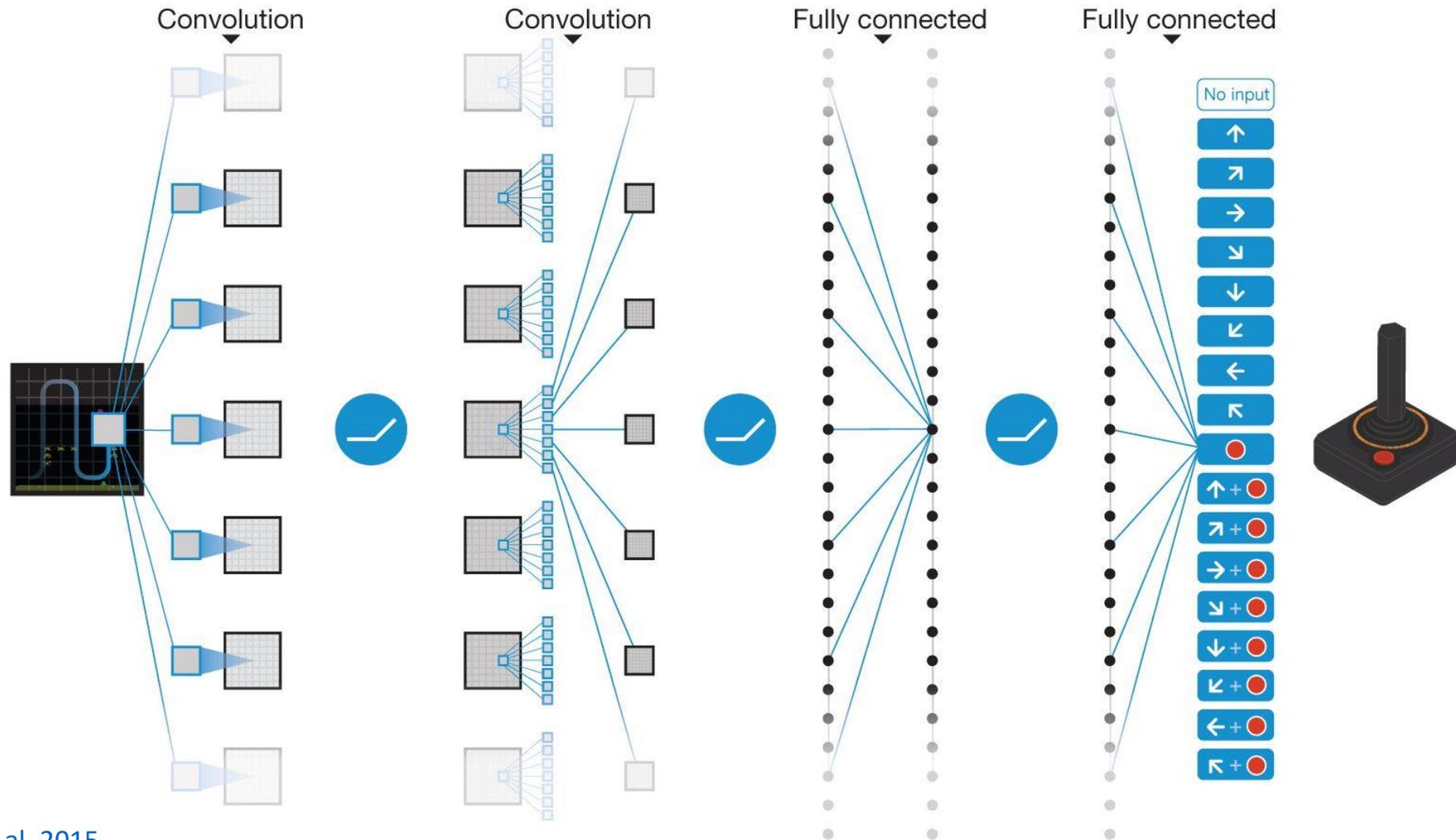
# Deep Q-Networks (DQN)

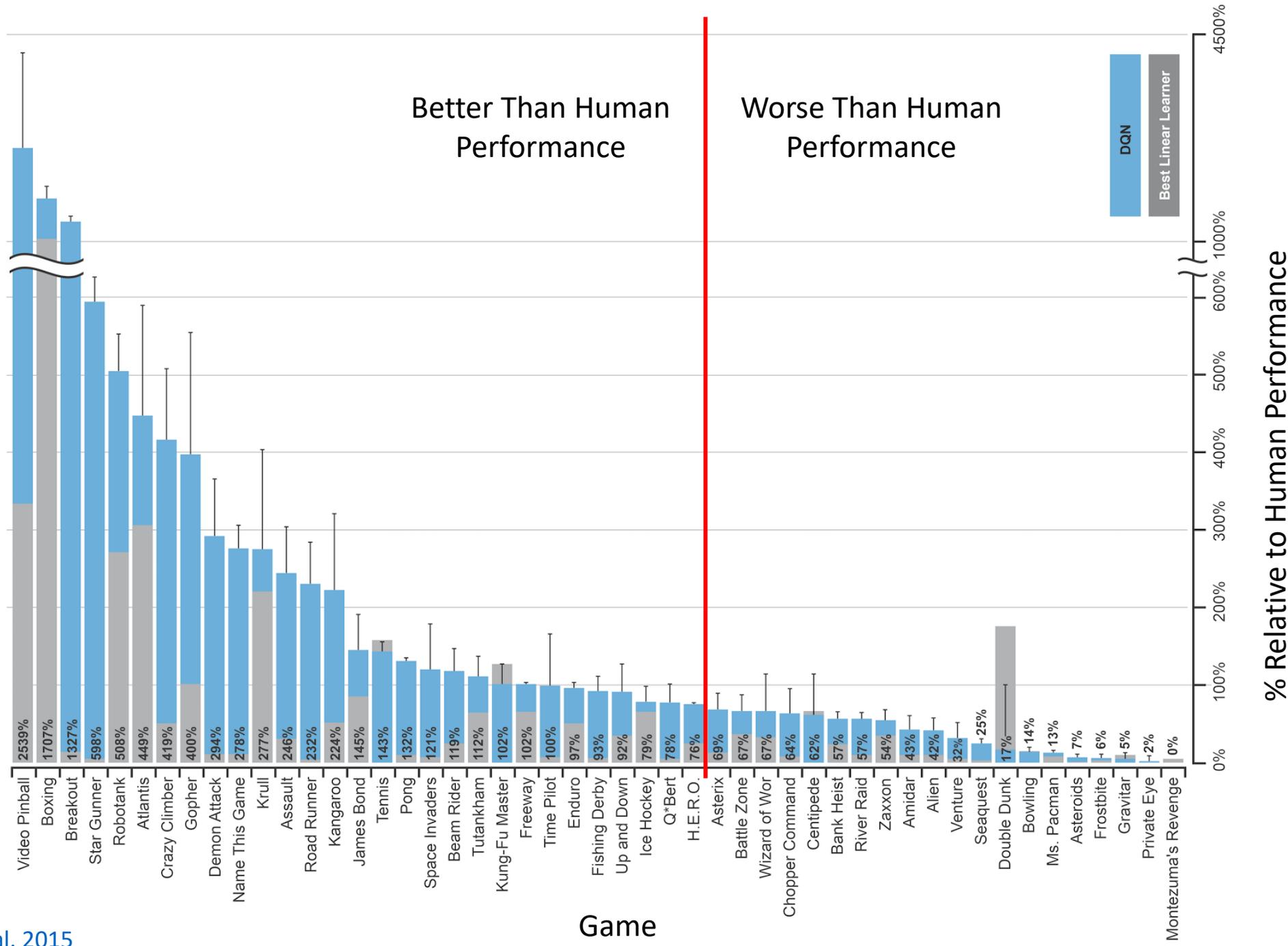


What should our neural network look like?

- We could have state-action pairs as input, and a single action-value as output.
- We could have only the state as input, then one output per available action as outputs.

# Deep Q-Networks (DQN)







# Ridiculously Naïve DQN Issues

- Data Highly Temporally Correlated
  - Samples are taken sequentially – distributional shift over time.
- The agent's policy could change rapidly.
  - Large errors will lead to large gradients and large updates.
  - Even small changes in Q-values may drastically change the policy.
  - This may cause sudden changes to the sampling distribution.
- We're using the same function to compute both our prediction ( $\hat{y}$ ) and our target ( $y$ ).
  - Whenever we perform an update, our target will change.
  - We're chasing a moving target – can lead to unstable training.

# DQN Fixes: Experience Replay

- Use **experience replay**.
  - Store a buffer  $D$  of experiences of the form  $(s_t, a_t, r_{t+1}, s_{t+1})$ .
  - Perform updates using a minibatch of experiences from the buffer.
- Decorrelates experiences temporally.
  - Recent experiences are not emphasised during training.
  - Old experiences are not forgotten.
- Requires a large amount of memory.
  - In excess of 1,000,000 recent experiences may need to be stored.
- Extensions exist to prioritise experiences and remove old ones.

# DQN Fixes: Reward Clipping/Normalising

- Reward magnitude not known ahead of time.
- Rewards with large magnitudes may cause dramatic weight updates, impacting stability.
- To deal with this, we could clip rewards to lie in the range  $[-1,1]$ .
- We can also clip the gradients of the loss function each before performing an update.
- We could also track the rewards we observe and normalise them as we go along.

# DQN Fixes: Fixed Target Network

- We currently use the same network  $\hat{q}(s, a, \theta)$  to select actions and compute our target value.
  - As we update our network, the target value will also change.
  - The network is updating towards a moving target!
  - This can cause instability when training, and decrease overall performance.
- To fix this, we can use a **fixed target network**.
  - Select actions using  $\hat{q}_1(s, a, \theta)$  with parameters  $\theta_1$ .
  - Calculate target values using  $\hat{q}_2(s, a, \theta_2)$  with parameters  $\theta_2$ .
  - Every  $C$  steps, update  $\theta_2 = \theta_1$ .
  - Allows us to perform update towards a (mostly) stationary target.

## Algorithm: Deep Q-Learning with Experience Replay and Fixed Target Network

Initialise replay memory  $D$  to capacity  $N$  ← Initialise Replay Buffer

Initialise action-value network  $\hat{q}_1$  with arbitrary weights  $\theta_1$

Initialise target action-value network  $\hat{q}_2$  with weights  $\theta_2 = \theta_1$  ← Initialise Target Network

**For** episode = 1,  $M$  **do**

    Initialise initial state  $S_1$

**For**  $t = 1, T$  **do**

        With probability  $\epsilon$  select random action  $A_t$

        With probability  $1 - \epsilon$  select action  $A_t = \operatorname{argmax}_a \hat{q}_1(S_t, a, \theta_1)$  ← Select action using  $\hat{q}_1$ .

        Execute action  $A_t$ , observe reward  $R_{t+1}$  and next state  $S_{t+1}$

        Store transition  $(S_t, A_t, R_{t+1}, S_{t+1})$  in  $D$  ← Store Experience in Replay Buffer

        Sample random minibatch of transitions  $(S_j, A_j, R_{j+1}, S_{j+1})$  from  $D$

        Set  $y_j = \begin{cases} R_{j+1} + 0, & \text{if } S_{j+1} \text{ is terminal} \\ R_{j+1} + \gamma \max_{a'} \hat{q}_2(S_{j+1}, a', \theta_2), & \text{otherwise} \end{cases}$  ← Compute TD target using  $\hat{q}_2$  (target network).      Update based on experience sampled from replay buffer.

        Perform gradient descent step  $\nabla_{\theta} (y_j - \hat{q}(S_j, A_j, \theta))^2$

        Every  $C$  steps, update  $\theta_2 = \theta_1$  ← Periodically update target network.

**End For**

**End For**

# Types of Reinforcement Learning Methods

- **Value-Based** RL Methods
  - Approximate the **optimal action-value function**  $Q^*(s, a)$ .
- **Policy-Based** RL Methods
  - Directly search the policy-space for the **optimal policy**  $\pi^*(a|s)$ .

**Deep RL methods use deep neural networks to represent the value function or policy.**

# Why Policy-Based Methods?

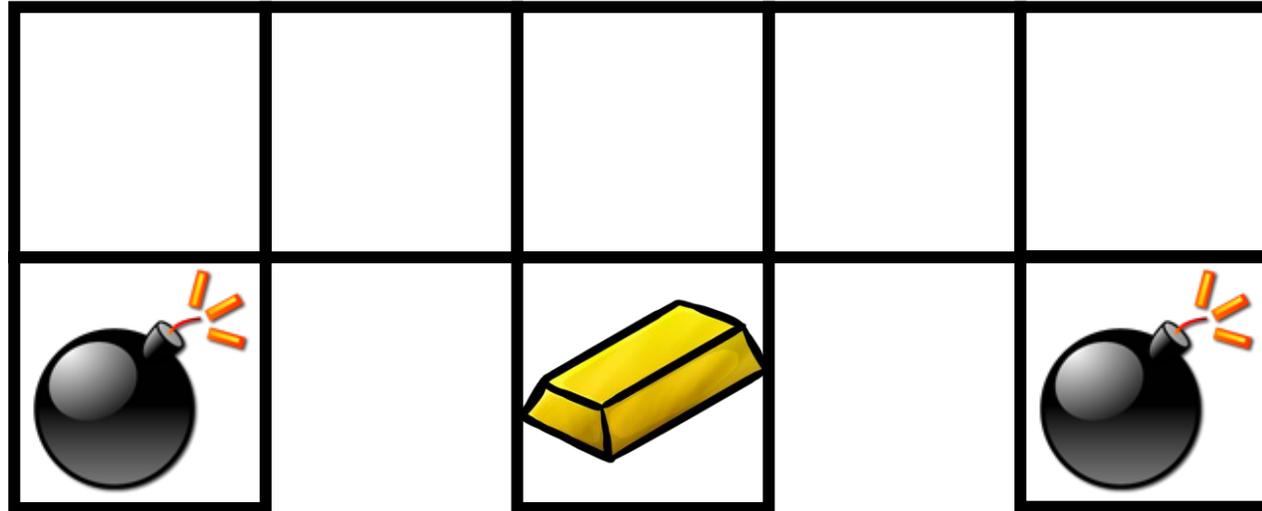
- So far, we have worked with **value-based** methods.
  - We'd learn the action-value function, then derive a policy (e.g.  $\epsilon$ -greedy).
- What if the optimal policy is **stochastic**?
  - Value-based methods have no natural way of dealing with this.



- Instead, could we learn the optimal policy directly?

# Aliased Gridworld

State = (Wall to North, Wall to South, Wall to East, Wall to West)



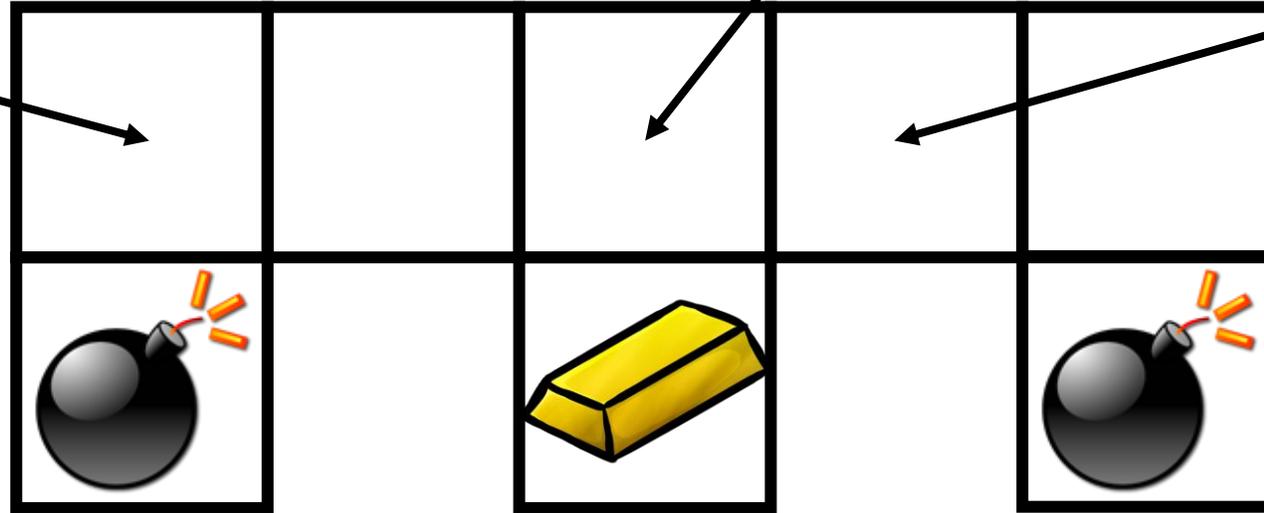
# Aliased Gridworld

State = (Wall to North, Wall to South, Wall to East, Wall to West)

(TRUE, FALSE, FALSE, FALSE)

(TRUE, TRUE, FALSE, FALSE)

(TRUE, FALSE, FALSE, TRUE)

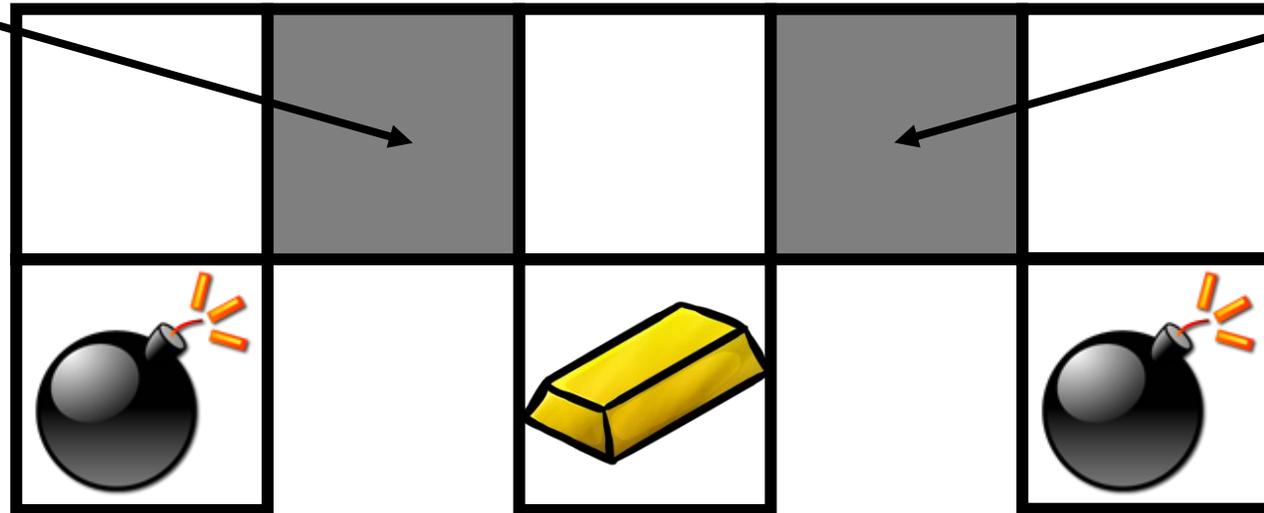


# Aliased Gridworld

State = (Wall to North, Wall to South, Wall to East, Wall to West)

(TRUE, TRUE, FALSE, FALSE)

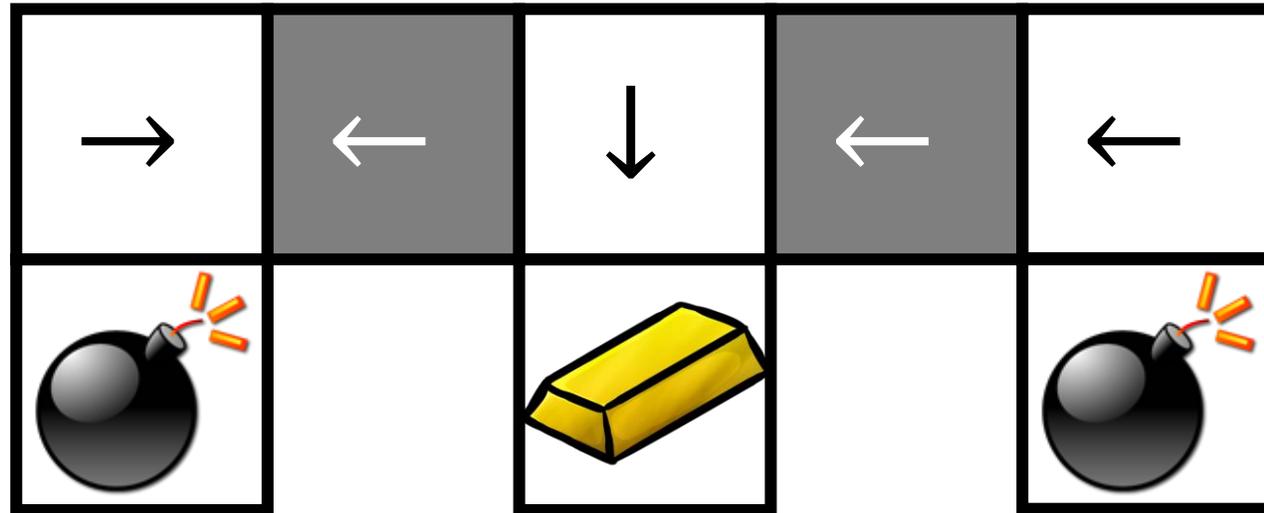
(TRUE, TRUE, FALSE, FALSE)



These two states have identical representations!

# Aliased Gridworld

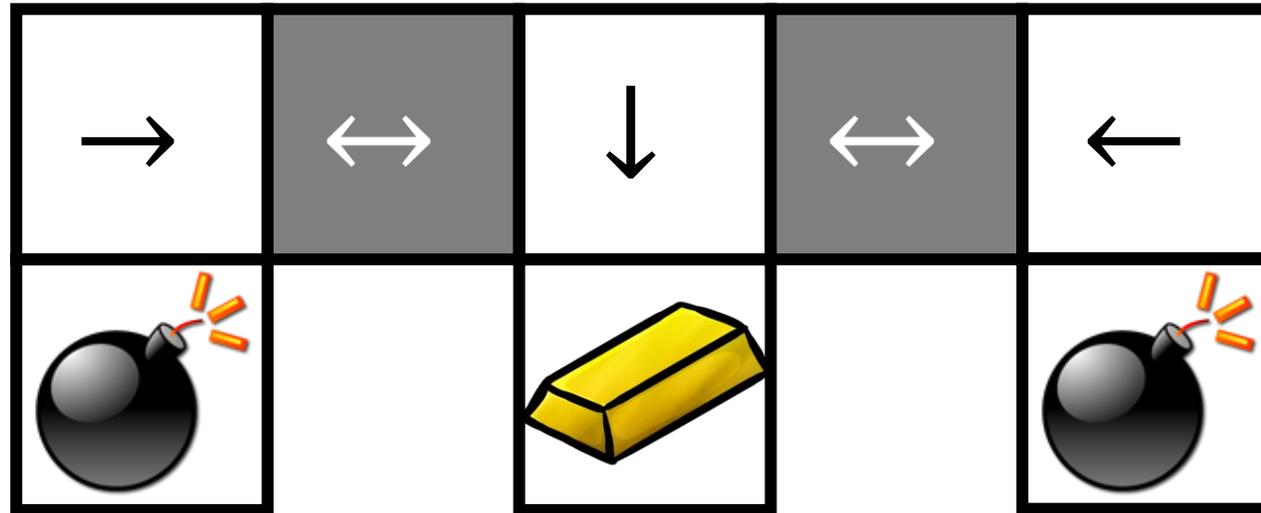
State = (Wall to North, Wall to South, Wall to East, Wall to West)



These two states have identical representations!  
A deterministic policy would get stuck.

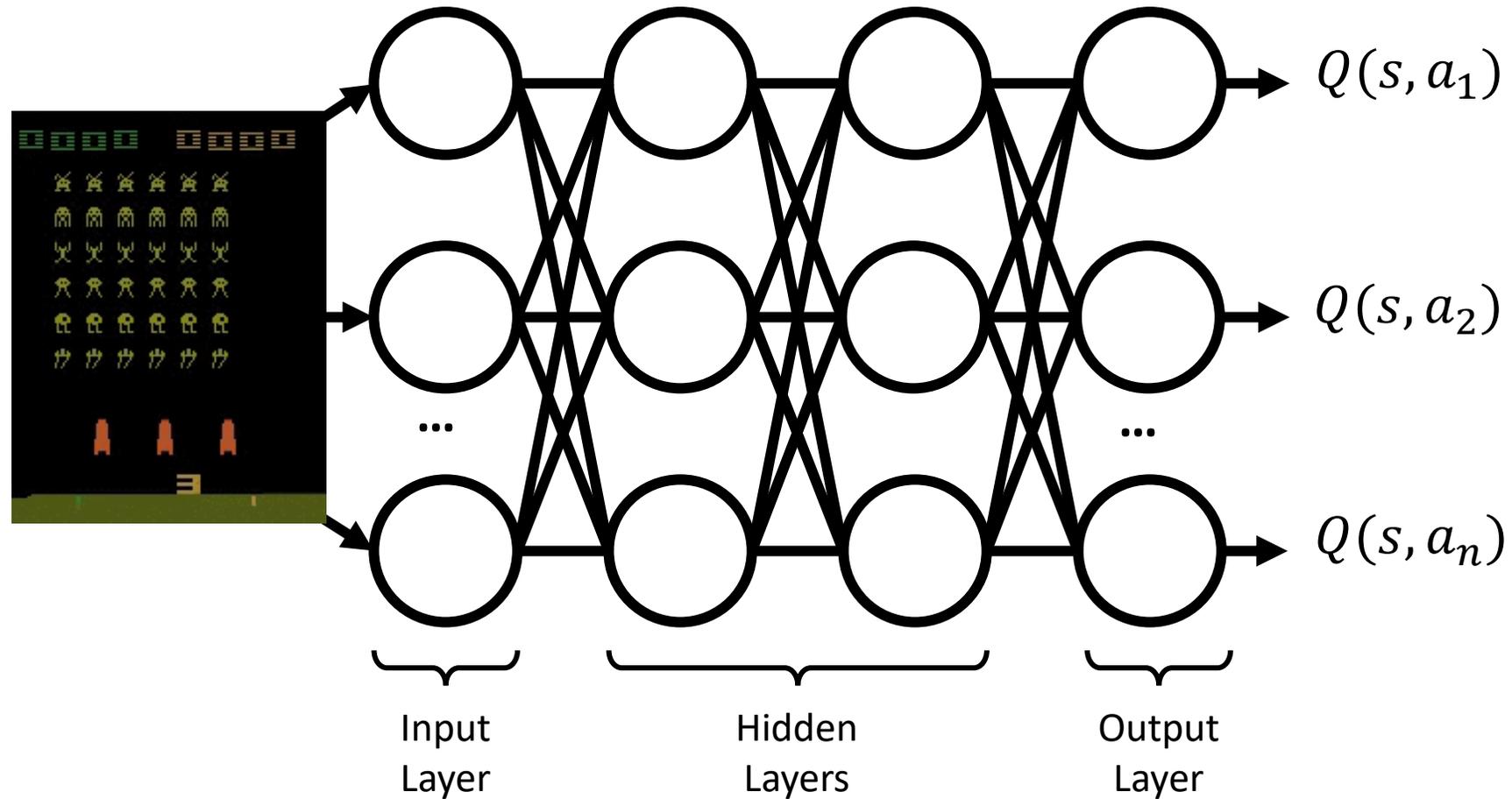
# Aliased Gridworld

State = (Wall to North, Wall to South, Wall to East, Wall to West)

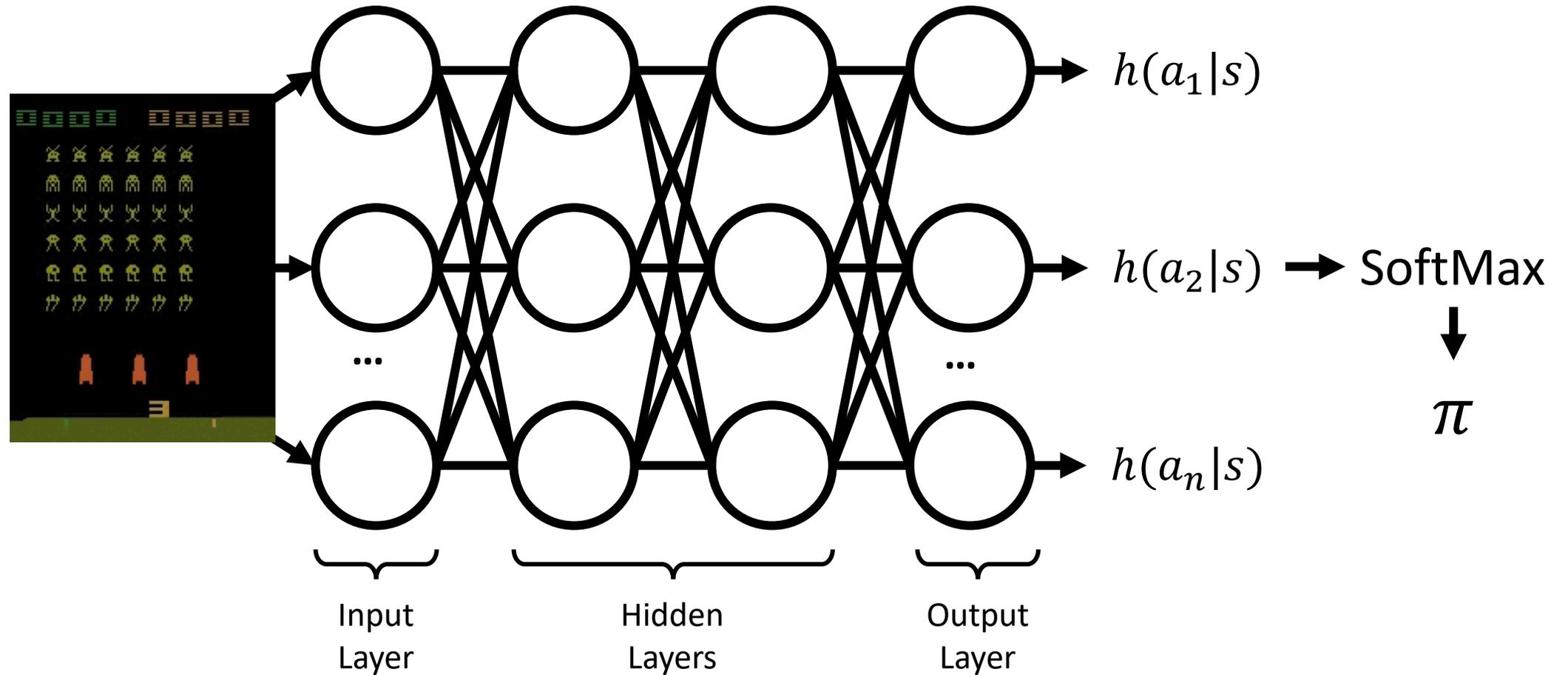


These two states have identical representations!  
A stochastic policy would work much better!

# Q-Network



# Policy Network



# Policy Gradient Methods

$\pi_{\theta}(a|s)$  = Probability of choosing action  $a$  given state  $s$  and policy parameters  $\theta$ .

- $J(\pi_{\theta})$  is some **objective function** for our policy, which we aim to maximise (e.g. expected return,  $E_{\pi_{\theta}}[G_t]$ ).
- Policy gradient update rule:  $\theta_{t+1} = \theta_t + \alpha \nabla_{\theta_t} J(\pi_{\theta_t})$
- $\nabla_{\theta} J(\pi_{\theta})$  is called the **policy gradient**.

# Policy Gradient Methods

- Policy gradient update:  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \nabla_{\boldsymbol{\theta}_t} J(\pi_{\boldsymbol{\theta}_t})$ 
  - Where our objective function,  $J(\pi_{\boldsymbol{\theta}_t})$ , is the discounted return,  $E_{\pi_{\boldsymbol{\theta}}} [G_t]$ .
- **Problem:** We don't have direct access to, the gradient of the discounted return w.r.t. our policy parameters!
- So, to actually use this update rule in an algorithm, we need an expression for the policy gradient which we can numerically compute.
  - This expression should be computable using only  $\pi, \nabla_{\boldsymbol{\theta}} \pi, \boldsymbol{\theta}_t, J \dots$
  - ...and a **trajectory**  $\tau$  of our agent's experience.

$$\begin{aligned}
\nabla_{\theta} J(\pi_{\theta}) &= \nabla_{\theta} E_{\tau \sim \pi_{\theta}} [G_t] \\
&= \nabla_{\theta} \int_{\tau} P(\tau | \theta) G_t && \text{Expand expectation.} \\
&= \int_{\tau} \nabla_{\theta} P(\tau | \theta) G_t && \text{Bring gradient under integral.} \\
&= \int_{\tau} P(\tau | \theta) \nabla_{\theta} \log P(\tau | \theta) G_t && \text{Log-derivative trick.} \\
&= E_{\tau \sim \pi} [\nabla_{\theta} \log P(\tau | \theta) G_t] && \text{Return to expectation form.} \\
&= E_{\tau \sim \pi} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi(A_t | S_t) G_t \right] && \text{Substitute expression for grad-log-prob.}
\end{aligned}$$

High-level, hand-wavy intuition:

- Push up the preferences of actions that lead to high returns.
- Push down the preferences of actions that lead to low returns.

We can approximate this using the sample mean over many trajectories  $\tau \in D$ :

$$\nabla_{\theta} J(\pi_{\theta}) \approx \frac{1}{|D|} \sum_{\tau \in D} \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t$$

## Algorithm: REINFORCE

Initialise parameters: step size  $\alpha \in (0,1]$

Initialise policy network  $\pi$  with parameters  $\theta$

**For** episode = 1,  $M$  **do**

    Generate an episode trajectory  $\tau \sim \pi_\theta$

**For**  $t = 1, T - 1$  **do**

$G \leftarrow \sum_{k=t+1}^T R_k$  ← Sample Return

$\theta \leftarrow \theta + \alpha \underbrace{\nabla_{\theta} \log \pi_{\theta}(A_t | S_t)}_{\text{Our Policy-Gradient!}} G$

**End For**

**End For**

Our Policy-Gradient!

# Actor-Critic Methods

- With REINFORCE, we are stuck with performing Monte Carlo updates.
- If we learn a **policy function** and a **value function**, we can use the value function to do bootstrapping, letting us perform TD updates.
- This gives us all the previous benefits that we've seen from bootstrapping, and the benefits of policy-based methods.
- We call bootstrapping methods that learn the policy function directly while using estimates from value functions **Actor-Critic Methods**.
  - The policy function is the “**Actor**”.
  - The value function is the “**Critic**”.

# Types of Reinforcement Learning Methods

- **Value-Based RL**
  - Approximates the **optimal action-value function**  $Q^*(s, a)$ .
- **Policy-Based RL**
  - Directly search the policy-space for the **optimal policy**  $\pi^*(a|s)$ .

Actor-Critic  
methods  
do both!

**Deep RL uses deep neural networks to represent the value function or policy**

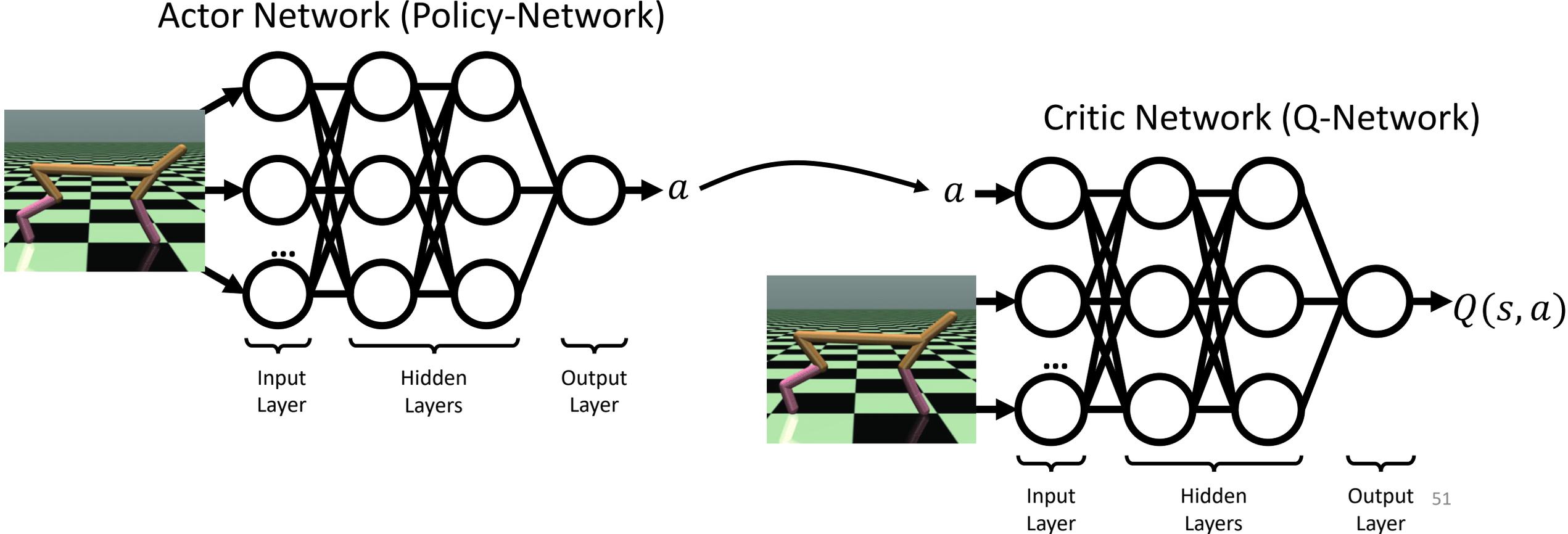
# Deep Deterministic Policy Gradients

- DDPG is an actor-critic algorithm for **continuous action-spaces**.
- It makes use of many of DQN's tricks, such as replay buffers and target networks.
- It is off-policy, so can make use of old experiences.
- It makes policy-gradient updates maximising  $E[Q(s, a)]$ .

Recall: We maximised  $E[G_t]$  earlier!



# DDPG Network Architecture



## Algorithm: Deep Deterministic Policy Gradients (DDPG)

Initialise replay memory  $D$  to capacity  $N$  ← Pre-Fill Replay Buffer

Initialise critic network  $Q$  with random weights  $\theta^Q$ , actor network  $\pi$  with weights  $\theta^\pi$

Initialise target critic network  $\hat{Q}$  with weights  $\hat{\theta}^Q = \theta^Q$ , target actor network  $\hat{\pi}$  with weights  $\hat{\theta}^\pi = \theta^\pi$

Initialise target network learning rate  $\beta \in (0,1]$

**For** episode = 1,  $M$  **do**

Initialise random process  $\mathcal{N}$  for action exploration ← We use random noise for exploration.

Initialise initial state  $s_1$

**For**  $t = 1, T$  **do**

Select action  $a_t = \pi(s_t) + \mathcal{N}_t$  ←

Execute action  $a_t$  and observe reward  $r_{t+1}$ , next state  $s_{t+1}$

Store transition  $(s_t, a_t, r_{t+1}, s_{t+1})$  in  $D$  ← Store Experience in Replay Buffer

Sample random minibatch of transitions  $(s_j, a_j, r_{j+1}, s_{j+1})$  from  $D$  ← Sample From Replay Buffer

Set  $y_j = r_{j+1} + \gamma \hat{Q}(s_{j+1}, \hat{\pi}(s_{j+1}))$  ← Generate Target Using Target Networks

Perform gradient descent step  $\nabla_{\theta^Q} (y_j - Q(s_j, a_j))^2$  on critic ← Critic update very similar to DQN's.

Perform gradient ascent step  $\nabla_{\theta^\pi} \mathbb{E} [Q(s_j, \pi(s_j))]$  on actor ← Actor update similar to REINFORCE's.

Update target networks  $\hat{\theta}^Q \leftarrow \beta \theta^Q + (1 - \beta) \hat{\theta}^Q$ ,  $\hat{\theta}^\pi \leftarrow \beta \theta^\pi + (1 - \beta) \hat{\theta}^\pi$  ← Update Target Networks

**End For**

**End For**

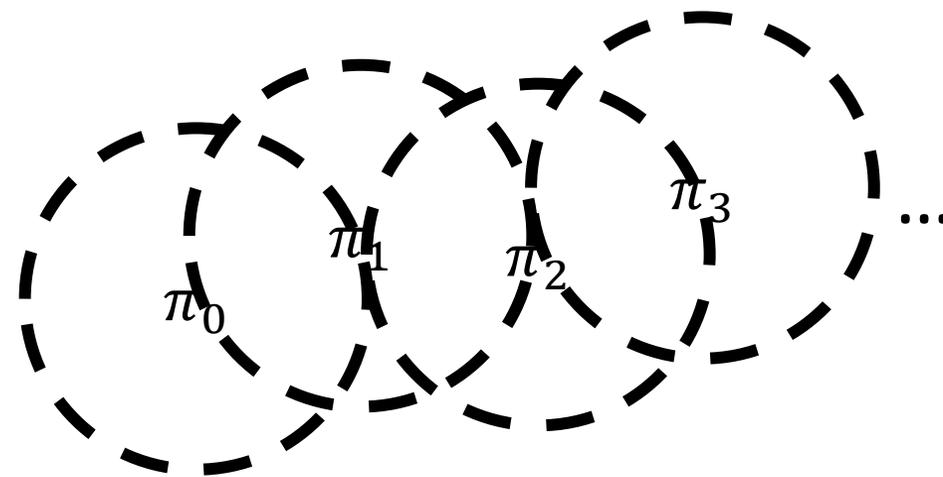
# Trust Region Policy Updates

- To avoid instability in training, we don't want to change our policy too much after any given update. **How can we enforce this?**
- **Idea:** Constrain the “distance” between our current and new policies.
  - The area of the policy space we can update within is called the **trust region**.

**How can we implement this? \***

$$\bar{D}_{KL}(\boldsymbol{\theta}_{k+1} || \boldsymbol{\theta}_k) = E_{s \sim \pi_{\boldsymbol{\theta}_k}} \left[ D_{KL} \left( \pi_{\boldsymbol{\theta}_{k+1}}(\cdot | s) || \pi_{\boldsymbol{\theta}_k}(\cdot | s) \right) \right]$$

During the optimisation step, constrain the KL divergence between the old and new policies across states visited by the old policy.



# Offline RL

- We've seen Off-Policy RL methods (e.g., Q-Learning) use data generated by one policy to learn about another.
- Offline RL methods take this idea to the extreme.
  - The agent is given a **fixed dataset of experience**. This experience could be generated by other agents, human demonstrated, etc.
  - The agent has to learn a policy using **only this fixed dataset** – it can't interact with the environment to generate any more data itself.
- This brings with it many unique challenges!
  - Exploring is not possible!
  - Distributional shifts between dataset experience and online experience.

# Inverse RL

- Regular RL:
  - What **policy** maximises the long-term reward for a given problem?
- Inverse RL:
  - What **reward function** is a given policy maximising?
- Useful when we have demonstrations from an agent completing a task, but we don't know the reward function being maximised.
  - We can use Inverse RL methods to derive the reward function, then train a reinforcement learning agent to solve the task.

# Intrinsically-Motivated RL

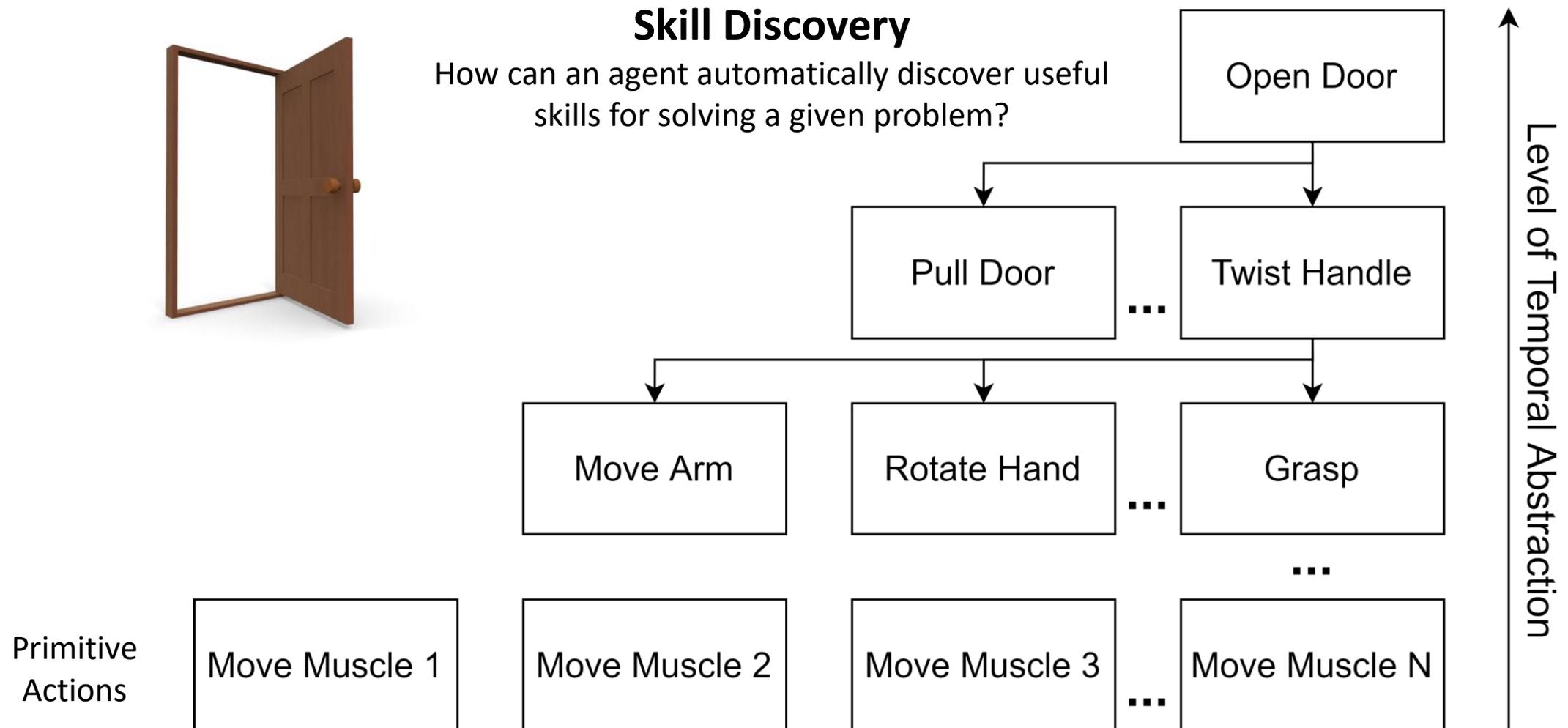
- Many difficult problems have very sparse rewards.
- Can we use **intrinsic** rewards to drive intelligent behaviour, even in the absence of **extrinsic** rewards from the environment?
  - Many examples of this in nature: novelty, curiosity, regularity, information-gain, empowerment, skill diversity...
- A closer look at one example: **Curiosity**
  - ~~Define intrinsic reward based on prediction errors.~~
  - Define intrinsic reward based on **prediction improvements**.

$$\begin{array}{ccccc} R_t & = & R_t^E & + & R_t^I \\ \uparrow & & \uparrow & & \uparrow \\ \text{Total} & & \text{Extrinsic} & & \text{Intrinsic} \\ \text{Reward} & & \text{Reward} & & \text{Reward} \end{array}$$



**Noisy TV Problem:** Our agent will never be learn to predict what's going to be shown next on this noisy TV screen.

# Hierarchical RL

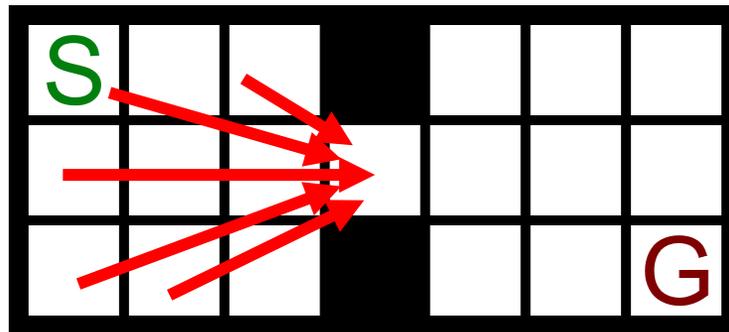


# Hierarchical RL

- How to represent skills? One approach: The **Options** framework.

$$o = \langle I_o, \pi_o, \beta_o \rangle$$

- $I_o$  - **Initiation Set**: Which states can I select this skill in?
- $\pi_o$  - **Option Policy**: What action does this skill select in each state?
- $\beta_o$  - **Termination Condition**: Which states does this skill terminate in?



- How to discover skills?
  - Many different approaches proposed. But still an open research question!

# Summary

