

# Reinforcement Learning

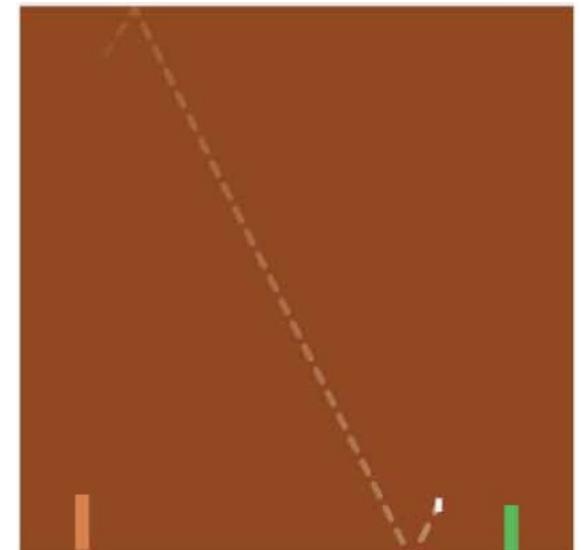
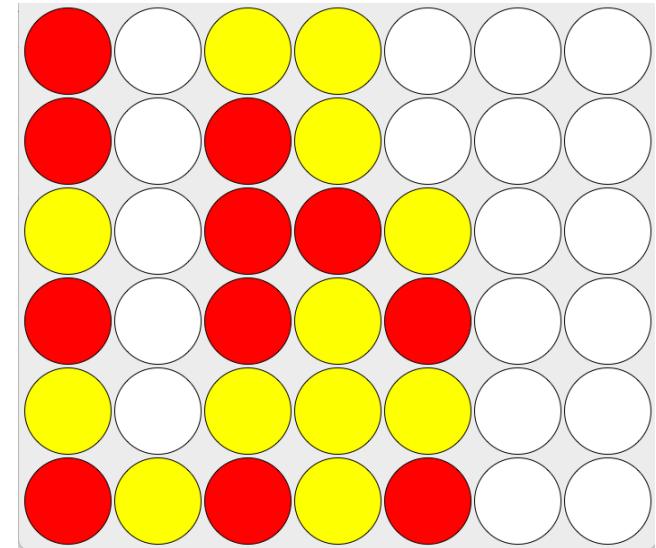
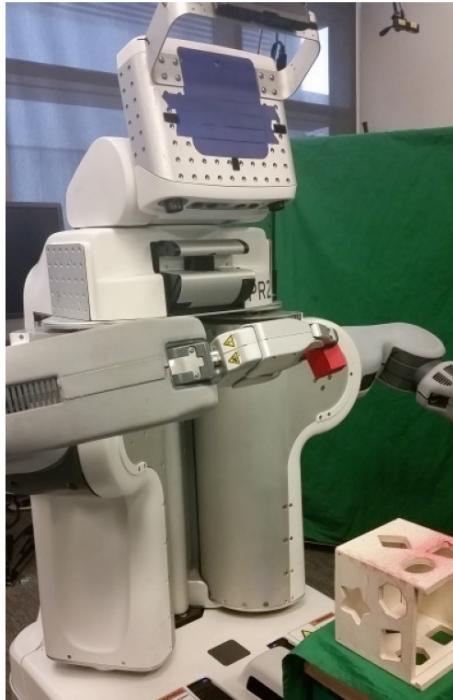
Forest Agostinelli  
University of South Carolina

# Outline

- Background
- Value functions
- Dynamic programming
- Model-free reinforcement learning

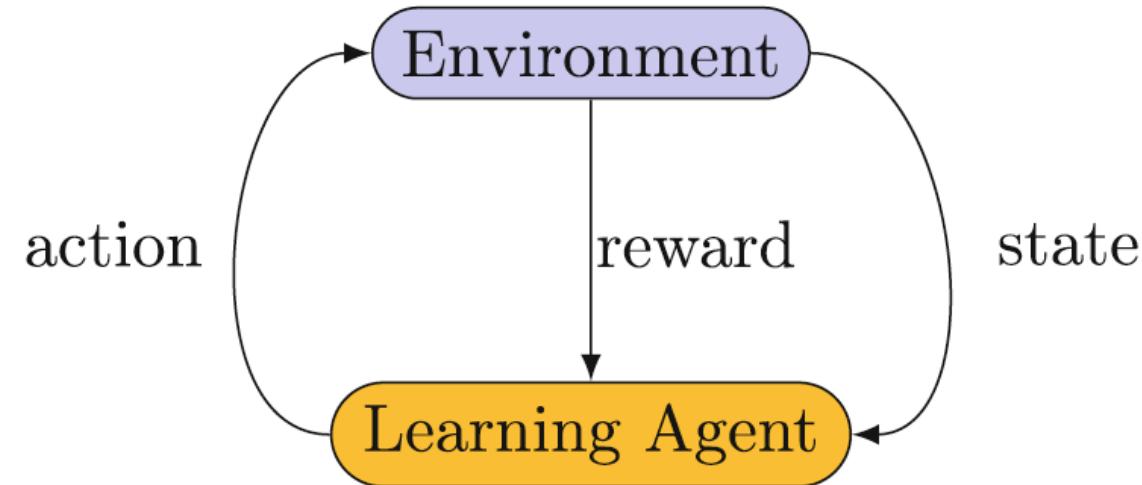
# Reinforcement Learning

- Learning to maximize reward in sequential decision-making problems
- Learning is done through experience
- Decisions affect experience and experience affects decisions
- There may be **stochasticity**
  - Inherent randomness in the outcome of actions
  - The probabilities associated with this randomness may be unknown



# Reinforcement Learning

- **Reinforcement learning:** learning to map **states** to **actions** so that we maximize the expected future **reward** we receive from the **environment**.
- This mapping of states to actions is called a **policy function**.
  - Deterministic:  $a = \pi(s)$
  - Stochastic:  $\pi(a|s) = P(A = a|S = s)$
- At each time step  $t$ 
  - In state  $S_t$ , agent takes action  $A_t$
  - Based on state  $s_t$  and action  $a_t$ , the environment transitions to state  $S_{t+1}$  and outputs reward  $R_{t+1}$



# Markov Decision Processes (MDPs)

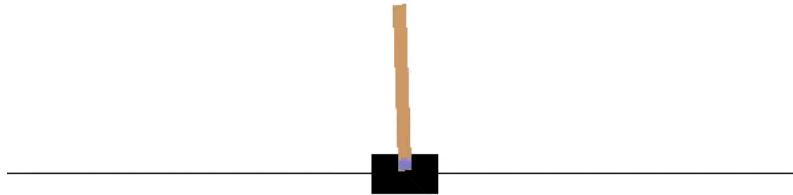
- To formally pose a problem as a reinforcement learning problem, we must define it as a Markov decision process
- **States**
  - The possible configurations of problem
- **Actions**
  - The actions that the agent can take
- **Transition Probabilities**
  - Given state,  $s$ , and action,  $a$ , this defines
    - The probability of transitioning to state,  $s'$
    - The probability of obtaining reward,  $r$

# The Markov Property

- In reinforcement learning, it is often assumed that the **Markov property** holds
- The Markov property states
  - Given the current state and action taken, we do not need any of the **history** to define the transition probabilities
  - The history consists of previous states, actions, and rewards
- In other words
  - The next state and reward is independent of the history given the current state and action

# Markov Property Example

- Cartpole
- Actions: apply force left, apply force right
- Rewards: +1 for every step
- Episode ends when the pole falls over
- The Markov property does **not** hold if we represent the state as the position of the cart and pole
  - Why?
- To ensure the Markov property holds, we must include the velocity of the cart and pole



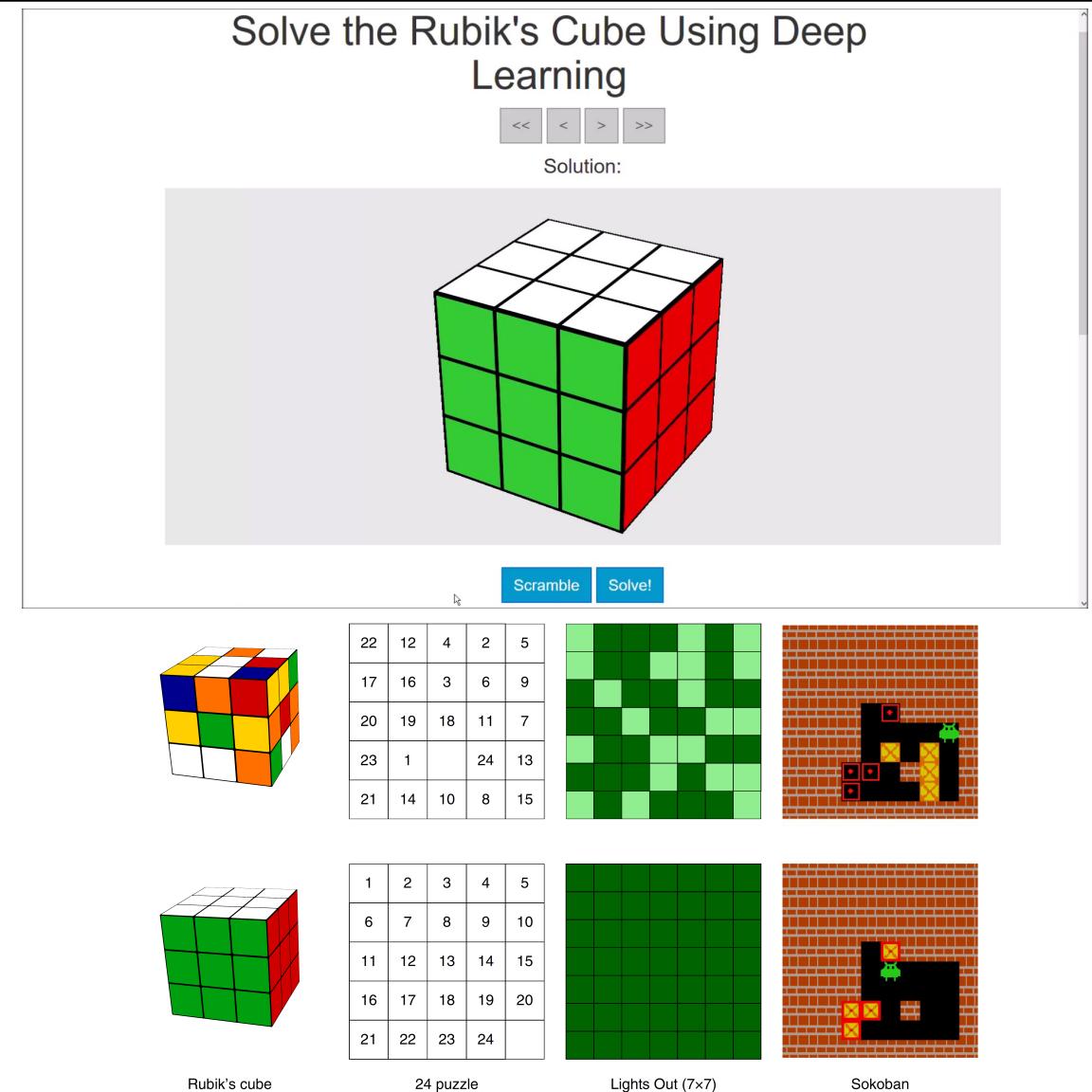
# RL Successes: Atari

- Using RL algorithms, a deep neural network is trained to play Atari games using raw pixels.
- Current work shows we can train deep neural networks to play better than humans on 57 different Atari games.
- RL was combined with **deep neural networks**



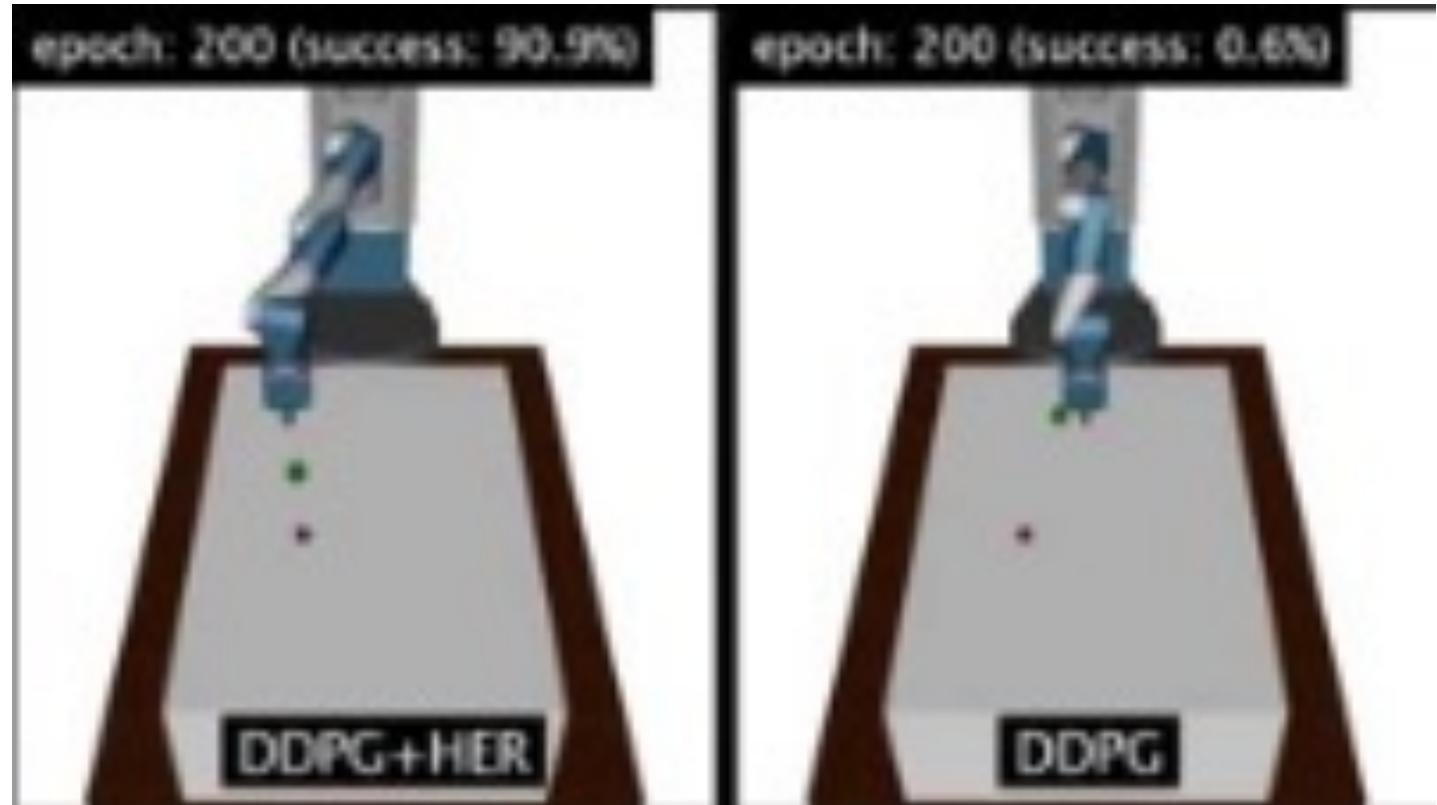
# RL Successes : Solving the Rubik's Cube

- $4.3 \times 10^{19}$  possible combinations
- No domain-specific knowledge
- DeepCubeA solves the Rubik's cube and other puzzles
  - Puzzles have up to  $3.0 \times 10^{62}$  possible combinations.
- Finds a shortest path in the majority of verifiable cases
- <http://deepcube.igb.uci.edu/>
- RL was combined with **deep neural networks** and **A\* search**



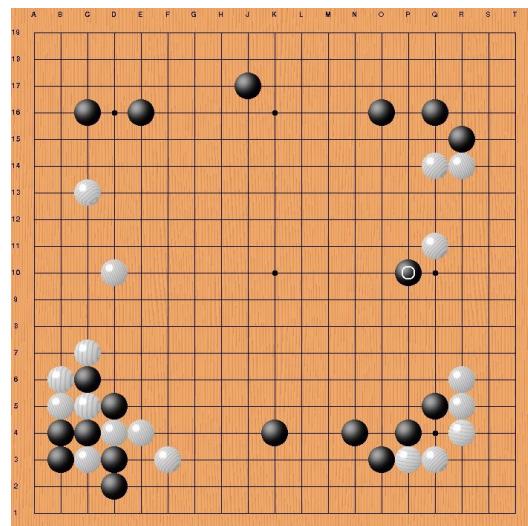
# RL Successes: Robotics

- Can solve problems in continuous environments
- Can train in simulation and transfer to the real world (Sim2Real)
- RL was combined with **deep neural networks**



# RL Successes: Go

- AlphaGo learned how to play Go from expert demonstrations data and from self-play
  - Defeated one of the best Go players, Lee Sedol, 4 to 1
  - Move 37
- AlphaGoZero builds on AlphaGo and learns only from self-play
- RL was combined with **deep neural networks** and **Monte Carlo tree search**

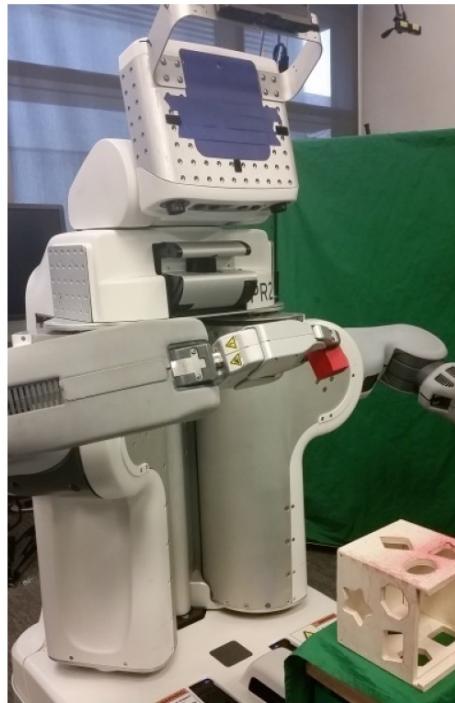
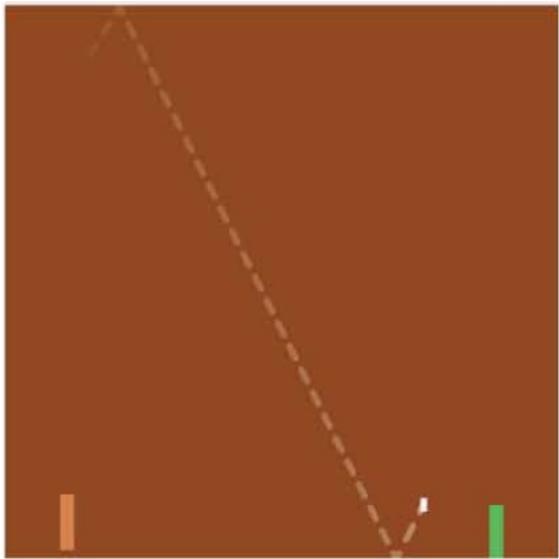


# Imitation Learning

- If we know how to solve the problem efficiently (or at all!), we can use **imitation learning** to train an agent to imitate the actions we want it to take
- Collect a dataset of states and the actions humans took to solve the problem
  - Input: states
  - Output: actions
  - Objective function: accurately predict actions given states
- We can now pose the problem as a **supervised learning** problem
  - We already have a plethora of tools to solve these problems

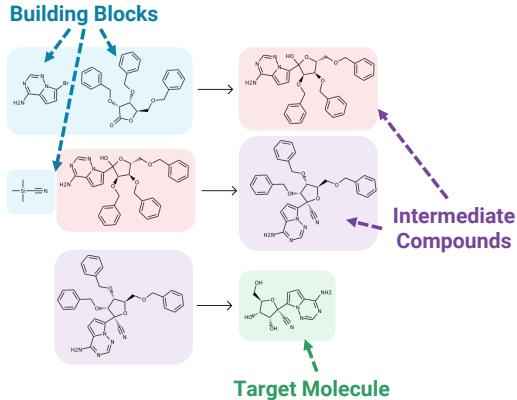
# Reinforcement Learning: Motivation

- Collecting data from humans may be very time-consuming
- Distribution shift may be an issue



# Reinforcement Learning: Motivation

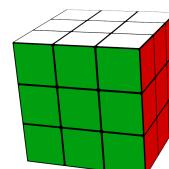
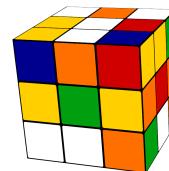
- Humans may not know how to solve the problem efficiently, or at all!



**Input:**  $n \times n$  skew-symmetric matrix  $\mathbf{A}$ , vector  $\mathbf{b}$ .

**Output:** The resulting vector  $\mathbf{c} = \mathbf{Ab}$  computed in  $\frac{(n-1)(n+2)}{2}$  multiplications.

```
(1) for  $i = 1, \dots, n-2$  do
(2)   for  $j = i+1, \dots, n$  do
(3)      $w_{ij} = a_{ij}(b_j - b_i)$                                 ▷ Computing the first  $(n-2)(n+1)/2$  intermediate products
(4)   for  $i = 1, \dots, n$  do
(5)      $q_i = b_i \sum_{j=1}^n a_{ji}$                                 ▷ Computing the final  $n$  intermediate products
(6)   for  $i = 1, \dots, n-2$  do
(7)      $c_i = \sum_{j=1}^{i-1} w_{ji} + \sum_{j=i+1}^n w_{ij} - q_i$ 
(8)    $c_{n-1} = -\sum_{i=1}^{n-2} \sum_{j=i+1}^{n-2} w_{ij} - \sum_{j=1}^{n-2} w_{jn} + \sum_{i=1, i \neq n-1}^n q_i$ 
(9)    $c_n = -\sum_{i=1}^{n-1} \sum_{j=i+1}^n w_{ij} + \sum_{i=1}^{n-1} q_i$ 
```

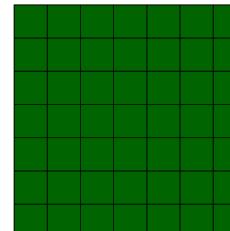
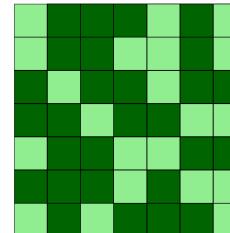
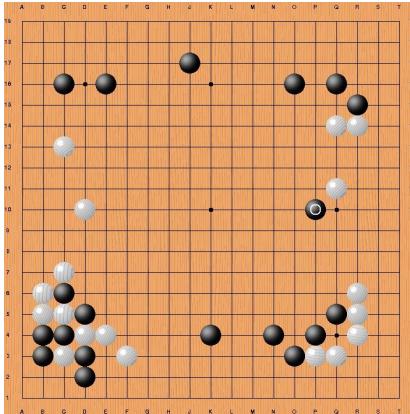


Rubik's cube

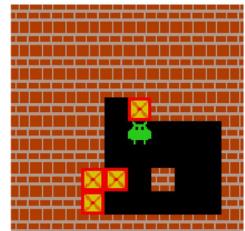
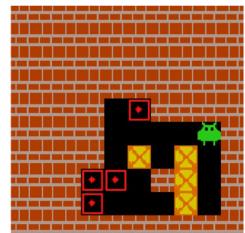
22	12	4	2	5
17	16	3	6	9
20	19	18	11	7
23	1		24	13
21	14	10	8	15

1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	

24 puzzle



Lights Out (7x7)



Sokoban

# Outline

- Background
- Value functions
- Dynamic programming
- Model-free reinforcement learning

# Value Functions

- **Episode:** Starts at some start state at timepoint 0 and ends at a special state, called the terminal state, at timepoint  $T$
- **Return:** the sum of rewards after timestep  $t$ 
  - $G_t = R_{t+1} + R_{t+2} + R_{t+3} \dots + R_T$
  - We seek to maximize the expected return
- **State-value function**
  - $v_\pi(s)$  is the expected return when in state  $s$  and following policy  $\pi$
- **Action-value function**
  - $q_\pi(s, a)$  is the expected return when taking action  $a$  in state  $s$  and then following policy  $\pi$
- Value functions are specific to a given policy  $\pi$

# Optimal Policy and Value Function

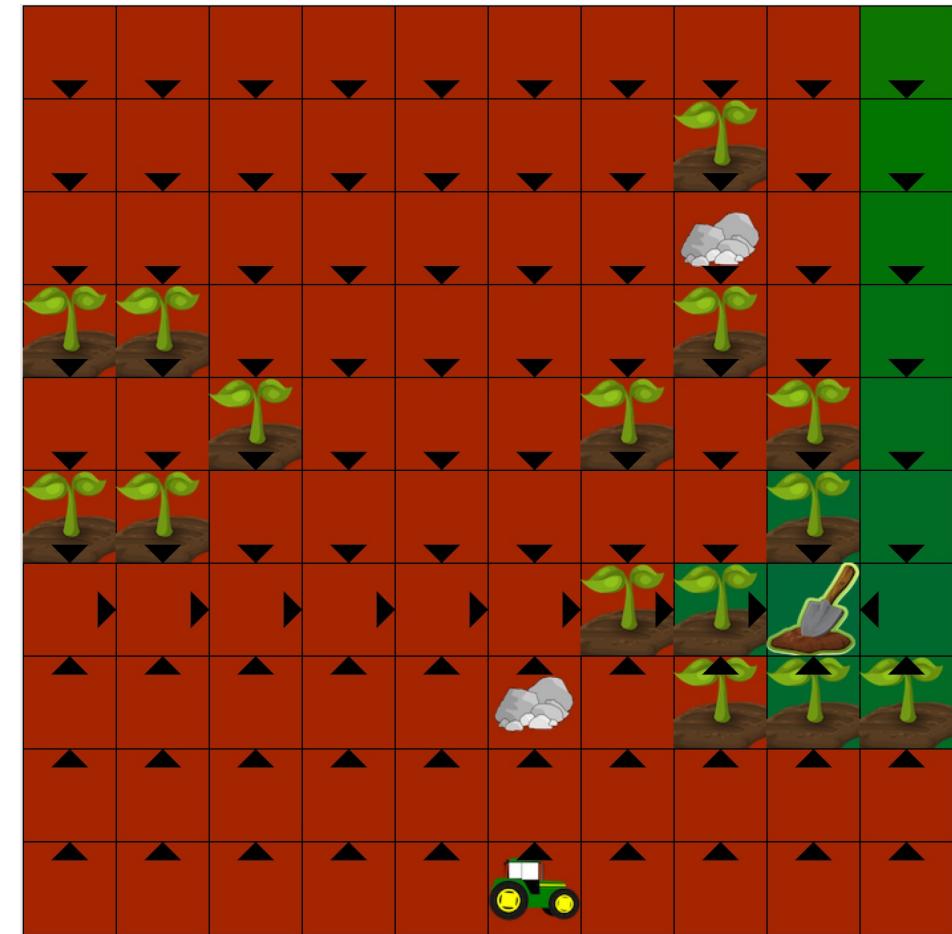
- Which policy is better?
  - $\pi \geq \pi'$  if and only if  $v_\pi(s) \geq v_{\pi'}(s)$  for all  $s \in S$
- A policy that achieves the greatest possible return from any state is an optimal policy
  - $\pi_* \geq \pi'$  for all  $\pi'$
- The optimal value function is the value function obtained when following the optimal policy
  - $v_*(s) = \max_\pi v_\pi(s)$
  - $q_*(s, a) = \max_\pi q_\pi(s, a)$
- Optimal policies can be obtained by behaving greedily with respect to the optimal value function
- Many RL methods first learn a value function and then **induce** a policy by behaving greedily with respect to the value function
- Is the optimal policy unique?
- Is the optimal value function unique?

# Value Functions



$$v_{\pi}(s)$$

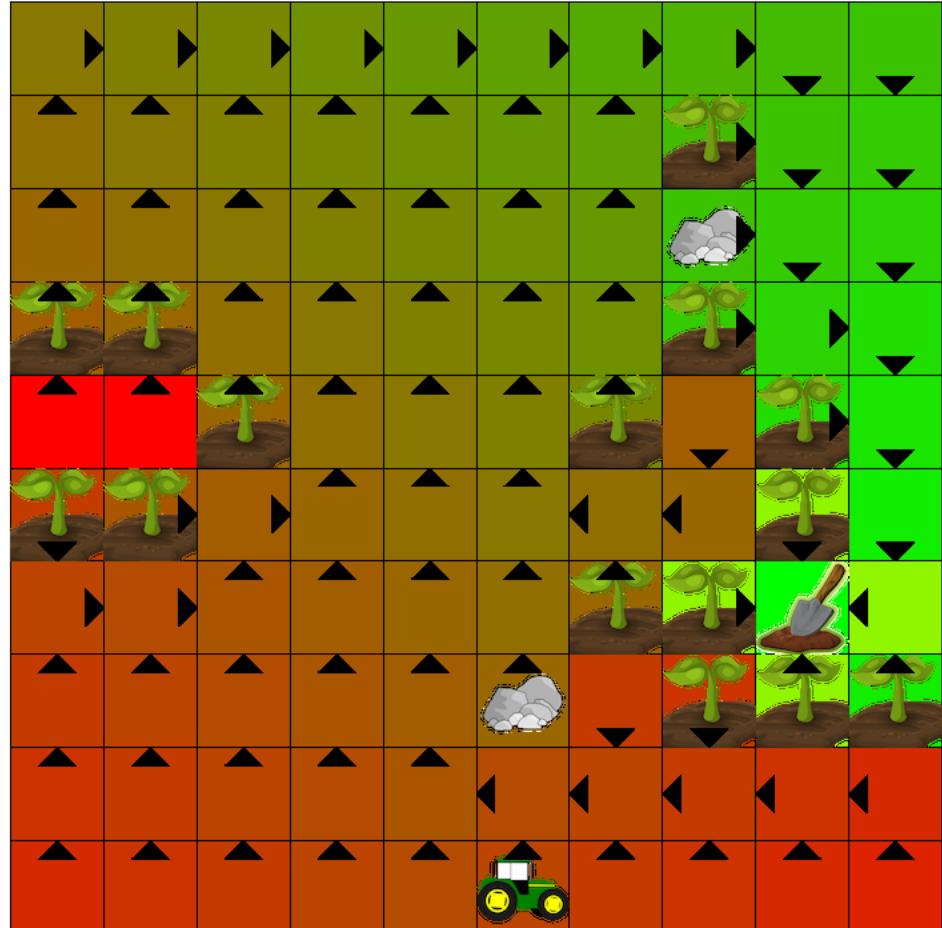
$\pi$  is uniform random for all states



$$v_{\pi}(s)$$

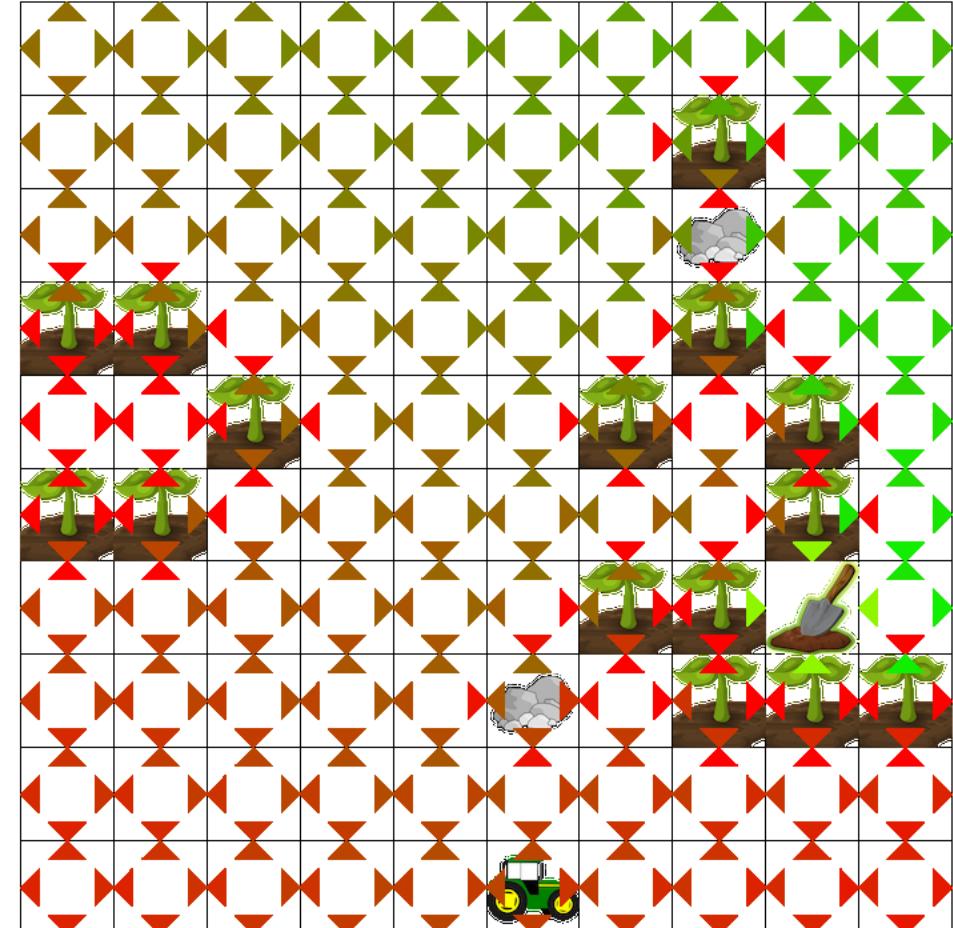
$\pi$  says to go up if below goal, go down if above goal, otherwise, go in the direction of the goal

# Optimal Value Function



$$v_*(s)$$

$$\pi_* = \underset{a}{\operatorname{argmax}}(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$$

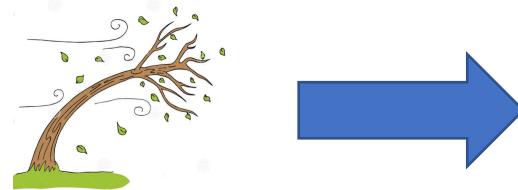
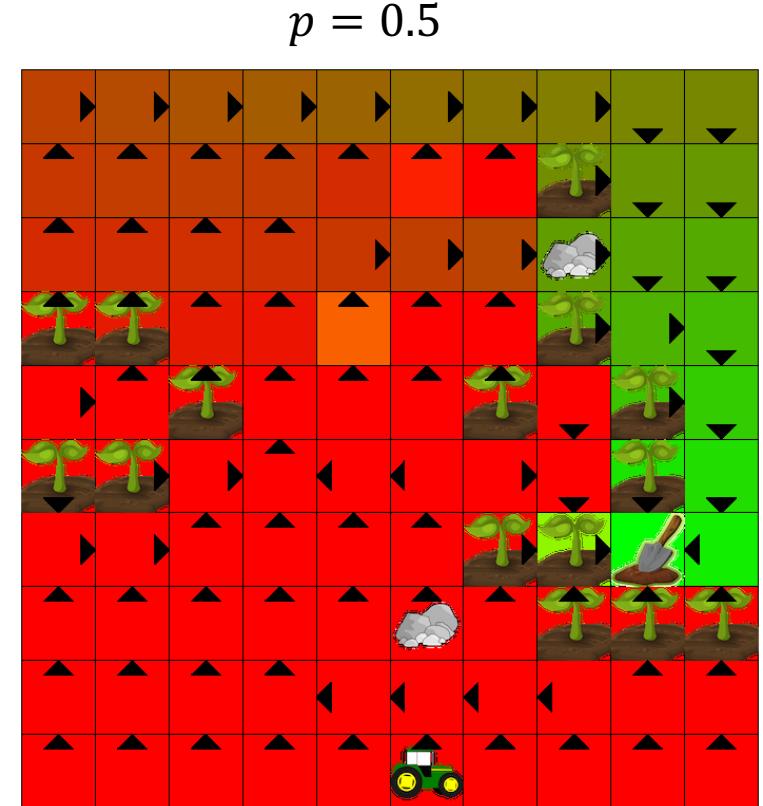
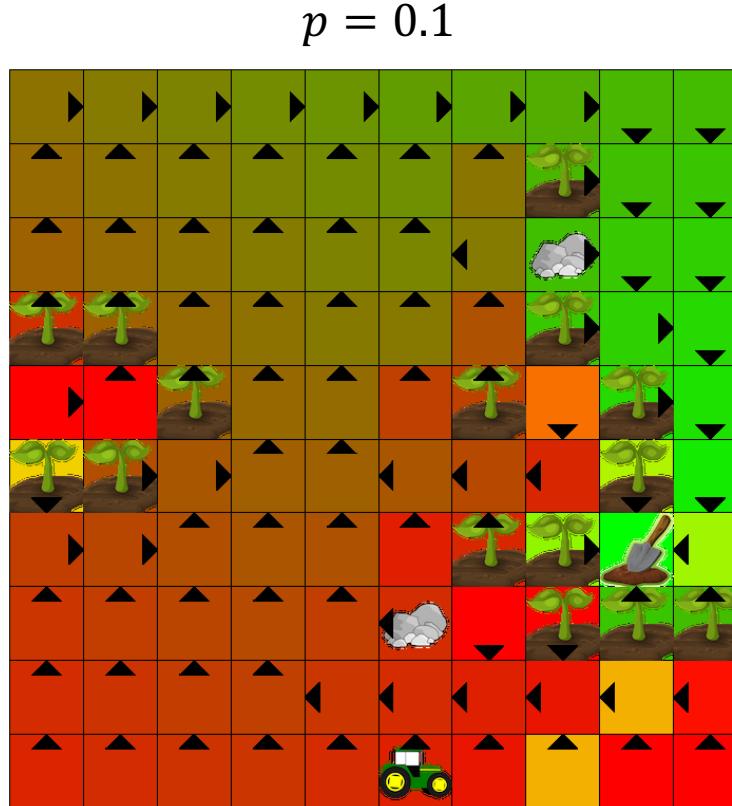
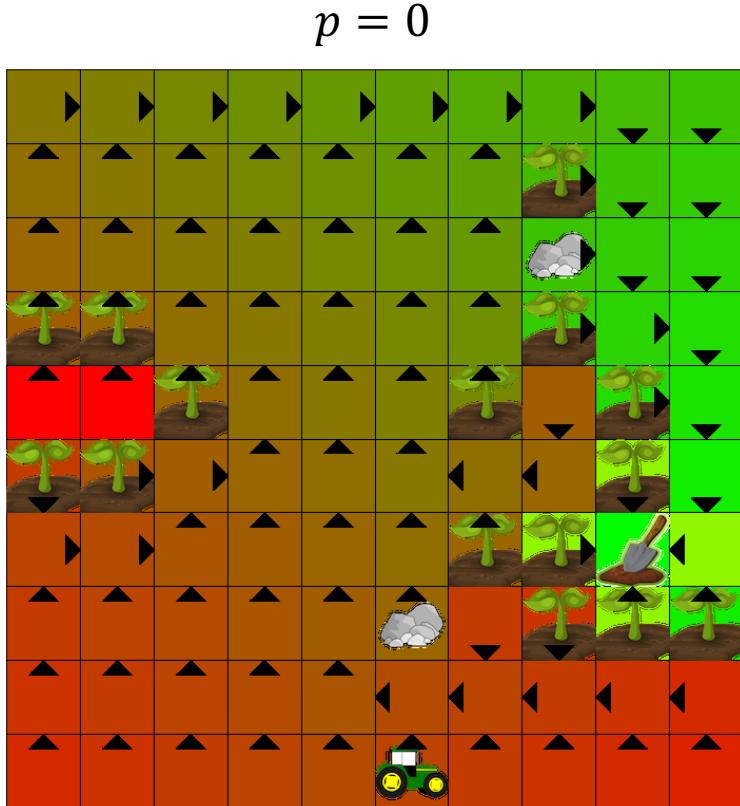


$$q_*(s, a)$$

$$\pi_* = \underset{a}{\operatorname{argmax}} q_*(s, a)$$

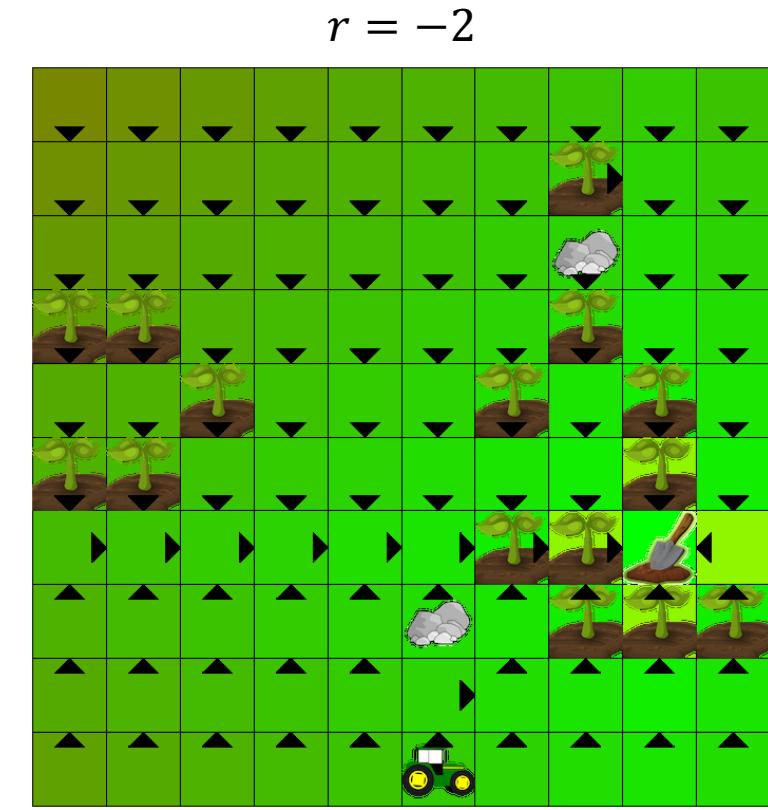
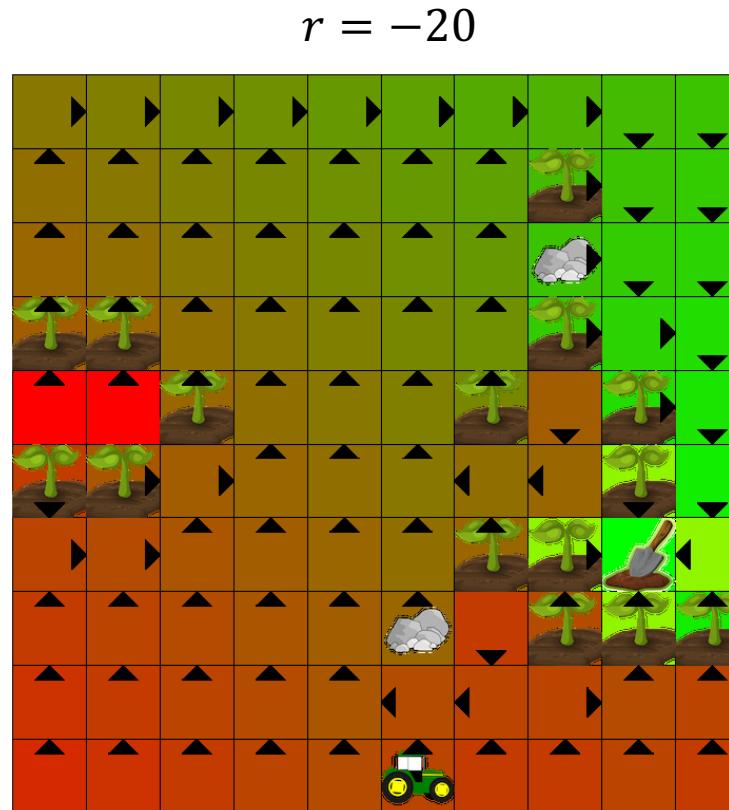
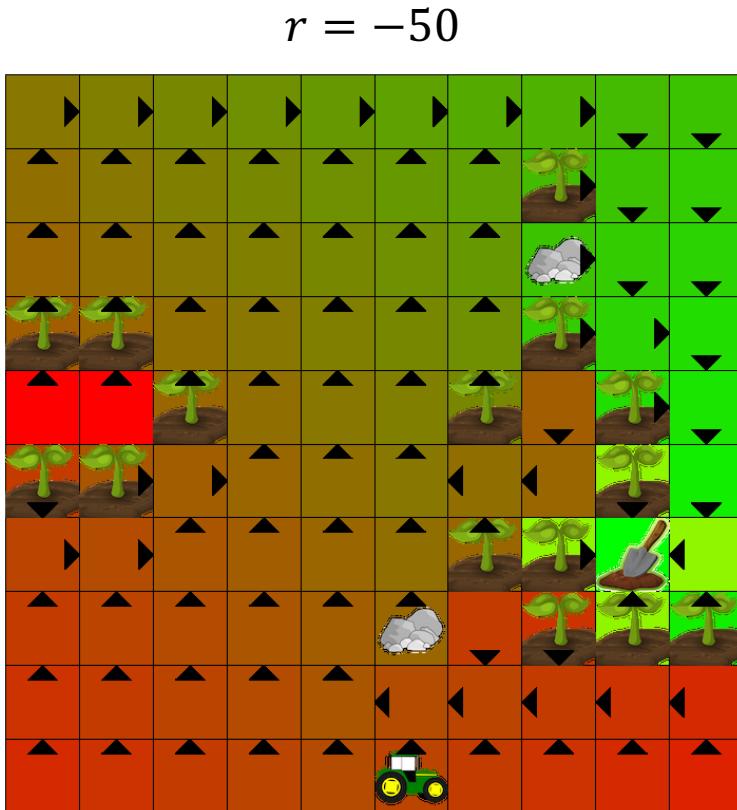
# Optimal Value Function: Strong Winds

- Wind blows you right with some probability  $p$



# Optimal Value Function: Effect of Rewards

- Reward  $r$  of driving on a plant

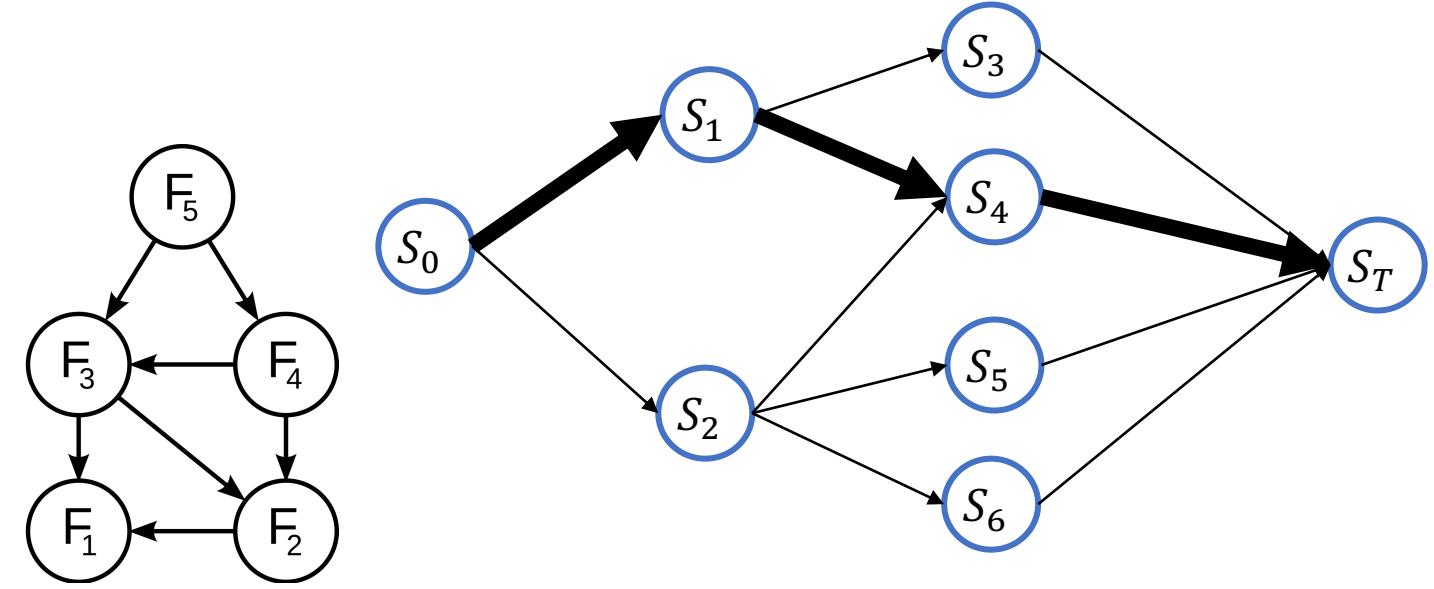
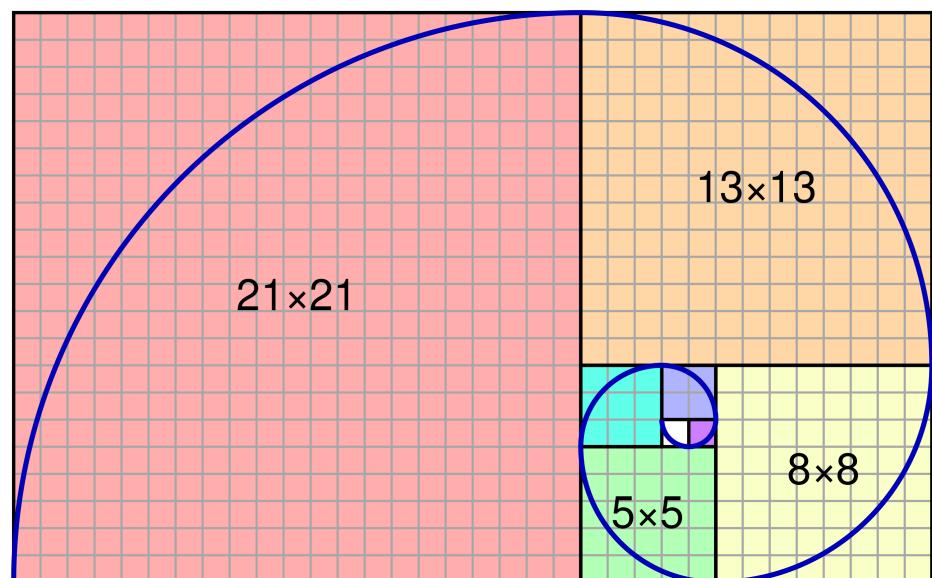


# Outline

- Background
- Value functions
- Dynamic programming
- Model-free reinforcement learning

# Dynamic Programming

- Solves problems by recursively breaking them down into simpler subproblems
- Requires
  - **Optimal substructure**: Can construct an optimal solution from optimal solutions of subproblems
    - Principle of optimality
  - **Overlapping subproblems**: Solutions to subproblems are re-used
    - Value functions
- Other applications: Fibonacci sequence, Scheduling, Sequence alignment (DNA)

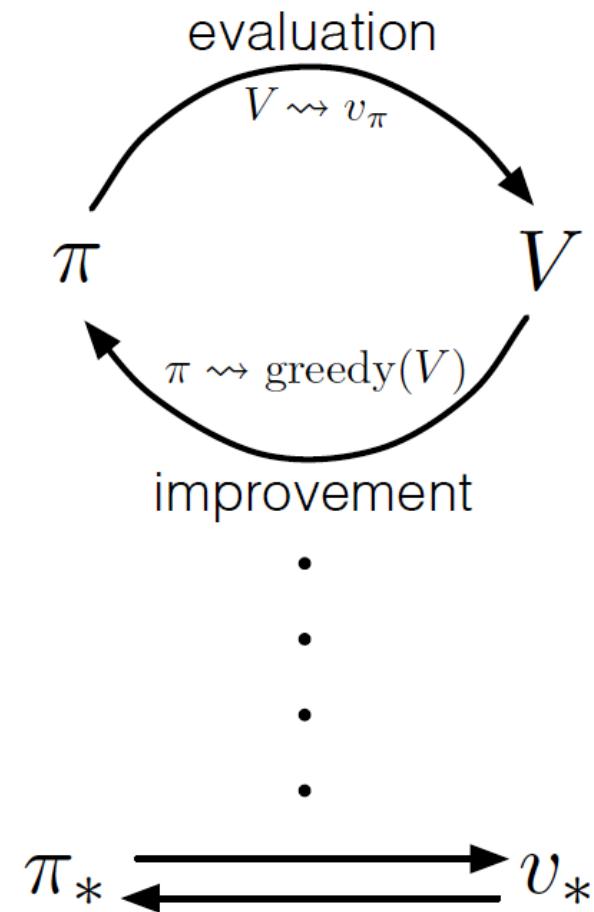
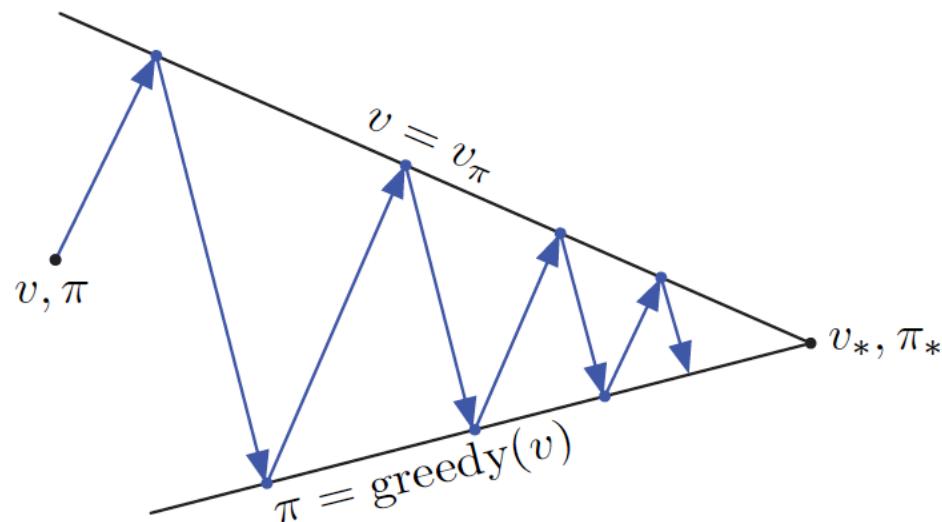


# Dynamic Programming

- We will use it to evaluate a policy and compute an optimal policy given a perfect model an MDP
- Foundational for reinforcement learning
- Using dynamic programming, we will do **policy iteration** by iterating between **policy evaluation** and **policy improvement**

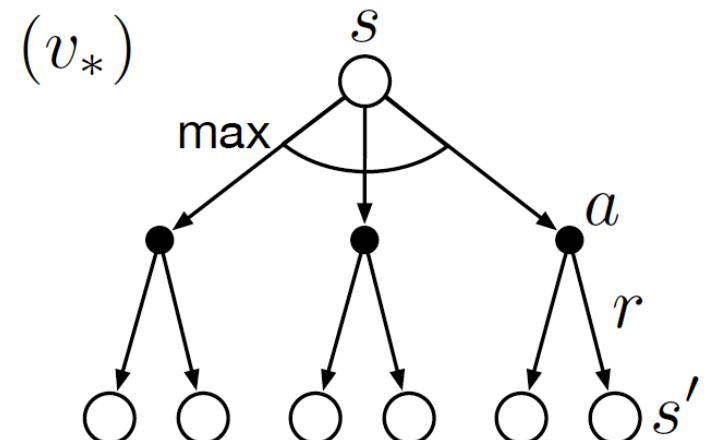
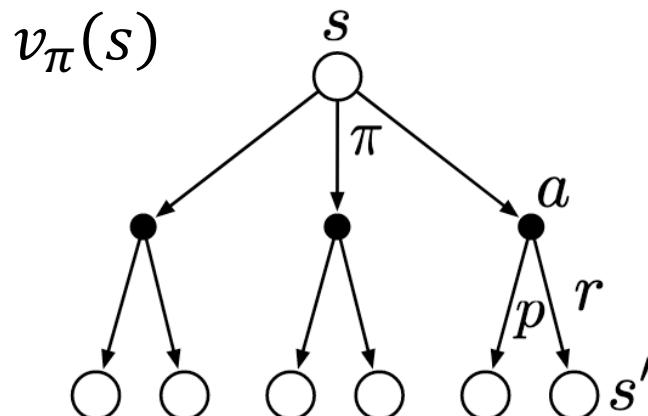
# Generalized Policy Iteration

- **Policy Evaluation:** Estimate the expected future reward when following policy  $\pi$
- **Policy Improvement:** Improve policy  $\pi$  so that it obtains a greater expected future reward
- We can obtain an optimal policy by iterating between **policy evaluation** and **policy improvement**



# Bellman Equation and Bellman Optimality Equation

- To determine the value of a state,  $s$ , we need only perform a one-step lookahead
  - Look at every possible action we can take
  - Look at the expected reward for taking those actions
  - Look at the probability of transitioning to a next state,  $s'$
  - Look at value of next state,  $s'$
- Bellman equation
  - Summarize based on probability of taking those actions
  - $v_\pi(s) = \sum_a \pi(a|s)(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s'))$
- Bellman optimality equation
  - Summarize based on the max of all possible actions
  - $v_*(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$

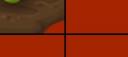


# Policy Evaluation

- Construct  $v_\pi$ , the expected future reward when following policy  $\pi$
- Use the Bellman equation as an update rule
- Theory shows that, in the tabular case, this is guaranteed to converge to  $v_\pi$

# Policy Evaluation AI Farm

- Policy: Uniform random

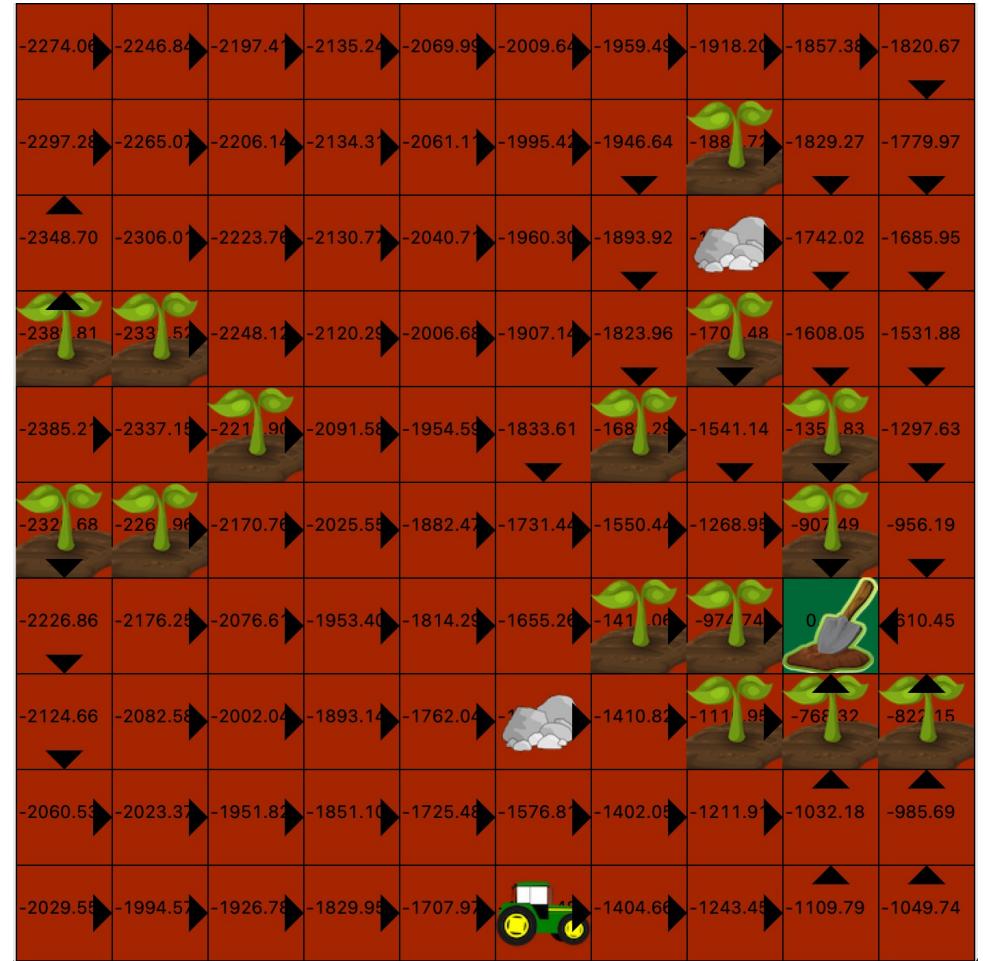
-2274.06	-2246.84	-2197.41	-2135.24	-2069.99	-2009.64	-1959.49	-1918.20	-1857.38	-1820.67	
-2297.28	-2265.07	-2206.14	-2134.31	-2061.11	-1995.42	-1946.64	-1881.72	-1829.27	-1779.97	
-2348.70	-2306.01	-2223.76	-2130.77	-2040.71	-1960.30	-1893.92	-1742.02	-1685.95		
-2381.81	-2331.52	-2248.12	-2120.29	-2006.68	-1907.14	-1823.96	-1704.48	-1608.05	-1531.88	
-2385.21	-2337.15	-2211.90	-2091.58	-1954.59	-1833.61	-1681.29	-1541.14	-1351.83	-1297.63	
-2321.68	-2261.96	-2170.76	-2025.55	-1882.47	-1731.44	-1550.44	-1268.95	-907.49	-956.19	
-2226.86	-2176.25	-2076.61	-1953.40	-1814.29	-1655.26	-1411.06	-974.74	0	-610.45	
-2124.66	-2082.58	-2002.04	-1893.14	-1762.04	-1410.82	-1111.95	-768.32	-822.15		
-2060.53	-2023.37	-1951.82	-1851.10	-1725.48	-1576.81	-1402.05	-1211.91	-1032.18	-985.69	
-2029.55	-1994.57	-1926.78	-1829.95	-1707.97	-1404.66	-1243.45	-1109.79	-1049.74		

# Policy Improvement

- We have evaluated policy  $\pi$ , how can we find a better policy?
- $\pi' \geq \pi$  if and only if  $v_{\pi'}(s) \geq v_\pi(s)$  for all  $s \in \mathcal{S}$
- We set the policy to be greedy with respect to  $v_\pi$ 
  - Perform a one-step lookahead and choose the action with the highest expected future reward
  - $\pi'(s) = \operatorname{argmax}_a(r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s'))$
- This policy will always be the same or better than the previous one
  - Policy improvement theorem.

# Policy Improvement

- Better than original uniform random policy
- Still not optimal



When acting greedily with respect to  $\nu_\pi$

# Policy Iteration

- Policy improvement improves a policy  $\pi$  and obtains a new policy  $\pi'$  such that,  $\pi' \geq \pi$
- If  $\pi' = \pi$ , then  $v_{\pi'} = v_\pi$
- Therefore,  $v_{\pi'}(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_{\pi'}(s'))$
- This is the same as the Bellman optimality equation!
- Does this mean  $\pi' = \pi_*$  and  $v_{\pi'} = v_*$ ?
- There is a proof showing that the Bellman optimality equation is a unique fixed point
  - Similar to that of the Bellman equation

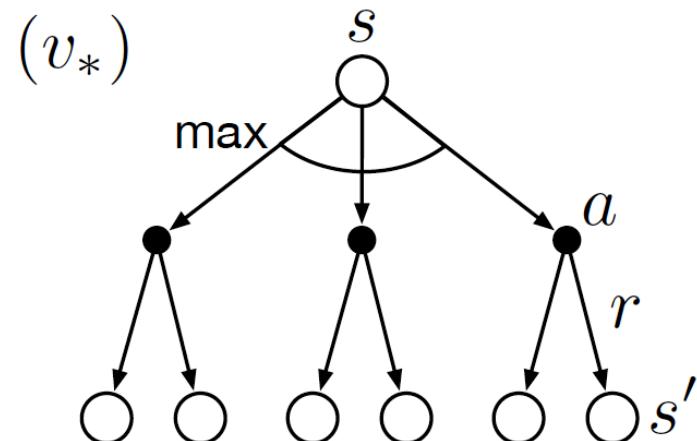
# Policy Iteration: AI Farm



```
(drl) forestagostinelli@Forests-MacBook-Pro Farm_Grid_World % python run_policy_iteration.py --map maps/map1.txt --wait 1.0 --wait_evaluation 0.0 --rand_right 0.0
```

# Value Iteration

- Find the optimal value function
- Recall the Bellman optimality equation
  - $v_*(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_*(s'))$
- Use this as an update rule
  - $V(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$
  - Guaranteed to converge to  $v_*$
- Combines policy evaluation and policy improvement into one step



# Value Iteration: AI Farm

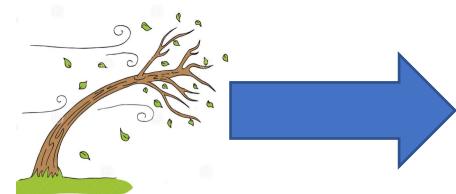
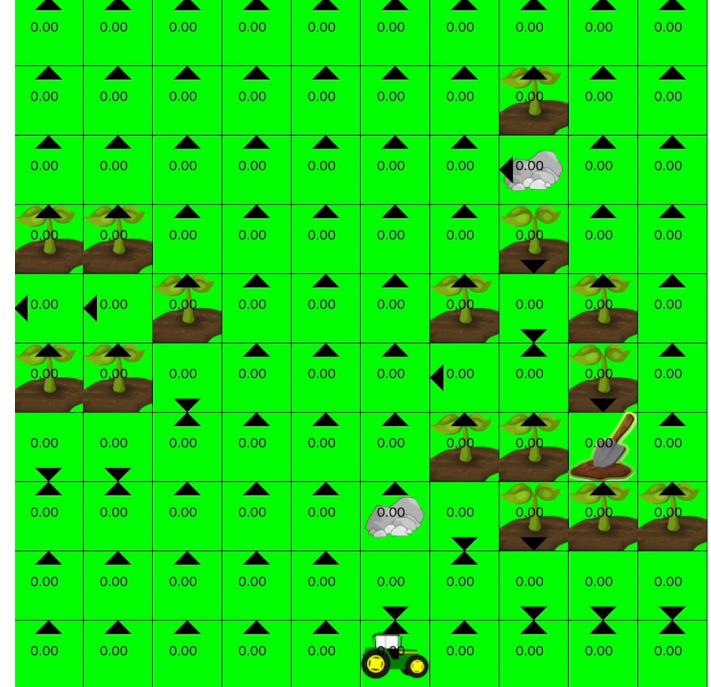
$p = 0$



$p = 0.1$

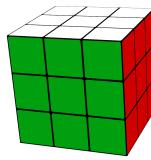
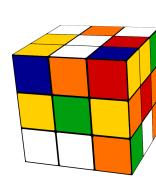


$p = 0.5$



# Approximate Value Iteration

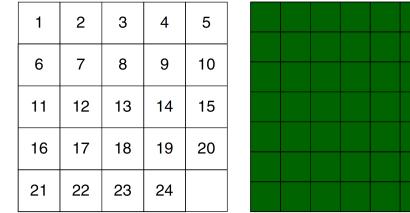
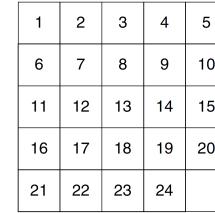
- Sometimes, we can accurately define the MDP, however, the state space is too large to store in a table
- Therefore, we will need to use an architecture to approximate the value
- The goal is to have an architecture whose number of parameters is far less than the size of the state space



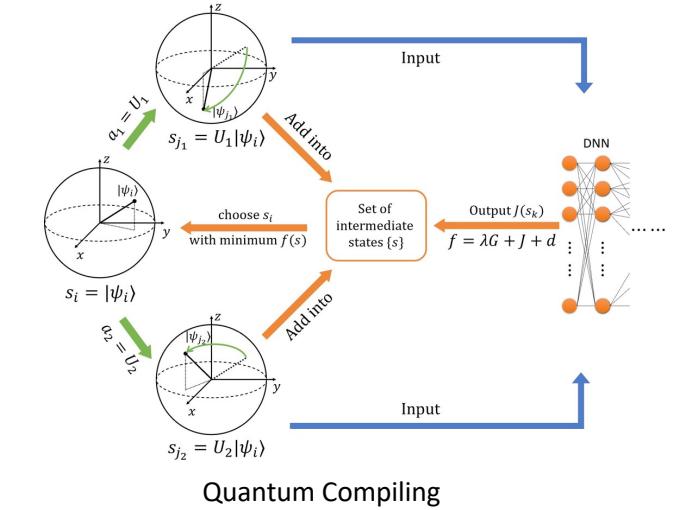
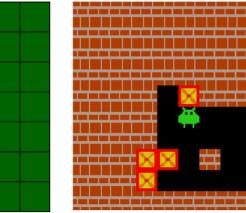
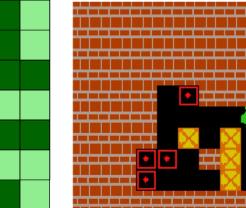
Rubik's cube

22	12	4	2	5
17	16	3	6	9
20	19	18	11	7
23	1		24	13
21	14	10	8	15

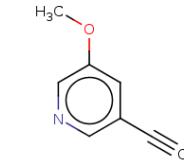
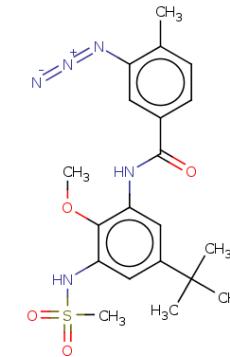
1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	



Puzzles



Quantum Compiling



Chemical Synthesis

# Approximate Value Iteration (AVI)

- Value Iteration
  - $V(s) = \max_a(r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$
- Approximate Value Iteration
  - $y = \max_a(r(s, a) + \gamma \sum_{s'} p(s'|s, a) \hat{v}(s', \mathbf{w}))$
  - $E(\mathbf{w}) = \frac{1}{2} (y - \hat{v}(s, \mathbf{w}))^2$
  - $\nabla_{\mathbf{w}} E(\mathbf{w}) = (y - \hat{v}(s, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(s, \mathbf{w})$
  - Even though  $y$  depends on  $\mathbf{w}$ , we do not differentiate  $y$  with respect to  $\mathbf{w}$ .
  - Only supervision is that  $y = 0$  for terminal states
- Unlike in the tabular case, we cannot guarantee convergence in the case of neural network function approximators
- However, there are many different methods that we can use to make deep reinforcement learning work in practice

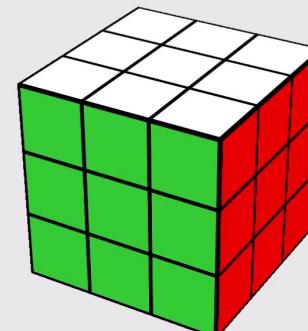
# AVI: Solving the Rubik's Cube

- $4.3 \times 10^{19}$  possible combinations
- No domain-specific knowledge
- DeepCubeA solves the Rubik's cube and other puzzles
  - Puzzles have up to  $3.0 \times 10^{62}$  possible combinations.
- Finds a shortest path in the majority of verifiable cases
- <http://deepcube.igb.uci.edu/>
- RL was combined with **deep neural networks** and **A\* search**

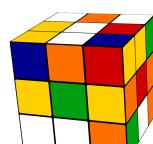
Solve the Rubik's Cube Using Deep Learning

<< < > >>

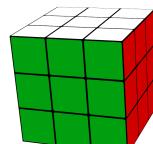
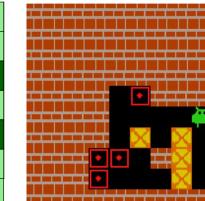
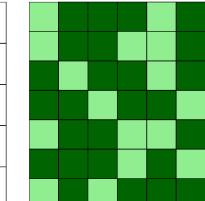
Solution:



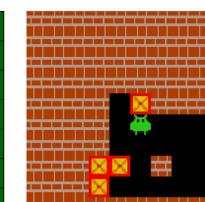
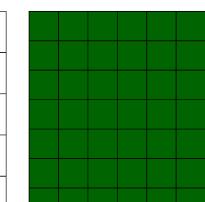
Scramble Solve!



22	12	4	2	5
17	16	3	6	9
20	19	18	11	7
23	1		24	13
21	14	10	8	15



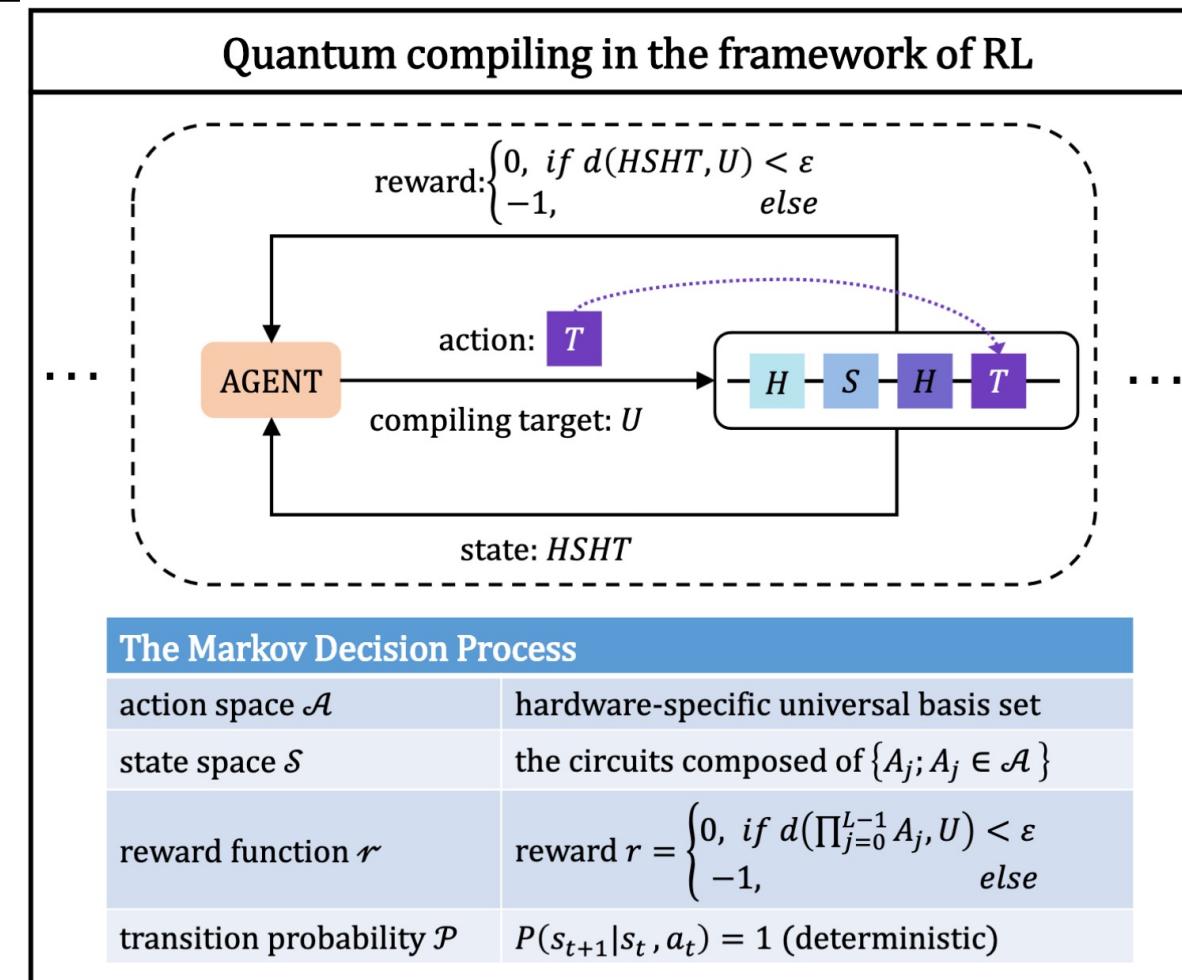
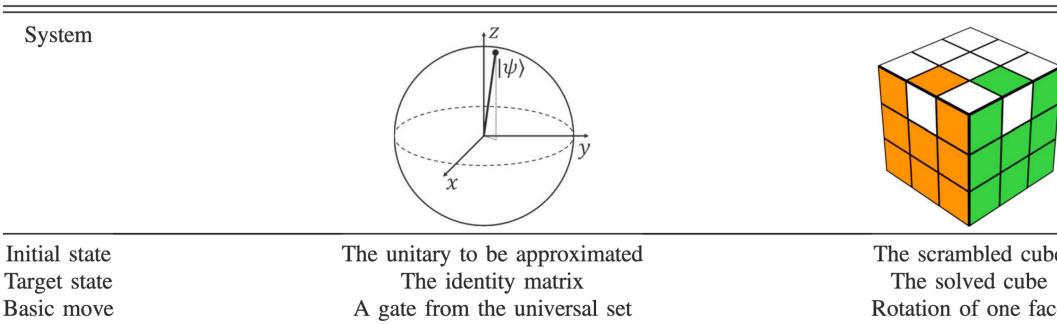
1	2	3	4	5
6	7	8	9	10
11	12	13	14	15
16	17	18	19	20
21	22	23	24	



Rubik's cube      24 puzzle      Lights Out (7x7)      Sokoban

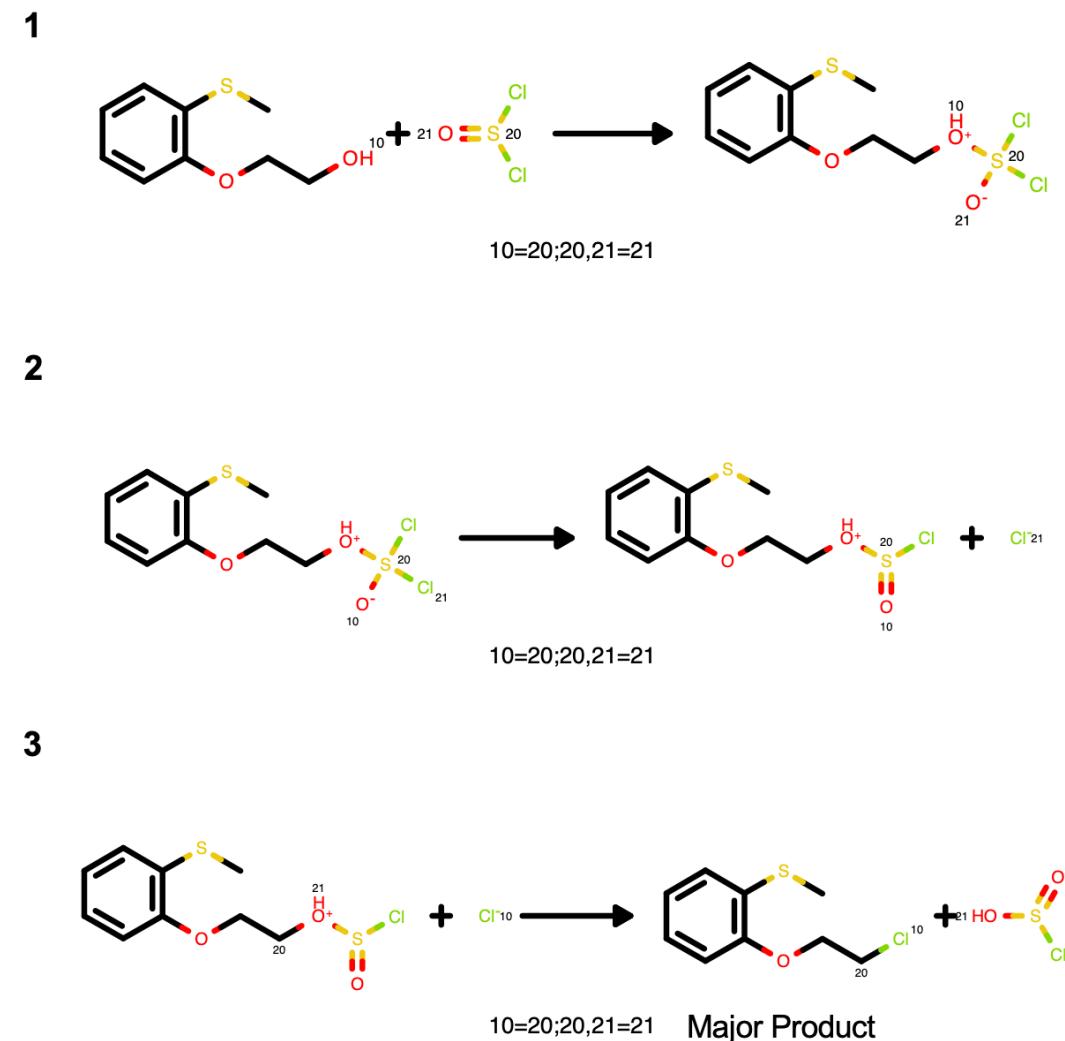
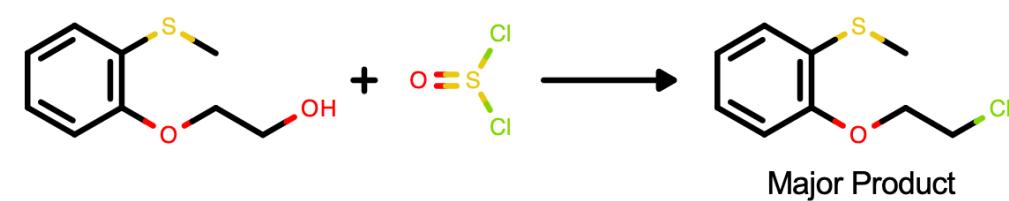
# AVI: Quantum Algorithm Compilation

- Given a quantum algorithm, a compiler must synthesize a quantum circuit for this algorithm from a given set of quantum gates
- If a given circuit is below an error threshold, then the problem is considered solved



# AVI: Reaction Mechanisms

- Chemical reactions are composed of smaller steps called **reaction mechanisms**
- Knowledge of the reaction mechanisms that compose a chemical reaction allows practitioners to
  - Validate reaction feasibility
  - Improve reaction efficiency
  - Predict reaction outcome under different conditions
- Most chemical reaction prediction methods skip reaction mechanisms and predict products directly from reactants



# Outline

- Background
- Value functions
- Dynamic programming
- Model-free reinforcement learning

# Crucial Assumption

- Assuming environment dynamics are known
  - $P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$
- Environment dynamics are unknown in many real-world scenarios
- Even if known, may be too costly to compute (i.e. physics)

$$V(s) = \max_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$$

The diagram consists of a mathematical equation for the value function  $V(s)$ . It features two black arrows pointing upwards from the word "Unknown" at the bottom to the red-colored terms  $r(s, a)$  and  $p(s'|s, a)$  in the equation. The first arrow is positioned under the term  $r(s, a)$ , and the second arrow is positioned under the term  $p(s'|s, a)$ .

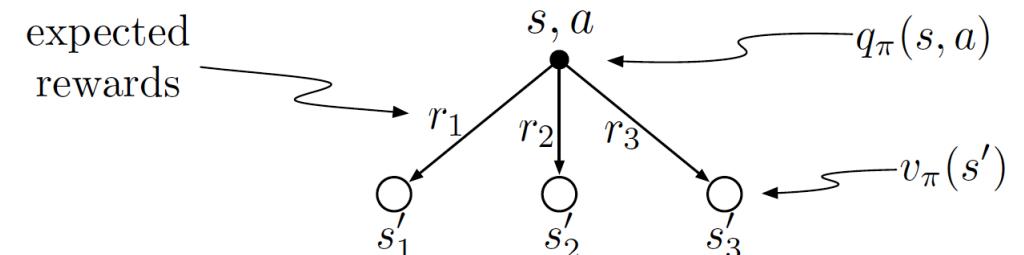
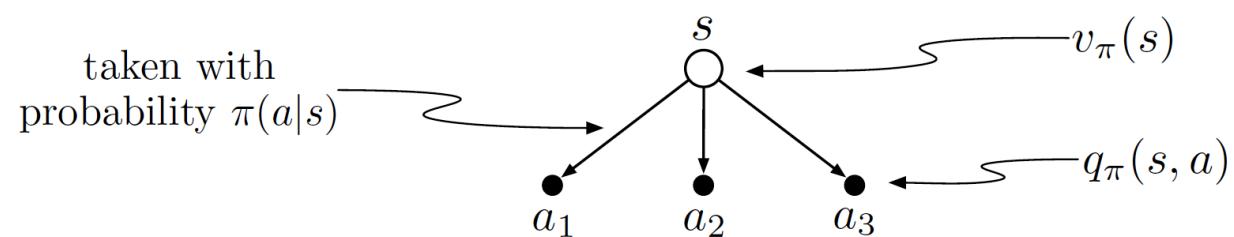
- Instead of using a model, we must learn from **experience**

# Model-Free RL Concepts

- Exploration vs Exploitation
  - **Exploration:** Learn more about the environment
  - **Exploitation:** Use what you have learned to obtain more reward
- On-policy vs Off-policy
  - Your policy determines your experience
  - We need to explore using randomness
    - May not be the best policy
  - Experience may be delicate and hard to obtain (i.e. a hospital)
  - **Behavior policy:** policy that we use to interact with the environment
  - **Target policy:** policy that we wish to evaluate and/or improve

# Model-Free Control

- In this dynamic programming, we induced a policy by doing a one step lookahead using the value function
  - $\pi(s) = \operatorname{argmax}_a (r(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s'))$
- However, we cannot do this in the model-free case because we do not have access to a model
- Therefore, we use an action-value function to induce a policy
  - $q_\pi(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) v_\pi(s')$
  - $\pi(s) = \operatorname{argmax}_a (Q(s, a))$

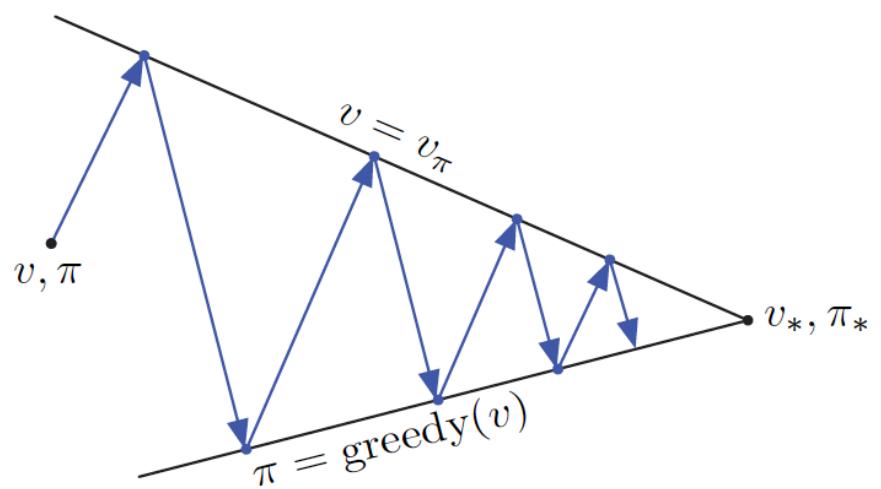


# Model-Free Control: Exploration

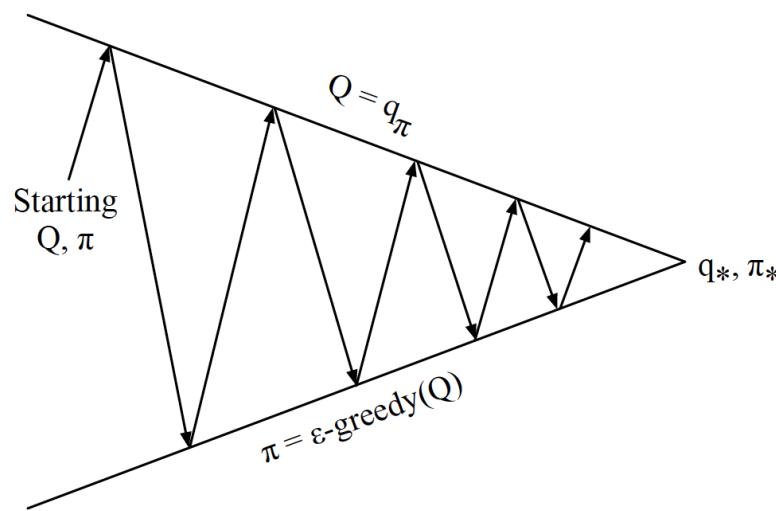
- How do we ensure that we explore our state space?
  - In dynamic programming, we assumed that we could just loop over every possible state
  - Cannot do this in the model-free case
- $\epsilon$ -greedy policy
  - Take a random action with probability  $\epsilon$
  - Take the greedy action,  $\operatorname{argmax}_a(Q(s, a))$ , with probability  $1 - \epsilon$
- While there are many more sophisticated exploration methods,  $\epsilon$ -greedy exploration can work well on some problems

# Model-Free Control

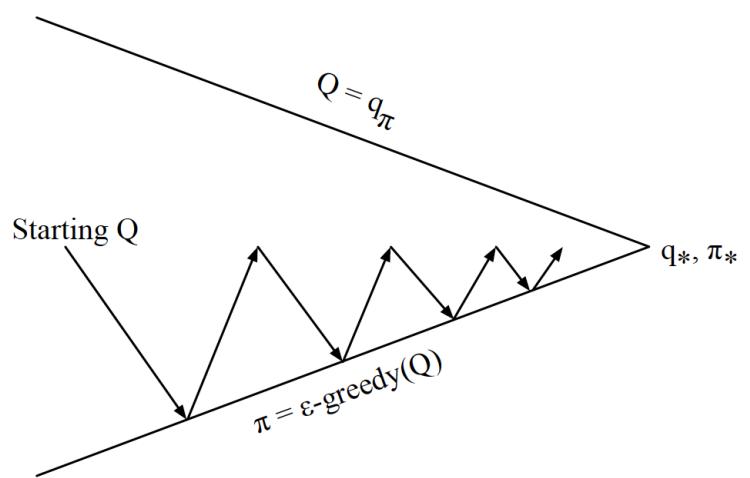
- **Policy Evaluation:** Learn an action-value function.
- **Policy Improvement:** Act epsilon greedily with respect to it.



Dynamic programming



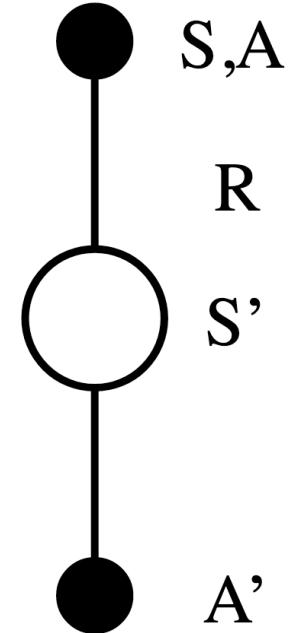
Model-free with infinite  
time to estimate  $q_\pi$



Model-free with finite time  
to estimate  $q_\pi$

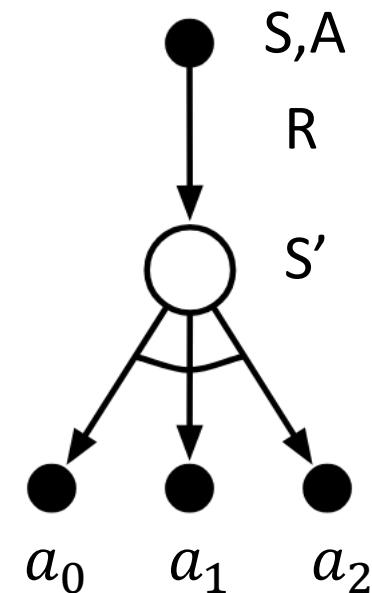
# Sarsa

- Model-free on-policy prediction (policy evaluation)
  - $V(S_t) = V(S_t) + \alpha[R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$
- Sarsa: model-free on-policy temporal-difference control
  - Sarsa: State, action, reward, state (next), action (next)
  - $Q(S_t, A_t) = Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$
  - Behavior policy: epsilon greedy
  - Target policy: epsilon greedy
  - Shown to converge to  $q_*$  if greedy in the limit with infinite exploration and if Robbins-Monro conditions hold for  $\alpha$



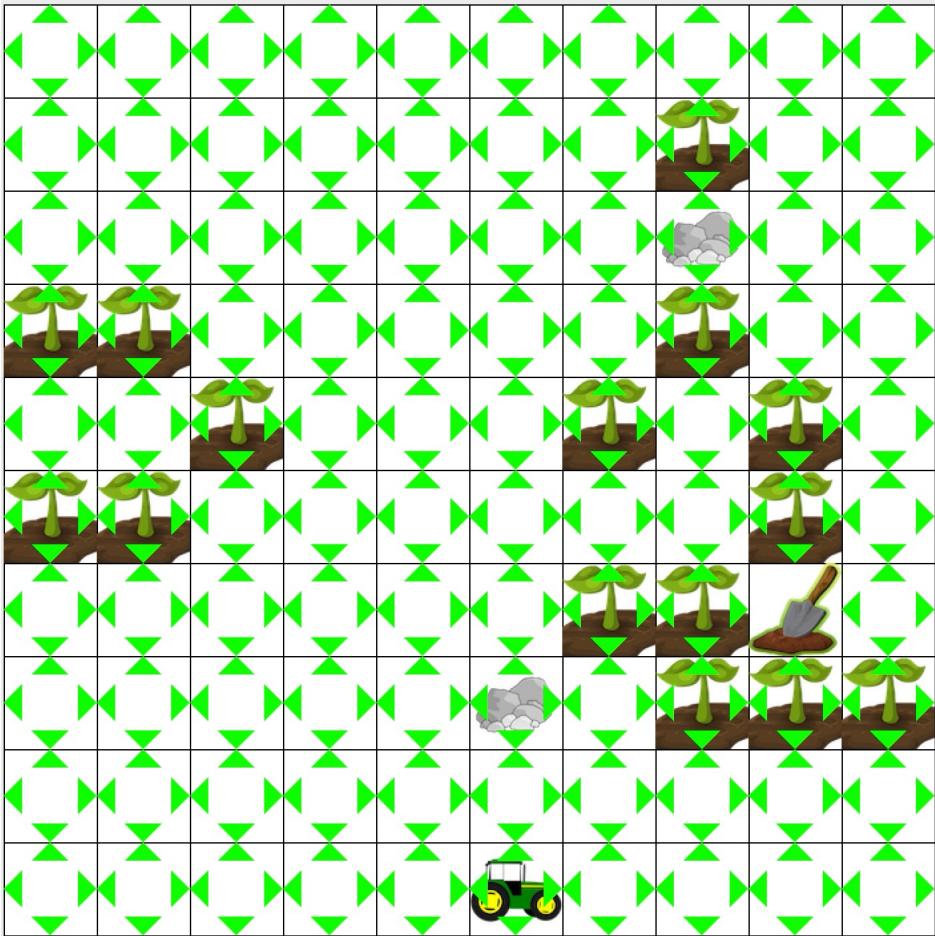
# Q-learning

- Q-learning: model-free off-policy temporal-difference control
  - $Q(S_t, A_t) = Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$
  - Behavior policy: epsilon greedy
  - Target policy: greedy
  - Converges to  $q_*$  if Robbins-Monro conditions hold for  $\alpha$

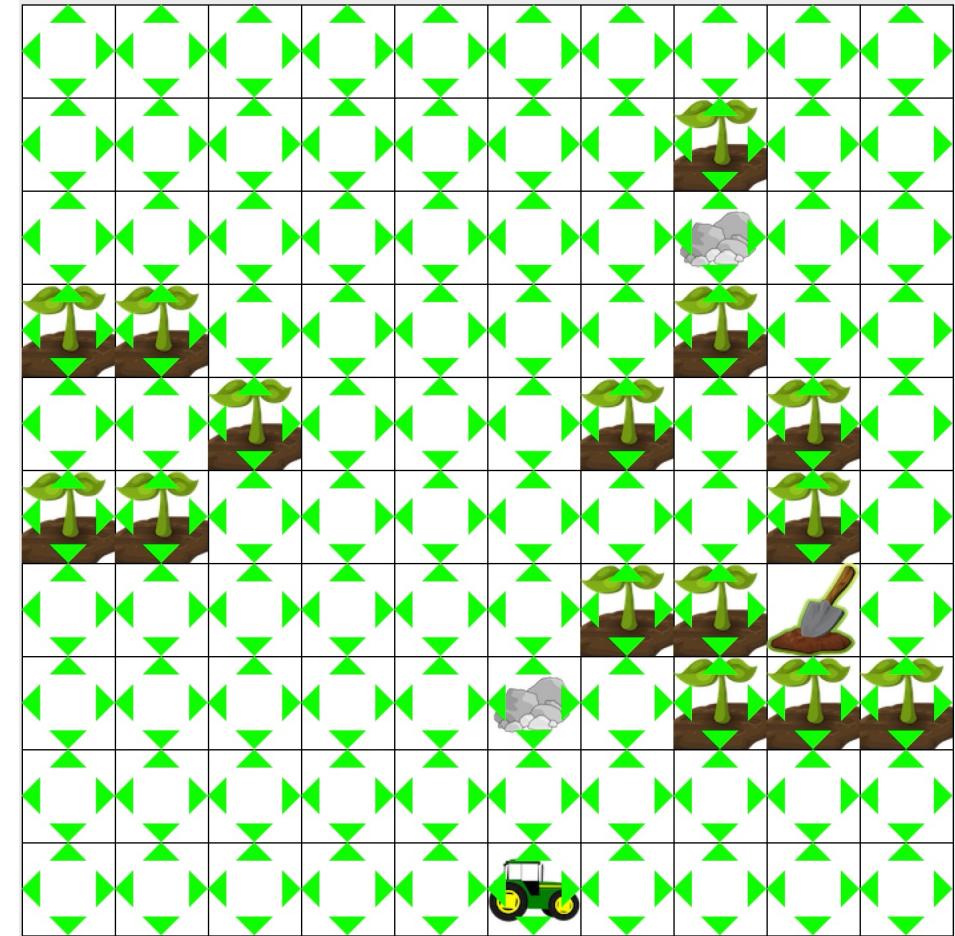


# Q-learning

- Why are the action-values for the area surrounded by plants not all red?



Q-learning step-by-step

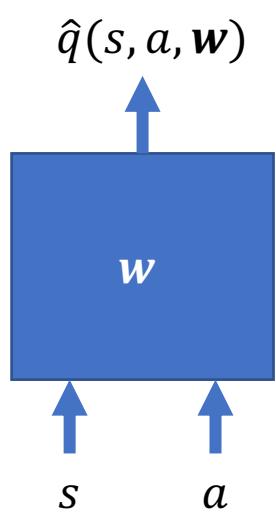


Q-learning. Showing greedy policy after every 100 episodes.

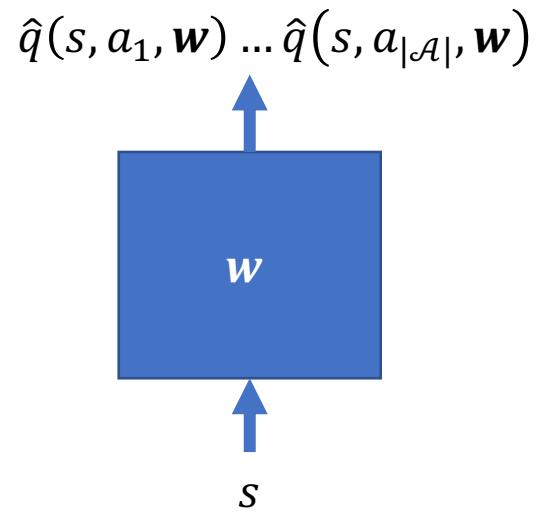
# Approximate Q-learning

- Q-learning
  - $Q(s, a) = Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$
- Approximate Q-learning
  - $y = r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w})$
  - $E(\mathbf{w}) = \frac{1}{2} (y - \hat{q}(s, a, \mathbf{w}))^2$
  - $\nabla_{\mathbf{w}} E(\mathbf{w}) = (y - \hat{q}(s, a, \mathbf{w})) \nabla_{\mathbf{w}} \hat{q}(s, a, \mathbf{w})$
  - Even though  $y$  depends on  $\mathbf{w}$ , we do not differentiate  $y$  with respect to  $\mathbf{w}$ .
  - Only supervision is that  $y = 0$  for terminal states

# Deep Q-Networks (DQNs)

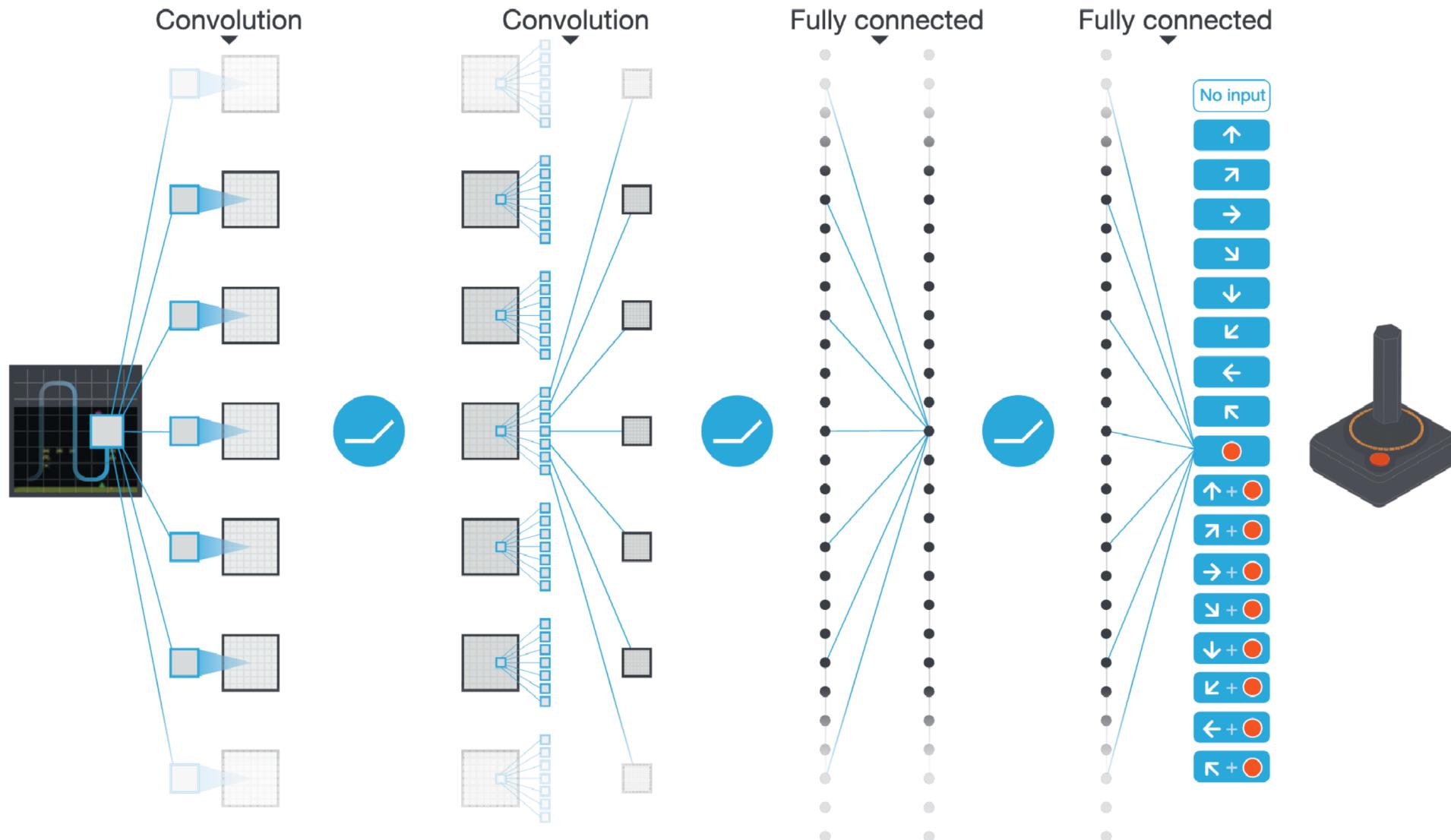


$|\mathcal{A}|$  forward  
passes needed



1 forward pass  
needed

# DQNs: Atari



# DQN: Practical Challenges

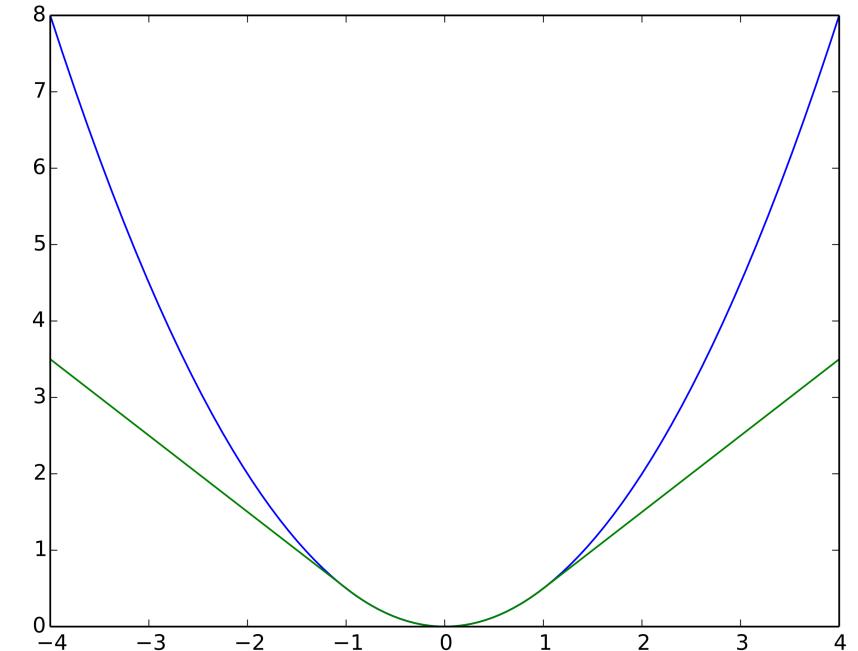
- Training instability due to non-stationary target  $y$ 
  - $y = r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w})$
  - $E(\mathbf{w}) = \frac{1}{2} (y - \hat{q}(s, a, \mathbf{w}))^2$
- Training examples are highly correlated
  - Neural networks are prone to “catastrophic forgetting”
- Training instability due to very large positive or negative rewards

# DQNs: Practical Challenges

- Instability in target  $y$ 
  - Solution: Maintain a “**target network**” with weights  $\mathbf{w}^-$
  - $y = r + \gamma \max_{a'} \hat{q}(s', a', \mathbf{w}^-)$
  - $E(\mathbf{w}) = \frac{1}{2} (y - \hat{q}(s, a, \mathbf{w}))^2$
  - Update  $\mathbf{w}^-$  to equal  $\mathbf{w}$  every N training iterations
  - Can do this for approximate value iteration, as well
- Training examples are highly correlated
  - Solution: do **experience replay**
  - Keep all of your experience in a replay buffer
  - Each element is a tuple  $(s, a, r, s')$
  - During training, randomly sample a batch from your replay buffer
  - Decorrelates training examples

# DQNs: Huber Loss

- Training instability due to very large positive or negative rewards
- Clips loss after reaching a certain magnitude

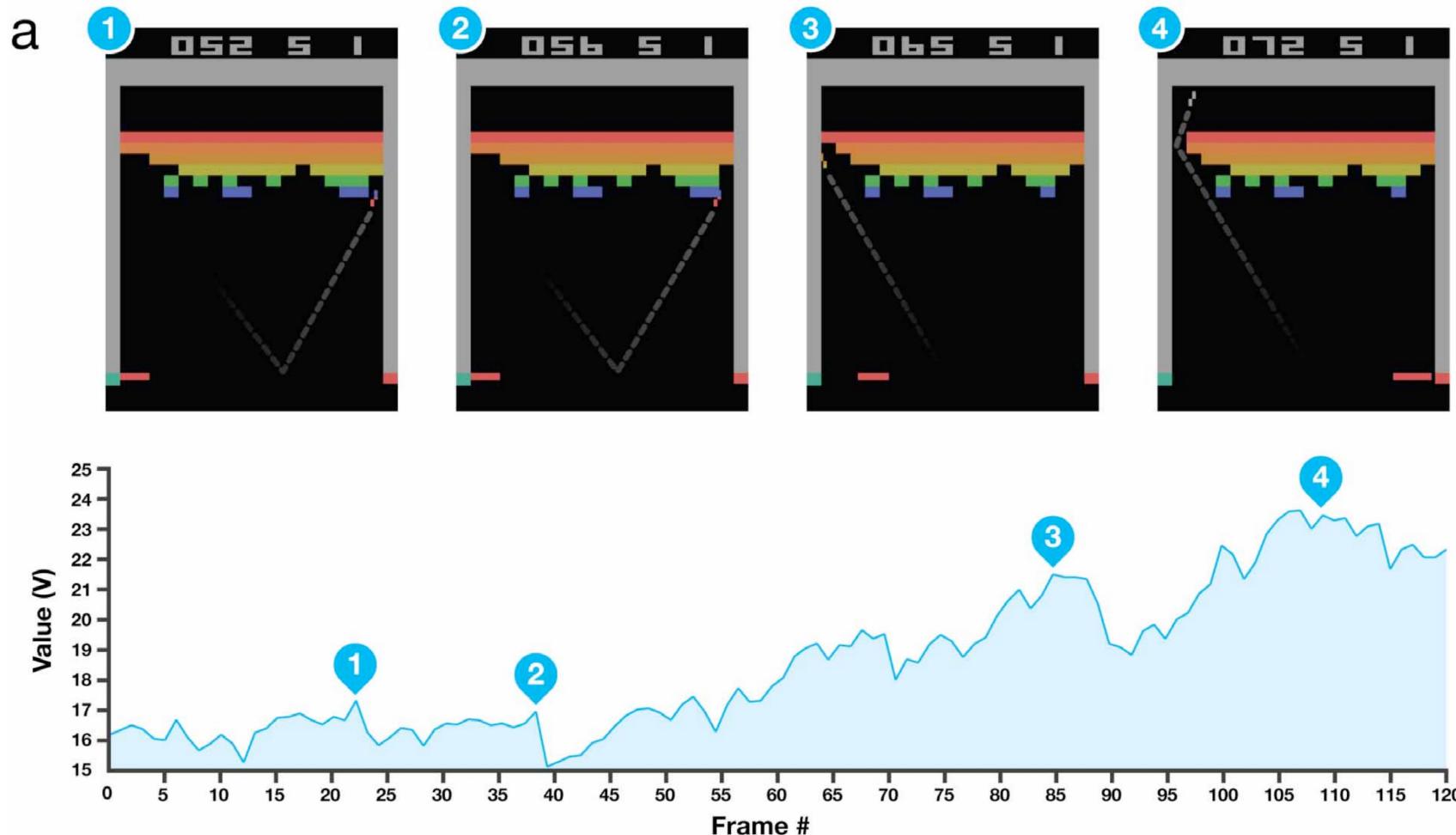


# DQNs: Example on Atari



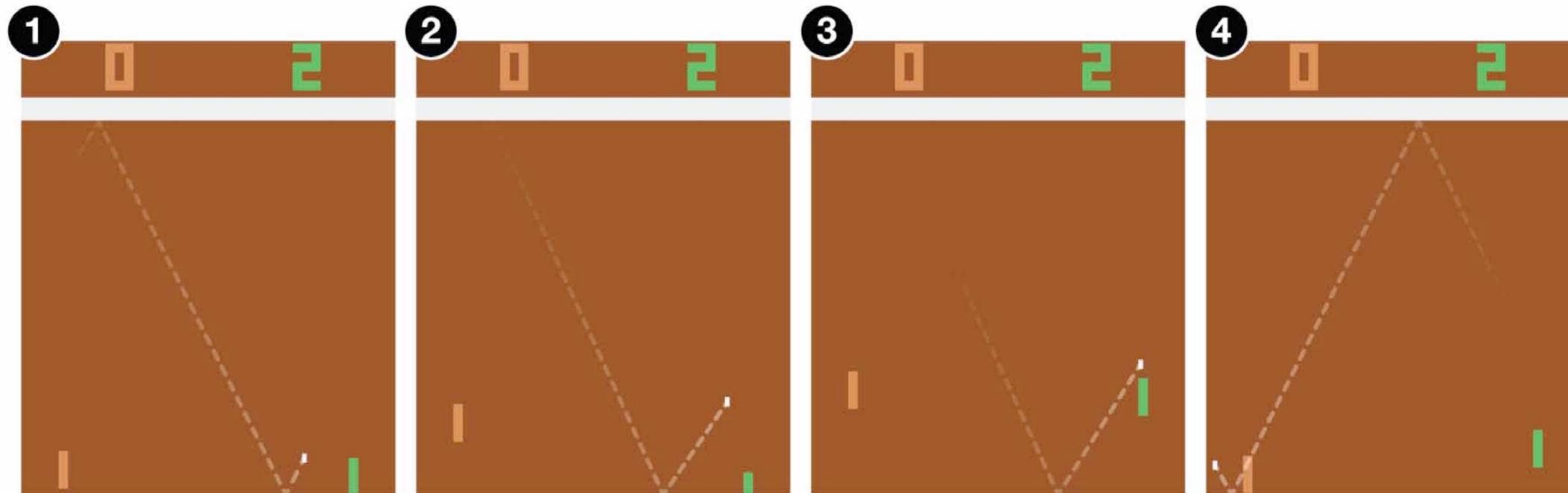
# DQN: State Value Behavior

- Value fluctuation due to discounting as well as nature of imperfect function approximation

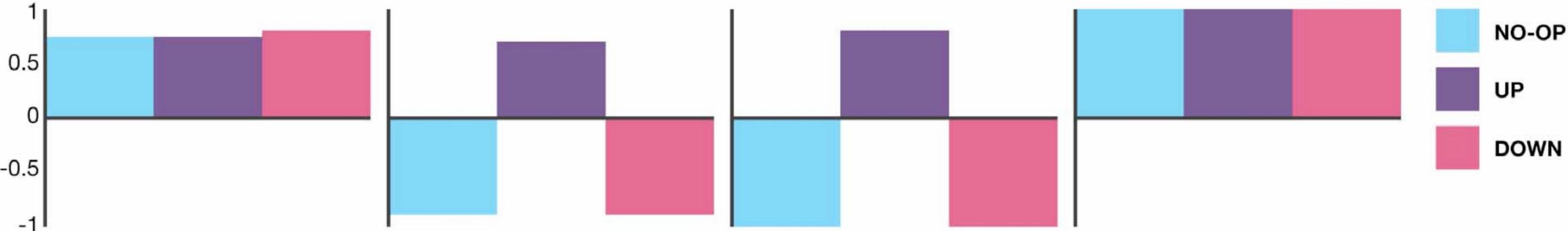


# DQN: Action-Value Behavior

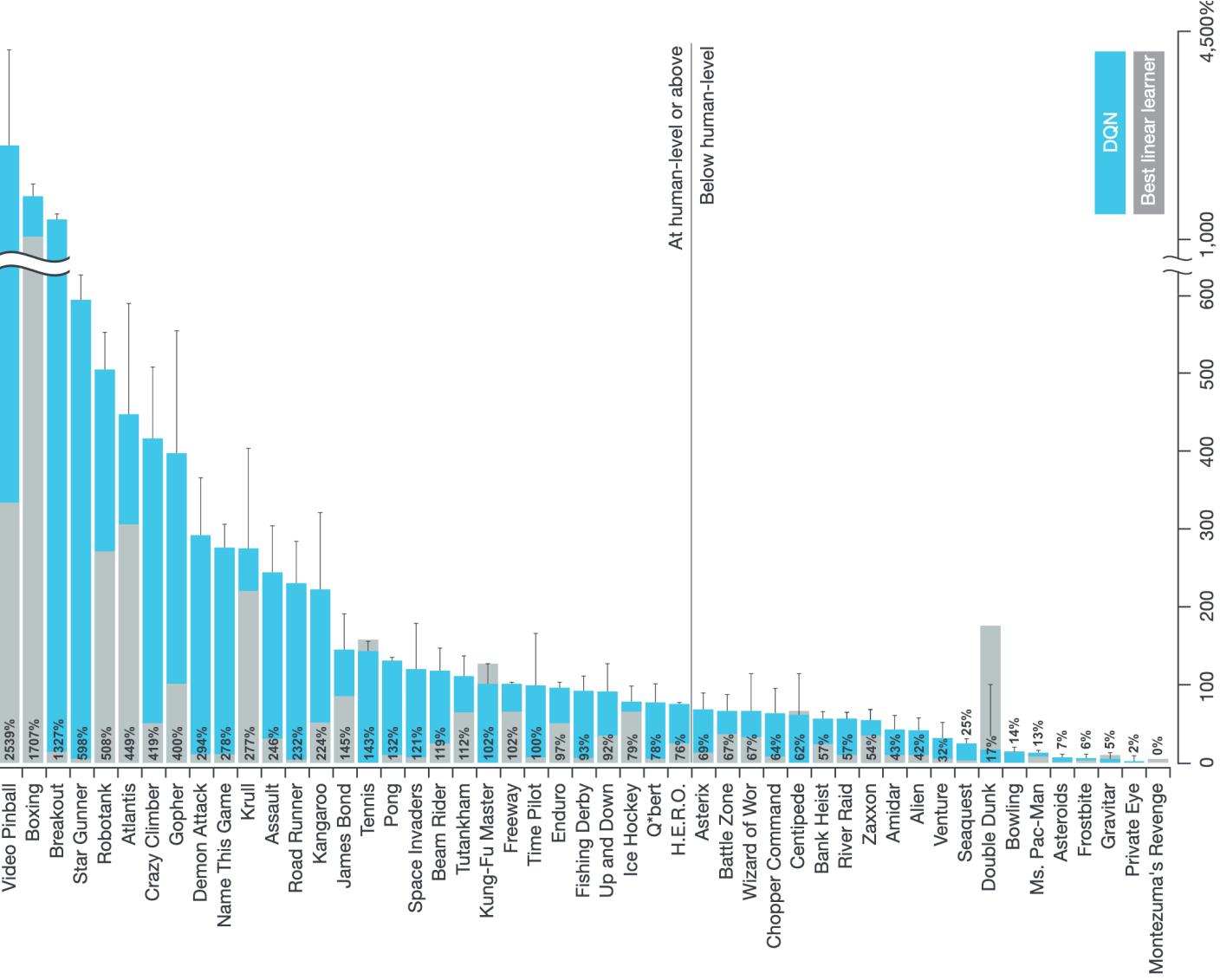
b



Action-Values (Q)



# DQN: Performance on Atari



# DQN: Ablation Study

<b>Game</b>	<b>With replay, with target Q</b>	<b>With replay, without target Q</b>	<b>Without replay, with target Q</b>	<b>Without replay, without target Q</b>
Breakout	316.8	240.7	10.2	3.2
Enduro	1006.3	831.4	141.9	29.1
River Raid	7446.6	4102.8	2867.7	1453.0
Seaquest	2894.4	822.6	1003.0	275.8
Space Invaders	1088.9	826.3	373.2	302.0

# DQN: Ablation Study

- How do deep neural networks compare with a linear function approximator?

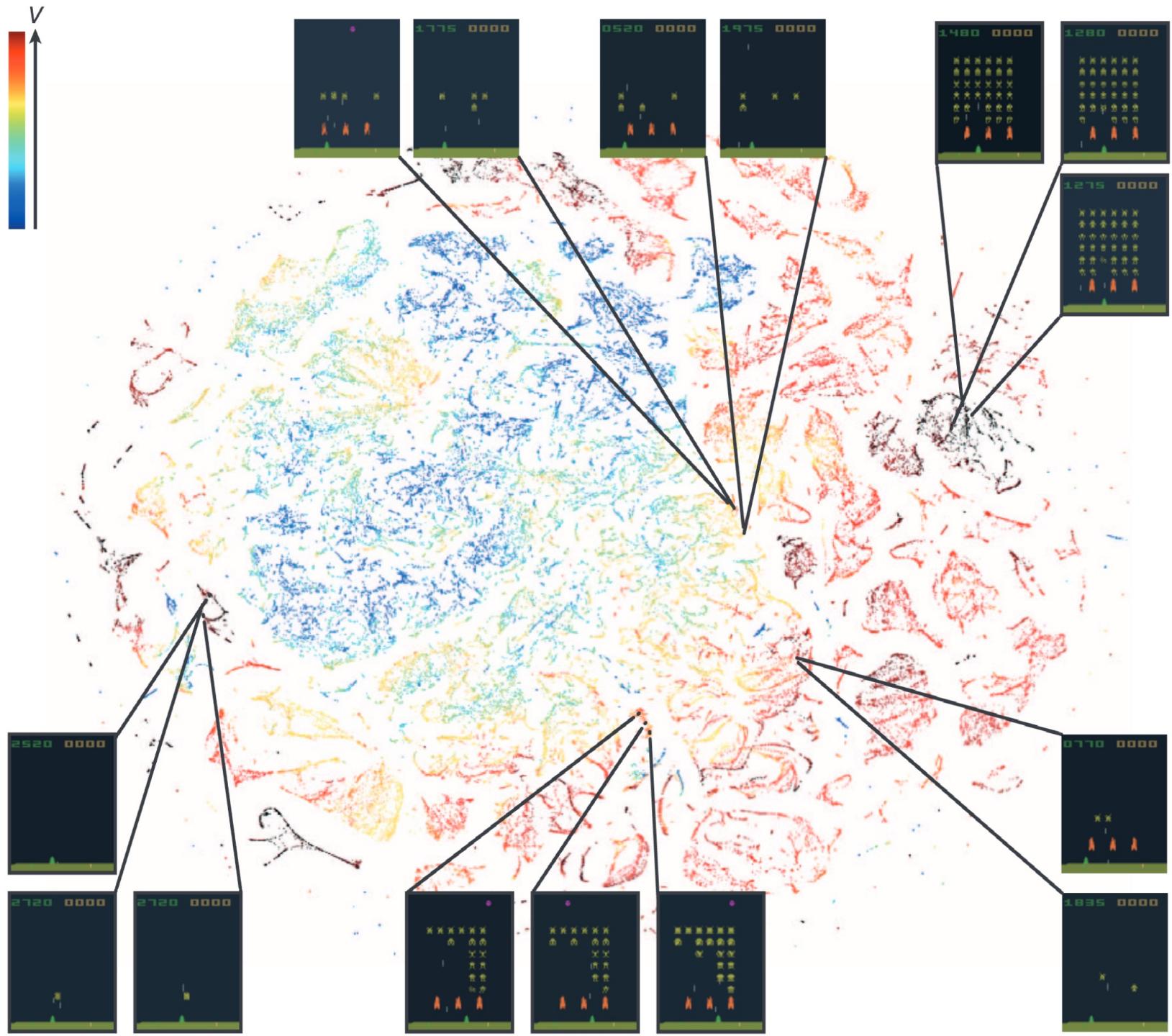
Game	DQN	Linear
Breakout	316.8	3.00
Enduro	1006.3	62.0
River Raid	7446.6	2346.9
Seaquest	2894.4	656.9
Space Invaders	1088.9	301.3

# DQN: Hyperparameters

Hyperparameter	Value	Description
minibatch size	32	Number of training cases over which each stochastic gradient descent (SGD) update is computed.
replay memory size	1000000	SGD updates are sampled from this number of most recent frames.
agent history length	4	The number of most recent frames experienced by the agent that are given as input to the Q network.
target network update frequency	10000	The frequency (measured in the number of parameter updates) with which the target network is updated (this corresponds to the parameter C from Algorithm 1).
discount factor	0.99	Discount factor gamma used in the Q-learning update.
action repeat	4	Repeat each action selected by the agent this many times. Using a value of 4 results in the agent seeing only every 4th input frame.
update frequency	4	The number of actions selected by the agent between successive SGD updates. Using a value of 4 results in the agent selecting 4 actions between each pair of successive updates.
learning rate	0.00025	The learning rate used by RMSProp.
gradient momentum	0.95	Gradient momentum used by RMSProp.
squared gradient momentum	0.95	Squared gradient (denominator) momentum used by RMSProp.
min squared gradient	0.01	Constant added to the squared gradient in the denominator of the RMSProp update.
initial exploration	1	Initial value of $\epsilon$ in $\epsilon$ -greedy exploration.
final exploration	0.1	Final value of $\epsilon$ in $\epsilon$ -greedy exploration.
final exploration frame	1000000	The number of frames over which the initial value of $\epsilon$ is linearly annealed to its final value.
replay start size	50000	A uniform random policy is run for this number of frames before learning starts and the resulting experience is used to populate the replay memory.
no-op max	30	Maximum number of "do nothing" actions to be performed by the agent at the start of an episode.

# DQN: Last Hidden Layer Activations

- Done for space invaders
- Two-dimensional projection using t-SNE
- Semantically similar states are clustered next to each other



# DQN: Last Hidden Layer Activations

- How does network generalize to states obtained under other policies?
  - Human: Orange
  - AI Agent: Blue

