

Reinforcement Learning

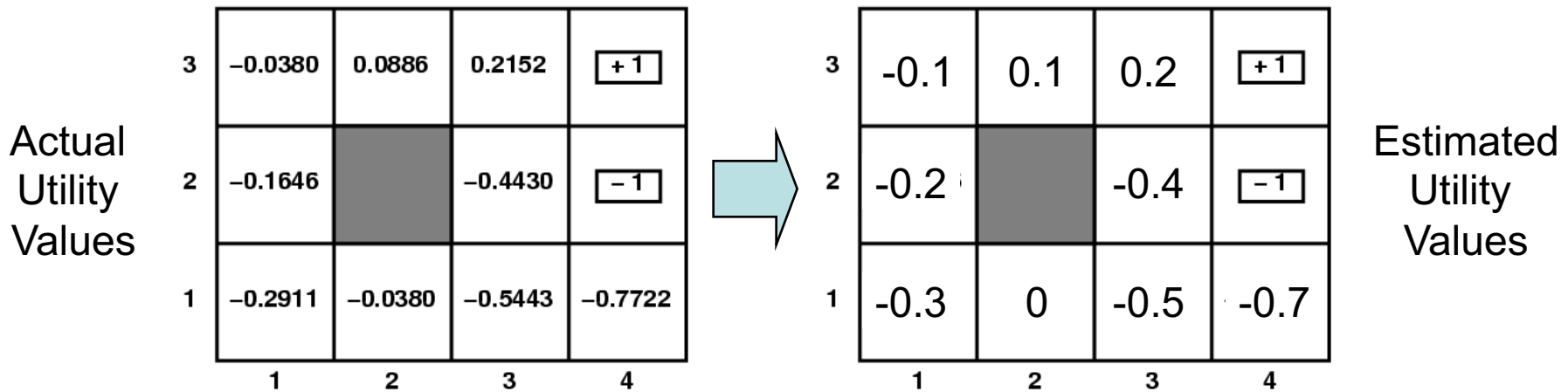
CPSC 470 – Artificial Intelligence

Brian Scassellati

From Supervised Learning to Reinforcement Learning

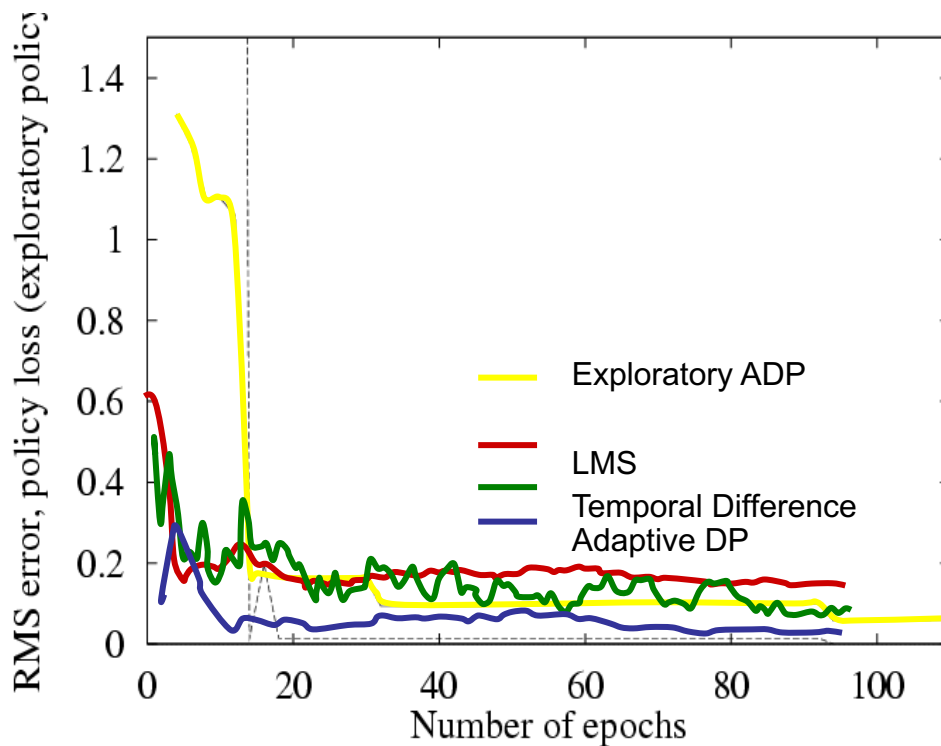
- In **supervised learning**, when the agent makes a mistake, it is immediately given feedback in the form of the **correct response**
- In **reinforcement learning**, when the agent makes a mistake, it will later be given feedback in the form of a **punishment or reward**

Passive Learning in Known Environments



- Given a set of training sequences that end in a terminal state (with a reward)
 - $(1,1) \rightarrow (2,1) \rightarrow (3,1) \rightarrow (3,2) \rightarrow (3,3) \rightarrow (4,3) \rightarrow +1$
 - $(1,1) \rightarrow (2,1) \rightarrow (3,1) \rightarrow (3,2) \rightarrow (4,2) \rightarrow -1$
 - $(1,1) \rightarrow (1,2) \rightarrow (1,3) \rightarrow (2,3) \rightarrow (3,3) \rightarrow (4,3) \rightarrow +1$
 - $(1,1) \rightarrow (2,1) \rightarrow (3,1) \rightarrow (4,1) \rightarrow (3,1) \rightarrow (3,2) \rightarrow (4,2) \rightarrow -1$
- Determine the expected utility $U(i)$ associated with each non-terminal state i

Comparison of Reinforcement Learning (Utility) Techniques



What if we have no good
model of the environment?

Two Basic Types of Reinforcement Learning

Utility Learning

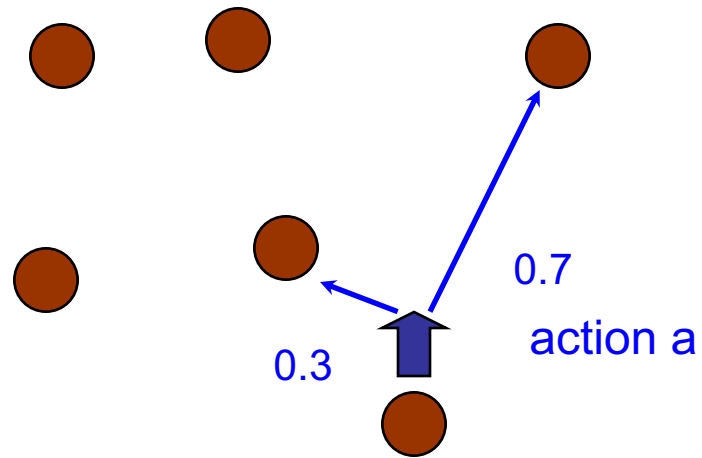
- Learn a utility function that maps states to utilities and select an action by maximizing the expected value
- Needs a model of the environment (needs to know which state an action will result in)
- Predictive

Action-Value Learning

- Learn an action-value function that gives the expected utility of taking a given action in a given state
- No need for an environment model
- Do not know where actions lead, so it cannot look ahead

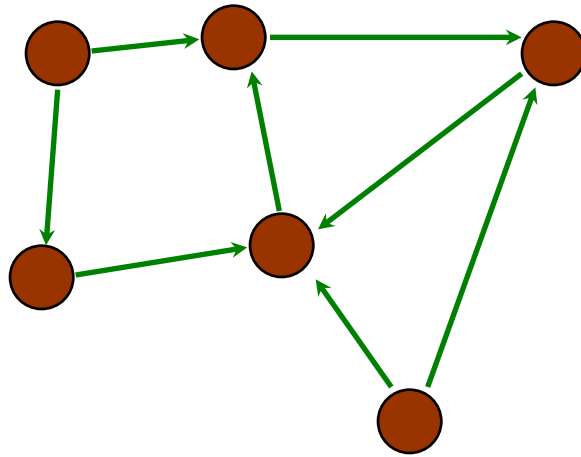
States and Actions

Environment = states



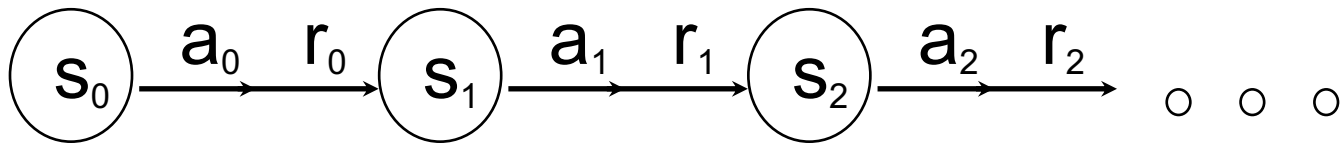
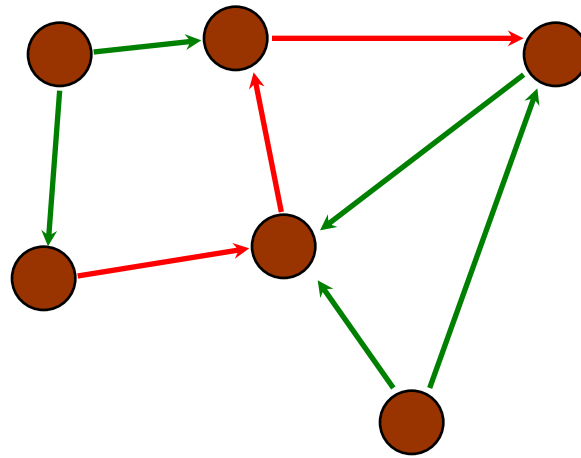
Actions = transitions $\delta(s, a, s')$

Rewards

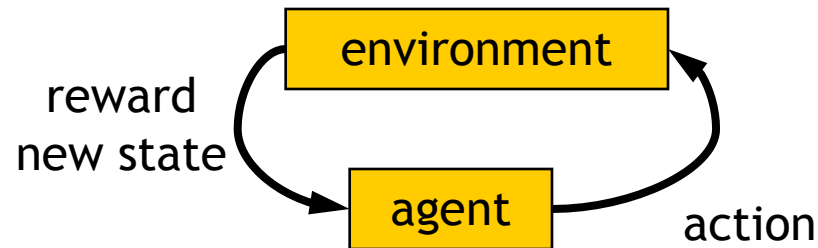


$R(\mathbf{s}, \mathbf{a})$ = reward at state \mathbf{s}
for doing action \mathbf{a}

Trajectories



Markov Decision Process (MDP)



- set of states S , set of actions A , initial state S_0
- transition model $P(s'|s,a)$
 - Markov assumption
- Reward function $R(s,a)$
- policy: mapping from S to A
 - $\pi(s)$ or $\pi(s,a)$

Agent Learns a Policy

Policy at step t , π_t :

a mapping from states to action probabilities

$\pi_t(s, a) =$ probability that $a_t = a$ when $s_t = s$

- Reinforcement learning methods specify how the agent **changes its policy** as a result of experience.
- Roughly, the agent's goal is to get as much reward as it can over the long run.

Goals and Rewards

- Is a scalar reward signal an adequate notion of a goal?
 - Maybe not, but it is surprisingly flexible.
- A goal should specify **what** we want to achieve, not **how** we want to achieve it.
- A goal must be outside the agent's direct control—thus outside the agent.
- The agent must be able to measure success:
 - explicitly;
 - frequently during its lifespan.

Returns

Suppose the sequence of rewards after step t is :

$$r_{t+1}, r_{t+2}, r_{t+3}, \dots$$

What do we want to maximize?

In general,

we want to maximize the **expected return**, $E\{R_t\}$, for each step t .

Episodic tasks: interaction breaks naturally into episodes, e.g., plays of a game, trips through a maze.

$$R_t = r_{t+1} + r_{t+2} + \dots + r_T,$$

where T is a final time step at which a **terminal state** is reached, ending an episode.

Returns for Continuing Tasks

Continuing tasks: interaction does not have natural episodes.

Discounted return:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1},$$

where γ , $0 \leq \gamma \leq 1$, is the **discount rate**.

shortsighted $0 \leftarrow \gamma \rightarrow 1$ farsighted

(Learning sequences of actions that generate rewards)

Q-LEARNING

Learning Sequences

- Q-Learning allows an agent to learn *chains* of actions.
- Even though the agent lives only in the present, it acts as if it can see into the future
 - Cannot even predict the next state
- Propagation of credit from the *consummatory* behavior back through the chain of *appetitive* behaviors

Temporal Assignment Problem

- How do you determine which actions successfully produced the goal state?
- Q-learning solves this problem without needing to explicitly remember the sequence of states that lead to a reward.
- The price you pay is that it depends upon repeated visits to each state/action combination.

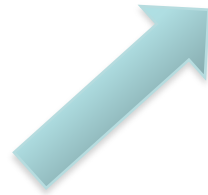
Q-Values

- Big idea:
 - Compute the quality value (or Q-value) for each possible action a in state x
 - Choose actions stochastically on the Q-values

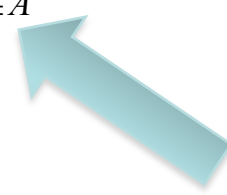
Q-Learning Approach

- Update the Q-values according to the following rule:

$$Q(x, a) \equiv (1 - \alpha)Q(x, a) + \alpha(r + \gamma \max_{b \in A} Q(y, b))$$



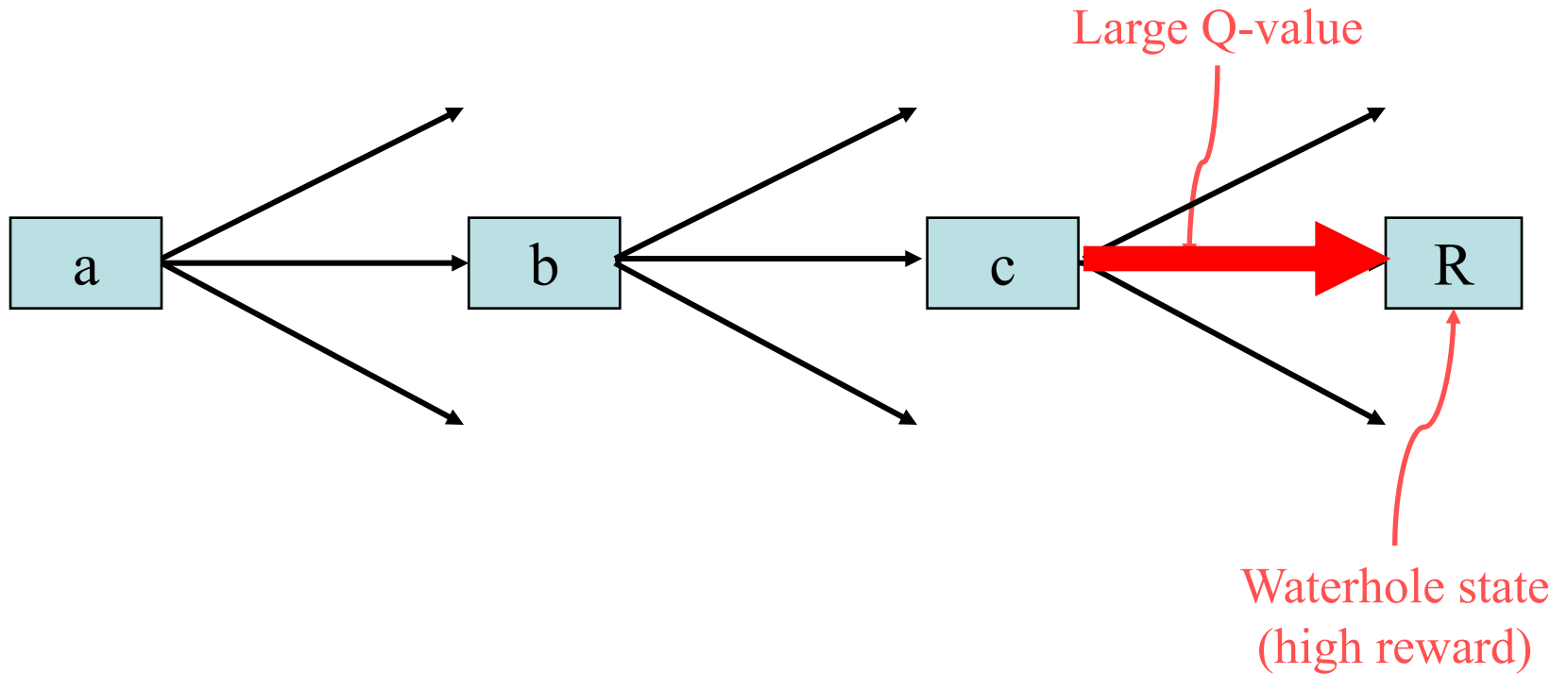
Discounted prior Q value



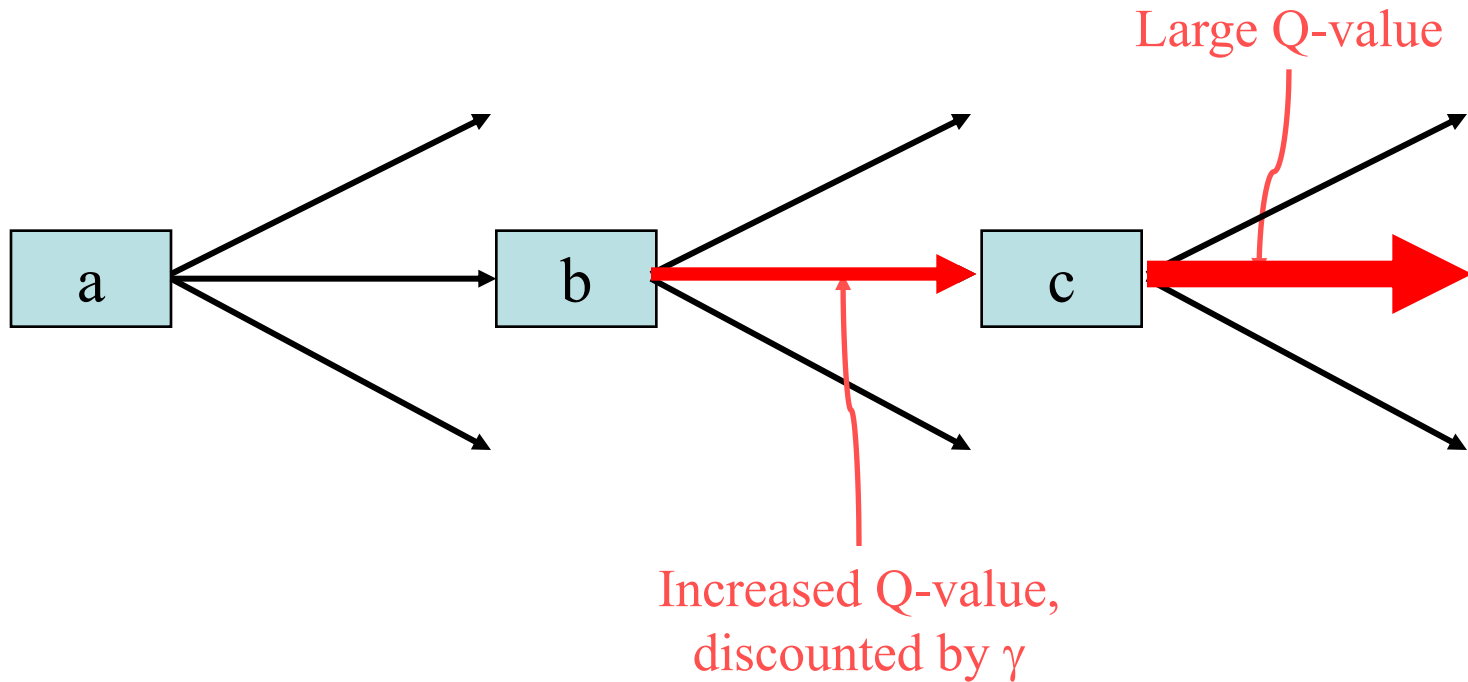
Updated with new Q value

- For a given learning rate α and a discount rate of γ

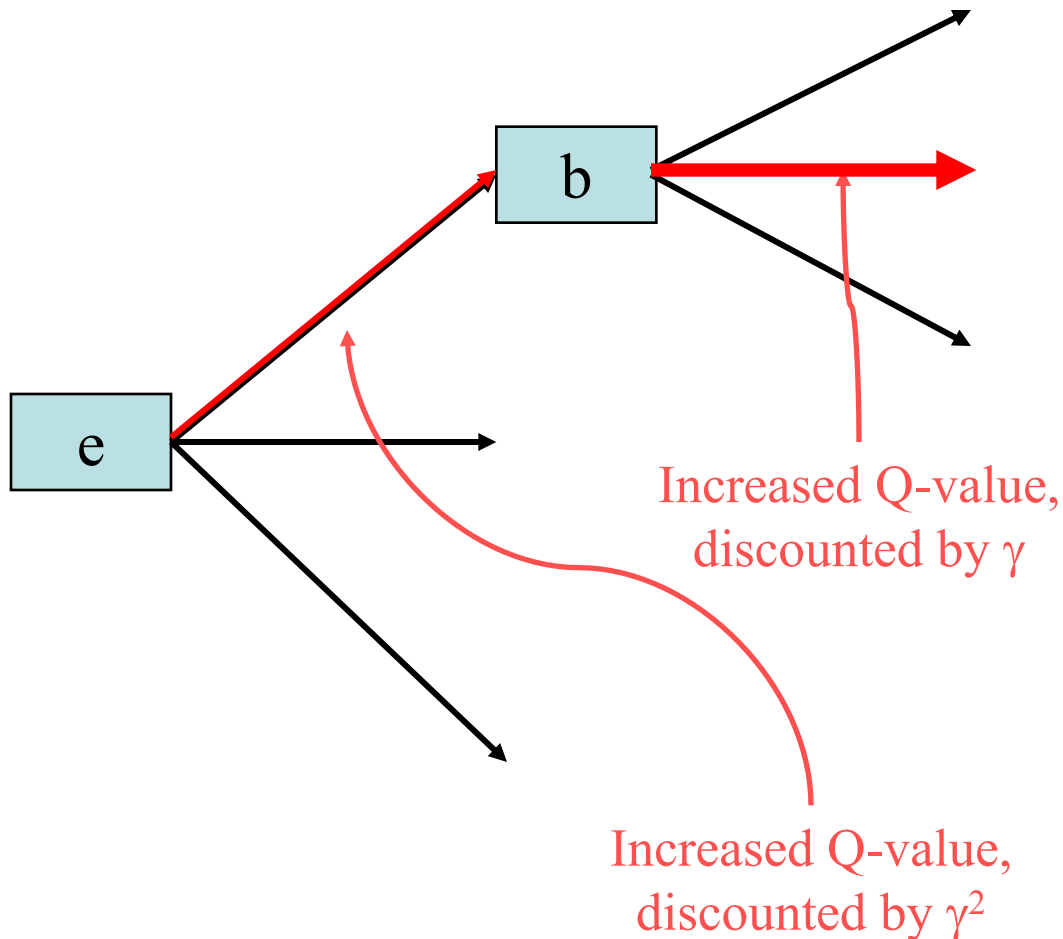
How Q-Learning Works



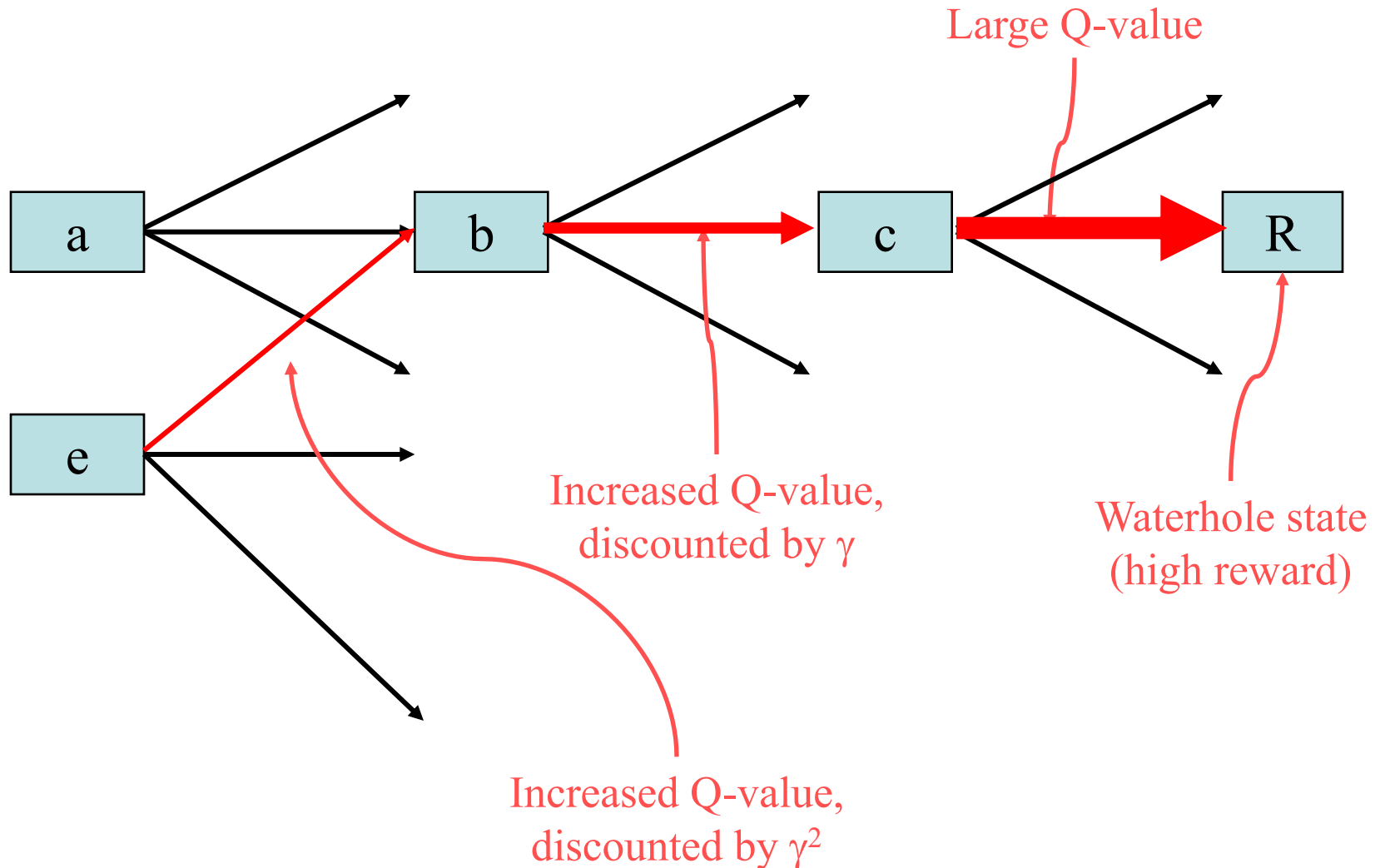
How Q-Learning Works



How Q-Learning Works



How Q-Learning Works



How to Change the Learning Rate

- What happens in the standard update rule if we just gradually drop the learning rate over time?

$$Q(x, a) \equiv (1 - \alpha)Q(x, a) + \alpha(r + \gamma \max_{b \in A} Q(y, b))$$

How to Change the Learning Rate

- What happens in the standard update rule if we just gradually drop the learning rate over time?

$$Q(x, a) \equiv (1 - \alpha)Q(x, a) + \alpha(r + \gamma \max_{b \in A} Q(y, b))$$

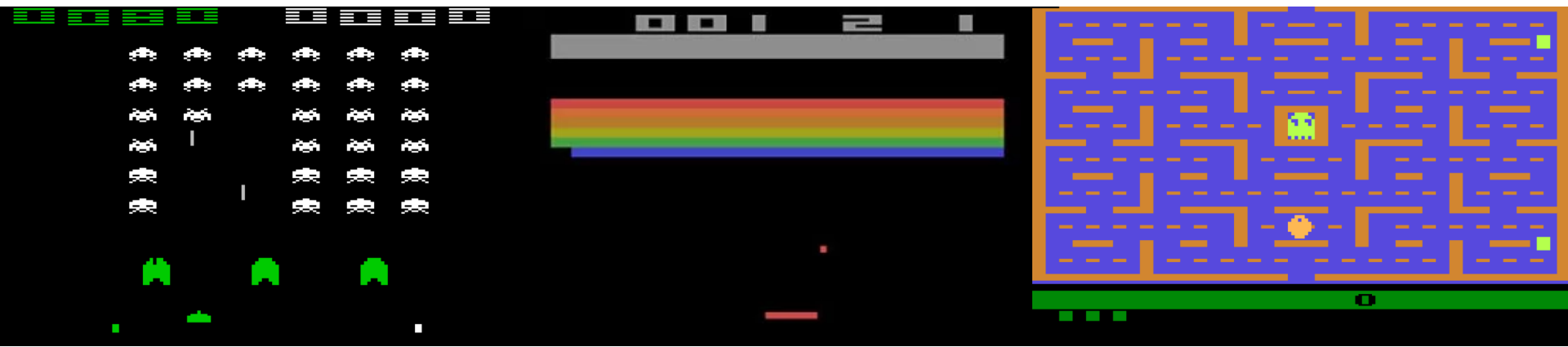
- Solution: Make a different learning rate for each state/action combination based on the number of times visited:

$$\alpha(x, a) = \frac{1}{n(x, a)}$$

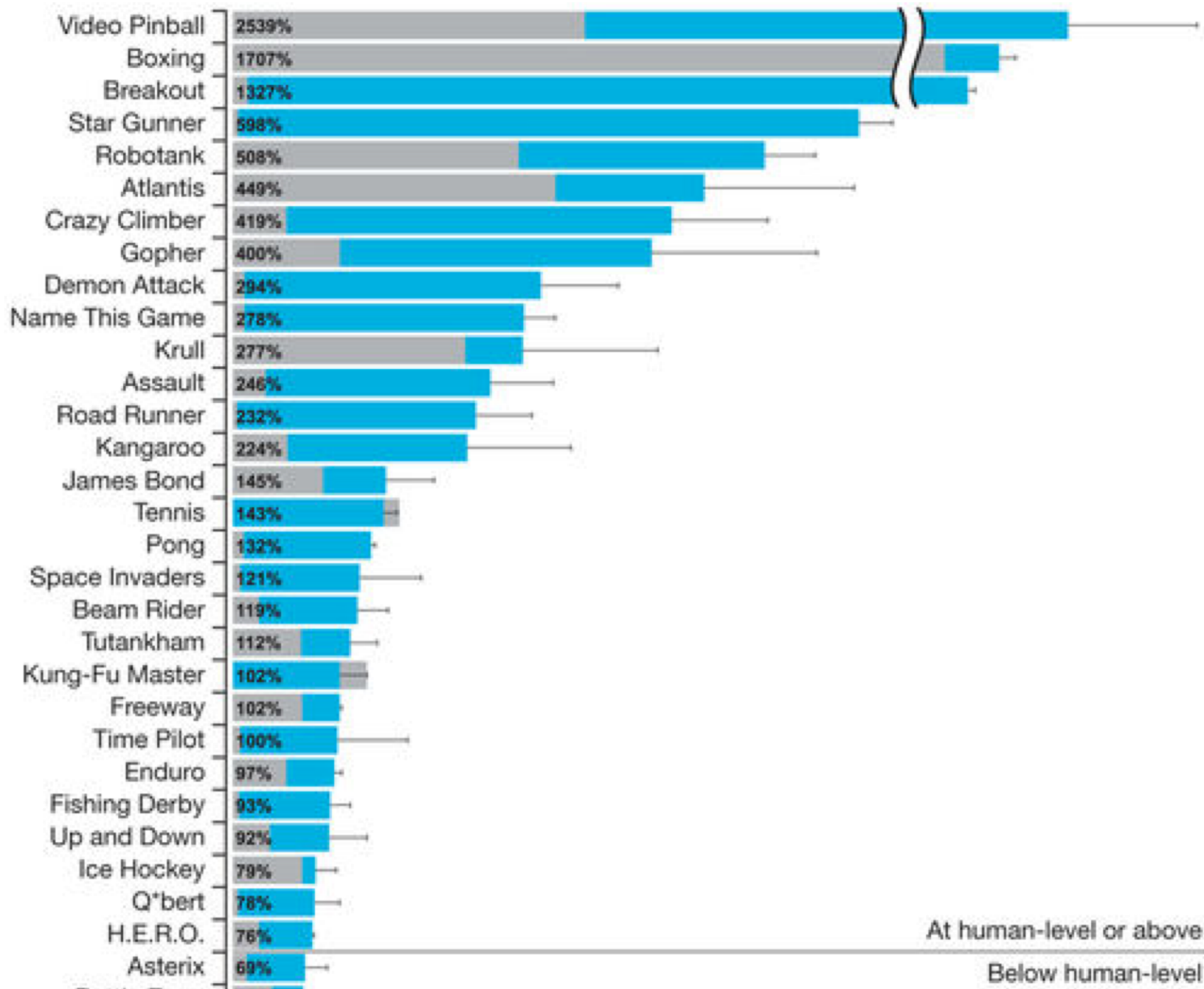
Problems with RL

- Clearly defined states and actions
- Infinite visits to each state
- No look-ahead
 - States with no exit

Reinforcement Learning for Classic Arcade Games

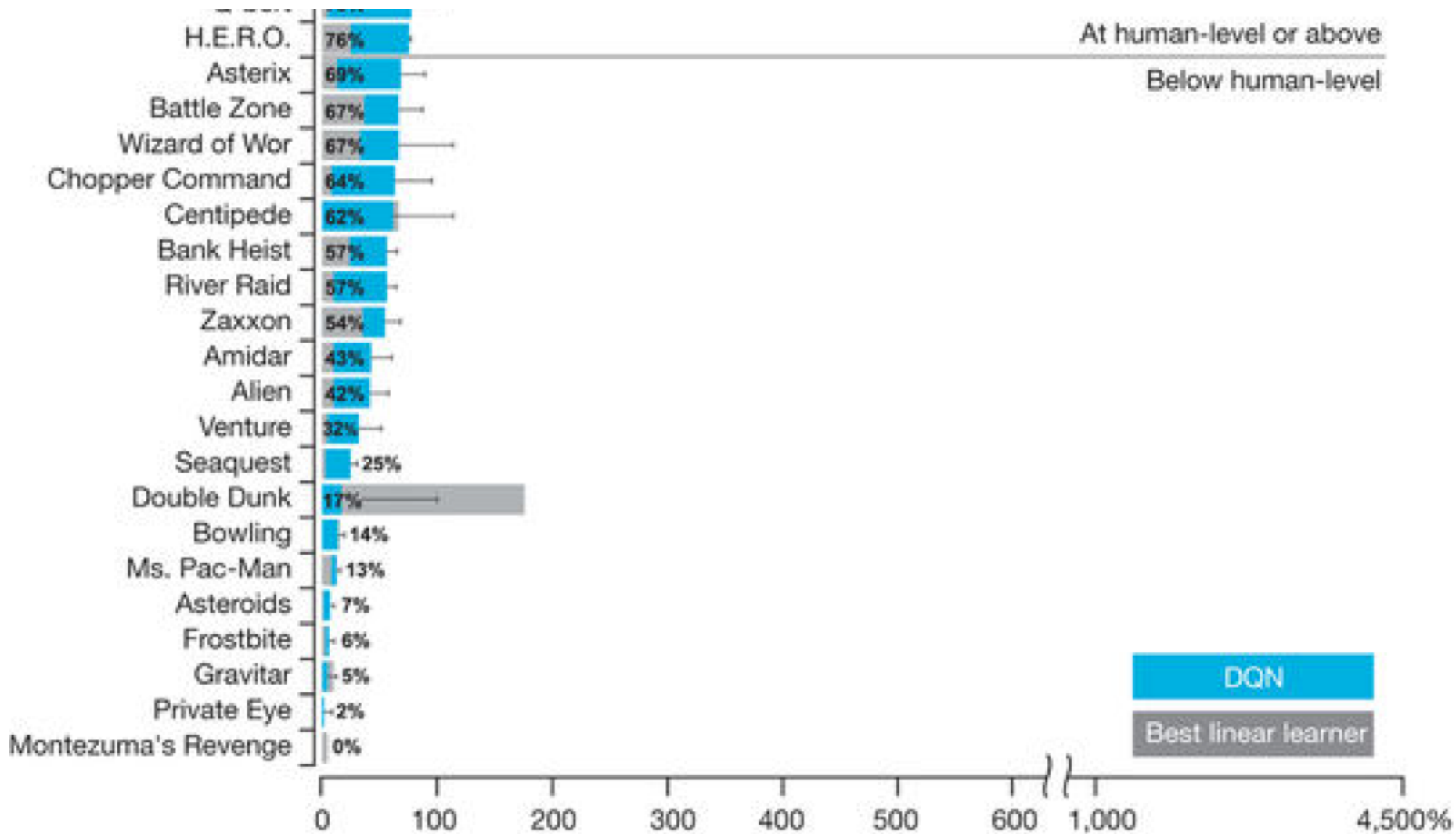


- 2014 Nature paper from Google's Deep Mind project
- Input: Pixels & score (reward)
- Output: joystick controls & button press



At human-level or above

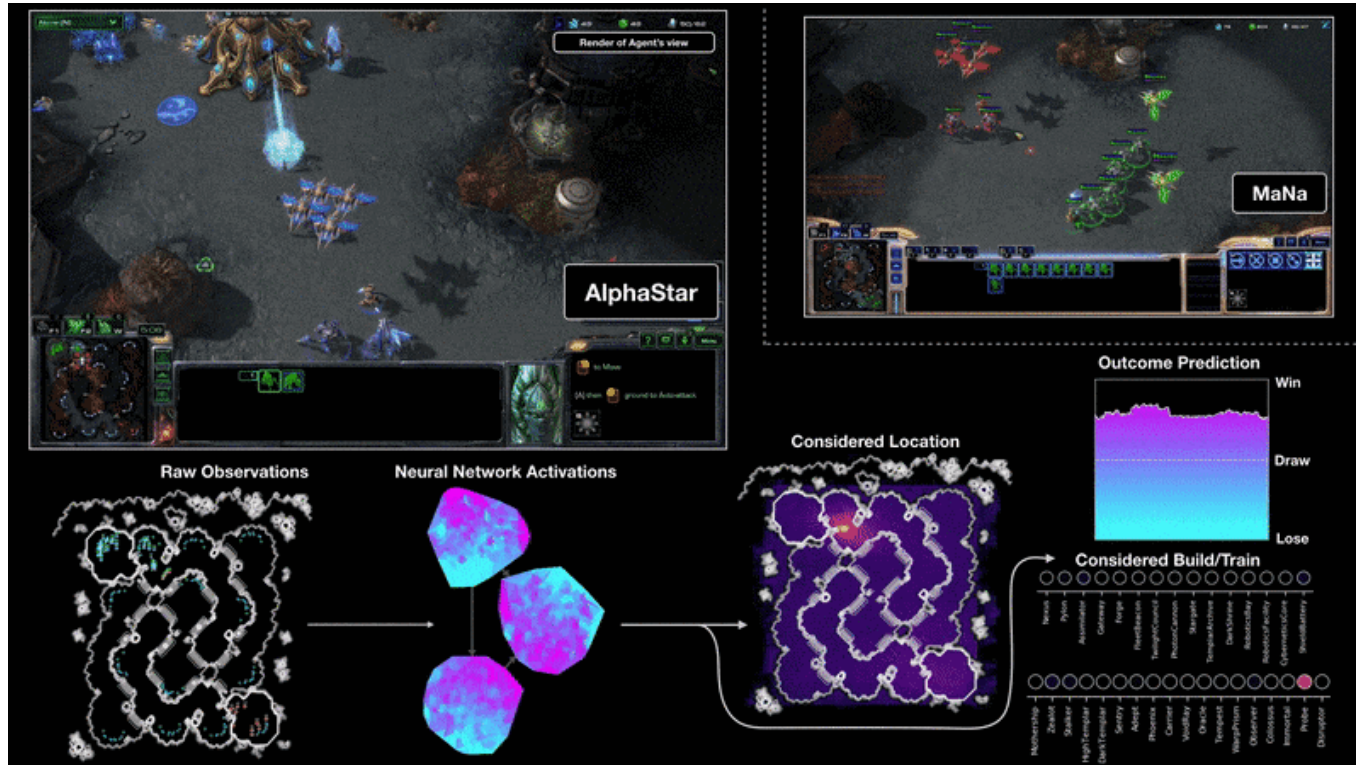
Below human-level



Why does RL fail on Montezuma's Revenge?



Starcraft II



- Real-time, imperfect information, long-term planning and reward, ~100 decisions per second
- Beat 4 top human players in Dec 2018, 10 to 1
- Trained with RL and Deep networks

Relative Complexity

Game	Board Size	State-Space Complexity	Year defeated
Tic Tac Toe	9	10^3	1952*
Connect 4	42	10^{13}	1995*
Backgammon	28	10^{20}	1979
Chess	64	10^{47}	1997
Go (19x19)	361	10^{170}	2015
Heads up NL Holdem	N/A	10^{180}	2017
StarCraft II	N/A	10^{1685}	???

Administrivia

- PS 5 due tonight
- PS 6 out today/tomorrow
- This week: Unsupervised Learning and Natural Language Processing