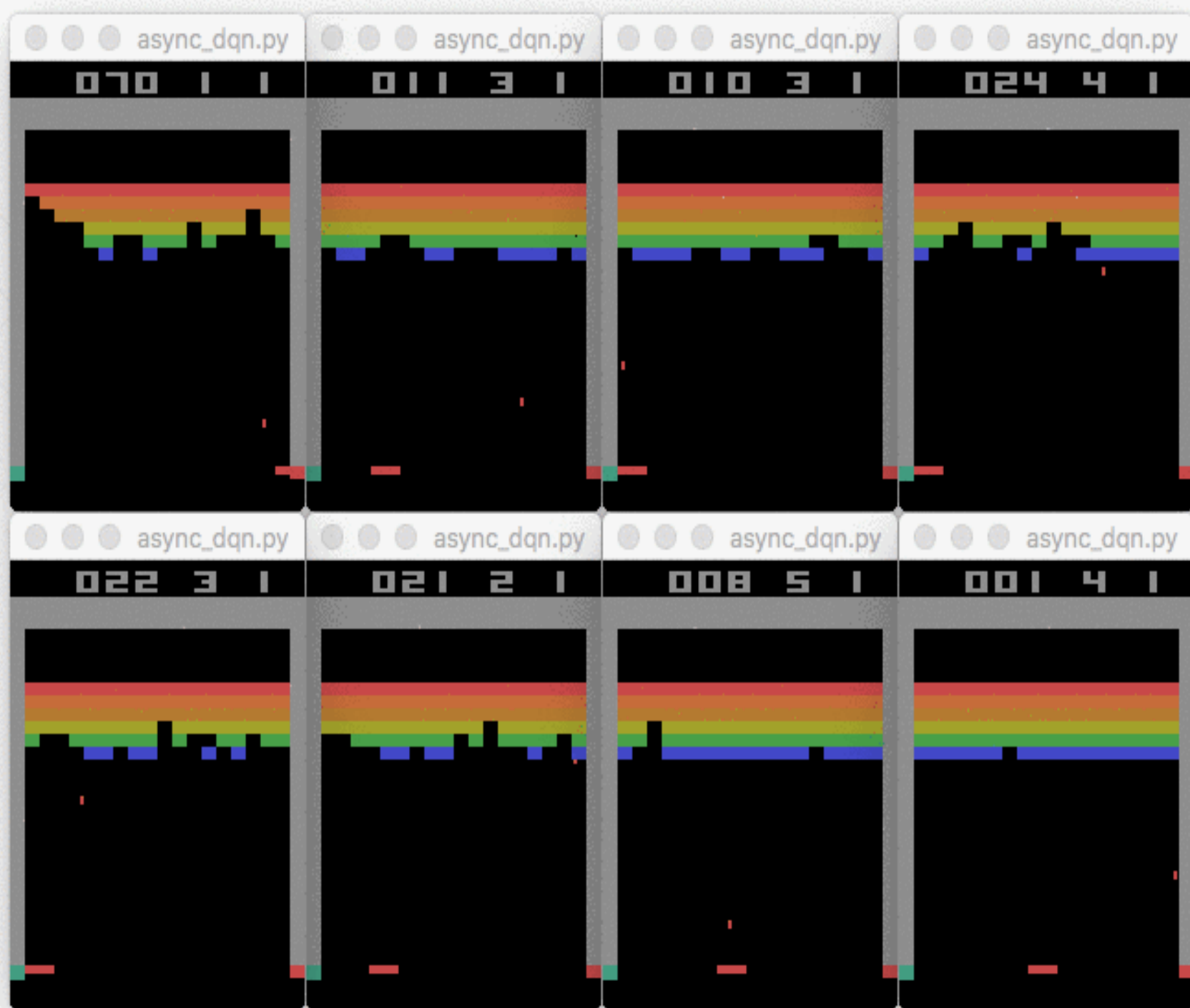# Lecture 13: Deep Reinforcement Learning

Deep Learning @ UvA

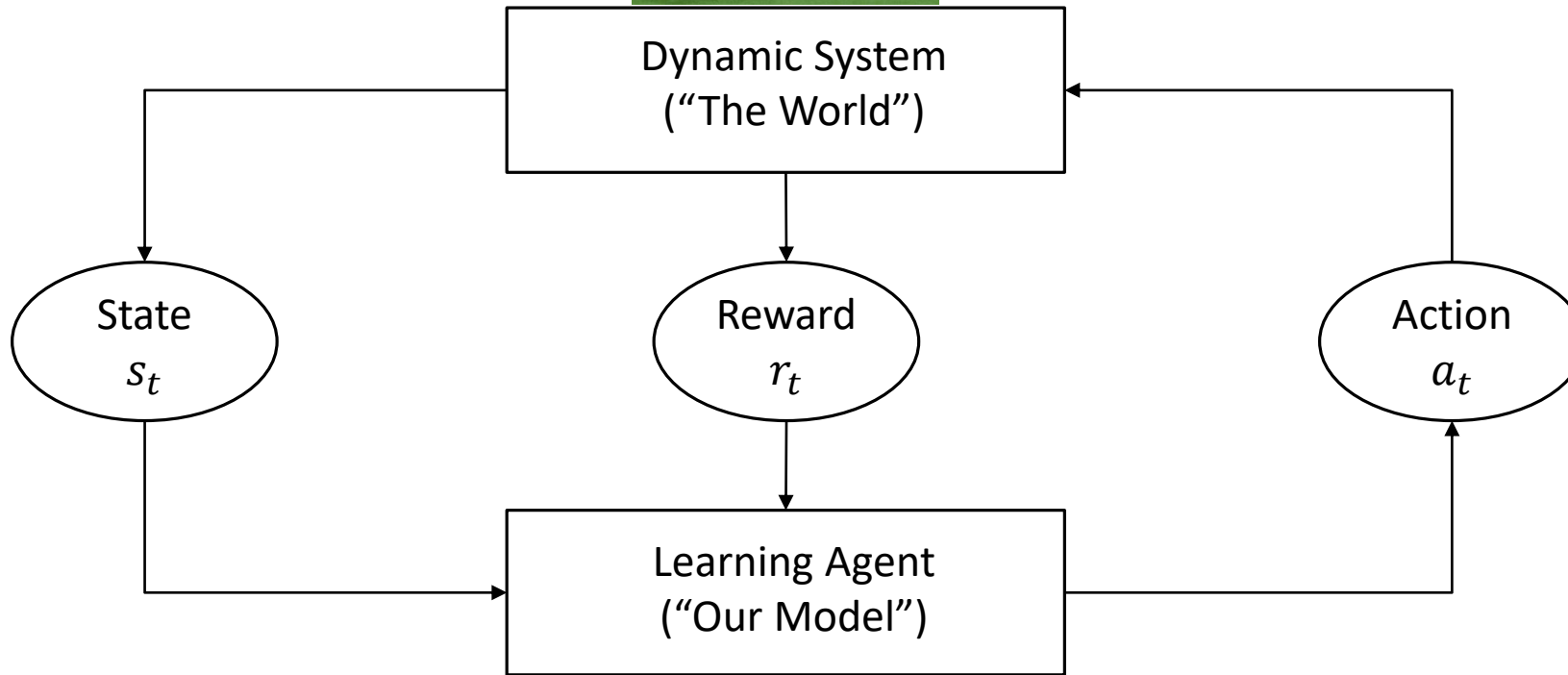# Reinforcement Learning

# What is Reinforcement Learning?

o General purpose framework for learning Artificial Intelligence models

o RL assumes that *the agent* (our model) can take *actions*

o These actions affect *the environment* where *the agent* operates, more specifically *the state* of the environment and *the state* of the agent

o Given the state of the environment and the agent, an action taken from the agent causes a reward (can be positive or negative)

o Goal: the goal of an RL agent is to learn how to take actions that maximize future rewards

# Some examples of RL

# Some examples of RL

o Controlling physical systems
  ◦ Robot walking, jumping, driving

o Logistics
  ◦ Scheduling, bandwidth allocation

o Games
  ◦ Atari, Go, Chess, Pacman

o Learning sequential algorithms
  ◦ Attention, memory

# Reinforcement Learning: An abstraction

*Slides inspired by P. Abbeel*

# How do we decide about actions, states, rewards?

o We model them as functions

o The *policy function* $a_t = \pi(s_t)$ selects an action given the current state

o The *value function* $Q^\pi(s_t, a_t)$ is the expected total reward that we will receive if we take action $a_t$ given state $s_t$

o What should our goal then be?

# How do we decide about actions, states, rewards?

o We model them as functions

o The **policy function** $a_t = \pi(s_t)$ selects an action given the current state

o The **value function** $Q^\pi(s_t, a_t)$ is the expected total reward that we will receive if we take action $a_t$ given state $s_t$

o What should our goal then be?

$$Q^\pi(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t, a_t)$$

o Learn to take actions $a_t$ that maximize the value function for different states

# Approaches to Reinforcement Learning

- Policy-based
  - Learn directly the optimal policy $\pi^*$
  - The policy $\pi^*$ obtains the maximum future reward

- Value-based
  - Learn the optimal value function $Q^*(s, a)$
  - This value function applies for any policy

- Model-based
  - Build a model for the environment
  - Plan and decide using that model

- Pros and cons?

# Bellman equation

o How can we rewrite the value function in more compact form

$$Q^\pi(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t, a_t) = ?$$

# Bellman equation

o How can we rewrite the value function in more compact form

$$Q^\pi(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t, a_t)$$
$$= \mathbb{E}_{s'}(r + \gamma Q^\pi(s', a') | s_t, a_t)$$

o This is the *Bellman equation*

o How can we rewrite the optimal value function $Q^*(s_t, a_t)$?

# Bellman equation

o How can we rewrite the value function in more compact form

$$Q^\pi(s_t, a_t) = \mathbb{E}(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \cdots | s_t, a_t)$$
$$= \mathbb{E}_{s'}(r + \gamma Q^\pi(s', a') | s_t, a_t)$$

o This is the **Bellman equation**

o How can we rewrite the optimal value function $Q^*(s, a)$?

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

# Q-Learning

o In the simplest case the value function $Q(s, a)$ is a table

o In the beginning of the learning the function $Q(s, a)$ is incorrect

o Still, to the limit value iteration algorithms solve the Bellman equation

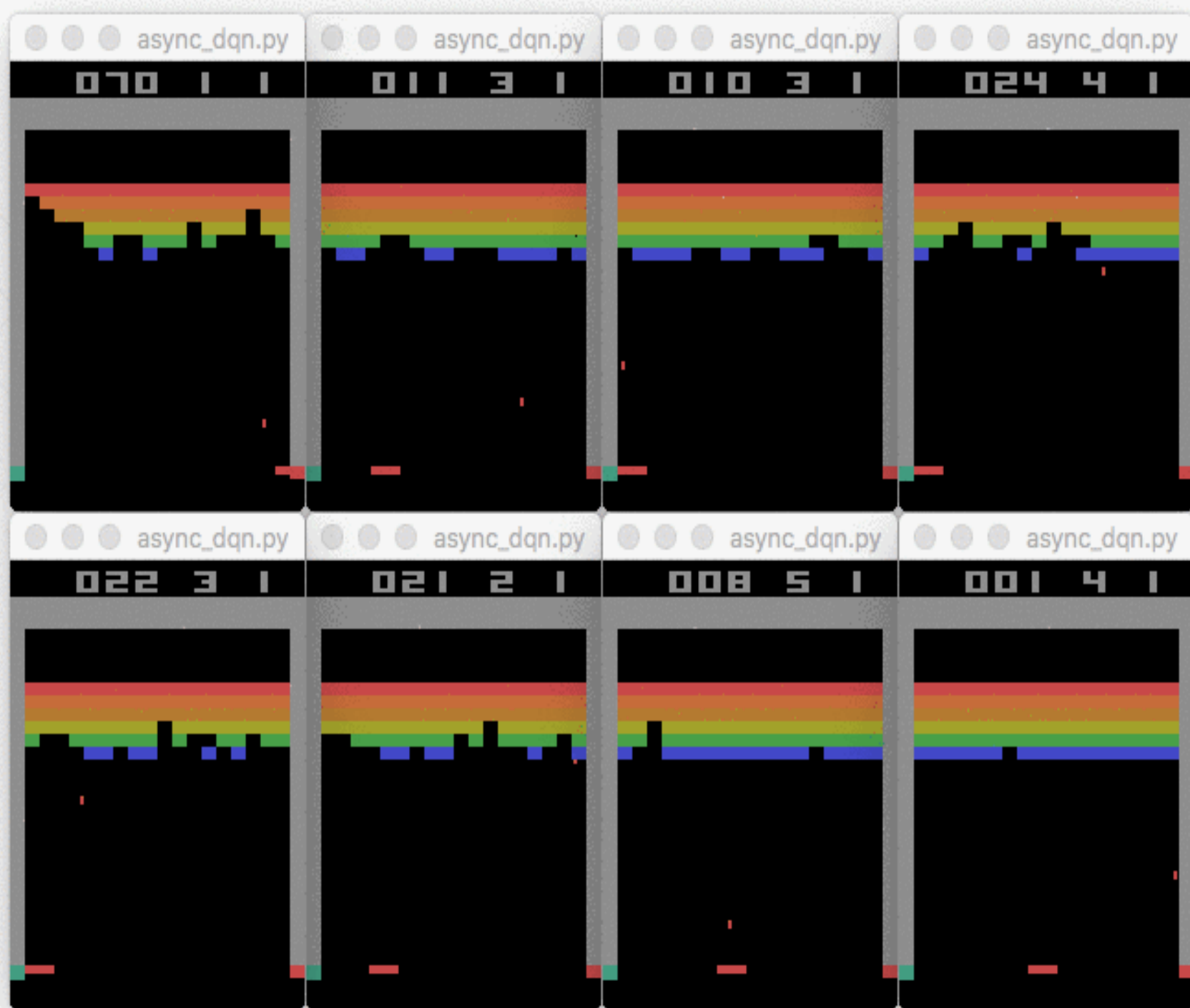$$Q_{i+1}(s, a) = \mathbb{E}_{s'}\left(r + \gamma \max_{a'} Q_i(s', a') \Big| s, a\right)$$

# Policy Optimization

o Computing the $Q$-value is often too expensive

  ◦ Hard to solve $\arg\max_{\alpha} Q_\theta(s, a)$

  ◦ Especially when having continuous or high-dimensional action spaces

o Often defining the policy $\pi_\theta(u|s)$ is easier than defining the $Q$-function

o Use a non-linear function approximator to model the action value function
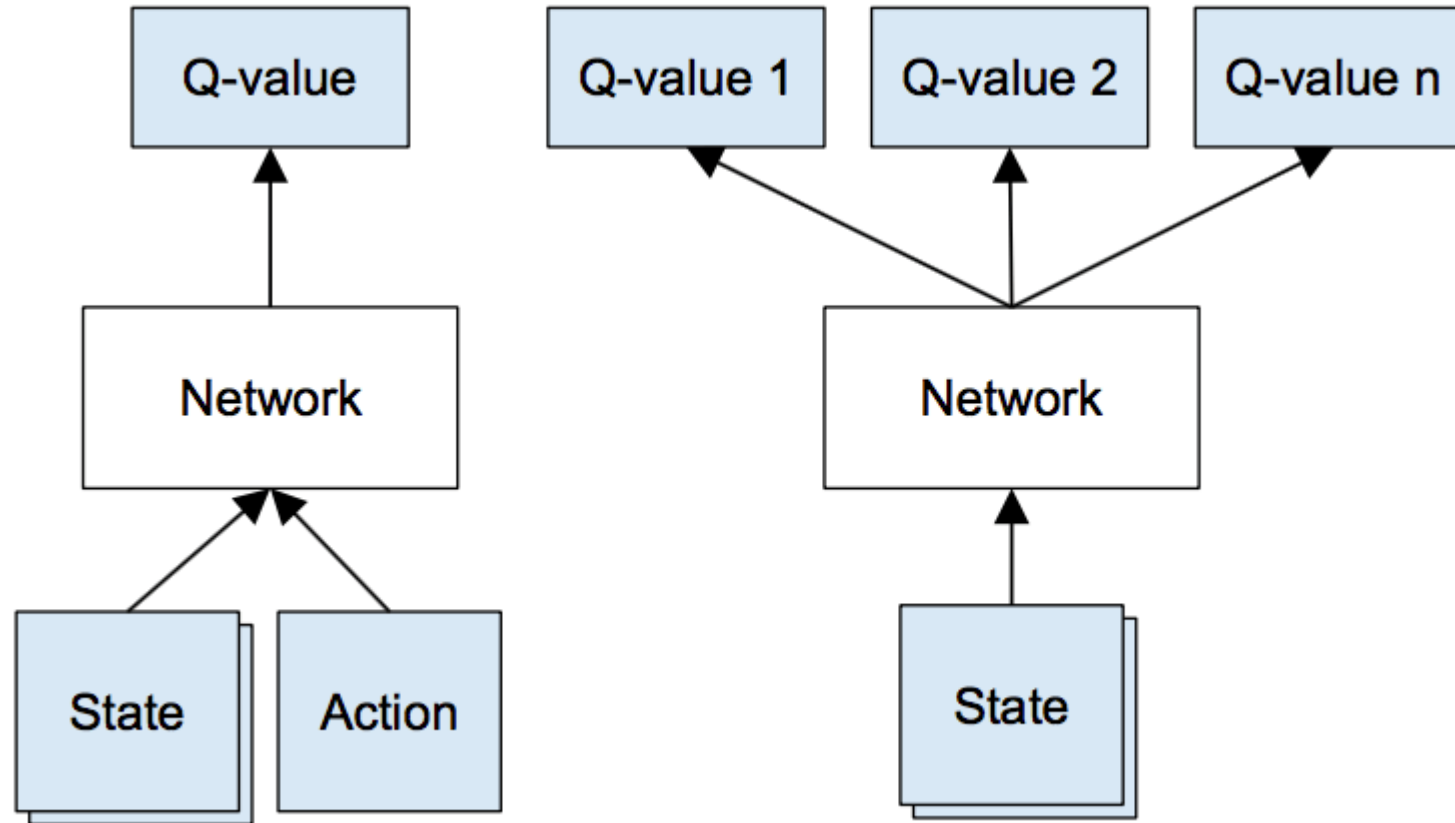
$$Q^*(s, a) \approx Q(s, a; \theta)$$

o Our deep network can be such a non-linear function approximator optimize the $\pi_\theta(u|s)$

# Deep Reinforcement Learning

# How to make RL deep?

# How to make RL deep?

# Deep Reinforcement Learning

o Rewards $r_t$ at time $t$

o Actions $\pi$ taken according to a policy $\pi = \mathrm{P}(a|s)$
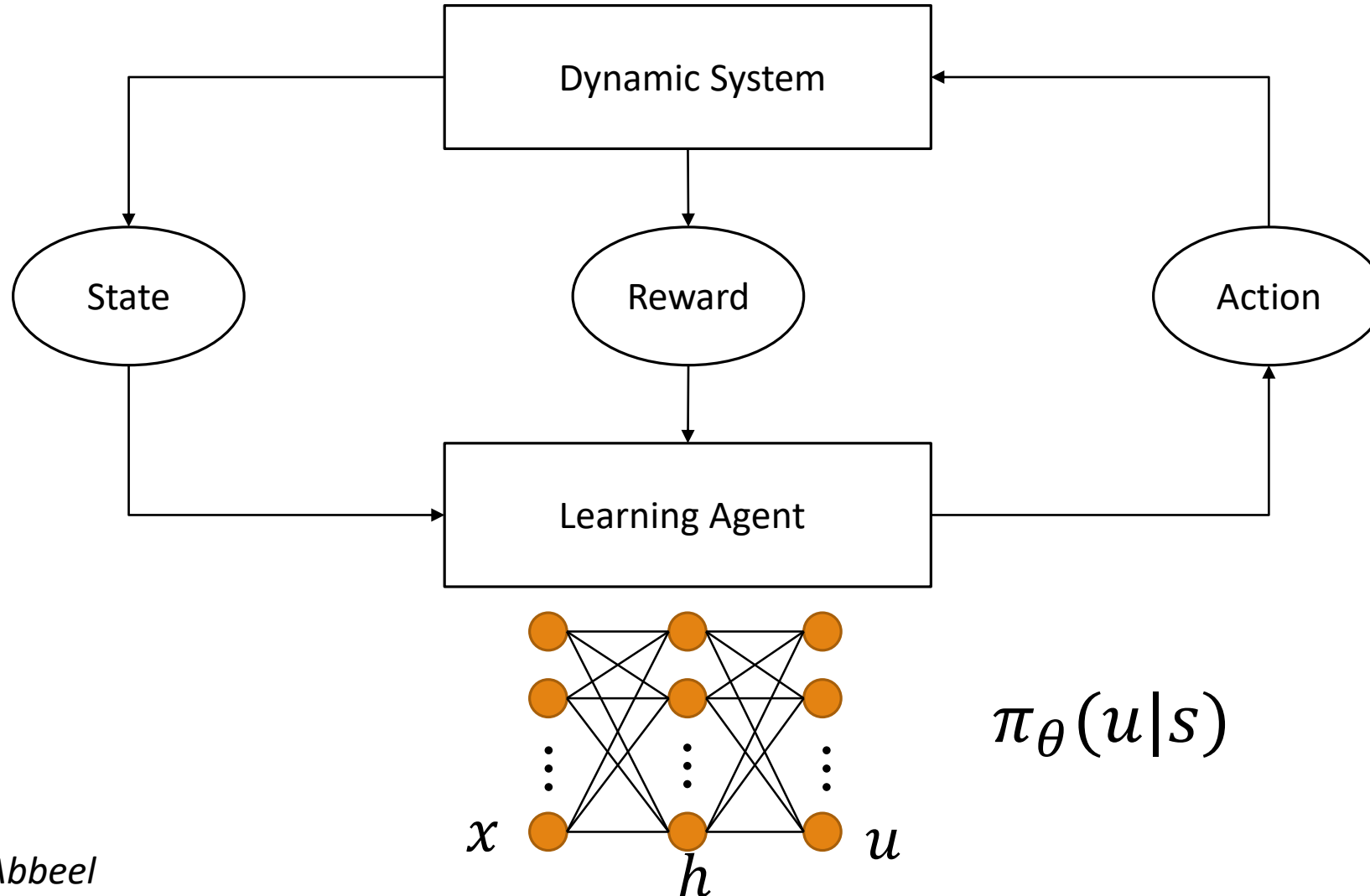
o Again, the action-value function

$$Q(s,a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots | s_t = s, a_t = a, \pi]$$

where $\gamma$ is a discount factor of the future rewards

  ◦ Future rewards should not be as important, because we do not know the future

o Use a non-linear function approximator to model the action value function
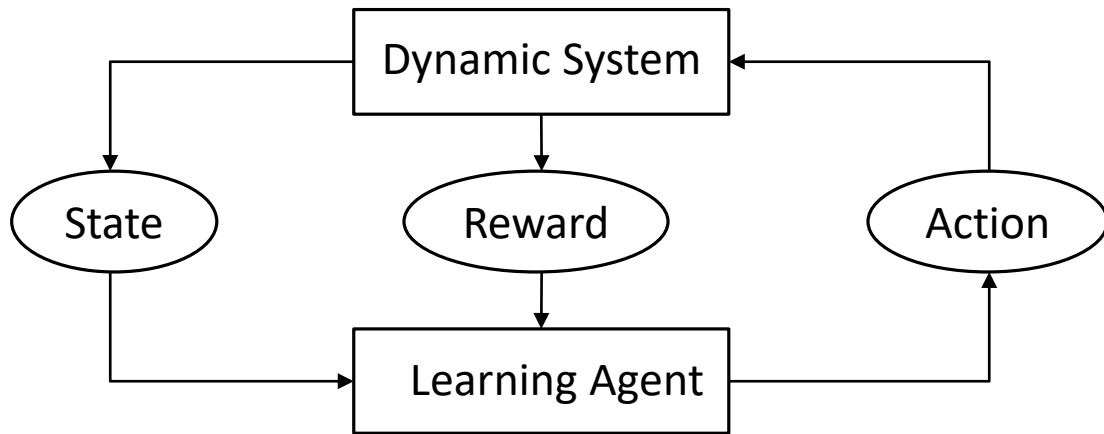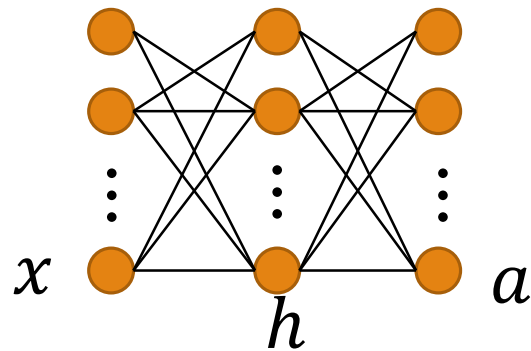
$$Q^*(s,a) \approx Q(s,a;\theta)$$

# Policy optimization



$$\pi_\theta(u|s)$$

*Slides inspired by P. Abbeel*

# Policy optimization

Dynamic System

State    Reward    Action

Learning Agent

$\pi_\theta(a|s)$

$x$    $h$    $a$

o Train learning agent for the optimal policy $\pi_\theta(a|s)$ given states $s$ and possible actions $a$

o The policy class can be either deterministic or stochastic

*Slides inspired by P. Abbeel*

# Deep Reinforcement Learning

o Non-linear function approximator: Deep Networks

o Input is as raw as possible, e.g. image frame
  ◦ Or perhaps several frames (When needed?)

o Output is the best possible action out of a set of actions for maximizing future reward

o **Important:** no need anymore to compute the actual value of the action-value function and take the maximum: $\arg \max_{\alpha} Q_{\theta}(s, a)$
  ◦ The network returns directly the optimal action

# How to optimize?

○ The objective is the mean squared-error in Q-values

$$\mathcal{L}(\theta) = \mathbb{E}[(\underbrace{r + \gamma \max_{a'} Q(s', a', \theta)}_{\text{target}} - Q(s, a, \theta))^2]$$

○ The Q-Learning gradient then becomes

$$\frac{\partial \mathcal{L}}{\partial \theta} = \mathbb{E}[(r + \gamma \max_{a'} Q(s', a', \theta) - Q(s, a, \theta)) \frac{\partial Q(s, a, \theta)}{\partial \theta}]$$

○ Optimize end-to-end with SGD

# In practice

1. Do a feedforward pass for the current state $s$ to get predicted Q-values for all actions

2. Do a feedforward pass for the next state $s'$ and calculate maximum overall network outputs $\max\limits_{a'} Q(s', a', \theta)$

3. Set Q-value target for action to $r + \gamma \max\limits_{a'} Q(s', a', \theta)$ (use the max calculated in step 2). For all other actions, set the Q-value target to the same as originally returned from step 1, <span style="color:red">making the error 0</span> for those outputs

4. Update the weights using backpropagation.

# Stability in Deep Reinforcement Learning

# Stability problems

- Naively, Q-Learning oscillates or diverges with neural networks

- Why?

# Stability problems

o Naively, Q-Learning oscillates or diverges with neural networks

o Why?

o Sequential data breaks i.i.d. assumption
  ◦ Highly correlated samples break SGD

o However, this is not specific to RL, as we have seen earlier

# Stability problems

- Naively, Q-Learning oscillates or diverges with neural networks

- Why?

# Stability problems

o The learning objective is

$$\mathcal{L}(\theta) = \mathbb{E}\left[\left(\textcolor{red}{r + \gamma \max_{a'} Q(s', a', \theta)} - \textcolor{green}{Q(s, a, \theta)}\right)^2\right]$$

o The target depends on the $Q$ function also. This means that if we update the current $Q$ function with backprop, the target will also change

o Plus, we know neural networks are highly non-convex

o Policy changes will change fast even with slight changes in the $Q$ function
  ◦ Policy might oscillate
  ◦ Distribution of data might move from one extreme to another

# Stability problems

- Naively, Q-Learning oscillates or diverges with neural networks
- Why?

# Stability problems

o Not easy to control the scale of the $Q$ values$\rightarrow$ gradients are unstable $Q$

o Remember, the $Q$ function is the output of a neural network

o There is no guarantee that the outputs will lie in a certain range
- Unless care is taken

o Naïve $Q$ gradients can be too large, or too small $\rightarrow$ generally unstable and unreliable

o Where else did we observe a similar behavior?

# Improving stability: Experience replay

o Replay memory/Experience replay

o Store memories $< s, a, r, s' >$

o Train using random stored memories instead of the latest memory transition

o Breaks the temporal dependencies – SGD works well if samples are roughly independent

o Learn from all past policies

# Experience replay

o Take action $a_t$ according to $\varepsilon$-greedy policy

o Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory $D$

o Sample random mini-batch of transitions $(s, a, r, s')$ from $D$

o Optimize mean squared error using the mini-batch

$$\mathcal{L}(\theta) = \mathbb{E}_{(s,a,r,s') \sim D}\left[\left(r + \gamma \max_{a'} Q(s', a', \theta) - Q(s, a, \theta)\right)^2\right]$$

o Effectively, update your network using random past inputs (experience), not the ones the agent currently sees

# Improving stability: Freeze target $Q$ network

o Instead of having "moving" targets, have two networks

 ◦ One Q-Learning and one Q-Target networks

o Copy the $Q$ network parameters to the target network every $K$ iterations

 ◦ Otherwise, keep the old parameters between iterations

 ◦ The targets come from another (Q-Target) network with slightly older parameters

o Optimize the mean squared error as before, only now the targets are defined by the "older" $Q$ function

$$\mathcal{L}(\theta) = \mathbb{E}\left[\left({\color{red}r + \gamma \max_{a'} Q(s', a', \theta_{old})} - {\color{green}Q(s, a, \theta)}\right)^2\right]$$

o Avoids oscillations

# Improving stability: Take care of rewards

o Clip rewards to be in the range $[-1, +1]$

o Or normalize them to lie in a certain, stable range

o Can't tell the difference between large and small rewards

# Results

| | Q-learning | Q-learning + Target Q | Q-learning + Replay | Q-learning + Replay + Target Q |
|---|---|---|---|---|
| Breakout | 3 | 10 | 241 | **317** |
| Enduro | 29 | 142 | 831 | **1006** |
| River Raid | 1453 | 2868 | 4103 | **7447** |
| Seaquest | 276 | 1003 | 823 | **2894** |
| Space Invaders | 302 | 373 | 826 | **1089** |

# Some extra tricks

o Skipping frames
  - Saves time and computation
  - Anyways, from one frame to the other there is often very little difference

o $\varepsilon$-greedy behavioral policy with annealed temperature during training
  - Select random action (instead of optimal) with probability $\varepsilon$
  - In the beginning of training our model is bad, no reason to trust the "optimal" action

o Alternatively: Exploration vs exploitation
  - early stages ≡ strong exploration
  - late stages ≡ strong exploitation

# Examples of Deep Reinforcement Learning

# Deep Reinforcement Learning in Atari
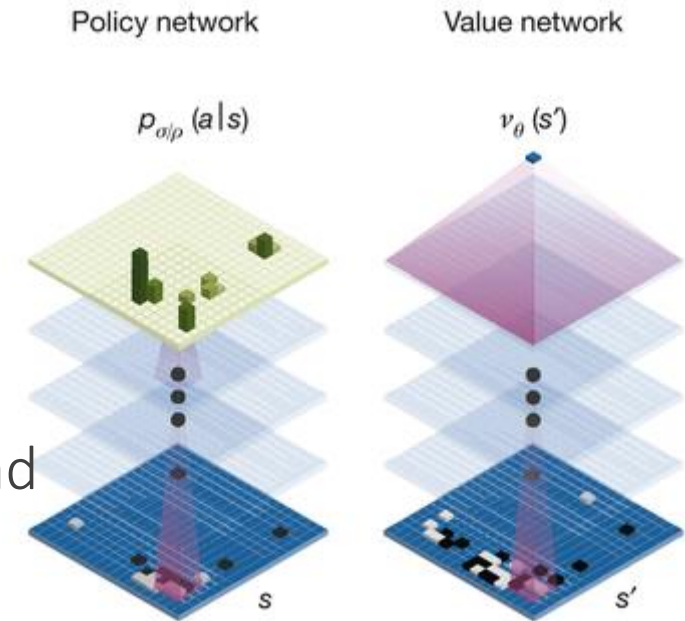
# AlphaGo

○ At least $10^{10^{48}}$ possible game states

  ◦ Chess has $10^{120}$

○ Monte Carlo Tree Search used mostly

  ◦ Start with random moves and evaluate how often they lead to victory

  ◦ Learn the value function to predict the quality of a move

  ◦ Exploration-exploitation trade-off
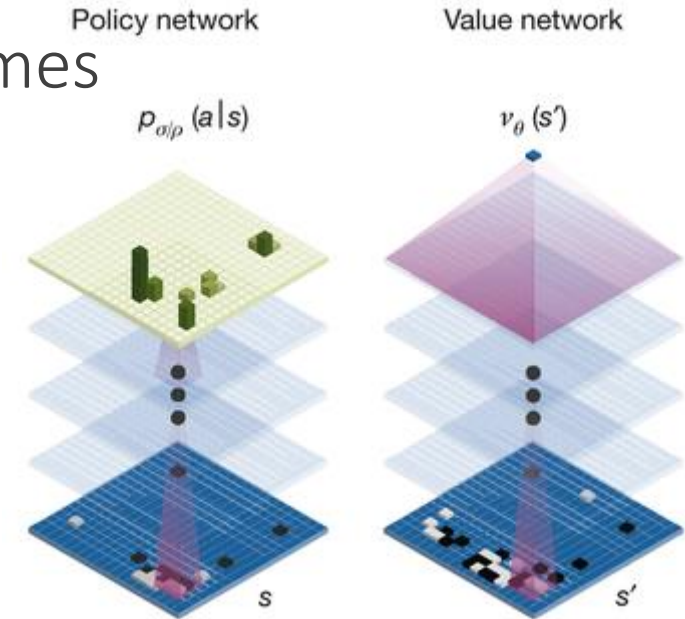
Tic-Tac-Toe possible game states

# AlphaGo

o AlphaGo relies on a tree procedure for search

o AlphaGo relies on ConvNets to guide the tree search

o A ConvNet trained to predict human moves achieved 57% accuracy

  ◦ Humans make intuitive moves instead of thinking too far ahead

o For Deep RL we don't want to predict human moves

  ◦ Instead, we want the agent to learn the optimal moves

o Two policy networks (one per side) + One value network

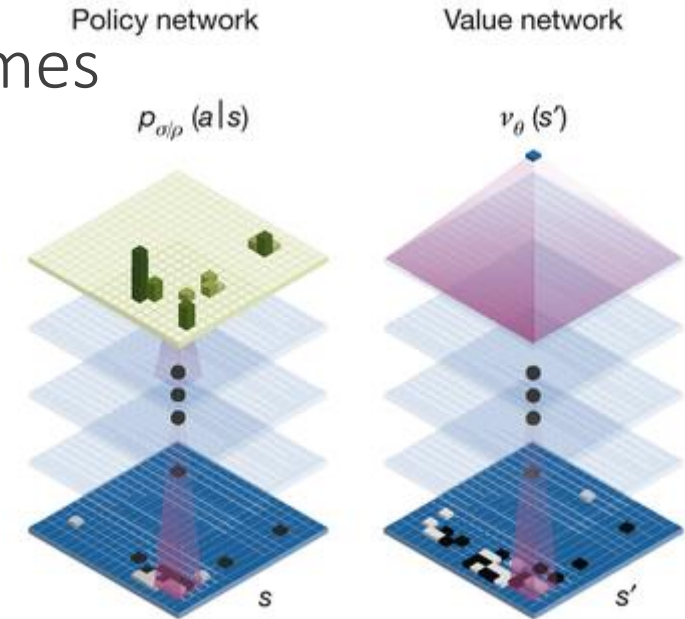o Value network trained on 30 million positions while policy networks play



Policy network

$p_{\sigma/\rho}(a|s)$

Value network

$v_\theta(s')$

# AlphaGo

o Both humans and Deep RL agents play better end games
  ◦ Maybe a fundamental cause?

o In the end the value of a state is computed equally from Monte Carlo simulation and the value network output
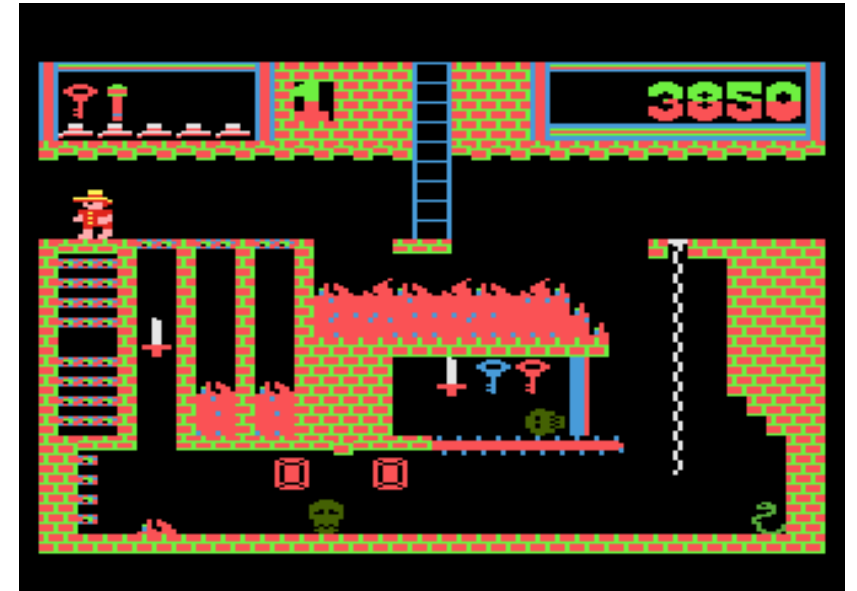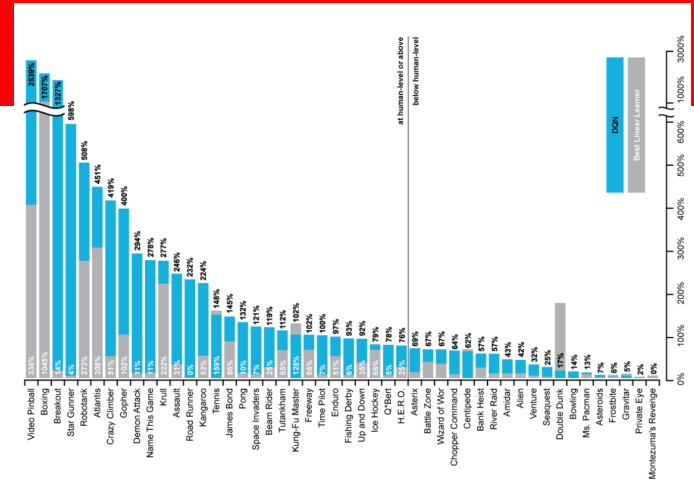  ◦ Combining intuitive play and thinking ahead

o Where is the catch?

Policy network

$p_{\sigma/\rho}(a|s)$

Value network

$v_\theta(s')$

$s$

$s'$

# AlphaGo

o Both humans and Deep RL agents play better end games
  ◦ Maybe a fundamental cause?

o In the end the value of a state is computed equally from Monte Carlo simulation and the value network output
  ◦ Combining intuitive play and thinking ahead

o Where is the catch?

o State is not the pixels but positions

o Also, the game states and actions are highly discrete



Policy network

$p_{\sigma/\rho}(a|s)$
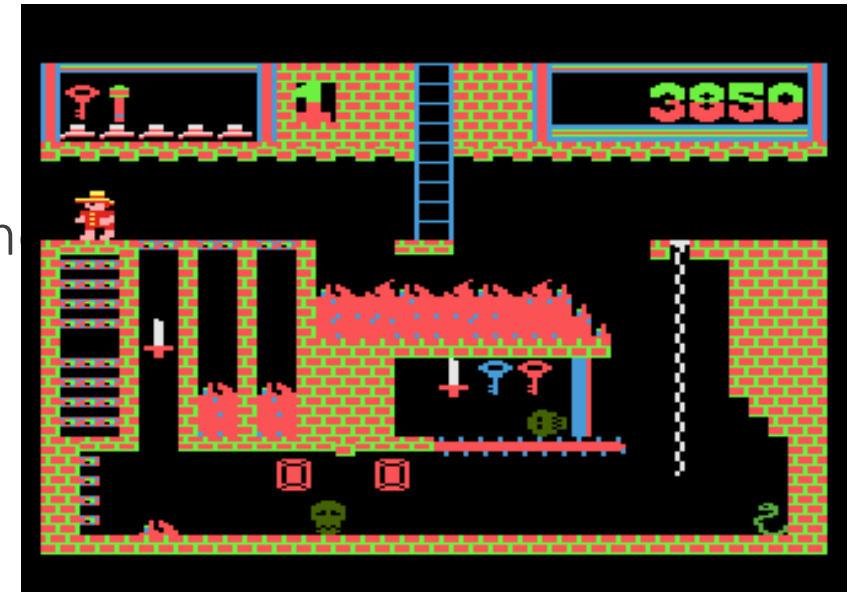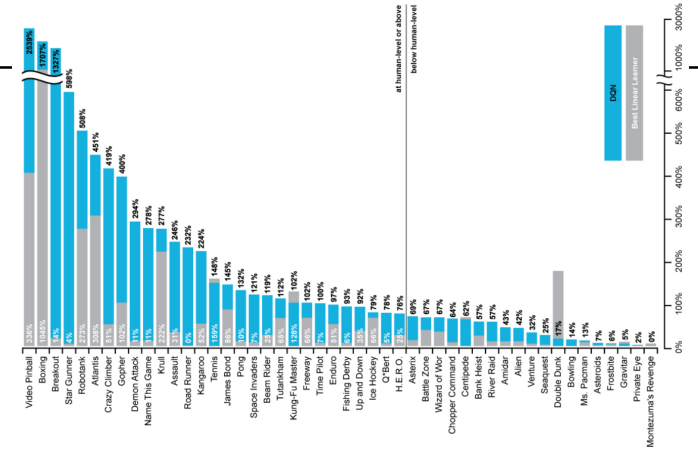
Value network

$v_\theta(s')$

# Montezuma's Revenge

- Hardest of the Atari games

- Why is it so hard?

# Montezuma's Revenge



- Hardest of the Atari games

- Why is it so hard?

- Very long-term dependencies
  - Must search in multiple rooms to find the "secret"

- Future rewards are too delayed
  - It takes a while to evaluate an action was good or n
  - Hard to optimize

- FeUdal Networks for Hierarchical Reinforcement Learning, A. Vezhnevets

# Starcraft II

o Dr. O. Vinyals in DeepMind tries to bring Starcraft II and Deep RL together
   ◦ There is a Machine Learning API from Blizzard
   ◦ There is a dataset of anonymized game replays
   ◦ An open source python toolkit from DeepMind
   ◦ And there is a paper

o More info
   ◦ https://deepmind.com/blog/deepmind-and-blizzard-open-starcraft-ii-ai-research-environment/
   ◦ https://www.youtube.com/watch?v=-fKUyT14G-8

o What are the possible difficulties?

# Summary

o Reinforcement Learning

o Q-Learning

o Deep Q-Learning

o Making Deep Q-Learning stable

o Examples of Deep Q-Learning