
Playing Atari with Deep Reinforcement Learning

**Volodymyr Mnih Koray Kavukcuoglu David Silver Alex Graves Ioannis Antonoglou
Daan Wierstra Martin Riedmiller**

DeepMind Technologies



Presented by Sung Kim <hunkim@gmail.com>

Deep Learning Workshop

NIPS 2013

[Call For Papers](#)[Submission](#)[Invited Speakers](#)[Accepted Papers](#)[Schedule](#)[Program Committee](#)

Accepted Papers

Oral presentations

- [Deep Supervised and Convolutional Generative Stochastic Network for Protein Secondary Structure Prediction](#)
Jian Zhou, Olga Troyanskaya
- **Playing Atari with Deep Reinforcement Learning**
Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, Martin Riedmiller

Poster presentations

Abstract

- First deep learning model using reinforcement learning
 - Successfully learn control policies directly
 - From high-dimensional sensory input (pixels)
- CNN model trained on a variant of Q-learning
 - Input: raw pixel, output: a value function estimating future reward
- Applied seven Atari 2600 games (no adjustment)
 - Outperforms ALL previous approaches on six games
 - Surpasses a human expert on three games

Problem to solve: motivation

- Learning directly from high dimensional sensory input
 - Long-standing challenges of RL
- Most successful RL relies on hand-crafted features

Insights: NN for RL?

- Recent advances in deep learning extract high-level features from raw sensory data
 - Breakthroughs in computer vision
 - Speech recognition
- Neural architectures on supervised/unsupervised learning
 - Convolutional networks
 - Multilayer perceptrons
 - Restricted Boltzmann matches
 - Recurrent Neural nets

Remaining challenges: cannot directly apply DL to RL

- Most DL requires hand labeled training data
 - RL must learn from a scalar reward signal
 - Reward signal is often sparse, noisy, and delayed
 - delay between actions and resulting rewards can be thousand time steps
- Most DL assumes data samples are independent
 - RL encounters sequences of highly correlated states

Solutions

- Most DL requires hand labeled training data
 - RL must learn from a scalar reward signal
 - Reward signal is often sparse, noisy, and delayed
 - delay between actions and resulting rewards can be thousand time steps
 - **CNN with a variant Q-learning**
- Most DL assumes data samples are independent
 - RL encounters sequences of highly correlated states
 - **Experience replay**

2 Background

We consider tasks in which an agent interacts with an environment \mathcal{E} , in this case the Atari emulator, in a sequence of actions, observations and rewards. At each time-step the agent selects an action a_t from the set of legal game actions, $\mathcal{A} = \{1, \dots, K\}$. The action is passed to the emulator and modifies its internal state and the game score. In general \mathcal{E} may be stochastic. The emulator's internal state is not observed by the agent; instead it observes an image $x_t \in \mathbb{R}^d$ from the emulator, which is a vector of raw pixel values representing the current screen. In addition it receives a reward r_t representing the change in game score. Note that in general the game score may depend on the whole prior sequence of actions and observations; feedback about an action may only be received after many thousands of time-steps have elapsed.

Since the agent only observes images of the current screen, the task is partially observed and many emulator states are perceptually aliased, i.e. it is impossible to fully understand the current situation from only the current screen x_t . We therefore consider sequences of actions and observations, $s_t = x_1, a_1, x_2, \dots, a_{t-1}, x_t$, and learn game strategies that depend upon these sequences. All sequences in the emulator are assumed to terminate in a finite number of time-steps. This formalism gives rise to a large but finite Markov decision process (MDP) in which each sequence is a distinct state. As a result, we can apply standard reinforcement learning methods for MDPs, simply by using the complete sequence s_t as the state representation at time t .

The goal of the agent is to interact with the emulator by selecting actions in a way that maximises future rewards. We make the standard assumption that future rewards are discounted by a factor of γ per time-step, and define the future discounted *return* at time t as $R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$, where T is the time-step at which the game terminates. We define the optimal action-value function $Q^*(s, a)$ as the maximum expected return achievable by following any strategy, after seeing some sequence s and then taking some action a , $Q^*(s, a) = \max_{\pi} \mathbb{E}[R_t | s_t = s, a_t = a, \pi]$, where π is a policy mapping sequences to actions (or distributions over actions).

The optimal action-value function obeys an important identity known as the *Bellman equation*. This is based on the following intuition: if the optimal value $Q^*(s', a')$ of the sequence s' at the next time-step was known for all possible actions a' , then the optimal strategy is to select the action a' maximising the expected value of $r + \gamma Q^*(s', a')$,

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') \middle| s, a \right] \quad (1)$$

The basic idea behind many reinforcement learning algorithms is to estimate the action-value function, by using the Bellman equation as an iterative update, $Q_{i+1}(s, a) = \mathbb{E}[r + \gamma \max_{a'} Q_i(s', a') | s, a]$. Such *value iteration* algorithms converge to the optimal action-value function, $Q_i \rightarrow Q^*$ as $i \rightarrow \infty$ [23]. In practice, this basic approach is totally impractical, because the action-value function is estimated separately for each sequence, without any generalisation. Instead, it is common to use a function approximator to estimate the action-value function, $Q(s, a; \theta) \approx Q^*(s, a)$. In the reinforcement learning community this is typically a linear function approximator, but sometimes a non-linear function approximator is used instead, such as a neural network. We refer to a neural network function approximator with weights θ as a Q-network. A Q-network can be trained by minimising a sequence of loss functions $L_i(\theta_i)$ that changes at each iteration i ,

$$L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot)} \left[(y_i - Q(s, a; \theta_i))^2 \right], \quad (2)$$

where $y_i = \mathbb{E}_{s' \sim \mathcal{E}} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ is the target for iteration i and $\rho(s, a)$ is a probability distribution over sequences s and actions a that we refer to as the *behaviour distribution*. The parameters from the previous iteration θ_{i-1} are held fixed when optimising the loss function $L_i(\theta_i)$. Note that the targets depend on the network weights; this is in contrast with the targets used for supervised learning, which are fixed before learning begins. Differentiating the loss function with respect to the weights we arrive at the following gradient,

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot), s' \sim \mathcal{E}} \left[\left(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i) \right) \nabla_{\theta_i} Q(s, a; \theta_i) \right]. \quad (3)$$

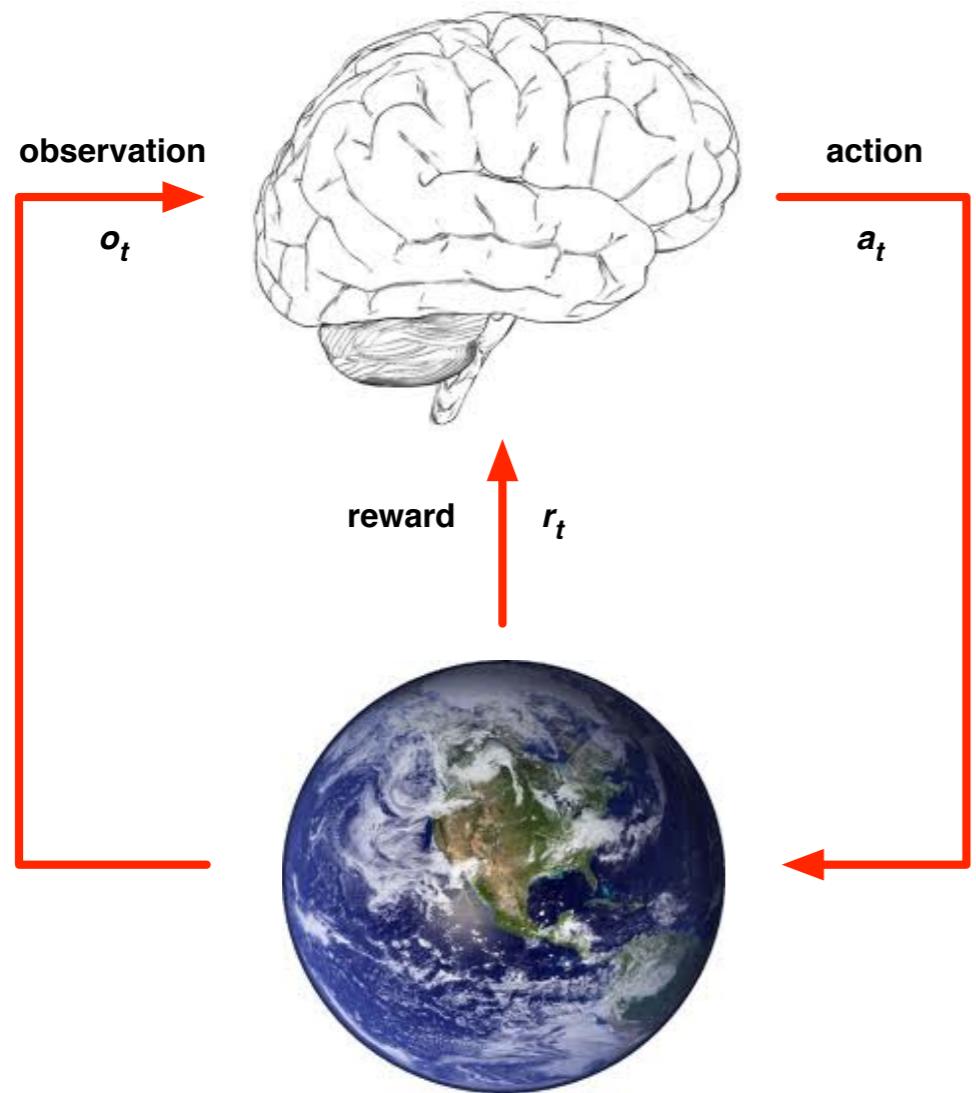
Rather than computing the full expectations in the above gradient, it is often computationally expedient to optimise the loss function by stochastic gradient descent. If the weights are updated after every time-step, and the expectations are replaced by single samples from the behaviour distribution ρ and the emulator \mathcal{E} respectively, then we arrive at the familiar *Q-learning* algorithm [26].

Note that this algorithm is *model-free*: it solves the reinforcement learning task directly using samples from the emulator \mathcal{E} , without explicitly constructing an estimate of \mathcal{E} . It is also *off-policy*: it learns about the greedy strategy $a = \max_a Q(s, a; \theta)$, while following a behaviour distribution that ensures adequate exploration of the state space. In practice, the behaviour distribution is often selected by an ϵ -greedy strategy that follows the greedy strategy with probability $1 - \epsilon$ and selects a random action with probability ϵ .

Tutorial: Deep Reinforcement Learning

David Silver, Google DeepMind

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 behaviour 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

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

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 Functions

- ▶ An optimal value function is the maximum achievable value

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

Optimal Value Functions

- ▶ 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 Functions

- ▶ 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:

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

Optimal Value Functions

- ▶ 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

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

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

Outline

Introduction to Deep Learning

Introduction to Reinforcement Learning

Value-Based Deep RL

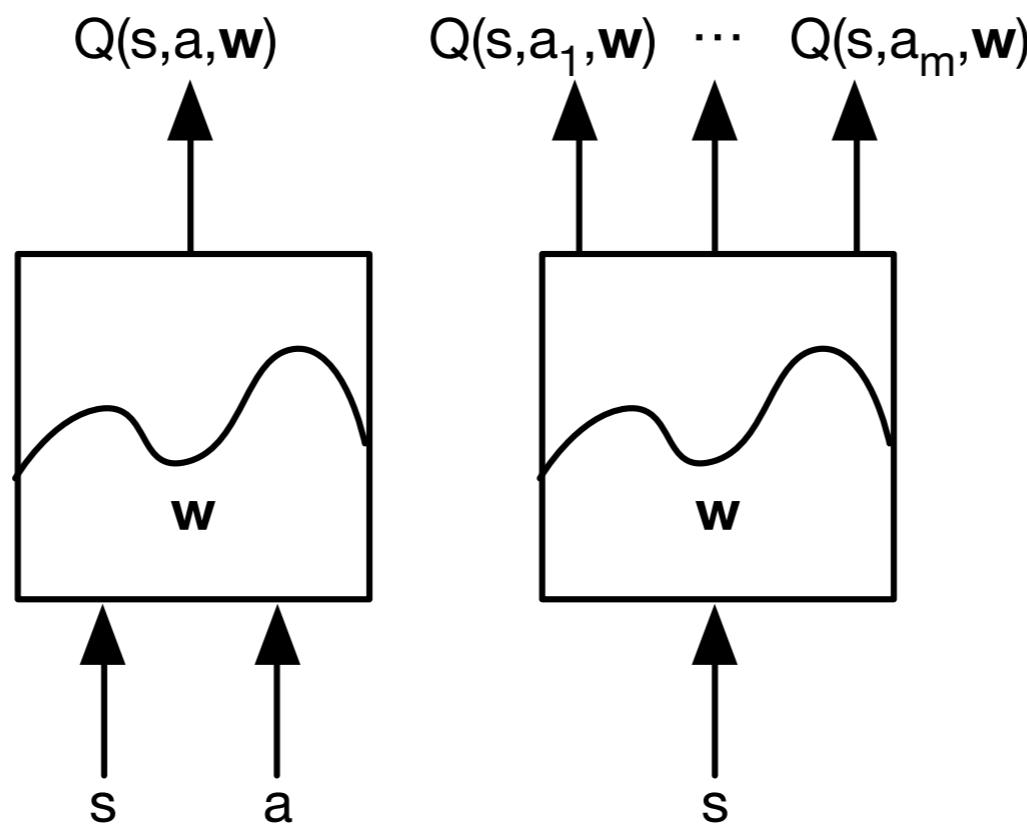
Policy-Based Deep RL

Model-Based Deep RL

Q-Networks

Represent value function by **Q-network** with weights \mathbf{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

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

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

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

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

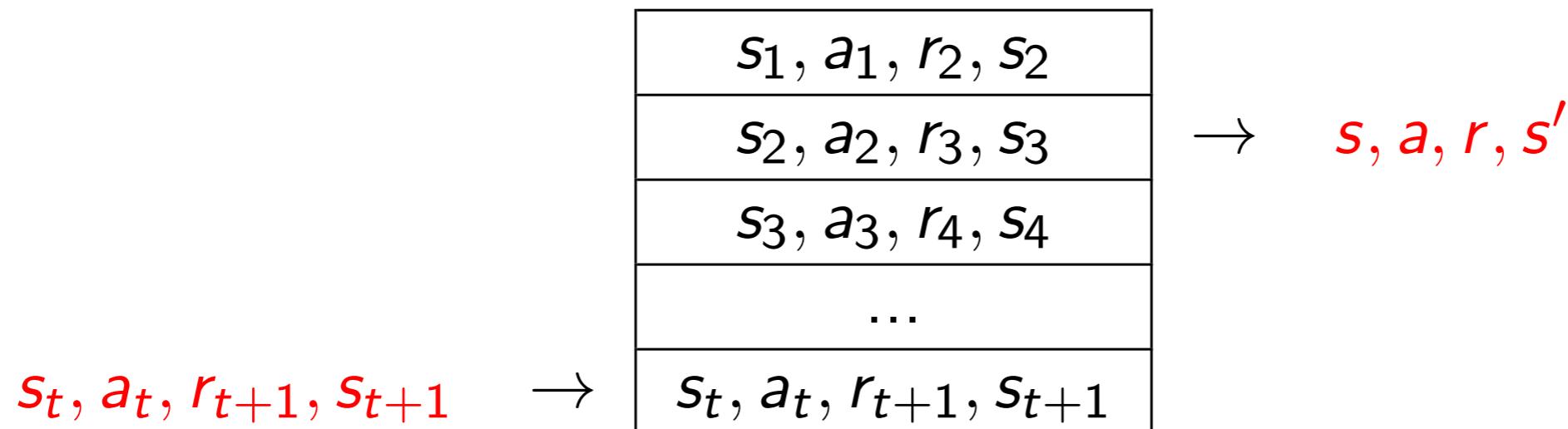
- ▶ Converges to Q^* using table lookup representation
- ▶ But **diverges** using neural networks due to:
 - ▶ Correlations between samples
 - ▶ Non-stationary targets

Challenges/Solutions

- Most DL requires hand labeled training data
 - RL must learn from a scalar reward signal
 - Reward signal is often sparse, noisy, and delayed
 - delay between actions and resulting rewards can be thousand time steps
 - **CNN with a variant Q-learning**
- Most DL assumes data samples are independent
 - RL encounters sequences of highly correlated states
 - **Experience replay**

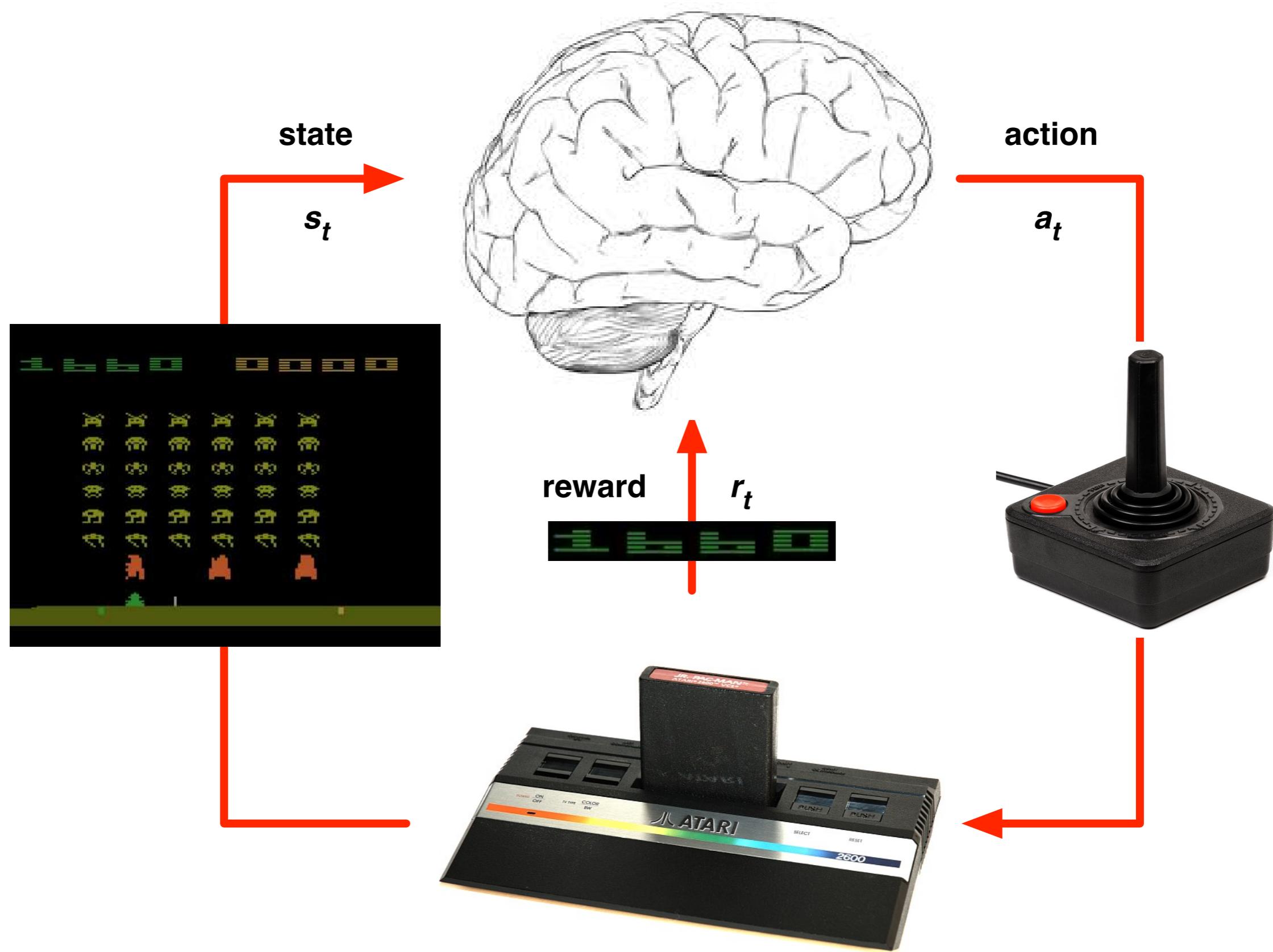
Deep Q-Networks (DQN): Experience Replay

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



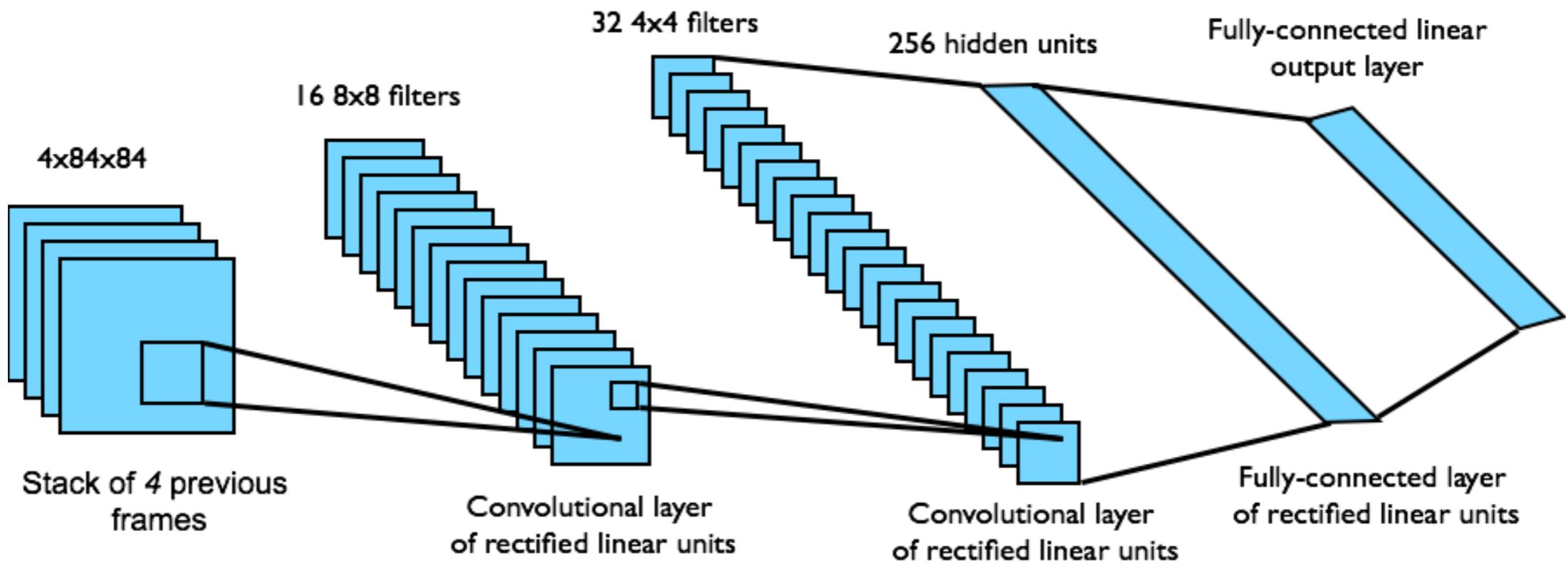
Sample experiences from data-set and apply update

Deep Reinforcement Learning in Atari



DQN in Atari

- ▶ End-to-end learning of values $Q(s, a)$ from pixels s
- ▶ Input state s is stack of raw pixels from last 4 frames
- ▶ Output is $Q(s, a)$ for 18 joystick/button positions
- ▶ Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

Results

	B. Rider	Breakout	Enduro	Pong	Q*bert	Seaquest	S. Invaders
Random	354	1.2	0	-20.4	157	110	179
Sarsa [3]	996	5.2	129	-19	614	665	271
Contingency [4]	1743	6	159	-17	960	723	268
DQN	4092	168	470	20	1952	1705	581
Human	7456	31	368	-3	18900	28010	3690
HNeat Best [8]	3616	52	106	19	1800	920	1720
HNeat Pixel [8]	1332	4	91	-16	1325	800	1145
DQN Best	5184	225	661	21	4500	1740	1075

Table 1: The upper table compares average total reward for various learning methods by running an ϵ -greedy policy with $\epsilon = 0.05$ for a fixed number of steps. The lower table reports results of the single best performing episode for HNeat and DQN. HNeat produces deterministic policies that always get the same score while DQN used an ϵ -greedy policy with $\epsilon = 0.05$.

Training and Stability

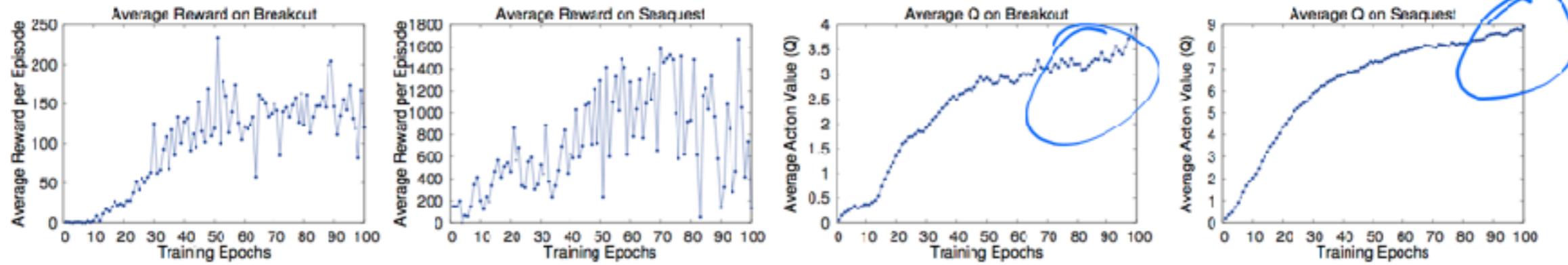


Figure 2: The two plots on the left show average reward per episode on Breakout and Seaquest respectively during training. The statistics were computed by running an ϵ -greedy policy with $\epsilon = 0.05$ for 10000 steps. The two plots on the right show the average maximum predicted action-value of a held out set of states on Breakout and Seaquest respectively. One epoch corresponds to 50000 minibatch weight updates or roughly 30 minutes of training time.

Training and Stability

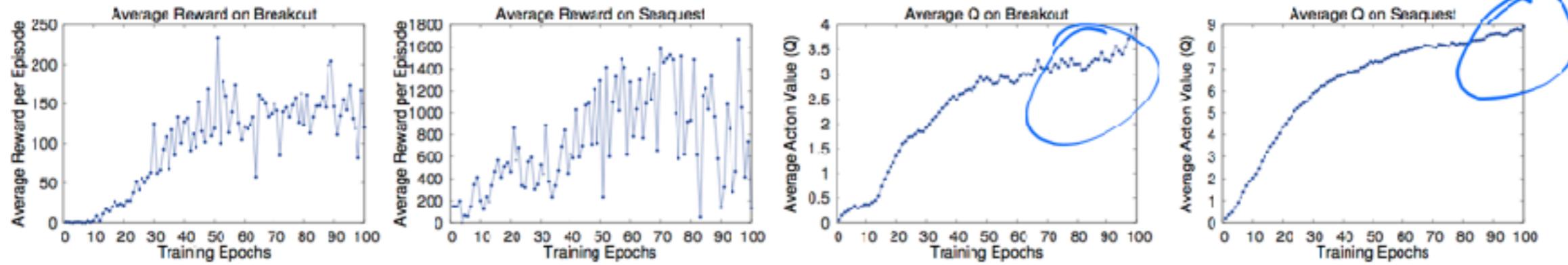


Figure 2: The two plots on the left show average reward per episode on Breakout and Seaquest respectively during training. The statistics were computed by running an ϵ -greedy policy with $\epsilon = 0.05$ for 10000 steps. The two plots on the right show the average maximum predicted action-value of a held out set of states on Breakout and Seaquest respectively. One epoch corresponds to 50000 minibatch weight updates or roughly 30 minutes of training time.

- Average total reward: noisy because small changes to the weights of a policy can lead to large changes in the distributions of states the policy visits.
- Q: smooth improvement despite lacking any theoretical convergence guarantees

Visualizing the value functions

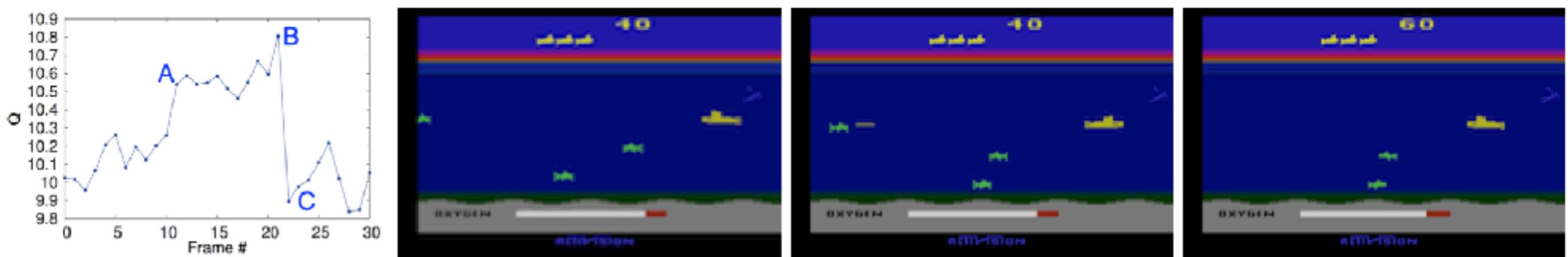
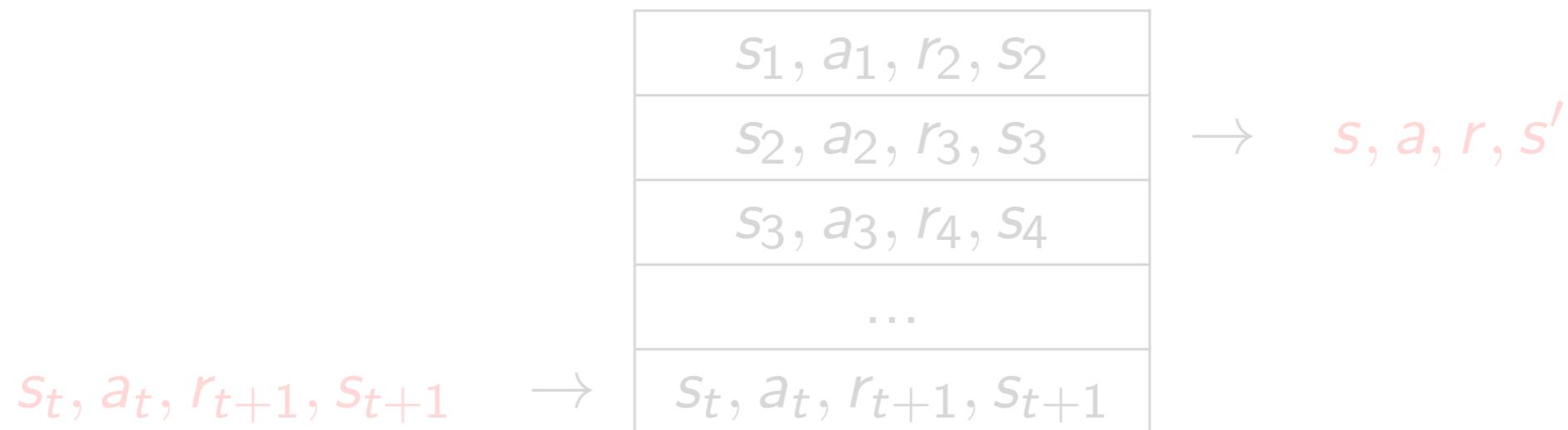


Figure 3: The leftmost plot shows the predicted value function for a 30 frame segment of the game Seaquest. The three screenshots correspond to the frames labeled by A, B, and C respectively.

Deep Q-Networks (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

DQN Atari Demo

DQN paper

www.nature.com/articles/nature14236

DQN source code:

sites.google.com/a/deepmind.com/dqn/



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

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

- ▶ **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

Conclusion

- ▶ General, stable and scalable RL is now possible
- ▶ Using deep networks to represent value, policy, model
- ▶ Successful in Atari, Labyrinth, Physics, Poker, Go
- ▶ Using a variety of deep RL paradigms

Model-free Methods

- Model-free: Learn a controller without learning a model.
- Model-based: Learn a model, and use it to derive a controller.

On Vs Off-policy

- On-policy ($\text{TD}(\lambda)$, SARSA)
 - Start with a simple soft policy
 - Sample state space with this policy
 - Improve policy
- Off-policy (Q-learning, R-learning)
 - Gather information from (partially) random moves
 - Evaluate states as if a greedy policy was used
 - Slowly reduce randomness

On-policy learning vs off-policy learning (according to S&B)

- On-policy methods
 - attempt to evaluate or improve the policy that is used to make decisions
 - often use soft action choice, i.e. $\pi(s, a) > 0 \forall a$
 - commit to always exploring and try to find the best policy that still explores
 - may become trapped in local minima
- Off-policy methods
 - evaluate one policy while following another, e.g. tries to evaluate the greedy policy while following a more exploratory scheme
 - the policy used for behavior should be soft
 - policies may not be sufficiently similar may be slower (only the part after the last exploration is reliable), but remains more flexible if alternative routes appear
 - May lead to the same result (e.g. after greedification)