# Lesson 5:  Deep Q-Network

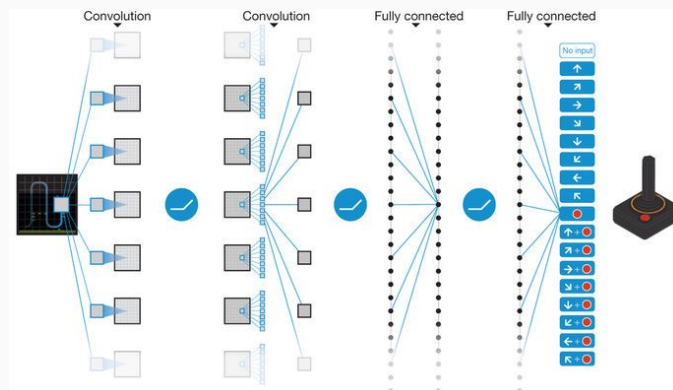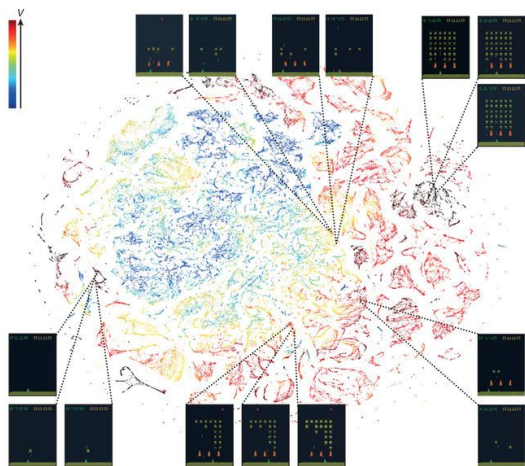# LETTER

## Human-level control through deep reinforcement learning

Volodymyr Mnih[1]*, Koray Kavukcuoglu[1]*, David Silver[1]*, Andrei A. Rusu[1], Joel Veness[1], Marc G. Bellemare[1], Alex Graves[1], Martin Riedmiller[1], Andreas K. Fidjeland[1], Georg Ostrovski[1], Stig Petersen[1], Charles Beattie[1], Amir Sadik[1], Ioannis Antonoglou[1], Helen King[1], Dharshan Kumaran[1], Daan Wierstra[1], Shane Legg[1] & Demis Hassabis[1]

**3 Main innovations:**
1. ConvNet Value approximator
.
.
.

CONVNET VALUE APPROXIMATOR!!??
2. Experience Replay
3. Target Networks

LETTER

## Human-level learning

Volodymyr Mnih[1]*, Koray Kavuko
Martin Riedmiller[1], Andreas K. Fi
Helen King[1], Dharshan Kumaran[1],

**3 Main innovations:**
1. ConvNet Value approximator
.
.
.
CONVNET VALUE APPROXIMATOR!!??
2. Experience Replay
3. Target Networks

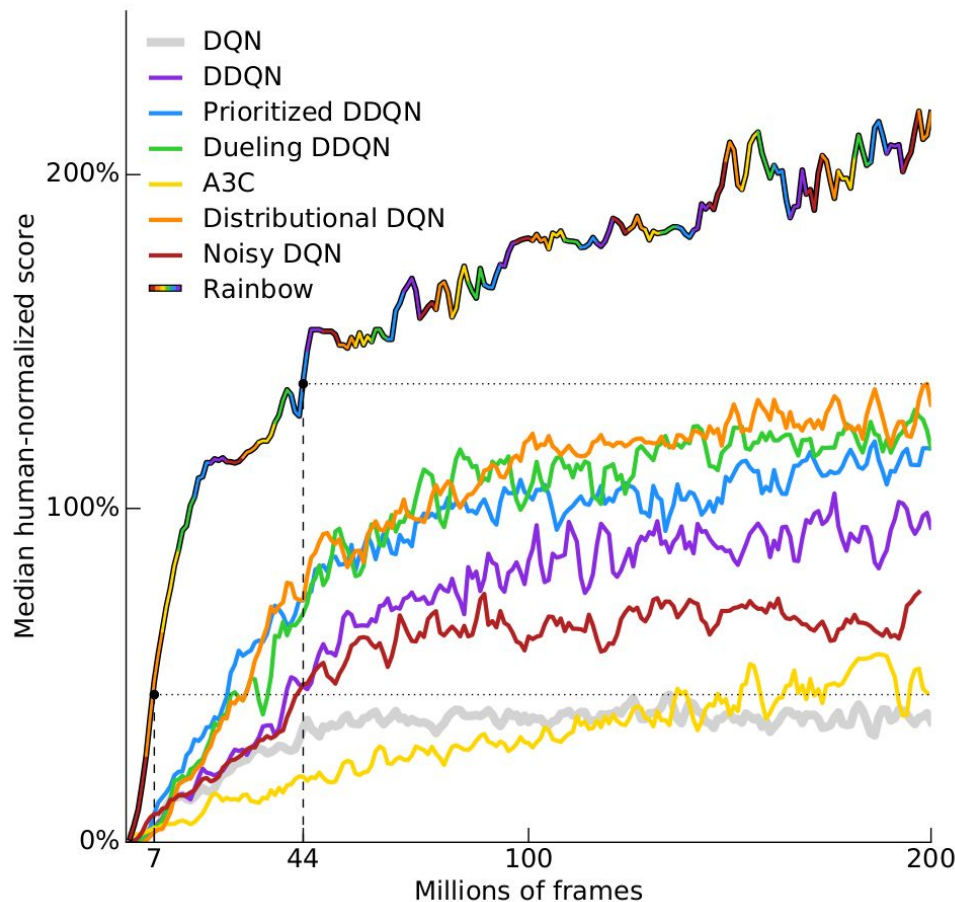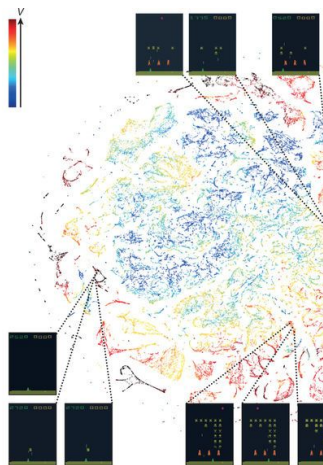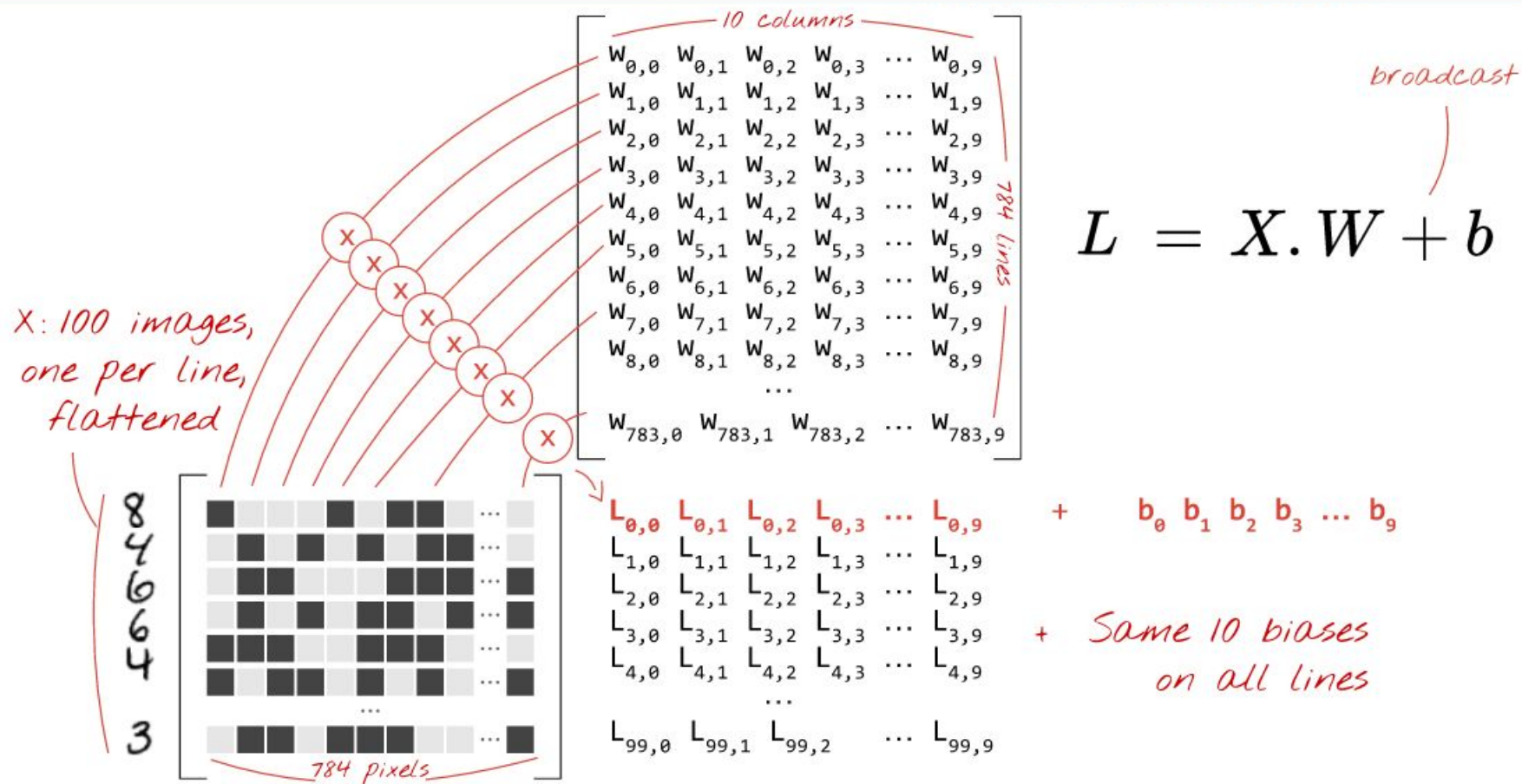$$L = X.W + b$$

broadcast

10 columns

$$\begin{matrix}
W_{0,0} & W_{0,1} & W_{0,2} & W_{0,3} & \cdots & W_{0,9} \\
W_{1,0} & W_{1,1} & W_{1,2} & W_{1,3} & \cdots & W_{1,9} \\
W_{2,0} & W_{2,1} & W_{2,2} & W_{2,3} & \cdots & W_{2,9} \\
W_{3,0} & W_{3,1} & W_{3,2} & W_{3,3} & \cdots & W_{3,9} \\
W_{4,0} & W_{4,1} & W_{4,2} & W_{4,3} & \cdots & W_{4,9} \\
W_{5,0} & W_{5,1} & W_{5,2} & W_{5,3} & \cdots & W_{5,9} \\
W_{6,0} & W_{6,1} & W_{6,2} & W_{6,3} & \cdots & W_{6,9} \\
W_{7,0} & W_{7,1} & W_{7,2} & W_{7,3} & \cdots & W_{7,9} \\
W_{8,0} & W_{8,1} & W_{8,2} & W_{8,3} & \cdots & W_{8,9} \\
& & \cdots & & & \\
W_{783,0} & W_{783,1} & W_{783,2} & & \cdots & W_{783,9}
\end{matrix}$$

784 lines

X: 100 images, one per line, flattened

784 pixels

$$\begin{matrix}
L_{0,0} & L_{0,1} & L_{0,2} & L_{0,3} & \cdots & L_{0,9} \\
L_{1,0} & L_{1,1} & L_{1,2} & L_{1,3} & \cdots & L_{1,9} \\
L_{2,0} & L_{2,1} & L_{2,2} & L_{2,3} & \cdots & L_{2,9} \\
L_{3,0} & L_{3,1} & L_{3,2} & L_{3,3} & \cdots & L_{3,9} \\
L_{4,0} & L_{4,1} & L_{4,2} & L_{4,3} & \cdots & L_{4,9} \\
& & \cdots & & & \\
L_{99,0} & L_{99,1} & L_{99,2} & & \cdots & L_{99,9}
\end{matrix}$$

$+ \quad b_0 \ b_1 \ b_2 \ b_3 \ \cdots \ b_9$

+ Same 10 biases on all lines

**Algorithm 1:** Online-Q

Initialize the value function $q$

**while** *not converged* **do**
    Get the initial state $s$
    **while** *s is not the terminal state* **do**
        Select an action $a$ according to a $\epsilon$-greedy policy derived from $q$
        Execute the action $a$, get the reward $r$ and the next state $s'$
        Update the value function $q$ with $(s, a, r, s')$ following the Q-learning update rule
        $s = s'$
    **end**
**end**

**Algorithm 2:** Buffer-Q

Initialize the value function $q$
Initialize the replay buffer $\mathcal{M}$

**while** *not converged* **do**
    Get the initial state $s$
    **while** *s is not the terminal state* **do**
        Select an action $a$ according to a $\epsilon$-greedy policy derived from $q$
        Execute the action $a$, get the reward $r$ and the next state $s'$
        Store the transition $(s, a, r, s')$ into the replay buffer $\mathcal{M}$
        Sample a batch of transitions $\mathcal{E}$ from $\mathcal{M}$
        Update the value function $q$ with $\mathcal{E}$ following the Q-learning update rule
        $s = s'$
    **end**
**end**

**Key Point:**
We don't sample from just recent experiences but from all.
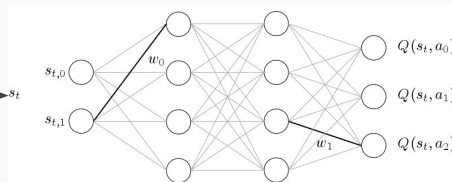
**Batch:**

Q-Update

vs.

**Replay:**

Q-Update

# Target Networks

**1. Exploration/Data Collection**

State Vector
(s(t))

$s_t$



$s_{t,0}$ $w_0$ $Q(s_t, a_0)$
$s_{t,1}$ $Q(s_t, a_1)$
$w_1$ $Q(s_t, a_2)$

Q-Values
(current
state)

e-greedy

max(Q)
or
random

Action
(a)

Update
State

next state, reward
s(t+1)     (r)

Environment
(gym.step(a))

**2. Network Training**

s(t+1)

$s_t$

$s_{t,0}$ $w_0$ $Q(s_t, a_0)$
$s_{t,1}$ $Q(s_t, a_1)$
$w_1$ $Q(s_t, a_2)$

Calculate loss

$$\left[ r_{t+1} + \gamma \max_{a} Q(s_{t+1}, a) - Q(s_t, a_t) \right]^2$$

Do training step

Freeze this network, update every N steps

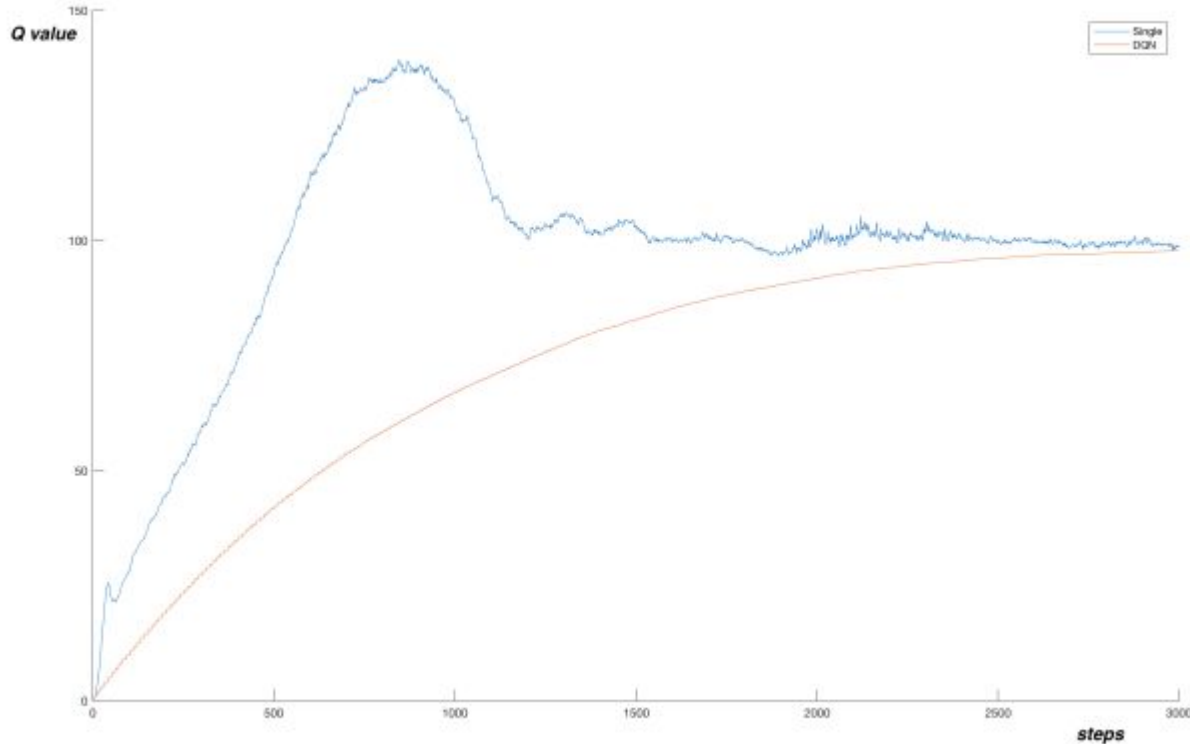**New Loss:** $\left[ r_{t+1} + \gamma \max_{a} \boxed{Q(s_{t+1}, a)} - Q(s_t, a_t) \right]^2$

1. **Explora**

Up
Sta

2. **Netwo**

s(t+1) ——→

Free
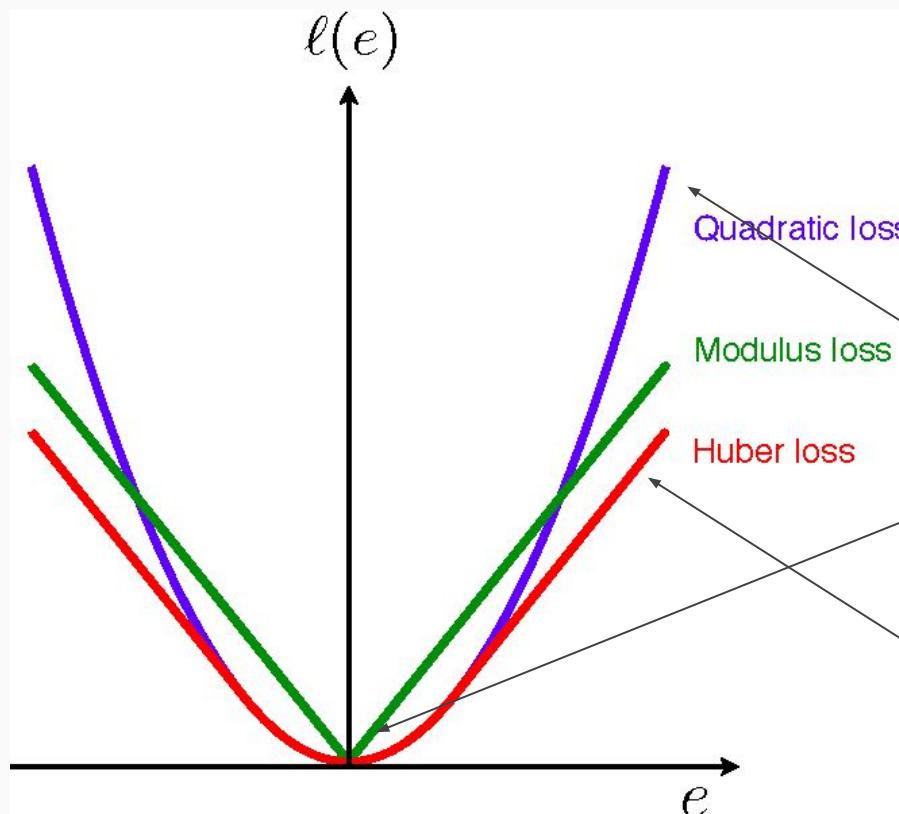
on
a)

o training step

$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \le \delta, \\ \delta(|a| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$$

For our purposes:
`tf.losses.huber_loss()`

Outliers have disproportionately strong effect
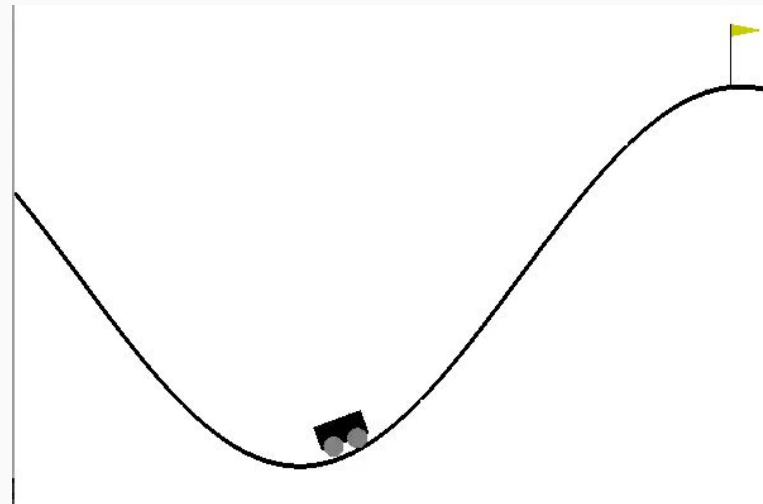
Sensitive to small perturbations

Best of both

**If you want to continue on this code-base:**
- Visualize training with tensorboard
- Try to implement some extensions: Double-DQN, Dueling Networks etc.
- How would any of this work in continuous action spaces?
- Try some of the other environments, MountainCar, Inverted pendulum etc.
- Can you get it working with a Convnet?
    - Warning: training times get out of hand quickly.
    - Test on the ATARI Envs

**Want to learn about other types of RL?**
- Policy Gradients
    - A3C, VPG, TRPO, PPO etc.
- Hybrid/Integrated Methods:
    - UNREAL, NEC,  Successor Learning etc.
- Applications/Misc.
    - Inverse RL/Imitation Learning,

**Resources:**
- To many to list here: UC Berkeley has tons of great stuff, start with their course or RL book by Richard Sutton