# DQN Walkthrough

Dan Neil

16 April 2015

# Quick Paper Overview

# History

## 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]

## Playing Atari with Deep Reinforcement Learning

Volodymyr Mnih    Koray Kavukcuoglu    David Silver    Alex Graves    Ioannis Antonoglou

Daan Wierstra    Martin Riedmiller

DeepMind Technologies

{vlad,koray,david,alex.graves,ioannis,daan,martin.riedmiller} @ deepmind.com
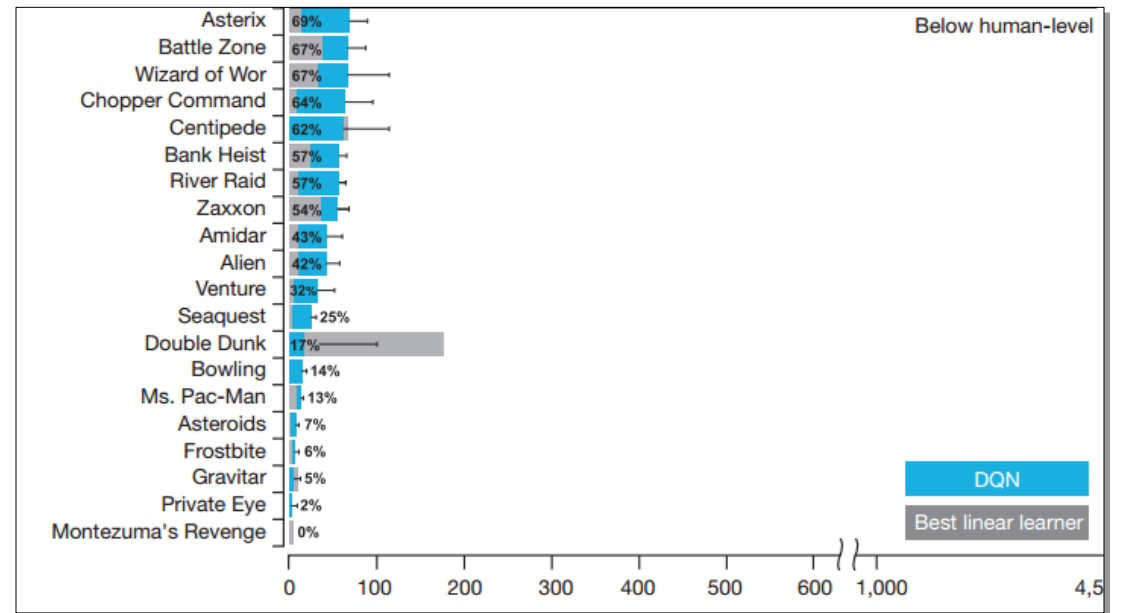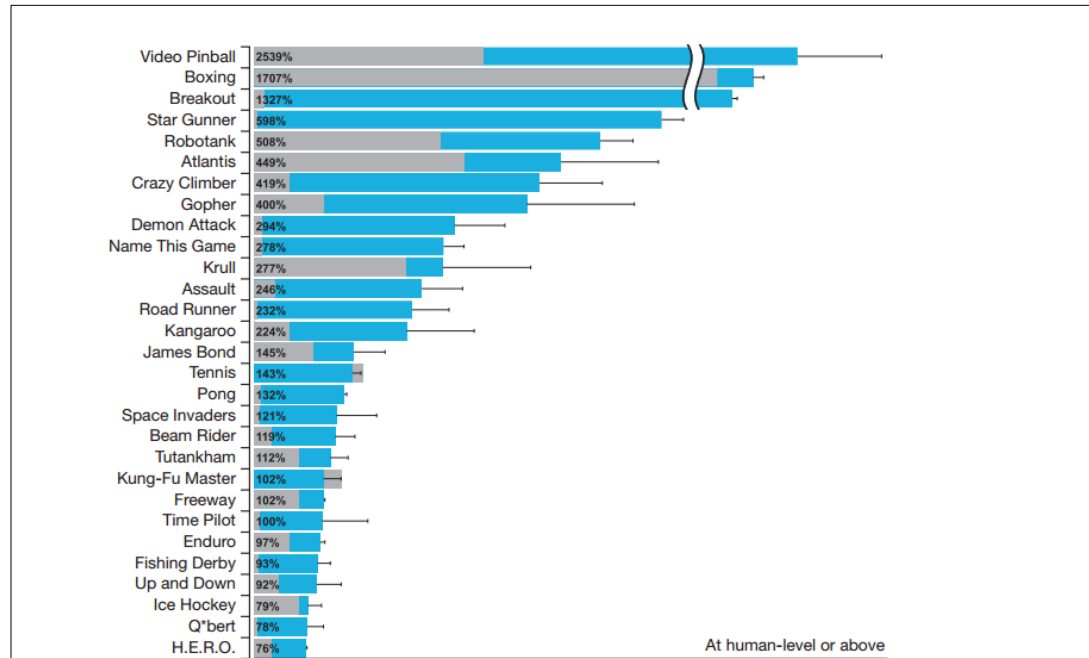
### Abstract

We present the first deep learning model to successfully learn control policies directly from high-dimensional sensory input using reinforcement learning. The model is a convolutional neural network, trained with a variant of Q-learning, whose input is raw pixels and whose output is a value function estimating future rewards. We apply our method to seven Atari 2600 games from the Arcade Learning Environment, with no adjustment of the architecture or learning algorithm. We find that it outperforms all previous approaches on six of the games and surpasses a human expert on three of them.
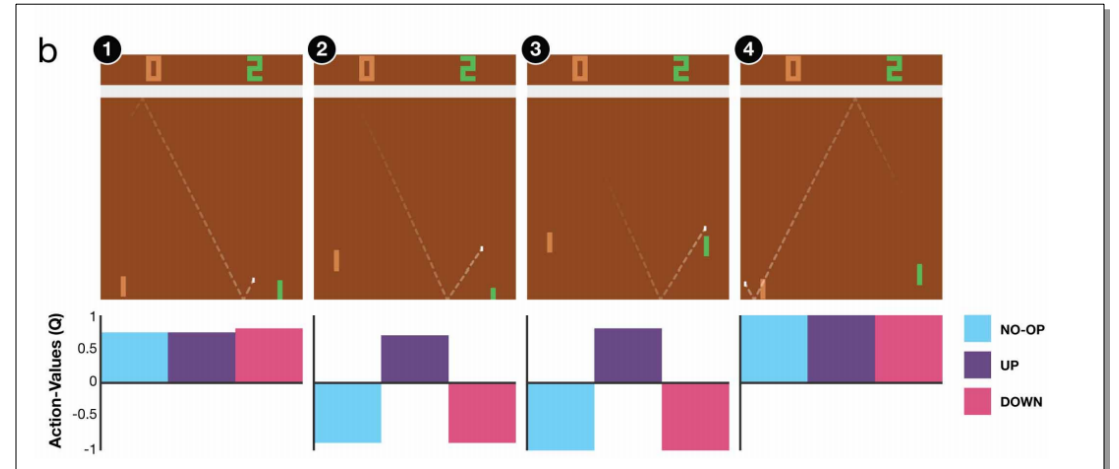
# Abstract (Following from Jakob Buhmann)

## Abstract

We present the first deep learning model to successfully learn control policies directly from high-dimensional sensory input using reinforcement learning. The model is a convolutional neural network, trained with a variant of Q-learning, whose input is raw pixels and whose output is a value function estimating future rewards. We apply our method to seven Atari 2600 games from the Arcade Learning Environment, with no adjustment of the architecture or learning algorithm. We find that it outperforms all previous approaches on six of the games and surpasses a human expert on three of them.
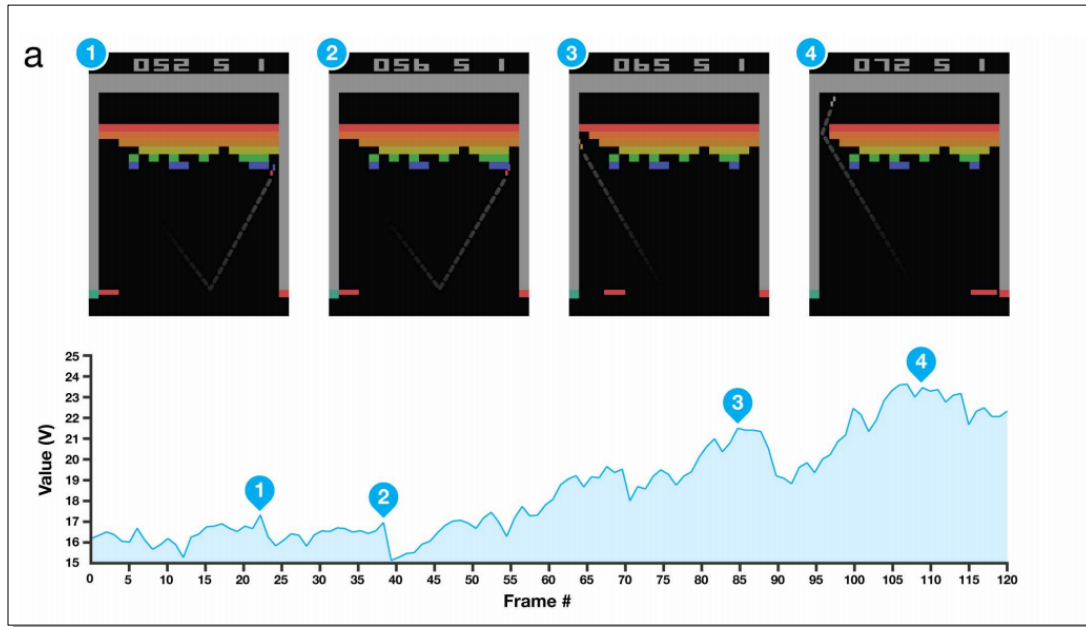
# Why DQNs?

# Show it working

# Reinforcement Learning Quick Recap

- Actor chooses an action
- Action is performed and influences the environment
- Based on the new state, the system receives a reward $r_t$
- System should learn to choose the action that maximizes the expected discounted reward
- Challenge is to learn in spite of sparse, noisy, or temporally shifted rewards

- **Q-Learning:**

  - State represented by a sequence $s_t = x_1, a_1, x_2, ..., a_t, x_t,$

  - Approximate action value function $Q(s, a)$

  - Find maximizing action $Q^*(s, a) = \max_a \mathbb{E}[R_t | s_t = s, a_t = a]$

  - Optimal solution fulfills Bellman equation

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

  - Approximate iteratively, s.t. $Q_i \to Q^*$ , $i \to \infty$

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

# Why function approximation?

- Learn a mapping of states to true rewards
  - Sounds like neural networks (model-free generic function approximators)
- Solution: neural networks
  - Fuzzy mapping from arbitrary features to rewards
  - Many tools exist to optimize learning

# Convolutional Networks

# Convolutional Network Embedding

# Challenges in NN as Function Approximators

- Correlations in successive states
- Small changes to Q-estimator may significantly change the policy and the subsequent data distribution
- Strong correlations exist between the action-values and the target values
- ✉Leads to feedback loops and instability.  How to solve?
- One previous approach: lots of randomly initialized networks
  - Very time costly

# Enter DQNs

- Two (and a half) key ideas:
- Replay memory
  - Store some memories of the past, and randomly sample from these instead of just the current state
  - Each memory can influence training more than once (good, efficient)
  - Random sampling from past breaks consecutive correlations
  - Behaviour distribution is distributed over history
- Target is only updated periodically, allowing the action-value estimator to diverge and reduce correlations
  - No longer does an increase of $Q(s_t, a_t)$ increase $Q(s_{t+1}, a)$, avoiding feedback loop
- Clip the error rate to (-1, 1)

# Code Deep Dive

# Algorithm Overview

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1$, $M$ **do**
   Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
   **For** $t = 1, T$ **do**
      With probability $\varepsilon$ select a random action $a_t$
      otherwise select $a_t = \text{argmax}_a\, Q(\phi(s_t), a; \theta)$
      Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
      Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
      Store transition $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$
      Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$

$$\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma\, \max_{a'} \hat{Q}\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$$

      Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the
      network parameters $\theta$
      Every $C$ steps reset $\hat{Q} = Q$
   **End For**
**End For**

# Algorithm Overview - Preprocess

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
   Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
   **For** $t = 1, T$ **do**
      With probability $\varepsilon$ select a random action $a_t$
      otherwise select $a_t = \mathrm{argmax}_a Q(\phi(s_t), a; \theta)$
      Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
      Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
      Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$
      Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$
$$\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma\, \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$
      Perform a gradient descent step on $\left(y_j - Q(\phi_j, a_j; \theta)\right)^2$ with respect to the network parameters $\theta$
      Every $C$ steps reset $\hat{Q} = Q$
   **End For**
**End For**

```
170 ▼ function nql:preprocess(rawstate)
171 ▼     if self.preproc then
172           return self.preproc:forward(rawstate:float())
173                       :clone():reshape(self.state_dim)
174       end
175
176       return rawstate
177   end
```

```
 7    require "nn"
 8    require "image"
 9
10    local scale = torch.class('nn.Scale', 'nn.Module')
11
12
13 ▼  function scale:__init(height, width)
14        self.height = height
15        self.width = width
16    end
17
18 ▼  function scale:forward(x)
19        local x = x
20        if x:dim() > 3 then
21            x = x[1]
22        end
23
24        x = image.rgb2y(x)
25        x = image.scale(x, self.width, self.height, 'bilinear')
26        return x
27    end
28
29    function scale:updateOutput(input)
30        return self:forward(input)
31    end
32
33    function scale:float()
34    end
```

# Algorithm Overview

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
  Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
  **For** $t = 1, \mathrm{T}$ **do**
    With probability $\varepsilon$ select a random action $a_t$
    otherwise select $a_t = \mathrm{argmax}_a Q(\phi(s_t), a; \theta)$
    Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
    Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
    Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$
    Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$
    Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$
    Perform a gradient descent step on $\left(y_j - Q(\phi_j, a_j; \theta)\right)^2$ with respect to the network parameters $\theta$
    Every $C$ steps reset $\hat{Q} = Q$
  **End For**
**End For**

```
315    self.transitions:add_recent_state(state, terminal)
316
317    local currentFullState = self.transitions:get_recent()
318
319    --Store transition s, a, r, s'
320 ▼  if self.lastState and not testing then
321        self.transitions:add(self.lastState, self.lastAction, reward,
322                             self.lastTerminal, priority)
323    end
324
325    if self.numSteps == self.learn_start+1 and not testing then
326        self:sample_validation_data()
327    end
328
329    curState= self.transitions:get_recent()
330    curState = curState:resize(1, unpack(self.input_dims))
331
332    -- Select action
333    local actionIndex = 1
334    if not terminal then
335        actionIndex = self:eGreedy(curState, testing_ep)
336    end
337
338    self.transitions:add_recent_action(actionIndex)
339
340    --Do some Q-learning updates
341 ▼  if self.numSteps > self.learn_start and not testing and
342        self.numSteps % self.update_freq == 0 then
343        for i = 1, self.n_replay do
344            self:qLearnMinibatch()
345        end
346    end
347
348    if not testing then
349        self.numSteps = self.numSteps + 1
350    end
351
352    self.lastState = state:clone()
353    self.lastAction = actionIndex
354    self.lastTerminal = terminal
355
356    if self.target_q and self.numSteps % self.target_q == 1 then
357        self.target_network = self.network:clone()
358    end
359
360    if not terminal then
361        return actionIndex
362    else
363        return 0
364    end
```

# Algorithm Overview

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
  Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
  **For** $t = 1,\mathrm{T}$ **do**
    With probability $\varepsilon$ select a random action $a_t$
    otherwise select $a_t = \mathrm{argmax}_a Q(\phi(s_t),a;\theta)$
    Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
    Set $s_{t+1} = s_t,a_t,x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
    Store transition $(\phi_t,a_t,r_t,\phi_{t+1})$ in $D$
    Sample random minibatch of transitions $(\phi_j,a_j,r_j,\phi_{j+1})$ from $D$
    Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \, \max_{a'} \hat{Q}(\phi_{j+1},a';\theta^-) & \text{otherwise} \end{cases}$
    Perform a gradient descent step on $(y_j - Q(\phi_j,a_j;\theta))^2$ with respect to the
    network parameters $\theta$
    Every $C$ steps reset $\hat{Q} = Q$
  **End For**
**End For**

```
315    self.transitions:add_recent_state(state, terminal)
316
317    local currentFullState = self.transitions:get_recent()
318
319    --Store transition s, a, r, s'
320 ▼  if self.lastState and not testing then
321        self.transitions:add(self.lastState, self.lastAction, reward,
322                             self.lastTerminal, priority)
323    end
324
325    if self.numSteps == self.learn_start+1 and not testing then
326        self:sample_validation_data()
327    end
328
329    curState= self.transitions:get_recent()
330    curState = curState:resize(1, unpack(self.input_dims))
331
332    -- Select action
333    local actionIndex = 1
334    if not terminal then
335        actionIndex = self:eGreedy(curState, testing_ep)
336    end
337
338    self.transitions:add_recent_action(actionIndex)
339
340    --Do some Q-learning updates
341 ▼  if self.numSteps > self.learn_start and not testing and
342        self.numSteps % self.update_freq == 0 then
343        for i = 1, self.n_replay do
344            self:qLearnMinibatch()
345        end
346    end
347
348    if not testing then
349        self.numSteps = self.numSteps + 1
350    end
351
352    self.lastState = state:clone()
353    self.lastAction = actionIndex
354    self.lastTerminal = terminal
355
356    if self.target_q and self.numSteps % self.target_q == 1 then
357        self.target_network = self.network:clone()
358    end
359
360    if not terminal then
361        return actionIndex
362    else
363        return 0
364    end
```

# Algorithm Overview – Epsilon Greedy

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
  Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
  **For** $t = 1, \text{T}$ **do**
    With probability $\varepsilon$ select a random action $a_t$
    otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$
    Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
    Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
    Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$
    Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$
    Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$
    Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters $\theta$
    Every $C$ steps reset $\hat{Q} = Q$
  **End For**
**End For**

```
368    function nql:eGreedy(state, testing_ep)
369        self.ep = testing_ep or (self.ep_end +
370            math.max(0, (self.ep_start - self.ep_end) * (self.ep_endt -
371            math.max(0, self.numSteps - self.learn_start))/self.ep_endt))
372        -- Epsilon greedy
373        if torch.uniform() < self.ep then
374            return torch.random(1, self.n_actions)
375        else
376            return self:greedy(state)
377        end
378    end
```

# Algorithm Overview – Greedy

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1$, $M$ **do**
   Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
   **For** $t = 1,\mathrm{T}$ **do**
      With probability $\varepsilon$ select a random action $a_t$
      otherwise select $a_t = \mathrm{argmax}_a\, Q(\phi(s_t),a; \theta)$
      Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
      Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
      Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$
      Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$

$$\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$

      Perform a gradient descent step on $\left(y_j - Q(\phi_j, a_j; \theta)\right)^2$ with respect to the network parameters $\theta$
      Every $C$ steps reset $\hat{Q} = Q$
   **End For**
**End For**

```
379
380
381 ▼  function nql:greedy(state)
382        -- Turn single state into minibatch.  Needed for convolutional nets.
383 ▼      if state:dim() == 2 then
384            assert(false, 'Input must be at least 3D')
385            state = state:resize(1, state:size(1), state:size(2))
386        end
387
388        if self.gpu >= 0 then
389            state = state:cuda()
390        end
391
392        local q = self.network:forward(state):float():squeeze()
393        local maxq = q[1]
394        local besta = {1}
395
396        -- Evaluate all other actions (with random tie-breaking)
397 ▼      for a = 2, self.n_actions do
398 ▼          if q[a] > maxq then
399                besta = { a }
400                maxq = q[a]
401            elseif q[a] == maxq then
402                besta[#besta+1] = a
403            end
404        end
405        self.bestq = maxq
406
407        local r = torch.random(1, #besta)
408
409        self.lastAction = besta[r]
410
411        return besta[r]
412    end
413
```

# Algorithm Overview – Store in Replay

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
    Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
    **For** $t = 1, T$ **do**
        With probability $\varepsilon$ select a random action $a_t$
        otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $\left( \phi_t, a_t, r_t, \phi_{t+1} \right)$ in $D$
        Sample random minibatch of transitions $\left( \phi_j, a_j, r_j, \phi_{j+1} \right)$ from $D$

$$\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}\left( \phi_{j+1}, a'; \theta^- \right) & \text{otherwise} \end{cases}$$

        Perform a gradient descent step on $\left( y_j - Q\left( \phi_j, a_j; \theta \right) \right)^2$ with respect to the
        network parameters $\theta$
        Every $C$ steps reset $\hat{Q} = Q$
    **End For**
**End For**

```
--Store transition s, a, r, s'
if self.lastState and not testing then
    self.transitions:add(self.lastState, self.lastAction, reward,
                         self.lastTerminal, priority)
end
```

# Algorithm Overview – Get Minibatch

**Algorithm 1: deep Q-learning with experience replay.**
Initialize replay memory $D$ to capacity $N$
Initialize action-value function $Q$ with random weights $\theta$
Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$
**For** episode $= 1, M$ **do**
   Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$
   **For** $t = 1, \text{T}$ **do**
      With probability $\varepsilon$ select a random action $a_t$
      otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$
      Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
      Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
      Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$
      Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$
$$\text{Set } y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$$
      Perform a gradient descent step on $\left(y_j - Q(\phi_j, a_j; \theta)\right)^2$ with respect to the
      network parameters $\theta$
      Every $C$ steps reset $\hat{Q} = Q$
   **End For**
**End For**

```lua
118  function trans:sample_one()
119      assert(self.numEntries > 1)
120      local index
121      local valid = false
122      while not valid do
123          -- start at 2 because of previous action
124          index = torch.random(2, self.numEntries-self.recentMemSize)
125          if self.t[index+self.recentMemSize-1] == 0 then
126              valid = true
127          end
128          if self.nonTermProb < 1 and self.t[index+self.recentMemSize] == 0 and
129              torch.uniform() > self.nonTermProb then
130              -- Discard non-terminal states with probability (1-nonTermProb).
131              -- Note that this is the terminal flag for s_{t+1}.
132              valid = false
133          end
134          if self.nonEventProb < 1 and self.t[index+self.recentMemSize] == 0 and
135              self.r[index+self.recentMemSize-1] == 0 and
136              torch.uniform() > self.nonTermProb then
137              -- Discard non-terminal or non-reward states with
138              -- probability (1-nonTermProb).
139              valid = false
140          end
141      end
142
143      return self:get(index)
144  end
```

# Algorithm Overview – Get Rewards

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory $D$ to capacity $N$

Initialize action-value function $Q$ with random weights $\theta$

Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$

**For** episode $= 1, M$ **do**

    Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

    **For** $t = 1, \text{T}$ **do**

        With probability $\varepsilon$ select a random action $a_t$

        otherwise select $a_t = \text{argmax}_a Q(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $\left(\phi_t, a_t, r_t, \phi_{t+1}\right)$ in $D$

        Sample random minibatch of transitions $\left(\phi_j, a_j, r_j, \phi_{j+1}\right)$ from $D$

        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}\left(\phi_{j+1}, a'; \theta^-\right) & \text{otherwise} \end{cases}$

        Perform a gradient descent step on $\left(y_j - Q\left(\phi_j, a_j; \theta\right)\right)^2$ with respect to the network parameters $\theta$

        Every $C$ steps reset $\hat{Q} = Q$

    **End For**

**End For**

```lua
-- The order of calls to forward is a bit odd in order
-- to avoid unnecessary calls (we only need 2).

-- delta = r + (1-terminal) * gamma * max_a Q(s2, a) - Q(s, a)
term = term:clone():float():mul(-1):add(1)

local target_q_net
if self.target_q then
    target_q_net = self.target_network
else
    target_q_net = self.network
end

-- Compute max_a Q(s_2, a).
q2_max = target_q_net:forward(s2):float():max(2)

-- Compute q2 = (1-terminal) * gamma * max_a Q(s2, a)
q2 = q2_max:clone():mul(self.discount):cmul(term)

delta = r:clone():float()

if self.rescale_r then
    delta:div(self.r_max)
end
delta:add(q2)

-- q = Q(s,a)
local q_all = self.network:forward(s):float()
q = torch.FloatTensor(q_all:size(1))
for i=1,q_all:size(1) do
    q[i] = q_all[i][a[i]]
end
delta:add(-1, q)

if self.clip_delta then
    delta[delta:ge(self.clip_delta)] = self.clip_delta
    delta[delta:le(-self.clip_delta)] = -self.clip_delta
end

local targets = torch.zeros(self.minibatch_size, self.n_actions):float()
for i=1,math.min(self.minibatch_size,a:size(1)) do
    targets[i][a[i]] = delta[i]
end
```

# Algorithm Overview – Update Target Network

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory $D$ to capacity $N$

Initialize action-value function $Q$ with random weights $\theta$

Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$

**For** episode = 1, $M$ **do**

    Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

    **For** $t = 1$,T **do**

        With probability $\varepsilon$ select a random action $a_t$

        otherwise select $a_t = \text{argmax}_a Q(\phi(s_t),a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t,a_t,x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t,a_t,r_t,\phi_{t+1})$ in $D$

        Sample random minibatch of transitions $(\phi_j,a_j,r_j,\phi_{j+1})$ from $D$

        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1},a'; \theta^-) & \text{otherwise} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j,a_j; \theta))^2$ with respect to the network parameters $\theta$

        Every $C$ steps reset $\hat{Q} = Q$

    **End For**

**End For**

```
356        if self.target_q and self.numSteps % self.target_q == 1 then
357            self.target_network = self.network:clone()
358    end
```

# What We Can Do

# DQNs for Game Playing

- For now, train a model using their default architecture
- Their parameters are identical across 49 games, indicating some robustness of these parameters
- In the future:
  - Pretrain visual model and reuse between games
  - Use model selection to pre-initialize game depending on type (puzzle, maze hack-and-slash, Mario-like, racing game, etc.)

# Possible Labor Division

- Input and Preprocessing
  - Interface to the game simulator
  - Represent game state as simplified image
- DQN Network without the deep part
  - Setup epsilon-greedy framework, replay memory, and Q-target copying
  - Build generalized DQN without the q-estimator part
- ConvNet Function Approximation
  - Set up a simple interface to use the convnet to estimate the rewards

# Appendix

# Alg. Overview ConvNet

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory $D$ to capacity $N$

Initialize action-value function $Q$ with random weights $\theta$

Initialize target action-value function $\hat{Q}$ with weights $\theta^- = \theta$

**For** episode $= 1, M$ **do**

    Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

    **For** $t = 1, \mathrm{T}$ **do**

        With probability $\varepsilon$ select a random action $a_t$

        otherwise select $a_t = \mathrm{argmax}_a Q(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $D$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $D$

        Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters $\theta$

        Every $C$ steps reset $\hat{Q} = Q$

    **End For**

**End For**

```lua
 9    function create_network(args)
10
11        local net = nn.Sequential()
12        net:add(nn.Reshape(unpack(args.input_dims)))
13
14        --- first convolutional layer
15        local convLayer = nn.SpatialConvolution
16
17        if args.gpu >= 0 then
18            net:add(nn.Transpose({1,2},{2,3},{3,4}))
19            convLayer = nn.SpatialConvolutionCUDA
20        end
21
22        net:add(convLayer(args.hist_len*args.ncols, args.n_units[1],
23                          args.filter_size[1], args.filter_size[1],
24                          args.filter_stride[1], args.filter_stride[1],1))
25        net:add(args.nl())
26
27        -- Add convolutional layers
28        for i=1,(#args.n_units-1) do
29            -- second convolutional layer
30            net:add(convLayer(args.n_units[i], args.n_units[i+1],
31                              args.filter_size[i+1], args.filter_size[i+1],
32                              args.filter_stride[i+1], args.filter_stride[i+1]))
33            net:add(args.nl())
34        end
35
36        local nel
37        if args.gpu >= 0 then
38            net:add(nn.Transpose({4,3},{3,2},{2,1}))
39            nel = net:cuda():forward(torch.zeros(1,unpack(args.input_dims))
40                    :cuda()):nElement()
41        else
42            nel = net:forward(torch.zeros(1,unpack(args.input_dims))):nElement()
43        end
44
45        -- reshape all feature planes into a vector per example
46        net:add(nn.Reshape(nel))
47
48        -- fully connected layer
49        net:add(nn.Linear(nel, args.n_hid[1]))
50        net:add(args.nl())
51        local last_layer_size = args.n_hid[1]
52
53        for i=1,(#args.n_hid-1) do
54            -- add Linear layer
55            last_layer_size = args.n_hid[i+1]
56            net:add(nn.Linear(args.n_hid[i], last_layer_size))
57            net:add(args.nl())
58        end
59
60        -- add the last fully connected layer (to actions)
61        net:add(nn.Linear(last_layer_size, args.n_actions))
62
63        if args.gpu >=0 then
64            net:cuda()
65        end
66        if args.verbose >= 2 then
67            print(net)
68            print('Convolutional layers flattened output size:', nel)
69        end
70        return net
71    end
72
```