

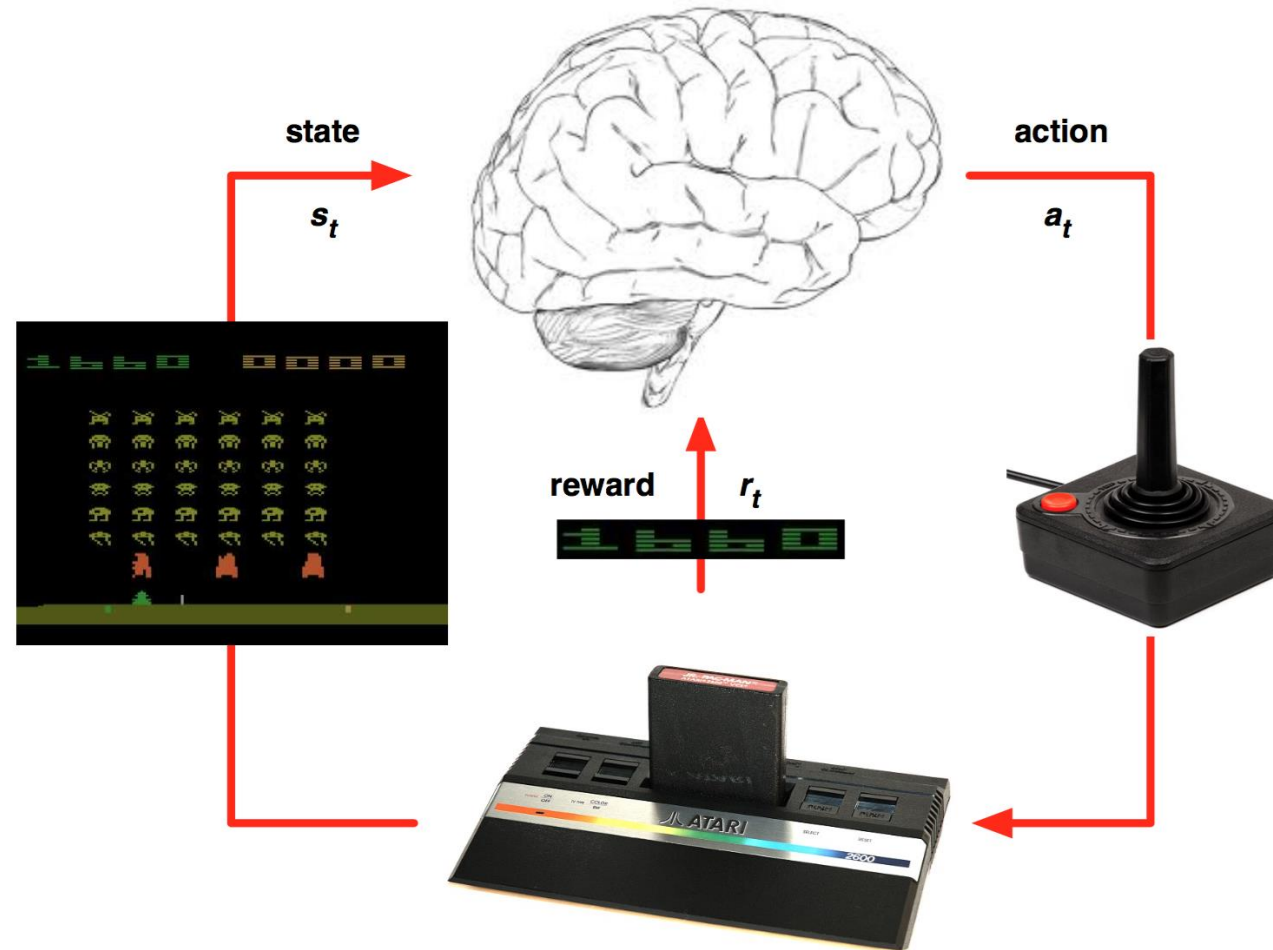
# DEEP Q LEARNING

By Comdet Phaudphut, [fb.com/comdet](https://fb.com/comdet), [comdet.p@gmail.com](mailto:comdet.p@gmail.com)

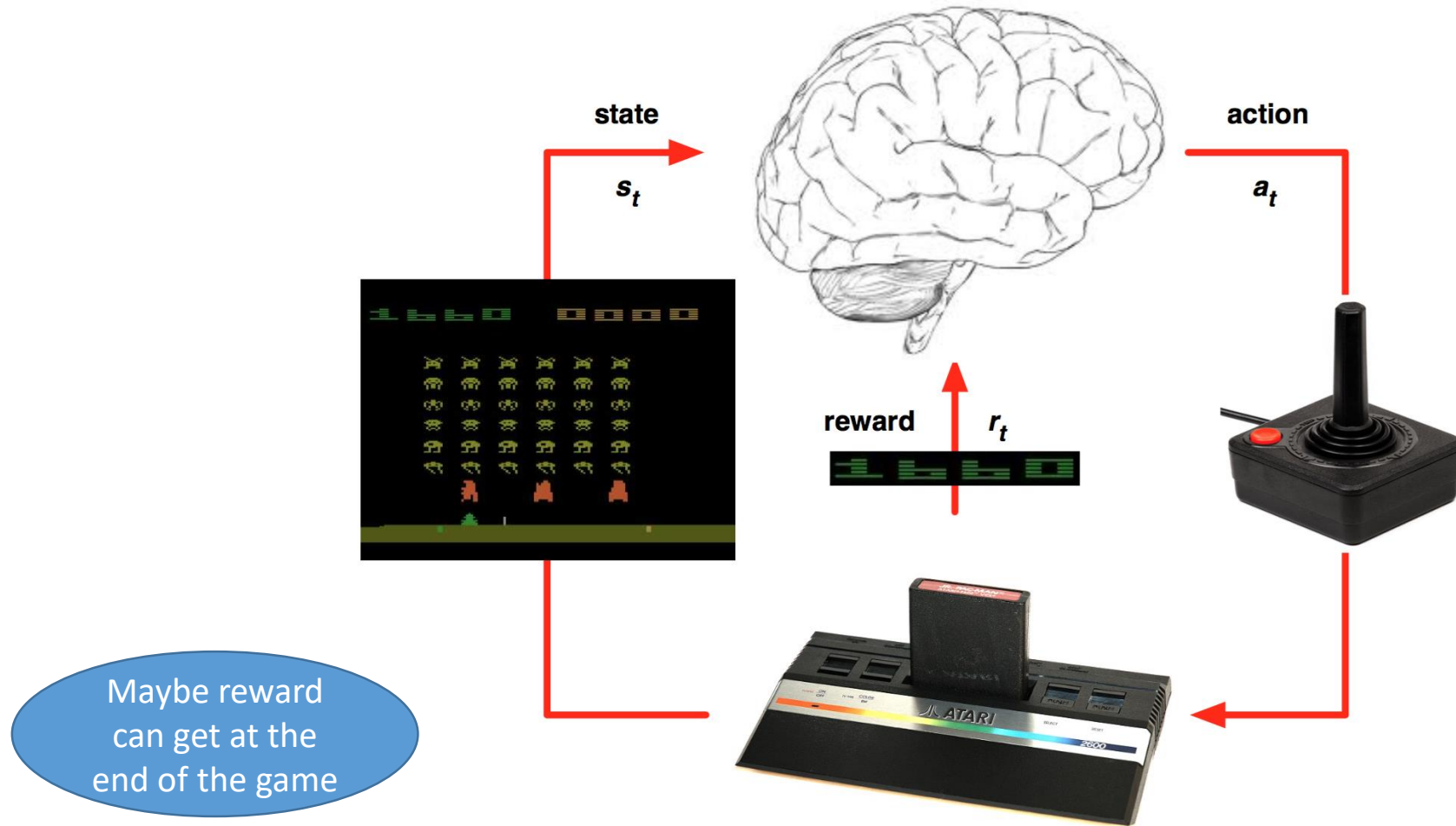
# Review Prev Lab

- We use CNN and gain more over 99% of accuracy.
- We know how convolution and pooling work.
- BUT
  - In real life data didn't represent as dataset ready to train

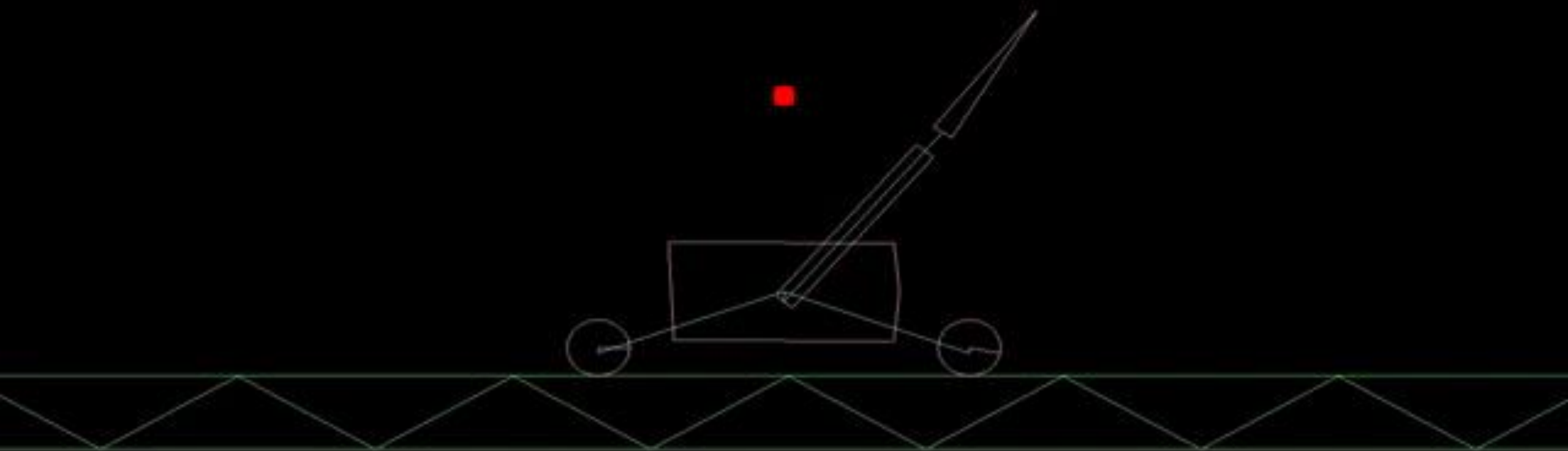
# Reinforcement Learning

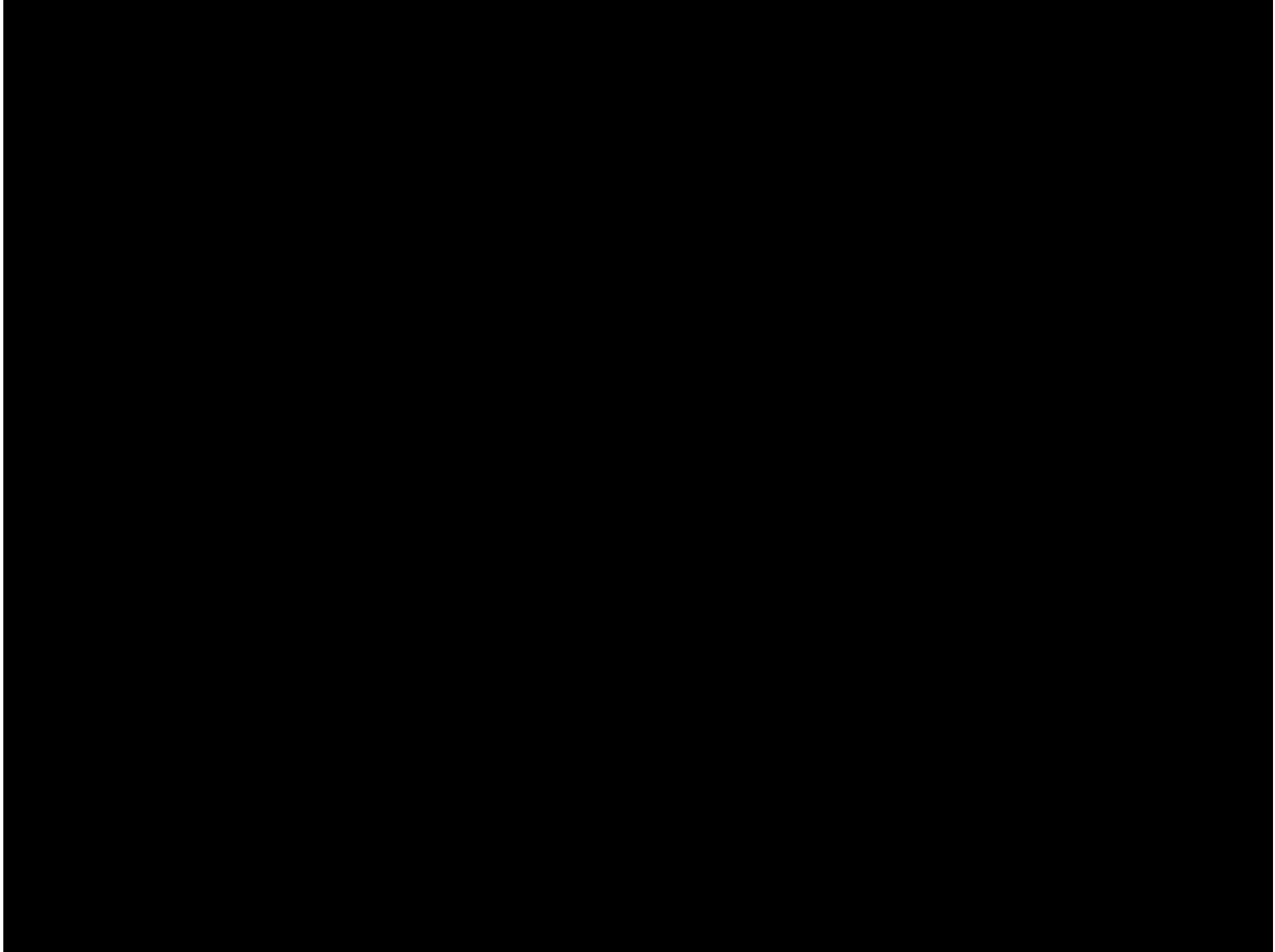


# Reinforcement Learning

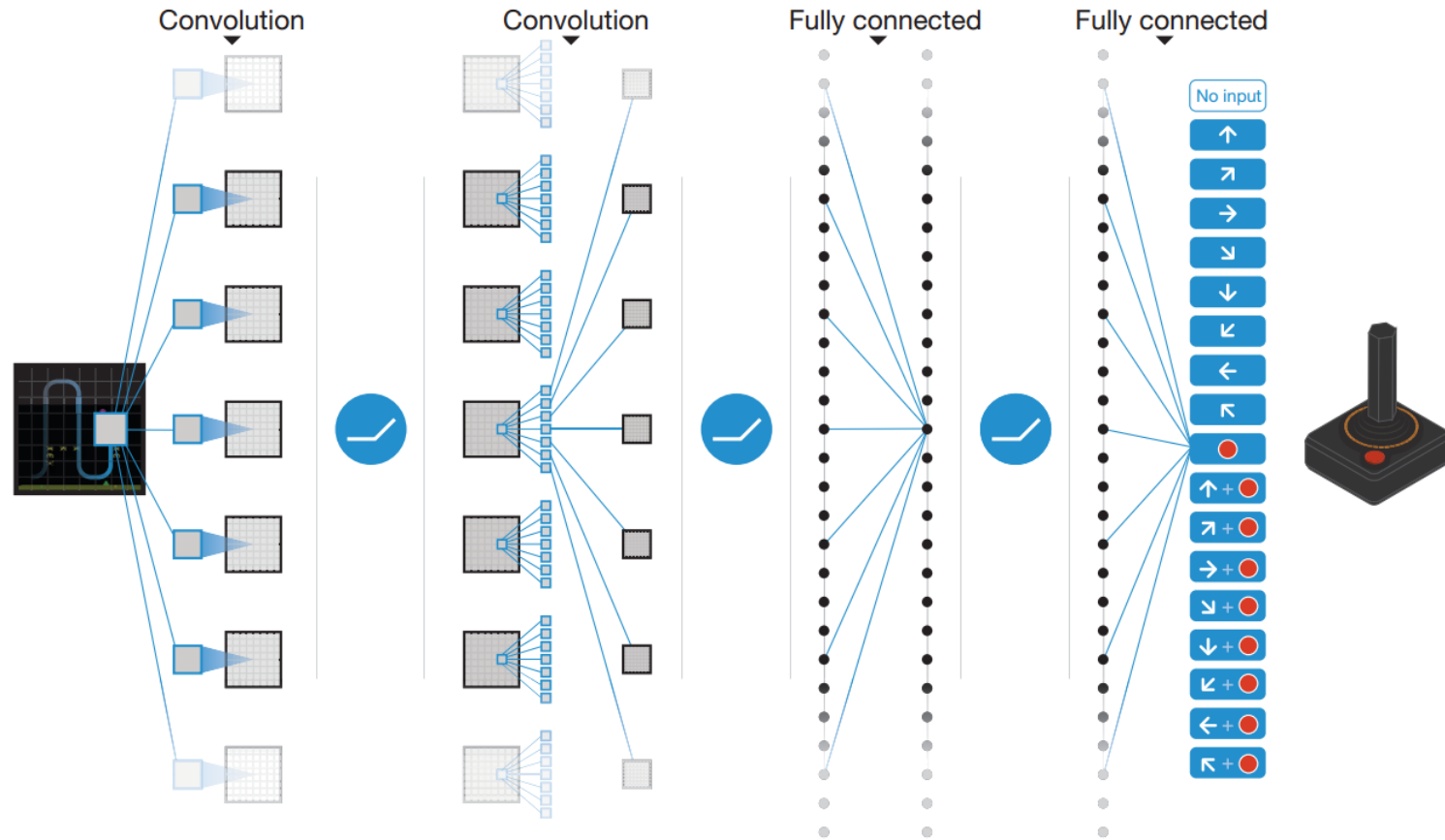


Goal: Move Right  
Q-Mode: Q-Table  
X: -0  
Speed: 0  
Acceleration: -0  
previousValue: -0.06  
valueDelta: -0  
timeSinceGoodValue: 13.6  
bestValue: 0.08  
worstValue: -0.07  
E: 0.1  
Alpha: 0.175  
Ipatience: 0.0043  
Max Randomness: 0.0999  
Min Randomness: 0.0597





# Algorithm



# Formulation

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

**2**  $Q'_t(s, a \mid \theta) \approx Q(s, a)$

**3**  $L(\theta) = \mathbf{E}_{s,a,r} \left[ \left( \mathbf{E}_{s'} \left[ r + \gamma \max_{a'} Q'_t(s', a' \mid \theta) \right] - Q_{t+1}(s, a \mid \theta) \right)^2 \right]$

Target



Variable





# Algorithm

---

**Algorithm 1** Deep Q-learning with Experience Replay

---

Initialize replay memory  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

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

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

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

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

        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        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  $\mathcal{D}$

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

        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3

**end for**

**end for**

---

# Our parameter

```
13  
14 ACTIONS = 2 # number of valid actions  
15 GAMMA = 0.99 # decay rate of past observations  
16 OBSERVE = 1000. # timesteps to observe before training  
17 EXPLORE = 2000000. # frames over which to anneal epsilon  
18 FINAL_EPSILON = 0.0001 # final value of epsilon  
19 INITIAL_EPSILON = 0.1 # starting value of epsilon  
20 REPLAY_MEMORY = 1000 # number of previous transitions to remember  
21 BATCH = 128 # size of minibatch  
22
```

# We create function that easy to use

```
1 def weight_variable(shape):
2     initial = tf.truncated_normal(shape, stddev = 0.01)
3     return tf.Variable(initial)
4
5 def bias_variable(shape):
6     initial = tf.constant(0.01, shape = shape)
7     return tf.Variable(initial)
8
9 def conv2d(x, W, stride):
10     return tf.nn.conv2d(x, W, strides = [1, stride, stride, 1], padding = "SAME")
11
12 def max_pool_2x2(x):
13     return tf.nn.max_pool(x, ksize = [1, 2, 2, 1], strides = [1, 2, 2, 1], padding = "SAME")
```

# We create variables

```
1  #create model
2  W_conv1 = weight_variable([8, 8, 4, 32])
3  b_conv1 = bias_variable([32])
4
5  W_conv2 = weight_variable([4, 4, 32, 64])
6  b_conv2 = bias_variable([64])
7
8  W_conv3 = weight_variable([3, 3, 64, 64])
9  b_conv3 = bias_variable([64])
10
11 W_fc1 = weight_variable([1600, 512])
12 b_fc1 = bias_variable([512])
13
14 W_fc2 = weight_variable([512, ACTIONS])
15 b_fc2 = bias_variable([ACTIONS])
16
```

# We create model

```
18 # input layer
19 s = tf.placeholder("float", [None, 80, 80, 4])
20
21 # hidden layers
22 h_conv1 = tf.nn.relu(conv2d(s, W_conv1, 4) + b_conv1) ### ==> 80x80x4 conv 4 ==> 20x20x32
23 h_pool1 = max_pool_2x2(h_conv1) ### ==> 20x20x32 maxpool => 10x10x64
24
25 h_conv2 = tf.nn.relu(conv2d(h_pool1, W_conv2, 2) + b_conv2) ###====10x10 conv2 ==>5x5x64
26 #h_pool2 = max_pool_2x2(h_conv2)
27
28 h_conv3 = tf.nn.relu(conv2d(h_conv2, W_conv3, 1) + b_conv3) ## 5x5 conv 1 padsame =>5x5x64
29 #h_pool3 = max_pool_2x2(h_conv3)
30
31 #h_pool3_flat = tf.reshape(h_pool3, [-1, 256])
32 h_conv3_flat = tf.reshape(h_conv3, [-1, 1600]) ##5x5x64 flatten =>1600
33
34 h_fc1 = tf.nn.relu(tf.matmul(h_conv3_flat, W_fc1) + b_fc1)
35
36 # readout layer
37 readout = tf.matmul(h_fc1, W_fc2) + b_fc2
```

# We create and here loss and game and D?

```
1 # define the cost function
2 a = tf.placeholder("float", [None, ACTIONS])
3 y = tf.placeholder("float", [None])
4
5 readout_action = tf.reduce_sum(tf.multiply(readout, a), reduction_indices=1)
6 cost = tf.reduce_mean(tf.square(y - readout_action)) #rms root mean square for cost function
7 train_step = tf.train.AdamOptimizer(1e-6).minimize(cost)
8
9 # open up a game state to communicate with emulator
10 game_state = game.GameState()
11
12 # saving and loading networks
13 saver = tf.train.Saver()
14 checkpoint = tf.train.get_checkpoint_state("saved_networks")
15 #if checkpoint and checkpoint.model_checkpoint_path:
16 #     saver.restore(sess, checkpoint.model_checkpoint_path)
17 #     print("Successfully loaded:", checkpoint.model_checkpoint_path)
18 #else:
19 #     print("Could not find old network weights")
20 |
21 # store the previous observations in replay memory
22 D = deque()
```

# Play game and saved in D

```
1 # start training
2 epsilon = INITIAL_EPSILON
3 rp = 0
4 while len(D) < REPLAY_MEMORY:
5
6     # choose an action epsilon greedily
7     readout_t = readout.eval(feed_dict={s : [s_t]})[0]
8     a_t, action_index = play_action(epsilon, readout_t)
9
10    ##### play ! #####
11    s_t1, r_t, terminal = get_stage(s_t, a_t)
12
13    # store the transition in D
14    D.append((s_t, a_t, r_t, s_t1, terminal))
15
16    s_t = s_t1
17    rp += 1
18    if(rp % 100 == 0):
19        print("TRY PLAY and RECORD : %d max readout %.4f" % (rp, np.max(readout_t)))
```

# Learning!

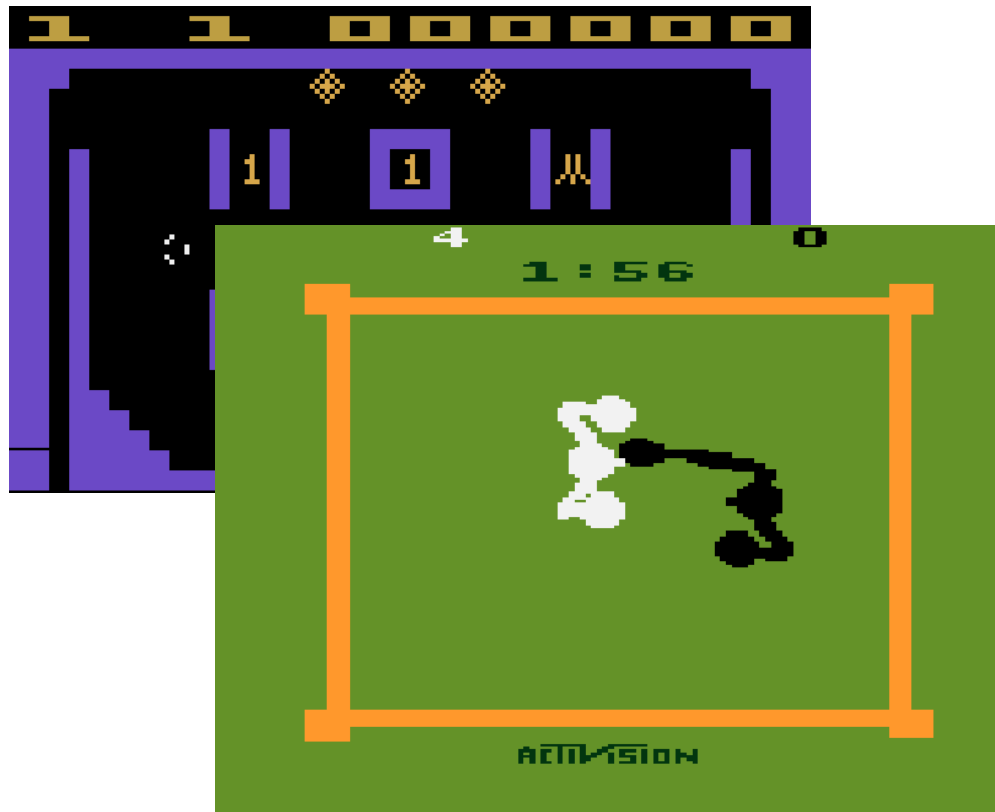
```
0 # run the selected action and observe next state and reward
1 s_t1,r_t,terminal = get_stage(s_t,a_t)
2 # store the transition in D
3 D.append((s_t, a_t, r_t, s_t1, terminal))
4
5 if len(D) > REPLAY_MEMORY:
6     D.popleft()
7 # sample a minibatch to train on
8 minibatch = random.sample(D, BATCH)
9
0 # get the batch variables
1 s_j_batch = [d[0] for d in minibatch]
2 a_batch = [d[1] for d in minibatch]
3 r_batch = [d[2] for d in minibatch]
4 s_j1_batch = [d[3] for d in minibatch]
5
6 y_batch = []
7 readout_j1_batch = readout.eval(feed_dict = {s : s_j1_batch})
8 for i in range(0, len(minibatch)):
9     terminal = minibatch[i][4]
0     # if terminal, only equals reward
1     if terminal:
2         y_batch.append(r_batch[i])
3     else:
4         y_batch.append(r_batch[i] + GAMMA * np.max(readout_j1_batch[i]))
5
6 # perform gradient step
7 train_step.run(feed_dict = {
8     y : y_batch,
9     a : a_batch,
0     s : s_j_batch}
1 )
2
3 # update the old values
4 s_t = s_t1
5 rp += 1
6
```



Let skip to final!

# Results Analysis

DQN is good at ...



DQN is bad at ...

