# Playing Atari with Deep Reinforcement Learning Paper Review
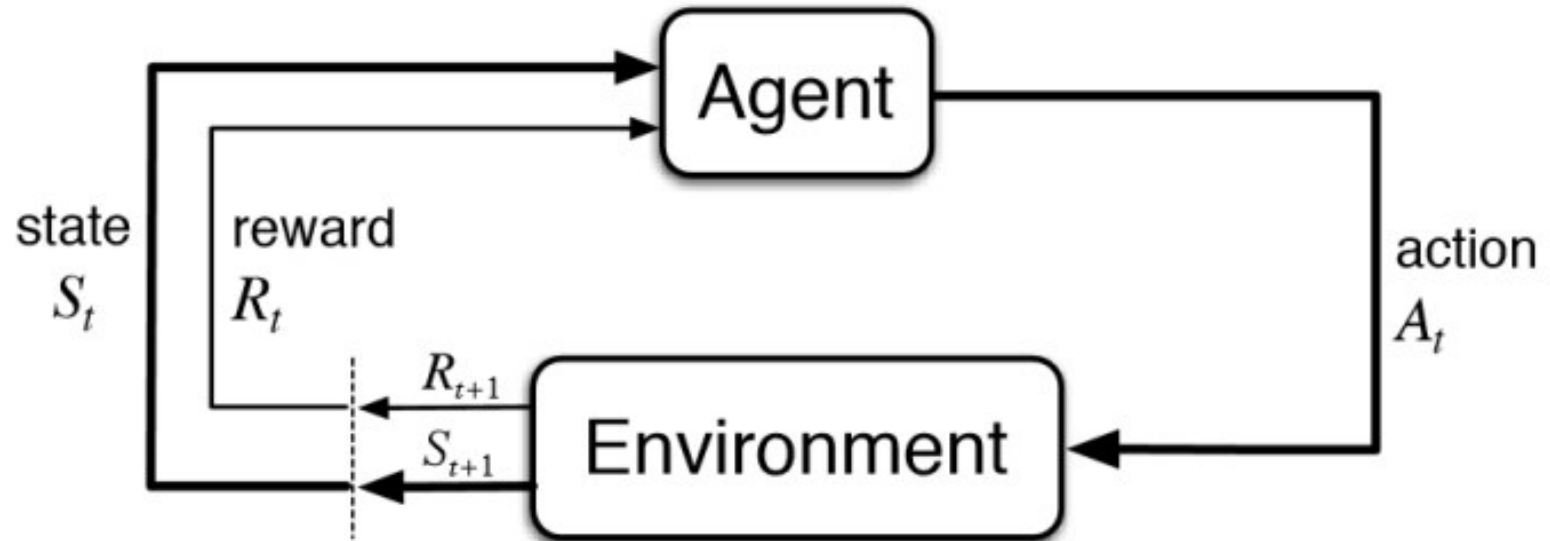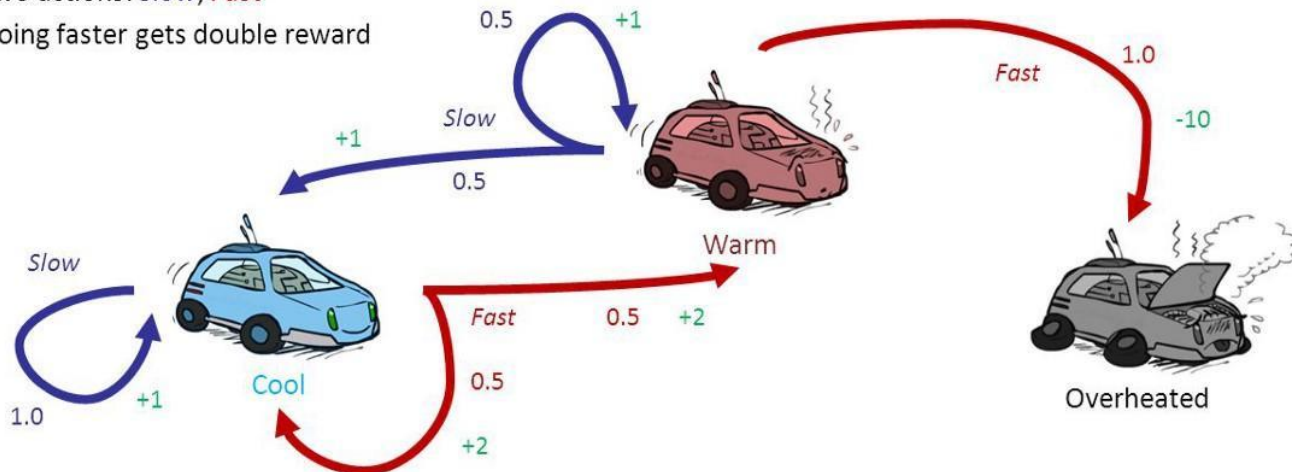
Tunahan PARLAYICI - 29913
Barış TEMEL - 19233

# Table Of Contents

# Reinforcement Learning Framework

# Markov Decision Process

## Example: Racing

- A robot car wants to travel far, quickly
- Three states: Cool, Warm, Overheated
- Two actions: Slow, Fast
- Going faster gets double reward

Slow
0.5    +1

Fast    1.0

+1    -10

0.5

Slow

Warm

Slow

Fast    0.5    +2

Cool    0.5

1.0    +1    +2

Overheated

# Bellman Equation, Q-learning

$$\text{New } Q(s,a) = Q(s,a) + \alpha [R(s,a) + \gamma \max Q'(s',a') - Q(s,a)]$$

New Q Value for that state and the action

Learning Rate

Reward for taking that action at that state

Current Q Values

Maximum expected future reward given the new state (s') and all possible actions at that new state.

Discount Rate

# Deep Q Network vs Q-Learning



Q Learning

Deep Q Learning

# SARSA vs Q-Learning

Cliff Walking Example - TD Learning On-Policy (SARSA) & Off-Policy (Q Learning)

R = -1

Safe path

Optimal path

S

The Cliff

G

R = -100

# Experience Replay

- Save transitions in a replay Memory D
- Random sample previous transitions to update parameters

# Experience Replay Advantages

- Greater Data Efficiency
- Prevents frames to be highly correlated
- Average the behaviour distribution

# Preprocessing

- Originally 210 × 160 image with 128 color palette
- Gray-scale and downsampled to 110 × 84, cropped to 84 × 84
- Frame stacking
- Stack 4 preprocessed frames as input
- Need multiple frames for velocity, etc
- Frame Skipping



endtoend. ai

# Model Architecture

## Model Architecture: Deep Q-Networks (DQN)

- Return $Q$ values for all 10 actions



4x84x84

16 8x8 filters

32 4x4 filters

256 hidden units

Fully-connected linear output layer

Stack of 4 previous frames

Convolutional layer of rectified linear units

Convolutional layer of rectified linear units

Fully-connected layer of rectified linear units

endtoend.ai

# Environment

- Same Algorithm implemented to 7 Atari Games
- Single Agent Problem
- Algorithm is model-free
  - No information to agents
  - No hand-crafted features
  - Same Neural Network structure and hyperparameters used in every game



Figure 1: Screen shots from five Atari 2600 Games: (*Left-to-right*) Pong, Breakout, Space Invaders, Seaquest, Beam Rider

# Pseudocode 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**
    Initialise 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**
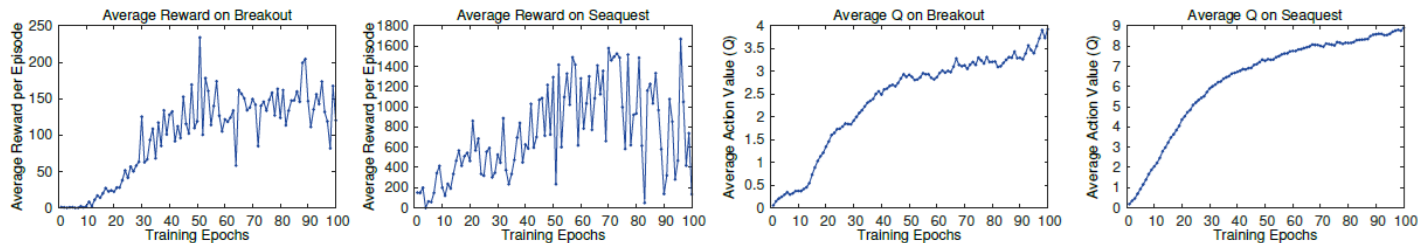**end for**

# Evaluation Metric



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 $\epsilon$-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.
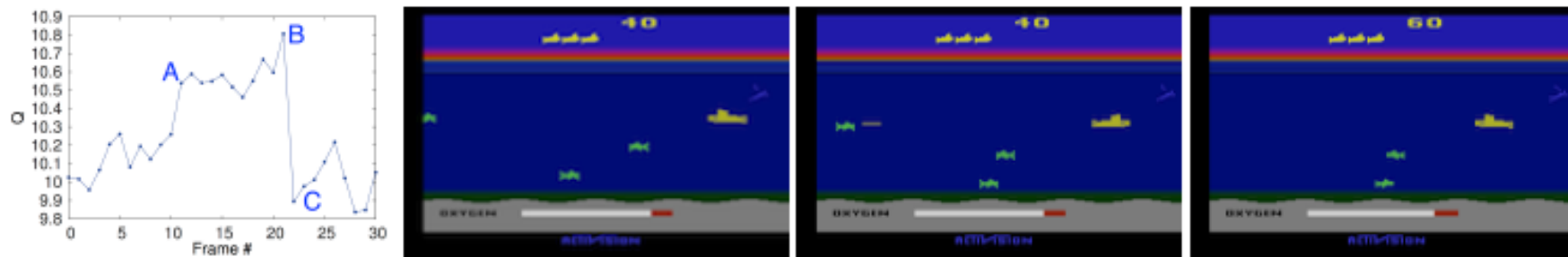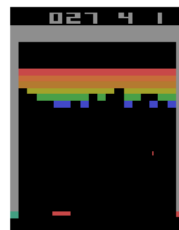
# Experiments



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.

# Results



Breakout game

| | 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 $\epsilon$-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 $\epsilon$-greedy policy with $\epsilon = 0.05$.