



The University of Manchester

Department of Computer Science  
Project Report 2020

**Reinforcement learning for a learnable agent  
in classic arcade games**

Author: Sean J Parker

Supervisor: Dr. Konstantin Korovin

## **Abstract**

### **Reinforcement learning for a learnable agent in classic arcade games**

**Author: Sean J Parker**

The aim of the project is to investigate the performance of Gismos and to design and construct a super multi-functional Gismo.

The novel aspects of the new Gismo are described. The abstract should perhaps be about half a page long.

The results of testing, which show the abject failure of the Gismo, are presented.

In the conclusions proposals for rectifying the deficiencies are outlined.

**Supervisor: Dr. Konstantin Korovin**

## **Acknowledgements**

I would like to thank my parents, my school teachers, my friends, my wonderful supervisor, and all my wonderful fellow students for their unswerving support during my project. Without your help none of this would have been possible.

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Motivation . . . . .	5
1.2	Objectives . . . . .	6
1.3	Report structure . . . . .	6
<b>2</b>	<b>Background</b>	<b>7</b>
2.1	Reinforcement learning . . . . .	7
2.1.1	Deep reinforcement learning . . . . .	7
2.2	DQN on Atari 2600 . . . . .	8
2.3	CNN Visualisation . . . . .	8
<b>3</b>	<b>Design</b>	<b>9</b>
3.1	Markov decision process . . . . .	9
3.2	Reinforcement learning . . . . .	9
3.2.1	Exploration vs Exploitation . . . . .	9
3.3	Q-Learning . . . . .	9
3.4	Q-Learning enhancements . . . . .	9
3.4.1	Double Q-Learning . . . . .	9
3.4.2	Duelling Q-Learning . . . . .	9
	<b>References</b>	<b>10</b>

# List of Figures

1.1	Screenshots of Pong, Breakout, Space Invaders (left to right). . . . .	6
2.1	OpenAI Robot solving a Rubik's cube . . . . .	8
2.2	Representation of end-to-end RL architectures . . . . .	8

# List of Tables

# Chapter 1

## Introduction

My project aims to replicate some of the reinforcement learning algorithms that can be used to play classic Atari 2600 games. It also compares the results of different tests with these algorithms such as varying the hyperparameters of the network. By observing the effects on the trained agents<sup>1</sup> when the hyperparameters are changed we can deduce a set of optimal values such that the networks can play three different Atari 2600 games. Overall, the main features of the project are the following:

- Agents with raw pixel game data as input, outputting a set of values for the best action.
- Agents attempt to find an optimal model of the environment without any prior knowledge.
- Visualization of the agent “brain” to provide insight into what information the agents is learning.

### 1.1 Motivation

Over the past 10 years there has been significant improvement in the RL (reinforcement learning) algorithms. One reason is that the computing power has become cheaply available by using discrete graphics cards. For example, for my project I used a Nvidia GeForce GTX 1070 that provides 1920 CUDA cores that can be used to accelerate training of neural networks. Despite this, RL algorithms are massively computationally expensive and hence take a long time to train.

Over recent years one of the pioneers in this area is DeepMind, which was acquired by Google in 2014, and they developed the DQN (deep q-network) algorithm in 2013 which they demonstrated could learn directly from the raw pixel data of games in order to achieve either human-level or super-human level performance.

This research was expanded upon by DeepMind and OpenAI which is based on the original DQN by DeepMind. This research focused on trying to approximate a Q-function and thereby infer the optimal policy. On the otherhand, there has recently been a focus on other methods such as A3C and PPO which instead seek to directly optimise in the policy space of the environment.

---

<sup>1</sup> Agent. In this case, agent refers to a trained neural network that takes actions in a chosen environment.

## 1.2 Objectives

Further to what was described in section 1 there was a few main objectives of the project. Firstly, I chose three games on which I decided to train the agents, Pong, Breakout, and Space Invaders which are shown below in Figure 1.1.

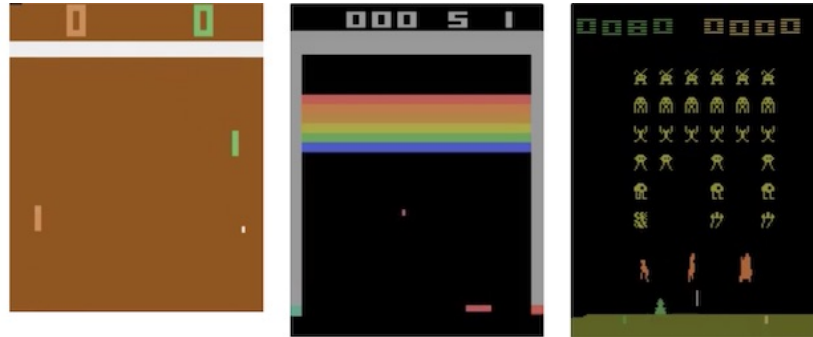


Figure 1.1: Screenshots of Pong, Breakout, Space Invaders (left to right).

Secondly, I wanted to find a way to explore the internals of a trained agent, in order to give further insight into what the agent is trying to learn. The reason for this is a researcher could use this information to determine, for example, where the agent has learnt to focus on the frame. Additionally, it provides a insight into how to optimally tune the hyperparameters which is described in section TODO: include section.

TODO: include gantt chart

## 1.3 Report structure

My report is divided into three main sections. Firstly, describing the background of the problem, then going onto giving details of my implementation and finally project evaluations/conclusions.



# Chapter 2

## Background

This chapter will cover some of the background material required for the following sections, it will cover the history of reinforcement learning (RL) and it's evolution to the current state-of-the-art. Additionally, it will cover the related work to this project and also cover some details of the past research papers for which this project has been based upon.

### 2.1 Reinforcement learning

Reinforcement learning is an area of machine learning has has been under active research since the late 1980s (TODO: ref watkins phd). The first defining algorithm of RL was called '*temporal-difference learning*' often referred to as TD-Learning. This algorithm learns by bootstrapping from the current value function in order iteratively converge towards a optimal policy (i.e. the agent's strategy for taking actions in the environment).

Further work by R. Sutton led to the development of TD-Lambda, an algorithm that was applied to the game of Backgammon, in 1992, by Gerald Tesauro to create TD-Gammon (TODO: ref Tesauro paper). It was a computer program that was shown to compete at expert-human level. The program also found novel strategies that were either unexplored, or dismissed in error as poor strategies. This was the first example of RL aiding in discovery or reconsideration of board game strategies. This trend of RL algorithms helping to improve human play would prove to only continue with DeepMind's AlphaZero program mastering the games of Chess (beating the strongest Chess programs such as Stockfish<sup>1</sup>) and Go.

#### 2.1.1 Deep reinforcement learning

Following from section 2.1 on RL, this section talks about the combination of two areas, deep learning and reinforcement learning methods, called deep reinforcement learning (DRL). Deep learning (DL) is a common class of machine learning methods that has been of much research focus over the past decade and can deal with high-dimensional sensory input; for the case of Atari this is 84x84 greyscale images after pre-processing of the raw Atari frames. On the other hand, reinforcement learning allows us to create an agent which can learn an optimal policy to navigate some environment in order to optimise its reward.

Through the combination of these methods, it has proven to provide solutions to previously intractable problems [1] in areas such as robotics, computer vision and healthcare. For example,

---

<sup>1</sup>Stockfish. One of the strongest Chess programs based on the CCRL ratings list.

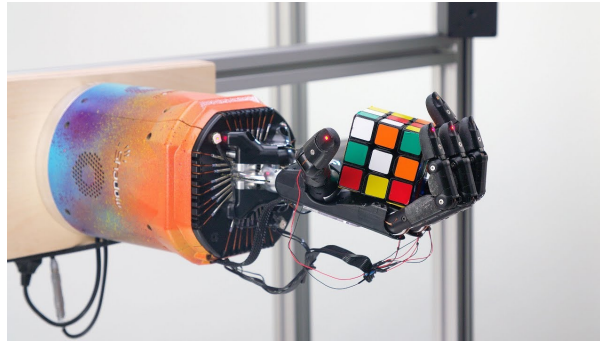


Figure 2.1: OpenAI Robot solving a Rubik's cube

in 2019 OpenAI developed a robotic hand that could solve a Rubik's cube [2.1](#), trained using deep reinforcement learning. As an extension to DRL, end-to-end reinforcement learning is a method for single layered neural network, trained by reinforcement learning. Figure [2.2](#) shows, diagrammatically, how DL and RL are used together in order to produce a single end-to-end model. In this simple architecture, there are two main components, the agent and the environment. This is a key feature of all DRL methods, an agent observes some state and reward from the environment after taking an action.

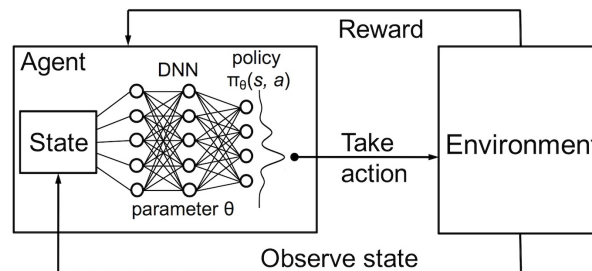


Figure 2.2: Representation of end-to-end RL architectures

Following on from Section [2.1](#)

## 2.2 DQN on Atari 2600

sdgsdsdfg

## 2.3 CNN Visualisation

sdfsdfsdfsdfsdf

# **Chapter 3**

## **Design**

### **3.1 Markov decision process**

### **3.2 Reinforcement learning**

#### **3.2.1 Exploration vs Exploitation**

### **3.3 Q-Learning**

### **3.4 Q-Learning enhancements**

#### **3.4.1 Double Q-Learning**

#### **3.4.2 Duelling Q-Learning**

# References

- [1] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath. A brief survey of deep reinforcement learning. *CoRR*, abs/1708.05866, 2017.