

Learning to Play Computer Games  
with Deep Learning and Reinforcement Learning

Final Report

COMP4801 Final Year Project  
UID: 3035207956

Student  
Mak Jeffrey Kelvin

Supervisor  
Dr. Dirk Schneiders

April 15, 2018

## **Abstract**

With the integration of deep learning into the traditional field of reinforcement learning in the recent decades, the spectrum of applications that artificial intelligence caters is currently very broad. As using AI to play games is a traditional application of reinforcement learning, the project's objective is to implement a deep reinforcement learning agent that can defeat a video game. Since it is often difficult to determine which algorithms are appropriate given the wide selection of state-of-the-art techniques in the discipline, proper comparisons and investigations of the algorithms are a prerequisite to implementing such an agent. As a result, this paper serves as a platform for exploring the possibility and effectiveness of using conventional state-of-the-art reinforcement learning methods for playing Pacman maps. In particular, this paper demonstrates that Combined DQN, a variation of Rainbow DQN, is able to attain high performance in small maps such as 506Pacman, smallGrid and mediumGrid. It was also demonstrated that the trained agents could also play Pacman maps similar to training with limited performance. Nevertheless, the algorithm suffers due to its data inefficiency and lack of human-like features, which may be remedied in the future by introducing more human-like features into the algorithm, such as intrinsic motivation and imagination.

# Acknowledgements

I would like to take this opportunity to thank my supervisor Dr. Dirk Schneiders from the Department of Computer Science for his advice and careful guidance in this project, as well as permission to use the GPU from the department. I would also like to express my gratitude towards Cezar Cazan from Center for Applied English Studies for various suggestions when drafting this paper. I also thank HKXF for providing financial support to this project through the FYP+ Support Scheme. Finally, a final thanks to HKU for providing a variety of quiet learning environments for me to write this paper without hindrance.

# Contents

<b>List of Figures</b>	<b>4</b>
<b>List of Tables</b>	<b>4</b>
<b>1 Abbreviations</b>	<b>6</b>
<b>2 Introduction</b>	<b>7</b>
<b>3 Background</b>	<b>8</b>
3.1 Deep Q Networks . . . . .	8
3.2 Variants of Deep Q Networks . . . . .	9
3.2.1 Double DQN . . . . .	9
3.2.2 Dueling Network architecture . . . . .	10
3.2.3 Prioritized Experience Replay . . . . .	10
3.2.4 Multi-step Q learning . . . . .	11
3.2.5 Noisy Networks . . . . .	11
3.2.6 Distributional DQN . . . . .	11
3.2.7 Rainbow DQN . . . . .	12
3.2.8 Count-Based Exploration and Intrinsic Motivation . . . . .	12
3.2.9 Deep Recurrent Q Network . . . . .	12
3.2.10 Other variants of Deep Q Network . . . . .	12
3.3 Other Deep Reinforcement Algorithms . . . . .	13
<b>4 Related Works</b>	<b>13</b>
4.1 History of Deep Reinforcement Learning . . . . .	13
4.2 Mrs. Pacman . . . . .	14
<b>5 Methodology</b>	<b>14</b>
5.1 Exploration phase . . . . .	15
5.2 Game . . . . .	15
5.3 Exploration Phase . . . . .	17
5.3.1 Deep Reinforcement Learning Algorithms . . . . .	17
5.3.2 Software and Hardware Requirements . . . . .	17
5.3.3 Neural Network . . . . .	17
5.3.4 Deep Q Networks and Its Variants . . . . .	18
5.3.5 Benchmark Algorithms . . . . .	19
5.3.6 Data Collection . . . . .	19
5.4 Implementation Phase . . . . .	20
5.4.1 Map . . . . .	20
5.4.2 Algorithm . . . . .	21
5.4.3 Experiment . . . . .	22
<b>6 Results and Discussion</b>	<b>23</b>
6.1 Exploration Phase . . . . .	23
6.2 Implementation Phase . . . . .	26
6.2.1 Generalization . . . . .	28
6.3 Limitation and Difficulties Encountered . . . . .	30
<b>7 Future Work</b>	<b>32</b>

8	Conclusion	34
9	References	34

## List of Figures

1	An illustration of the interaction between the agent and the environment through states and actions. Extracted from [1]. . . . .	7
2	Image extracted from [2]. CNN architecture used in DQN (left) and Dueling DQN (right). Note that the convolutional layers are followed by two separate fully connected layers, which splits the network into two parts, and leads to the value function and the advantage function as outputs of the CNN. The Q function is reconstructed from the the value function and advantage function at the end of the network. . . . .	10
3	A Pacman map consisting of all game elements. Pacman, represented by a yellow circle with a mouth, is enclosed by walls, as shown in a blue outline. The ghosts are represented by characters with two googly eyes. The dots and capsules are represented by small and large white dots. Image produced in the Pacman software by UC Berkeley[3]. . . . .	15
4	Examples of start game states for 506Pacman, a $3 \times 3$ map consisting of Pacman, one dot and one ghost that are randomly placed. Image produced in the Pacman software by UC Berkeley[3]. . . . .	16
5	Map of smallGrid. The map consists of pacman, one ghost, two dots and walls in a $7 \times 7$ map. Image produced in the Pacman software by UC Berkeley[3]. . .	16
6	Start game states for Pacman maps mediumGrid and smallClassic. Image produced in the Pacman software by UC Berkeley[3]. . . . .	16
7	$32 \times 32$ greyscale image used to observe change in Q values during training. Image produced in the Pacman software by UC Berkeley[3]. . . . .	19
8	Number of game state visited over the when playing a random policy in 506Pacman, smallGrid, mediumGrid, miniClassic and smallClassic. Note that game states are also classified as different in the Pacman API if the sprites' locations are the same but the direction Pacman faces is different. . . . .	20
9	Deep Q Networks training results with default ghost AI in 506Pacman . . . . .	22
10	Deep Q Networks training results with default ghost AI in 506Pacman . . . . .	24
11	Deep Q Networks training results with minimax ghost AI in 506Pacman . . . . .	24
12	Deep Q Networks training results with default ghost AI in 506Pacman . . . . .	25
13	Deep Q Networks training losses in Pacman . . . . .	25
14	Training loss of DQN with proportional-based PER in 506Pacman for 6000 games	26
15	Results obtained from playing 506Pacman with random ghost AI for 10000 games	27
16	Results obtained from playing smallGrid with random ghost AI for 5000 games .	28
17	Results obtained from playing mediumGrid with random ghost AI for 5000 games	29
18	Results obtained from playing smallClassic with random ghost AI for 10000 games	30
19	Sum of noisy training parameters during training in 506Pacman, smallGrid and mediumGrid. No data was gathered due to the policy having 0% win rate for smallClassic. . . . .	31

## List of Tables

1	Reported Ms. Pacman scores from [4] . . . . .	14
2	Collation of Ms. Pacman scores from [5], [6], [2], [7], [8], and [9] . . . . .	14
3	Score modification for various events in the game Pacman. Ghosts become scared for 40 moves when Pacman eats a capsule. . . . .	15

4	Description of the convolutional neural network used in the DQN, DDQN, Duel DQN and DQN with PER implementation for Pacman. The input of the network is an array of $32 \times 32$ preprocessed grayscale images for the three implementations, and the output of the network is an array, where each element consists of five Q-values, one for each possible action from the input image's state. The network is trained using Adam [10] as the optimizer, where the huber loss is minimized. The total number of trainable parameters in this network is 87205. . . . .	18
5	List of DQN training hyperparameters and their values in Pacman . . . . .	18
6	Shaped reward for agent in Pacman. The rewards are kept between -1 and 1 in order to prevent from obtaining large gradients, which may hinder learning. . . .	21
7	List of Combined DQN training hyperparameters and their values in Pacman . .	22
8	Win rate and average score obtained after running a trained 506Pacman agent using different ghost AIs in 506Pacman. . . . .	29
9	Win rate and average score obtained after running a trained 506Pacman agent in 506Pacman containing different number of dots. . . . .	30
10	Results from testing smallGrid agent on different variants of smallGrid . . . . .	31
11	Project Schedule . . . . .	33

# 1 Abbreviations

<b>AI</b>	Artificial Intelligence
<b>CNN</b>	Convolutional Neural Network
<b>DRL</b>	Deep Reinforcement Learning
<b>DQN</b>	Deep Q Networks
<b>DDQN</b>	Double Deep Q Networks
<b>Duel DQN</b>	Dueling Deep Q Networks
<b>DRQN</b>	Deep Recursive Q Networks
<b>MDP</b>	Markov Decision Process
<b>PER</b>	Prioritized Experience Replay
<b>PIL</b>	Python Image Library
<b>POMDP</b>	Partially Observable Markov Decision Process
<b>RL</b>	Reinforcement Learning



## 2 Introduction

Reinforcement learning, a subdiscipline of artificial intelligence, is becoming increasingly popular as a method of creating AI agents. Despite of such popularity, the idea of reinforcement learning is not a recent one, and originated from experiments in behavioural psychology relating to how animals learn tasks by perceiving rewards. In the context of computer science, reinforcement learning deals with how an agent learns an optimal policy  $\pi^*$ , ie. a series of actions, that maximizes the expected cumulative reward  $R_t = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$  by interacting with the environment, where  $r_t$  is the reward received at time step  $t$  and  $\gamma$  is the discount rate.

In order for an agent to learn, the manner in which it interacts with the environment is key for facilitating learning. In the field of reinforcement learning, there are two major approaches to modelling agent-environment interaction in reinforcement learning, namely Markov Decision Process (MDP) and Partially Observable Markov Decision Process (POMDP). In an MDP, the action  $A_t$  taken by the agent is based on only the reward  $r_t$  and state  $s_t$  observed from the environment at time step  $t$ , as illustrated in Figure 1. However, it is assumed that the agent always observes full information of the current state in each time step. To compensate for such a limitation, one can adopt a Partially Observable Markov Decision Process (POMDP) instead, where the agent only observes partial information of its current state. For example, the interaction would be modelled as a MDP for the game Pacman, since the player can see the entire map in Pacman during gameplay. Conversely, a POMDP would be used as a model for interactions in FPS games [11], since the player's view is often obscured by obstacles. This makes certain game information hidden, such as the position of enemies.

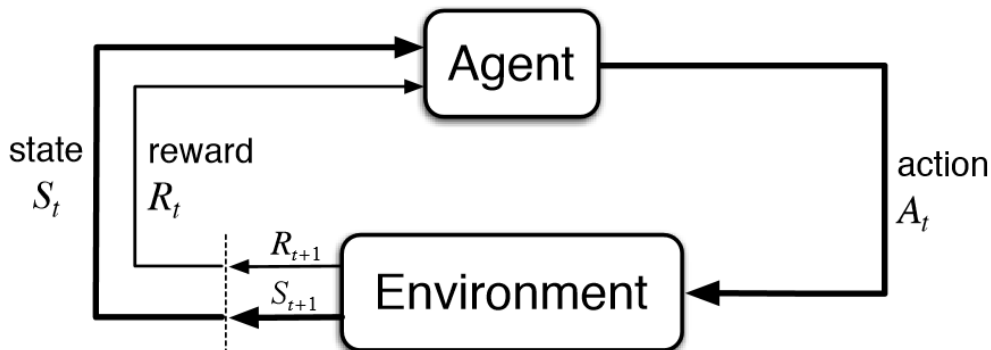


Figure 1: An illustration of the interaction between the agent and the environment through states and actions. Extracted from [1].

One may question the rationale of using deep reinforcement learning to play games, since the primary purpose of games is for entertainment. Nonetheless, a major reason for this is due to the relatively cheap training cost in a game simulation in comparison to training in the physical world. In particular, training in the real world often requires regular hardware maintenance of physical agents such as robots. For instance, a robot attempting to learn to walk would fall constantly during training, which would become unfavorable due to constant repairing of the robot. Another reason is that replication of agent training is much easier in games than in the real world, since the physical environment is often unpredictable. A consequence of this is that games serve as a convenient platform for benchmarking and comparing various reinforcement learning techniques. Finally, games are traditionally considered as a method of measuring human intelligence. This is especially apparent in chess, where logical reasoning

plays a key component for the construction of a winning strategy. As a result, games can be used for comparison between artificial intelligence and human intelligence.

In light of deep reinforcement learning being an active research field and game playing as convenient medium for data collection, this project’s objective is to construct a deep reinforcement learning agent in an unexplored game in literature. Nevertheless, proper comparison and investigation of existing deep reinforcement learning algorithms is required. As a result, this paper explore the possibilities of using Deep Q Network and its variants to play the game Pacman through small-scale experiments, while creating software implementations of the aforementioned techniques. With the results in hand, this paper serves as a gateway for an improved understanding on pros and cons of well-studied techniques, as well as their applicability in small-scale problems in game playing.

The remainder of the paper is organized as follows. Chapter 3 surveys existing state-of-the-art deep reinforcement learning algorithms in the literature, and chapter 4 summarizes the list of games that have been investigated in the literature. Chapter 5 details the phases in the project, followed by implementation details of Deep Q Network and its variants and experiments used to evaluate the techniques. Chapter 6 then compares the implemented methods by analyzing experimental results, and concludes by discussing the future work to be completed in Chapter 7.

## 3 Background

### 3.1 Deep Q Networks

Similar to tabular Q learning in traditional reinforcement learning, Deep Q Networks [12, 4] also uses Q values in order to predict the expected sum of future discounted rewards. However, a problem with tabular Q learning is that the Q values for state-actions pairs are encoded using a table, which uses a large amount of memory in the case of an environment involving a large state-action space. Moreover, each Q-value for different state-action pairs needs to be learned separately, thereby making improvement in training performance slow. As a result, the replacement of the Q table by a convolutional neural network in the form of a nonlinear function approximator in DQN that minimizes the loss function

$$L_i(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim U(D)} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta'_i) - Q(s, a; \theta_i) \right)^2 \right]$$

, where  $U(D)$  is the experience replay training minibatch,  $\theta_i$  is the network parameters for the online neural network and  $\theta'_i$  is the network parameters for the target neural network. The use of a convolutional neural network greatly reduces the memory required to encode the Q-values, as well as allowing the agent to generalize encoded information learned from the past to unfamiliar scenarios.

Apart from the use of a convolutional neural network, a prominent feature in DQN is the use of experience replay. Such a concept originated from the observed hippocampus activity in the brain. In particular, experience replay involves the use of an experience replay buffer, where the agent’s past experience is stored in the form of experience tuples  $e_t = (s_t, a_t, r_t, s_{t+1})$ .

During gameplay the agent learns from past experience by sampling mini-batches from the experience replay buffer in order to train the CNN. Such sampling allows for increased use of past memory when compared to the on-policy variant, which reduces the time taken to learn a policy. Moreover, experience replay prevents the CNN from being trained on consecutive experience tuples, which could cause convergence to a suboptimal policy [12].

To solve the exploration-exploitation dilemma, the DQN algorithm uses an epsilon-greedy algorithm. In particular, the agent picks actions uniformly at random with probability  $\epsilon$  and greedily with probability  $1 - \epsilon$ . To encourage exploration during the start of training and exploitation near the end of training, epsilon is set to 1 and linearly annealed to 0.1 over training games [12, 1].

To improve the training of DQN, the original paper proposed several implementation tricks used during coding [12]. First, the raw image is preprocessed in order to reduce the size of the CNN while still retaining sufficient information to solve the desired task. Various techniques can be used for preprocessing, such as downscaling of images and conversion of RGB image to grayscale. Owing to the large variation in reward ranges between different Atari 2600 games, the original papers uses clipping rewards to the range  $[-1, 1]$  in order to limit the gradient during training. Lastly, a target network is used in order to stabilize the Q value during training [12].

An advantage of Deep Q Networks is that minimal prior game-specific information can be used to develop an end-to-end framework, where the input and output of the network would be the preprocessed image from the game and the Q values of the state-action pairs of the current state respectively. In other words, the algorithm is agnostic to domain knowledge. This idea is also reinforced by the fact that the agent only learns by playing the game itself. In particular, the agent begins with minimal information about the game’s environment (eg. controller inputs), and training data is generated and learned by the agent as training occurs. The use of an end-to-end framework is unlike previous approaches, where game-specific feature extraction to construct an internal grid world representation of the game state coupled with tabular Q learning was required [13] in order for learning to occur, which may limit the policies that the agent could develop.

## 3.2 Variants of Deep Q Networks

After the development of DQN in 2012, many variations of DQN with improvements were developed subsequently. Some of the significant variations are detailed in the following subsections. Note that the improvements can be combined together, since they each focus on a different component of the algorithm.

### 3.2.1 Double DQN

One limitation of DQN is that Q values are often overestimated during training, resulting in overly optimistic policies. To remedy this situation, Double DQN, ie. DDQN, [5] is introduced where the target value for the CNN to train with is modified to

$$r + \gamma \max_{a'} Q(s', \operatorname{argmax}_a Q(s', a; \theta_i); \theta'_i)$$

Such modification increases the stability of the Q values during learning, thus allowing the algorithm to attain improved policies, as shown by the substantially better scores attained in various Atari 2600 games [5].

### 3.2.2 Dueling Network architecture

Another limitation of DQN is that Q values for different actions for the same state are learned separately. To allow the network to generalize learning of state value across state-action pairs with the same states, one may consider the following definition of the Q value function

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + A(s, a; \theta, \alpha)$$

where  $V$  is the value function,  $A$  is the advantage function, and  $\theta$ ,  $\alpha$ , and  $\beta$  are parameters of the dueling CNN. However, the paper describes the use of such an equation for estimating Q values as inefficient, since the value and advantage functions developed are not unique to the problem. To fix this issue, the paper proposes to use the definition of the Q value function

$$Q(s, a; \theta, \alpha, \beta) = V(s; \theta, \beta) + \left( A(s, a; \theta, \alpha) - \max_{a' \in |A|} A(s, a'; \theta, \alpha) \right)$$

instead [2]. The use of the above equation then leads to the modification of the CNN architecture as shown in Figure 2, which allows the dueling variant of DQN, ie. Duel DQN, to learn the two functions separately. When compared to DDQN, Duel DQN outperforms it in 46 out of 57 Atari 2600 games [2]. It is also mentioned that the Duel DQN with PER outperforms DDQN significantly.

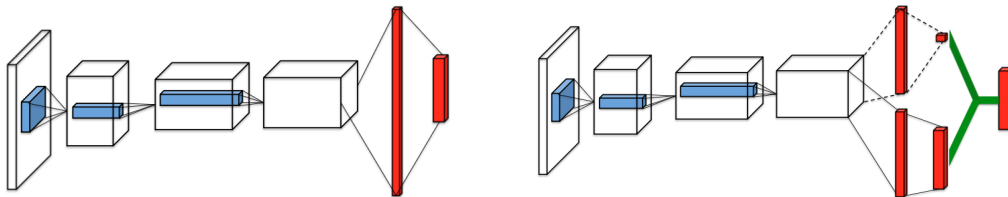


Figure 2: Image extracted from [2]. CNN architecture used in DQN (left) and Dueling DQN (right). Note that the convolutional layers are followed by two separate fully connected layers, which splits the network into two parts, and leads to the value function and the advantage function as outputs of the CNN. The Q function is reconstructed from the the value function and advantage function at the end of the network.

### 3.2.3 Prioritized Experience Replay

At the end of [4] and [12], Minh et al. described that it may be possible to improve the algorithm via prioritized sweeping, where training would be biased towards significant events. Such idea

was subsequently further developed by Schaul et al. [6], where two types of prioritized experience replay (PER) was investigated, namely rank-prioritized PER and proportional-prioritized PER. In particular, prioritized experience replay involves a non-uniform distribution when sampling for experience tuples, where the probability of sampling experience tuple  $i$  is

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

In particular,  $p_i = |\delta_i| + \epsilon$  in proportional prioritization, where  $\delta$  is the training error and epsilon is a small constant, and  $p_i = \frac{1}{rank(i)}$  in rank-based prioritization, where  $rank(i)$  is the rank of  $i$  when sorted by  $|\delta_i|$  in descending order. Results from the paper showed that DQN with PER outperformed DQN in 41 out of 49 Atari 2600 games [6].

### 3.2.4 Multi-step Q learning

Multi-step Q learning uses multi-step targets as opposed to single targets, leading to better Q value estimates [citation]. In particular, the truncated n-step return can be defined as  $R^{(n)}_t = \sum_{k=0}^{n-1} \gamma^k R_{t+k+1}$ , and the n-step target as  $Q(s, a) = R_t^n + \gamma^n \max_{a'} Q(s_{t+n}, a')$  [1].

### 3.2.5 Noisy Networks

Contrary to the inefficient state-independent exploration in epsilon-greedy, NoisyNet provides an efficient state-dependent exploration strategy without bias by adding noise into fully connected layers in the CNN in order to perturb Q-values and the corresponding policy. The output of the noisy dense layer is described by  $y = (\mu_w + \sigma_w \odot \epsilon_w)x + \mu_b + \sigma_b \odot \epsilon_b$  where  $\mu_w$ ,  $\sigma_w$ ,  $\mu_b$  and  $\sigma_b$  are trainable weights, and  $\epsilon_w$  and  $\epsilon_b$  is the Gaussian noise injected. Factorised Gaussian noise is preferred, as it requires less random number generations, thus lowering computational time. In an ideal situation, the noisy training weight parameters should converge to zero when the agent attains an optimal policy during training [8].

### 3.2.6 Distributional DQN

In reinforcement learning, it is often the case that the Q value is not a single value, but rather a distribution due to stochasticity in agent-environment interaction. In fact, this model could even be multimodal in nature. As a result, Distributional DQN models Q values using a distribution and allows the distributional Bellman equation to be satisfied by minimizing the Wasserstein metric, which measures the distance between two distributions. In particular, c51 [14] models the distribution using 51 bins, whereas QR-DQN [7] models the distribution using quantiles. In either case the output of the CNN is replaced with a Q value distribution for each action. Results show that QR-DQN outperforms c51.

### 3.2.7 Rainbow DQN

Since many algorithmic improvements are complementary in nature, the effect of combining different improvements were investigated. The result showed that combining DDQN, PER, Dueling DDQN, c51 and Noisy DQN, termed as Rainbow DQN [9], resulted in a significantly higher median human-normalized score when compared to any individual augmentation.

### 3.2.8 Count-Based Exploration and Intrinsic Motivation

Another drawback of using DQN to play games is that the algorithm requires frequent gratification in order to learn efficiently, meaning that such a limitation becomes problematic in games such as Montezuma’s Revenge, where rewards are sparse. One solution to this problem is the use of pseudo-counts and bonus rewards with DQN in order to increase the agent’s curiosity about the environment and encourage in-depth explorations. Bellemare et al. (2016) showed that the agent with exploration bonus was able to explore 15 rooms in Montezuma’s Revenge during training, as compared to 2 rooms for an agent without exploration bonus [15].

### 3.2.9 Deep Recurrent Q Network

A requirement for the use of MDP for modelling the agent-environment interaction is that the agent must be able to observe all information about the current state. However, this requirement is often not satisfied in games such as Doom where partial information is observed at each frame, resulting POMDP being used to model the agent-environment interaction instead. In order to preserve state information across states, Hausknecht and Stone [16] proposed to attach a Long Short Term Memory (LSTM) [17] to the end of DQN’s architecture in order to make the neural network recurrent. Although results showed no significant improvement over DQN in Atari 2600 games if stacked frames are used as input to the neural network, Deep Recurrent Q Networks (DRQN) are still able to play Atari 2600 games even in the presence of a flickering screen [16], and is effective for playing first player shooting (FPS) games such as Doom [11].

### 3.2.10 Other variants of Deep Q Network

Apart from the variants above, there exists many other modifications that can be found in recent literature regarding improvement on the DQN algorithm. For example, a Bayesian exploration approach can be taken by using the uncertainty Bellman equation (UBE) to obtain a tight bound for the true Q value interval. Together with the use of Thompson sampling, the paper reports better performance for UBE in 51 out of 57 games when compared to that of DQN [18]. Another example is the use of optimality tightening through constrained optimization in order for the network to converge faster to a policy [19].

### 3.3 Other Deep Reinforcement Algorithms

Apart from Deep Q Networks and its variants, other deep reinforcement algorithms have also been developed and experimented in the context of games. For example, A3C [20] and ACER [21] are actor-critic algorithms that were shown have comparable performance to DQN. Policy gradient algorithms, where the network outputs action probabilities instead of Q values, such as DDPG [22], TRPO [23] and PPO [24], can also be used to play Atari 2600 games. In the context of board games, AlphaGo was able to achieve superhuman performance either with [25] or without [26] the presence of prior training data in the game Go through a Monte Carlo Tree Search-based approach.

## 4 Related Works

### 4.1 History of Deep Reinforcement Learning

The history of using reinforcement learning and neural networks to play games dates back to the 1990's, where agents were developed to play traditional board games such as Backgammon and chess. In particular, TD-Gammon [13] was developed to play backgammon through the use of temporal difference algorithm and a multilayer perceptron network. The incorporation of handcrafted features further improved TD-Gammon's performance. However, progress in the field stagnated subsequently due to lack of success in replicating similar performances in other board games [1]. It was not until the development of Deep Q Networks by Google DeepMind that a wider variety of games could be played, notably video games such as those in Atari 2600 [12, 4]. The concept of agents playing highly complex video games was furthered by using a CNN-LSTM network architecture to play Doom, an FPS game. The recent spark in deep reinforcement learning inspired development of DRL agents in many video games, including the popular games Super Mario world and Flappy Bird. Open-source APIs for games such as Minecraft [27] and Starcraft II [28] were also developed to encourage research on the use of deep reinforcement learning in such open-ended games. Finally, OpenAI Gym [29] provides a simple platform for playing a wide variety of games or simulations, and is often used for establishing baselines for comparing DRL algorithms.

## 4.2 Mrs. Pacman

The above papers all use the Atari Learning Environment in order to play Atari 2600 games. One of the games in Atari 2600 is Ms. Pacman. Since the nature of Ms. Pacman is similar to the Pacman game used in this paper, a review of the scores is performed below.

Table 1: Reported Ms. Pacman scores from [4]

Random Play	Best linear Learner	Contingency (SARSA)	Human	DQN
307.3	1692	1227	15693	2311525

Table 2: Collation of Ms. Pacman scores from [5], [6], [2], [7], [8], and [9]

Reference	Condition	Random	Human	DQN	A3C	DDQN	Duel DDQN	PER DDQN	PER-Duel	c51	QRDQN	NoisyNet-DQN	NoisyNet-A3C	NoisyNet-Dueling	Rainbow
DDQN	no-op starts	307.3	15693.4	2311.0	N/A	3210.0	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
DDQN	human starts	196.8	15375.0	763.5	N/A	1401.8	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
						1241.3 (tuned)									
PER	human start	197.8	15375.0	763.5, 1263.0 (Gorilla), 964.7 (rank-b)	N/A	1241.3, 1865.9 (rank-b)	N/A	1824.6	N/A	N/A	N/A	N/A	N/A	N/A	N/A
Duel	no-op starts	307.3	6951.6	3085.6	N/A	2711.4	6283.5	6518.7	3327.3	N/A	N/A	N/A	N/A	N/A	N/A
Duel	human starts	197.8	15375.0	1092.3	N/A	1241.3	2250.6	1865.9	1007.8	N/A	N/A	N/A	N/A	N/A	N/A
QR-DQN	no-op starts	307.3	6951.6	3085.6	N/A	N/A	N/A	N/A	3327.3	3415	5822.5821 (huber)	N/A	N/A	N/A	N/A
NoisyNet	no-op starts	307	6952	2674±43	2436±249	N/A	3650±445	N/A	N/A	N/A	N/A	2722±148	3401±761	5546±367	N/A
Rainbow	no-op starts	N/A	N/A	1092.3	653.7	1241.3	2250.6	1824.6	N/A	2064.1	N/A	1012.1	N/A	N/A	2570.2
Rainbow	human starts	N/A	N/A	3085.6	N/A	2711.4	6235.5	4751.2	N/A	3769.2	N/A	2501.6	N/A	N/A	5380.4

Table 1 and 2 shows the raw scores obtained from evaluating various agents in Ms. Pacman using either human starts or no-op starts. From the scores, it can be seen that DQN and its variants reach higher scores generally in comparison to other DRL algorithms, such as A3C. It can also be seen that each improvement leads to better scores. In addition, it can be seen that the combination of techniques in general leads to a higher score. In particular, Rainbow, which combine five improvements to DQN, scores highest when compared to any other variation of DQN. Nevertheless, any agent’s performance does not supercede human expert performance.

However, the first DQN agent that is able to achieve super-human performance and even attain the highest score in Ms. Pacman was designed by Microsoft ???. In particular, it uses a hybrid reward architecture, where the reward function is decomposed and subsequently used to train different Q value heads in the architecture. Nevertheless, this project will not adopt such an algorithm, as decomposing the reward function involves a significant amount of reward shaping. This would likely introduce domain knowledge of the game into the algorithm, which is not desirable for the purpose of this project.

## 5 Methodology

The project is broadly separated into two phases, namely an exploration phase and an implementation phase. In particular, the objective of the exploration phase is to construct one or more deep reinforcement learning agents that can beat the game Pacman. Various RL methods are to be used in this phase, including but not limited to Deep Q Networks and its variants. The results gathered from the selected implemented algorithms are to be evaluated against benchmark algorithms, namely traditional AI agents, ie. Minimax and Expectimax.

The implementation phase involves further improving the algorithm using state-of-the-art techniques, and subsequently using it to play larger Pacman maps. Since DRL techniques provide a domain-agnostic approach for training a task, this phase also investigates how agents inter-



act with similar but unseen environments. In case of Pacman, how the agent would react to different maps of the same size.

## 5.1 Exploration phase

## 5.2 Game

Pacman, a popular arcade game, is the game studied in the exploration phase of this project. In particular, the game consists of a map containing dots, capsule, Pacman, and one or more ghosts (see Figure 3), and the objective of the game is to consume all the dots. A score is counted throughout the game, and is modulated as Pacman receives positive and negative rewards (see Table 3). This game was first studied, as the associated game implementation is open-sourced by UC Berkeley, meaning that implementing deep reinforcement algorithms would not be too challenging [3]. In terms of writing code, the game implementation also allows for map customization and control over behavior of both Pacman and ghosts, thus providing a convenient platform for algorithm testing on custom designed maps of varying difficulty.

Table 3: Score modification for various events in the game Pacman. Ghosts become scared for 40 moves when Pacman eats a capsule.

Event	Score
Win	+500
Lose	-500
Eats dot	+10
Eats capsule	+10
Collides with scared ghost	+200
Each turn	-1

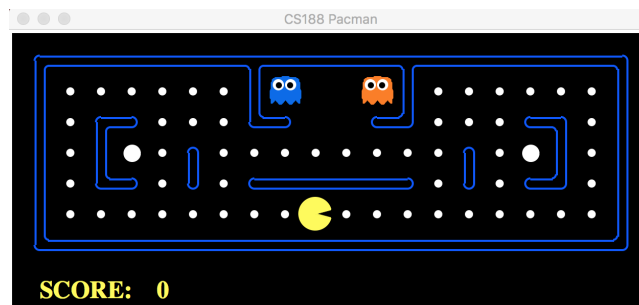


Figure 3: A Pacman map consisting of all game elements. Pacman, represented by a yellow circle with a mouth, is enclosed by walls, as shown in a blue outline. The ghosts are represented by characters with two googly eyes. The dots and capsules are represented by small and large white dots. Image produced in the Pacman software by UC Berkeley[3].

In the exploration phase, algorithmic implementations are tested on a simple map named 506Pacman (see Figure 4). In particular, the map consists of Pacman, one dot and one ghost in a  $3 \times 3$  grid, and the sprite positions are randomized uniformly each time the map is initialized. Specifically, the map's name originates from the fact that there are 506 possible game states, with 72 accessible states in one game. One of the main motivation for using this map is

its randomness, which helps to demonstrate whether the algorithm has performed generalized learning, as opposed to rote memorization of the path to a win the game. Furthermore, the implementations are also tested on a simple realistic map named smallGrid, as shown in Figure 5 in order to see whether the agent is also able to successfully learn in this map as well. The paper starts with 506Pacman and smallGrid, namely small maps, in order to make training of the DRL algorithms feasible in a reasonable amount of time.

In the implementation phase, the algorithmic implementations are trained and tested in both small and large maps, where small maps consist of 506Pacman and smallGrid, and large maps include mediumGrid and smallClassic, as shown in Figure 6.

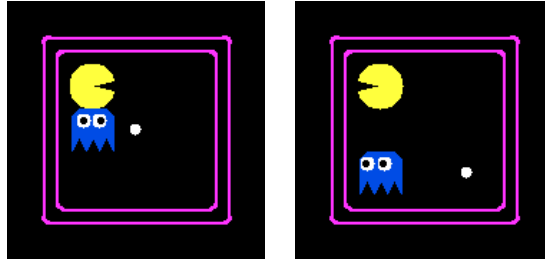


Figure 4: Examples of start game states for 506Pacman, a  $3 \times 3$  map consisting of Pacman, one dot and one ghost that are randomly placed. Image produced in the Pacman software by UC Berkeley[3].

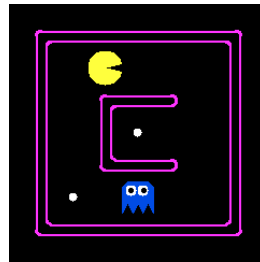
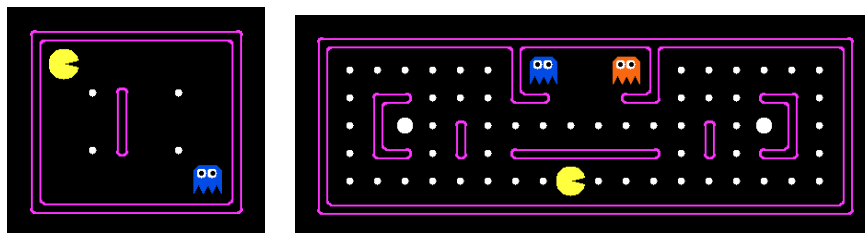


Figure 5: Map of smallGrid. The map consists of pacman, one ghost, two dots and walls in a  $7 \times 7$  map. Image produced in the Pacman software by UC Berkeley[3].



(a) mediumGrid

(b) smallClassic

Figure 6: Start game states for Pacman maps mediumGrid and smallClassic. Image produced in the Pacman software by UC Berkeley[3].

## 5.3 Exploration Phase

### 5.3.1 Deep Reinforcement Learning Algorithms

Since deep reinforcement learning is currently an active research field, there is a wide selection of algorithms to choose from. As a result, the current phase focuses on simple and effective methods that can solve Pacman maps. Since Deep Q Networks and its variants, ie. Double Deep Q Network, Dueling Deep Q Network, and Deep Q Networks with prioritized experience replay, are well studied in deep reinforcement learning, these algorithms are implemented first. In particular, they are implemented as detailed in [4, 12], [5], [2] and [6] respectively. Due to time limitations, this paper only considers proportional-based prioritized experience replay, as results from Schaul et al. indicated that the two variants of PER had similar performance [6].

### 5.3.2 Software and Hardware Requirements

Algorithms used in this paper are implemented in Python in view of its conciseness. In addition, neural networks are implemented either in Keras [30] with Tensorflow [31] backend, since Keras provides a platform for convenient and efficient prototyping and experimentation. Since the experiments are conducted in a small scale, they are run on a 2.8 GHz Intel Core i5 CPU on a Macbook Pro if possible. If the time taken to run the experiment is too long, an NVIDIA GPU can be used to speed up training.

Since the original Pacman software does not offer an API for capturing images, an API must be written in order to run the previously described algorithms. There are two approaches for capturing each frame of gameplay in Pacman, namely through direct screen capture by ImageGrab or to redraw the frame through Python Image Library (PIL). Although ImageGrab contains a simple API, it does not exactly capture the correct area properly. This is important since only pixels within the Pacman should be fed into the program. Moreover, using PIL for the frame capture implementation allows training to be done in the absence of a GUI and thereby reducing training time. As a result, this paper adopts the use of PIL for capturing gameplay frames in Python as opposed to ImageGrab.

### 5.3.3 Neural Network

In light of the effectiveness of the neural network used by Minh. et. al. in playing a wide range of Atari 2600 games[4, 12], similar architectures are adopted for DQN and its variants, as detailed in Table 4. Since the input image to the network has size  $32 \times 32$ , the filter size for the last convolutional layer is adjusted to size 2. Similar to the original paper, the RGB input image is converted into a grayscale image before feeding it into the network to reduce the size of the input layer. The parameter weights of the neural networks are also initialized uniformly at random with zero bias in order to break the symmetry of the networks and to increase the number of features learned during training.

Table 4: Description of the convolutional neural network used in the DQN, DDQN, Duel DQN and DQN with PER implementation for Pacman. The input of the network is an array of  $32 \times 32$  preprocessed grayscale images for the three implementations, and the output of the network is an array, where each element consists of five Q-values, one for each possible action from the input image’s state. The network is trained using Adam [10] as the optimizer, where the huber loss is minimized. The total number of trainable parameters in this network is 87205.

Layer	Input	Filter Size	Stride	No. of filters/Nodes	Activation	Output	No. of parameters
conv2d.1 (input)	(#samples,32,32,1)	8 x 8	4	32	Relu	(#sample, 7, 7, 32)	2080
conv2d.2	(#samples,7,7,32)	4 x 4	2	64	Relu	(#samples,2,2,32)	32832
conv2d.3	(#samples,2,2,32)	2 x 2	1	64	Relu	(#samples,1,1,64)	16448
flatten.1	(#samples,1,1,64)	-	-	-	-	(#samples, 64)	0
dense.1	(#samples,512)	-	-	512	Relu	(#samples, 512)	33280
dense.2 (output)	(#samples, 512)	-	-	5	-	(#samples, 5)	2565

### 5.3.4 Deep Q Networks and Its Variants

Since DQN, DDQN, Dueling DQN and DQN with PER are very similar value-based RL methods, they are implemented using the same training hyperparameter for ease of implementation. In particular, Table 5 shows the hand-optimized hyperparameters used in DQN and DDQN. Similar to the original publications [12, 4, 5], several trick are used to accelerate network training for both algorithms. Firstly, an experience replay buffer is used to store the last 128 steps of gameplay in order to decorrelate consecutive training data when updating the network. Such storage of game history is especially useful, since training correlated experiences induces a large variance when training a neural network [12]. Secondly, the rewards for each step is clipped to the range [-10,10] in order to increase speed of convergence to optimal Q values. Lastly, the gradient during training is clipped to [-1,1] in order to stabilize network training [12].

To train the network, Adam [10] is used as the neural network’s optimizer, and to facilitate convergence of Q values, and the Huber loss function provided by Tensorflow is used as the training loss function in order to clip the gradient of the network.

Table 5: List of DQN training hyperparameters and their values in Pacman

Hyperparameter	Value
No. of frames per input	1
Experience buffer size	2000
Batch size	256
Initial training epsilon	1
Final training epsilon	0.1
No. of exploration frames	2000
Discount factor	0.99
Target network update frequency	20
Replay start size	100
Optimizer learning rate	0.00025

### 5.3.5 Benchmark Algorithms

In this experiment, three baseline algorithms are chosen, namely Minimax, Expectimax, and DQN agent. In particular, minimax and expectimax agents are traditional optimal AI algorithms that have access to game state information. The algorithmic implementations will also be compared against an existing implementation of DQN in Pacman that uses feature extraction in order to confirm whether an end-to-end architecture is better than the use of feature extraction for learning in DQN [32]. The establishment of these benchmark algorithms serve as a medium for showing the optimality and learning progress that the deep reinforcement learning techniques have achieved.

### 5.3.6 Data Collection

All four implementations were trained for 7000 games for 506Pacman, where the epsilon value was linearly decreased from 1 to 0.1 over 2000 frames and kept at 0.1 for the rest of training. After training, the trained agents were tested on the same map for 500 games but with a constant test epsilon value of 0.05 for 506Pacman. The epsilon value was set to 0.05 in order to prevent the CNN from overfitting. To vary the map difficulty, the agent is trained and tested with random ghost and Minimax ghost behaviour.

In order to compare the effect of using map features or raw greyscale image as input data to the CNN on the agent's performance, the two versions of DQN were trained for 5000 games and tested for 1000 games on smallGrid using the same training and testing parameters as previously described.

In terms of metrics measured for both maps, the average score, percentage of games won and maximum Q values were measured over the course of training and testing. In addition, average scores and win percentages are measured as an indication of whether the agent was successful in completing the map. To gain insight on the policy learned, the maximum Q value obtained for the map shown in Figure 7 are recorded.

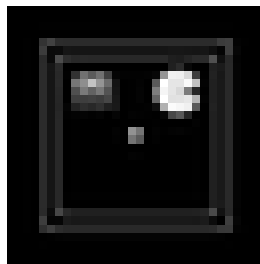
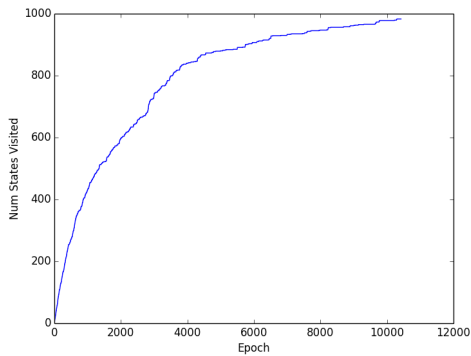


Figure 7:  $32 \times 32$  greyscale image used to observe change in Q values during training. Image produced in the Pacman software by UC Berkeley[3].

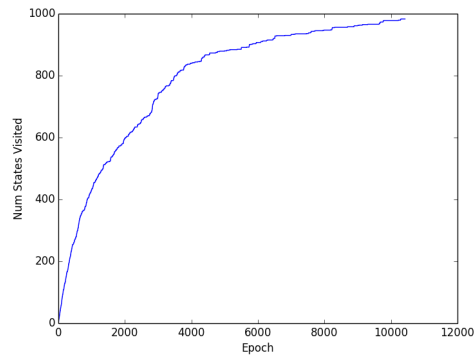
## 5.4 Implementation Phase

### 5.4.1 Map

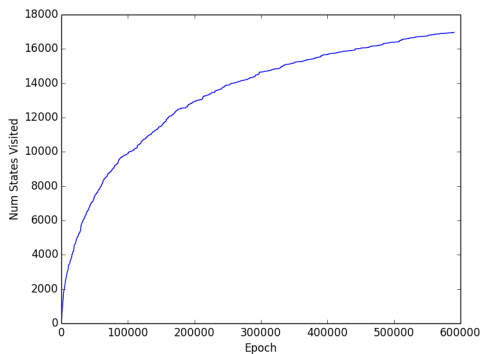
In order to see how Pacman handles tasks of varying difficulty, five maps were tested in this phase, namely 506Pacman, smallGrid, mediumGrid, miniClassic and smallClassic. For simplicity, a random ghost AI is To compare the difficulty between the maps, an agent with random was run on 506Pacman and smallGrid for 2000 games and 10000 games for mediumGrid, miniClassic and smallClassic, and the number of states visited was recorded, as shown in Figure 8. It can be seen that 506Pacman and smallGrid have approximately the same amount of states, and is fastest to converge in terms of number of states visited, reaching approximately 1000 and 800 states respectively (see Figures 8(a) and 8(b)). Figures 8(c) and (d) shows that mediumGrid and smallClassic, with approximately 20000 and more than 200000, reaches substantially more number of states. From this, it can readily be determined that the map difficulty is in the following order: 506Pacman, smallGrid, mediumGrid and smallClassic.



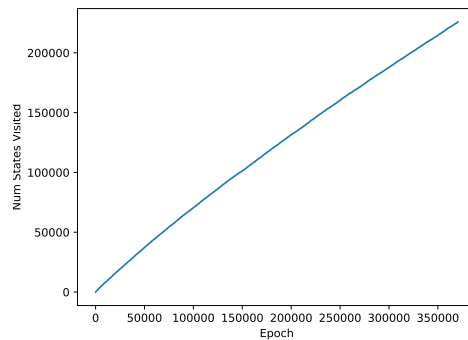
(a) 506Pacman



(b) smallGrid



(c) mediumGrid



(d) smallClassic

Figure 8: Number of game state visited over the when playing a random policy in 506Pacman, smallGrid, mediumGrid, miniClassic and smallClassic. Note that game states are also classified as different in the Pacman API if the sprites' locations are the same but the direction Pacman faces is different.

### 5.4.2 Algorithm

During the implementation phase, the algorithm was improved to include NoisyNet [citation] and Quantile Regression, thus forming a variant of Rainbow DQN, termed Combined DQN. The neural network architecture, as well as some hyperparameters were modified.

Similar to the previous phase, the neural network architecture used in this phase consists of three convolutional layers and two fully connected layers. However, the three convolutional layers uses  $8 \times 8$  filter with stride 3,  $4 \times 4$  filter with stride 3 and  $3 \times 3$  with stride 2. Such modification allows for larger image input sizes, which is required for playing larger maps in Pacman, as well as to minimize the number of network parameters. As different maps have different sizes, neural networks with different input dimensions are required for different maps. In particular, the input image size is  $64 \times 64$  for 506Pacman, smallGrid,  $64 \times 56$  mediumGrid and  $120 \times 45$  for smallClassic.

The algorithm implemented in this phase is Rainbow DQN, with a few modifications. First, c51 is replaced with QR-DQN since QR-DQN directly minimizes the Wasserstein metric and provides theoretical convergence to an optimal policy. In addition, multi-step learning is not used, since Pacman is a reward-rich game. In order to reduce training time, the CNN is trained using a batch size of 256 every 4 steps taken by the agent. The game reward function was also modified in order to speed up convergence of the Q value (see Table 6).

Table 6: Shaped reward for agent in Pacman. The rewards are kept between -1 and 1 in order to prevent from obtaining large gradients, which may hinder learning.

Event	Score
Win	+1
Lose	-1
Eats dot	1
Eats capsule	1
Collides with scared ghost	0.5
Each turn	-0.1

Since the algorithm uses quantile regression, a quantile regression huber loss function is used. For prioritized experience replay, the  $\beta$  is annealed from 0.4 to 1 over 90% of the number of training games, where  $\beta$  is annealed per game episode [6]. Since Pacman is a simple game, 10 quantiles were used to model the Q value distribution [7]. In terms of the noisy dense layers used in the CNN, factorized normal distribution was used to reduce computational cost and sigma was set to 0.5. Note that NoisyNet replaces epsilon-greedy algorithm, meaning that the agent's policy is greedy with respect to noisy Q-values and does not involve making random actions [8]. Table 7 summarizes the hyperparameters used during training.

Table 7: List of Combined DQN training hyperparameters and their values in Pacman

Hyperparameters	value
No. of frames per input	1
Experience buffer size	50000
Batch size	256
No.of steps per training step	4
Discount	0.95
Target network update frequency	20
Replay start size	10000
Adam learnin grate	0.0042
Adam epsilon	0.00001953125
PEP alpha	0.6
PEP initial beta	0.4
QR no. of targets	10
NoisyNet sigma	0.5

### 5.4.3 Experiment

The experimentation in this phase is separated into two parts. Similar to the exploration phase, the first part involves training the algorithm for sufficient amount of games until the agent's policy converges. The agent is subsequently testing for 500 games on the Pacman maps without initial human or random actions. To track the agent's performance during training, the win rate, game scores, loss function and Q values are recorded. During testing, the average score and win rate is recorded. To further examine the policy generated from training, Q values of individual states are also recorded.

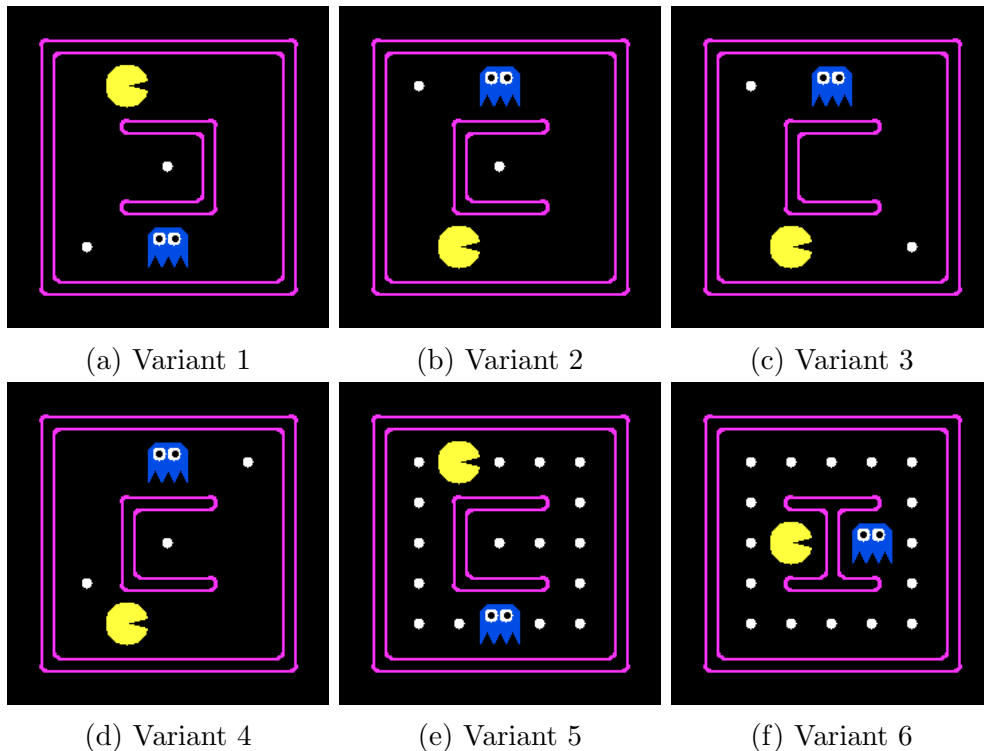


Figure 9: Deep Q Networks training results with default ghost AI in 506Pacman



The second part involves testing the developed agents from the first part on different maps or conditions. Specifically, a Combined DQN agent trained using 506Pacman is tested using variants of the 506Pacman map, where either the number of dots in the map or the ghost AI is altered. To determine how the agent interacts with an unfamiliar environment, a trained smallGrid agent is tested on six variants of smallGrid, as shown in Figure 9. Note that the map used during training and testing must have identical map dimensions, as the sprites' sizes would change, which would disadvantage the network.

## 6 Results and Discussion

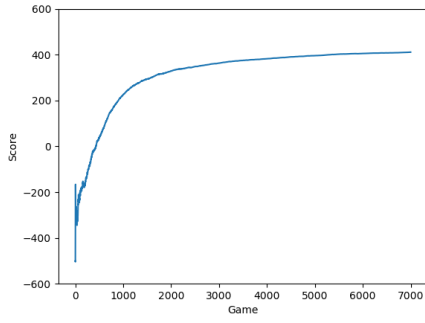
### 6.1 Exploration Phase

Figures 10 and 11 show results from training in 506Pacman using default and minimax ghost AI respectively, Figure 12 shows the results obtained from training in smallGrid, Figure 13 show the training losses for training in the two maps and Figure 14 shows the training loss when using DQN with proportional-based prioritization. In particular, the metrics are plotted either against training games or training epochs, where a single epoch corresponds to one neural network batch update in the algorithms.

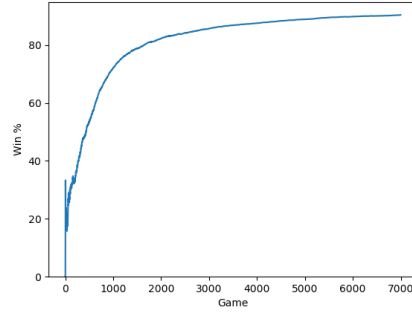
From Figures 10(a), 10(b), 11(a), and 11(b), a gradual increase in average score and win rate is observed, indicating that the agent does perform learning during training. This can be confirmed in Figures 10(c) and 11(c), where the maximum Q value converges to a value close to 10 during training. During training, DQN won 6330 and 6141 out of 7000 games, ie. a win rate of 90.4% and 87.7%, in 506Pacman with default and Minimax ghost behaviour respectively. The amount of training steps taken were around 16000-17000, taking approximately 3-4 hours to train. During testing, DQN won 455 and 451 out of 500 games (91% win rate) for default and ghost AI respectively. These scores are lower than the win rate of a Minimax or Expectimax agent in 506Pacman, which achieved over 98% win rate for both ghost behaviors, implying that the policy attained by DQN is not optimal.

It is interesting to note that by analysing the Q values obtained with the image in Figure 7 as input, it is observed that the action corresponding to the maximum Q value in 10(d) and 11(d) is south, which is consistent with what a human player would do, since moving south in the map would increase Pacman's distance from the ghost and decrease Pacman's distance from the dot.

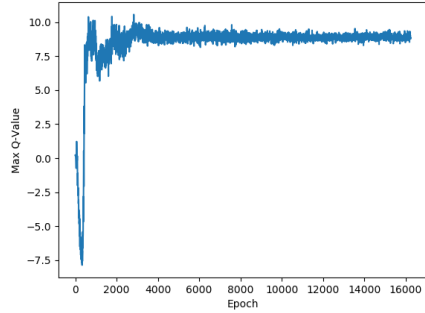
In addition, Figure 12 shows the results obtained from running DQN in smallGrid. Similar to 506Pacman, Figures 11(a) and 11(b) show an increase in average score and win rate respectively, which also shows that the agent has converged to a winning strategy in smallGrid. In particular, 4007 out of 5000 games (80.4%) was achieved during training, and 891 out of 1000 games (89.1%) was achieved during testing. This is in contrast to the results of the feature extraction DQN implementation from [32], where 1267 out of 5000 and 499/1000 games were won during training and testing respectively. From the results above, it can be concluded that . The explanation for such an observation may be due to extra information encoded by the input pixel image that is not conveyed via feature extraction, such as the direction that Pacman is facing, thereby allowing better policies to be made.



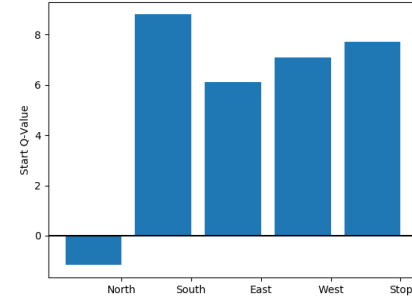
(a) Average scores during training



(b) Win percentage during training

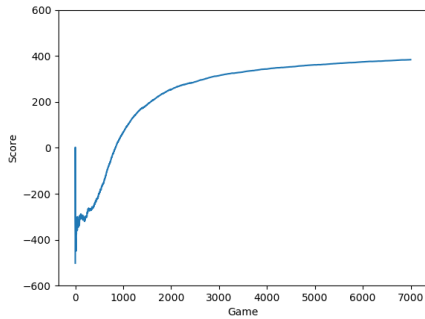


(c) max Q values for a specific start state during training

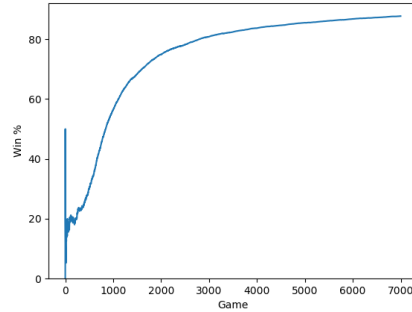


(d) start Q values from map in Figure 7

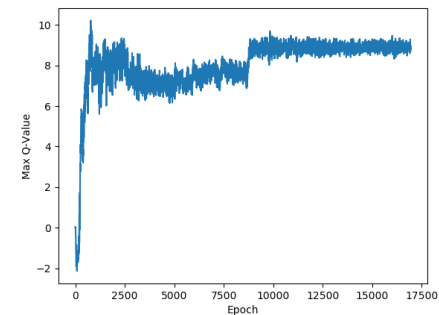
Figure 10: Deep Q Networks training results with default ghost AI in 506Pacman



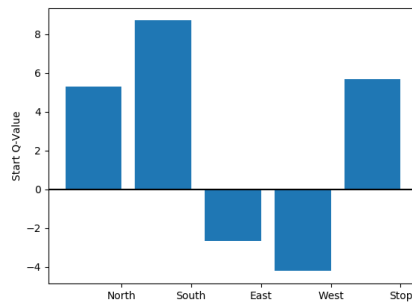
(a) Average scores during training



(b) Win percentage during training

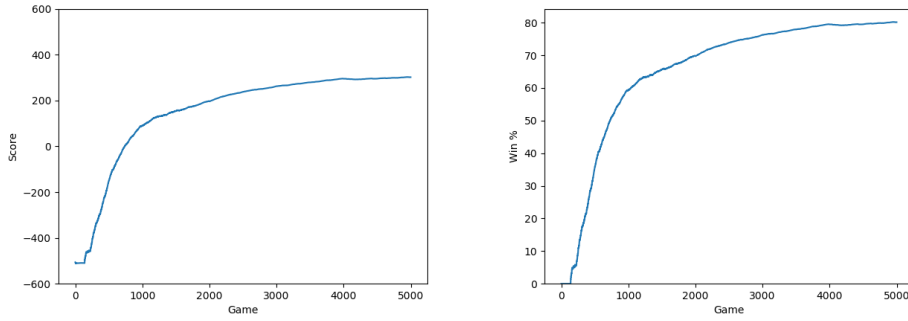


(c) max Q values for a specific start state during training



(d) Last training game's start Q values

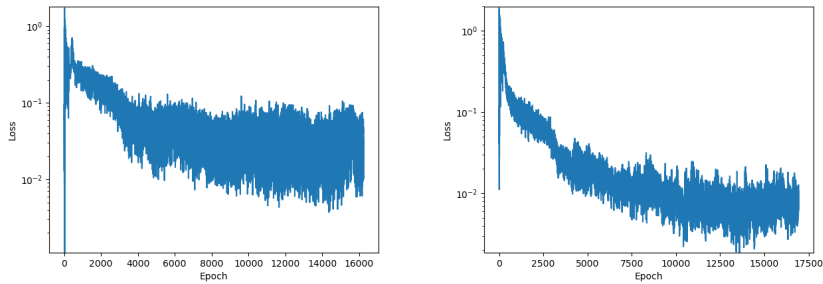
Figure 11: Deep Q Networks training results with minimax ghost AI in 506Pacman



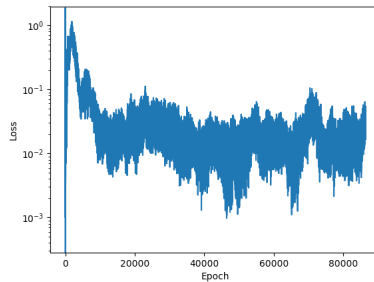
(a) Average scores during training      (b) Win percentage during training

Figure 12: Deep Q Networks training results with default ghost AI in 506Pacman

Figure 13(a), 13(b) and 13(c) show the training losses during training in smallGrid and 506Pacman. All three figures show a decline in training losses, indicating learning performed by the agent. However, it is worthwhile to note that the decline stops with training loss at around  $10^{-2}$ , indicating that the agent does not converge to an optimal policy. The training loss not converging to zero may be an indication of the CNN either under-fitting. As a result, further hyper-parameter tuning may be required in order to further improve the agent's performance.



(a) Training loss in 506Pacman with default ghost behavior      (b) Training loss in 506Pacman with default ghost behavior



(c) Training loss in smallGrid

Figure 13: Deep Q Networks training losses in Pacman

Apart from performing the aforementioned experiments, there were many preliminary experimentation performed during the process. In particular, there were also attempts to use Double DQN and Duel DQN, as well as combination of the DQN variants in playing 506Pacman. However, the subsequent results indicated that the agent often scored below -550 after training for approximately 1000 games, indicating that the agent learned how to avoid the ghost but did not learn to eat the dot in order to win the game. This may be likely due to the CNN having trouble distinguishing between the ghost and dot sprite. It is also observed that the use of prioritized experience replay leads to a faster reduction in training loss, as well as a substantial decrease in training loss to the magnitude of  $10^{-6}$ , as shown in Figure 14. Moreover, testing of the resulting agent results in a 96% win rate, which agrees with the. Therefore, it can be concluded that DQN with proportional-based PER leads to a better agent than DQN, which is in line with the literature [6].

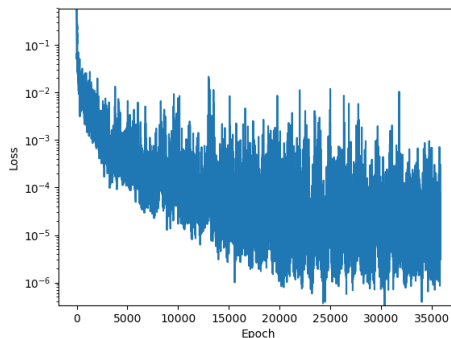


Figure 14: Training loss of DQN with proportional-based PER in 506Pacman for 6000 games

## 6.2 Implementation Phase

Figure 15 shows the results obtained after training for 10000 games on 506Pacman using random ghost AI. From Figure 15, it can be seen that the agent has converged to a policy to a score over 500 (Figure 15(a)). Figures 15(b), 15(c) and 15(d) shows that the policy is converging, although the loss function (Figure 15(c)) shows that the convergence is not complete. Nevertheless, during testing, the agent obtains an average score of 507.67 and a win rate of 100%, indicating an optimal policy is attained. It is interesting to note that the policy converges early on after approximately 4000 games, but still makes errors infrequently (see figure (b)).

Figure 16 shows the data obtained after training for 4000 games on smallGrid using random ghost AI. As shown in Figure 16(e), the Q value function has not converged yet, indicating that the agent has not completed learning. As a result, during testing, the agent attains an average score of -21.35 and 48% win rate. During testing, a spike in number of states visited was observed at around epoch 40000, which is likely due to the agent being able to consistently eat one dot in the map (see Figure ??). Similar to 506Pacman, informal experimentation suggests that the converged agent will attain 100% after approximately 7000 games, although testing has yet to be done due to time limitaiton.

Figure 17 shows the data obtained after training 2500 games on mediumGrid using random Ghost AI. The agent is able to converge to an optimal policy, attaining 478.34 and 95.2% win rate during testing (see Figure 17(b)). Similar to the previous two maps, the Q value

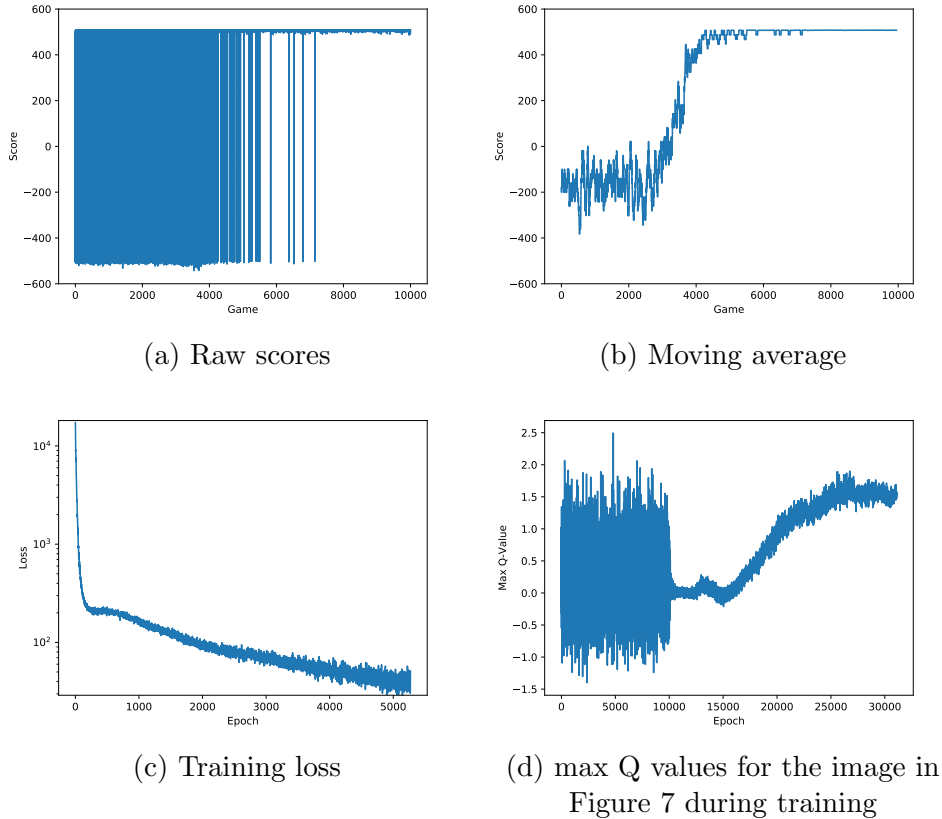


Figure 15: Results obtained from playing 506Pacman with random ghost AI for 10000 games

function converges (see Figures 17(c) and 17(d)), but not completely yet. Upon informal experimentation, it was noted that an agent with 100% win rate can be attained by training for at least 7000 games.

Figure 18 shows the data obtained after training 1000 games on smallClassic using random Ghost AI. It can be observed that the agent has reached a suboptimal policy (see Figures 18(a) and 18(b)) even though the training loss and Q value functions appear to have converged (see Figures 18(c) and 18(d)). During testing, the average score and win rate obtained are -261.84 and 0% respectively. A reason for the long time taken for the agent's policy to converge is due to the large number of dots on the map. Such a large amount forces the agent to perform implicit path planning, which DQN and its variants are not very good at currently [9].

In summary, although Combined DQN is able to attain high scores and win rates in general, the number of steps taken to attain such a result increases substantially as the map difficulty increases.

In general, the training weights associated with noise decreases as the agent learns (see Figure 19). However, the weights do not decrease to a value near zero even after an optimal policy is attained in the respective maps. This means that the policy relies on a degree of randomness even after an optimal policy was attained. Alternatively, the non-zero weights could be due to the fact that the amount of noise does not affect which action the agent chooses.

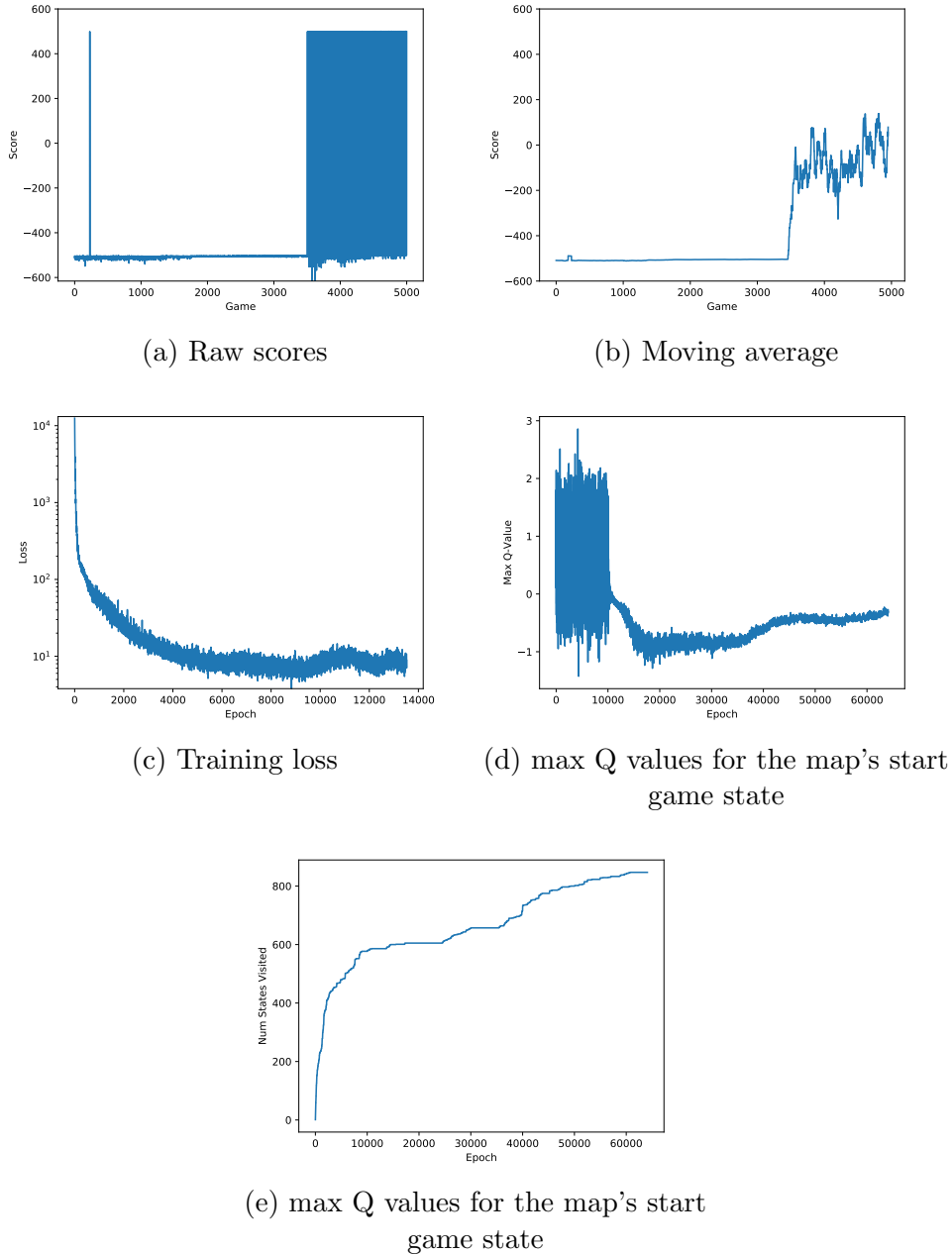


Figure 16: Results obtained from playing smallGrid with random ghost AI for 5000 games

### 6.2.1 Generalization

To investigate the extent to which trained agents can generalize, some trained agents from the previous part were tested on maps with different conditions. Table 8 shows the win rate and average score after testing a 506Pacman Combined DQN agent on a 506Pacman map with either random, direct, Minimax or Expectimax ghost behavior. In particular, it is observed that changing the ghost behavior does not hurt the agent's performance substantially. The consistent high win rate is likely due to the random ghost AI, which promotes the agent to pick an action based on the worst-case move that the ghost could pick. As a result, this phenomenon prevents the agent from losing games even if the ghost acts more aggressively.

Table 9 shows how the agent's performance is affected from increasing the number of dots in

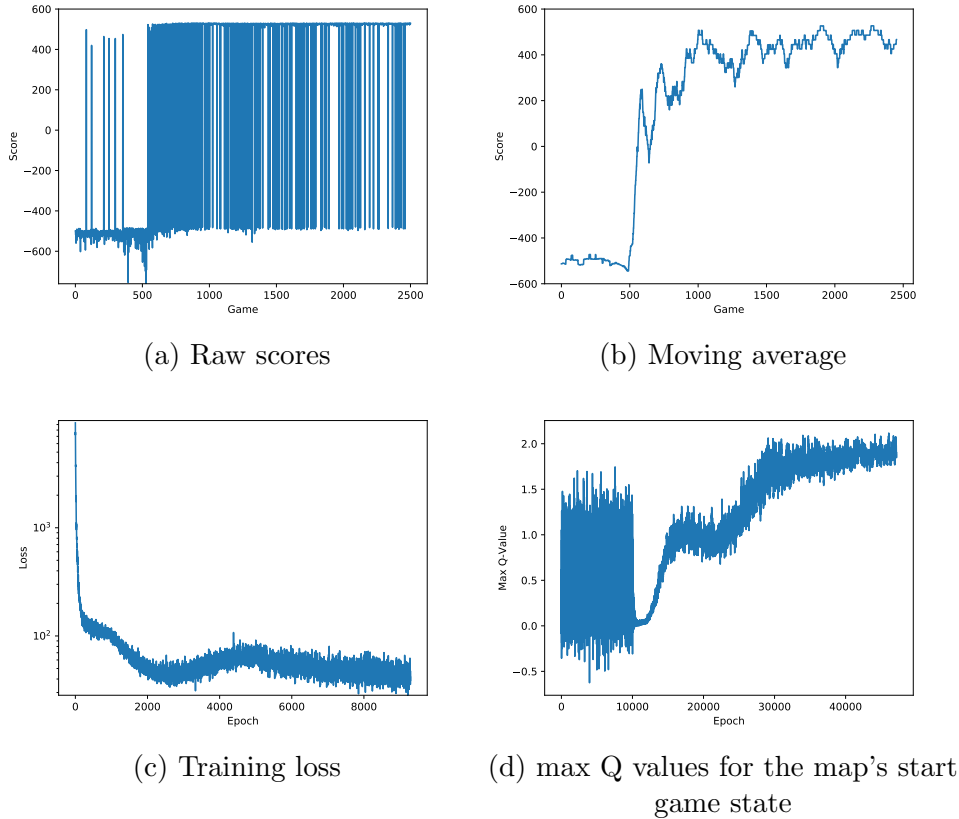


Figure 17: Results obtained from playing mediumGrid with random ghost AI for 5000 games

Table 8: Win rate and average score obtained after running a trained 506Pacman agent using different ghost AIs in 506Pacman.

Ghost AI	Win Rate (%)	Average Score
random	100	507.67
direct	100	507.52
Minimax	100	507.27
Expectimax	100	507.60

the 506Pacman map. In general, the win rate and average score decreases as the number of dots increases. The main reason to this is likely because Pacman did not learn to perform path finding during training, which leads it to act impulsively during testing.

Table 10 shows the results obtained from running a smallGrid agent on variations of the small-Grid map. In general, the performance is quite poor, although maps containing characteristics that resemble the original map tend to have higher scores and win rates. In particular, the agent in the first map only eats one dot, and even gets stuck in the lower left corner of the map in the third variant. From such observations, it appears that the policy attained from training on one map only is unable to generalize to unfamiliar situations.

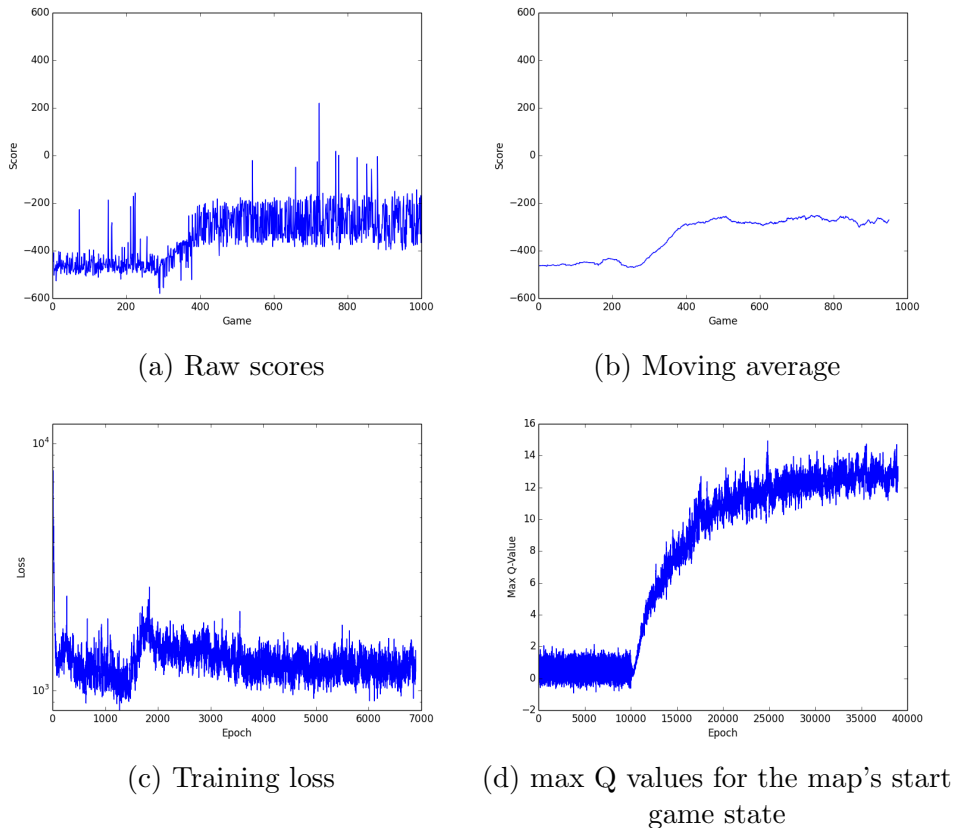


Figure 18: Results obtained from playing smallClassic with random ghost AI for 10000 games

Table 9: Win rate and average score obtained after running a trained 506Pacman agent in 506Pacman containing different number of dots.

Number of dots	Win Rate (%)	Average Score
2	93.4	448.18
3	85.6	375.96
4	82.4	351.06
5	68.8	217.56
6	60.2	135.32

### 6.3 Limitation and Difficulties Encountered

\*\*hyperparameter tuning completed \*\*bottleneck on computation instead of Pacman API now ( eg. computation of target values for training)

Several difficulties were encountered during implementation and experimentation of the algorithms used in this paper. In particular, it was challenging to optimize the algorithms, since there seems to be no fast and systematic method of optimizing the hyperparameters for the above algorithms. This was especially the case for tuning epsilon since a balance between exploration and exploitation is required in order for the reinforcement agents to learn. Since hand optimization of hyperparameters is expensive and time consuming, the project progress was slowed by a significant amount. Since little literature regarding this issue has been read, this problem may have to be addressed in the future through further literature search.



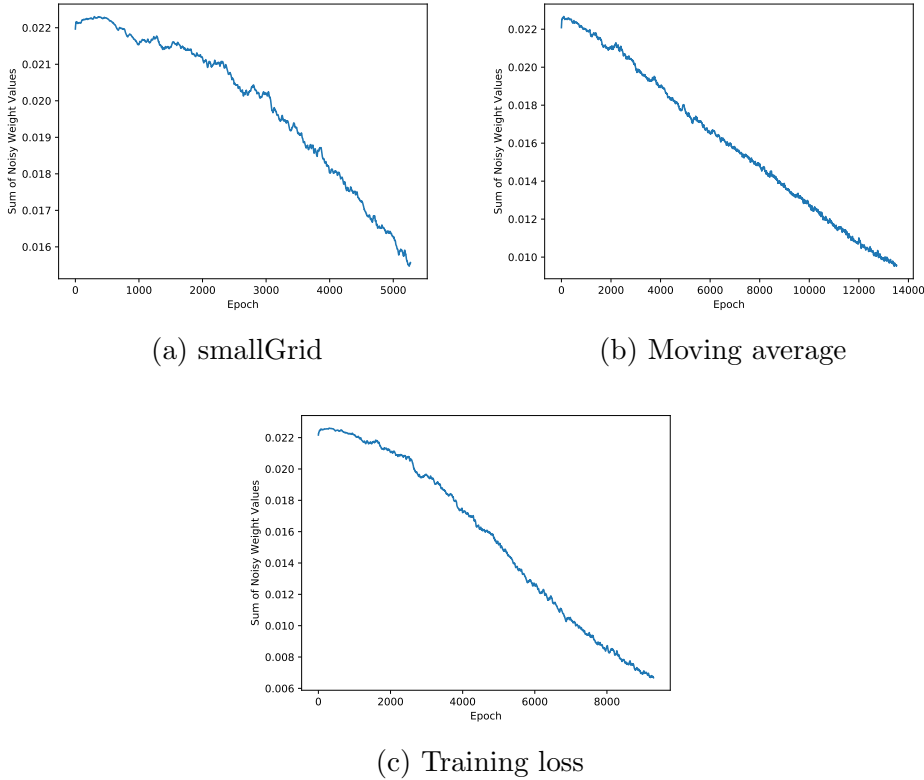


Figure 19: Sum of noisy training parameters during training in 506Pacman, smallGrid and mediumGrid. No data was gathered due to the policy having 0% win rate for smallClassic.

Table 10: Results from testing smallGrid agent on different variants of smallGrid

Variant	win rate	average score
smallGrid_var1	0	-510.41
smallGrid_var2	42	-82.37
smallGrid_var3	0	-510.37
smallGrid_var4	34	-159.46
smallGrid_var5	0	-481.46
smallGrid_var6	0	-492.92

Apart from hyperparameter tuning, another problem was that the values of the hyperparameters were dependent on the map being played. As a result, only two map was studied thus far. Nevertheless, similar experiments involving much larger maps with more ghosts are to be performed at a later stage in the exploration phase.

In terms of limitations of the algorithms, one of the problems is that an increase the state-action space causes the deep reinforcement learning agent to take longer to converge to an optimal policy. Such a characteristic of deep reinforcement learning algorithms can be problematic, due to time and resource constraints. Even when a GPU is used, training time can still be lengthy due to image drawing and processing being the bottleneck. As a result, further investigations regarding such bottleneck may be required in order to minimize the use of computational resource and time in the future. One possible way to do this would be to parallize to, although a thread-safe implementation of DQN or Combined DQN would be rather tricky.

Numerous technical problems appeared when implementing the algorithms. One such issue was that the provided Pacman software did not provide a direct API method for observing the visual output of a game state. As a result, unexpected time was spent on implementing such a feature in addition to the original software, which slowed the project's progress by a small amount.

Another challenge was on the mechanism when handling illegal moves. Since the Pacman API does not accept illegal moves, the algorithms were programmed so that any illegal action would cause Pacman to execute the stop action. Nonetheless, the intended action that Pacman wanted to take is recorded in the experience replay buffer so that the action taken was propagated back to the associated nodes in the network.

## 7 Future Work

In view of the major obstacle being the inability for the deep reinforcement learning agent to approach 100% win rate, the next step of the project involves the modification of the existing implementation to mitigate this issue. For example, additional experiments may be devised in order to determine the impact of neural architecture on DQN for a better understanding of the optimal architecture to use in deep reinforcement learning. Apart from tackling with the problems previously mentioned, various benchmark algorithms may also be explored in order to better compare experimental results obtained. Lastly, experiments conducted in this paper will also be repeated with Pacman maps of higher complexity in order to better compare the computational differences between the studied deep reinforcement learning techniques.

Since the survey of deep reinforcement learning algorithms on different games is still incomplete, it is rather difficult to plan in a detailed manner in the exploration phase. Nevertheless, it is projected with a high probability that the project will do some kind of AI agent implementation in either Minecraft or StarCraft, since Microsoft and Deepmind have respectively provided an API for programmers to create agents in these games as part of an initiative to encourage community implementation and exploration of deep reinforcement learning through such games [27, 28].

There are many things that can be done for this project. First, the DQN agents developed in this project are unable to play large-sized maps, such as the originalClassic map in Ms. Pacman. In addition, the current agents have the limitation of only being able to play in maps of the same size as that during training due to a fixed input in the CNN. As a result, it would be interesting to see whether a general Pacman agent can be developed that can play maps of varying sizes. Lastly, one limitation of current DQN algorithms is that the agent lacks imagination or curiosity, which may help in ameliorating the agent in exploring the map to find rewards that are more sparse.

The project is completed, and is summarized in Table 11, where bold items represent strict deadlines and deliverables. Note that significant amendments were made over the course of the project in view of the change of the project's scope to narrow the project's scope to only consider DQN and its variants. In terms of met deadlines, all items up to Apr 15, 2017 are currently completed.

Table 11: Project Schedule

Date	Task
Early October 2017	Preliminary Research <ul style="list-style-type: none"> <li>• Perform preliminary literature research on existing games that use reinforcement learning and deep learning techniques.</li> <li>• Experiment with various classic RL methods, including policy iteration, value iteration and tabular Q-Learning</li> </ul>
October 1, 2017	<b>Phase 1 Deliverables (Inception): Project Scheme, Detailed Project Plan</b>
Mid October 2017	Stage 1: Algorithmic Exploration <ul style="list-style-type: none"> <li>• Implement and evaluate the performance of DQN in Pacman</li> <li>• Read up on variants of DQN algorithm</li> </ul>
October 31 2017	<ul style="list-style-type: none"> <li>• Read up on variants of DQN algorithm</li> <li>• Implement DDQN, Duel DQN and DQN with proportional-based PER</li> </ul>
Mid November 2017	<ul style="list-style-type: none"> <li>• Implement the 506Pacman map</li> <li>• Experiment with DRL algorithms in different Pacman maps</li> </ul>
December 2017	<ul style="list-style-type: none"> <li>• Train and Test DQN and its variants on 506Pacman and establish results</li> </ul>
January 22, 2018	<b>First Presentation</b>
January 29, 2018	<b>Phase 2 Deliverables (Elaboration): Pacman RL Implementations, Interim Report</b>
Jan –April 2018	Stage 2: RL implementation and Optimization <ul style="list-style-type: none"> <li>• Improve upon existing reinforcement learning algorithms</li> <li>• Play on larger Pacman maps</li> </ul>
April 15, 2018	<b>Phase 3 (Construction): Game RL Implementation, Final Report</b>
April 19, 2018	<b>Final Presentation</b>
May 2, 2018	<b>Project Exhibition</b>

## 8 Conclusion

This project aims to explore the applicability of deep reinforcement learning methods in the context of playing games. Starting with an implementation of DQN in the exploration phase, the project saw success in attaining a high score and win rate in the simple map 506Pacman. However, it was determined that the current paper’s DQN implementation is not optimal, and further work is required to determine whether it can be further improved. In terms of network architecture, it was confirmed that the use of an end-to-end CNN architecture is more advantageous than the use of feature extraction and a CNN in the context of smallGrid. Several observations were also made, leading to additional motivation to improve the use of algorithms such as DDQN and DuelDQN in Pacman. Nevertheless, the exploration phase saw development of Combined DQN, which led to significantly better

Despite of the limited selection of techniques and games explored in this paper, steps are to be taken in order to explore other techniques and tricks that ameliorate learning of a task in a wide variety of games, such as Atari 2600 games, in the exploration phase. This is an important part of the project in view of the rapid progress in the field, which could easily render previous surveys outdated. With proper comparisons and investigation completed in the future, the implementation phase will subsequently use the results to make informed decisions on the selection of algorithms to use and apply in a sophisticated video game or to gain insight on drawbacks of existing algorithms and subsequently improvements to them.

## 9 References

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning : An Introduction*. MIT Press, 1998.
- [2] Z. Wang, N. de Freitas, and M. Lanctot, “Dueling network architectures for deep reinforcement learning,” *CoRR*, vol. abs/1511.06581, 2015. [Online]. Available: <http://arxiv.org/abs/1511.06581>
- [3] (2017, Oct). [Online]. Available: <http://ai.berkeley.edu/reinforcement.html>
- [4] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 02 2015. [Online]. Available: <http://dx.doi.org/10.1038/nature14236>
- [5] H. van Hasselt, A. Guez, and D. Silver, “Deep reinforcement learning with double q-learning,” *CoRR*, vol. abs/1509.06461, 2015. [Online]. Available: <http://arxiv.org/abs/1509.06461>
- [6] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, “Prioritized experience replay,” *CoRR*, vol. abs/1511.05952, 2015. [Online]. Available: <http://arxiv.org/abs/1511.05952>

- [7] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos, “Distributional reinforcement learning with quantile regression,” *CoRR*, vol. abs/1710.10044, 2017. [Online]. Available: <http://arxiv.org/abs/1710.10044>
- [8] M. Fortunato, M. G. Azar, B. Piot, J. Menick, I. Osband, A. Graves, V. Mnih, R. Munos, D. Hassabis, O. Pietquin, C. Blundell, and S. Legg, “Noisy networks for exploration,” *CoRR*, vol. abs/1706.10295, 2017. [Online]. Available: <http://arxiv.org/abs/1706.10295>
- [9] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, “Rainbow: Combining improvements in deep reinforcement learning,” *CoRR*, vol. abs/1710.02298, 2017. [Online]. Available: <http://arxiv.org/abs/1710.02298>
- [10] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *CoRR*, vol. abs/1412.6980, 2014. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [11] G. Lample and D. S. Chaplot, “Playing FPS games with deep reinforcement learning,” *CoRR*, vol. abs/1609.05521, 2016. [Online]. Available: <http://arxiv.org/abs/1609.05521>
- [12] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” in *NIPS Deep Learning Workshop*, 2013.
- [13] G. Tesauro, “Temporal difference learning and td-gammon,” *Commun. ACM*, vol. 38, no. 3, pp. 58–68, Mar. 1995. [Online]. Available: <http://doi.acm.org/10.1145/203330.203343>
- [14] M. G. Bellemare, W. Dabney, and R. Munos, “A distributional perspective on reinforcement learning,” *CoRR*, vol. abs/1707.06887, 2017. [Online]. Available: <http://arxiv.org/abs/1707.06887>
- [15] M. G. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos, “Unifying count-based exploration and intrinsic motivation,” *CoRR*, vol. abs/1606.01868, 2016. [Online]. Available: <http://arxiv.org/abs/1606.01868>
- [16] M. J. Hausknecht and P. Stone, “Deep recurrent q-learning for partially observable mdps,” *CoRR*, vol. abs/1507.06527, 2015. [Online]. Available: <http://arxiv.org/abs/1507.06527>
- [17] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997. [Online]. Available: <http://dx.doi.org/10.1162/neco.1997.9.8.1735>
- [18] B. O’Donoghue, I. Osband, R. Munos, and V. Mnih, “The uncertainty bellman equation and exploration,” *CoRR*, vol. abs/1709.05380, 2017. [Online]. Available: <http://arxiv.org/abs/1709.05380>
- [19] F. S. He, Y. Liu, A. G. Schwing, and J. Peng, “Learning to play in a day: Faster deep reinforcement learning by optimality tightening,” *CoRR*, vol. abs/1611.01606, 2016. [Online]. Available: <http://arxiv.org/abs/1611.01606>
- [20] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and

- K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” *CoRR*, vol. abs/1602.01783, 2016. [Online]. Available: <http://arxiv.org/abs/1602.01783>
- [21] Z. Wang, V. Bapst, N. Heess, V. Mnih, R. Munos, K. Kavukcuoglu, and N. de Freitas, “Sample efficient actor-critic with experience replay,” *CoRR*, vol. abs/1611.01224, 2016. [Online]. Available: <http://arxiv.org/abs/1611.01224>
- [22] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, “Continuous control with deep reinforcement learning,” *CoRR*, vol. abs/1509.02971, 2015. [Online]. Available: <http://arxiv.org/abs/1509.02971>
- [23] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, “Trust region policy optimization,” *CoRR*, vol. abs/1502.05477, 2015. [Online]. Available: <http://arxiv.org/abs/1502.05477>
- [24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *CoRR*, vol. abs/1707.06347, 2017. [Online]. Available: <http://arxiv.org/abs/1707.06347>
- [25] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, pp. 484–489, January 2016.
- [26] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, “Mastering the game of go without human knowledge,” *Nature*, vol. 550, pp. 354 EP –, 10 2017. [Online]. Available: <http://dx.doi.org/10.1038/nature24270>
- [27] M. Johnson, K. Hofmann, T. Hutton, and D. Bignell, “The malmo platform for artificial intelligence experimentation,” in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence, IJCAI 2016, New York, NY, USA, 9-15 July 2016*, 2016, pp. 4246–4247. [Online]. Available: <http://www.ijcai.org/Abstract/16/643>
- [28] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schrittwieser, J. Quan, S. Gaffney, S. Petersen, K. Simonyan, T. Schaul, H. van Hasselt, D. Silver, T. P. Lillicrap, K. Calderone, P. Keet, A. Brunasso, D. Lawrence, A. Ekermo, J. Repp, and R. Tsing, “Starcraft II: A new challenge for reinforcement learning,” *CoRR*, vol. abs/1708.04782, 2017. [Online]. Available: <http://arxiv.org/abs/1708.04782>
- [29] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai gym,” *CoRR*, vol. abs/1606.01540, 2016. [Online]. Available: <http://arxiv.org/abs/1606.01540>
- [30] F. Chollet *et al.*, “Keras,” <https://github.com/fchollet/keras>, 2015.
- [31] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado,

A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, “TensorFlow: Large-scale machine learning on heterogeneous systems,” 2015, software available from tensorflow.org. [Online]. Available: <https://www.tensorflow.org/>

[32] T. van der Ouderaa, “Deep reinforcement learning in pac-man,” 2016.