



INVESTIGATING STATE REPRESENTATIONS IN DEEP REINFORCEMENT LEARNING FOR PELLET EATING IN AGAR.IO

Bachelor's Project Thesis

Nil Stolt Ansó, s2705338, nilstoltanso@gmail.com,
 Supervisor: Dr M. Wiering

Abstract: The online game Agar.io has become massively popular on the internet due to its intuitive game design and its ability to instantly match players with others around the world. The game has a continuous input and action space and allows to have diverse agents with complex strategies compete against each other. This paper first investigates how different state representations influence the learning process of a Q-learning algorithm. The representations examined range from raw pixel values to extracted handcrafted feature vision grids. Secondly, we investigate how different value function network architectures compare in performance. The architectures examined are two convolutional Deep Q-networks (DQN) of varying depth and one smaller multilayer perceptron (MLP). The results show that the Q-learning algorithm, together with prioritized experience replay, is able to play quite well. Handcrafted feature vision grids seem to require minimal resolution and network complexity, and outperform raw pixel input for all conditions and tasks tested.

1 Introduction

Reinforcement learning (RL) is a machine learning paradigm which uses a reward function that assigns a value to a specific state an agent is in as a supervision signal [1]. The agent attempts to learn what actions to take in an environment to maximize this reward signal. These agents are usually trained in simulations or games. This is for multiple reasons: the level of noise can be directly controlled, the researcher has access to all relevant information, and the simulation can be sped up and parallelized. Once agents achieve optimal performances in complex and noisy simulated environments, they can also be employed for real-world tasks.

The environment for this research is based on the game of `Agar.io`, which is itself inspired by the behaviour of biological cells. The player controls circular cells in a 2D plane (as if laid out on a Petri dish) which follow the player's mouse cursor. The player can signal their cells to split or eject little vesicles of mass. Cells of a player can eat small food pellets scattered in the environment or other smaller enemy player-controlled cells to grow in size. The

environment of Agar.io is therefore very interesting for RL research, as it is mainly formed by the behavior of other (larger) players on the same plane, it is stochastic and constantly changing. Also the output space is continuous, as the cells of the player move towards the exact position of the mouse cursor. On top of that, the complexity of the game can be scaled by introducing or removing additional features. This paper therefore studies how to use RL to build an intelligent agent for this game, especially focusing on how to represent the game state for the agent.

1.1 Previous Research

Ever since it was shown that a multilayer perceptron (MLP) could be trained through back-propagation to store internal representations of the provided input [2], the use of artificial neural networks (ANN) has demonstrated great promise at learning representations of complex environments. Given a large enough hidden layer, an ANN has been shown to be able to approximate any continuous function on the input space to any degree of accuracy [3].

Tesauro was among the first to show that optimal decision making could be learned in the large state space of the game Backgammon through the use of an MLP as a decision maker [4]. Over the years, this principle has been extended through the use of convolutional neural networks (CNNs). This has been shown to achieve human level performance by learning from solely pixel values in a variety of Atari games [5], and even first-person perspective 3D games like Doom [6]. Despite much of the success of deep learning coming with great computational requirements, the combination of deep learning and reinforcement learning has been finding increasing successes [5] [7] [6].

One approach to overcoming the issue of large state spaces is the preprocessing of the state representation in order to extract features that will boost performance and reduce the amount of potentially irrelevant information required for the network to process. The use of vision grids is one such approach that has been employed in games such as Starcraft [8] and Tron [9] by extracting hand-crafted features into grids. These methods can greatly simplify the state space and allow for a decreased network complexity. Despite this benefit, feature extraction might introduce biases and has no guarantee to achieve the same performance as a network being fed the raw game representation (given enough training time).

A widely successful algorithm for reinforcement learning that acts on state representations is Q-Learning [10]. This algorithm can be combined with a function approximator to estimate the Quality, or long-term reward prospect, of a state-action pair. This estimation is improved by observing the actual reward received after taking a specific action in a specific state. By combining the value of the reward received and the value of the discounted estimation of the best state-action pair in the new state, the algorithm slowly propagates expected reward values backwards through the state-action space. Which action to take in a given state is chosen by taking the action with the highest predicted Q-value. Such approaches are used in the research mentioned above on Atari games [7], Doom [6], and Tron [9]. In Atari games, in Doom, and in Tron the possible actions in each state are equivalent to the buttons that the player can press. In Agar.io, the relative position of the mouse cursor on the screen is used to direct the player. This has a range of continuous

values, similarly to real-life robotic actuators.

1.2 Contributions of this Paper

This paper explores how the complexity of the state space affects the convergence and final performance of the algorithm. This is explored through a core task of the game: pellet collection. In this task, the agent has to navigate in the environment and eat as many food pellets as possible.

More importantly, this research focuses on how different state representations, varying resolutions of such state representations, and varying the structure for the function approximators affect the performance of the Q-learning algorithm. First, different kinds of low-level information used in the state representation are compared, each one providing a different kind of information. This includes grayscale pixel values, RGB pixel values, and a semantic vision grid for pellets in the screen. The effect of the resolution of these state representations is also explored. Finally the ability of Q-learning to achieve a good playing performance is examined when using two different CNN structures (which differ in the number of layers) and compared to the use of an MLP.

1.3 Paper Outline

This paper has the following structure. In section 2, the game of Agar.io and the core behaviours required to play the game are described. Furthermore, this paper explains why the game is interesting from the perspective of reinforcement learning. In section 3 the fundamental principles behind Q-learning and the techniques used to enhance its performance are outlined. Section 4 then describes why state representations are important for the training the algorithm. This is followed by a description of the types of state representations explored in this paper. The experimental setup follows in section 5, where the network structures and parameters used are described. Next, section 6 shows and discusses the results obtained. Finally, an outline of the conclusions is provided and ideas for future research are described in section 7.

2 The Game

Agar.io is a multiplayer online game in which the player controls one or more cells. The game has a top-down perspective on the map of which the size of the visible area of the player is based on the mass and count of their cells. The goal of the game is to grow as much as possible. This can be done by having the player’s cell absorb food pellets, viruses, or other smaller enemy player’s cells. The game itself has no end. Players can join an ongoing game at any point in time. Players start the game as a single small cell in an environment with other player’s cells of all sizes. When all the cells of a player are eaten, that player loses and can choose to re-enter the game.

Every cell in the game loses a small percentage of its mass in every time step. This makes it harder for large cells to grow quickly and it punishes inaction or hesitation. The game has simple controls. The cursor’s position on the screen determines the direction all of the player’s cells move towards. The player also has the option to ‘split’, in which case every player cell (given the cell has enough mass) will split into two cells of the same mass, both with half the mass of the original cells. One of these cells will be shot in the direction of the cursor with a given momentum. Furthermore, the player has an option to have every cell ‘eject’ a small mass blob, which can be eaten by other cells or viruses.

For the purpose of this research, the game was simplified to fit the computational resources available. The version of the game used has viruses disabled and is run with only one player. Furthermore, ejecting and splitting actions were disabled for the experiments in this paper. This is because ejecting is only useful for very advanced strategies, and splitting requires tracking of when the player’s cells are able to merge back together over long time intervals. The use of these actions would require recurrent neural networks such as LSTMs [11] which are outside of the scope of this research. Figure 2.1 shows a screenshot of the clone of Agar.io used for this research.

We introduce a ‘Greedy’ bot to the game to compare the RL agents against. This bot is preprogrammed to move towards the cell with the highest cell mass to distance ratio. It ignores cells with a mass above its biggest own cell’s absorption threshold. The bot also has no splitting or ejecting be-

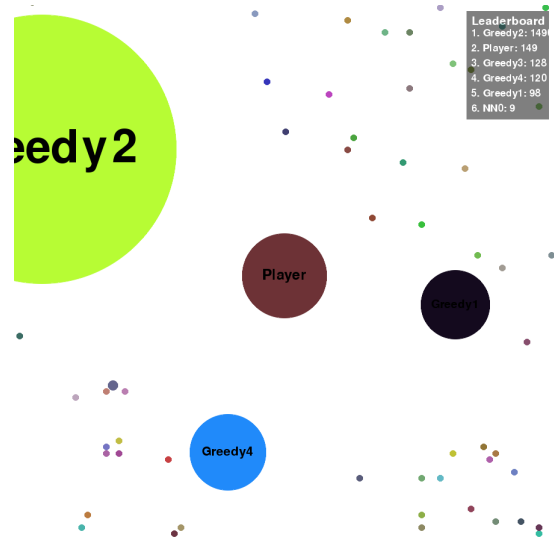


Figure 2.1: A clone of the game Agar.io used for this research. The player has one cell in the center of the screen. This player is in danger of being eaten by the Greedy bot seen on the top left, the other cells have a similar size as the player’s cell and therefore pose no danger. The little colored dots are pellets that can be consumed to grow in mass.

havior. This relatively naive heuristic, outperforms human players at early stages of the game. On the other hand, the heuristic can be outperformed by abusing its lack of path planning and general world knowledge later in the game.

3 Reinforcement Learning

This paper follows the general conventions [1] to model the reinforcement learning (RL) problem as a Markov decision process (MDP). In a Markov Decision Process an agent can take an action in a state to get to a new state. Most importantly the transition from the state to the new state has the Markov property: the stochastic transition probabilities between the states are only dependent on the current state and selected action.

To model the RL problem as an MDP, it must be defined what a state constitutes of. In short, the state consists of the properties the environment has and how the agent perceives the available relevant information. The transition between a state and an

action to a new state is handled by the game engine. This research applies frame-skipping to the MDP. In frame skipping a certain number of frames, or states, are skipped and the action that the agent chose is applied during all of these skipped frames. Also the rewards during these skipped frames are summed up until the next non-skipped state where the sum total is used as the reward. Frame skipping offers a direct computational advantage, as it allows the agent to not have to calculate the best action in every single frame of the game. More importantly, frame skipping leads to successive states in the MDP to be more different from each other than without frame skipping and leads to higher rewards, simply because more steps happened in between states. Making successive states more different from each other makes it easier for a function approximator to differentiate states. Larger rewards also have a positive effect on the training speed.

3.1 The Reward Function

In RL there must be some function that maps a state to a reward, also called the reward function. The aim of the agent in RL is to maximize the total expected reward that the agent receives in the long run through this reward function, also called the gain (G):

$$G = \sum_{t=0}^{\infty} r_t \cdot \gamma^t \quad (3.1)$$

r_t indicates the reward the agent receives at time t and γ indicates the discount factor. This discount factor is number between 0 and 1 which controls how much future rewards are discounted and therefore how much immediate rewards are preferred. A value of 1 would mean that the agent takes for every action into consideration how much reward this action will yield over the episode, whereas a value of 0 would make the agent completely myopic and disregard any future rewards.

The aim in Agar.io is to grow as big as possible. That means the agent has the aim to maximize its combined overall mass of all its cells in the shortest amount of time possible. This leads to the idea of the reward being the change in mass between the previous state and the current state:

$$r_t = \begin{cases} 0, & \text{if } t = 0 \\ m_t - m_{t-1}, & \text{otherwise} \end{cases} \quad (3.2)$$

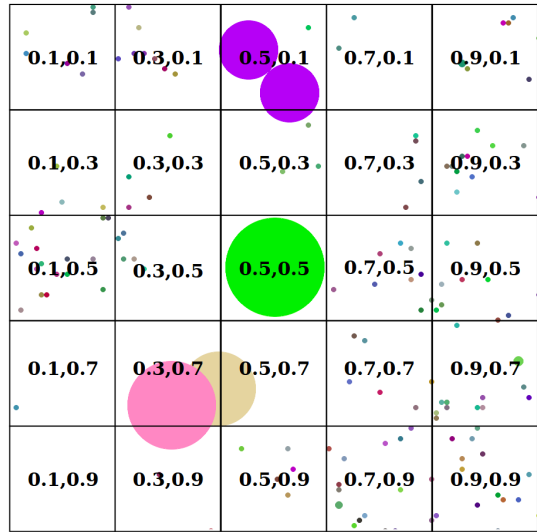


Figure 3.1: Possible action coordinates are laid out in a grid-like fashion. At a given state, the network will choose the square with the highest Q-value as an action. In total, 25 possible actions are used.

3.2 Q-Learning

Q-learning [10] predicts the quality (Q-value) of an action in a specific state. By iterating through all possible actions in a state, the algorithm picks the action with the highest Q-value as the action that the agent should take in that state. The Q-value indicates how much reward in the long term, or how much gain, the agent can expect to receive when choosing action a in state s . This prediction is updated over time by shifting it towards the sum of the rewards that the agent got for taking that action and the predicted value of the best possible action in the next state.

As Q-learning iterates over all possible actions in a state, the action space cannot be continuous. Therefore we discretize the action space by laying a grid of actions over the screen (Figure 3.1). Every center point of a square in the grid indicates a possible mouse position that the algorithm can choose.

To predict the Q-value for an action in a state, an artificial neural network (ANN) is used, which is trained through backpropagation. To construct the ANN to predict the Q-values we took inspiration from the network structure proposed in [5]. This

architecture feeds the state as an input to the network and has one output node per possible action. Additional details such as the learning rate or the number of layers can be found in the appendix. The tabular Q-learning update for a transition from state s_t after selecting action a_t with reward r_t and the new state s_{t+1} is:

$$Q(s_t, a_t) = Q(s_t, a_t) \cdot (1 - \alpha) + \alpha \cdot (r_t + \gamma \cdot \max_a Q(s_{t+1}, a)) \quad (3.3)$$

In this formula α indicates the learning rate. This formula is adapted so that it can be used to train an ANN by calculating the target for backpropagation for a specific state-action tuple (s_t, a_t) :

$$\text{Target}(s_t, a_t) = r_t + \gamma \cdot \max_a Q(s_{t+1}, a) \quad (3.4)$$

3.2.1 Exploration

Even though an optimistic initialization of the Q-values leads to a natural initial exploration [1], it is necessary to explore the action space throughout training to avoid being stuck in local optima. For Q-learning the ϵ -greedy exploration [1] was chosen due to its simplicity. The ϵ value indicates how likely it is that a random action is chosen, instead of choosing greedily the action with the highest Q-value. For this research the ϵ value is annealed exponentially from 1 to a specific value close to 0 over the course of training. The ϵ value should decrease over time, as this allows the agent to progress more in the game by taking more greedy actions. This causes the agent to progress steadily while exploring alternative actions over the course of training.

3.2.2 Target Networks

To stabilize Q-learning, Mnih et al. [5] introduced target networks. As the training of Q-learning maximizes over the possible actions taken in the next state, the combination of this training with function approximators can lead to the deadly triad [1]. This deadly triad gives a high probability of the Q-function to diverge from the true function over the course of training. A possible remedy to this problem is Double-Q-learning [12], which uses two Q-value networks. For the training of one network, the other network is used to calculate the Q-value of the action in the next state of a transition to avoid the positive feedback loop of the deadly triad. Mnih

et al. simplify this approach by introducing a target network in addition to the Q-value network. The parameters of the Q-value network are copied to the target network every time after a certain amount of steps. This requires no need to introduce a new separate network, but the maximization of the Q-values is still done by a slightly different network, therefore mitigating the unwanted effect.

3.2.3 Prioritized Experience Replay

Q-learning is an off-policy algorithm. This means that Q-learning can learn on transitions that are not directly generated by the Q-value network itself, but also by other policies or by an older version of the Q-value network. Lin [13] introduced a technique named experience replay to further stabilize and improve the performance of Q-learning. The technique has been shown to work well for DQN [7]. When using experience replay every transition tuple (s_t, a_t, r_t, s_{t+1}) is stored in a buffer instead of being trained on directly. If this buffer reaches its maximum capacity the oldest transitions in it get replaced. To train the value network using experience replay in every training step N random transitions from the replay buffer are sampled with replacement. For each of the transitions in the mini-batch the target for s_t is calculated and then the value network is trained on this mini-batch.

This form of experience replay offers a big advantage over pure online Q-learning. One assumption of using backpropagation to train an ANN is that the samples that are used to train in the mini-batches are independent and identically distributed. This assumption does not hold for online Q-learning, as each new transition is a partial result of the previous transition. Therefore random sampling from a large buffer of transitions partially restores the validity of this assumption. Furthermore, with experience replay, experiences are used more effectively, as the agent can learn multiple times from them.

As an enhancement to experience replay, Schaul et al. [14] developed prioritized experience replay (PER). PER does not sample uniformly from the replay buffer, but instead assigns the sampling probability to an experience i :

$$P(i) = \frac{TDE_i^\alpha}{\sum_k TDE_k^\alpha} \quad (3.5)$$

Here, the α coefficient determines how much pri-

oritization is used, $\alpha = 1$ would mean full prioritization. TDE stands for the temporal difference error of transition i , computed as:

$$TDE_i = r_t + \gamma \cdot \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \quad (3.6)$$

This implies that the badly predicted transitions will be replayed more often in the network, which was shown to lead to faster learning and better final performance [14].

Because more transitions with high TDEs will be trained on in PER, leading to proportionally larger changes in the weights of the network, Schaul et al. also introduce an importance sampling weight which decreases the magnitude of the weight change in the MLP for transition i anti-proportionally to its TDE_i . This is done to reduce the bias of training on average on more high TDE transitions. Therefore, a weight w_i is applied to the weight changes induced by each transition i of magnitude:

$$w_i = \left(\frac{1}{N} \cdot \frac{1}{TDE_i}\right)^\beta \quad (3.7)$$

In this formula N is the batch size and β controls the amount of applied importance sampling. In practice the weights are used in the Q-learning update by multiplying the prediction error for transition i , used in backpropagation, by w_i .

This research uses the OpenAI baselines repository [15] for prioritized experience replay to enhance reproducibility.

4 The State Representation

The information used in state representations can have varying levels of abstraction. The choice of a given state representation often brings positive and negative influences on the algorithm’s learning process, which the designer has to balance optimally. State representations with high levels of abstraction usually have the environment information preprocessed before it is fed to the algorithm. This has the advantage of allowing for a simpler network which will take less training time to converge. This is not without its downsides, as such approach has additional processing requirements. It also goes without saying that these state representations of hand-crafted features are inherently biased

due to being created with the programmer’s own heuristic in mind. An example of this is Bom et al.’s paper on learning to play Ms. PacMan [16], where a small network learns to play the game by using a representation that includes the distance to the closest collectable pills as determined by an A* search algorithm.

On the other end of the spectrum there are approaches where the unfiltered raw data of the environment is fed to the learning algorithm. The simplicity of this approach allows for agents to learn in complex state-action spaces for which humans might have non-optimal existing heuristics. The downside is that the large number of parameters the networks are required to have, brings issues with processing power and amount of training time before convergence.

One of the aims of this paper is to research how state representations of the same resolution, but with varying levels of preprocessing compare to one another. The base representation of the game is the raw state representation of the game, which comes in the form of RGB pixel values.

The second state representation uses the grayscale pixel values. A player in Agar.io aims to locate food pellets and cells in its view against a white background. Processing the RGB channels into a single grayscale channel will reduce the amount of weight tuning required for the network in order to extract non-white objects. This processing is performed by a pixel-wise averaging across the RGB channels.

The third state representation is a ‘semantic representation’ of objects in the environment. This consists of a vision grid in which every individual area unit has a value equal to the amount of food pellets contained in that area (see Figure 4.1). For the following experiments, the grid values at given areas were obtained from the game engine itself to reduce computational costs, but one could theoretically obtain these values from the RGB pixel values of the real game using preprocessing techniques.

Another aim of this paper is to explore how much the performance is influenced by the resolution of these representations. The DQN approach has shown success with state representation sizes of 84 by 84 [7], but one could hypothesize that as the representation resolution drops, it will be harder for the network to understand the semantics of the game. Given this, semantic representations should

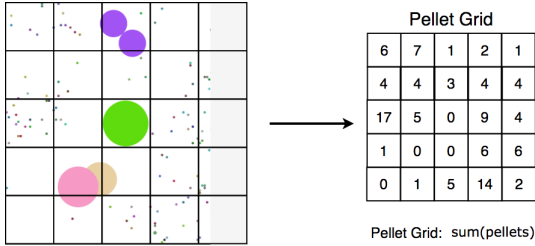


Figure 4.1: The semantic state representation consists a vision grid laid out on the player’s view. Values are then extracted from each area unit based on how many food pellets are present in it.

be expected to perform marginally better than pixel values at lower resolutions.

The last aim of this paper is to compare how different representations perform with different architectures. Every subsequent layer in a neural network can be thought of as providing recognition of more abstract concepts. Providing the network with a more semantically complex state representation to begin with, might relieve the network from the need to extract objects such as circles (for cells), as well as features such as the size of the circles (for estimating the mass of the cell). This hypothesis will be tested by comparing the performance of the pixel and semantic representations between a CNN with 3 convolutional layers to that of a CNN with 2 convolutional layers. The first network has the same structure as the one used in the 2015 DQN paper [7]. With the only difference being that the one used here only uses one single channel for the current representation of the game, whereas the ones used by Mnih et al. use convolution over the 4 last frames. The second CNN has a similar structure to the 2013 DQN paper [?], but with a slight difference in the number of filters.

Furthermore, to emphasize how the semantic representation can be used to achieve high performances with relatively small networks, the mentioned methods will also be compared to that of a small MLP without convolutional layers that uses a state representation of resolution 11 by 11.

5 Experimental Setup

5.1 General Experimental Setup

In the experiments of this paper, the various networks and state representations will be evaluated under a pellet collection task. In this task, one agent is placed into an environment with only food pellets present and the goal is to grow as large as possible.

In all experiments the environment is reset after 20,000 game steps. This is considered to be one episode. Upon reset all agents are reassigned a new cell with mass 10 at a random location and all pellet locations are randomized. This is done to avoid that the learning agents learn peculiarities of pellet locations on the map and to force the agents to also learn to deal with low cell mass strategies. Furthermore, to avoid the network from overfitting to one particular color in the pixel value representations, we also have the player cell colors be randomized every time an episode ends.

Each algorithm instance was trained with 300,000 state transitions. Given the network used a frame skip rate of 10, one state transition experience was generated every 11 in-game frames, giving a total of 3,300,000 game steps. These states were generated on-line as the network learned to play and stored into the experience replay buffer of size 20,000. Every training step, the network was trained on 32 experiences sampled (with replacement) into a single batch. On an Intel Xeon E5 2680v3 CPU @2.5Ghz it took between 32 to 84 hours to train each individual CNN run depending on the trial. State representations of 42 by 42 in resolution were at the lower end due to their smaller amount of network parameters, while resolutions of 84 by 84 took the longest to train on. On the other hand, the MLP runs took approximately 5.5 hours. Each condition was trained for a total of 10 independent runs and the mean across those runs was taken.

Every 5% of the training process the performance of one agent is tested five times. The noise factor of the agent (ϵ of ϵ -greedy) is set to zero. In this environment the agent can only collect pellets for 15,000 in-game steps. Furthermore, after training is completed the agent is placed in the environment 10 times to measure the final performance.

For the testing during and after training, the performance for each experimental condition is calculated by taking the mean across the testing runs of

the 10 independent training runs.

5.2 Network structures

We assembled two CNN architectures to be our value function networks. We also constructed a simple MLP to further test the semantic representation’s performance with low resolutions with a small network.

The simple MLP architecture consists of a variable input length, 3 fully connected layers, and an output layer. The input to the network consists of the grids of the semantic representation, which were first flattened into a 1D vector, and then had 2 extra values appended to it: the current mass of the player, and the ‘field of view’ (FoV) size of the player. These two extra values are information the human player has implicit access to in the real game through estimation of the total mass and FoV size by comparison to features such as the relative sizes of a food pellet, or the game background. This source of information is useful, as the optimal strategy in the game changes depending on size. For a state’s semantic representation resolution of 11 by 11 where a pellet vision grid are used, the 1D input of the network would be 123 in length. This input is then fed into 3 subsequent fully connected layers of 250 rectified linear units each. This is then followed by an output layer of 25 linear units. The output layer, as specified in the ‘Reinforcement Learning’ section, has units symbolizing a possible mouse position on the screen on a grid-like fashion.

The first CNN architecture has the same structure as that used for Atari games in the 2015 DQN paper [7], with the only difference being that the structures used here only use the current frame in the input for the convolution. The default input consists of 84 by 84 units in length, the number of channels is dependent on the type of representation used. The first convolutional layer uses a kernel size of 8 by 8 with stride 4 for a total of 32 filters. The second convolutional layer uses a kernel size of 4 by 4 with stride 2 for 64 filters. The third convolutional layer uses a kernel size of 3 with stride 1 for 64 filters. Every convolutional layer applies a rectified linear activation function. At this point in the network, the current layer’s output was flattened and, similar to the case of the MLP’s input, the values for the mass of the player and the FoV size were appended to it. Next, this 1D vector was fed to a fully con-

nected layer of 512 rectifier units, which was then followed by an output layer of 25 linear units.

Lastly, the second CNN architecture has a similar structure to the first one, but has only 2 convolutional layers. Again, the input by default consists of 84 by 84 units in length. The number of channels is dependent on the type of representation used. The first convolutional layer uses a kernel size of 8 by 8 with stride 4 and a total of 32 filters. The second convolutional layer uses a kernel size of 4 by 4 with stride 2 and 64 filters. Every convolutional layer applies a rectified linear activation function. Just like in the other CNN architecture, at this point the layer’s output is flattened into a 1D array and gets appended the mass and FoV player values. Next, a fully connected layer of 256 rectifier units is used, which is then followed by an output layer of 25 linear units.

The best hyper-parameters for all methods were coarsely searched for. Please refer to the appendix for all the parameter values used in these algorithms. All artificial neural networks were constructed using Keras 2.1.4 [17].

6 Experimental Results

6.1 General Results

As seen in Figure 6.1, for the set of resolutions tested for the vision grid representation on both CNN architectures, the 42 by 42 resolutions seem to perform marginally better than every other one. The 84 by 84 resolution performs the second best closely followed by the 63 by 63 resolution. Unsurprisingly, as seen in Figure 6.2, the 42 by 42 resolution also achieves good performances the fastest due to its lower number of trainable parameters. The 3 convolutional layer network seems to also achieve slightly better performances than the 2 convolutional layer CNN for resolutions of 42 by 42 and 63 by 63, but not for 84 by 84. The MLP architecture with the 11 by 11 resolution seems to achieve a higher performance than both CNNs using 84 by 84 resolutions, although not as high as CNNs using 42 by 42 resolutions. The MLP seems to learn at a similar rate as 84 by 84 CNN resolutions (see Figure 6.2).

The RGB pixel value representations seem to all have similar during-training performances as seen in Figure 6.3, suggesting that resolution does not have

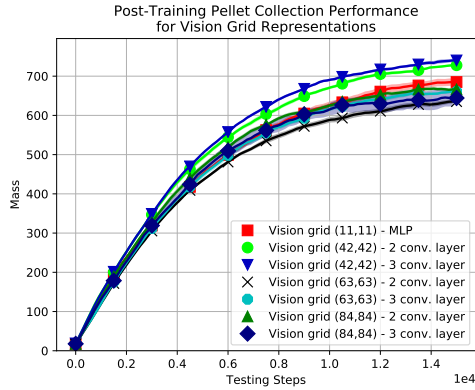


Figure 6.1: Post-training performance of vision grid representations with differing resolutions for the two CNN architectures, as well as for the MLP architecture with a 11 by 11 resolution. Each point represents the average of the 10 testing rounds and the shaded area denotes its 1 S.D. range. Results are averaged for 10 simulations.

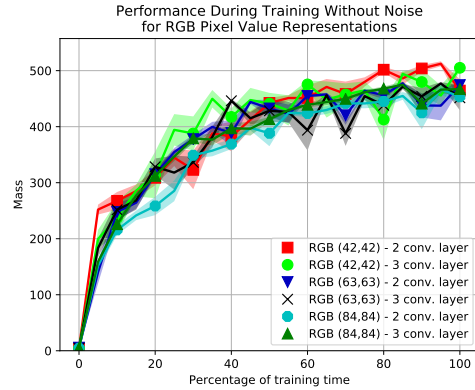


Figure 6.3: During-training performance of RGB pixel value representations with differing resolutions for the two CNN architectures. Each point represents the average of the 5 testing rounds and the shaded area denotes its 1 S.D. range.

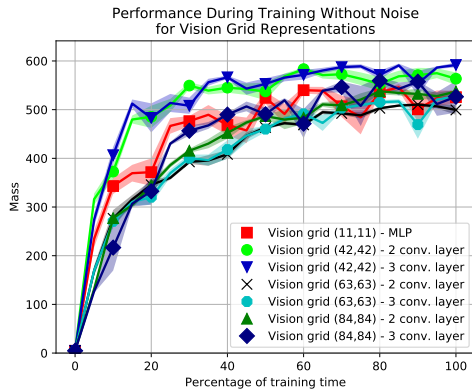


Figure 6.2: During-training performance of vision grid representations with differing resolutions for the two CNN architectures, as well as for the MLP architecture with a 11 by 11 resolution. Each point represents the average of the 5 testing rounds and the shaded area denotes its 1 S.D. range.

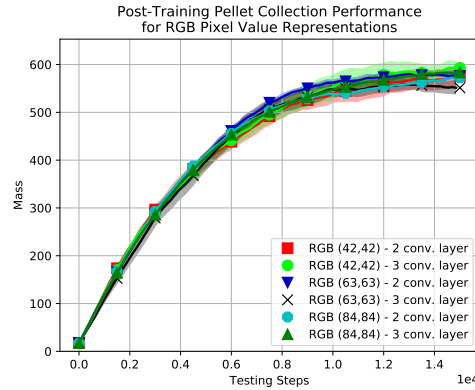


Figure 6.4: Post-training performance of RGB pixel value representations with differing resolutions for the two CNN architectures. Each point represents the average of the 10 testing rounds and the shaded area denotes its 1 S.D. range.

much of an effect on the learning of the networks for the resolutions tested. This is further emphasized by Figure 6.4, where there are no noticeable differences between the resolutions or architectures in post-training performance.

Lastly, the grayscale pixel value representation appears to have trends similar to those of vision

grids, but not to the same extent. The during training performance of the 42 by 42 resolutions seems to converge at a slightly higher performance than the other 2 resolutions (see Figure 6.5). The post-training performance seen in Figure 6.6 also seems to indicate that 42 by 42 resolutions have a higher performance after 300,000 training steps.

In order to observe how different kinds of state representations compare to one another, some of the highest performing runs were plotted together as seen in Figures 6.7 and 6.8. The 42 by 42 grayscale

representation on a 2 convolutional layer network seems to achieve a similar final performance as both of the 84 by 84 vision grid representations. This grayscale run is noticeably better than the plotted 42 by 42 RGB representation on a 3 convolutional layer network in terms of final performance.

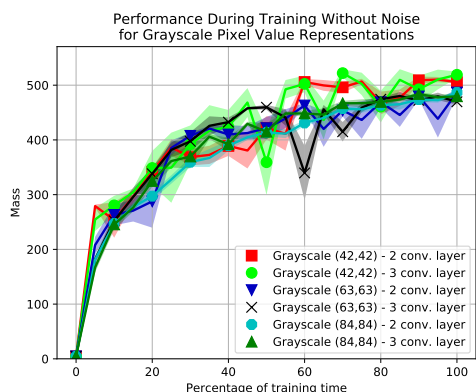


Figure 6.5: During-training performance of grayscale pixel value representations with differing resolutions for the two CNN architectures. Each point represents the average of the 5 testing rounds and the shaded area denotes its 1 S.D. range.

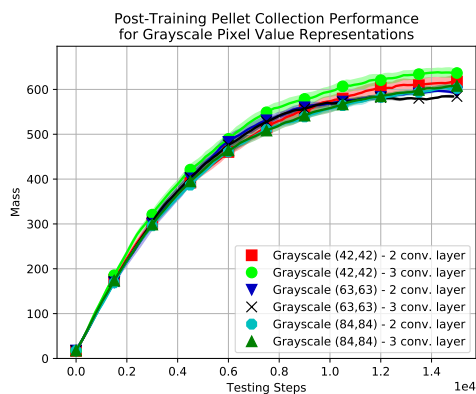


Figure 6.6: Post-training performance of grayscale pixel value representations with differing resolutions for the two CNN architectures. Each point represents the average of the 10 testing rounds and the shaded area denotes its 1 S.D. range.

The post-training performances for all conditions can be seen in Table 6.1. The column 'Mean Performance' shows the mean testing scores across all 10

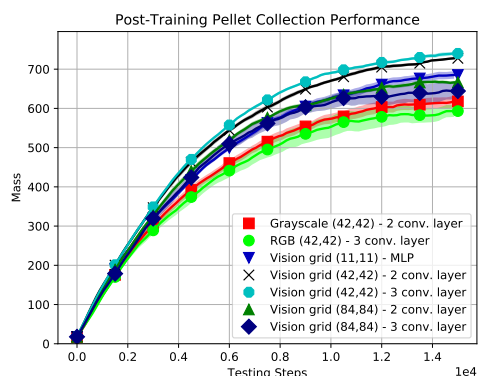


Figure 6.7: Post-training performance of various top-performing runs of various representations and network architectures. Each point represents the average of the 10 testing rounds and the shaded area denotes its 1 S.D. range.

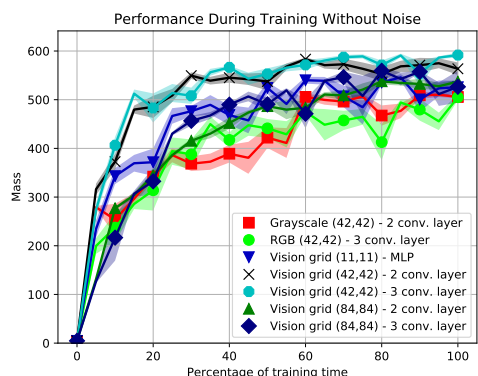


Figure 6.8: During-training performance of various top-performing runs of various representations and network architectures. Each point represents the average of the 10 testing rounds and the shaded area denotes its 1 S.D. range.

tests of all 10 simulations. The column 'Mean Max Performance' shows the average maximum scores across all 10 tests of all 10 simulations. The top scoring condition appears to be 'Vision grid (42,42) - 3 conv. layers' with a max score of 763. Comparing its max score to its 2 convolutional layer counterpart (which holds the second highest score) through the use of a t-test yields a p-value of 0.019, suggesting there is a significant difference between their maximum scores. All other conditions yield p-values below 0.001. Furthermore, the top scoring grayscale condition (Grayscale (42,42) - 3 conv. layer, hold-

Table 6.1: Post-training mean performances across 10 simulations. The 'Mean Performance' column contains the mean mass value for the post-training averaged performance curve (such as the ones seen in Figure 6.7). The 'Mean Max Performance' column is the max scores of the post-training averaged performance curve.

	Mean Performance	Std. Error Mean	Mean Max Performance	Std. Error Max
Random	18	0.1	31	0.2
Greedy Heuristic	527	0.3	693	0.6
Vision grid (11,11) MLP	481	2.8	704	5.9
Vision grid (42,42) 3 conv. layer	537	0.6	763	0.9
Vision grid (42,42) 2 conv. layer	526	1.1	750	1.1
Vision grid (63,63) 3 conv. layer	479	1.3	688	0.6
Vision grid (63,63) 2 conv. layer	461	1.4	662	1.8
Vision grid (84,84) 3 conv. layer	480	5.7	675	5.4
Vision grid (84,84) 2 conv. layer	494	1.2	696	1.1
Grayscale (42,42) 3 conv. layer	471	4.2	669	4.9
Grayscale (42,42) 2 conv. layer	449	4.5	648	6.8
Grayscale (63,63) 3 conv. layer	446	2.0	617	1.2
Grayscale (63,63) 2 conv. layer	448	3.1	625	1.1
Grayscale (84,84) 3 conv. layer	442	2.1	632	1.7
Grayscale (84,84) 2 conv. layer	439	1.5	627	2.2
RGB (42,42) 3 conv. layer	432	7.2	628	6.8
RGB (42,42) 2 conv. layer	425	7.0	609	9.9
RGB (63,63) 3 conv. layer	421	7.8	588	7.4
RGB (63,63) 2 conv. layer	436	2.3	609	2.7
RGB (84,84) 3 conv. layer	421	7.8	588	7.4
RGB (84,84) 2 conv. layer	427	0.9	598	1.6

ing a max score of 669) can be tested against the top scoring RGB condition (RGB (42,42) - 3 conv. layer, holding a max score of 628). Performing a t-test yields a p-value of 0.096, suggesting the difference is not significant.

6.2 Discussion

As seen in Figure 6.8, the best performances are achieved by the CNN networks using vision grid representations. Although the MLP network achieves a surprising performance despite its requirement of having a low resolution state representation, both CNN architectures using a 42 by 42 resolution input reach a higher performance at a faster pace.

The performance increase in relation to the MLP is likely due to CNNs' increased ability to process local changes in the environment, thus not having to evaluate potentially uncorrelated inputs far apart in the network's input. This also allows a CNN

to have a higher resolution input while keeping its number of parameters low, which helps explain why CNNs reach higher performances faster than the MLP. The deeper CNN architecture at 42 by 42 input resolution has 118,969 trainable parameters while the MLP architecture has 169,753.

The semantic representations yield a surprising performance in comparison to the RGB and grayscale pixel values. Even at resolutions of 11 by 11, the MLP yield a significantly higher performance. It should be noted that the pixel value representations have not fully converged after 300,000 training steps (see Figure 6.8), and that it could be the case that given enough training time, these could match the performance of vision grids. The same can be said about higher vision grid resolutions (particularly 84 by 84), which by the end of the training period have also not converged.

A reason that could explain why the 63 by 63 resolutions performed worse, is the input dimensions

are odd-numbered while the kernel stride are even. This causes the network to ignore 3 columns on the right of the input and 3 columns in the bottom of the input, leading to a loss of possibly important information.

7 Conclusion

This paper has researched the effect that different types of state representations have on the learning process of the Q-learning algorithm. Also, the effect that the resolution of these representations have was investigated. Furthermore, the performance and learning speed of 3 different artificial neural networks was explored.

The best performing resolutions for the CNN networks was 42 by 42, which outperformed 63 by 63 and 84 by 84 resolutions for the vision grid representations. For the pixel value representations, the change in resolution had little effect. For the resolutions of 42 by 42, the state representation that performed the best was the vision grid with a significant increase over the two pixel value representations. The grayscale pixel value representation performed somewhat better than the RGB pixel value representation.

As for the value function networks used, the CNN with 3 convolutional layers performed similar to the CNN with 2 convolutional layers, although for occasional runs, the 3 convolutional layer network showed slight increases in performance over the 2 layered one. The MLP architecture with an 11 by 11 resolution had a higher performance and faster convergence than most other experimental conditions. The MLP performance was only surpassed by the vision grids representation with 42 by 42 resolution. This might be explained by the low number of parameters these two conditions have. It could be possible that higher resolutions such as 84 by 84 could achieve higher performances given longer training times, since they seemed to not have fully converged after the given training time.

7.1 Future Work

First and foremost, an interesting goal for future work would be to expand upon the research done in this paper, so as to create an algorithm that can learn to play the full game of Agar.io. Following,

are possible strategies that could be employed to achieve such goal.

As reinforcement learning techniques are applied on more complex games of the current age, something to consider are design choices for algorithms to faster propagate rewards backwards in the state-action space. One such approach could be the use of multi-step algorithms, although these come at the disadvantage of increased computational requirements. One such example would be multi-step Q-learning, such as the one by Peng and Williams [18].

Another possibility could be the use of incremental discount factors. This would allow agents to first develop greedy short term behaviours, and as those establish the agent with a stable performance foundation, start gradually valuing longer-term strategies with distant rewards. One such example is annealing of the discount factor as implemented in OpenAI's 'OpenAI Five' Dota2 project [19].

Furthermore, another important necessity of reinforcement learning algorithms is the generation of interesting training data. Instead of the ϵ -greedy approach, other algorithms can provide smarter exploration methods through high-level decision making. One such approach is used in hierarchical actor critic methods, such as the h-DQN approach of Kulkarni et al. [20].

As processing requirements become a limiting factor in the field of deep reinforcement learning, the development of smarter techniques is accelerating. The use of introspection to gain inspiration of our own learning processes is opening doors to achieving human-level control in evermore complex systems, where even today, the applications of reinforcement learning remain unforeseen.

8 Acknowledgements

We would like to thank the Center for Information Technology of the University of Groningen for their support and for providing access to the Peregrine high performance computing cluster.

References

- [1] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2017.

- [2] David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533, 1986.
- [3] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*, 2(4):303–314, Dec 1989.
- [4] Gerald Tesauro. Temporal difference learning and TD-Gammon. *Communications of the ACM*, 38(3), 1995.
- [5] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- [6] Guillaume Lample and Devendra Singh Chaplot. Playing FPS games with deep reinforcement learning. In *AAAI*, pages 2140–2146, 2017.
- [7] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [8] A. Shantia, E. Begue, and M. A. Wiering. Connectionist reinforcement learning for intelligent unit micro management in Starcraft. In *Neural Networks (IJCNN), The 2011 International Joint Conference on*, pages 1794–1801. IEEE, 2011.
- [9] Stefan J. L. Knecht, Madalina M. Drugan, and Marco A. Wiering. Opponent modelling in the game of Tron using reinforcement learning. In *Proceedings of the 10th International Conference on Agents and Artificial Intelligence - Volume 2: ICAART*, pages 29–40. INSTICC, SciTePress, 2018.
- [10] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Cambridge, 1989.
- [11] Sepp Hochreiter and Jürgen Schmidhuber. Long Short-Term Memory. *Neural Computation*, 9(8):1735–1780, November 1997.
- [12] Hado V. Hasselt. Double Q-learning. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages 2613–2621. Curran Associates, Inc., 2010.
- [13] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8(3):293–321, May 1992.
- [14] T. Schaul, J. Quan, I. Antonoglou, and D. Silver. Prioritized Experience Replay. *ArXiv e-prints*, November 2015. arxiv:1511.05952.
- [15] P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu. OpenAI baselines. <https://github.com/openai/baselines>, 2017.
- [16] L. Bom, R. Henken, and M. Wiering. Reinforcement learning to train Ms. Pac-Man using higher-order action-relative inputs. In *2013 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL)*, pages 156–163, 2013.
- [17] François Chollet et al. Keras. <https://keras.io>, 2015.
- [18] Jing Peng and Ronald J Williams. Incremental multi-step Q-learning. In *Machine Learning Proceedings 1994*, pages 226–232. Elsevier, 1994.
- [19] OpenAI Five Blog. <https://blog.openai.com/openai-five/>. Accessed: 16-07-2018.
- [20] Tejas D Kulkarni, Karthik Narasimhan, Arda van Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. In *Advances in neural information processing systems*, pages 3675–3683, 2016.

A Appendix

Parameter	Value
Reset Environment After	20,000 training steps
Frame Skip Rate	10
Discount Factor	0.85
Total Training Steps	300,000
Optimizer	Adam
Loss Function	Mean-Squared Error
Weight Initializer	Glorot Uniform
Activation Function Hidden Layers	ReLU
Activation Function Output Layer	Linear
Prioritized Experience Replay Alpha	0.6
Prioritized Experience Replay Beta	0.4
Prioritized Experience Replay Capacity	20,000
Training Batch Length	32
Q-Learning Steps Between Target Network Updates	1500