
Learning Slither.io

Amol Kapoor

Department of Computer Science
Oxford University

Adam Cobb

Department of Computer Science
Oxford University

Abstract

Reinforcement learning aims to create artificial intelligence models that can dynamically learn how to complete various tasks without preexisting information and minimal guidance. While these agents have shown great success in solving simple tasks, only recently has the scale and complexity of reinforcement learning models rivalled human intelligence and creativity. Thus, few models exist for tasks where dozens of players may be competing at one time. This paper presents a model for learning the popular internet multiplayer game Slither.io. Our agent is trained using a combination of inverse reinforcement learning methods and the reinforcement learning A3C algorithm. Our work provides a jumping off point to learn more about solving the Slither.io environment for future research. After significant training, further techniques are still required to teach an agent to successfully learn Slither.io.

1 Introduction: Literature Review

In this paper, we present a model that can be trained to play Slither.io¹. In Section 1.1 we include an introduction to Reinforcement Learning (RL) and then describe relevant literature to our problem in Sections 1.2-1.5. In Section 2, we discuss the model and training scheme we use. In Section 3.1, we discuss the experiment, including a review of the Slither.io environment. In Section 3.2 we show results. Section 4 presents our conclusions and discusses ways to move forward and improve on our research.

1.1 RL Preliminaries

The goal of reinforcement learning is to create agents that can learn to operate in unknown environments. Because reinforcement learning agents can be applied in any system with a state and a reward system, reinforcement learning has applications in numerous fields. Reinforcement learning promises a generalized artificial intelligence. Such an artificial intelligence could help humans better understand and solve complex problems [1].

Reinforcement learning is a rich field with a long history that was originally inspired by neuropsychological research into how humans learn. However, though reinforcement learning has been studied for many years, the principles of the field have remained fairly consistent over time.

Reinforcement learning is an extension of the Markov Decision Process (MDP) framework. MDPs describe discrete time stochastic control processes. It is a mathematical framework where outcomes are decided partly by probability of state transitions, and partly by agent decision making. Thus, an MDP is composed of states, state transition probabilities, and possible actions given a state.

Reinforcement learning tasks add rewards to MDPs. They consist of a set of environment states \mathcal{S} , a set of actions available to the agent \mathcal{A} , policies defining transitions between states and actions,

¹Code located here: <https://github.com/theahura/Oxford-Reinforcement-Learning/tree/master/final-proj>

and rules defining the immediate reward of a transition from one state to the other. A reinforcement learning agent interacts with the environment in discrete time steps. Starting from state s in S , at each time t , the agent receives an observation \mathbf{o} and a reward \mathbf{r} . The agent then chooses an action \mathbf{a} from the set \mathbf{A} based on policy π , and moves the environment to a new state s' . The goal of an agent is to maximize reward by changing the policies it follows [1].

A natural development from the basic definition of reinforcement learning tasks is the development of internal state representations by the agent. This internal representation, called a value function, maps a state to the amount of reward expected from that state based on a policy: $v_\pi(s) = \mathbf{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$, where gamma is a parameter representing the discounted value of rewards from future states [1]. This definition for the value of a state naturally lends itself to an optimality equation, or the Bellman Equation, $v_\pi(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s, a)[r(s, a, s') + \gamma(v_\pi(s'))]$, defining the recurrence relationship between the value of a state and the value of its successors.

The definition for an optimal value function also naturally lends itself to a definition for the optimal policy - simply, take the action that leads to the state with the highest value, or greedy behavior: $\pi(s) = \operatorname{argmax}_a q_\pi(s, a)$, where q represents the state action value, or the value of taking a certain action in a certain state² [1]. Traditional reinforcement learning algorithms attempt to determine the value function through the iterative experience of the agent in the environment. Important algorithms include Monte-Carlo methods, TD(λ) methods, Q-Learning, and SARSA [1], though explaining each of them is out of the scope of this paper.

Because the agent must maximize reward, reinforcement learning algorithms need to determine a balance between exploration and exploitation. Without exploration, the agent may never find the optimal reward path; without exploitation, the agent may never actually maximize reward based on information that it has discovered. Mechanisms for balancing exploration and exploitation is an ongoing field of study within reinforcement learning [1].

1.2 Policy Gradient Methods and Actor Critic

Early reinforcement learning techniques focused on the value function approximation, with the default action policy being represented as greedy to the value function [2]. Value function approximation methods have some significant problems, including being oriented towards deterministic instead of stochastic policies and having discontinuous changes in action selection due to arbitrarily small value changes [2]. Further, many value function approximation methods do not have convergence guarantees. Policy gradient methods aim to solve these problems by estimating the policy directly. Actor critic methods are a subsection of policy gradient methods that combine value function approximation (the critic) to supply the policy (the actor) with low-variance knowledge of the performance of a given policy [3][4].

1.3 Deep Reinforcement Models

Traditional algorithms for reinforcement learning often store state values and policy actions in arrays representing each state. These storage mechanisms grow exponentially with new states, making traditional reinforcement learning algorithms untenable for complex problems that often have hundreds or thousands of states. Complex reinforcement learning tasks therefore require function approximation methods to determine optimal value functions and perform policy evaluation. Thanks to the explosion of computing power in the last decade, deep learning function approximators have seen significant success in solving complex reinforcement learning tasks [5][6].

Two common building blocks of deep learning models are particularly appealing for reinforcement learning algorithms. Convolutional networks have shown success in understanding position and identifying objects in images [7]. Long short-term memory modules (LSTMs) have shown success in understanding sequential data feeds over time [7]. Together, these elements of deep neural networks

²Finding the q function is trivial from the value function, and takes a similar form:

$$q_\pi(s, a) = \mathbf{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a]$$

81 allow a reinforcement learning agent to use an unprocessed input stream (e.g. pixels on a computer
82 screen) to understand how actions lead to outcomes.

83 1.4 A3C

84 While modern hardware allows for large neural networks, these systems are often limited to expensive
85 and massive GPU server clusters. A recent learning algorithm proposed by Mnih et. al. optimizes
86 deep reinforcement learning for easily available multicore CPUs [8]. The algorithm, known as
87 Asynchronous Actor-Critic (A3C), asynchronously updates a global network from the experiences
88 generated by numerous agents, dubbed workers, running simultaneously on different threads. These
89 workers then synchronize with the global network after some amount of time or at the end of an
90 episode. The A3C algorithm proposes the following update rule for parameters θ defining the value
91 function and policy for a given state: $\nabla_{\theta'} \log \pi(a_t | s_t, \theta') A(s_t, a_t, \theta, \theta_v)$, where A is an estimate of the
92 advantage function, a function that estimates how much better an action was to take over any other
93 action (or the advantage of taking that action). We use the advantage function proposed by Schulman
94 et. al. [9]: $\sum_{l=0}^{\infty} (\lambda \gamma)^l \delta_{t+l}^V$, where δ^V is defined as $r_t + \gamma V(s_{t+k}, \theta_v) - V(s_t, \theta_v)$. This advantage
95 function addresses the significant data requirements of previous policy gradient methods [9], and is
96 therefore suitable for a large state environment like Slither.io.

97 The asynchronous method accounts for exploration through the asynchronous action of numerous
98 simultaneous agents, each experiencing a different portion of the environment. In order to further
99 encourage exploration, we follow Mnih et. al. in adding the entropy of the policy π to the loss
100 function [8].

101 1.5 Inverse Reinforcement Learning

102 Inverse reinforcement learning, or apprenticeship learning, encapsulates methods and algorithms
103 that teach an agent a reward function from an expert who already knows the environment [10]. For
104 environments with numerous complex states, it can be difficult for an agent to ‘make the first jump’
105 to understanding how to play the game. We use inverse reinforcement learning methods interspersed
106 with self-learning to better train the Slither.io agent. Current research in the field will compare the
107 reward function proposed by the agent with the true reward function [11]. However, in our case the
108 true reward function is not known. We instead treat the expert action as the ‘true’ action and compare
109 action-by-action with a softmax cross entropy function. Gradients can then be backpropagated
110 through the model based on the ‘difference’ between the actual action and the predicted one.

111 2 Model

112 2.1 Network Construction

113 Our model consists of a series of convolutional network layers with an LSTM network on top. The
114 value function and the policy function are both determined by separate fully connected linear layers
115 on top of the LSTM network. The final model can be viewed in Figure 1. State of the art research
116 recommends using elu activation neurons [12], initialized using the He initialization function [13].

117 As per the A3C algorithm, at any given time there are multiple worker networks running in parallel.
118 Gradients are collected based on the experiences and rewards of these threaded workers. The gradients
119 are then used to update a global network located on the master thread. After an update, the global
120 network syncs network variables with the worker that provided the update. Thus, at any given
121 time, the global network contains the most recent updates from all of the parallel threads, while the
122 individual threads all eventually update to the most recent network [8]. This network structure can be
123 viewed in Figure 2.

124 2.2 Training

125 The model is trained with a combination of inverse reinforcement learning and reinforcement learning
126 techniques. In the inverse reinforcement learning stages, the model assumes the human action is the
127 perfect action for the given state. The inverse reinforcement learning is used as a corrective measure.
128 The model is allowed to self train for most of the training steps. If the model seems to get caught in a

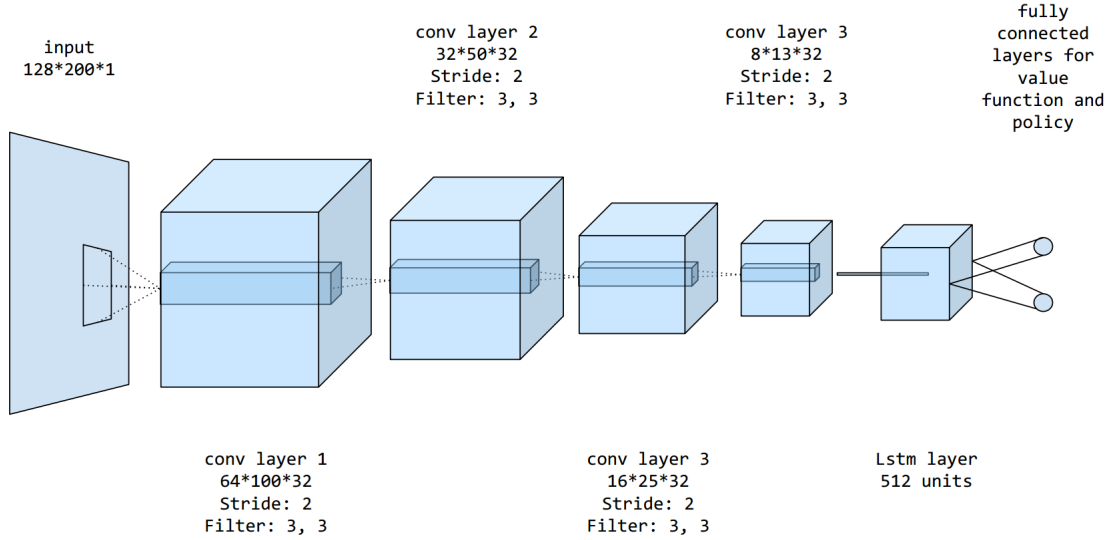


Figure 1: Visualization of the final model used for the Slither.io agent. The model consists of convolutional network layers each with a filter of 5x5 and a stride of 2, followed by an LSTM layer with 512 units, followed by two fully connected linear layers for the value function and the policy function respectively.

129 rut - for example, only choosing a single action - inverse reinforcement learning is used to correct the
130 model. In this way, the model learns with a mix of self taught and guided learning.

131 The model uses the softmax cross entropy function for its loss calculation. In the reinforcement
132 learning stages, the model implements the A3C algorithm stated above, and uses the A3C loss
133 calculation. To prevent exploding gradients and weights, the model also uses clipped gradients,
134 dropout, and L2 regularization. In both stages, the Adam Optimizer was used for propagating
135 gradients.

136 3 Experiment

137 3.1 Setup

138 3.1.1 Model Construction

139 The model was built using Tensorflow 1.0 without CUDA. Convolutional layers were built using the
140 implementation found in *tf.nn.conv2d*. The LSTM layers were built using the implementation found
141 in *tf.contrib.rnn.LSTMCell*. Training and testing was done on an ASUS ROG GL551J.

142 While there are numerous hyperparameters that can affect the model, due to lack of time and
143 computational ability we were unable to rigorously examine the optimal hyperparameter configuration.
144 Similarly, we were unable to rigorously search for the optimal network structure. The full list of
145 hyperparameters that we used for our final model can be found in Table 1.

146 3.1.2 Environment

147 Slither.io is a snake game where the goal is to create as long a snake as possible. The game takes
148 in any one of six actions - straight, left, right, fast-straight, fast-left, and fast-right - and returns as
149 reward the delta in the length of the agent. The game is episodic, in the sense that after some amount
150 of time the player agent will die (due to crashing into an opponent or the boundary of the stage) after
151 which the game will reset from the beginning. Further, any fast action will cause the snake to lose
152 length, therefore earning negative reward. If the action for speeding up is held until it is impossible to
153 lose length, the agent will go back to its normal speed and will no longer accrue negative reward even

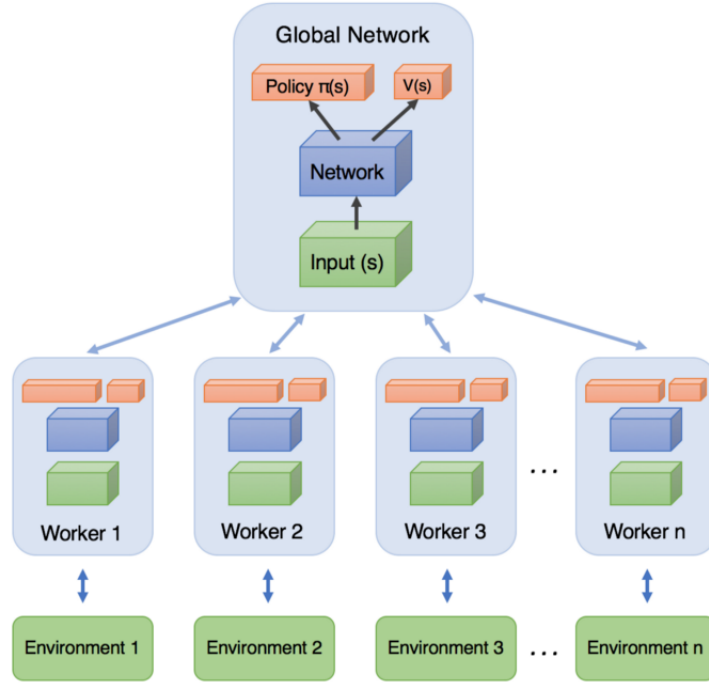


Figure 2: Visualization of A3C network[14]. The global network spins up n threads (generally n is less than or equal to the number of cores on the CPU), each of which runs a separate agent in its own version of the environment.

154 if the fast action is held. The game has the capacity to host 500 players in a single server at a single
 155 time. After each death, the agent would begin again in a new server.

156 The Slither.io environment is set up using OpenAI Universe package[15]. At each time step, Universe
 157 provides the model with a grey scale image representing the Slither.io environment. An example is
 158 shown in Figure 3. Thus, there are over eight million states. Universe also provides the reward earned
 159 in the last step. At each frame, the model inputs one of the six actions to the environment. Universe
 160 environment parameters can be found in Table 2. Due to limitations in computation, all training was
 161 done online - instead of storing and replaying memories, we had the network play in the Slither.io
 162 environment more often.

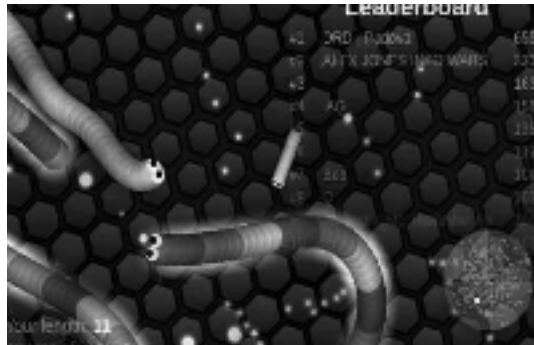


Figure 3: An example state observation given by Universe.

Table 1: Hyperparameters for Slither.io agent network.

Parameter	Value
Network Construction	
Convolution Layers	4
Filter Shape	(3, 3)
Stride	(2, 2)
Output Channels	32
LSTM Units	512
Loss Function	
Value Function Loss Constant	0.5
Entropy Constant	0.01
Discount Constant	0.99
Learning Rate	0.1
Regularization	
Max Gradient Clipping	60.0
Dropout Probability	0.5
L2 Constant	1e-5

Table 2: Hyperparameters for Universe Slither.io environment.

Parameter	Value
Input Height	128
Input Width	200
FPS	10.0
No Reward Value	0
Game Over Reward Value	0
Number of Actions	6

163 3.1.3 Data Collection

164 At the end of each episode or at the end of 1000 steps, the global network collects data from the agent.
 165 For quantitative results, this data was compared to data collected on a random agent that, for each
 166 frame, selected a random action of the six actions available. For qualitative results, we observed and
 167 recorded the behavior of the agent in game.

168 3.2 Results

169 3.2.1 Quantitative

170 Our model is unable to perform better on average than the random model. The loss for the trained
 171 model over all time steps can be seen in Figure 4. The total reward over time of the trained model
 172 compared to that of the random model over all time steps can be seen in Figure 5. These results
 173 suggest that the model was not learning, or that the training time was not a sufficient.

174 It is important to note the significant jumps that can be observed both in the loss and the total reward.
 175 These jumps show the unpredictable nature of Slither.io, and likely make it difficult for an agent to
 176 determine how its actions impact the environment.

177 3.2.2 Qualitative

178 For qualitative tests we examine how our agent performs at various points in its training. If the agent
 179 met the following criterion, we considered the environment solved:

- 180 1. Minimize use of the speed up action to avoid negative reward.
- 181 2. Move to avoid the outer boundary.

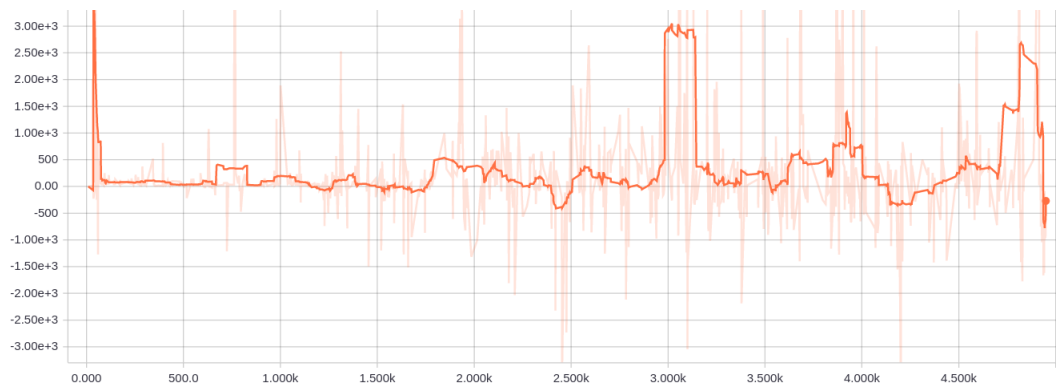


Figure 4: Loss at each time step for the trained model. The loss is smoothed with a rolling window average of about 500 steps. An agent that learned how to play the game from experience should have a decreasing average loss over time.

182 3. Move to avoid running into enemy agents.

183 Our agent was unable to meet any of these goals.

184 Our final agent was tested using a reward scheme that did not punish game overs or not growing.
 185 However, it is worth noting that we did test using alternative reward schemes before settling on
 186 the final version. Surprisingly, whenever the reward scheme had any kind of negative reward on
 187 game over, the agent would generalize to choosing only a single action³. Single action selection
 188 occurred even after the implementation of alternating inverse reinforcement and regular reinforcement
 189 learning stages. Notably, the selected action distribution did change significantly for each inverse
 190 reinforcement learning stage. This suggests that, as expected, the human intervention successfully
 191 pulled the model back into a more exploratory state.

192 Single action selection differs significantly from the random agent, which (obviously) showed no
 193 changes in its action selection. We expected that the agent would learn to avoid enemy agents due
 194 to the game over negative reward. We believe that this was caused due to an error in the Universe
 195 environment package. The Slither.io environment would continue polling for actions for a few
 196 moments after the agent had died before a game over was registered and the negative reward given.
 197 Thus, the agent was unable to tie the actions that led to its death with the negative reward, as the
 198 negative reward was too far away from the most relevant actions.

199 When the reward scheme did not punish game overs, we found that the agent had an evenly spread
 200 distribution of action choices⁴. This can be seen in Figure 6. Notably, we expected but did not observe
 201 a bias against selecting fast actions.

202 4 Conclusion and Future Work

203 In this paper we developed a reinforcement learning model to play Slither.io on the Universe platform,
 204 using state of the art reinforcement learning architectures and paradigms. Though our results left
 205 much to be desired, they indicated the potential for an active area of research.

206 Slither.io is a complex and chaotic game with state transitions that are uniquely difficult to learn
 207 due to the nature of multiple different human players simultaneously playing at any moment and the
 208 sheer number of possible states. The lack of consistency between enemy agents made predicting
 209 enemy behavior an extremely difficult task. Though we trained for over a week, this was likely
 210 nowhere near the amount of time necessary for the agent to properly generalize how its actions
 211 impacted the environment. Further, it is very likely that our model architecture was severely limited
 212 by computation - for example, attempting to increase the number of convolution layers or decrease the

³ We do not show the data for this model. However, an example recording can be seen here:
<https://youtu.be/M9YVTZ3aQqA>

⁴ A recording can be seen here: https://www.youtube.com/watch?v=_I59oRiANcg

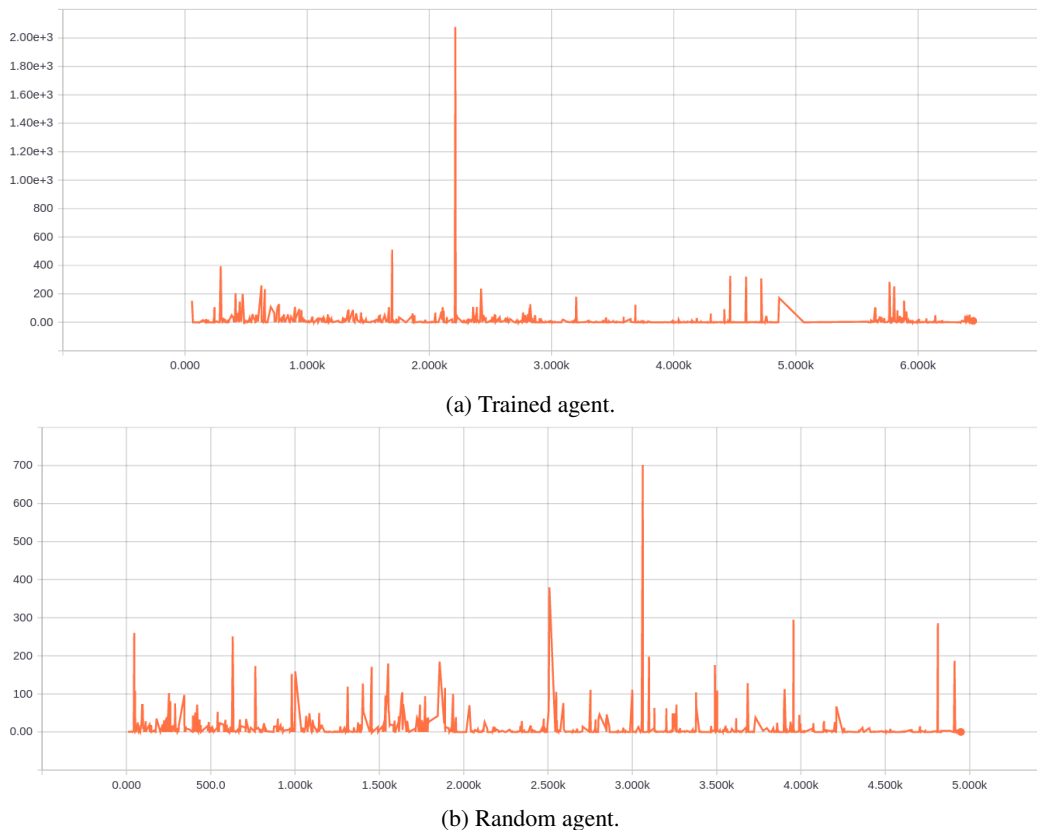


Figure 5: Comparison of total reward gathered by the trained agent and the total reward gathered by the random agent. Both agents have close to the same average score of about 20. An agent that learned how to play the game from experience should have an increasing total reward over time compared to the random control.

stride would cause the system to hang. These issues were further compounded by an inability to test for hyperparameters due to time constraints. With more training time and computational resources, we believe that our model will successfully improve. For future experiments, we also recommend rigorously examining various algorithms. Though A3C is state of the art for many RL tasks, each environment is different. We selected the A3C algorithm in part because it is optimized to multicore CPUs instead of GPUs. This hardware limitation prevented us from branching out and examining how other algorithms - for example, deep Q-Learning, would have performed on this task. It is also worth examining traditional algorithms - for example, tabular Q-Learning - to help establish and gain insight into the minimum viable baseline for this environment.

We found relative success in our alternating inverse reinforcement learning approach. Each inverse reinforcement learning stage led to a more exploratory agent. The ability to nudge an agent into a more exploratory state is useful for long term training, and can be viewed similar to a hands-off student-teacher relationship where the teacher intervenes to guide the student when stuck. While we did not formalize our human intervention schedule, we suspect that a more rigorous implementation of alternating between reinforcement learning and inverse reinforcement learning would speed up agent training significantly. We further recommend memory playback[5], a popular mechanism to aid reinforcement learning tasks by using memory replay as a simulated supervised learning environment.

The A3C algorithm - and, to the best of our knowledge, reinforcement learning in general - does not have a great solution to highly inconsistent state transitions. When combined with a large state size, it can be difficult for an agent to understand how its actions impact the environment. This is a significant problem in environments with multiple other agents, and makes learning games like Slither.io particularly challenging. Most reinforcement learning paradigms attempt to address probabilistic state changes through the resulting expected value of the state in the agent value function.

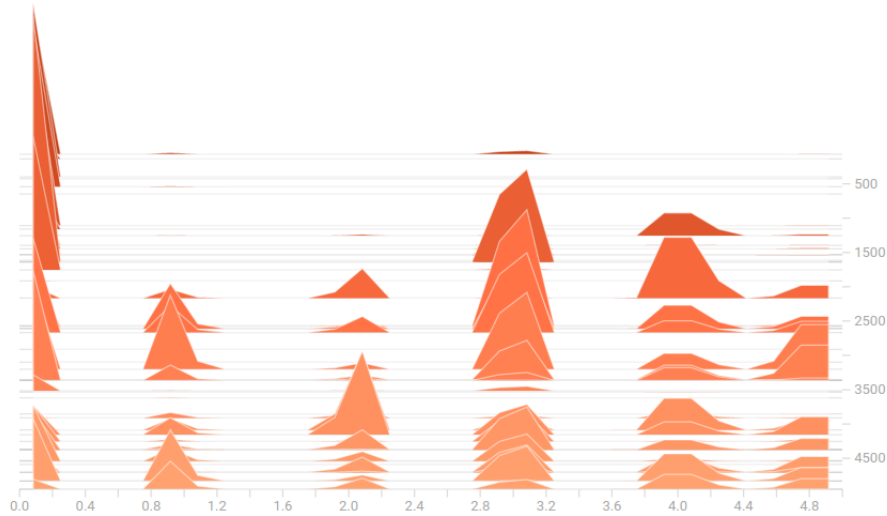


Figure 6: Frequency of selected actions over time for a no-punishment reward scheme. The right hand axis is the step. The distribution peaks (six peaks for six actions) indicate the objective number of times a certain action was chosen on a given step. Note that as time goes on, the agent begins to select across all actions evenly. Even selection occurred even though the agent began by favoring a single action.

236 Though this may eventually converge, training time can be prohibitive. Attempting to tackle this
 237 problem without introducing external (i.e. not self learned) information to the model proved difficult.

238 While it is hard to predict what a deep learning model will learn, we suspect that our current
 239 architecture was insufficient at predicting the possible state transitions of a given state. Due to
 240 a lack of pooling layer, our architecture was unable to generalize the shape of an enemy agent
 241 across the environment. This significantly increased the number of states our agent had to process
 242 before our agent would learn to generalize. It also preventing enemy behavior in one area of
 243 the environment from influencing predicted outcomes of enemy behavior in another area of the
 244 environment. However, pooling layers remove the geo-spatial relationships between detected objects
 245 in an image. It is important to preserve these relationships through the convolutional layers for the
 246 network to understand the agent's location in relation to enemy agents and the environment. We
 247 propose a modification to our architecture that replaces the single stacked convolutional networks
 248 with multiple such stacks, some with pooling layers and some without. We believe that this would
 249 allow later layers of the network to identify both the presence of an enemy agent as well as its location.
 250 Further, we believe that a single LSTM layer is insufficient to accurately predict all of the potential
 251 future behaviors of an enemy agent. We propose a modification to our architecture that replaces our
 252 single unit LSTM layer with multiple LSTM units across the same layer. We suspect that each LSTM
 253 unit would be able to predict a different outcome given an input state, similar to how lower layers of
 254 a convolutional network can specialize to identify specific shapes.

255 Acknowledgments

256 Thank you Adam for an amazing class. I learned more than I could have hoped, and definitely feel
 257 better equipped for further studies into this fascinating field. I look forward to applying what I have
 258 learned both over the summer and well into the future! I will keep you posted with what I am up to at
 259 Google this summer - I am in a machine learning group so hopefully I will get to apply a lot of what I
 260 learned this term!

References

- [1] Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, Cambridge, MA, USA, 1st edition, 2012.
- [2] Richard S Sutton, David A McAllester, Satinder P Singh, Yishay Mansour, et al. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, volume 99, pages 1057–1063, 1999.
- [3] Ivo Grondman, Lucian Busoniu, Gabriel AD Lopes, and Robert Babuska. A survey of actor-critic reinforcement learning: Standard and natural policy gradients. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 42(6):1291–1307, 2012.
- [4] Vijay R Konda and John N Tsitsiklis. Actor-critic algorithms. In *NIPS*, volume 13, pages 1008–1014, 1999.
- [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] 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):529–533, 2015.
- [7] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- [8] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International Conference on Machine Learning*, pages 1928–1937, 2016.
- [9] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- [10] Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the twenty-first international conference on Machine learning*, page 1. ACM, 2004.
- [11] Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Deep inverse reinforcement learning. *CoRR*, 2015.
- [12] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [14] Arthur Juliani. Simple reinforcement learning with tensorflow part 8: Asynchronous actor-critic agents (a3c), 2016.
- [15] OpenAI. Universe, 2016.