

## 1 A Additional Reward Plots

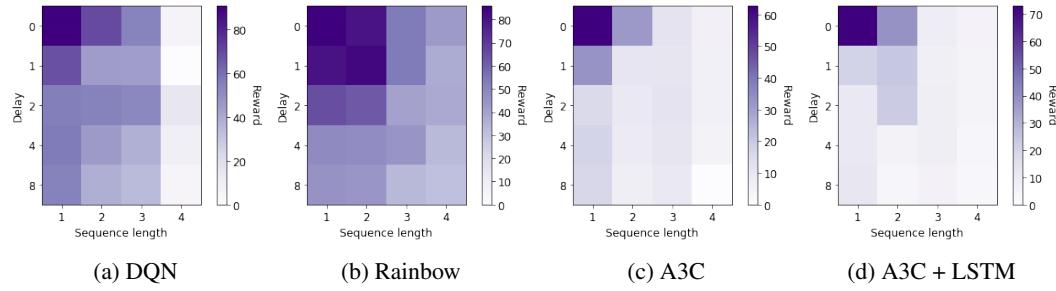


Figure 1: Mean episodic reward for evaluation rollouts (limited to 100 timesteps) at the end of training for the different algorithms **when varying delay and specific sequences**. Please note the different colorbar scales.

## 2 B Additional Learning Curves

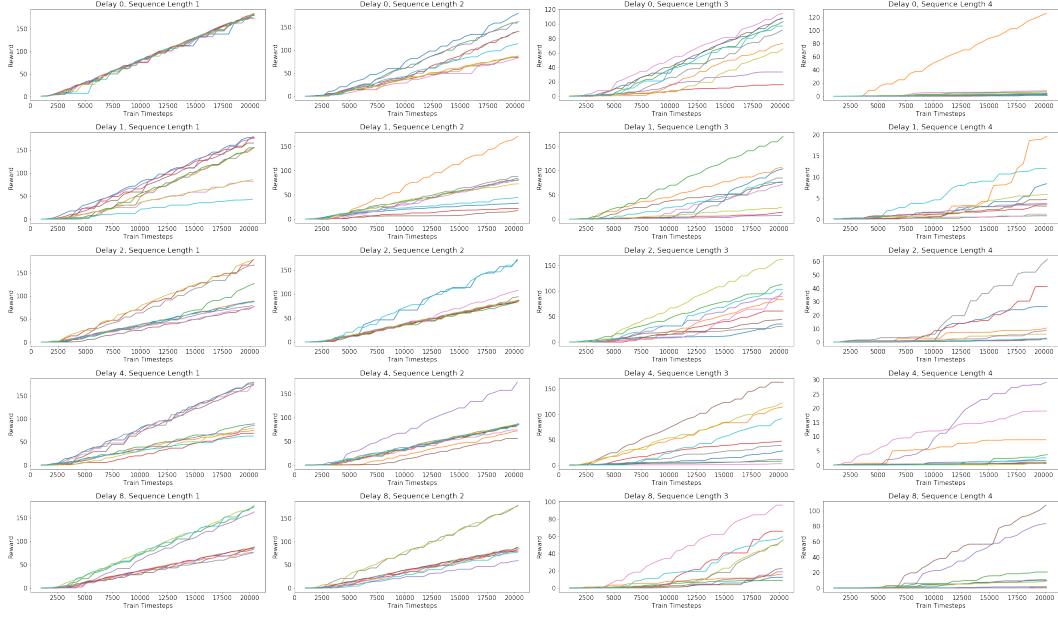


Figure 2: Training Learning Curves for DQN when varying delay and specific sequences. Please note the different colorbar scales.

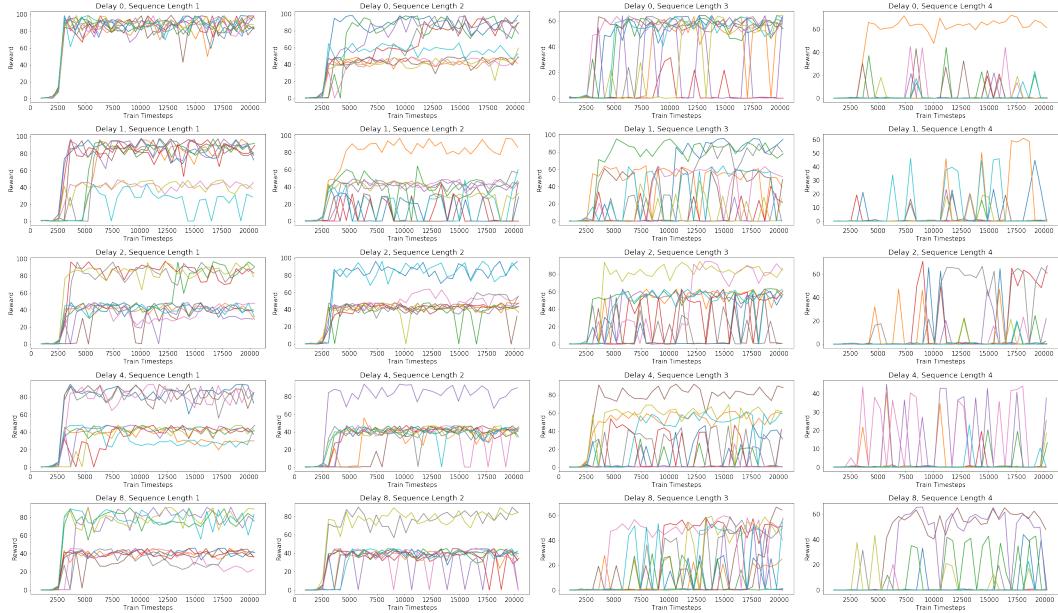
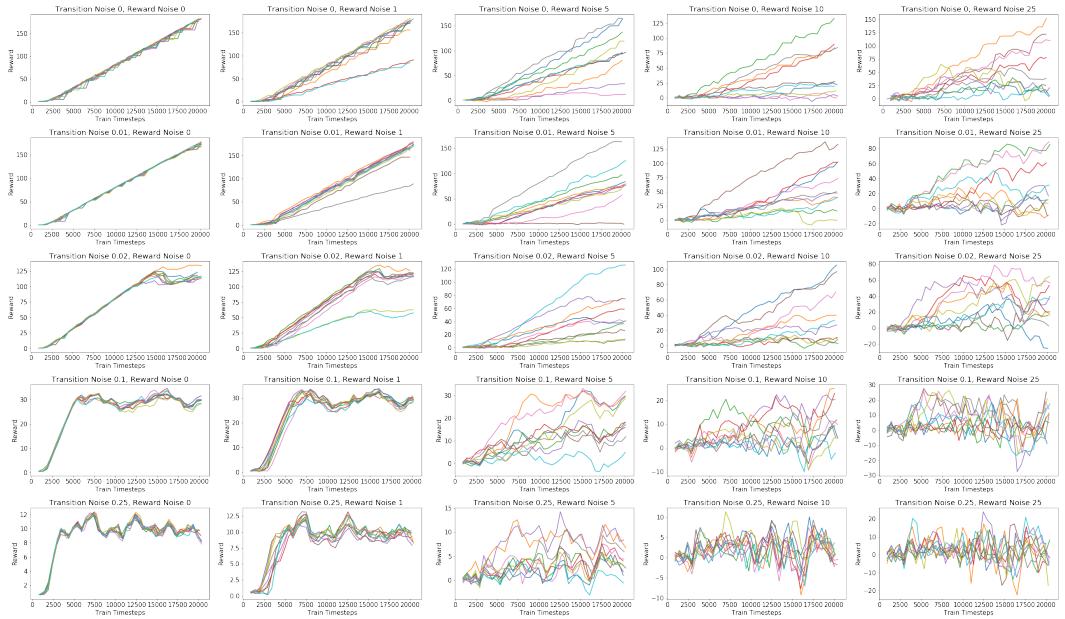
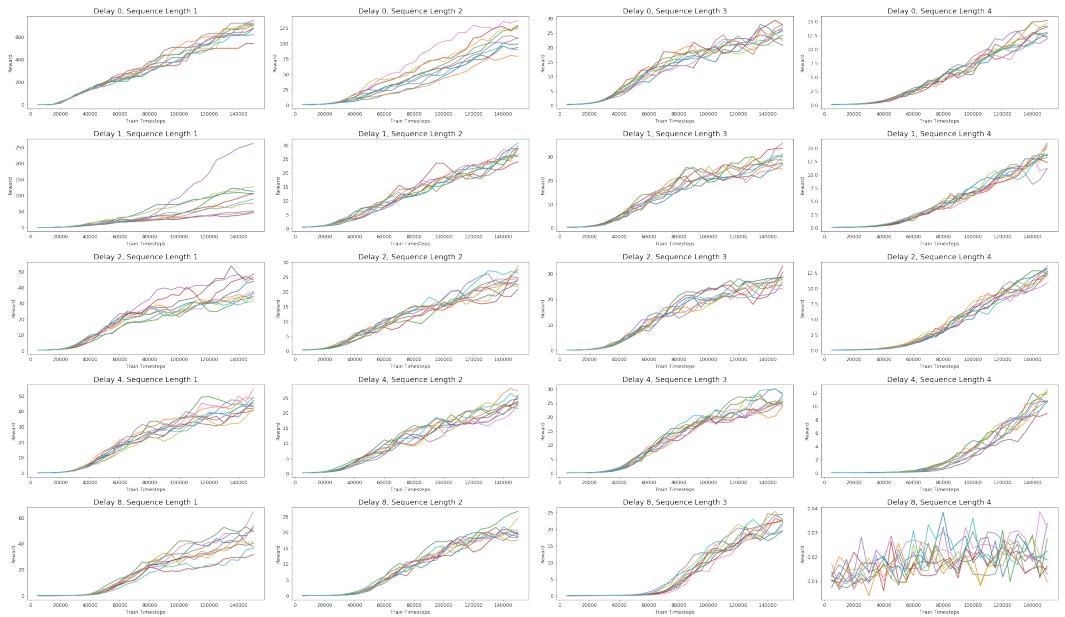


Figure 3: Evaluation Learning Curves for DQN when varying delay and specific sequences. Please note the different colorbar scales.



**Figure 4: Training Learning Curves for DQN when varying transition noise and reward noise.** Please note the different colorbar scales.



**Figure 5: Training Learning Curves for A3C when varying delay and specific sequences.** Please note the different Y-axis scales.

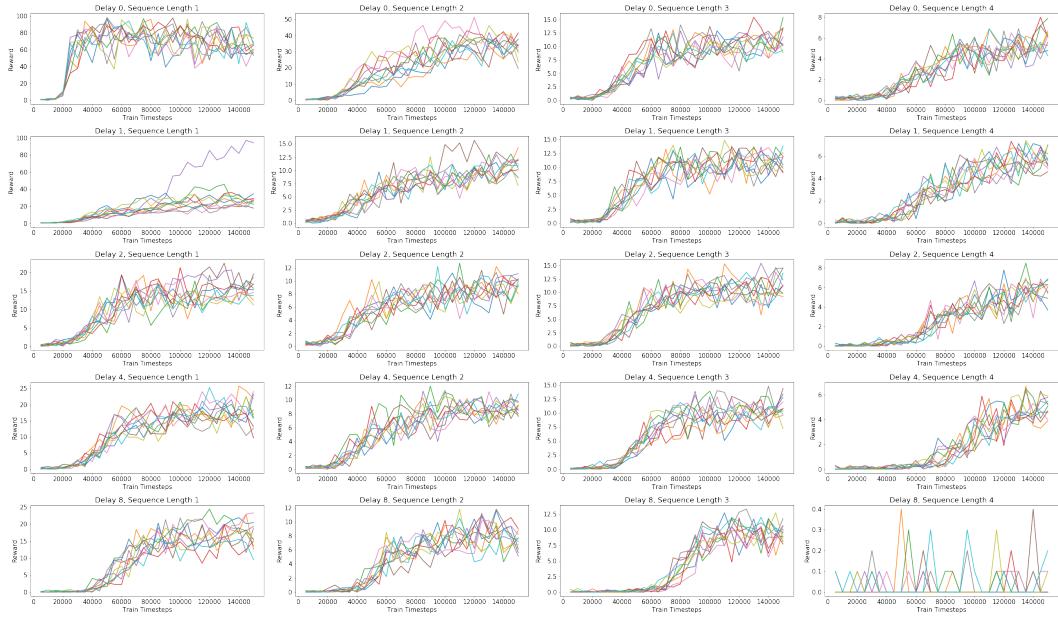


Figure 6: Evaluation Learning Curves for A3C when varying delay and specific sequences. Please note the different Y-axis scales.

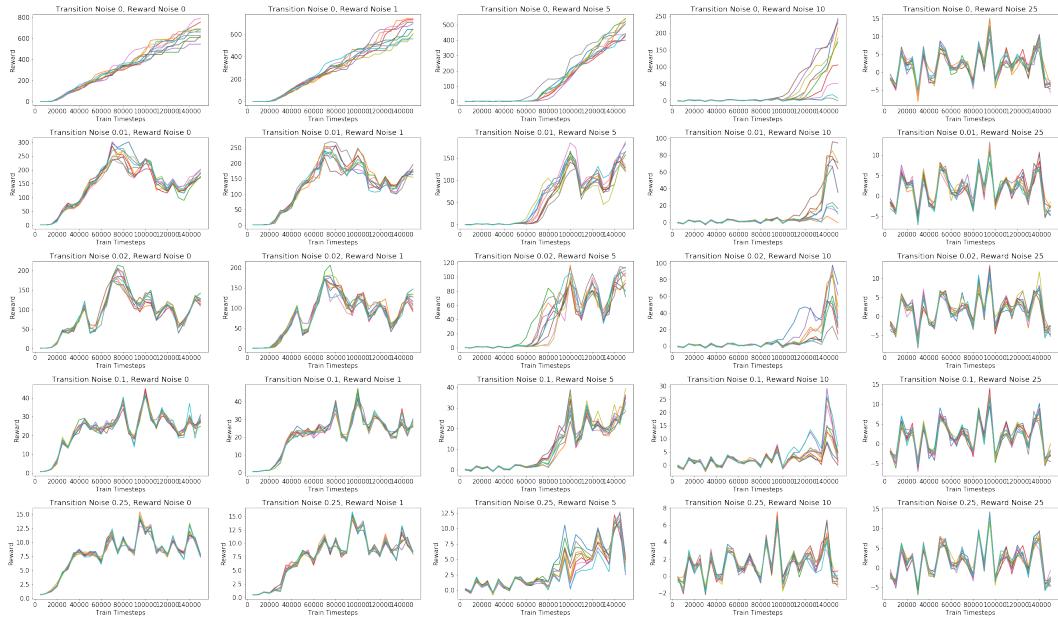


Figure 7: Training Learning Curves for A3C when varying noises. Please note the different Y-axis scales.

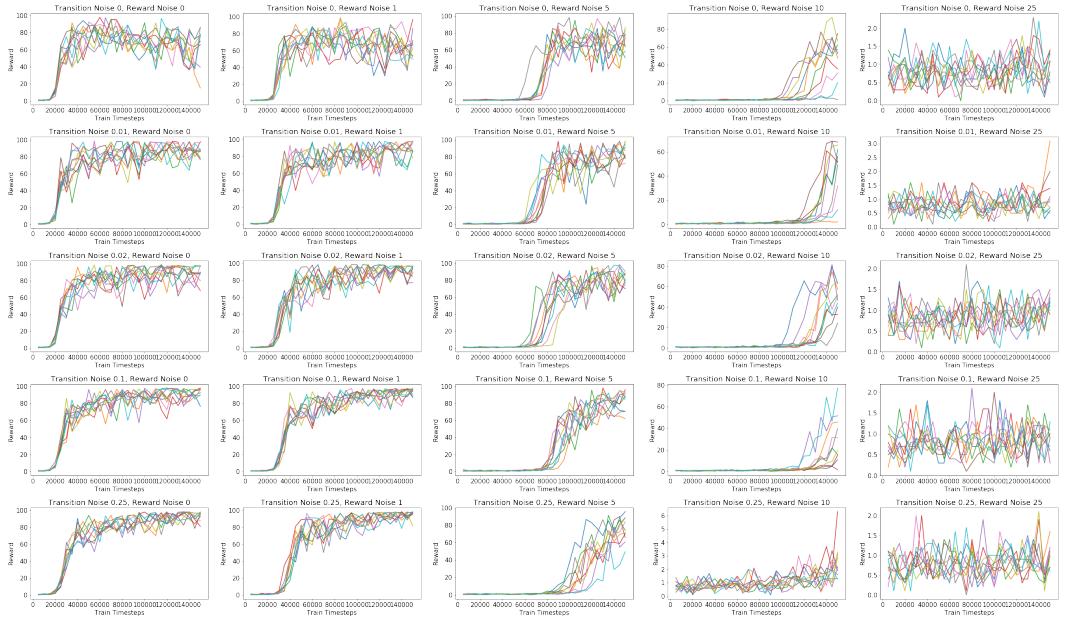


Figure 8: Evaluation Learning Curves for A3C **when varying noises**. Please note the different Y-axis scales.

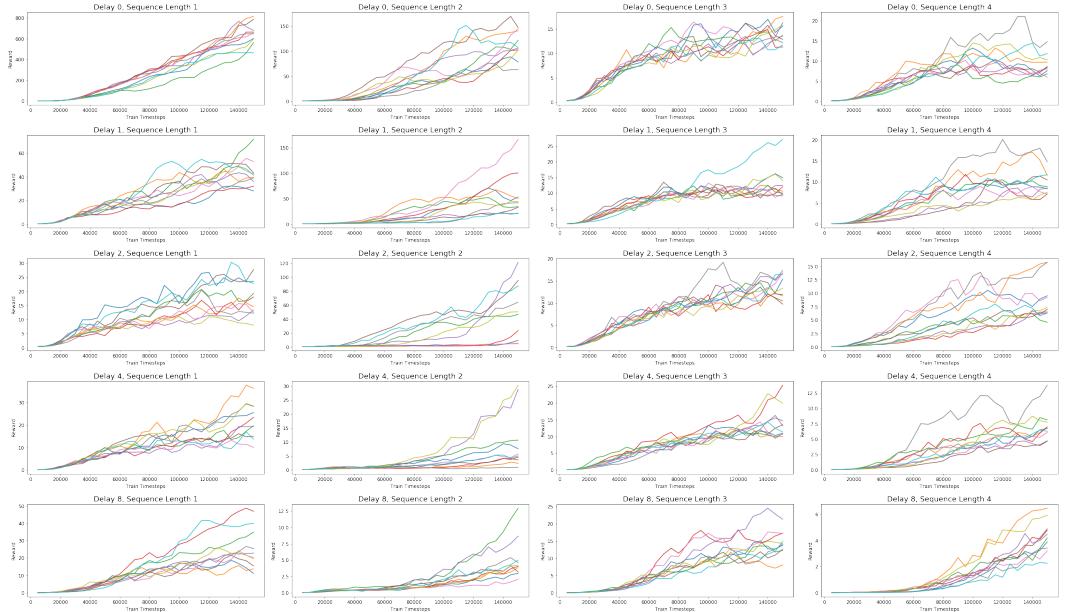
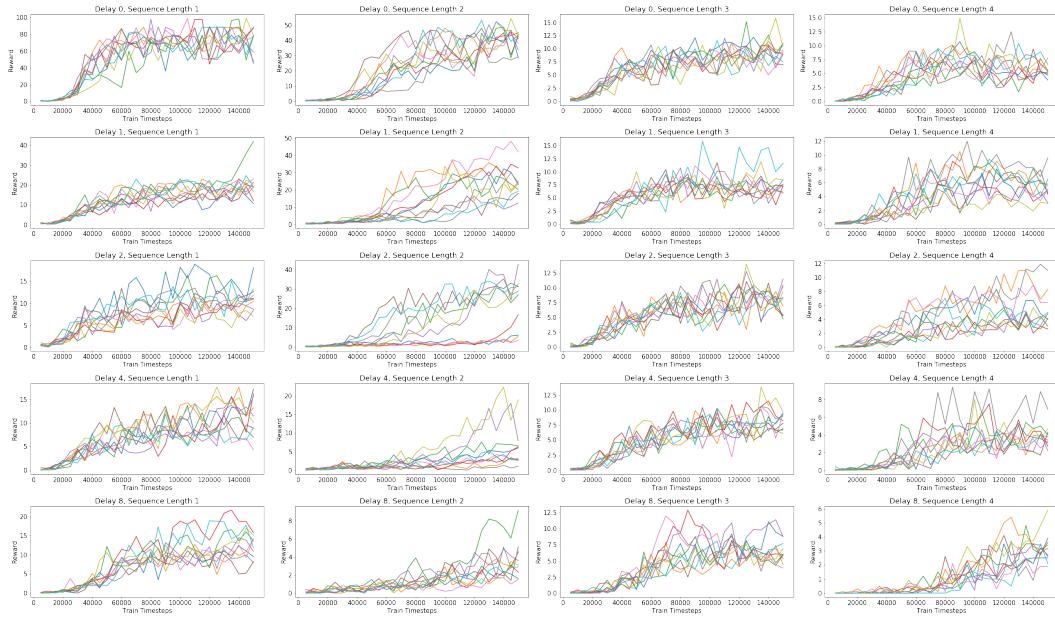
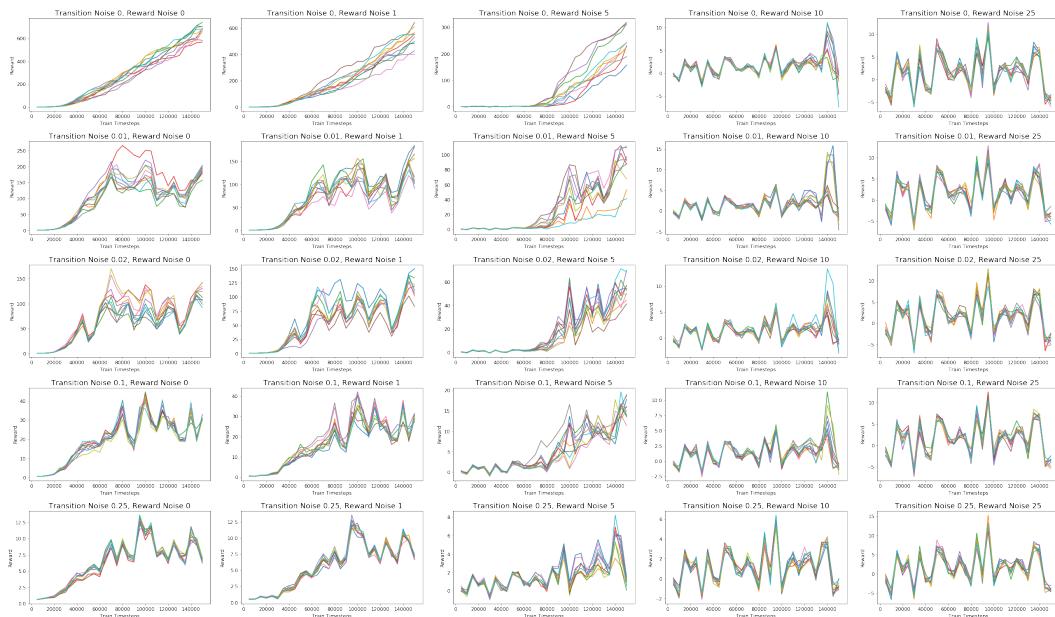


Figure 9: Training Learning Curves for A3C with LSTM **when varying delay and specific sequences**. Please note the different Y-axis scales.



**Figure 10: Evaluation Learning Curves for A3C with LSTM when varying delay and specific sequences.**  
Please note the different Y-axis scales.



**Figure 11: Training Learning Curves for A3C with LSTM when varying noises.** Please note the different Y-axis scales.

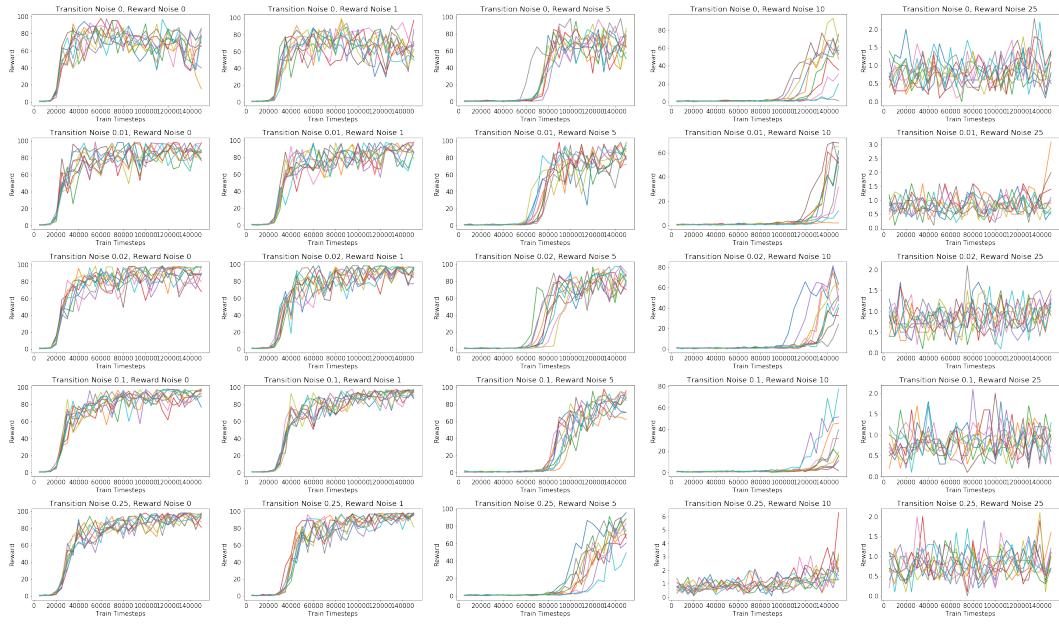


Figure 12: Evaluation Learning Curves for A3C with LSTM when varying noises. Please note the different Y-axis scales.

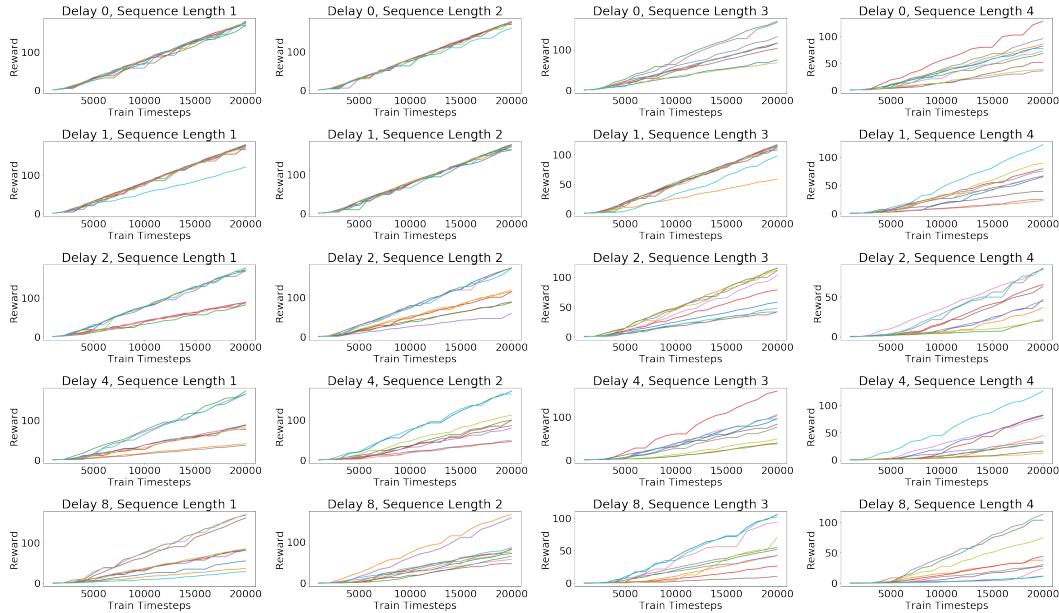


Figure 13: Training Learning Curves for Rainbow when varying delay and specific sequences. Please note the different Y-axis scales.

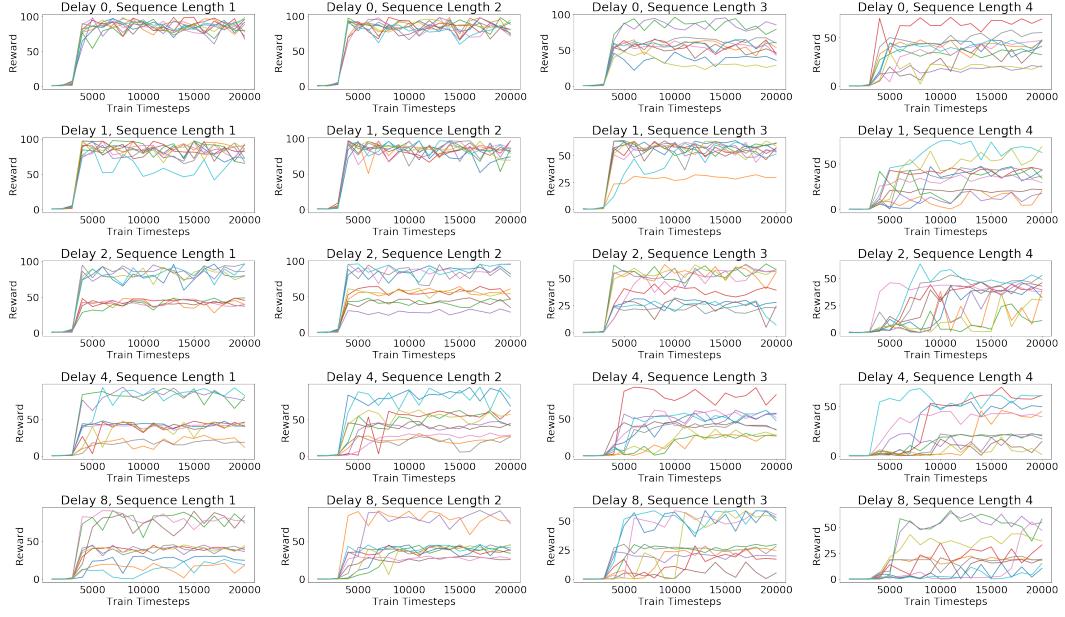


Figure 14: Evaluation Learning Curves for Rainbow when varying delay and specific sequences. Please note the different Y-axis scales.

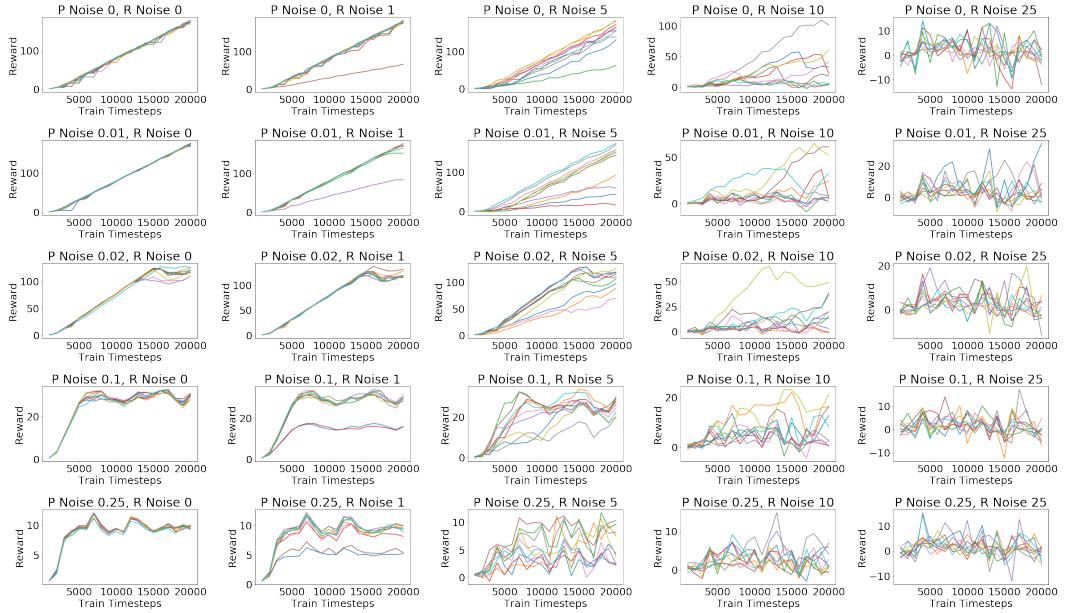
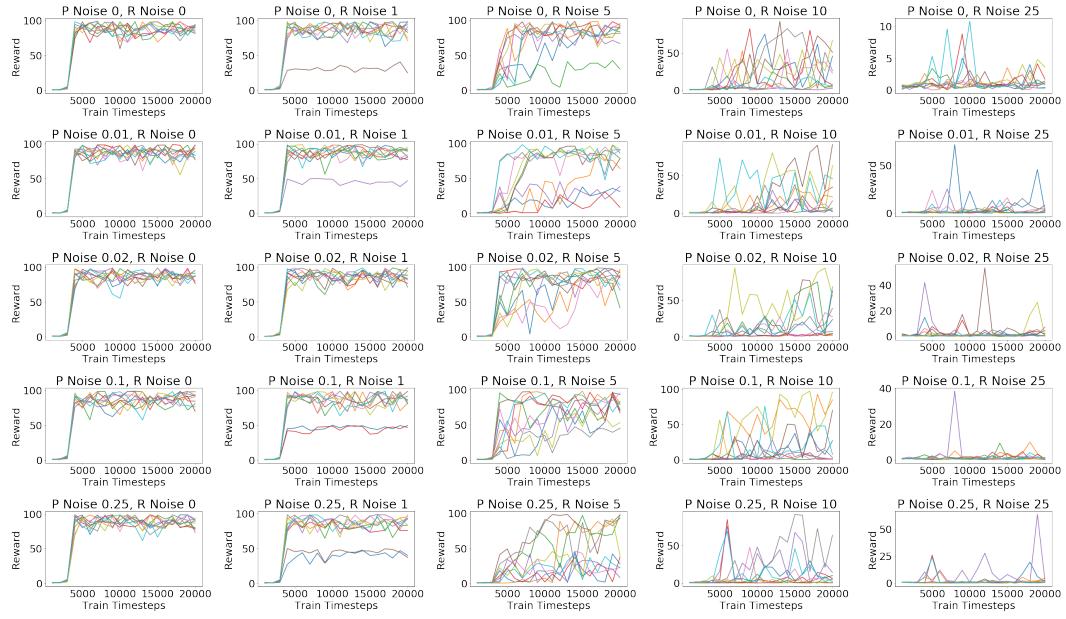


Figure 15: Training Learning Curves for Rainbow when varying noises. Please note the different Y-axis scales.



**Figure 16:** Evaluation Learning Curves for Rainbow **when varying noises**. Please note the different Y-axis scales.

3 C Plots for Making Rewards Denser

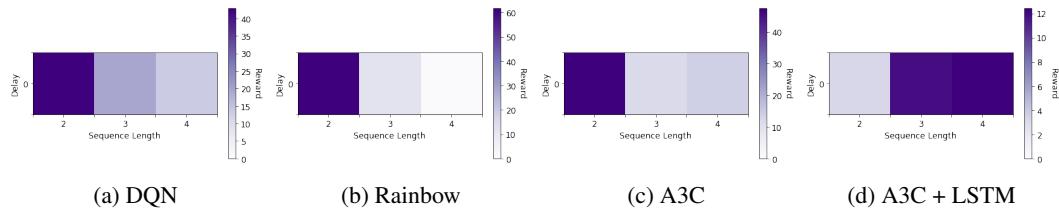


Figure 17: Mean episodic evaluation rollout reward at the end of training for the different algorithms **when making reward for specific sequences denser**. Please note the different colorbar scales.

4 **D More on Conclusion and Future Work**

5 Another key meta-feature, or rather meta-feature related to specific sequences, is manifolds. Let's say  
6 every specific sequence that was rewarded had a length  $n$  and the reward density was, say, 0.1. If  
7 we had no prior knowledge of the environment and did completely random exploration, we would  
8 need  $O(n^{|S|})$  sequences to obtain a reward signal which is exponential in the state space size (if we  
9 ignore the action space). We need to have some prior knowledge about the manifolds that exist in  
10 the environment to explore in a directed manner. This we intend to infuse through having a model  
11 of the environment. We will initially create simple manifolds for different tasks, such as rewarding  
12 following a circle in a 3-D continuous space and test directed exploration strategies to see how these  
13 choose to explore.

14 The states and actions contained in a specific sequence could just be a single *compound* state and  
15 *compound* action if we discretised time in a suitable manner. This brings us to the idea of learning at  
16 multiple timescales. HRL algorithms with formulations like the options framework (Sutton et al.,  
17 1999), could try to identify these specific sequences at the higher level and then carry out "atomic"  
18 actions at the lower level.

19 We also hope to benchmark other algorithms like PPO<sup>1</sup> (Schulman et al., 2017), Rudder (Arjona-  
20 Medina et al., 2018), MCTS (Silver et al., 2016) and table-based algorithms and to show theoretical  
21 results match with practice on toy benchmarks.

22 We also aim to promote reproducibility in RL as in (Henderson et al., 2017) and hope our benchmark  
23 helps with that goal.

24 We need different RL algorithms for different environments. Aside from some basic heuristics  
25 such as applying DDPG (Lillicrap et al., 2015) to continuous environments and DQN to discrete  
26 environments, it is not very clear when to use which RL algorithms. We hope this will be a first  
27 step to being able to identify from the environment what sort of algorithm to use and to help build  
28 adaptive algorithms which adapt to the environment at hand. Additionally, aside from being a great  
29 benchmark for RL algorithms, it is also a great didactic tool for teaching how RL algorithms work in  
30 different environments.

31 **References**

- 32 Arjona-Medina, J. A., Gillhofer, M., Widrich, M., Unterthiner, T., Brandstetter, J., and Hochreiter, S.  
33 (2018). Rudder: Return decomposition for delayed rewards.
- 34 Henderson, P., Islam, R., Bachman, P., Pineau, J., Precup, D., and Meger, D. (2017). Deep reinforcement  
35 learning that matters.
- 36 Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D.  
37 (2015). Continuous control with deep reinforcement learning.
- 38 Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. (2017). Proximal policy  
39 optimization algorithms.
- 40 Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J.,  
41 Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. (2016). Mastering the game of go with  
42 deep neural networks and tree search. *nature*, 529(7587):484.
- 43 Sutton, R. S., Precup, D., and Singh, S. (1999). Between mdps and semi-mdps: A framework for  
44 temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2):181–211.

---

<sup>1</sup>We tried PPO but could not get it to learn