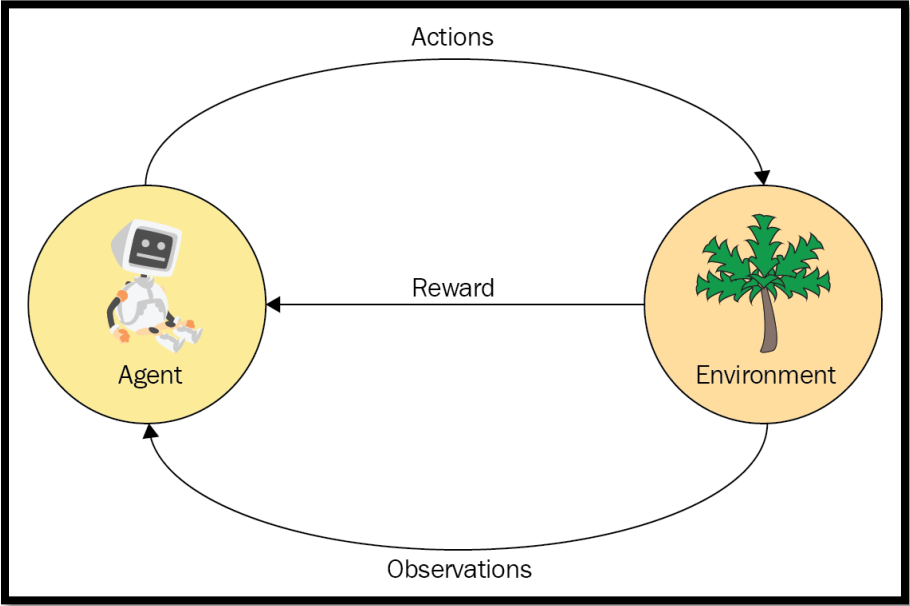
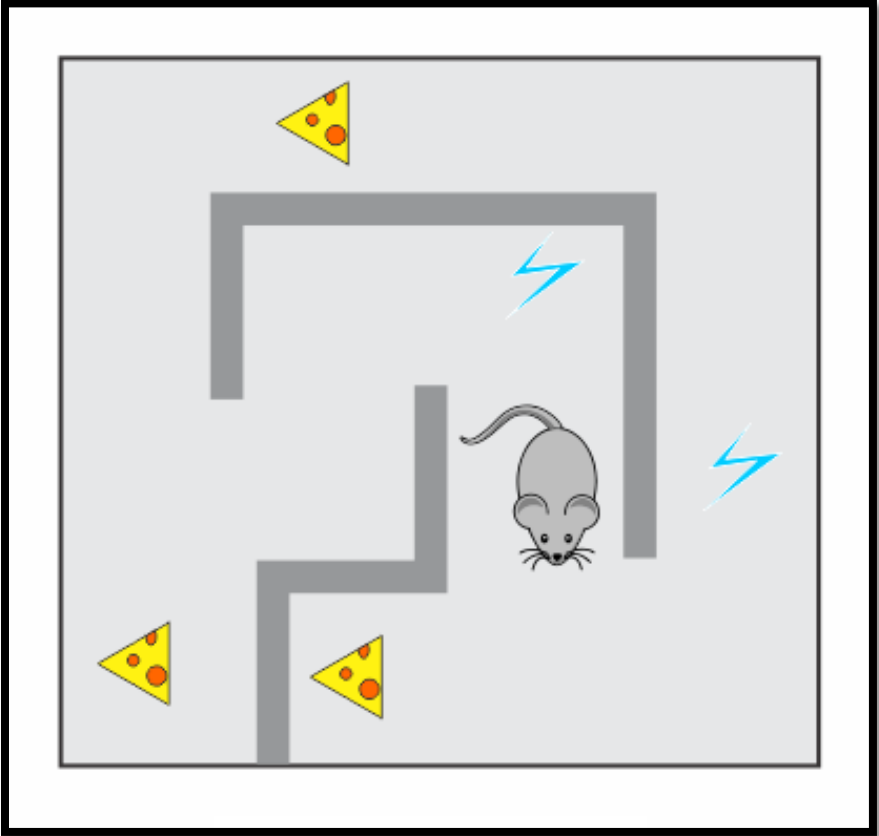
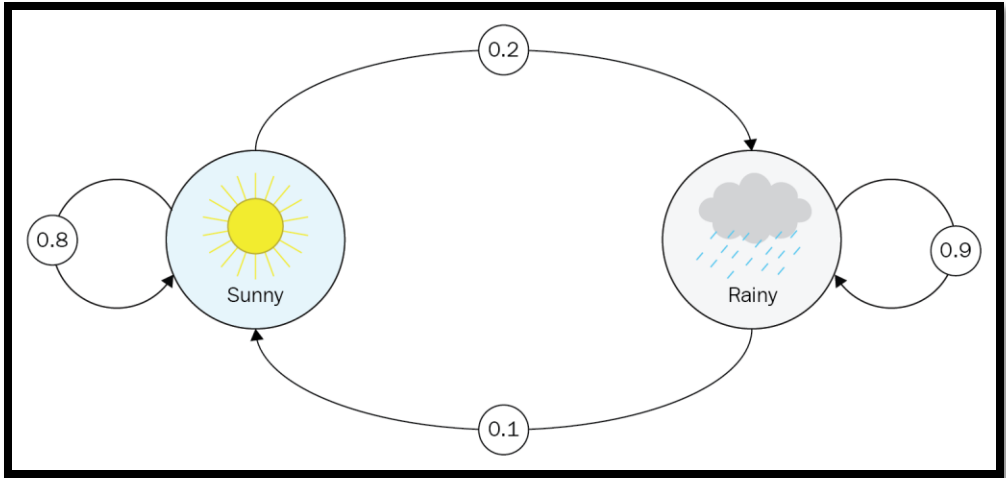
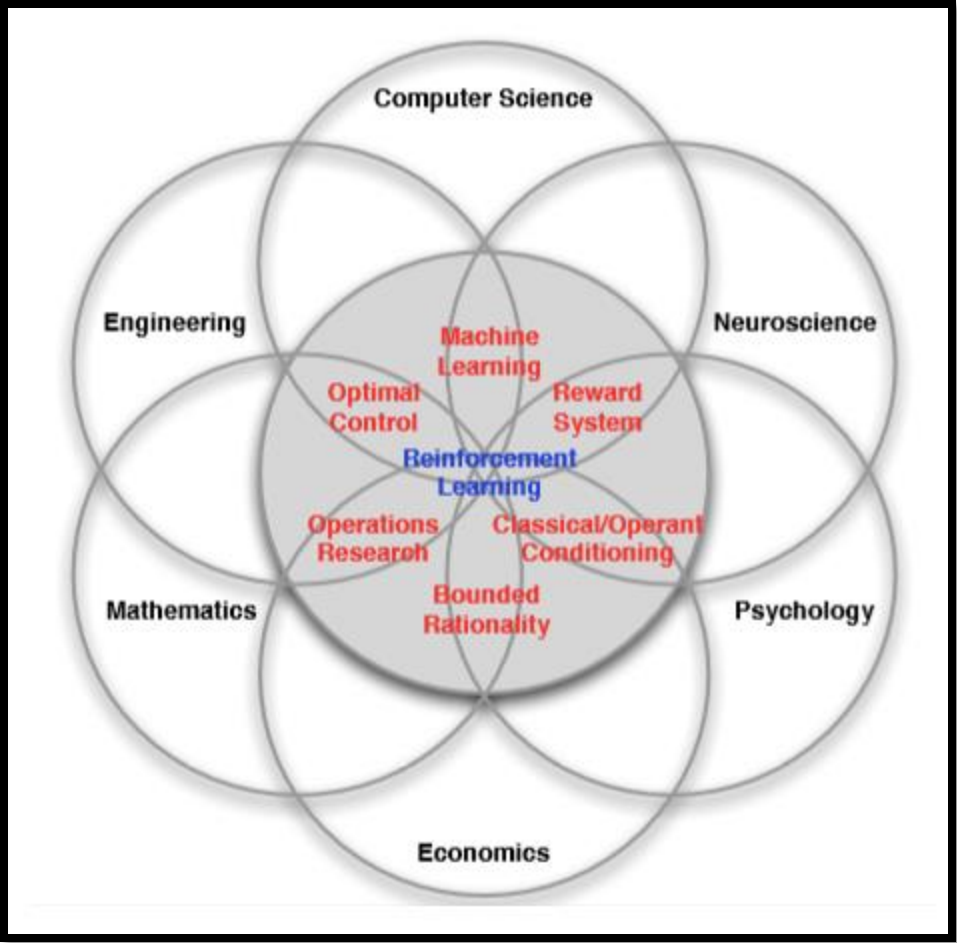
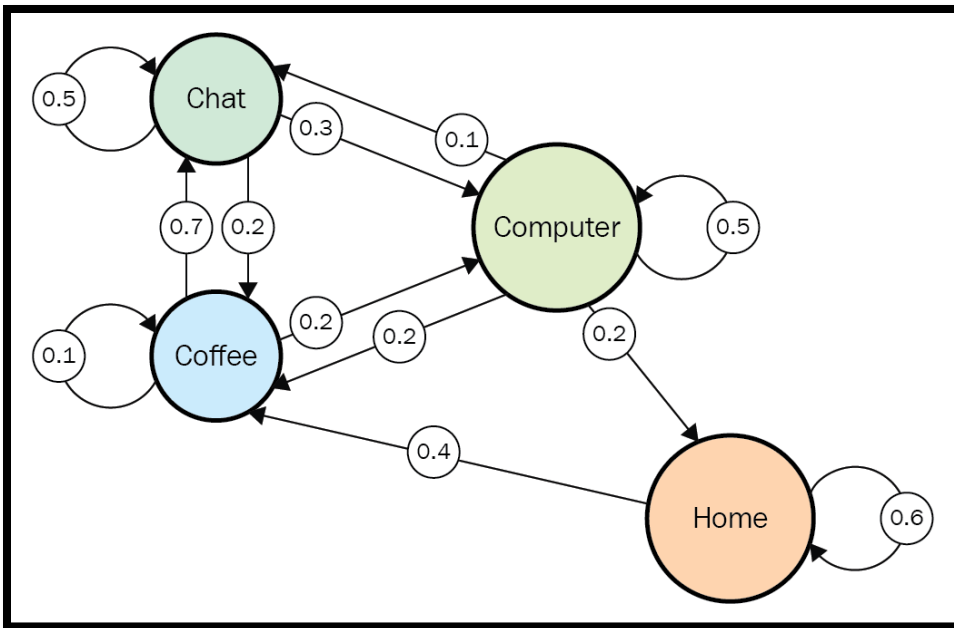
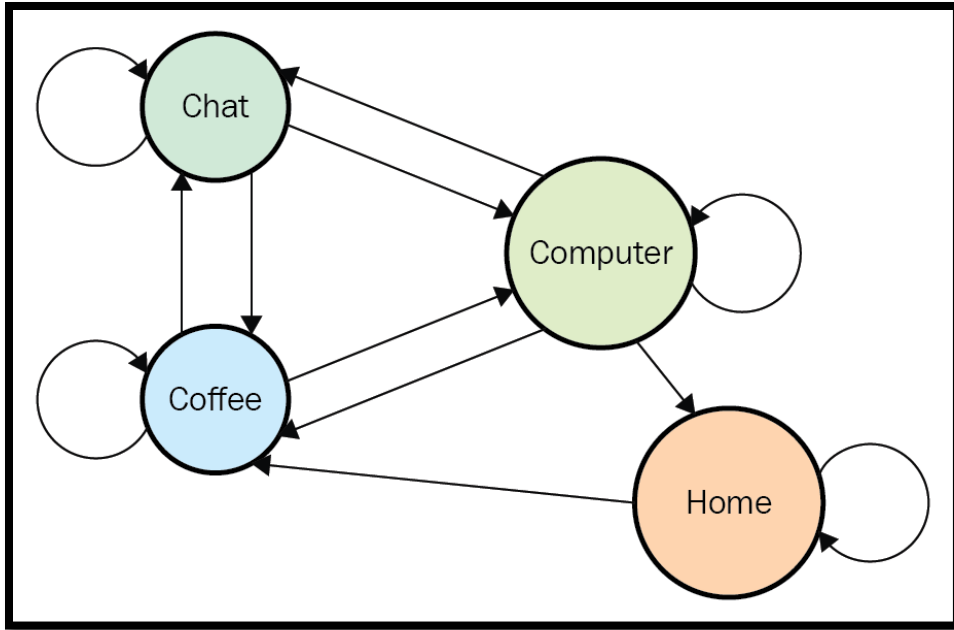
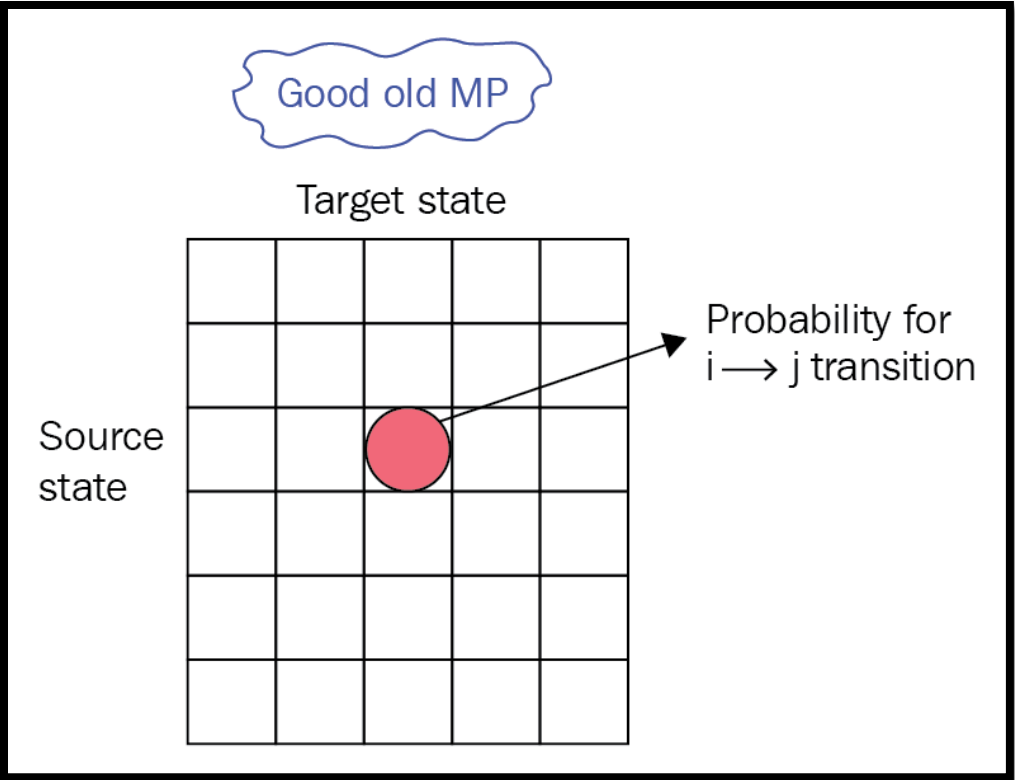
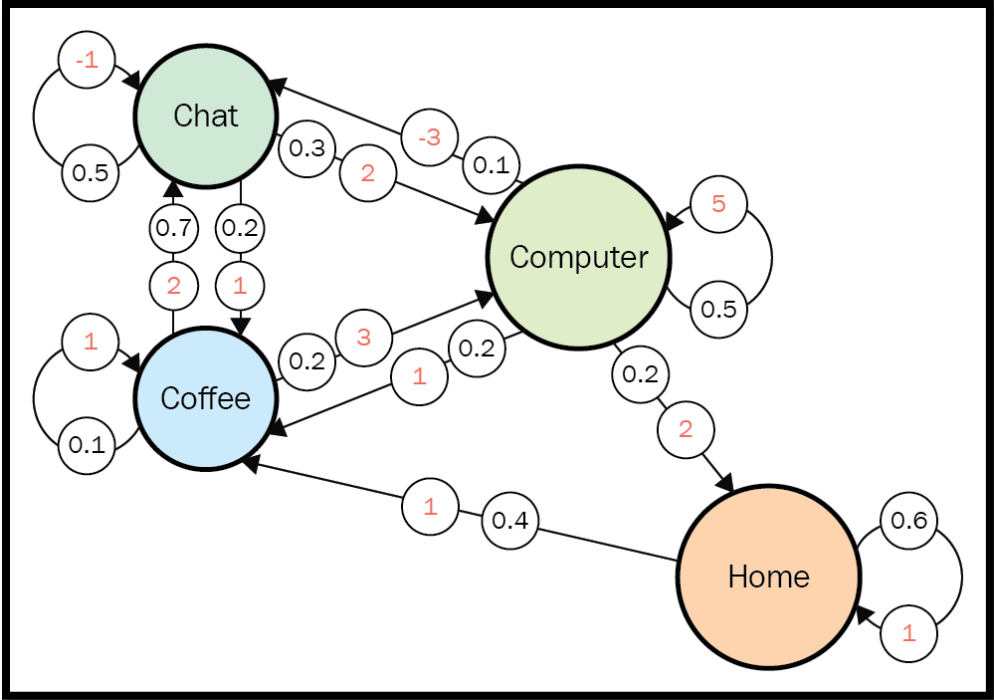


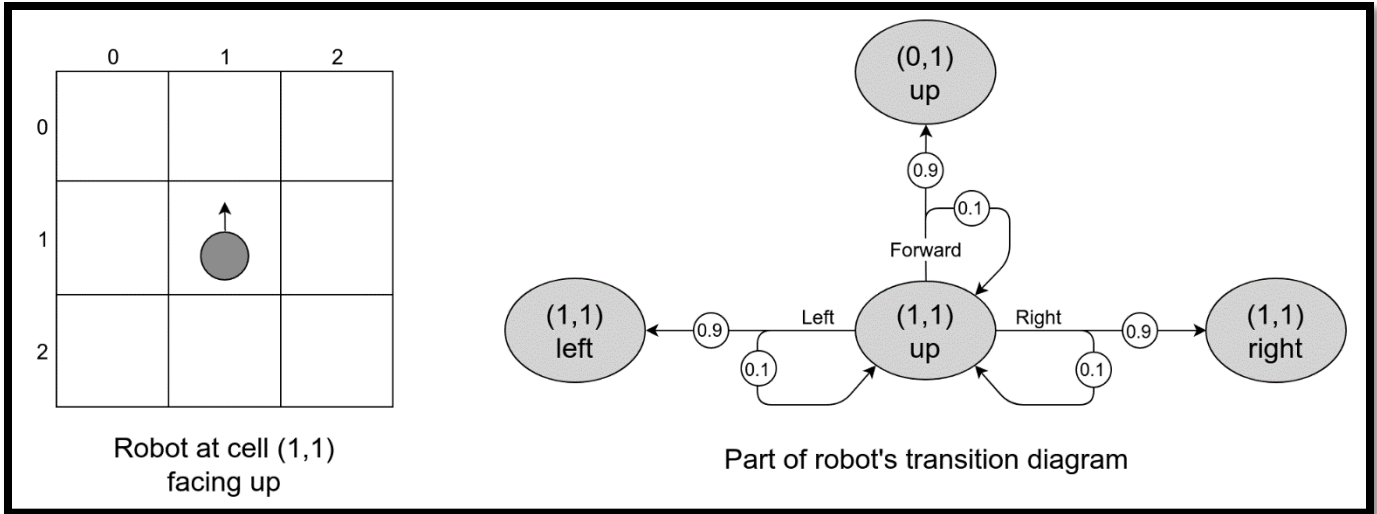
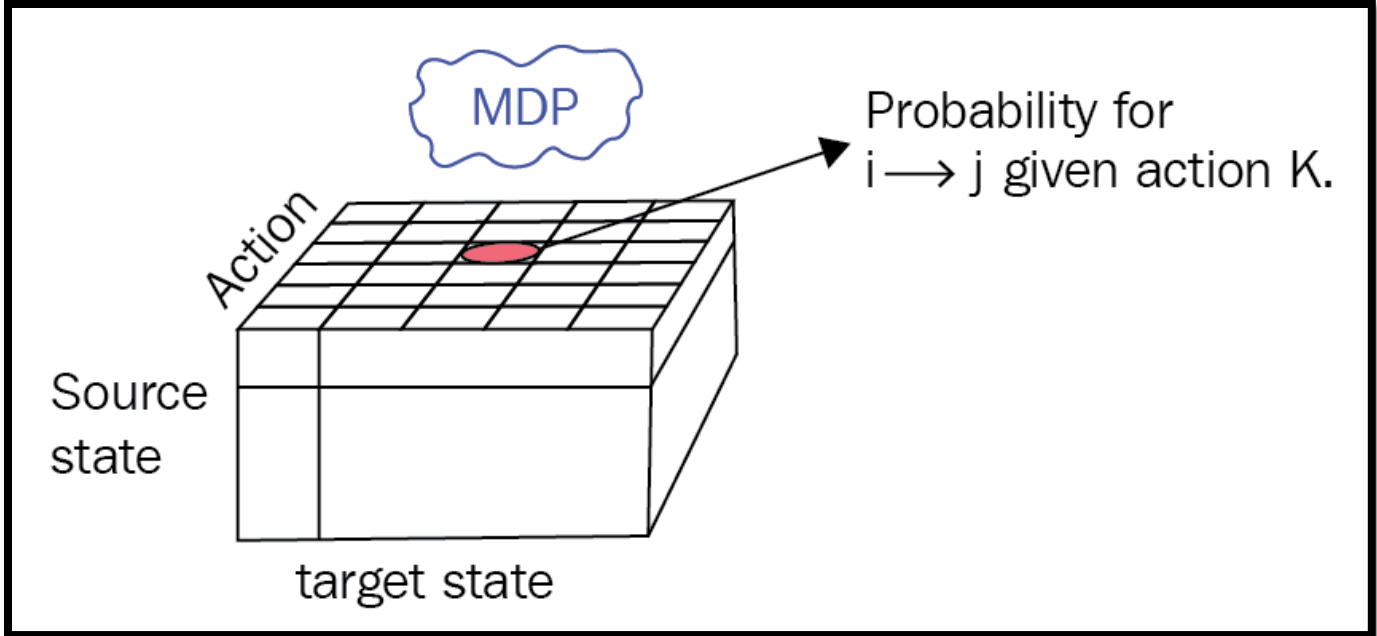
# Chapter 1: What is Reinforcement Learning?



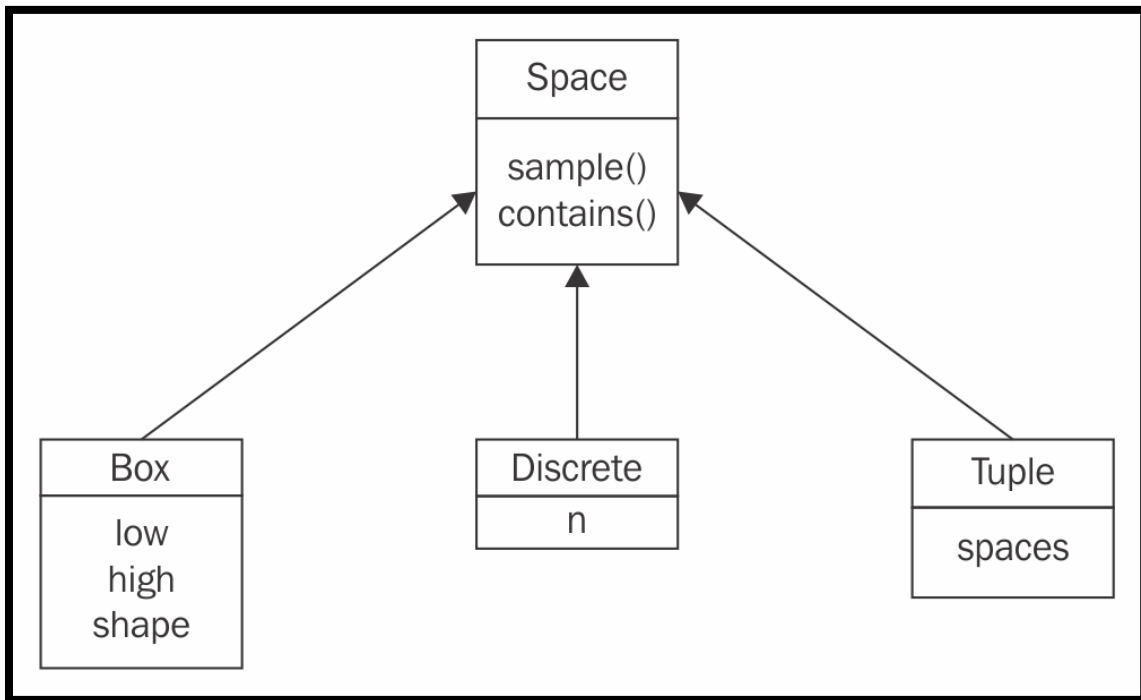


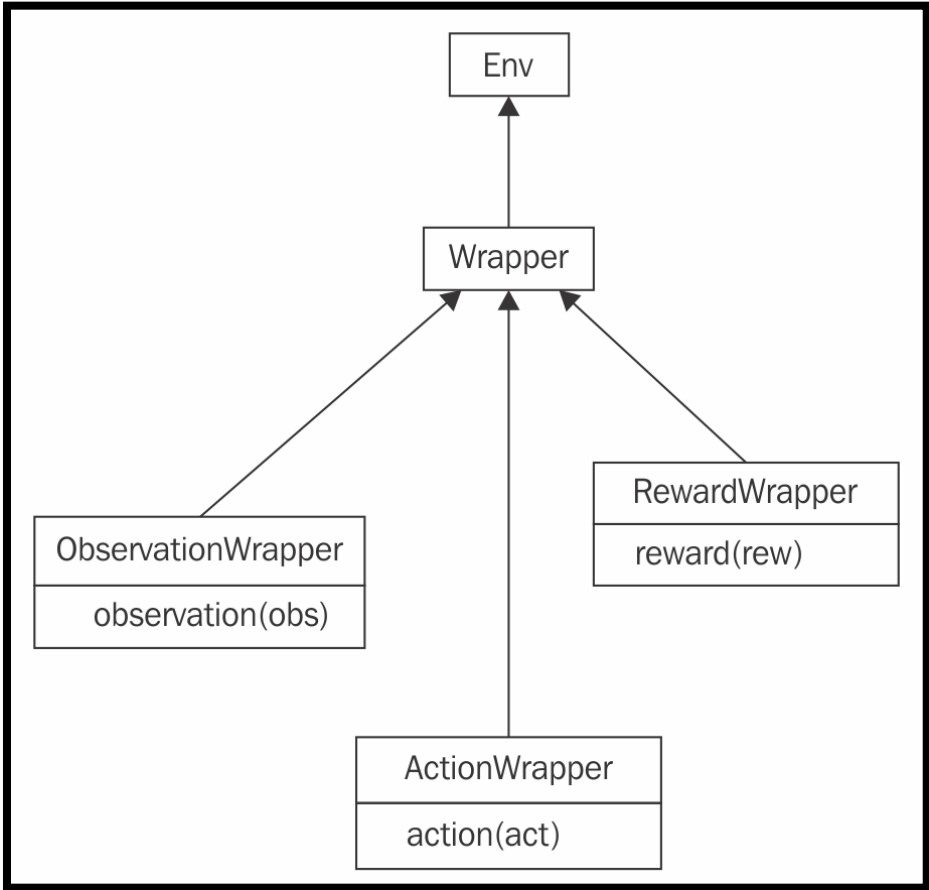
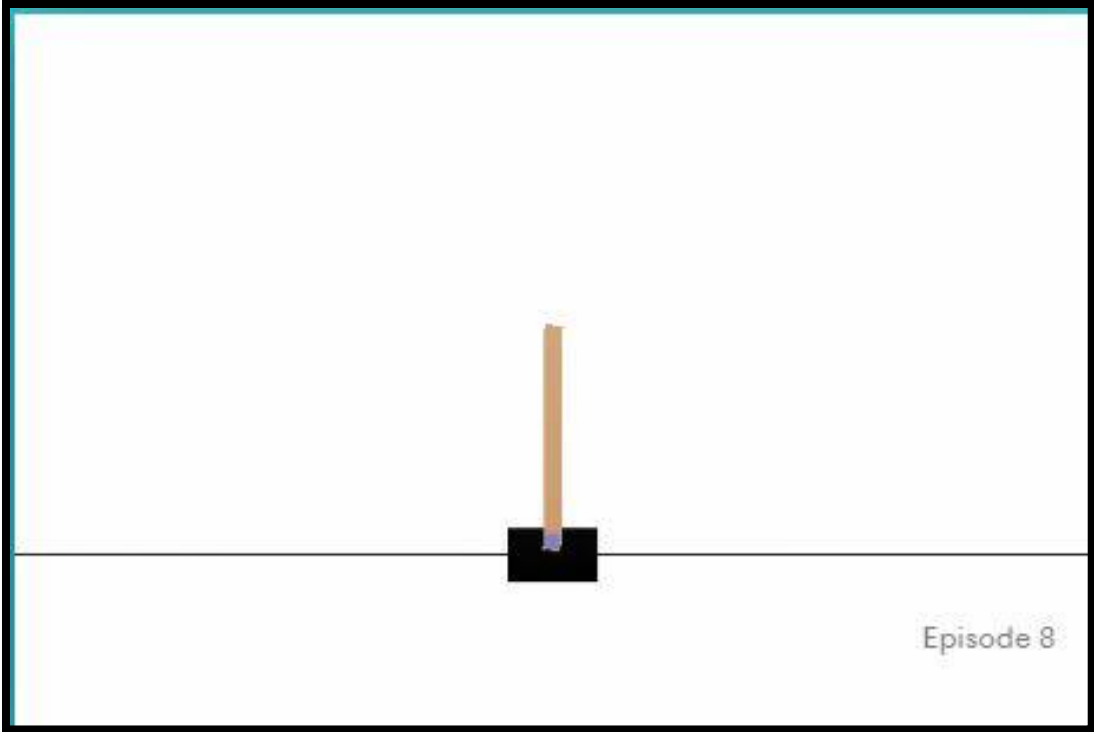






## Chapter 2: OpenAI Gym



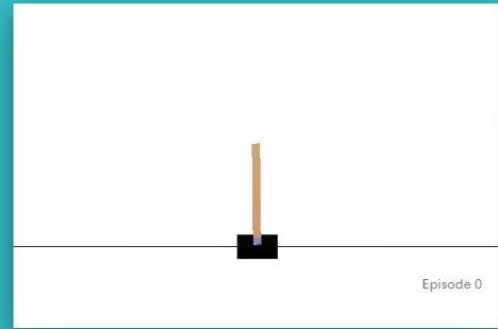


## CartPole-v0

A pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The system is controlled by applying a force of +1 or -1 to the cart. The pendulum starts upright, and the goal is to prevent it from falling over. A reward of +1 is provided for every timestep that the pole remains upright. The episode ends when the pole is more than 15 degrees from vertical, or the cart moves more than 2.4 units from the center.

*CartPole-v0 defines "solving" as getting average reward of 195.0 over 100 consecutive trials.*

*This environment corresponds to the version of the cart-pole problem described by Barto, Sutton, and Anderson [Barto83].*



karpathy's algorithm, Took 211 episodes to solve the environment. 195.27 ± 1.57 a year ago

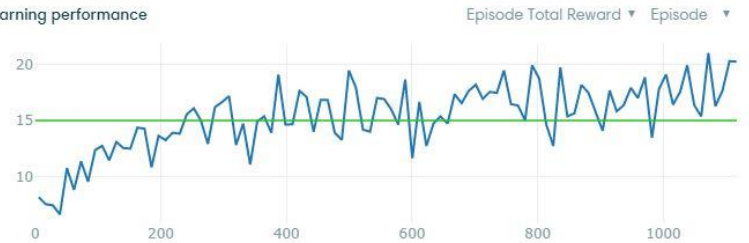
### CartPole-v0 Evaluations

ALGORITHM	EPISODES BEFORE SOLVE	SUBMITTED
n1try's algorithm <a href="#">writeup</a>	85.0	11 days ago
mbalunovic's algorithm <a href="#">writeup</a>	306.0	11 days ago
ruippeixotog's algorithm <a href="#">writeup</a>	933.0	12 days ago
ruippeixotog's algorithm <a href="#">writeup</a>	961.0	13 days ago





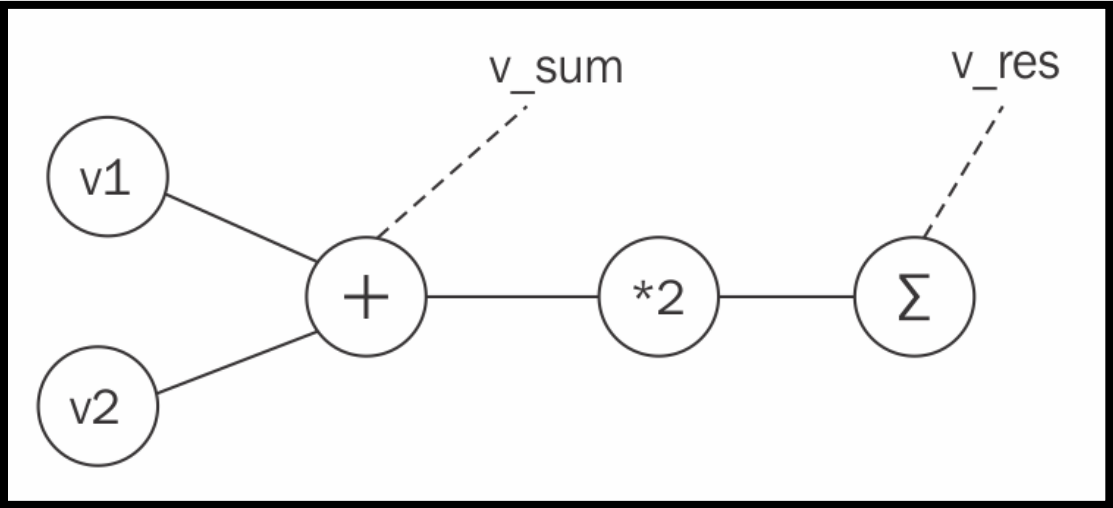
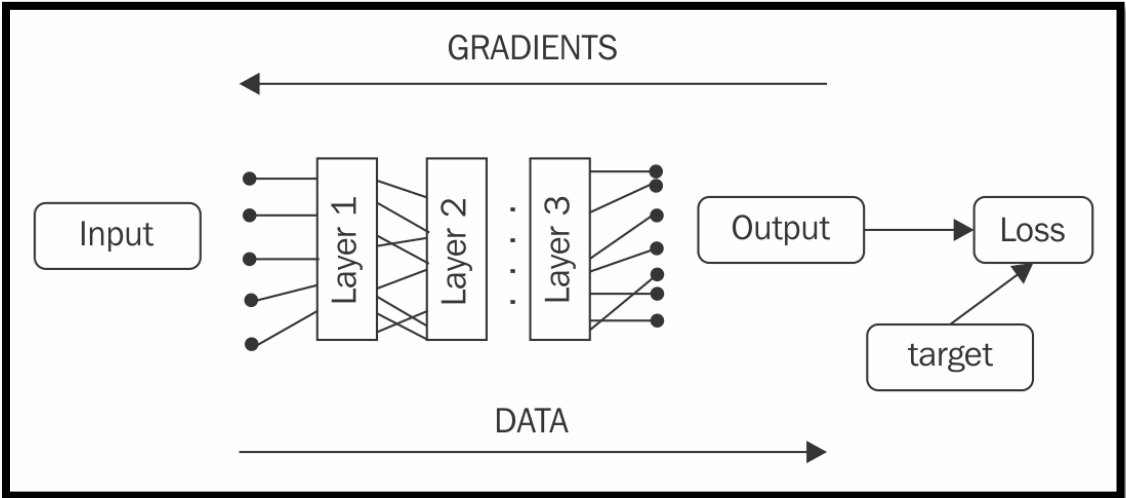
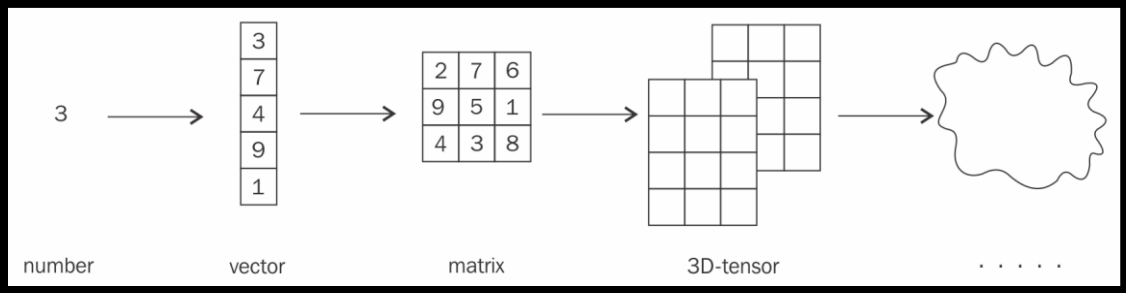
## Learning performance

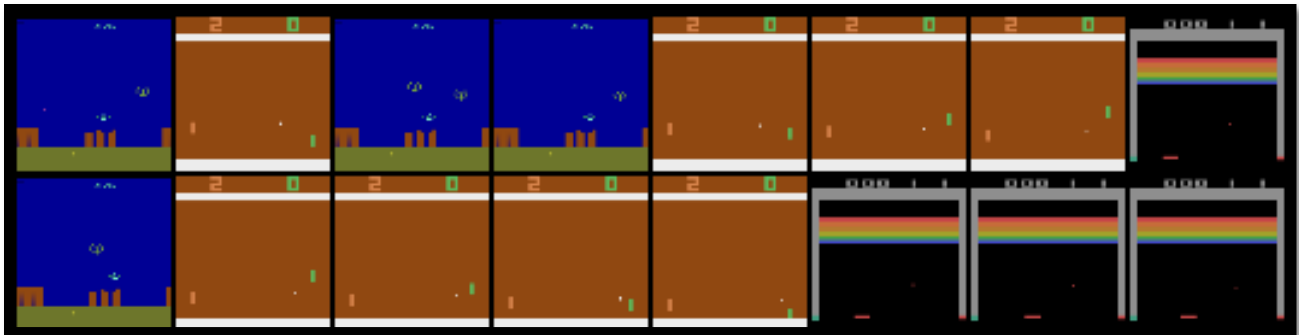
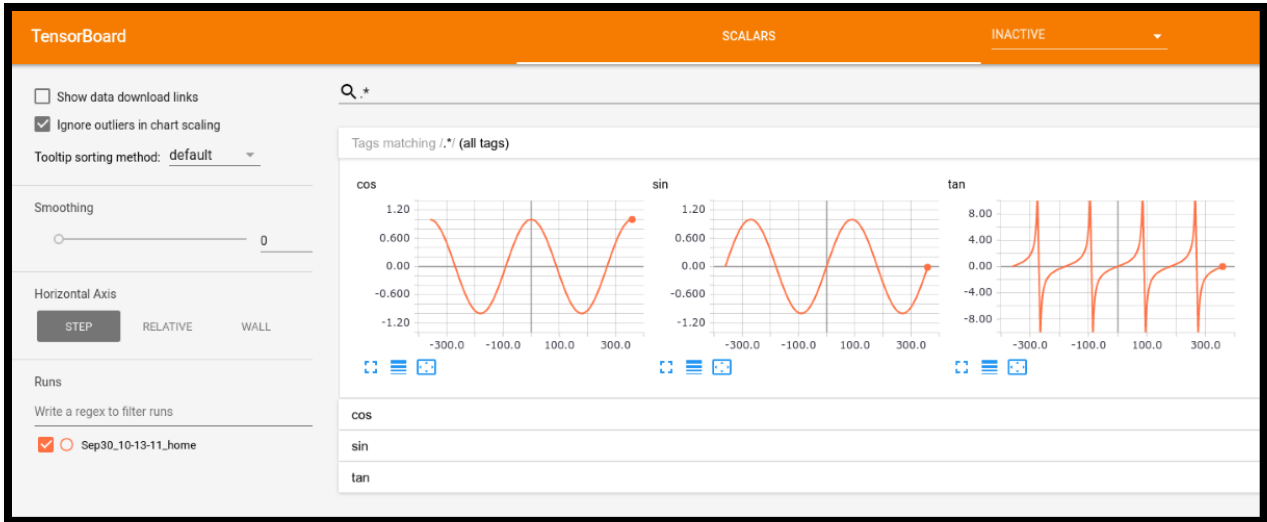
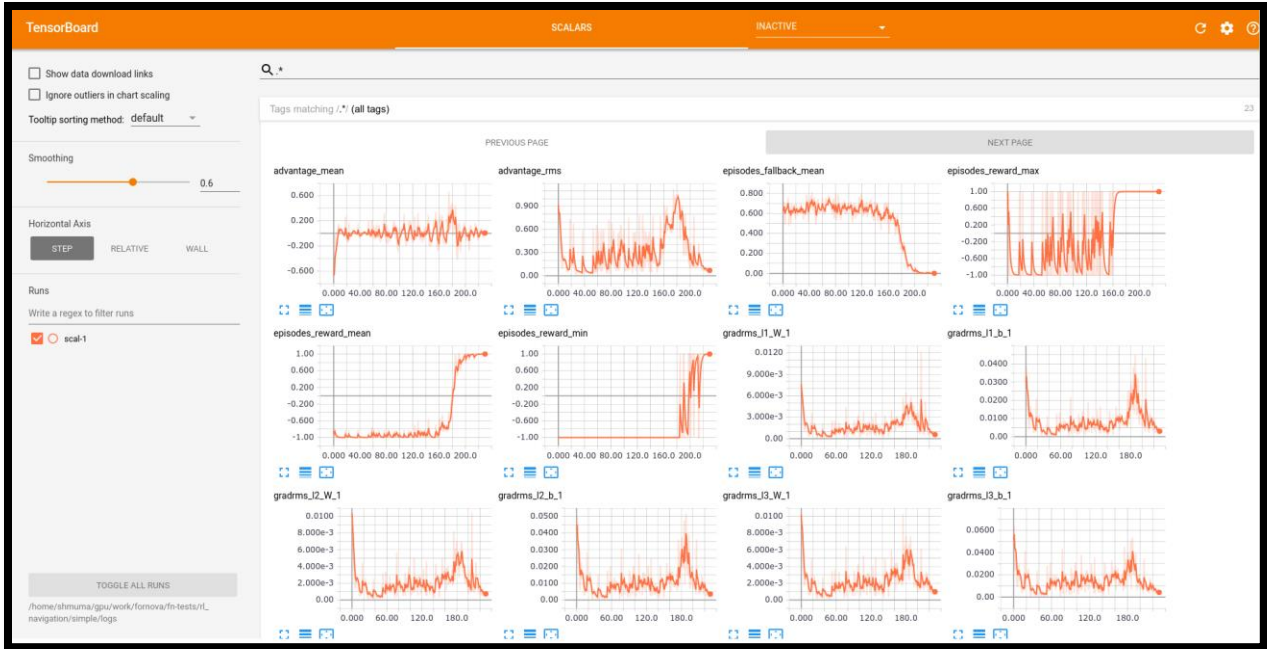


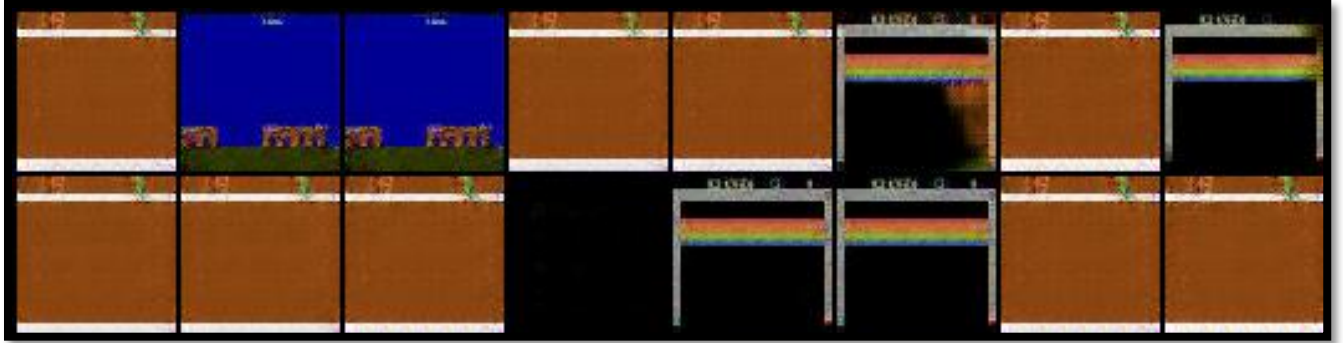
Solved after 205 episodes. Best 100-episode average reward was  $18.42 \pm 0.76$ .  
(DoomDefendLine-v0 is considered "solved" when the agent obtains an average reward of at least 15.0 over 100 consecutive episodes.)

205	1123	✓	40m	↓	🐦
Episodes to solve	Total episodes	Solved	Time to solve	Download	Tweet

# Chapter 3: Deep Learning with PyTorch



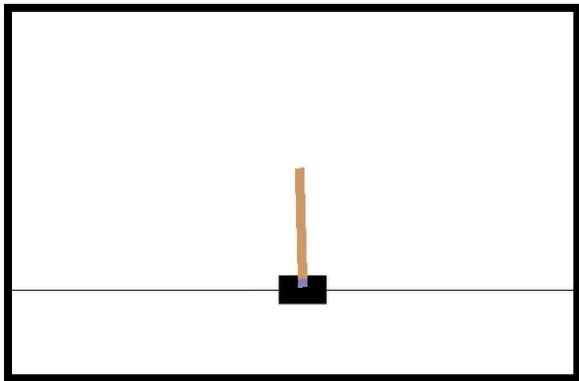
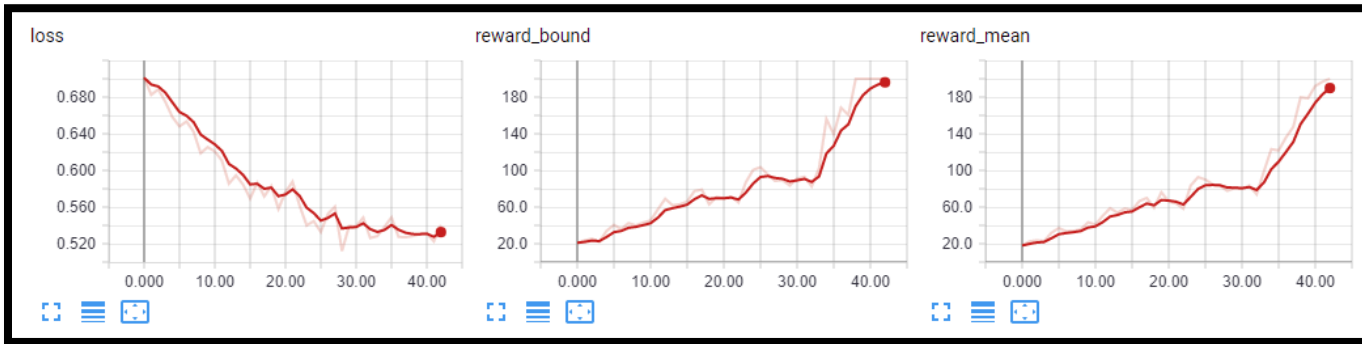


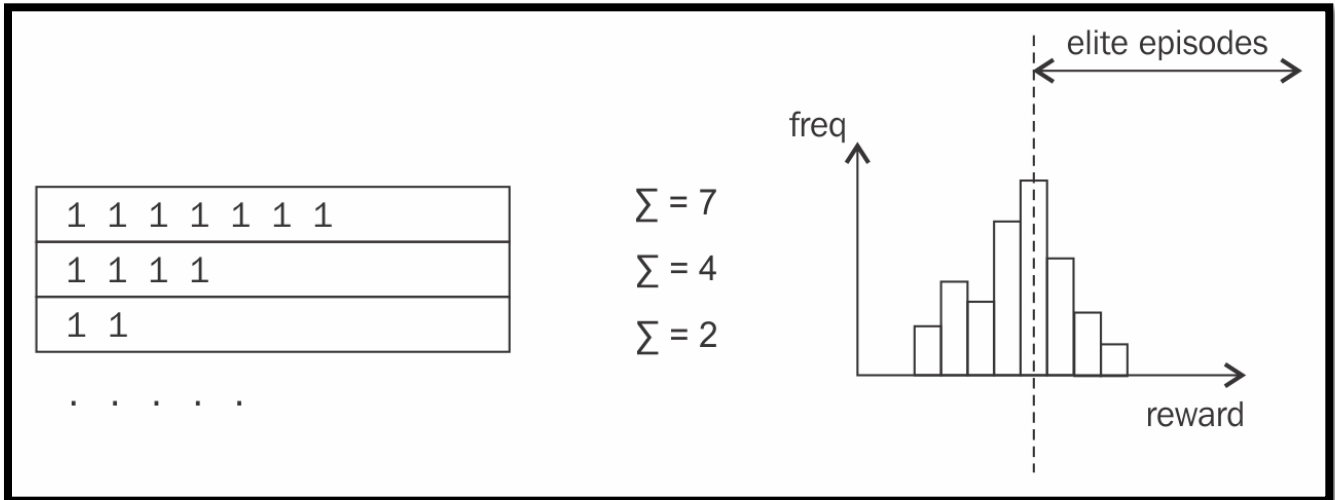
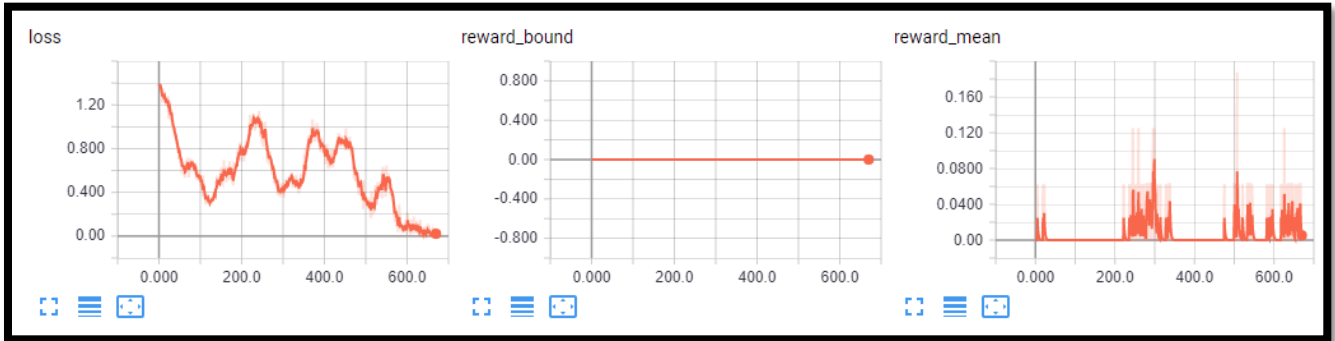
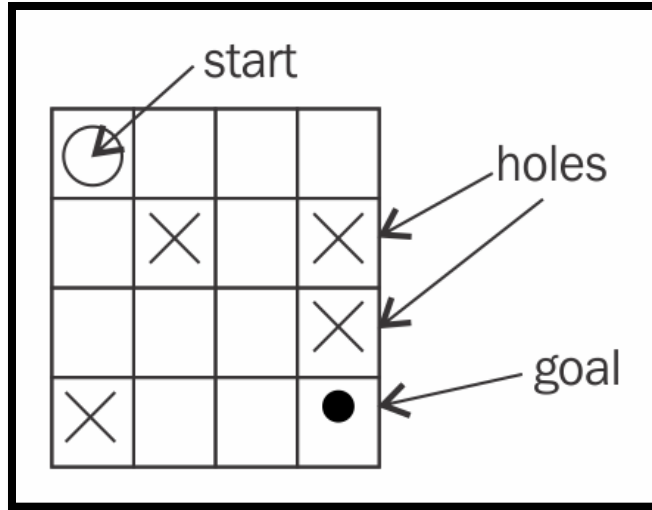


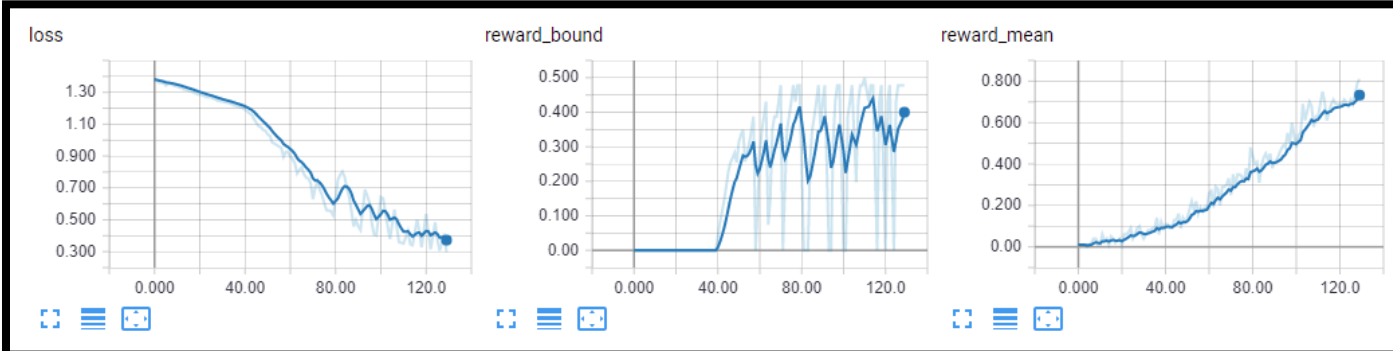
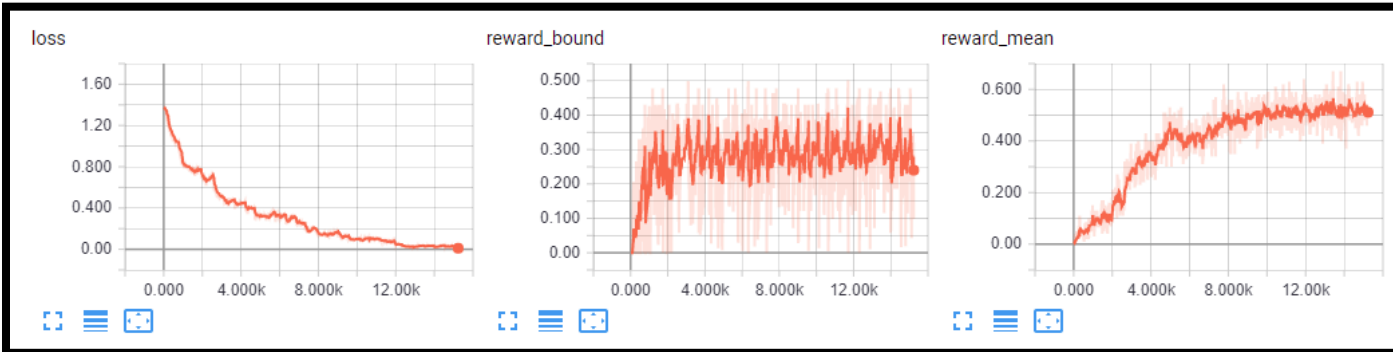
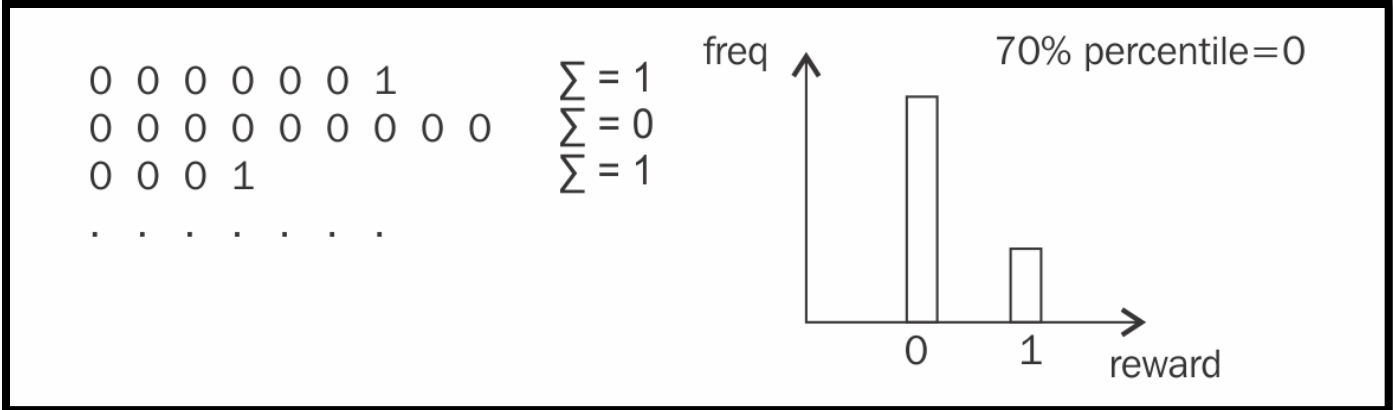
# Chapter 4: The Cross-Entropy Method

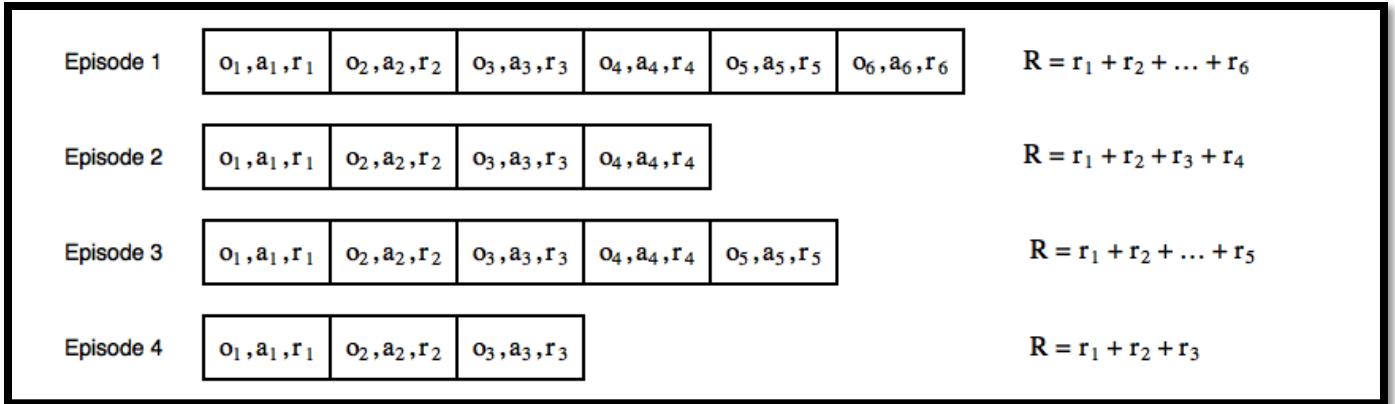
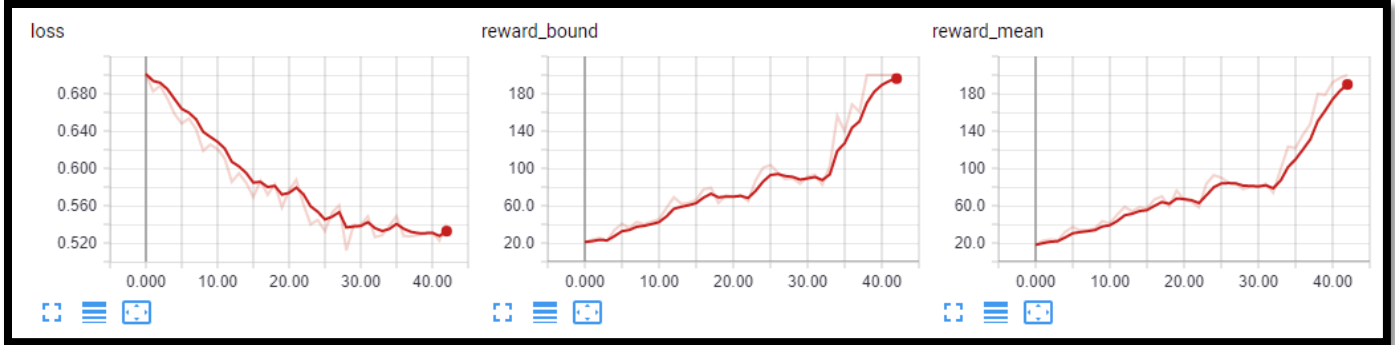


Episode 1	$o_1, a_1, r_1$	$o_2, a_2, r_2$	$o_3, a_3, r_3$	$o_4, a_4, r_4$	$o_5, a_5, r_5$	$o_6, a_6, r_6$	$R = r_1 + r_2 + \dots + r_6$
Episode 2	$o_1, a_1, r_1$	$o_2, a_2, r_2$	$o_3, a_3, r_3$	$o_4, a_4, r_4$			$R = r_1 + r_2 + r_3 + r_4$
Episode 3	$o_1, a_1, r_1$	$o_2, a_2, r_2$	$o_3, a_3, r_3$	$o_4, a_4, r_4$	$o_5, a_5, r_5$		$R = r_1 + r_2 + \dots + r_5$
Episode 4	$o_1, a_1, r_1$	$o_2, a_2, r_2$	$o_3, a_3, r_3$				$R = r_1 + r_2 + r_3$



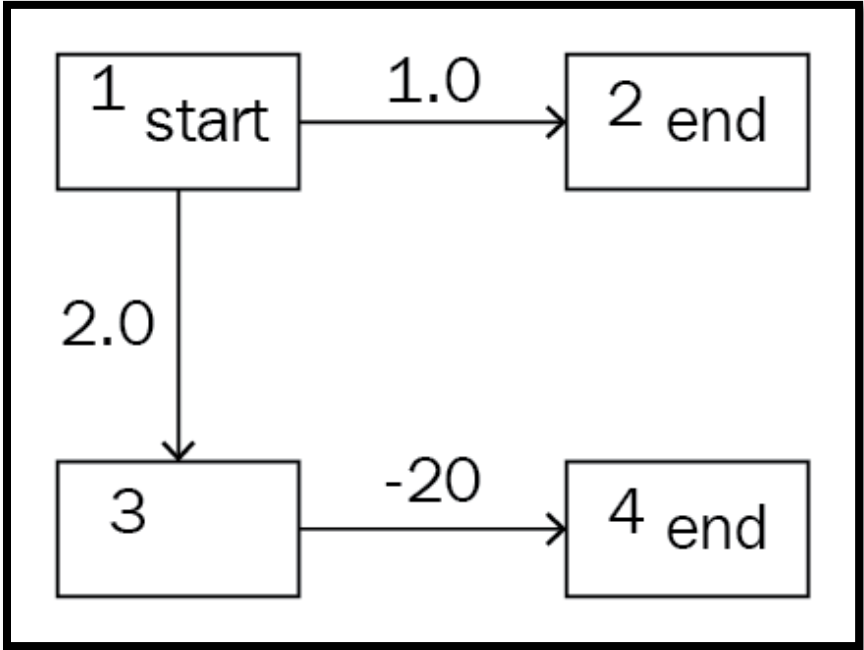
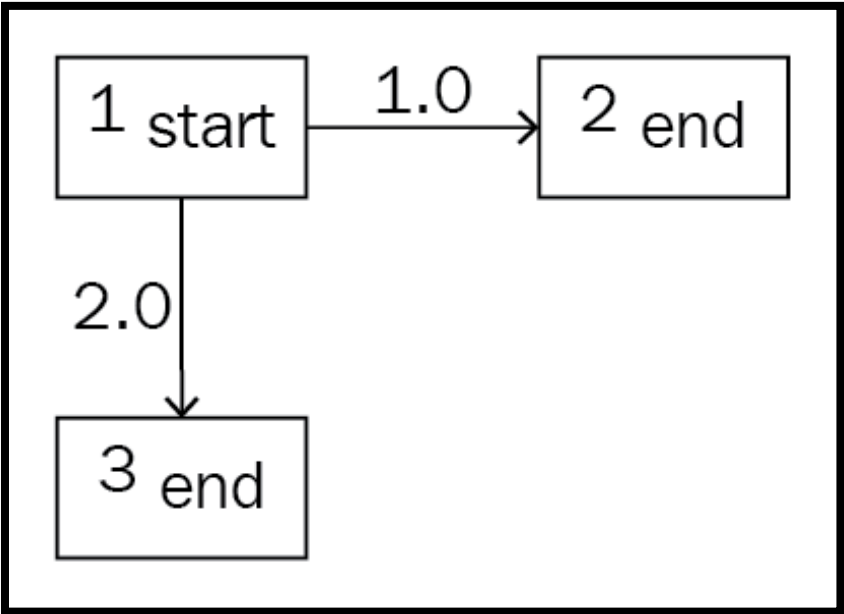


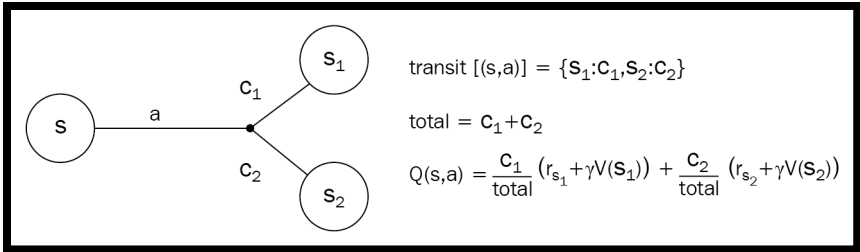
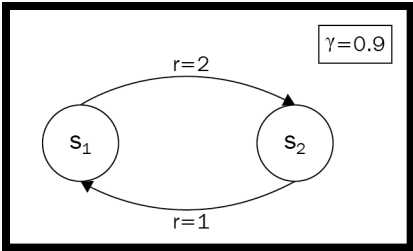
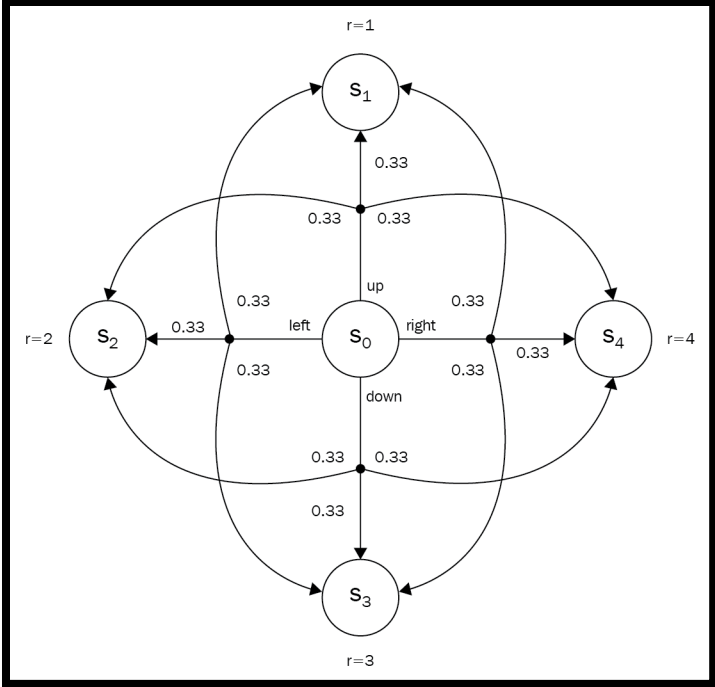
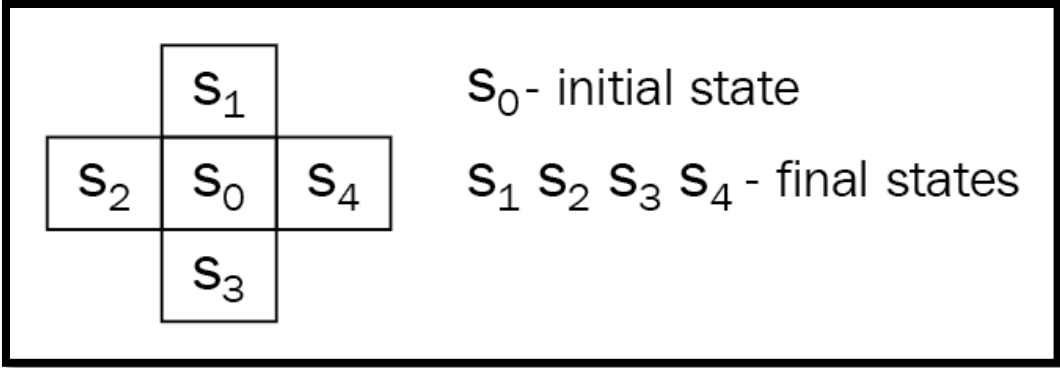


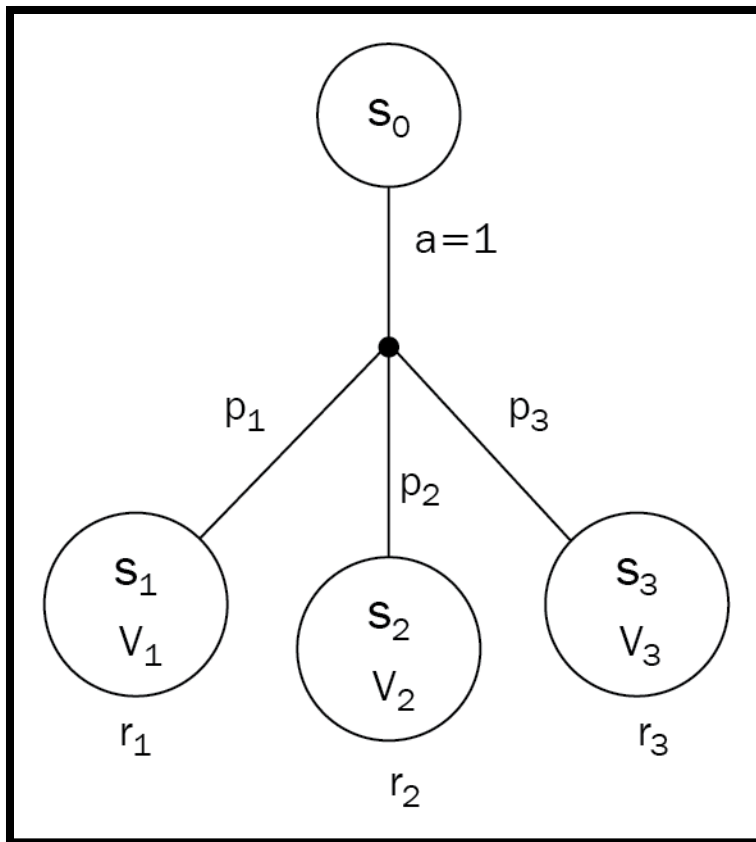
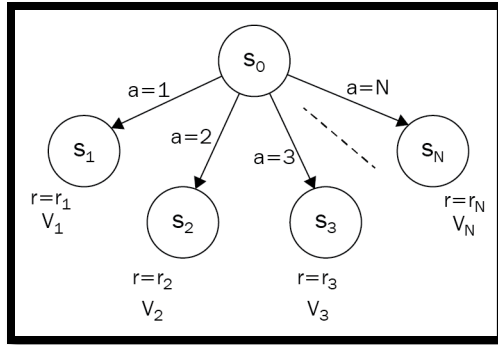




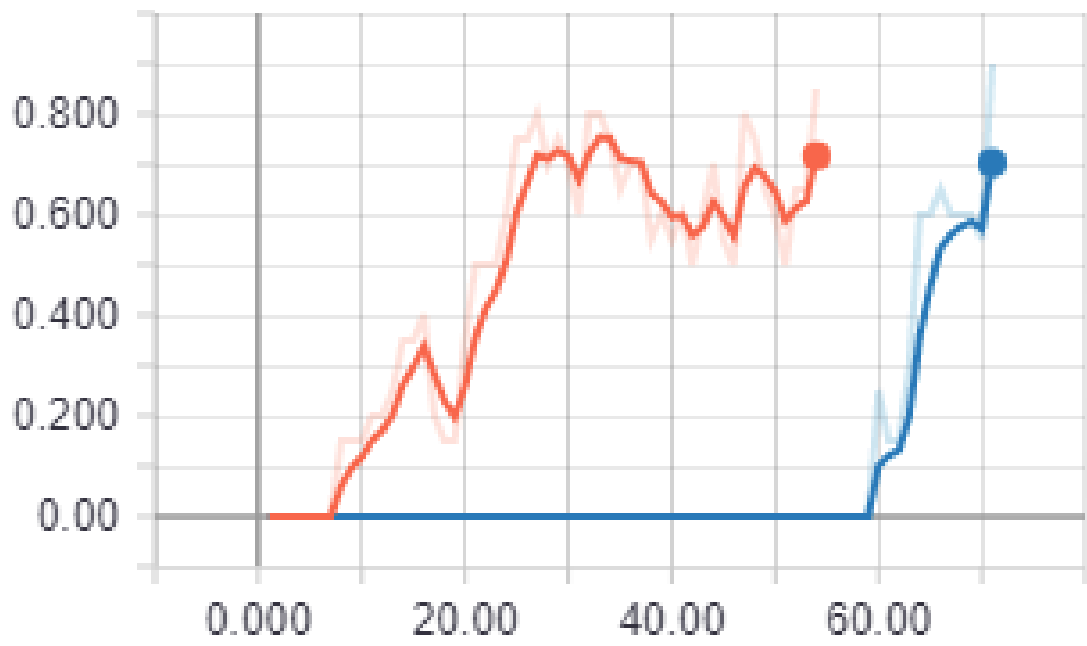
Chapter 5: Tabular Learning and the Bellman Equation



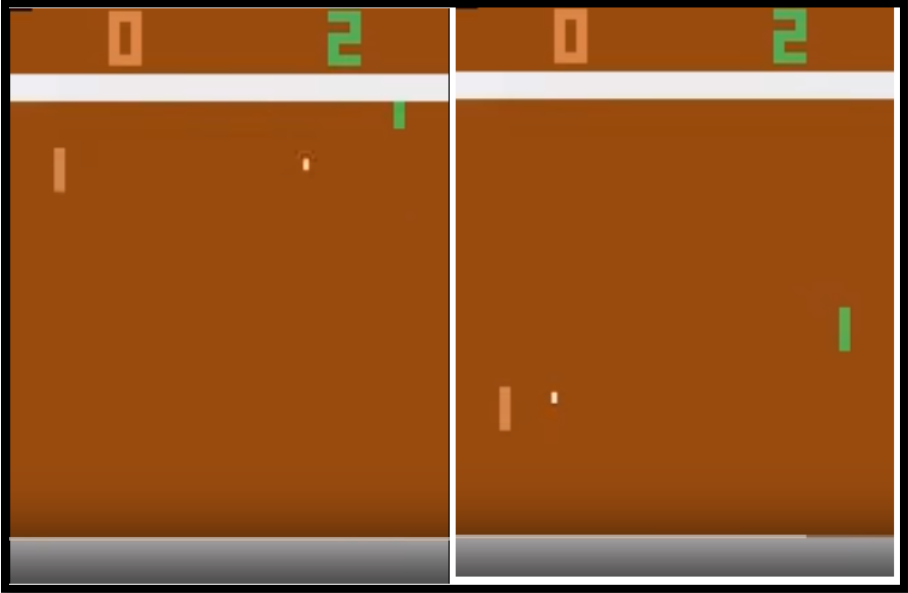
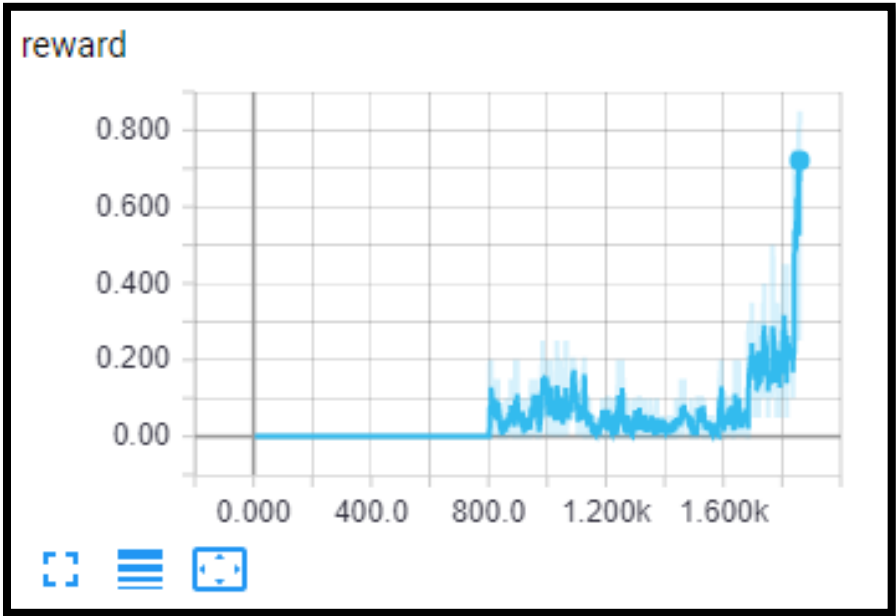


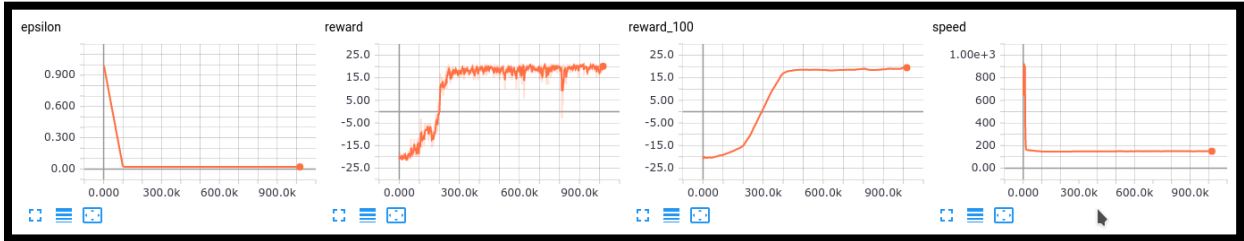
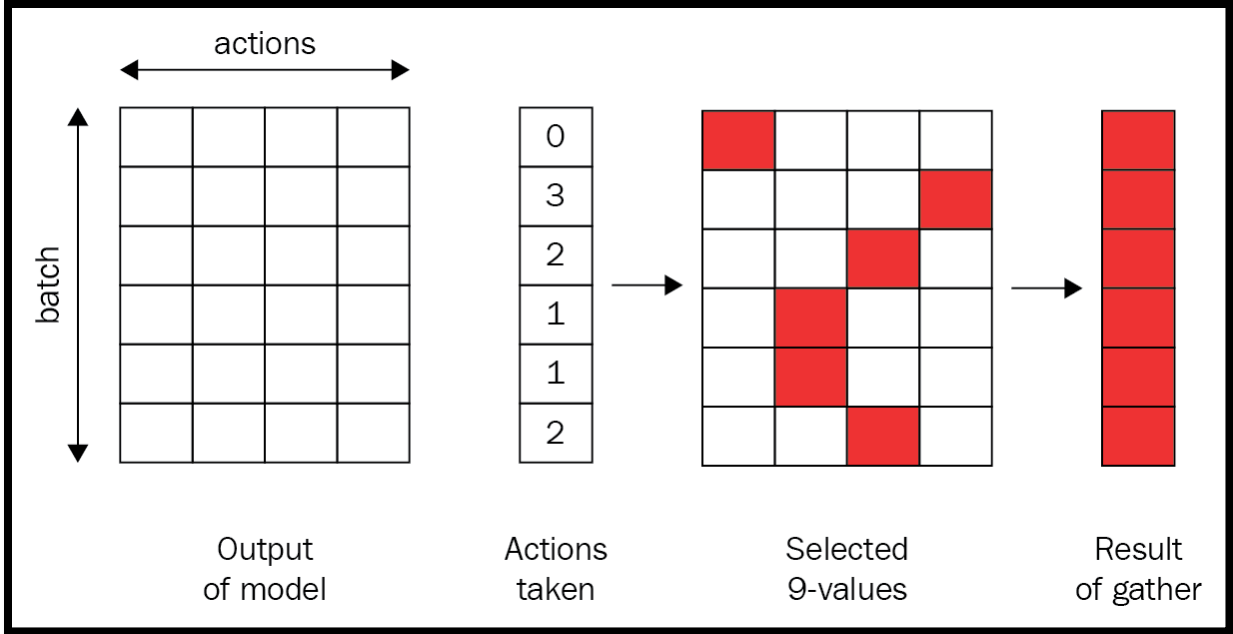


reward

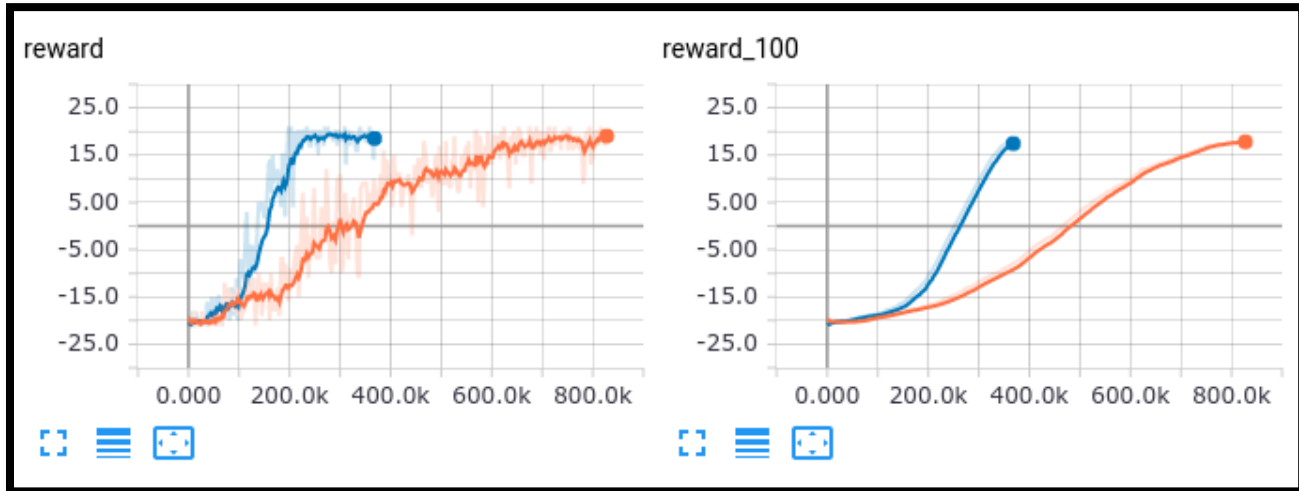
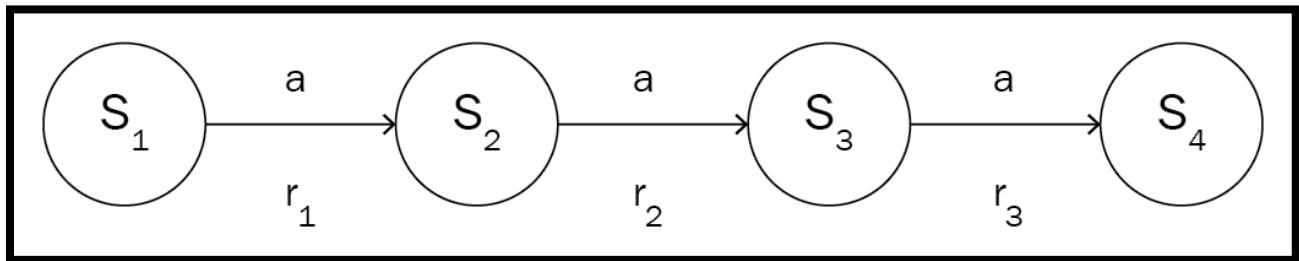
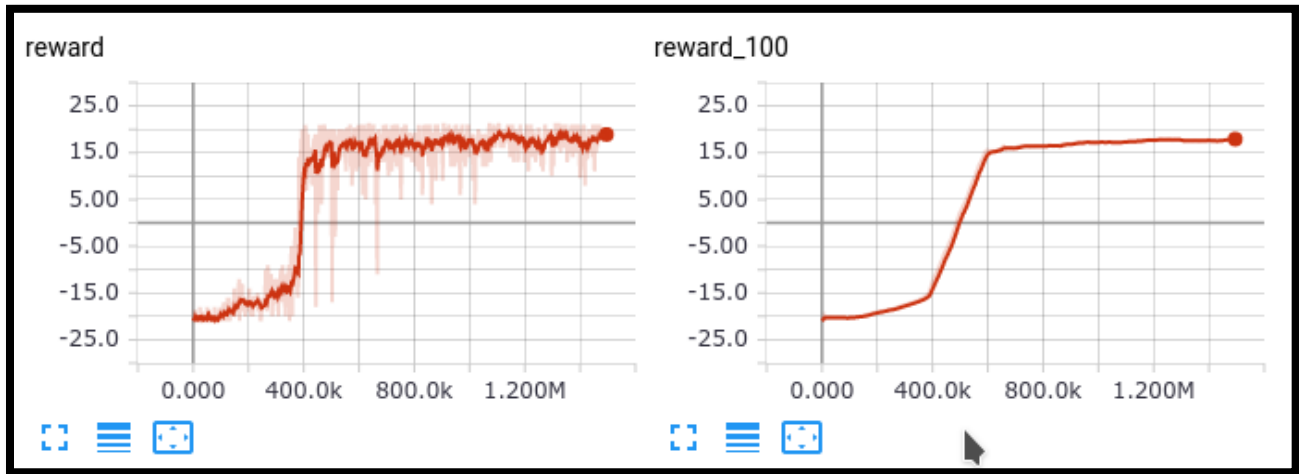


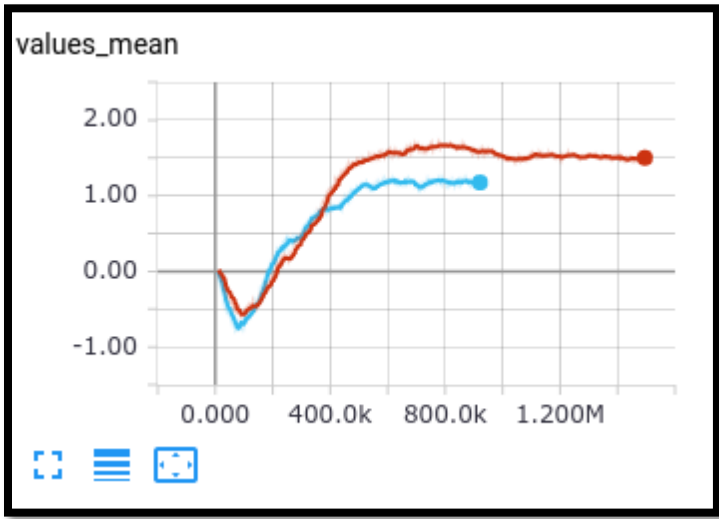
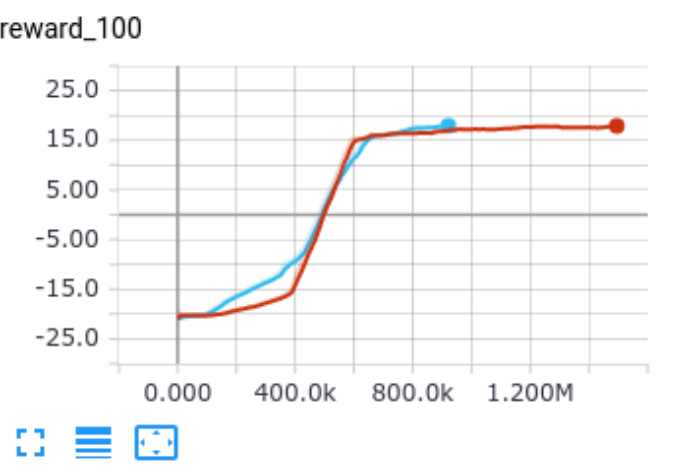
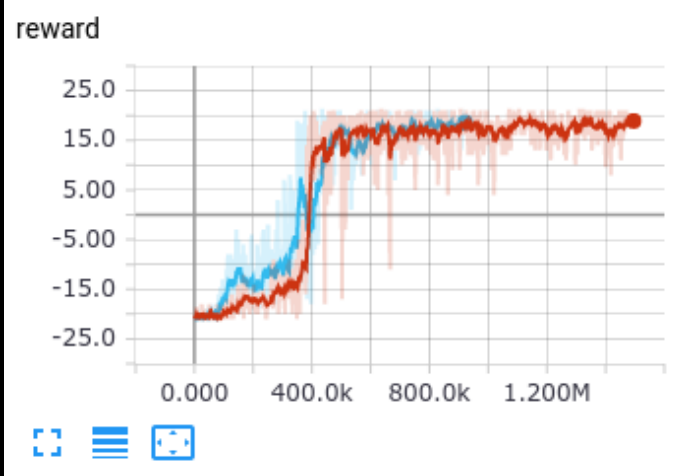
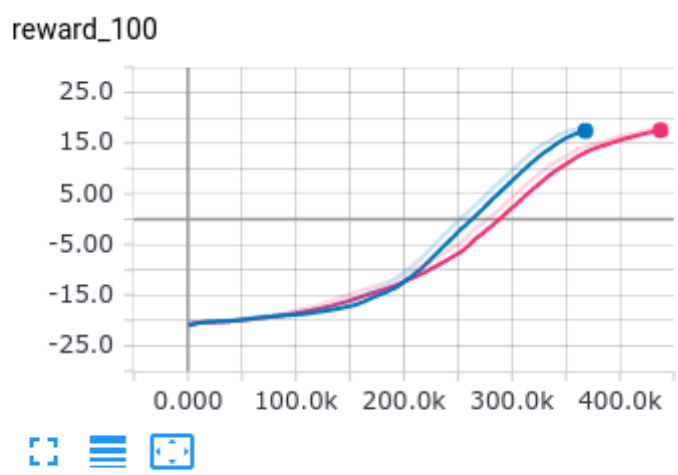
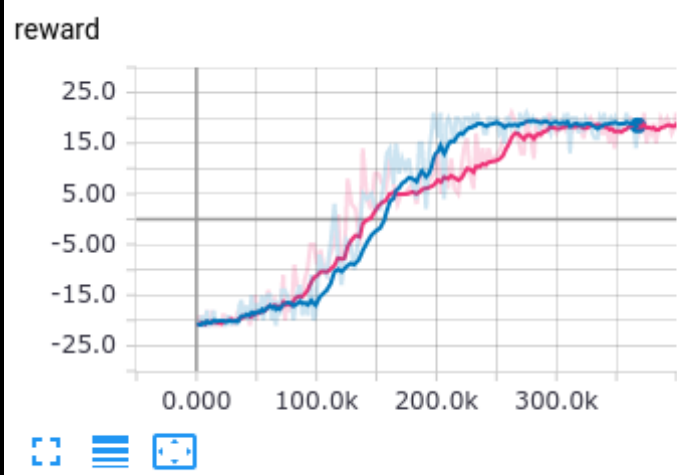
# Chapter 6: Deep Q-Networks





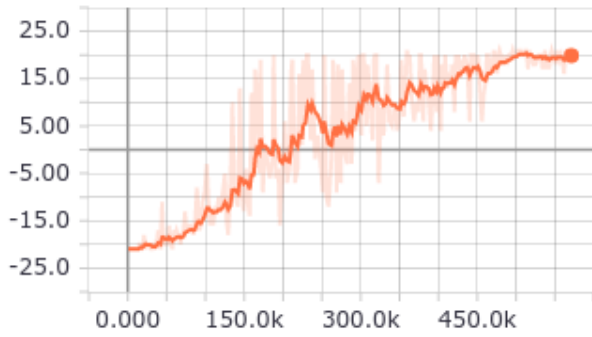
# Chapter 7: DQN Extensions



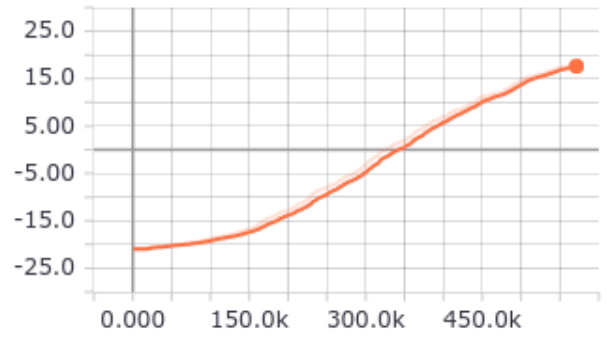




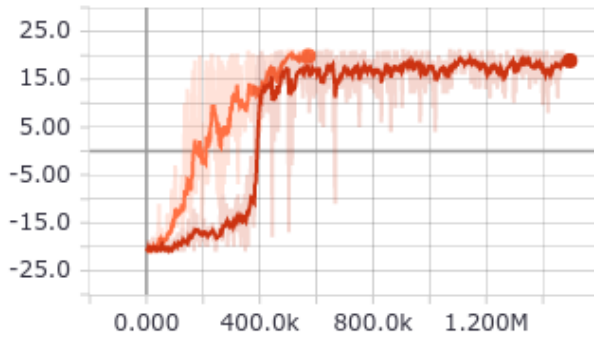
reward



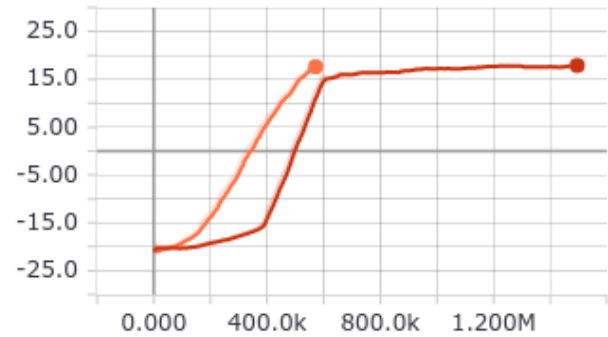
reward\_100



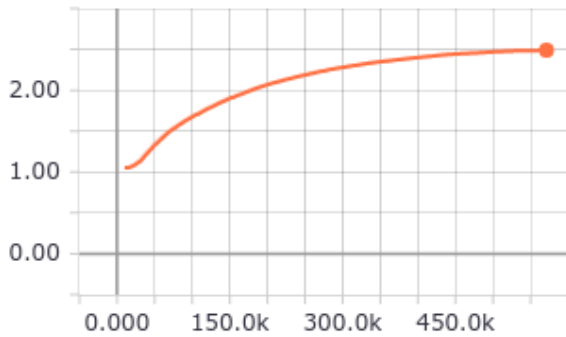
reward



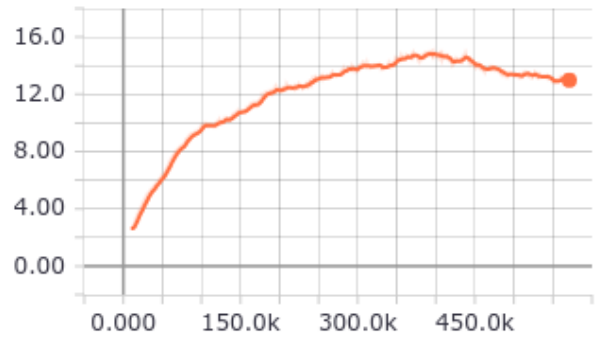
reward\_100



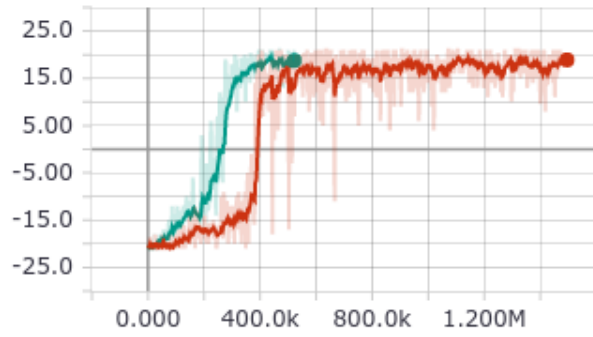
sigma\_snr\_layer\_1



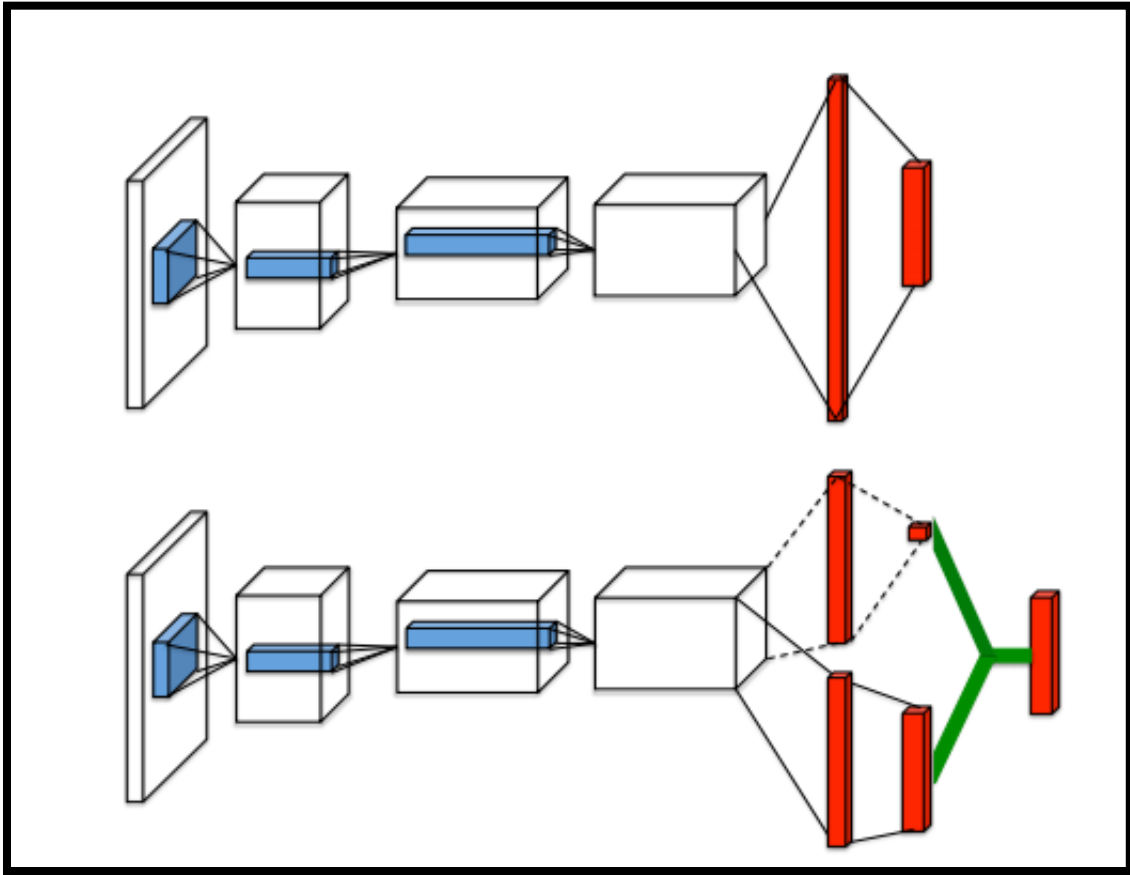
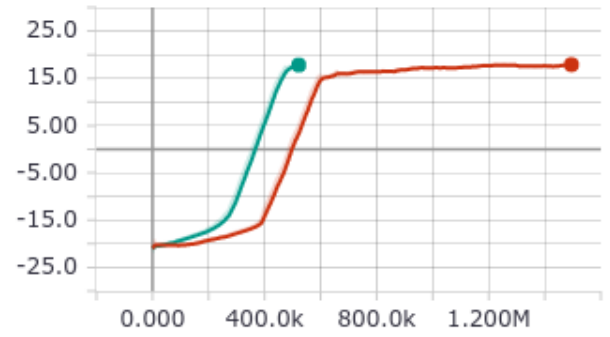
sigma\_snr\_layer\_2



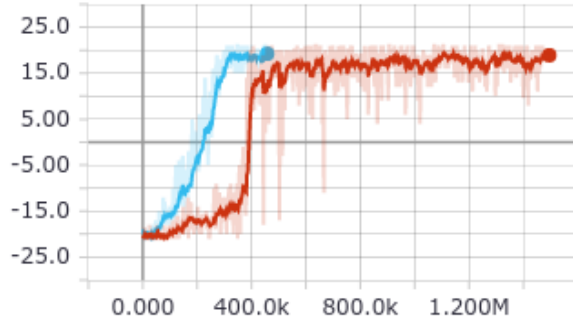
reward



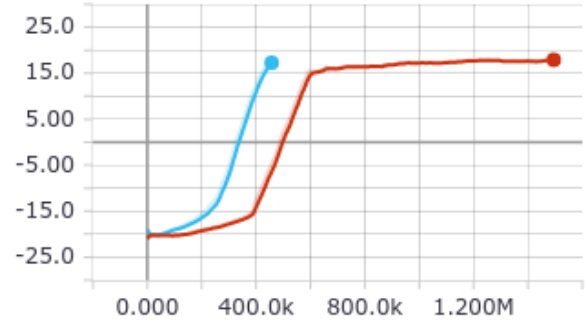
reward\_100



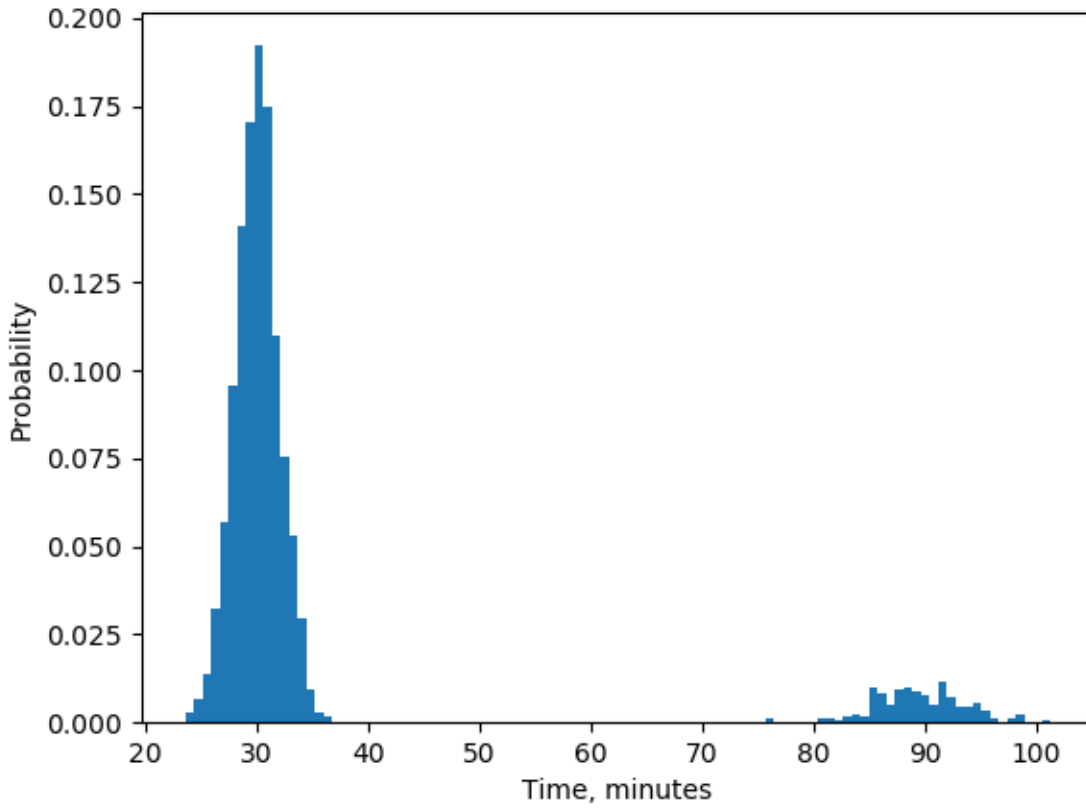
reward

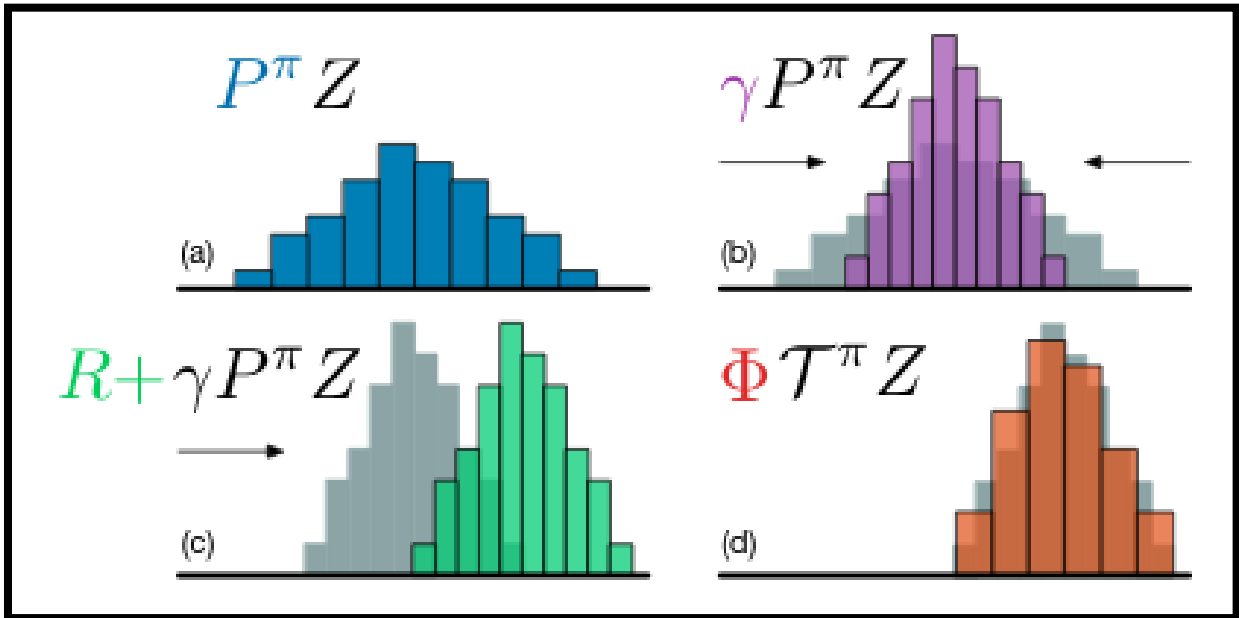
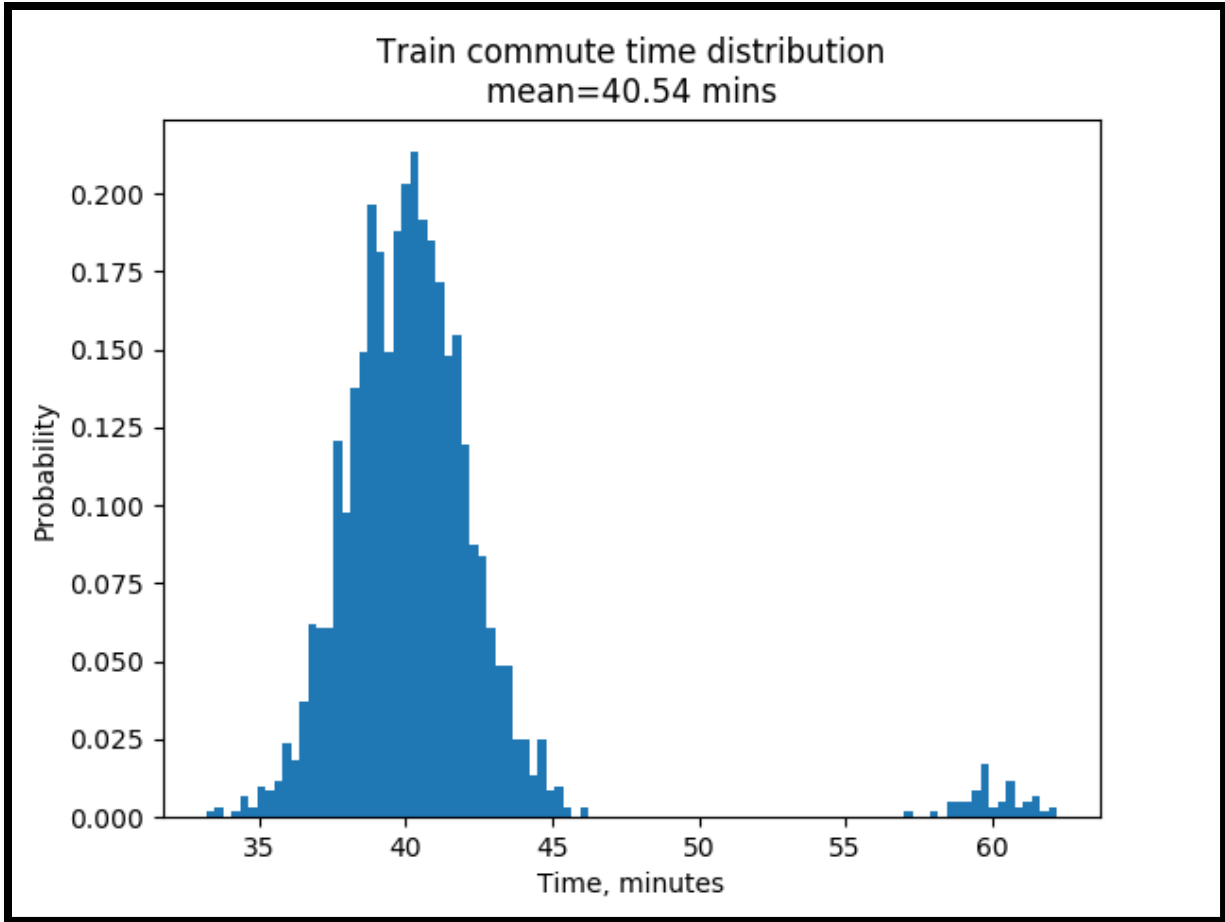


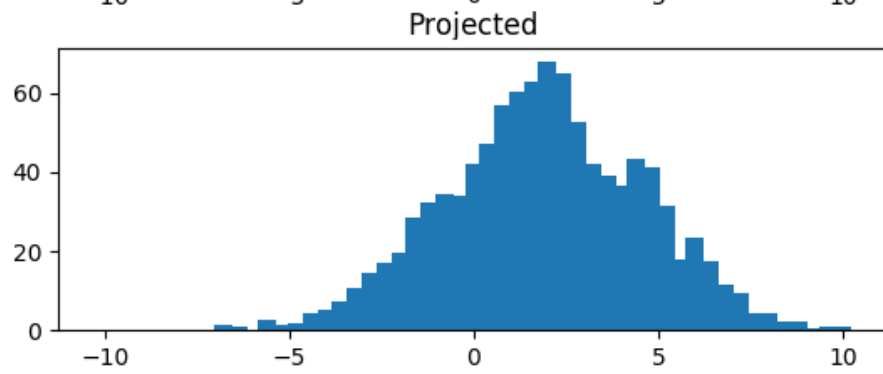
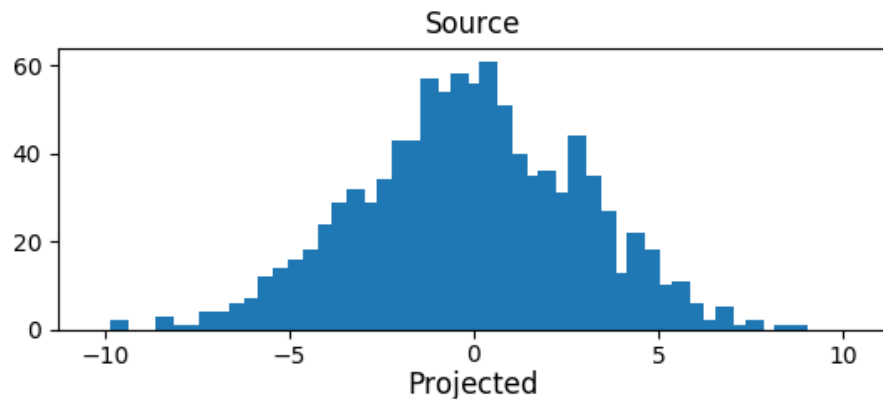
reward\_100

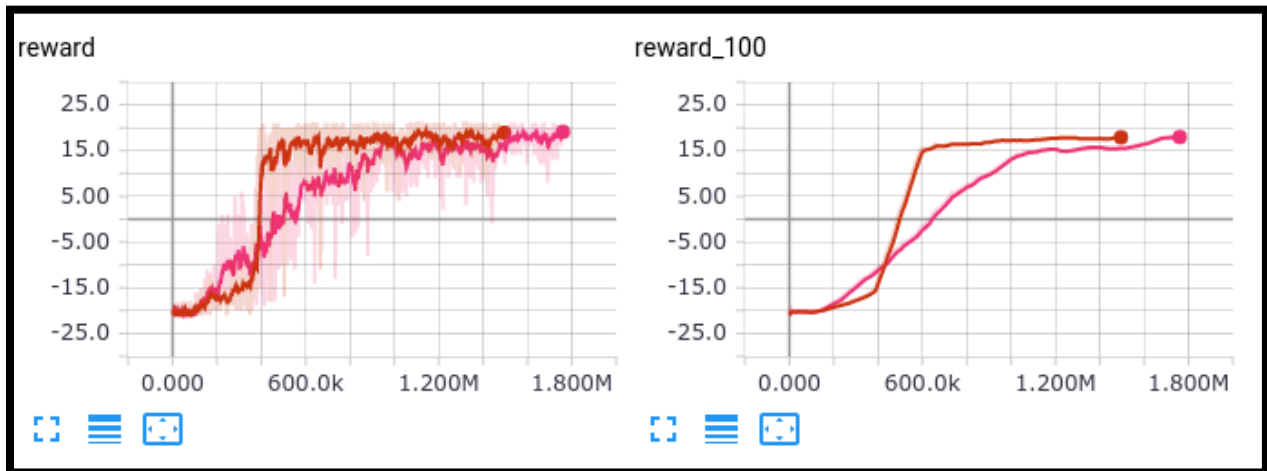
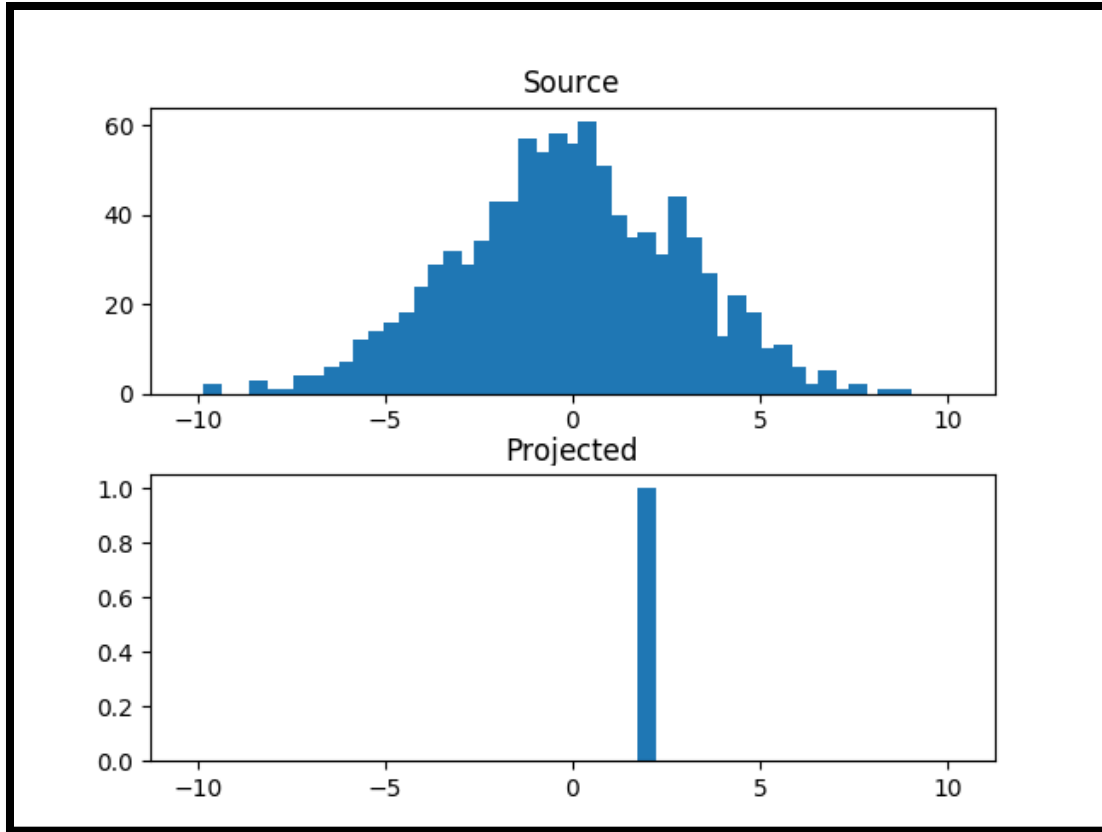


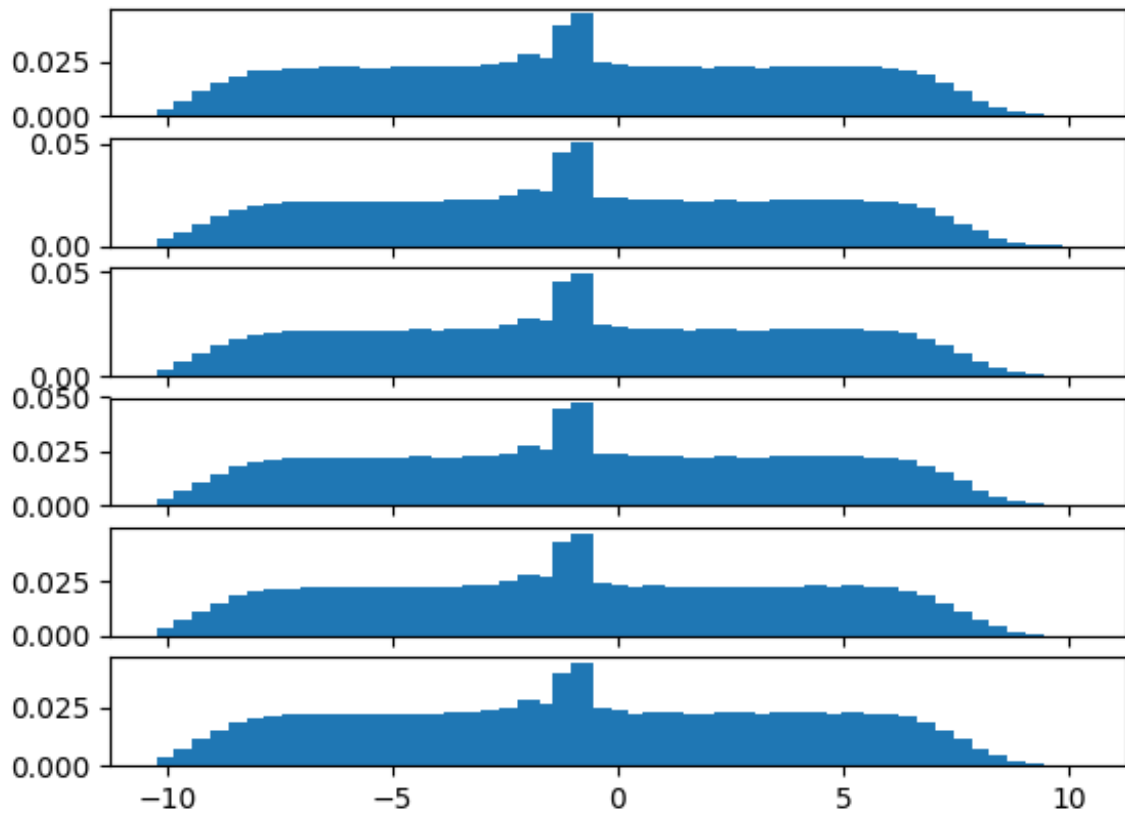
Car commute time distribution  
mean=35.43 mins

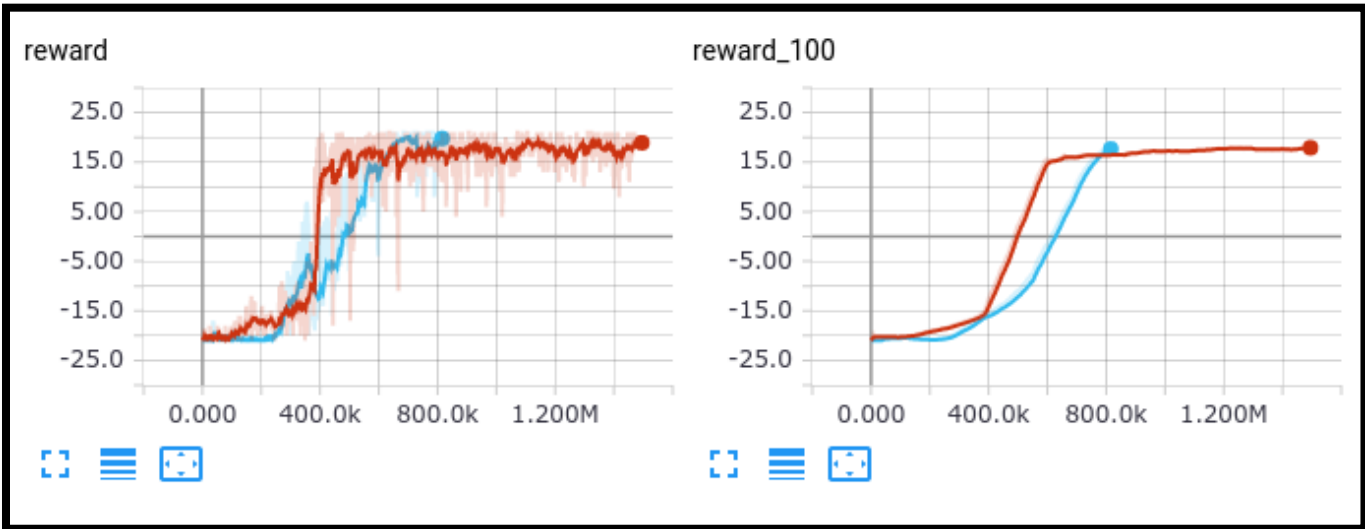
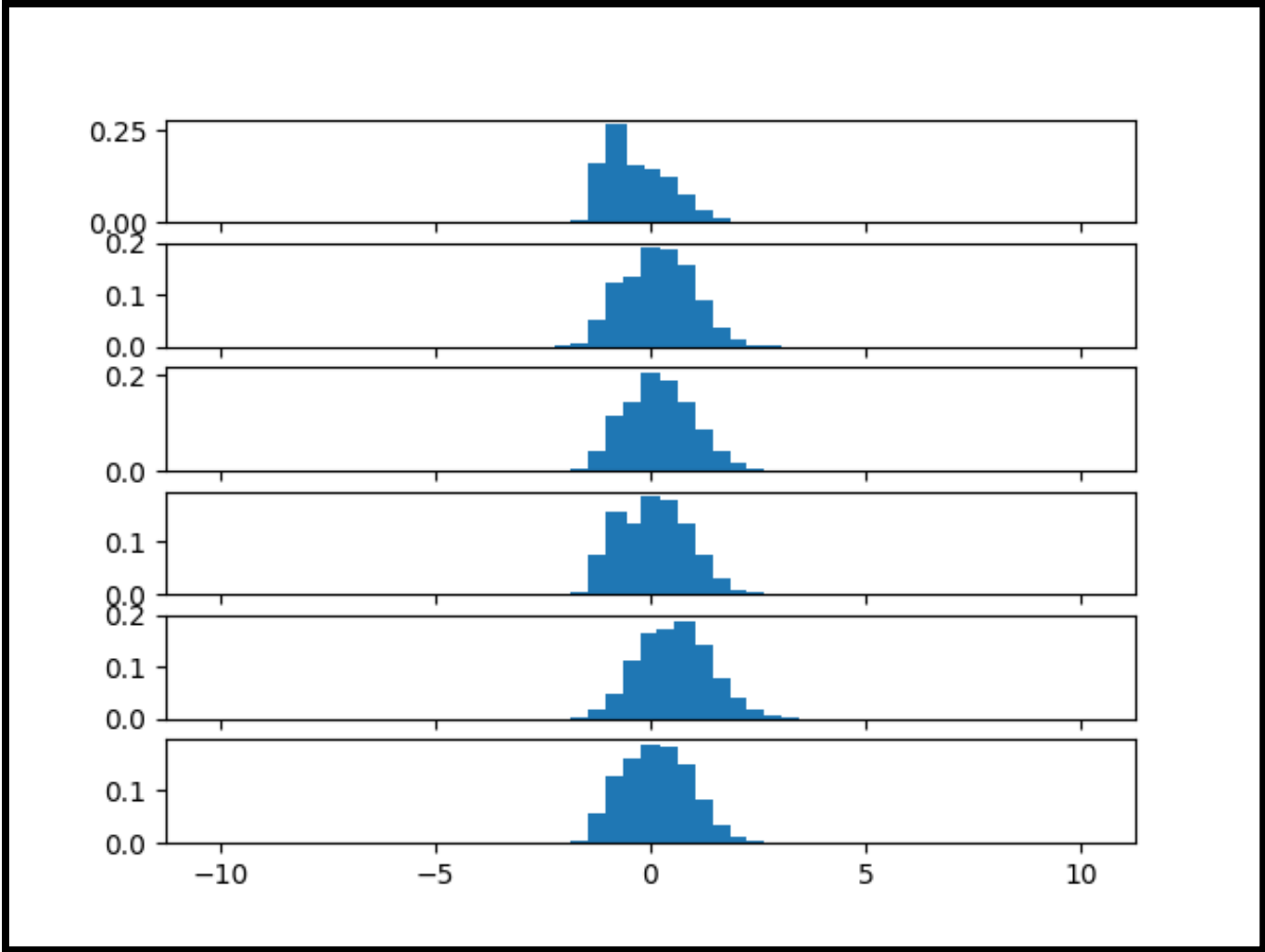










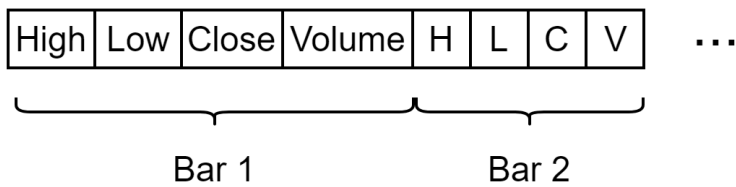




## Chapter 8: Stocks Trading Using RL

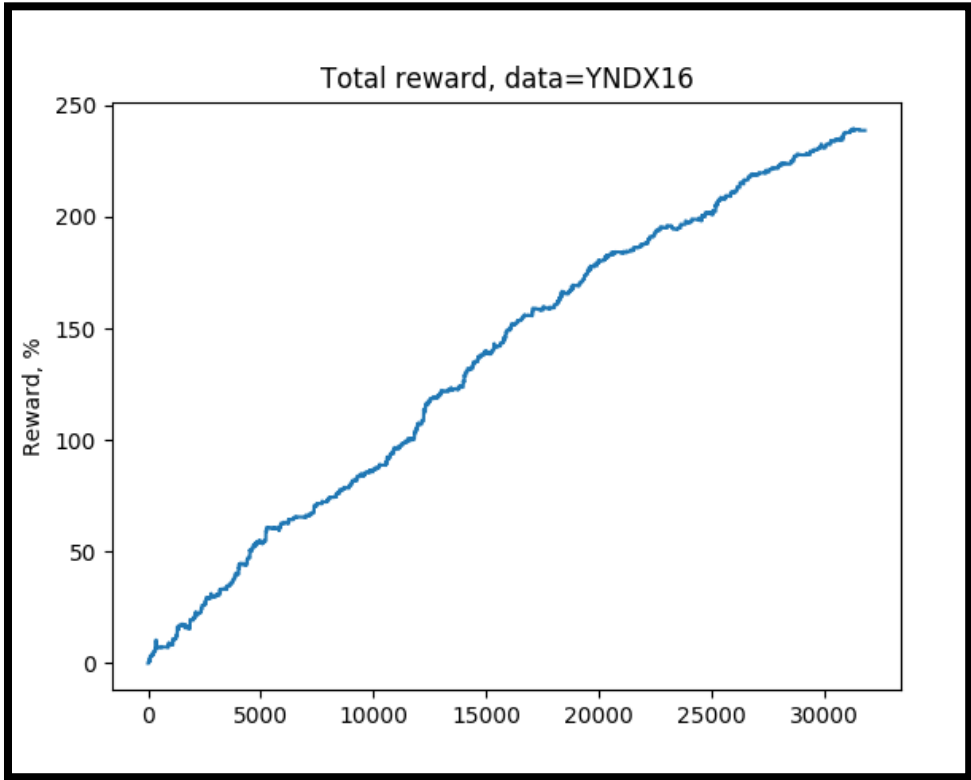
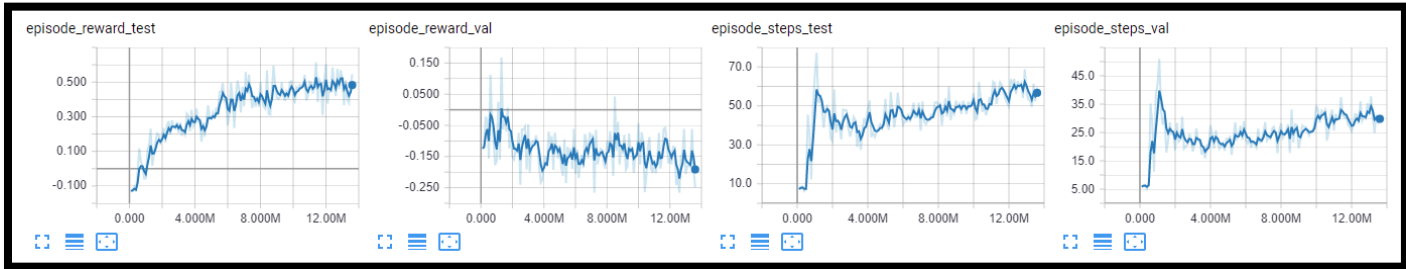
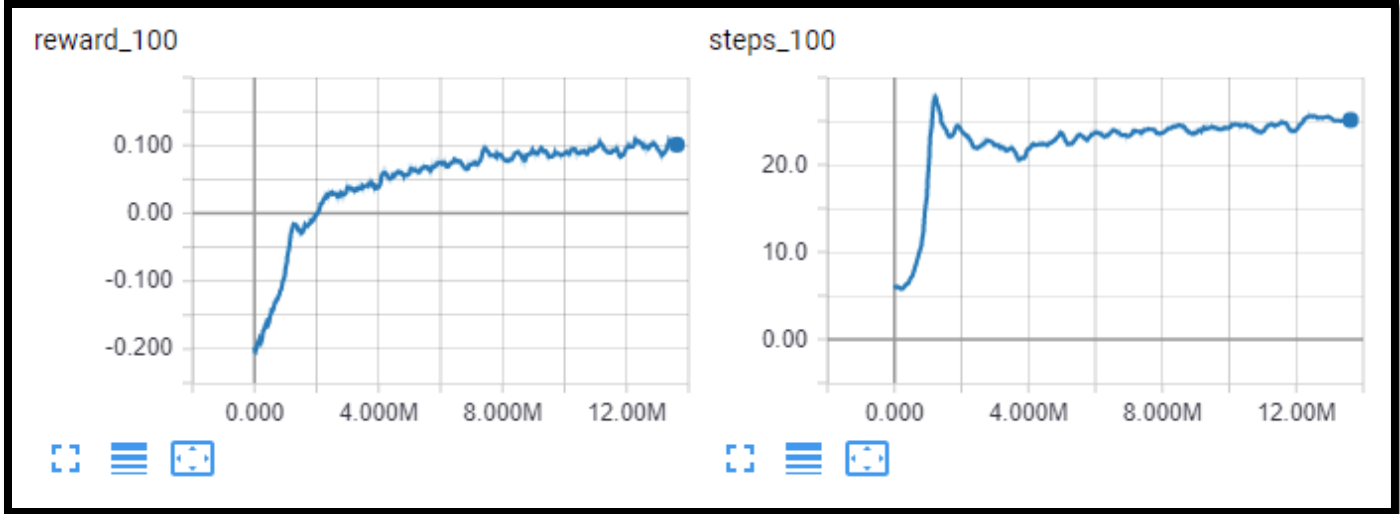


Data encoding: conv\_1d=False, vector

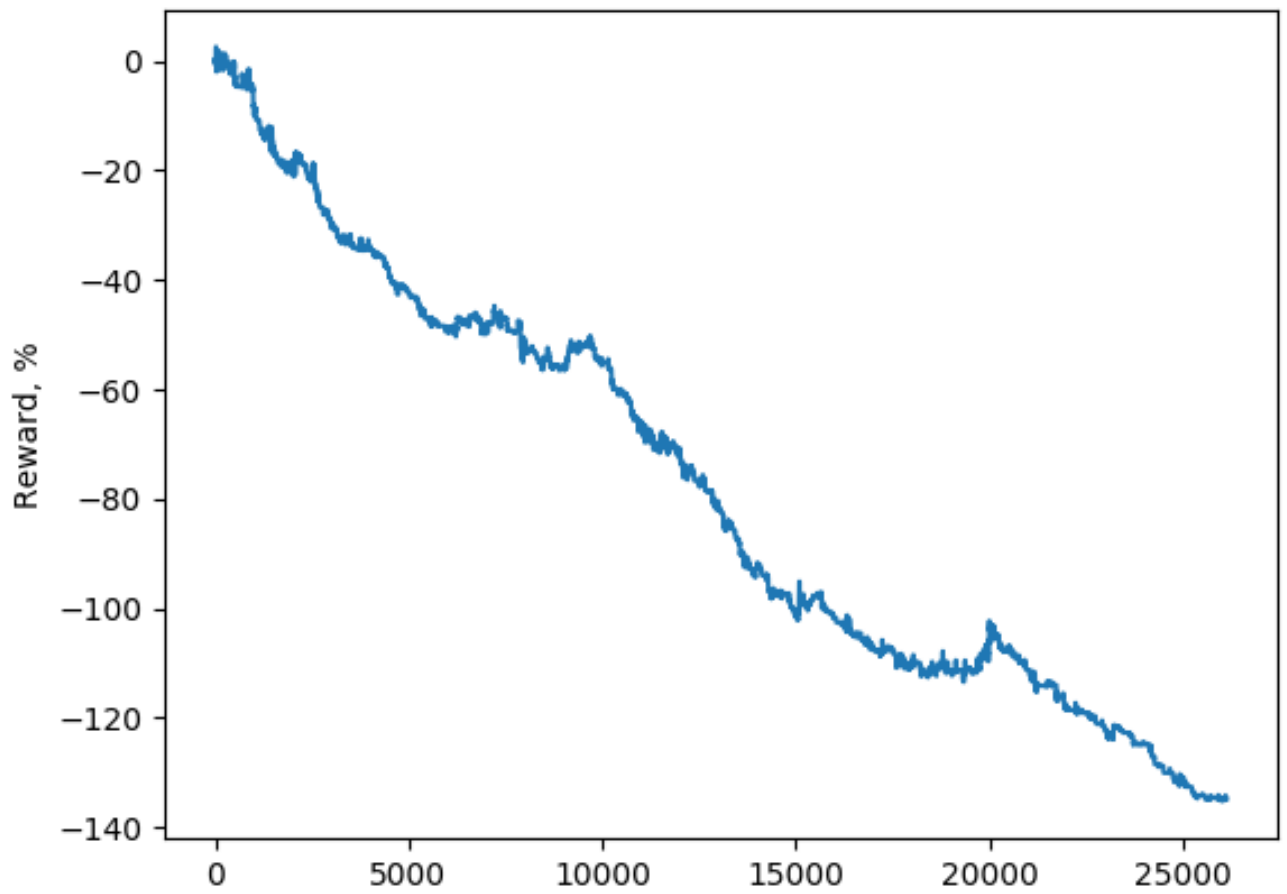


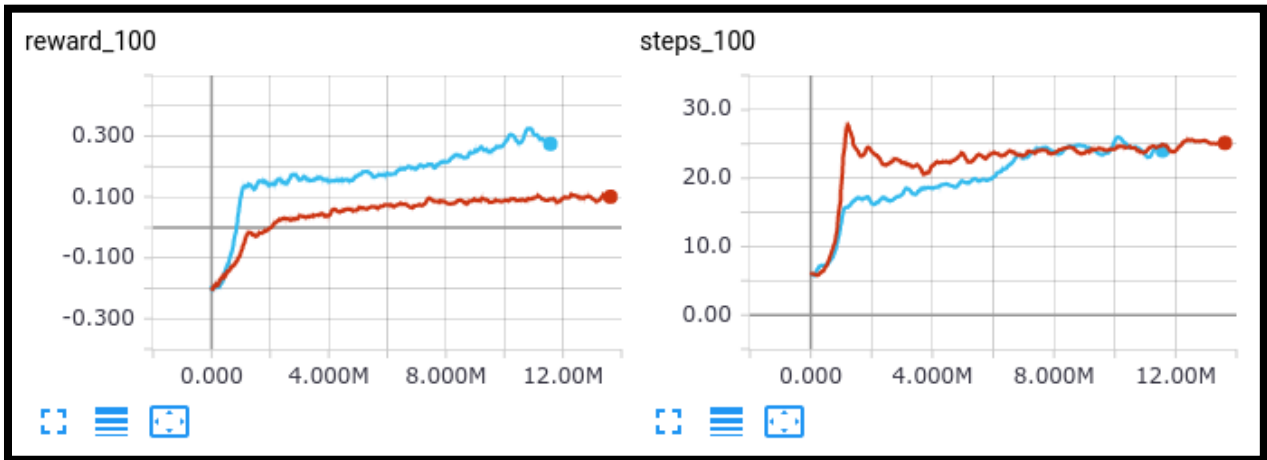
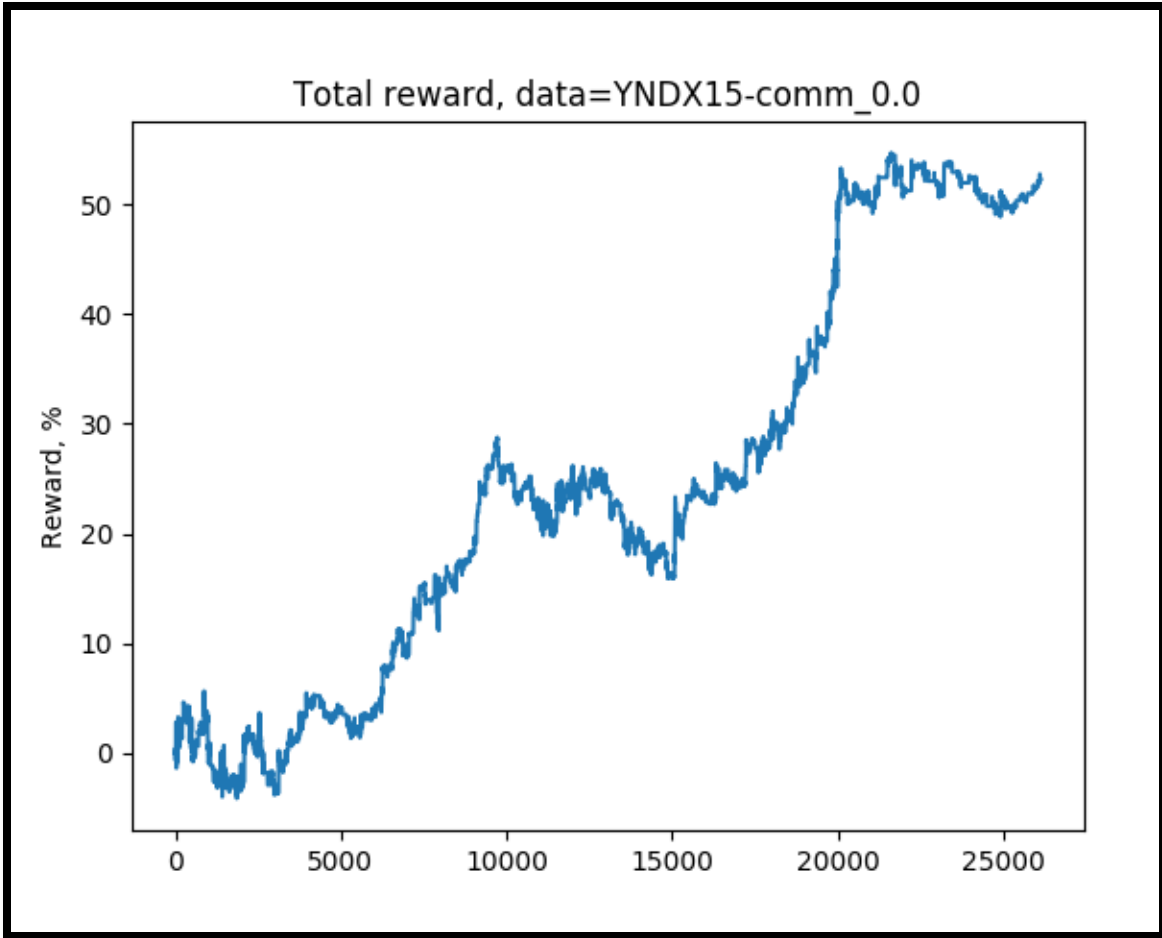
Data encoding: conv\_1d=True, matrix

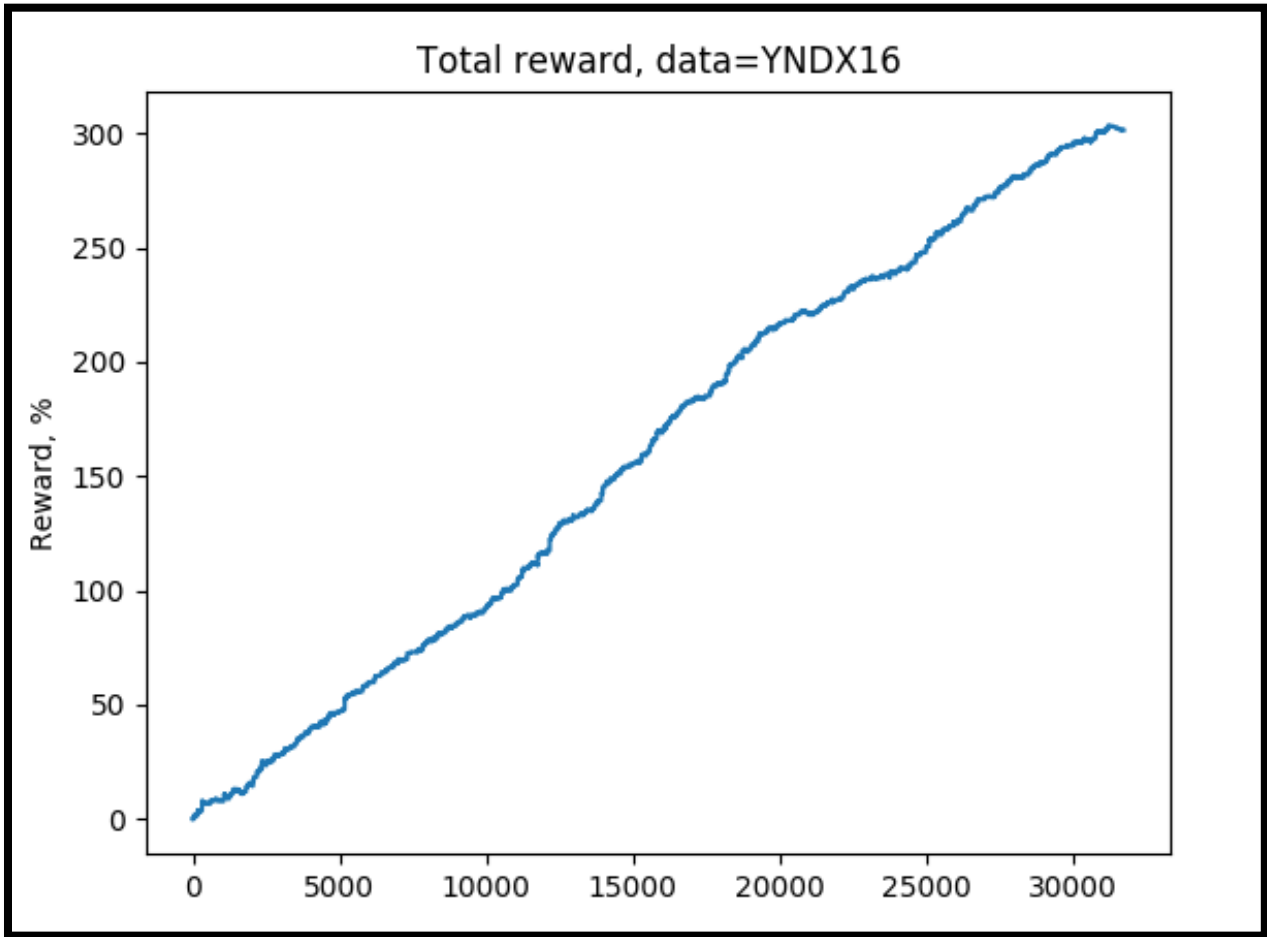
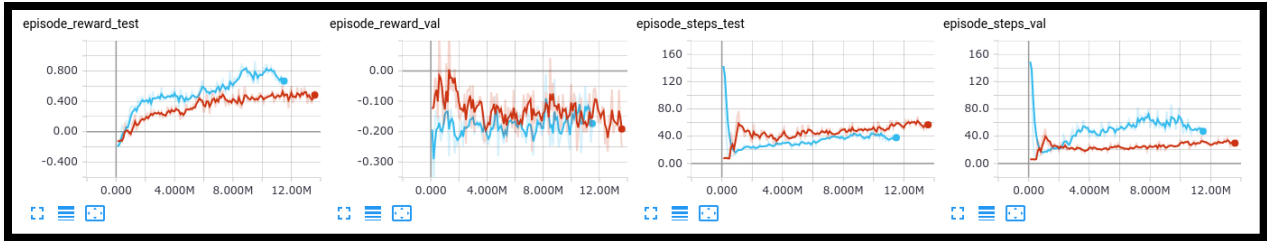
High prices for all bars
Low prices
Close prices
Volumes

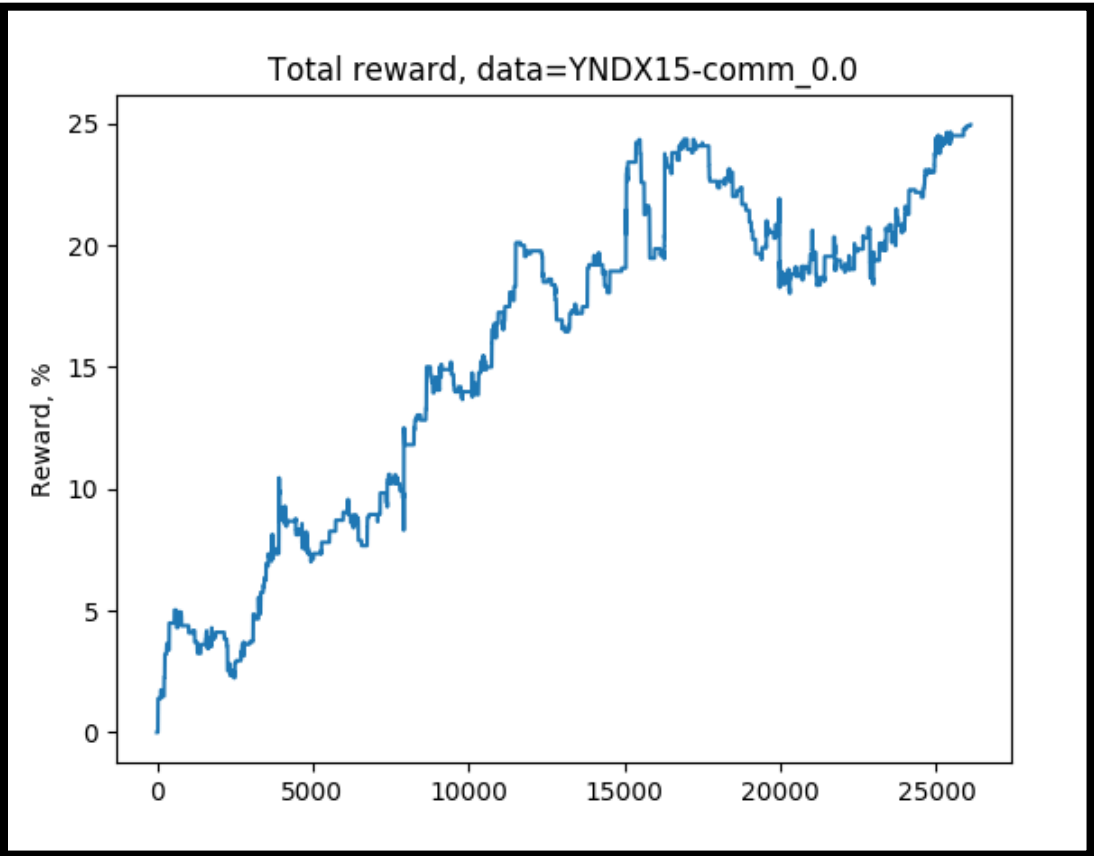
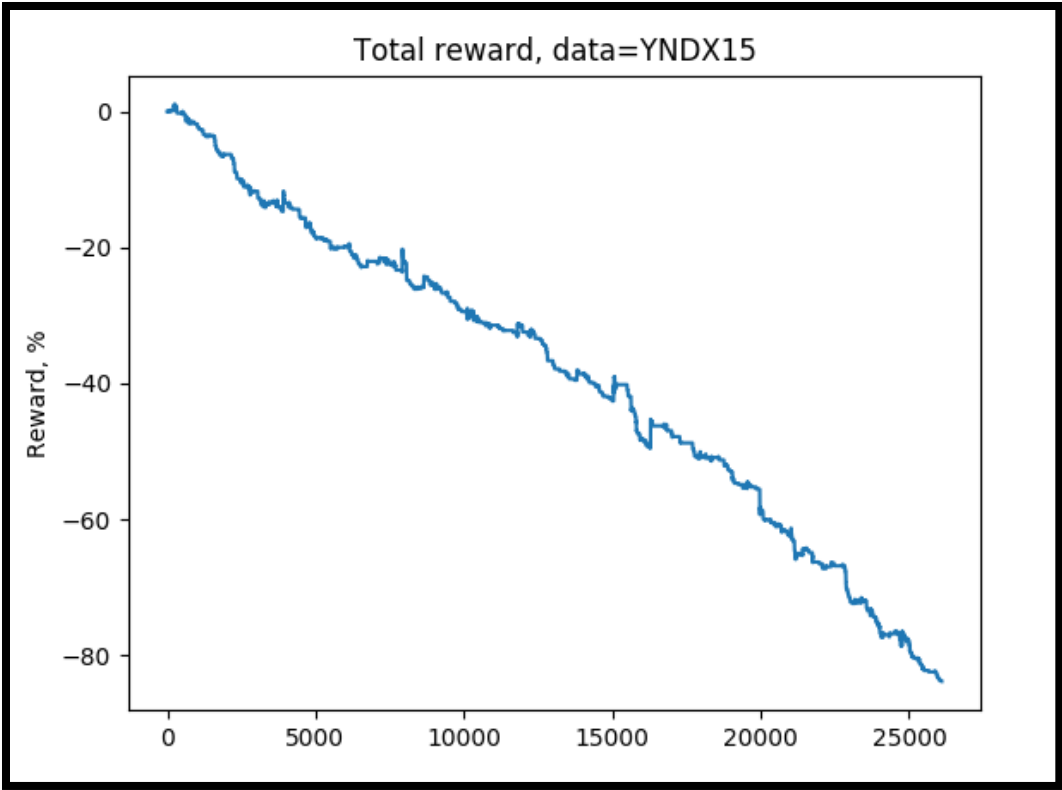


Total reward, data=YNDX15

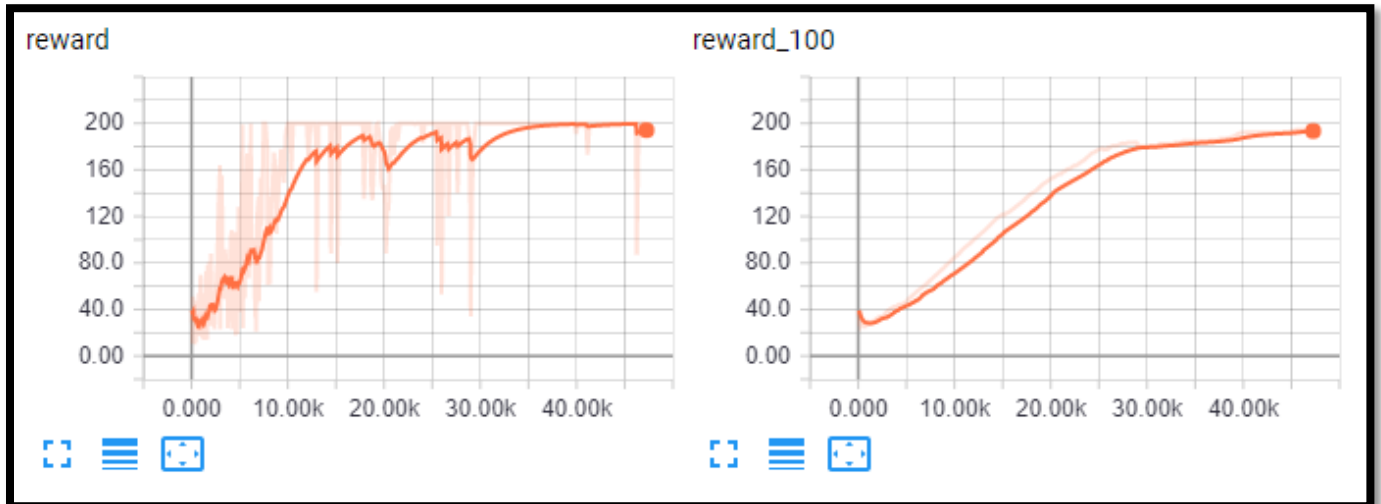
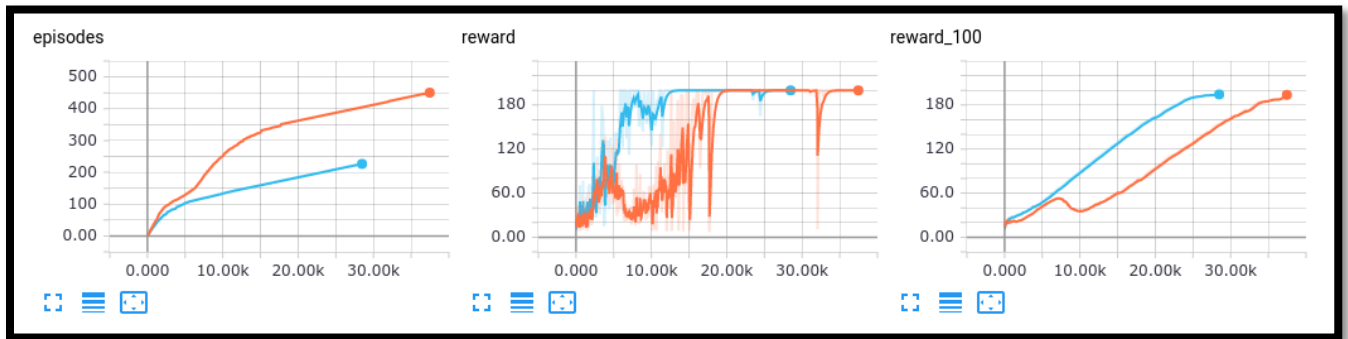
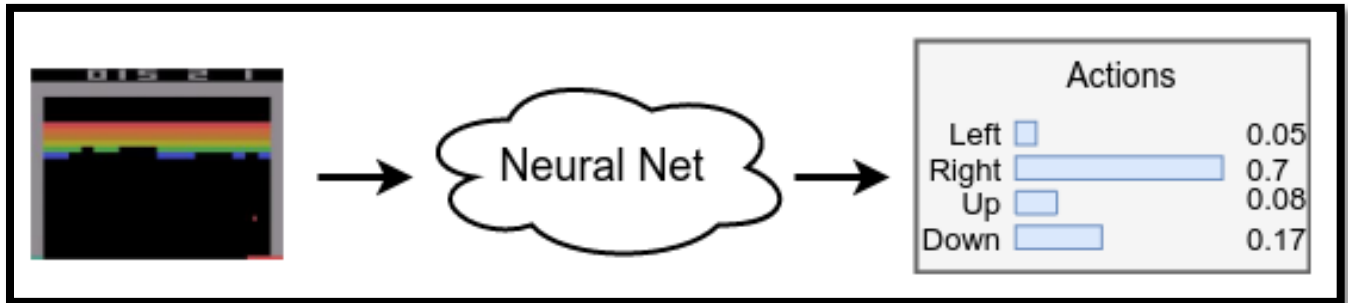


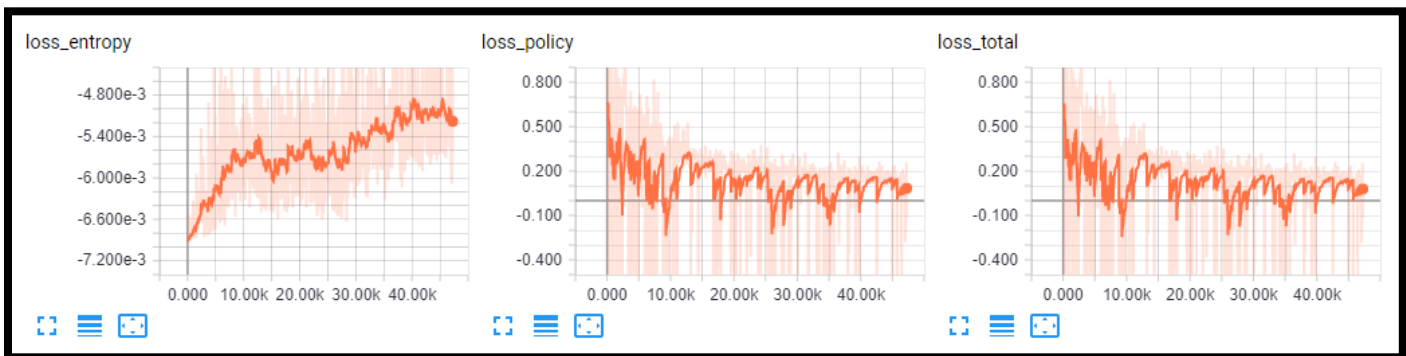
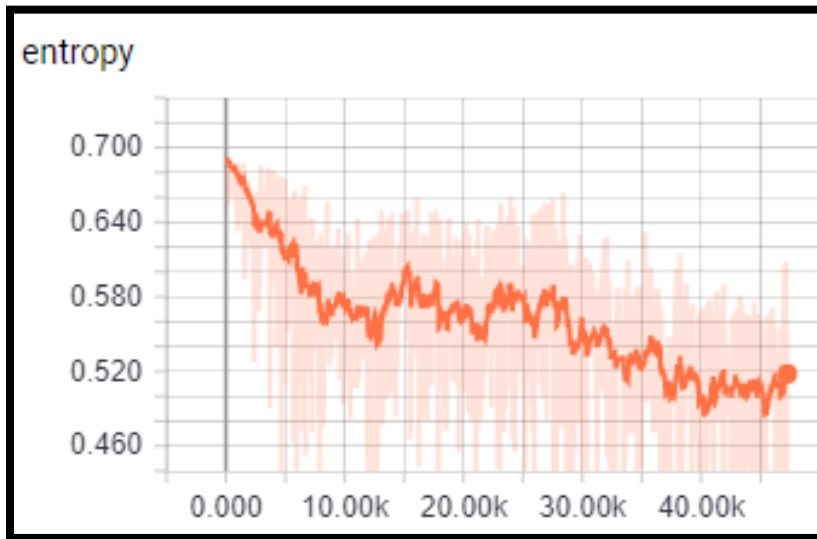
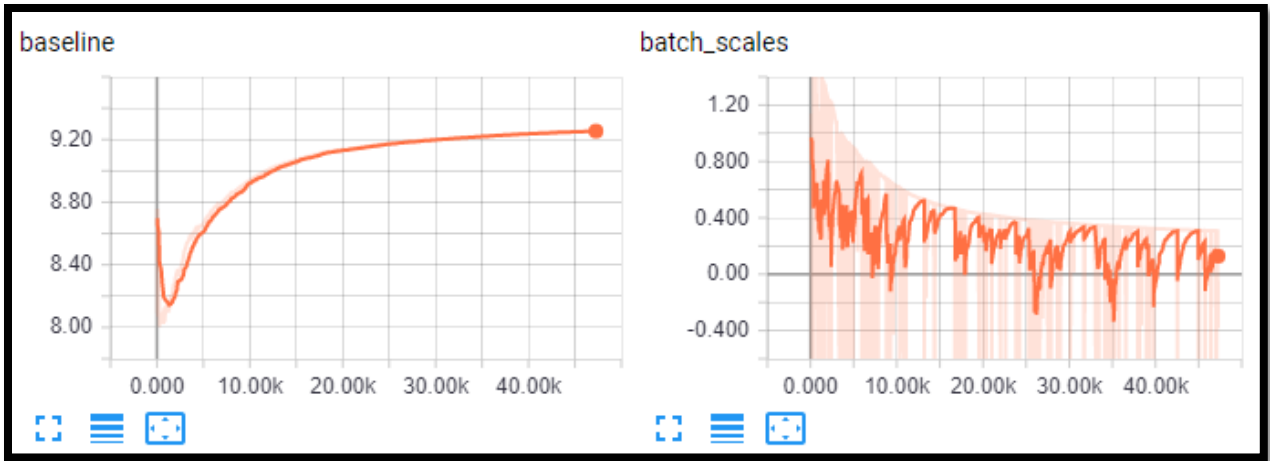




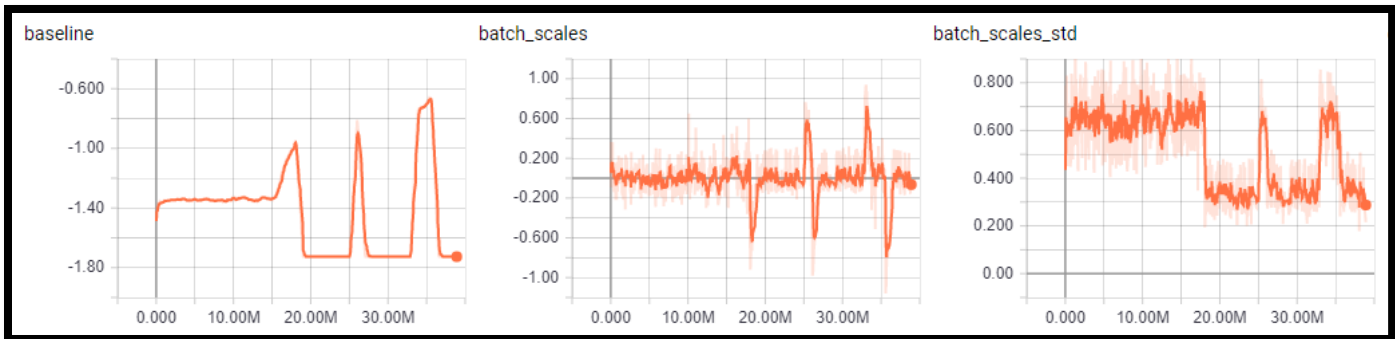
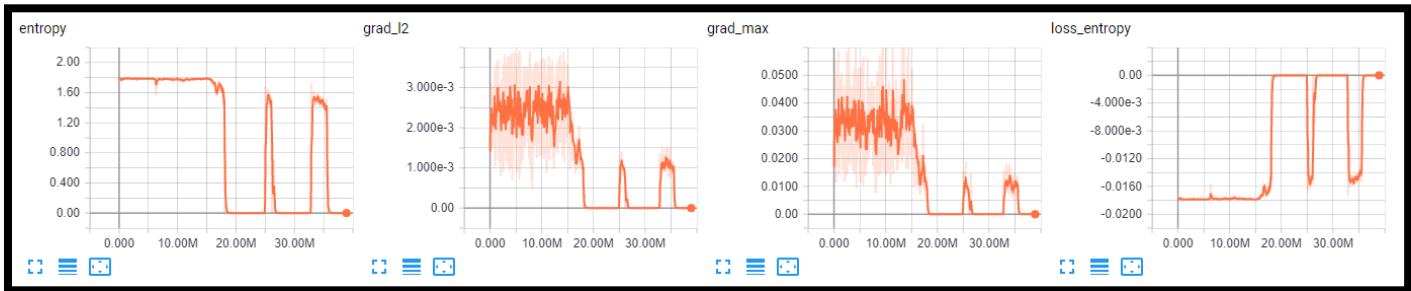
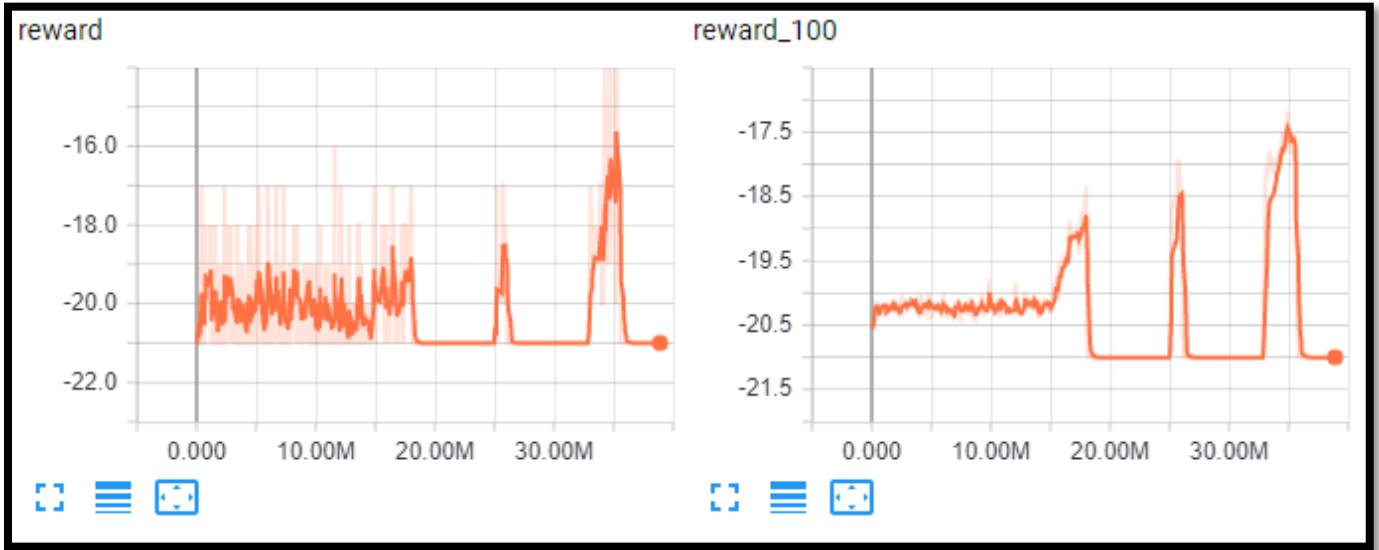
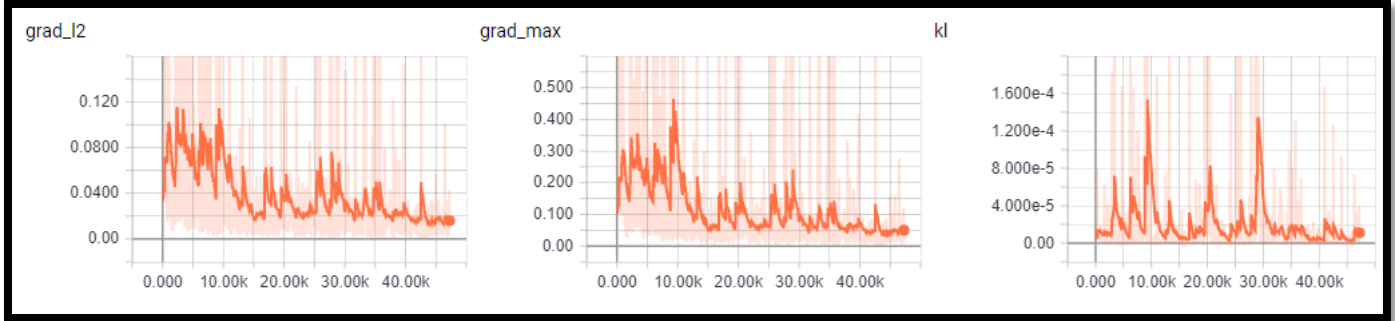


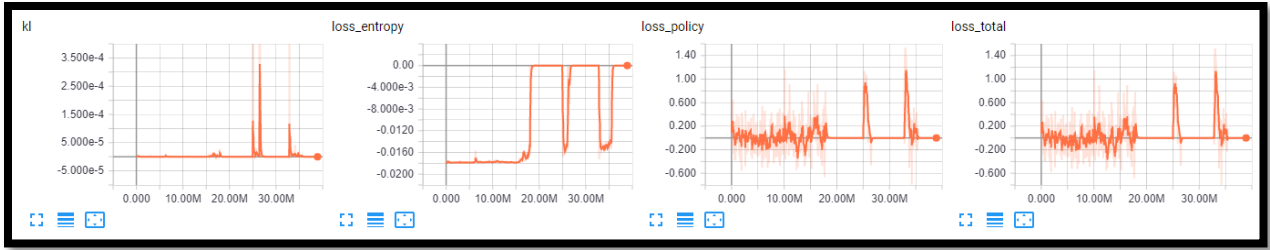
# Chapter 9: Policy Gradients – An Alternative



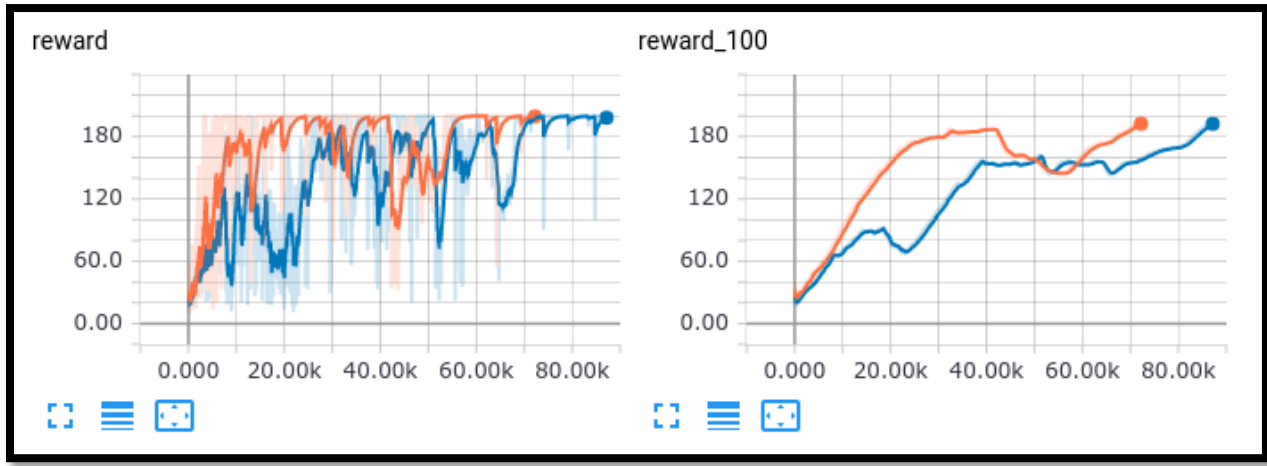
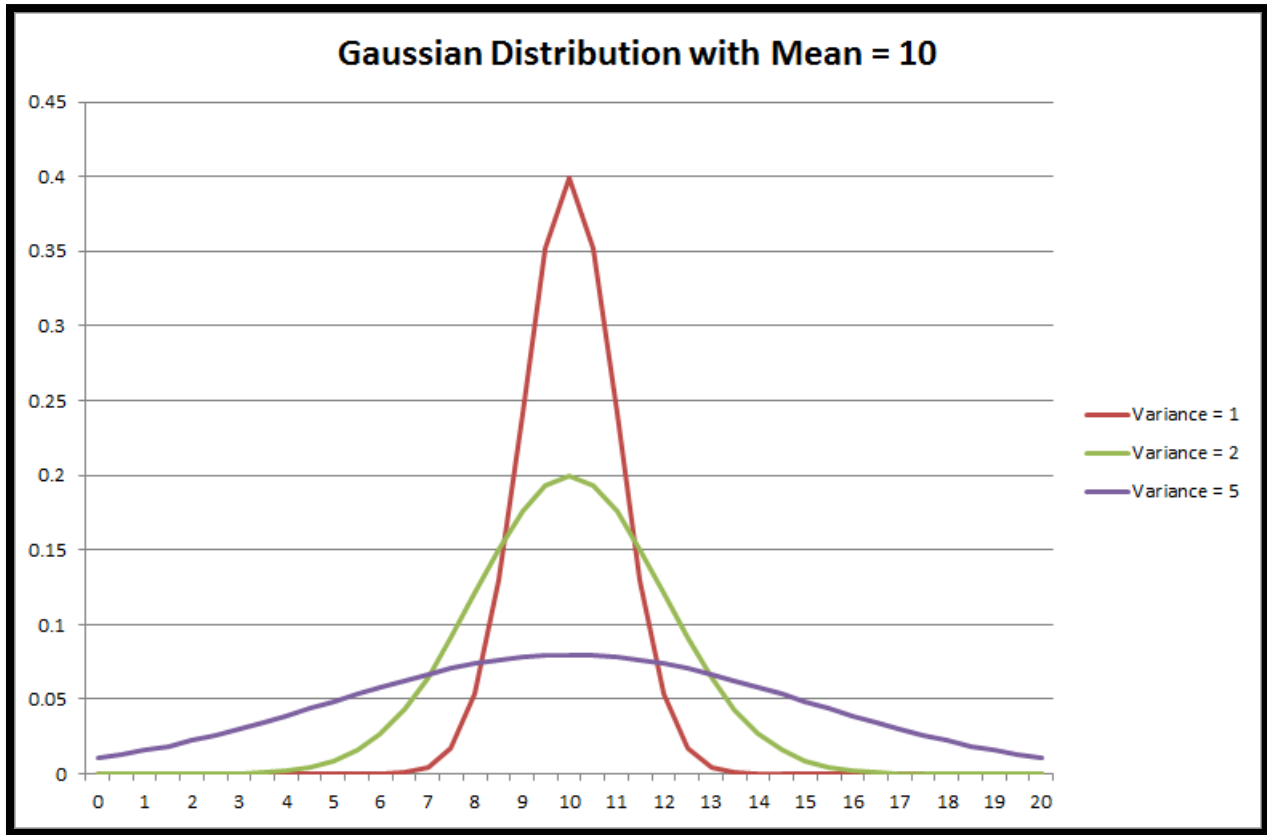


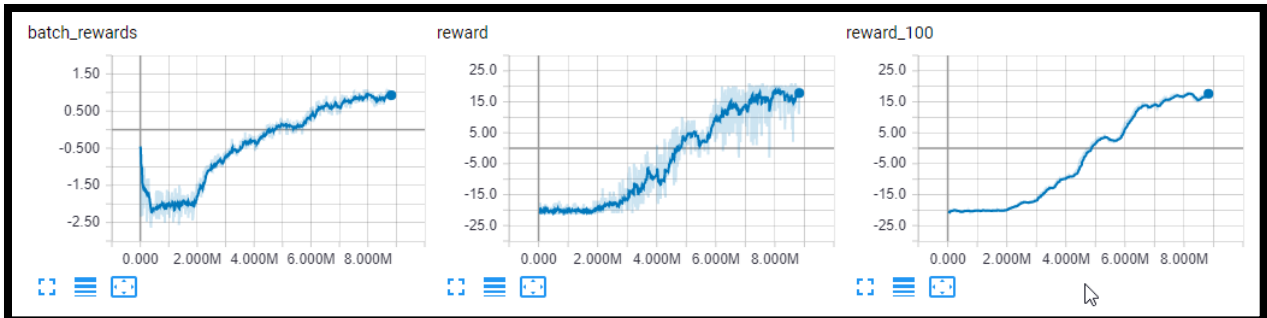
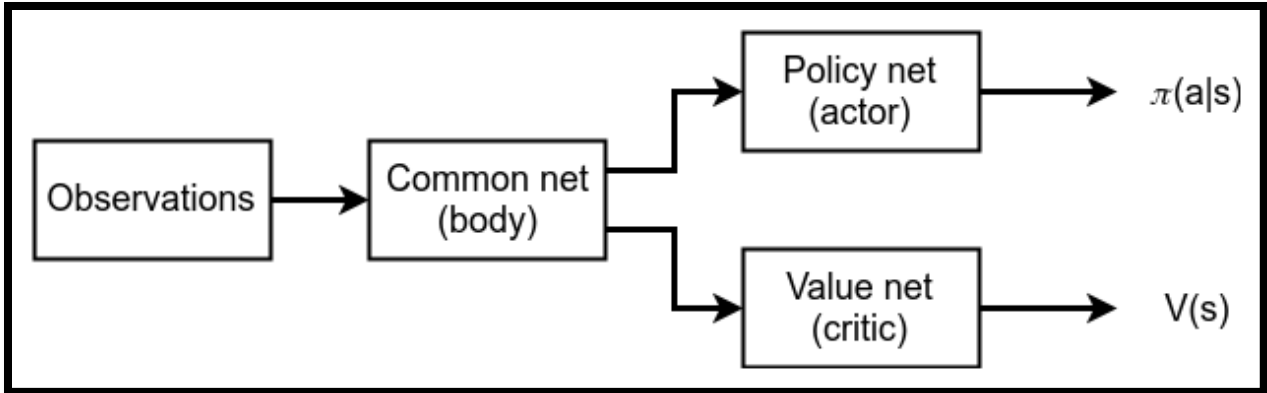
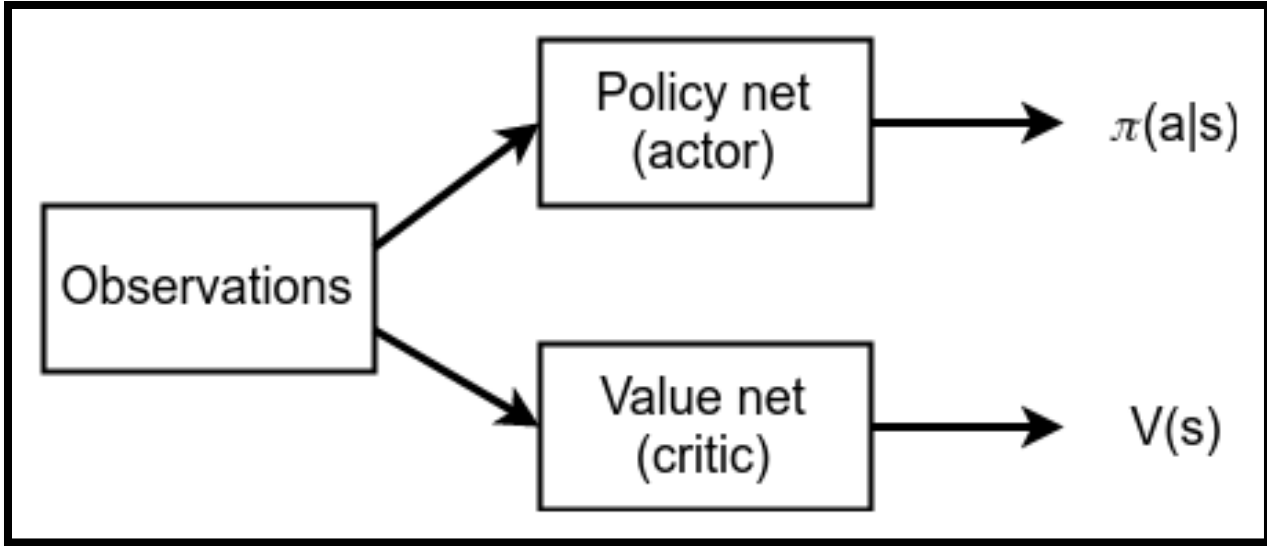
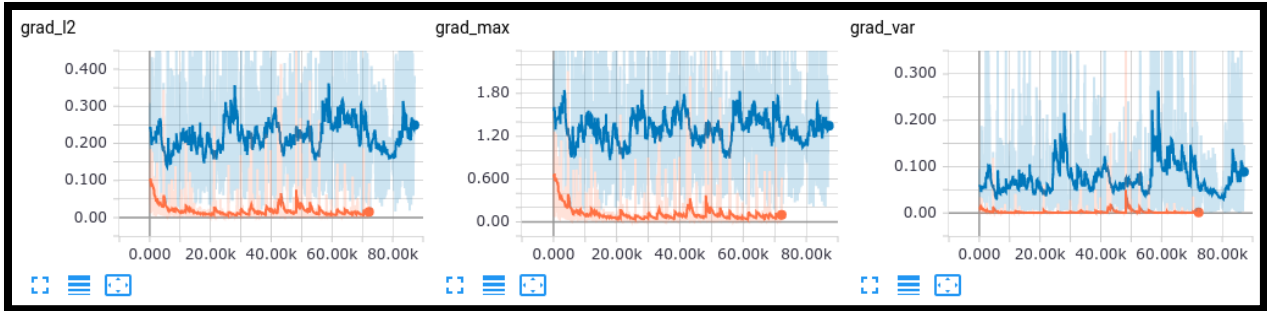


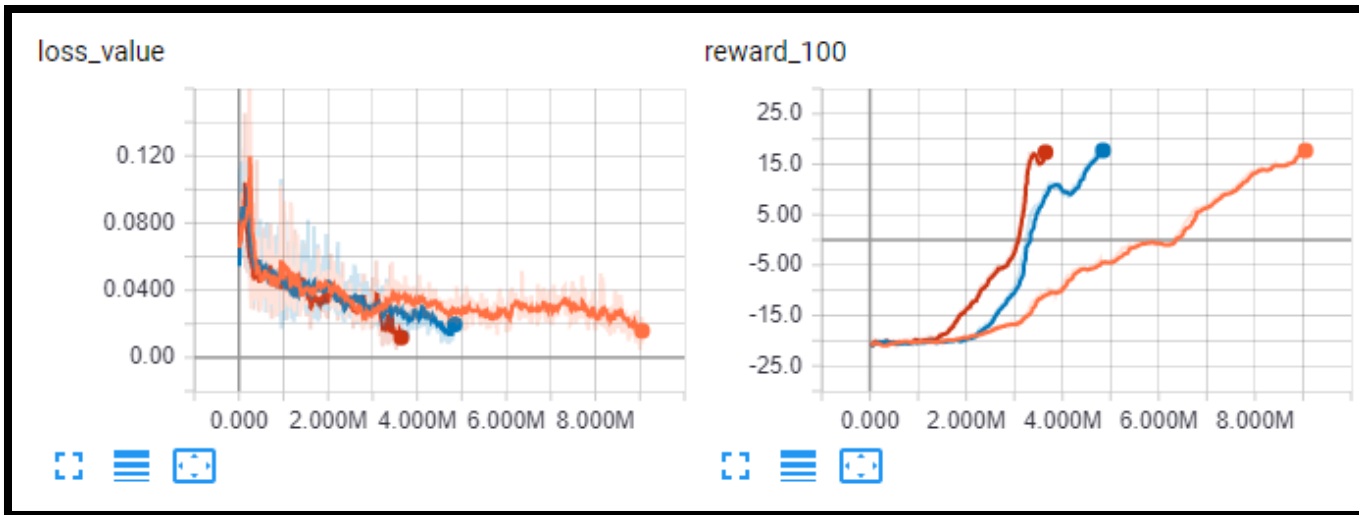
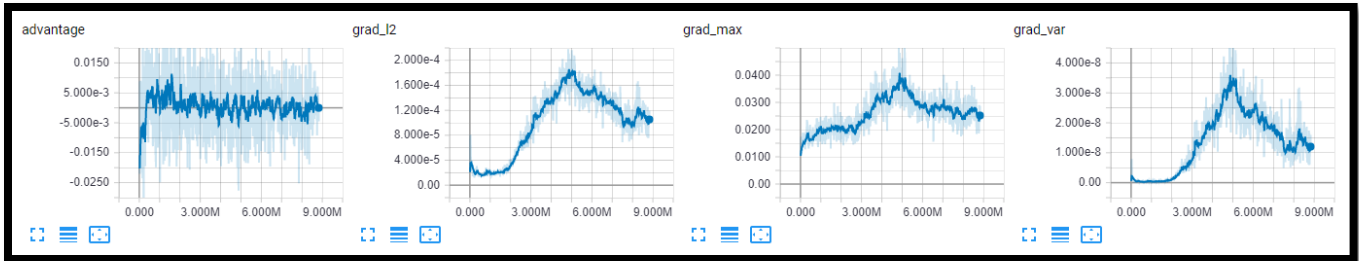
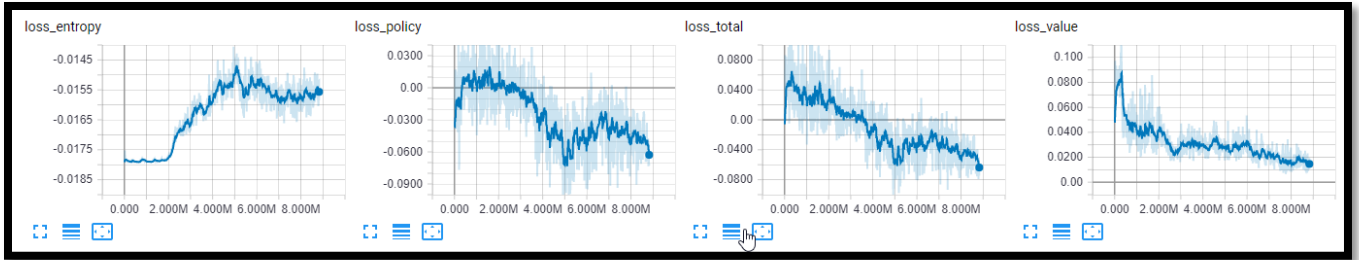




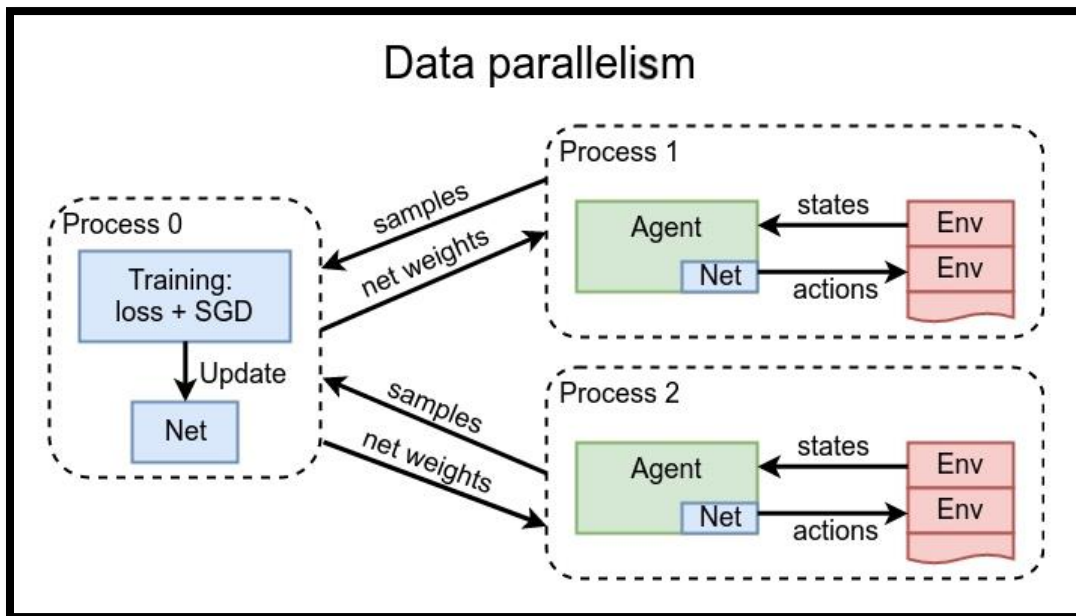
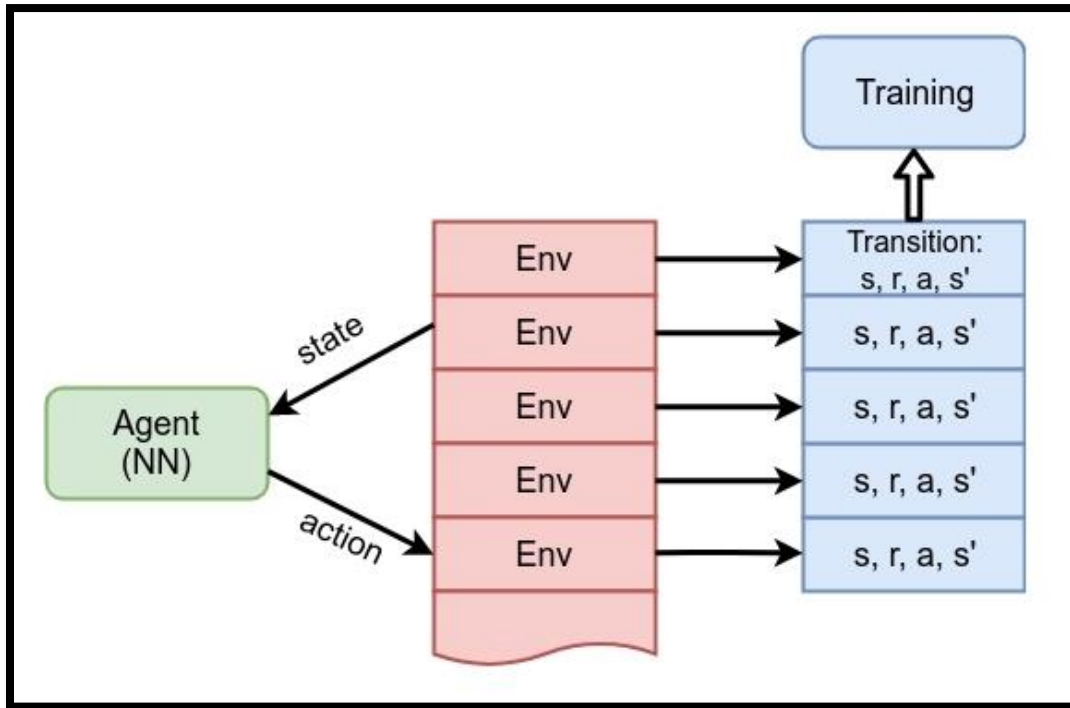
# Chapter 10: The Actor-Critic Method



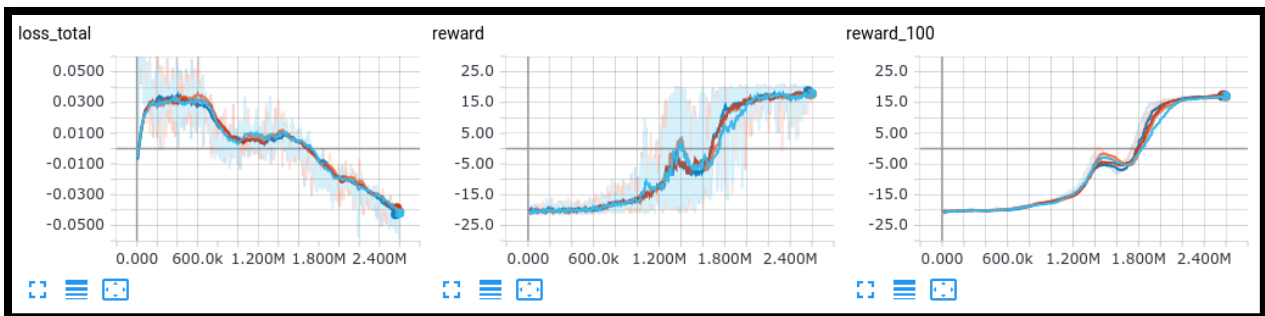
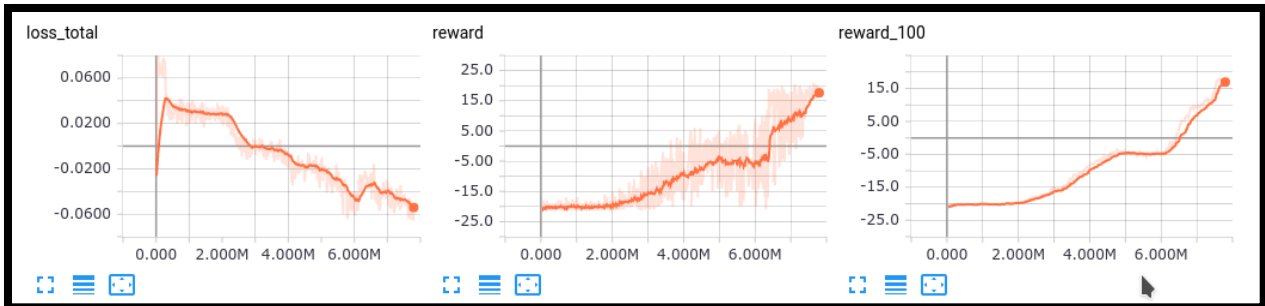
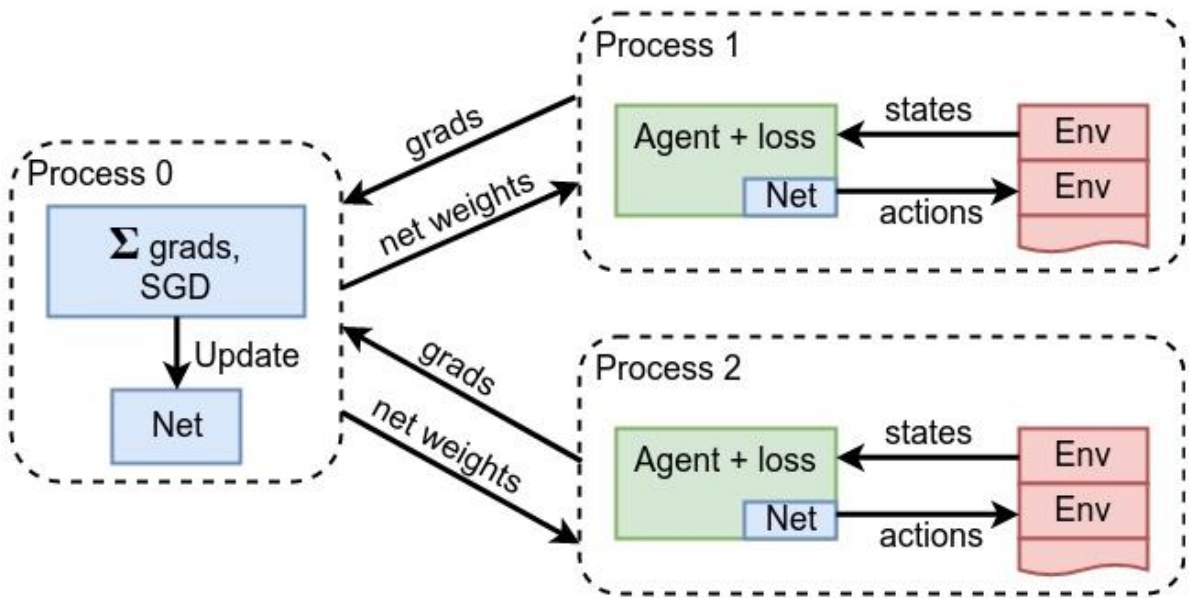




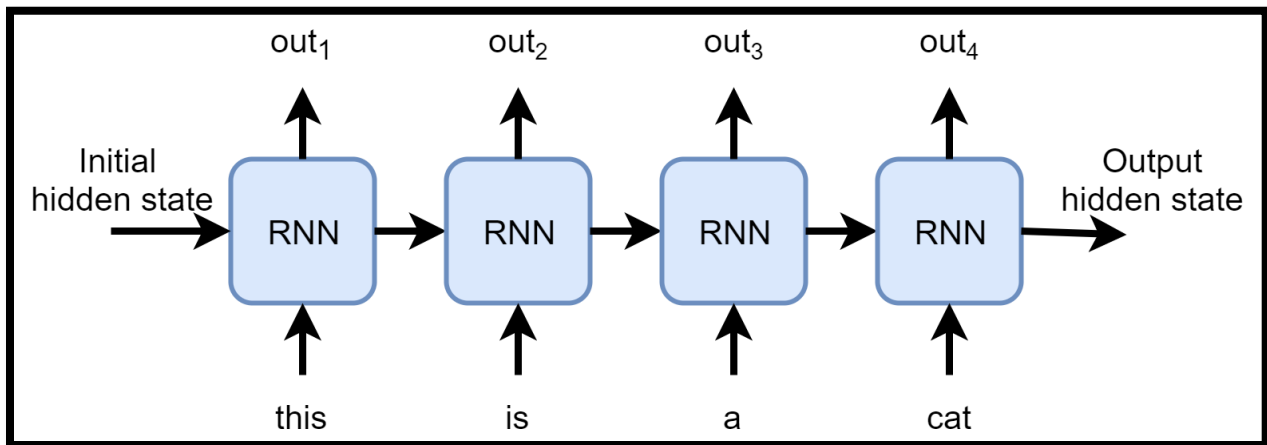
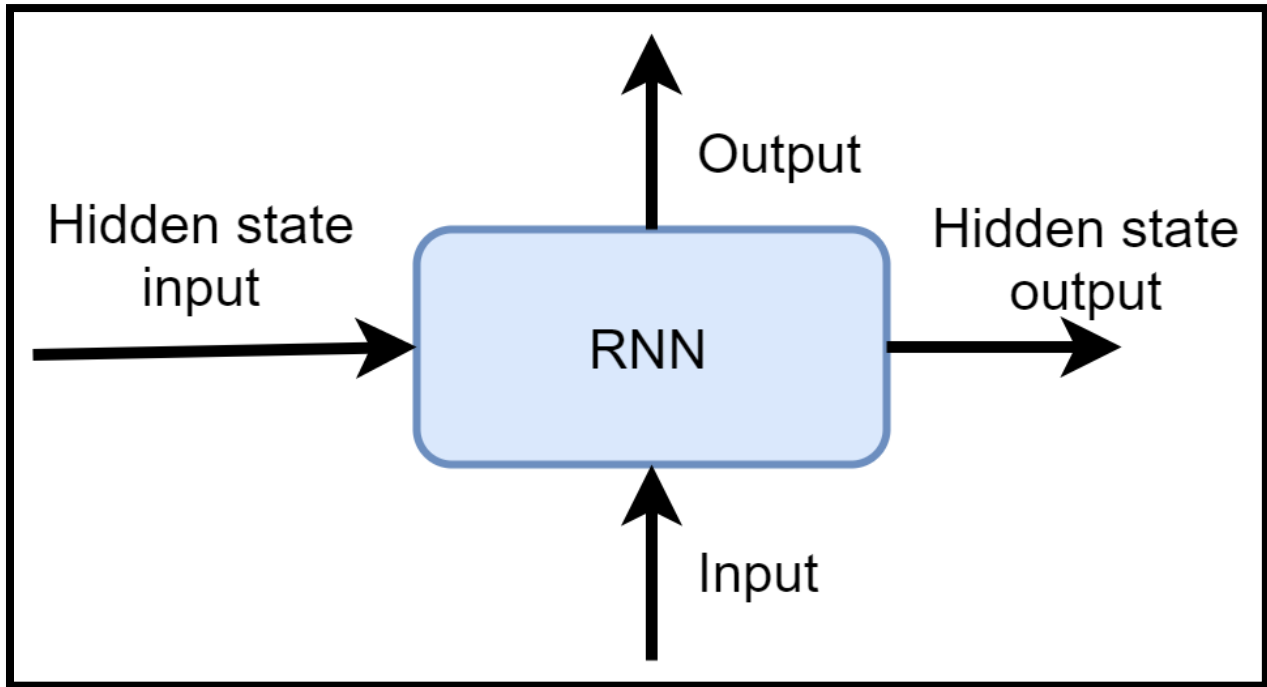
# Chapter 11: Asynchronous Advantage Actor-Critic



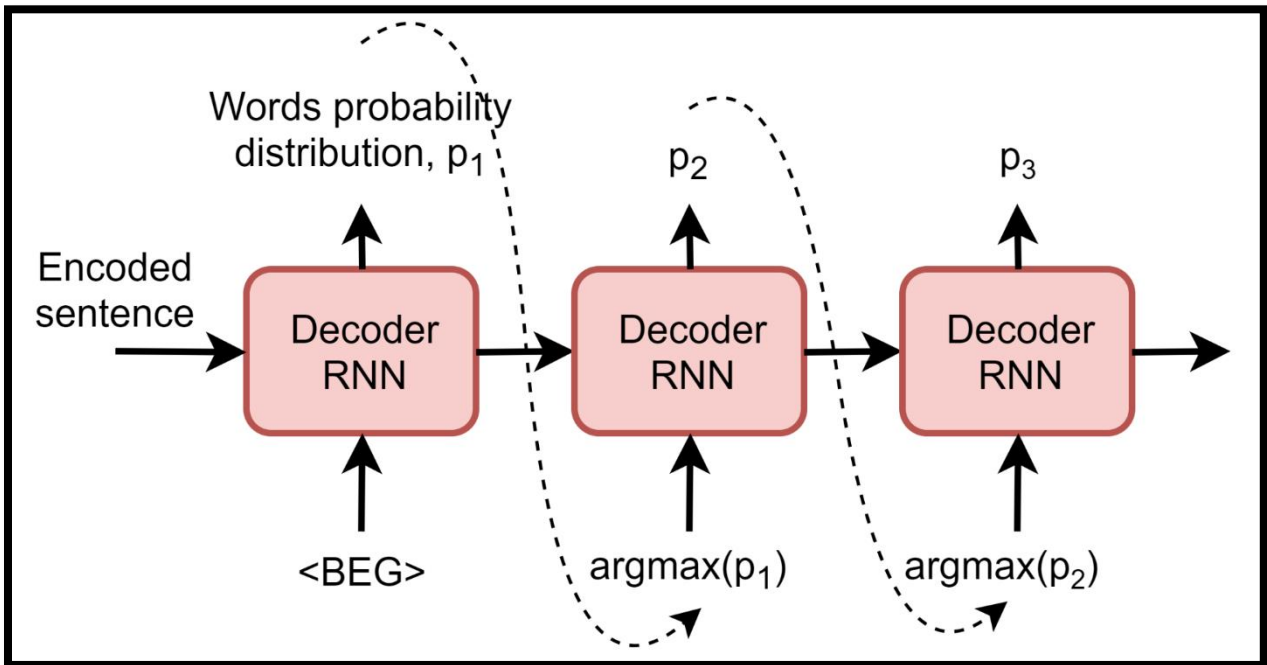
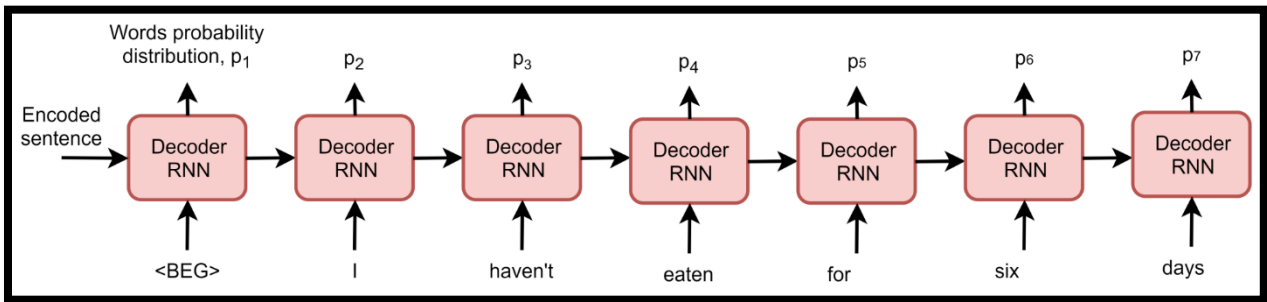
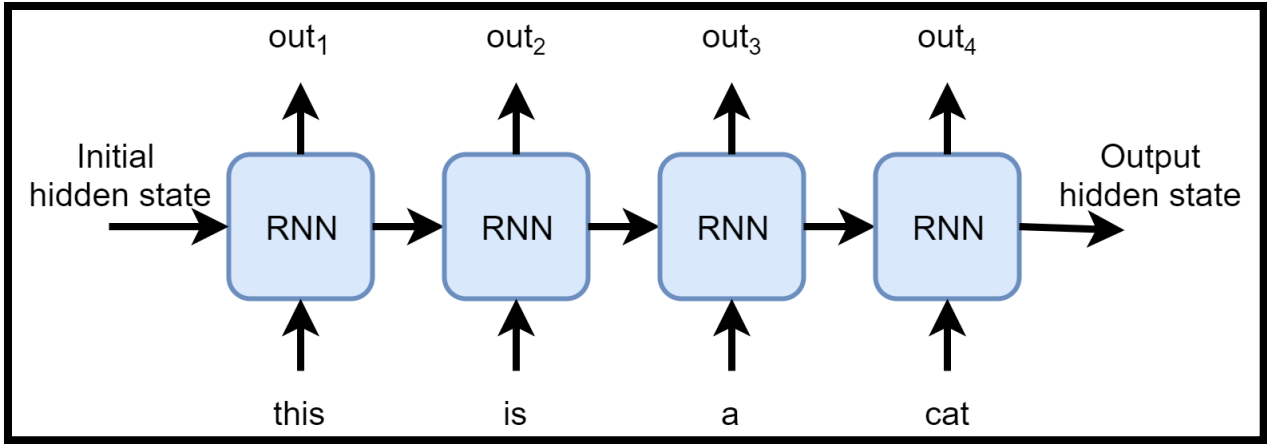
# Gradients parallelism

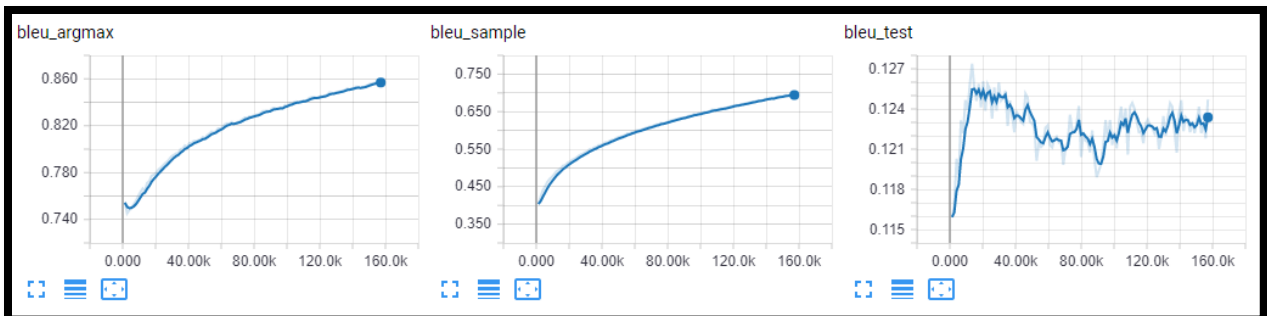
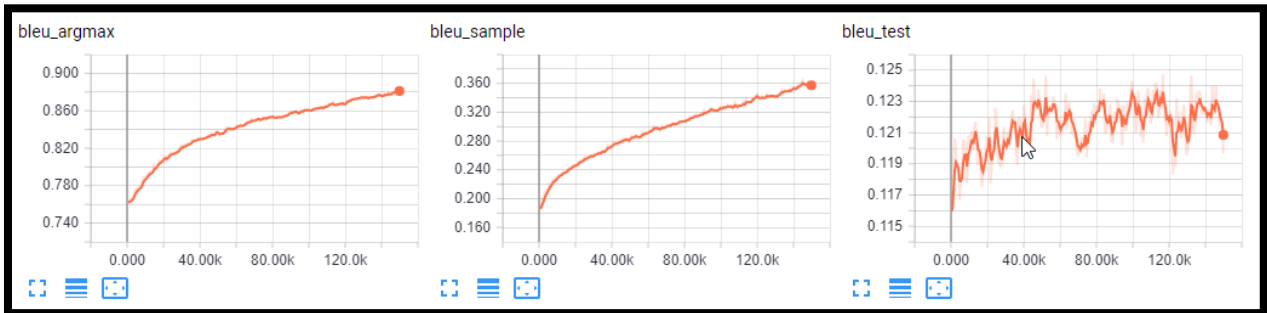
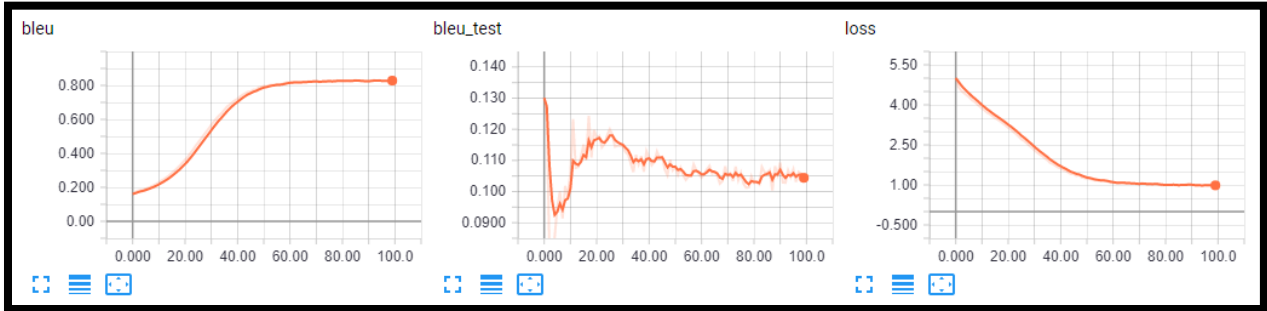
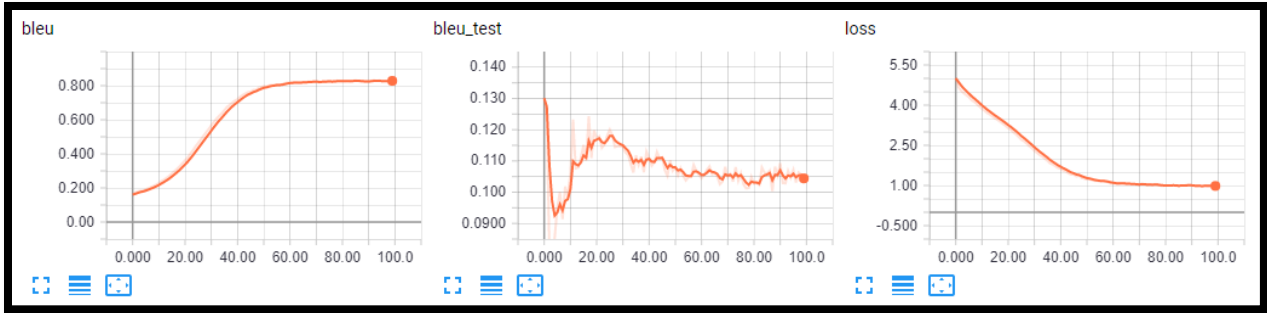


## Chapter 12: Chatbots Training with RL











**Shmuma**

12:32:19 PM

/bot how are you?



**rl\_bot\_ch12**

12:32:36 PM

very well. thank you.



**Shmuma**

12:32:29 PM

/bot are you going to enslave humanity?



**rl\_bot\_ch12**

12:32:46 PM

yes ... i don't know.



**Shmuma**

12:32:36 PM

/bot are you sure?



**rl\_bot\_ch12**

12:32:54 PM

yeah.



**Shmuma**

12:32:43 PM

/bot why?

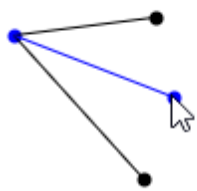

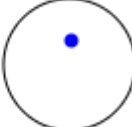
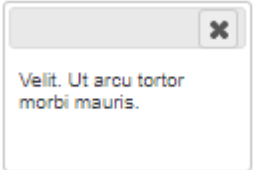


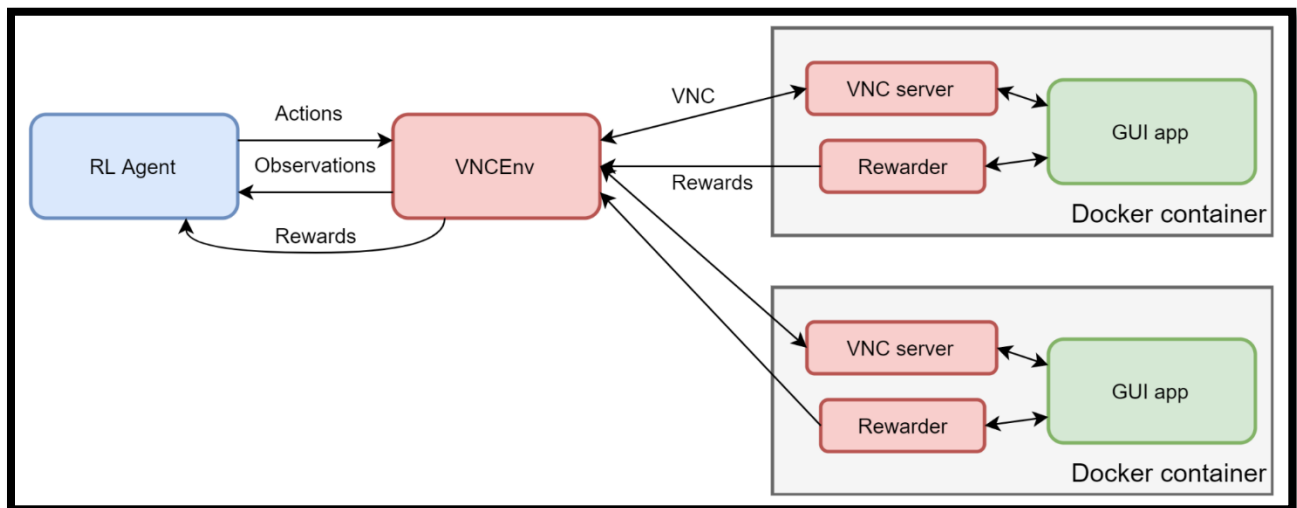
**rl\_bot\_ch12**

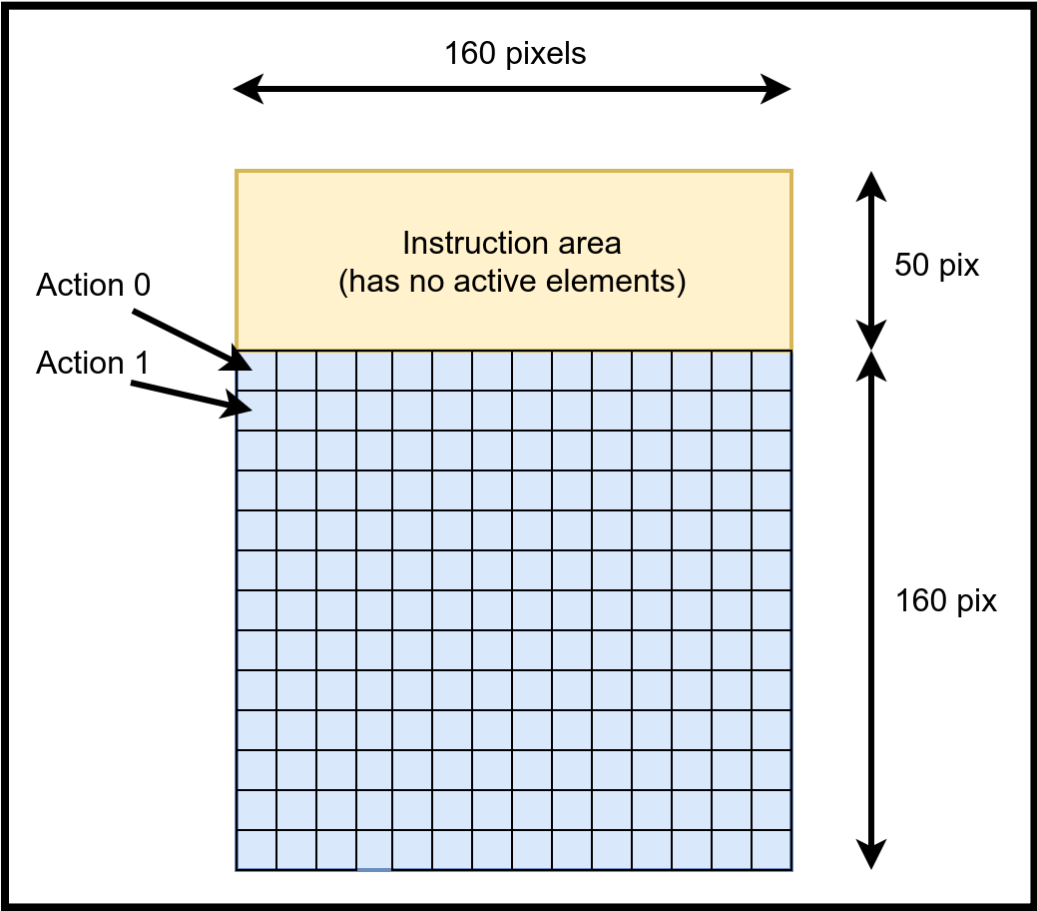
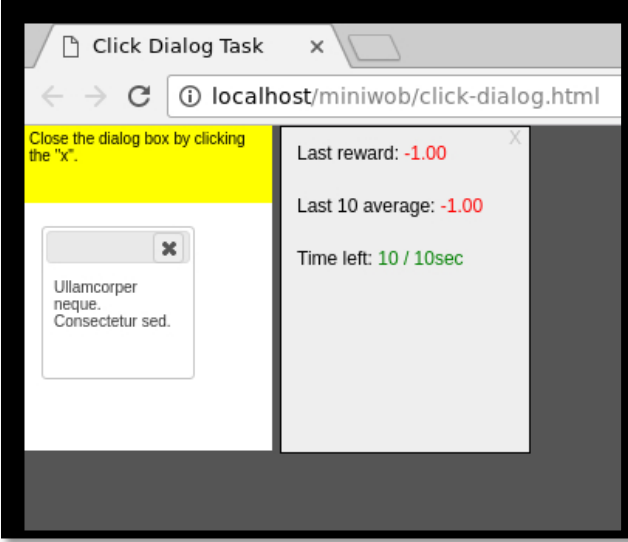
12:33:00 PM

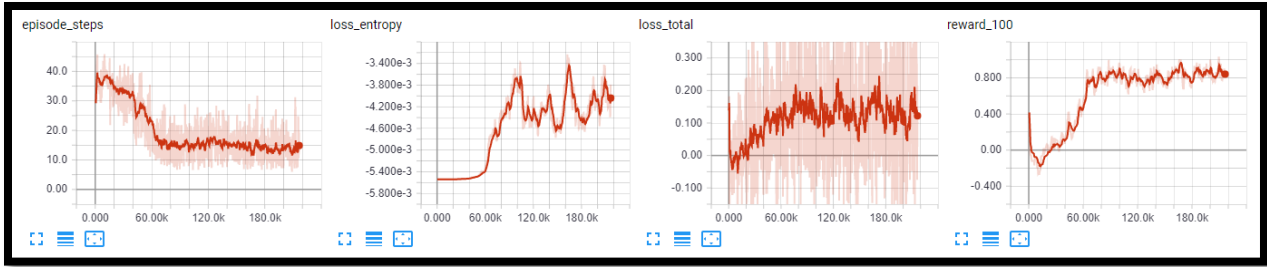
because i want to.

# Chapter 13: Web Navigation

<p>Create a line that bisects the angle evenly in two, then press submit.</p>  <p>Submit</p>	<p>Select 10/04/2016 as the date and hit submit.</p> <p>Date: <input type="text"/></p> 	<p>Select Caty from the list and click Submit.</p> <p>Caty</p> <ul style="list-style-type: none"> <li>Caty</li> <li>Leona</li> <li>Maud</li> <li>Merilyn</li> <li>Athena</li> <li>Dena</li> <li>Ilise</li> </ul>
<p>Find and click on the center of the circle, then press submit.</p>  <p>Submit</p>	<p>Expand the section below and click submit.</p> <p><b>Section #24</b></p> <p>Euismod. Sed tincidunt interdum. Interdum maecenas ut nibh risus massa facilisis elementum ullamcorper quisque quis porttitor. Metus feugiat lectus est.</p> <p>Submit</p>	<p>Close the dialog box by clicking the "x".</p> 







Close the dialog box by clicking the "x".

✕

Dolor. Sed. Nam  
interdum morbi  
turpis. Ac.

Last reward: -1.00

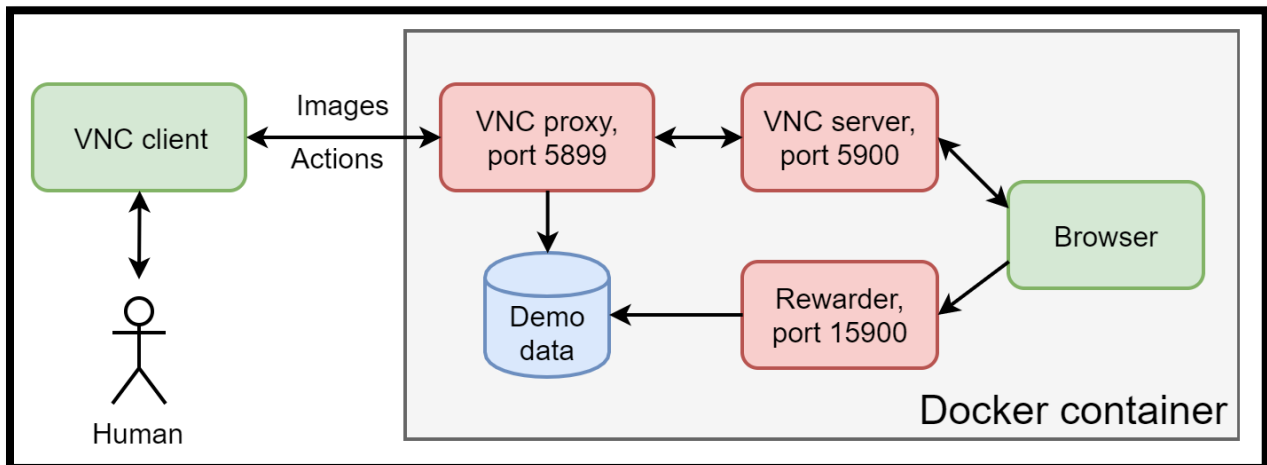
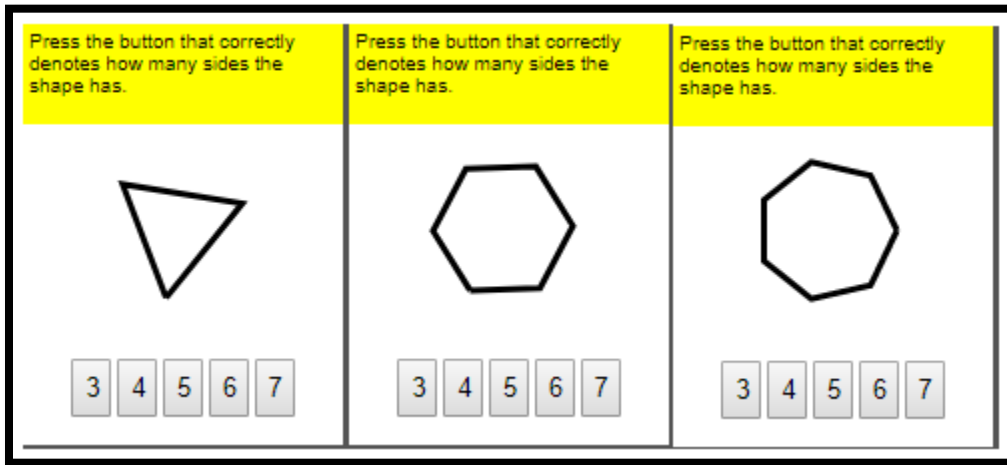
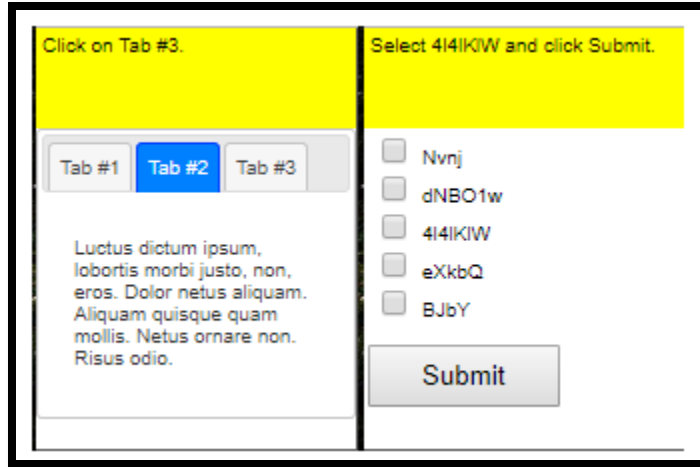
Last 10 average: -1.00

Time left: 2 / 10sec

Click button ONE, then click button TWO.

TWO

ONE



Press the button that correctly denotes how many sides the shape has.

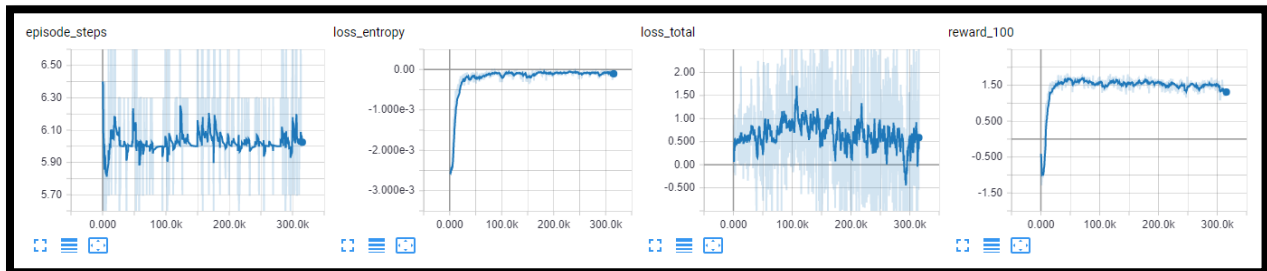
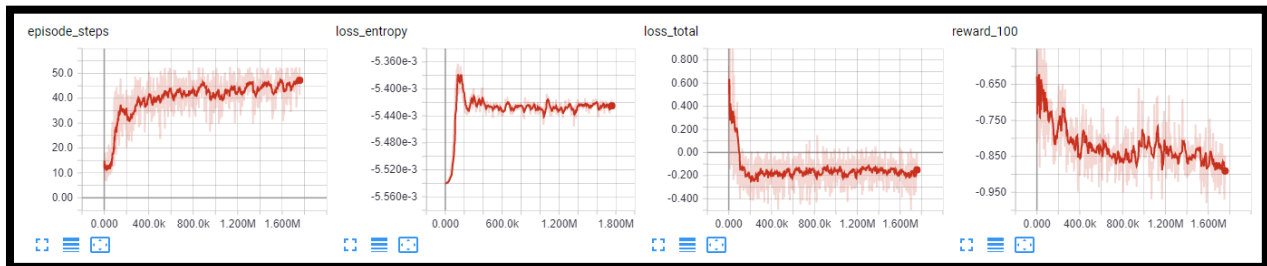
Press the button that correctly denotes how many sides the shape has.

Press the button that correctly denotes how many sides the shape has.

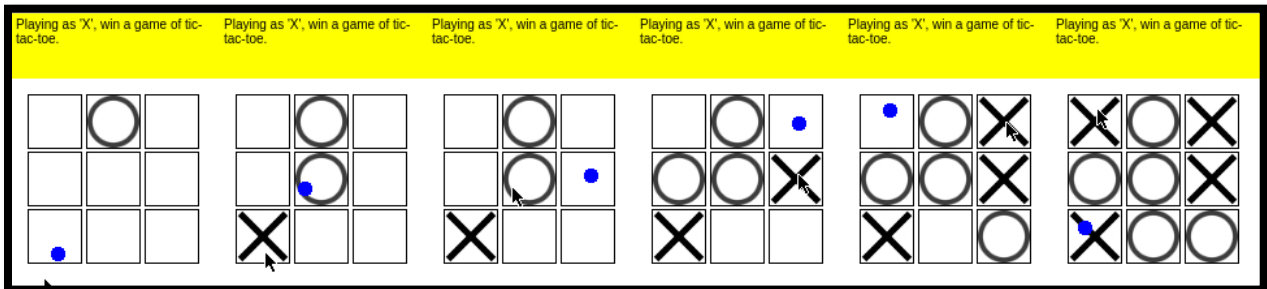
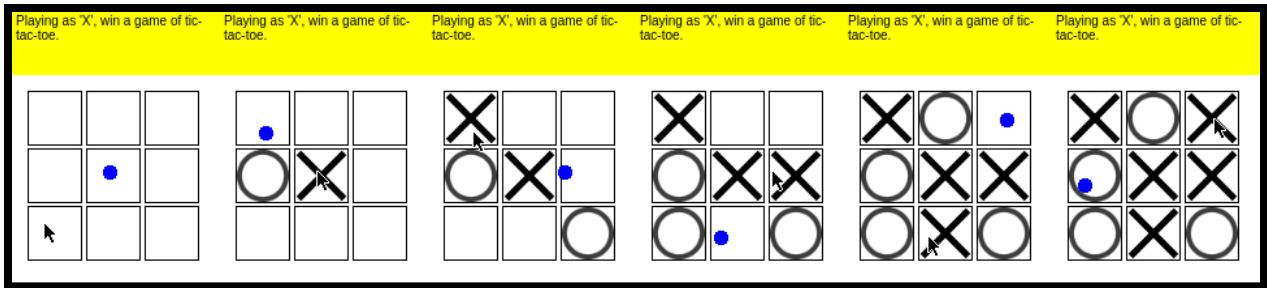
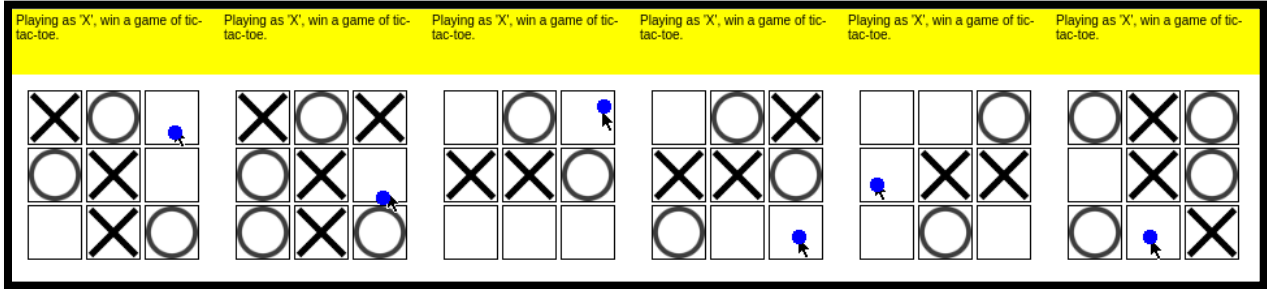
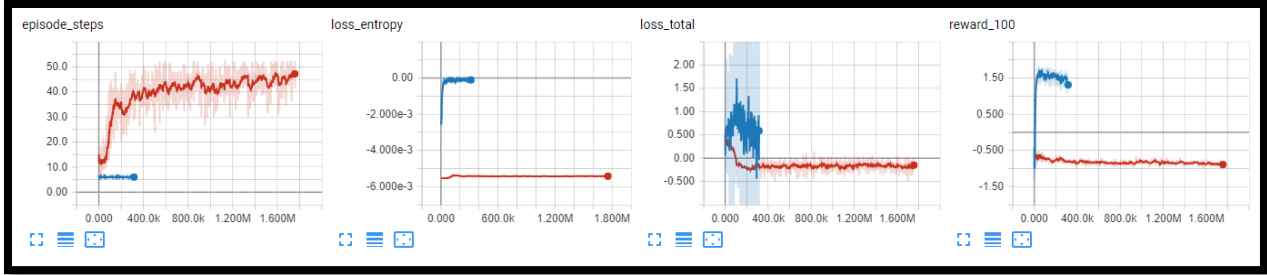
Press the button that correctly denotes how many sides the shape has.

Press the button that correctly denotes how many sides the shape has.

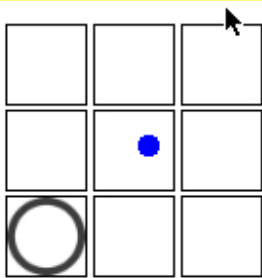
Press the button that correctly denotes how many sides the shape has.



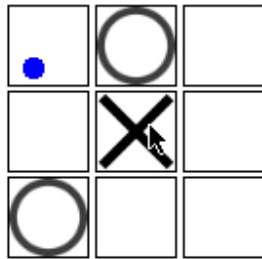




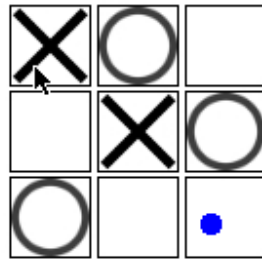
Playing as 'X', win a game of tic-tac-toe.



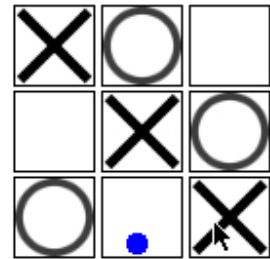
Playing as 'X', win a game of tic-tac-toe.



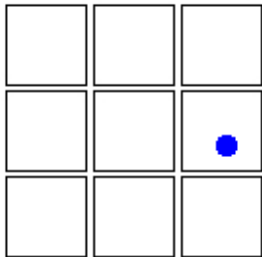
Playing as 'X', win a game of tic-tac-toe.



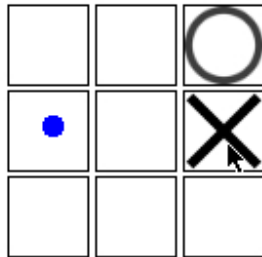
Playing as 'X', win a game of tic-tac-toe.



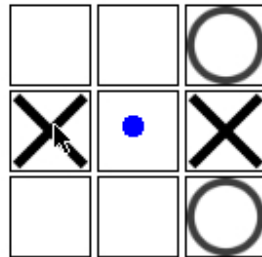
Playing as 'X', win a game of tic-tac-toe.



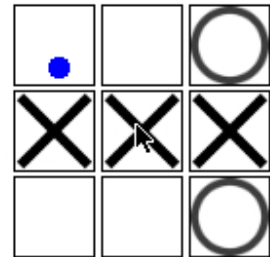
Playing as 'X', win a game of tic-tac-toe.



Playing as 'X', win a game of tic-tac-toe.



Playing as 'X', win a game of tic-tac-toe.

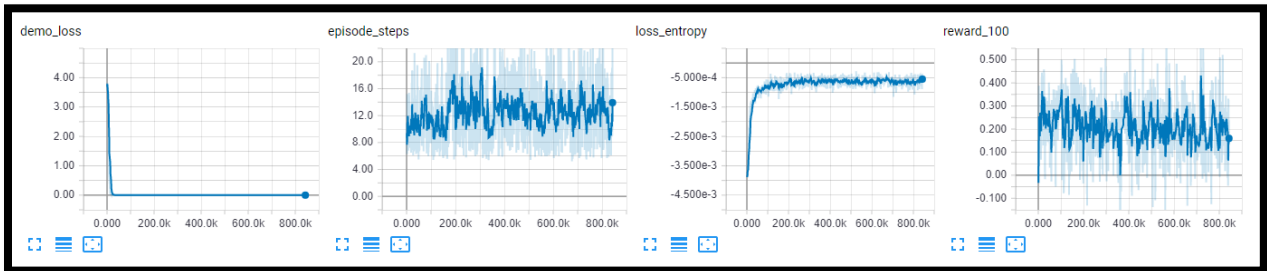
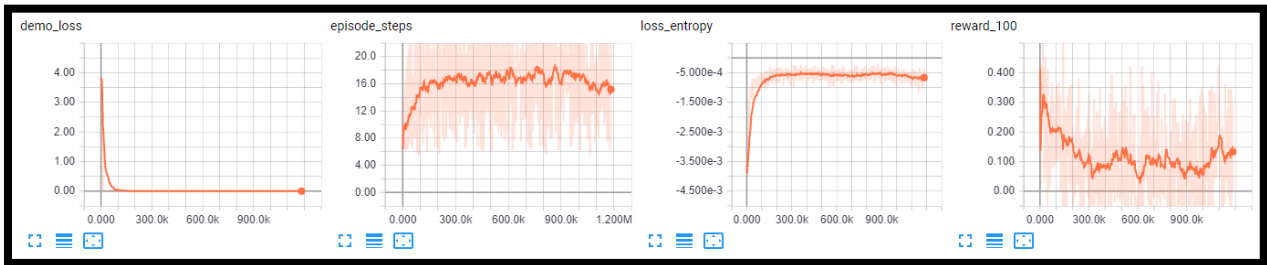


Click on the "Yes" button.      Click on the "Next" button.      Click on the "Previous" button.

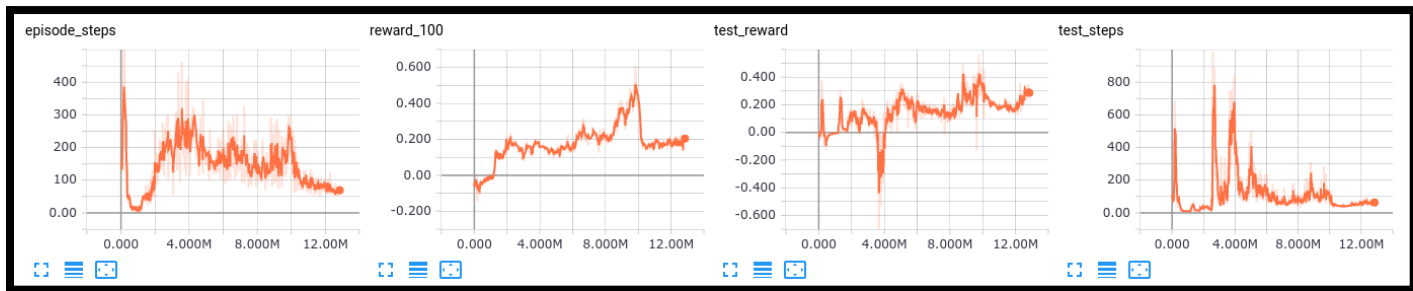
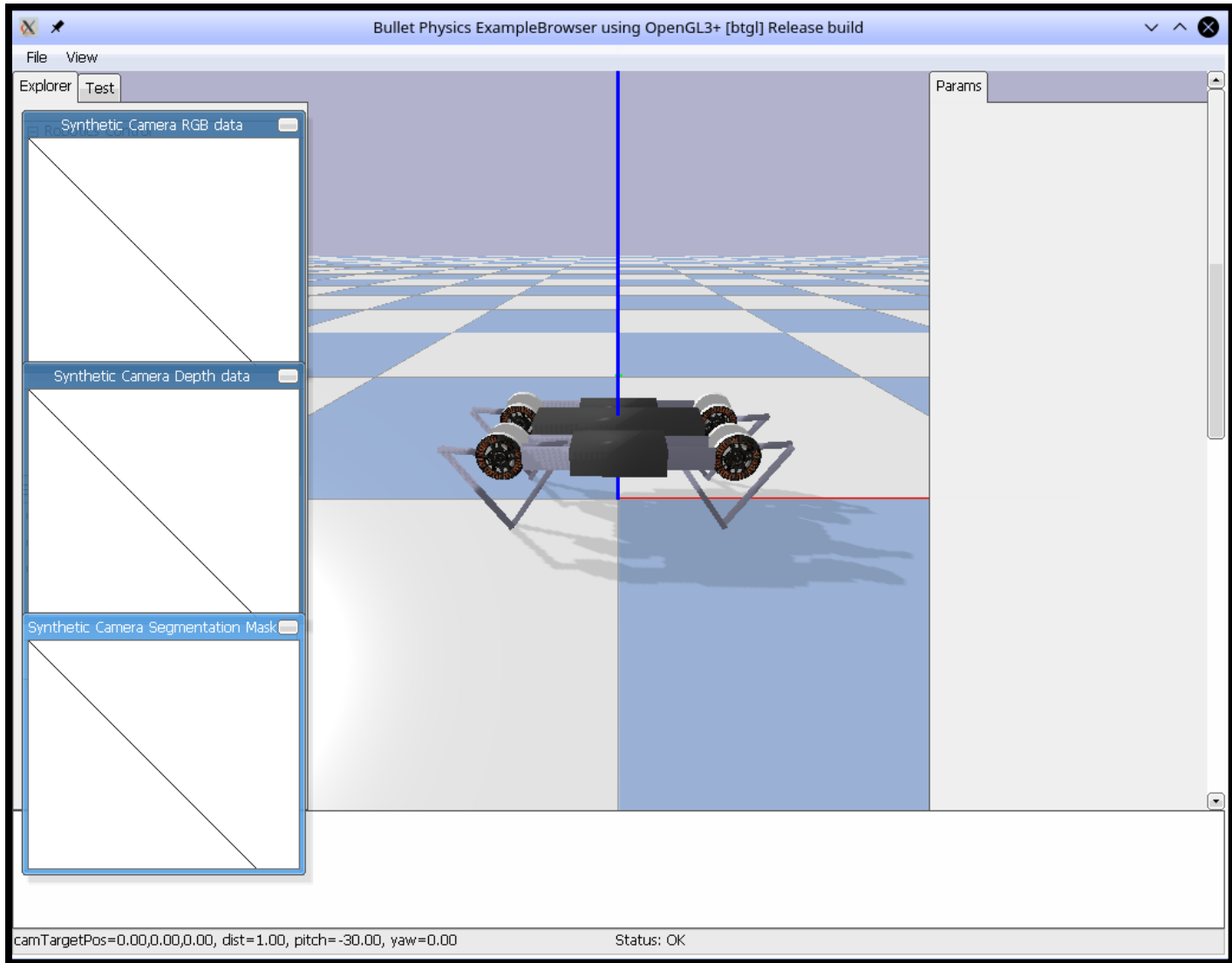
     ornare urna sit tempor, et, enim:       consequat massa auctor:   
      tellus fames tellus       yes venenatis scelerisque   
            aliquam:   
donec placerat aliquet       quis sociis vitae augue semper adipiscing      nunc, elit sapien   
urna id augue

Click on the "Okay" button.      Click on the "Submit" button.      Click on the "Submit" button.

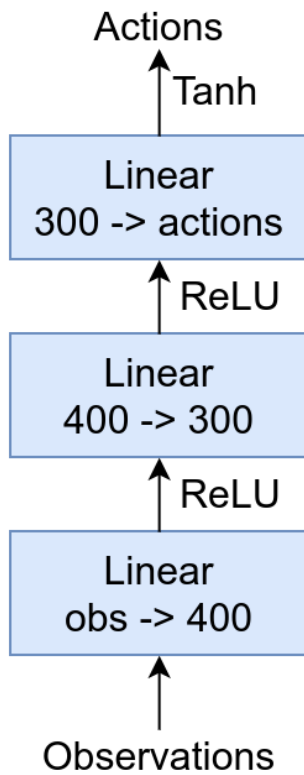
adipiscing adipiscing libero       sem nulla ac         
      ac imperdiet nulla:       No tortor sit in:   
duis odio vel:       suspendisse lacinia tempor tortor pharetra, porta         
      quam in tempus:         
egestas sit aliquet



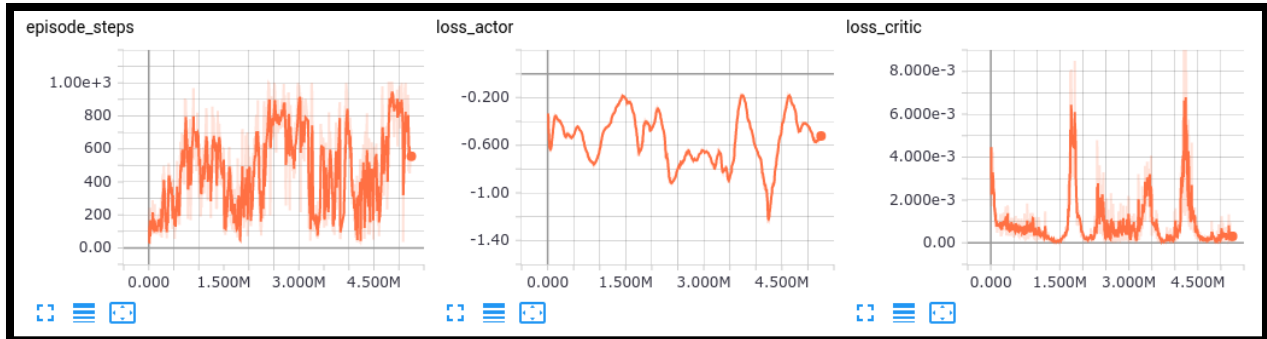
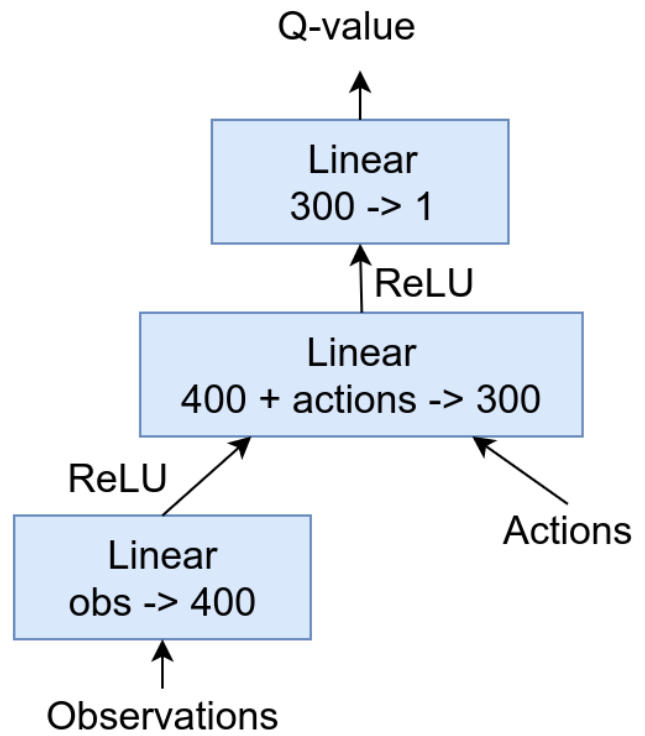
# Chapter 14: Continuous Action Space

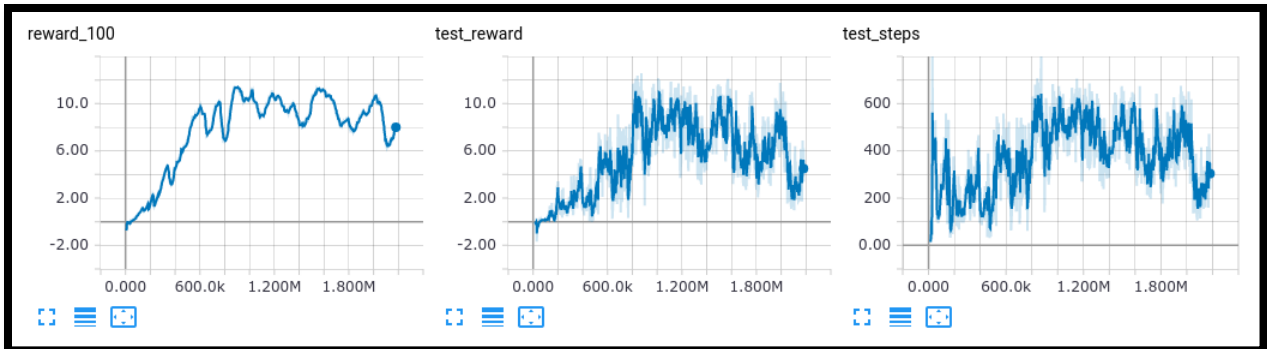
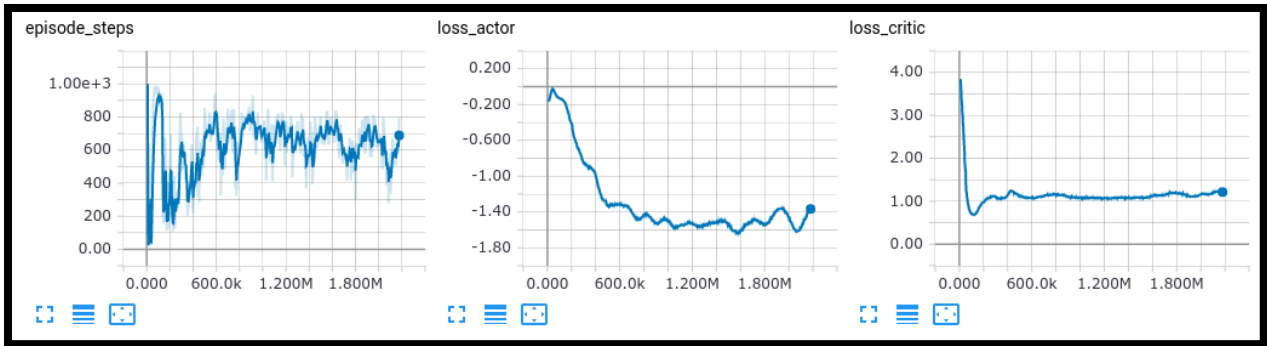
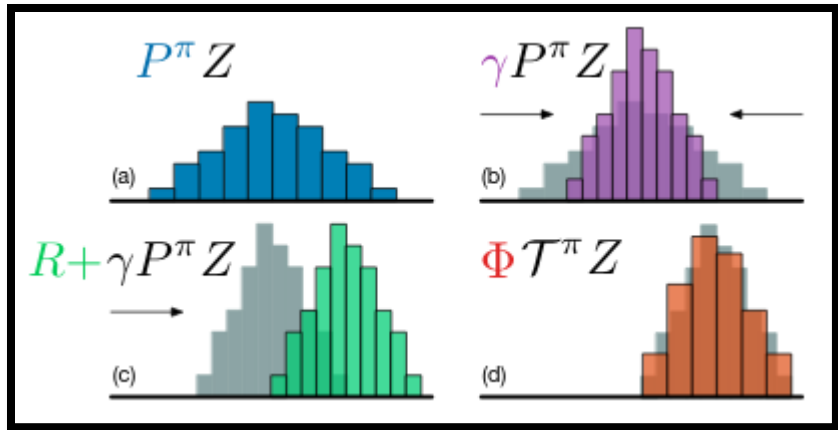
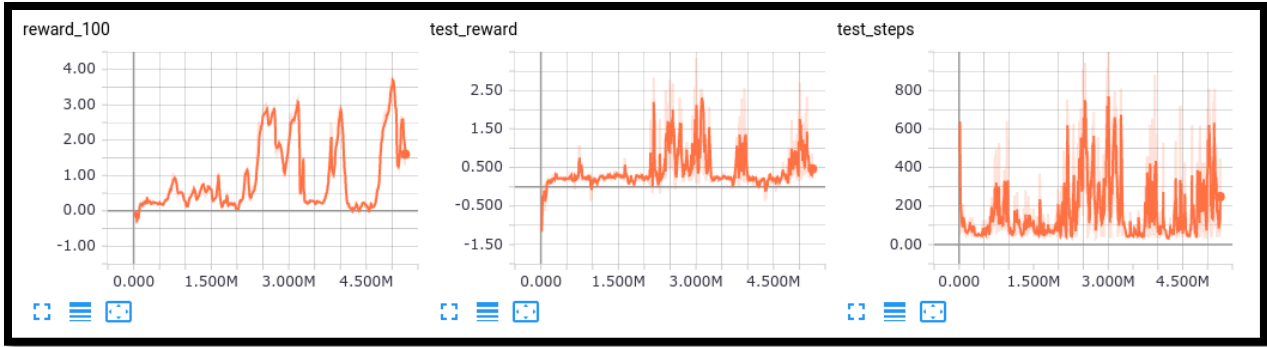


## Actor network

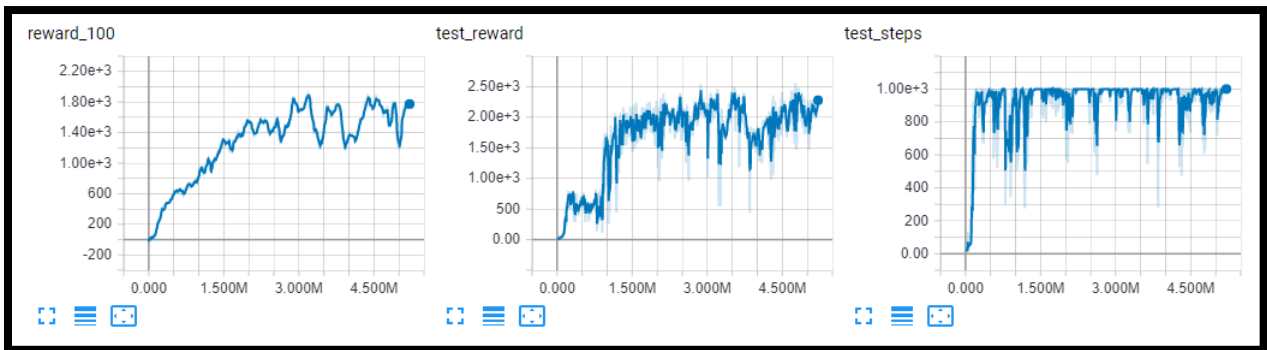
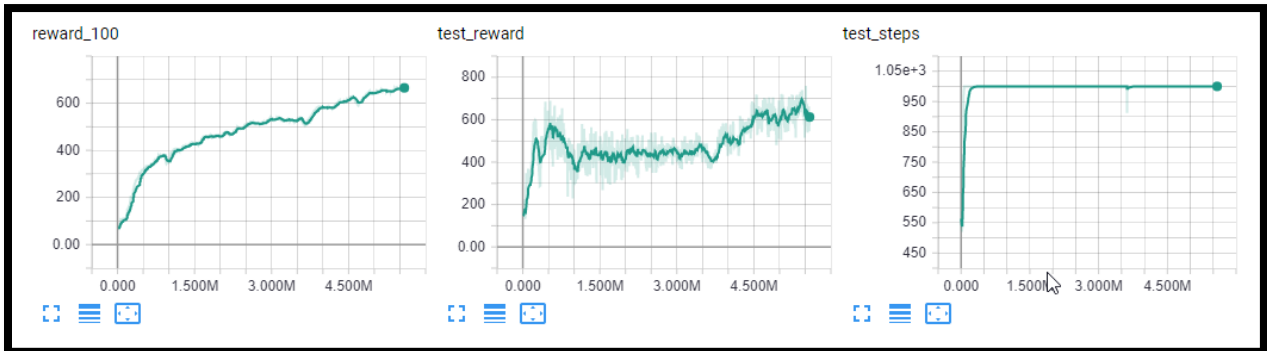
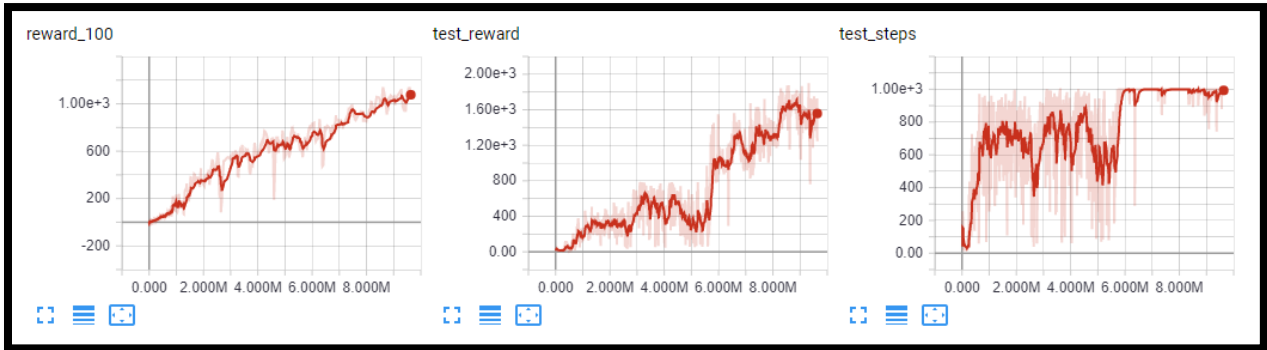
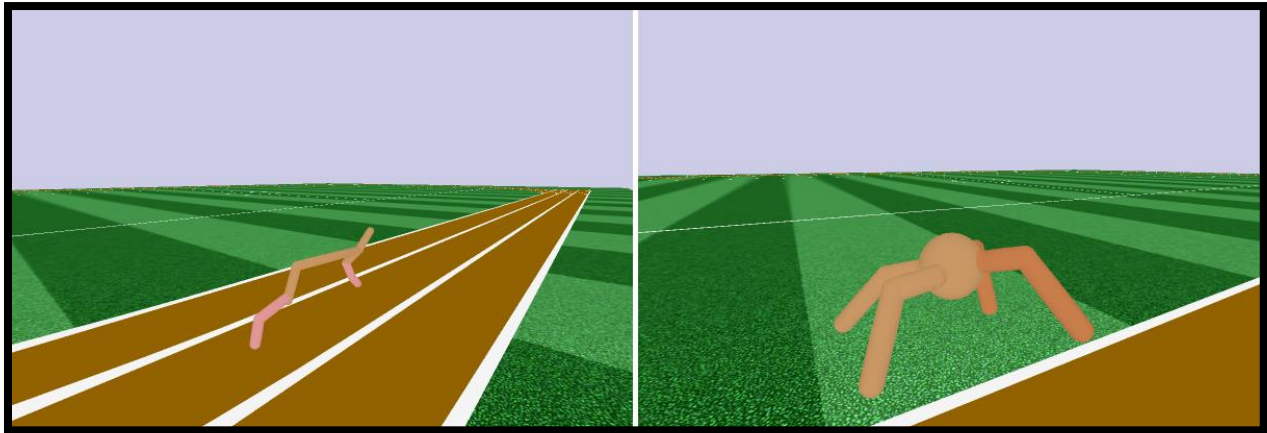


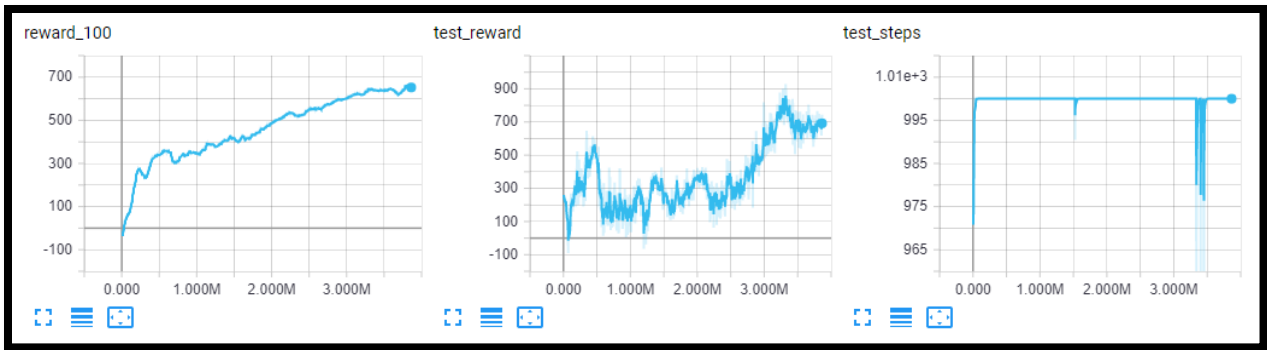
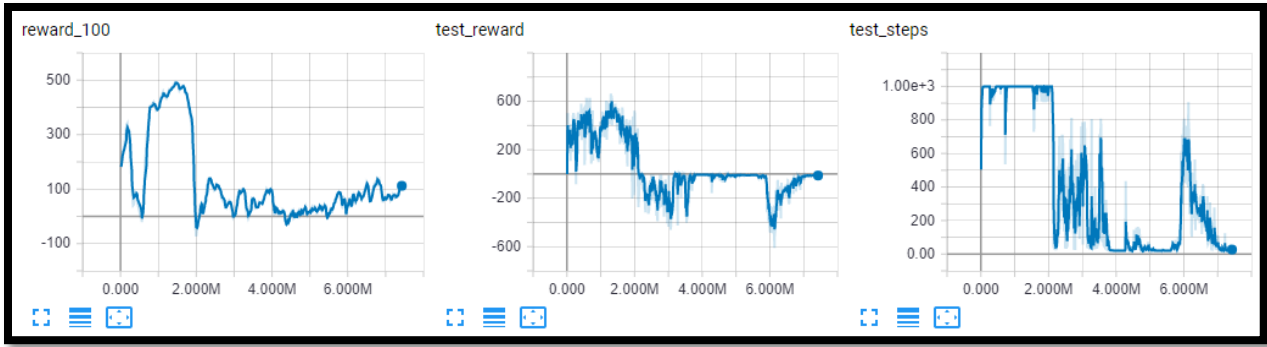
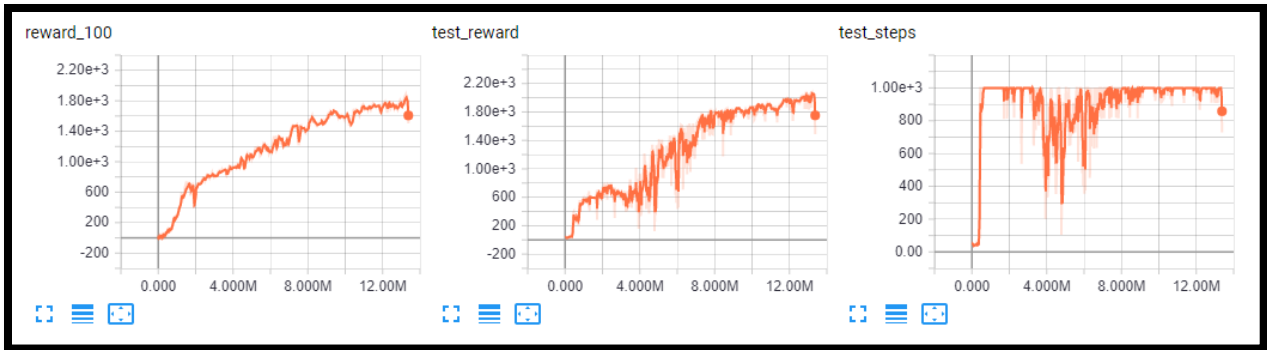
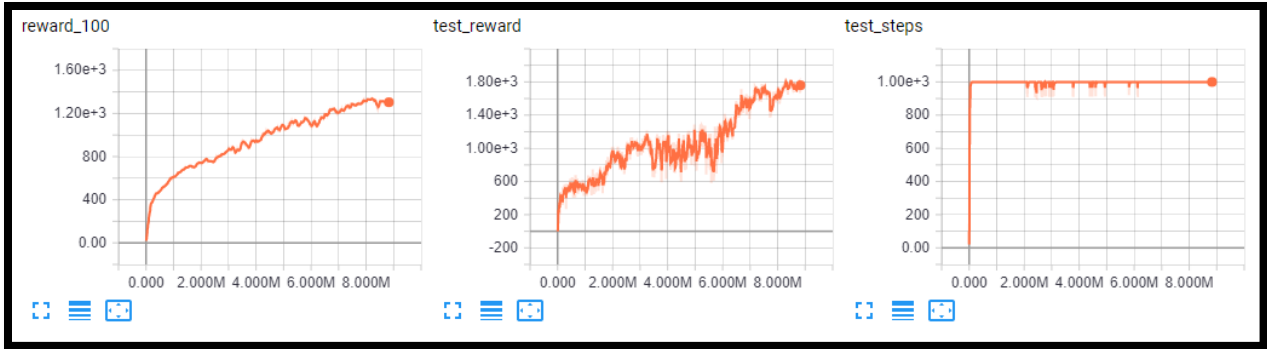
## Critic network





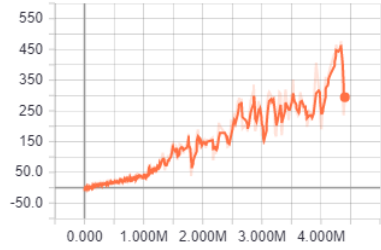
# Chapter 15: Trust Regions – TRPO, PPO, and ACKTR



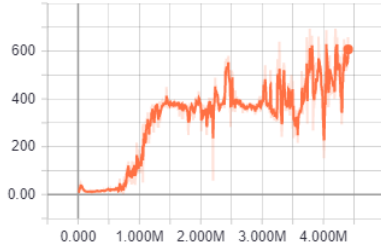




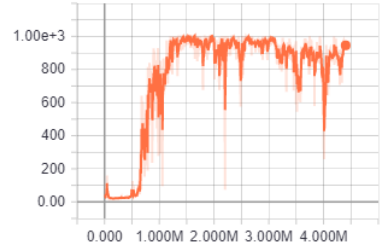
reward\_100



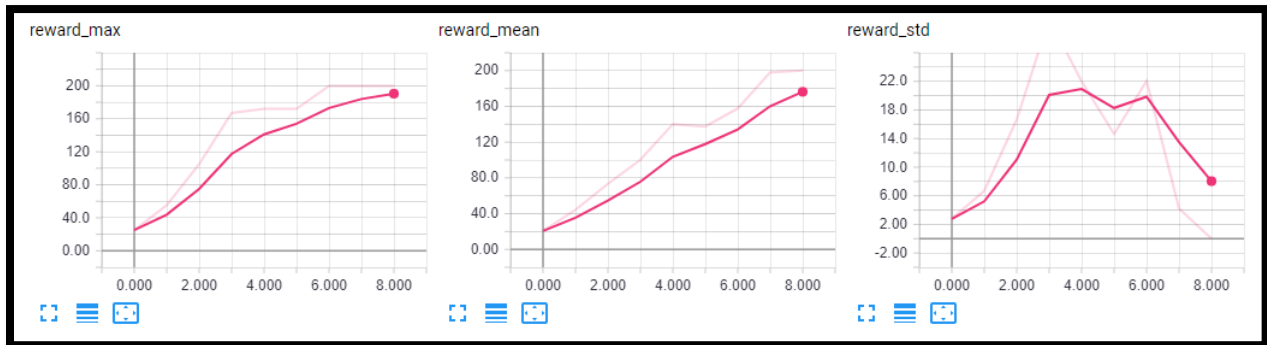
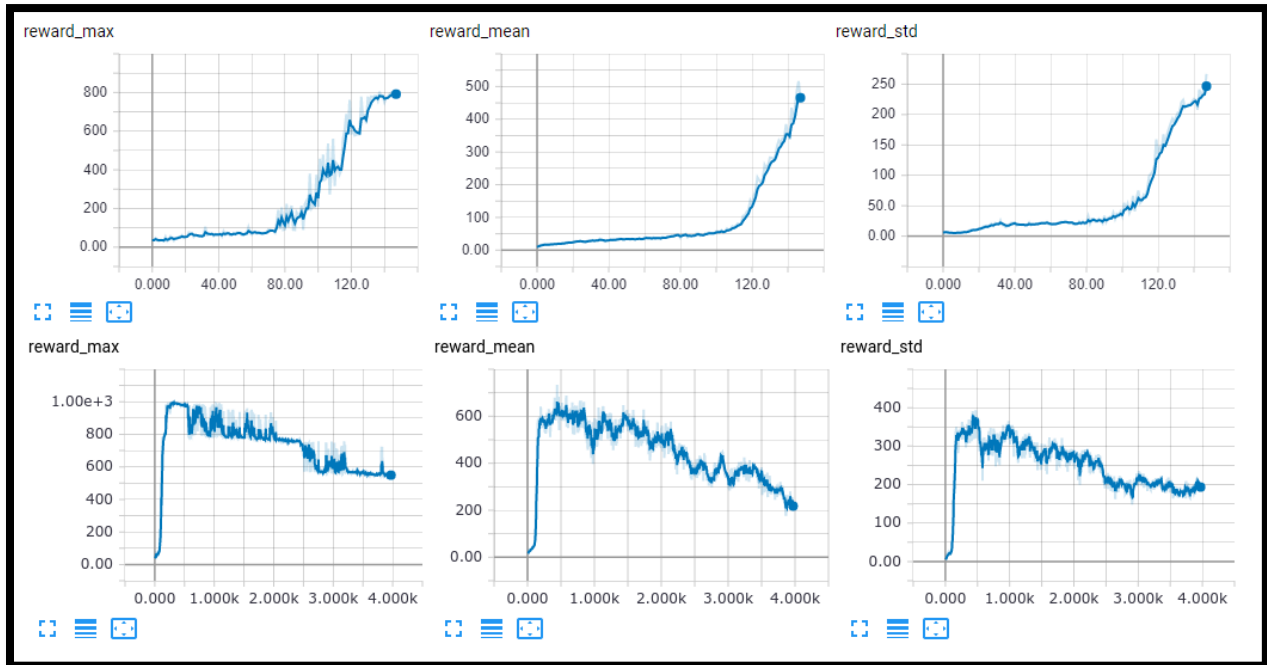
test\_reward

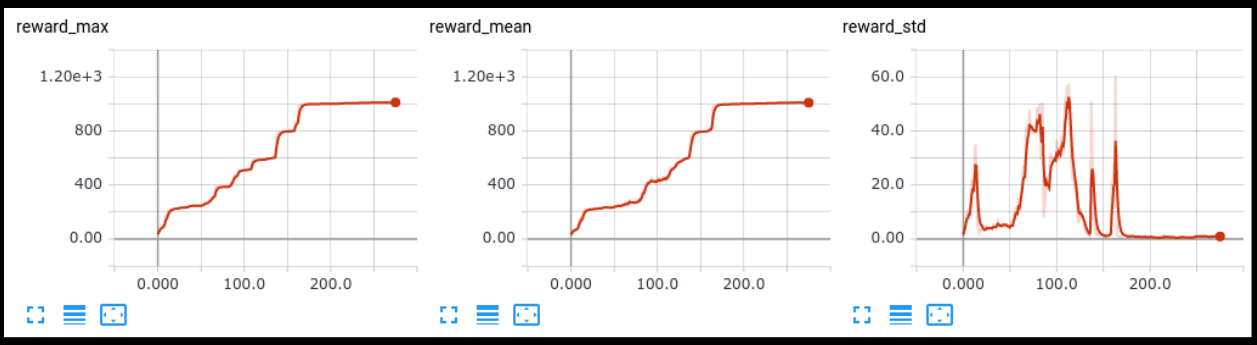


test\_steps

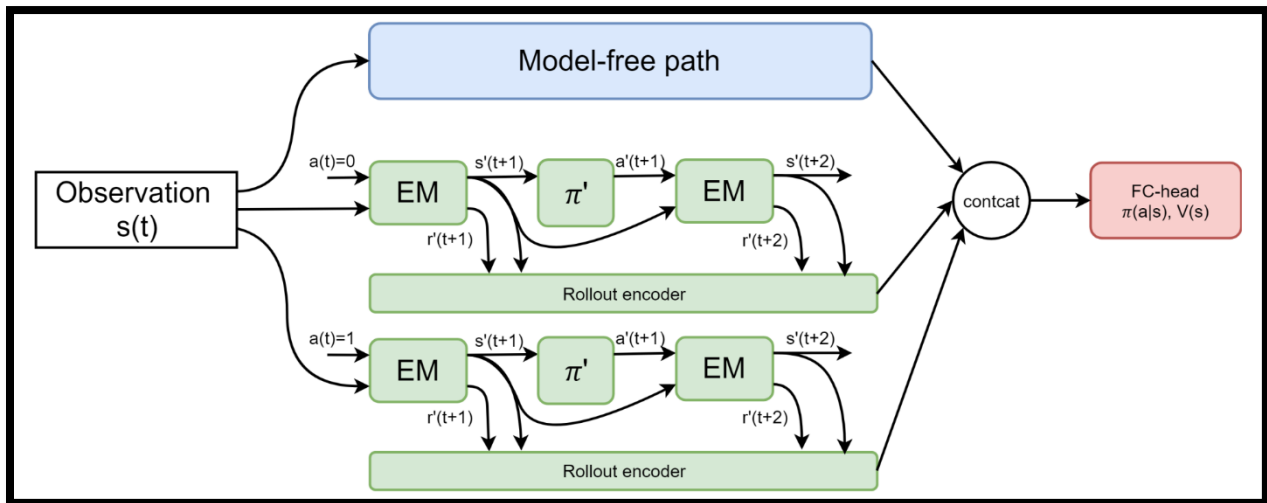
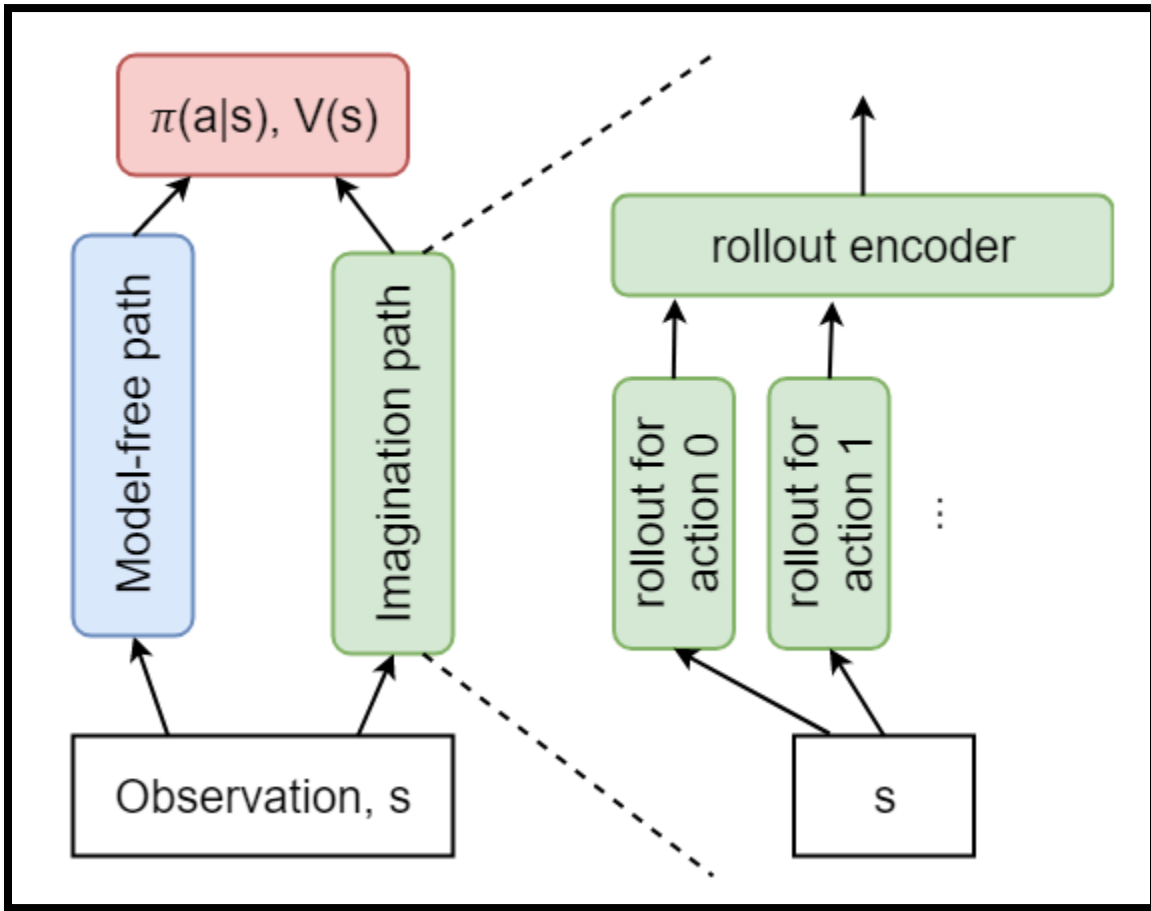


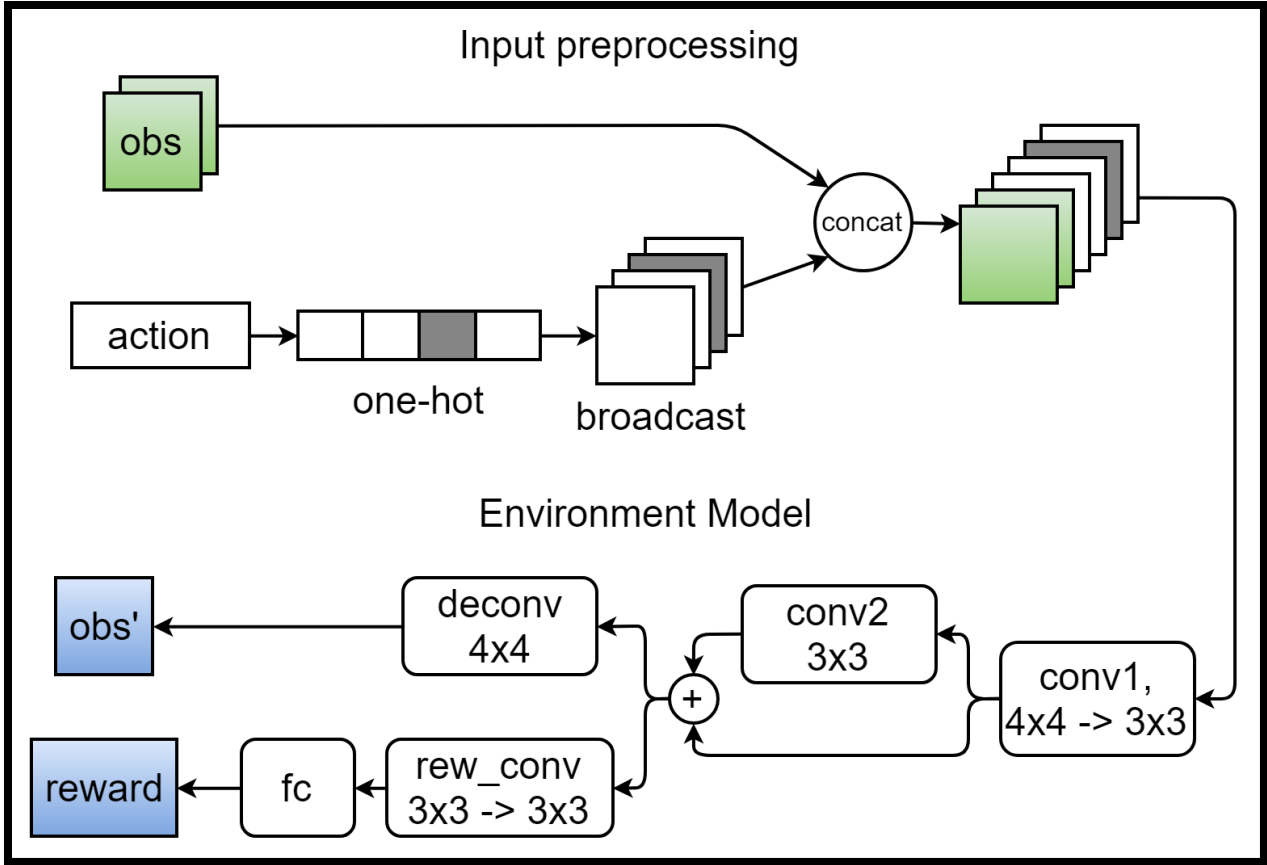
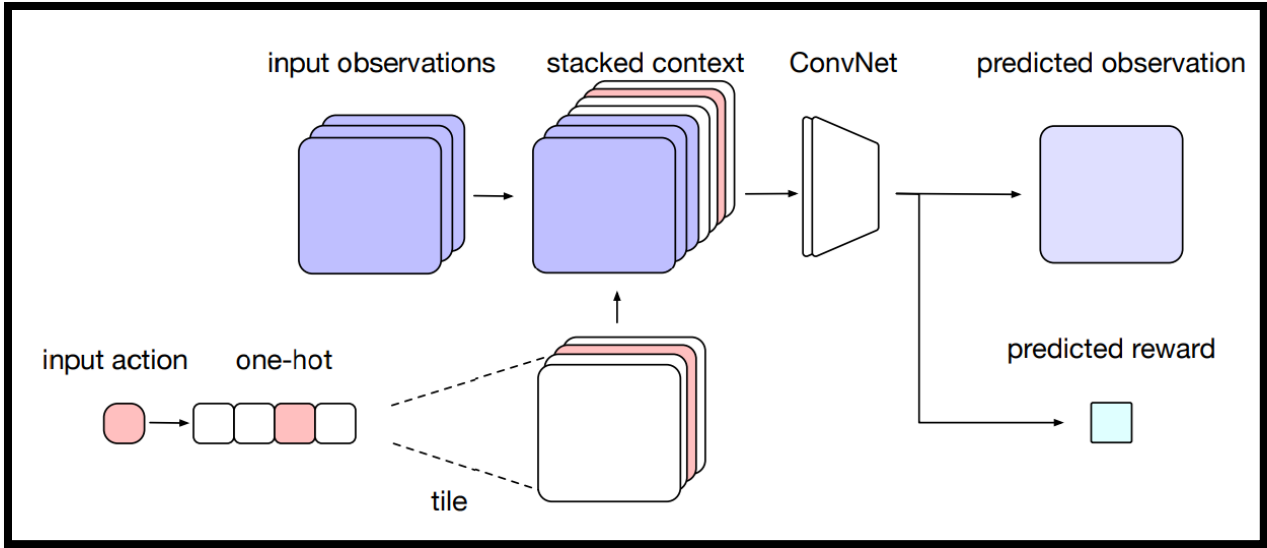
# Chapter 16: Black-Box Optimization in RL

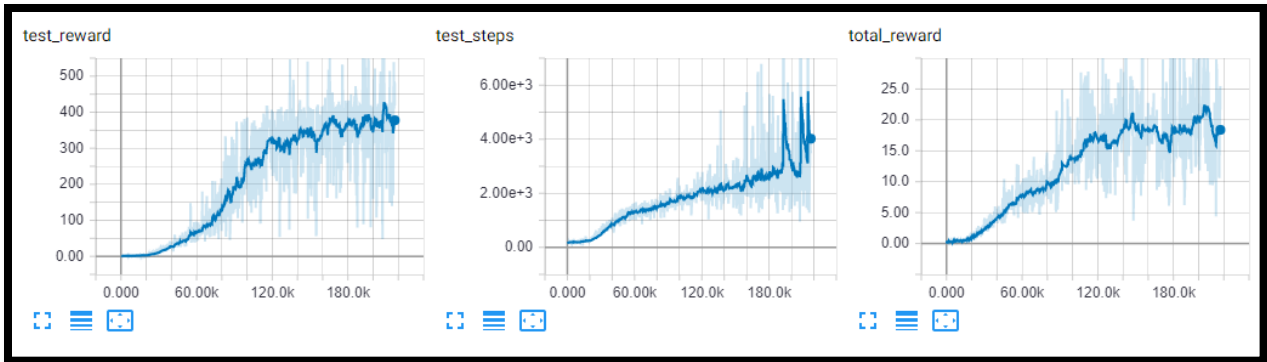
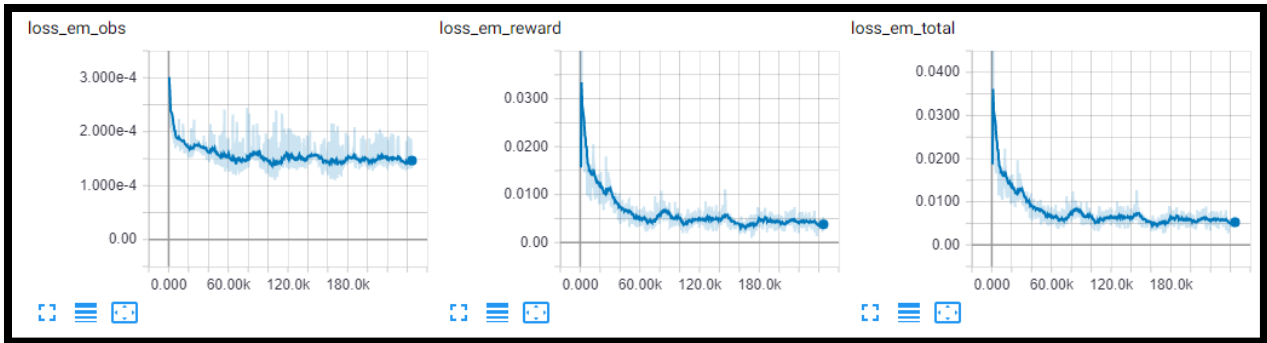
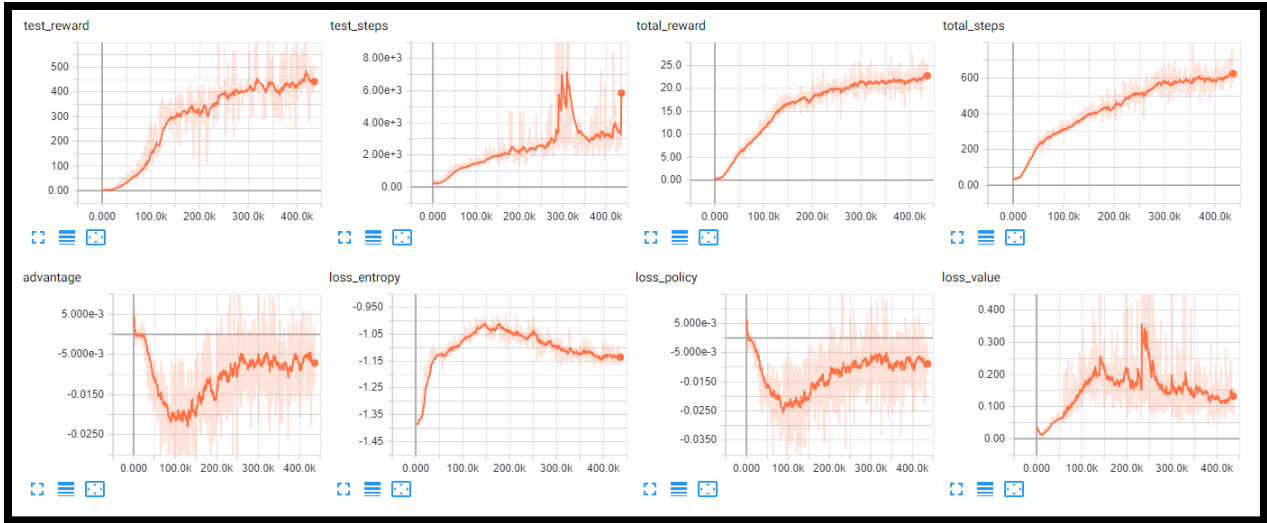


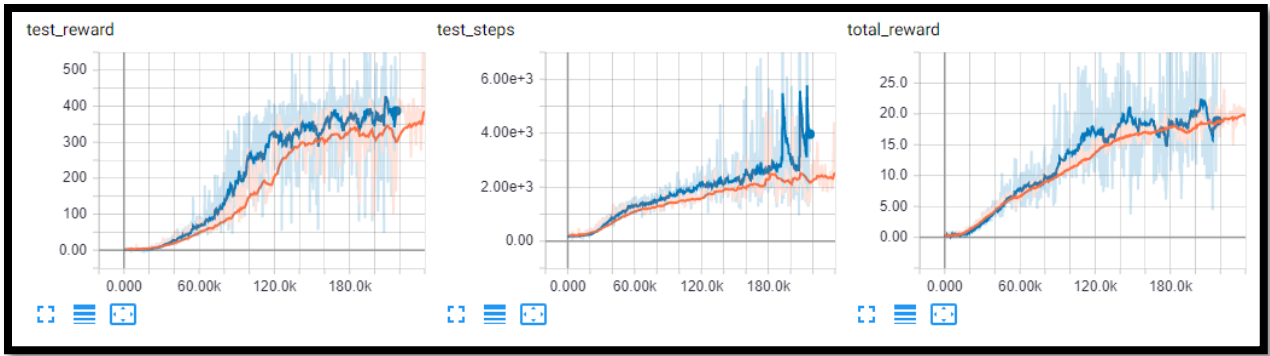


# Chapter 17: Beyond Model-Free – Imagination

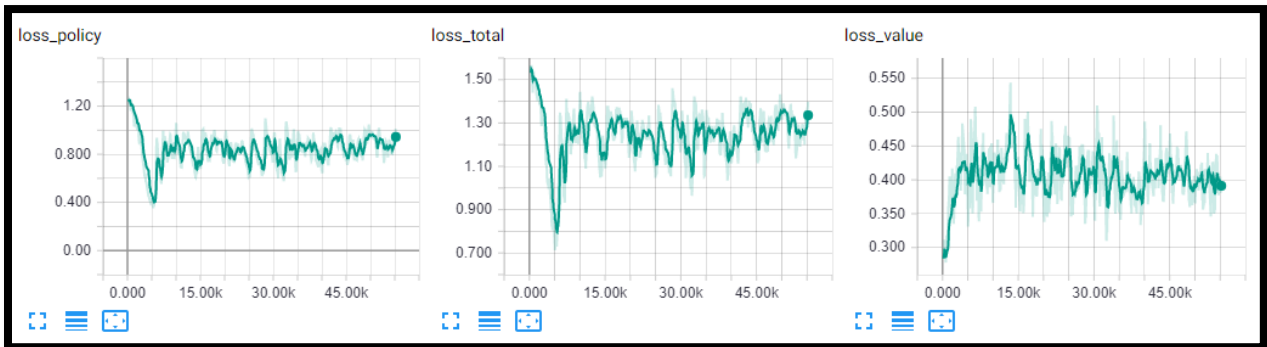
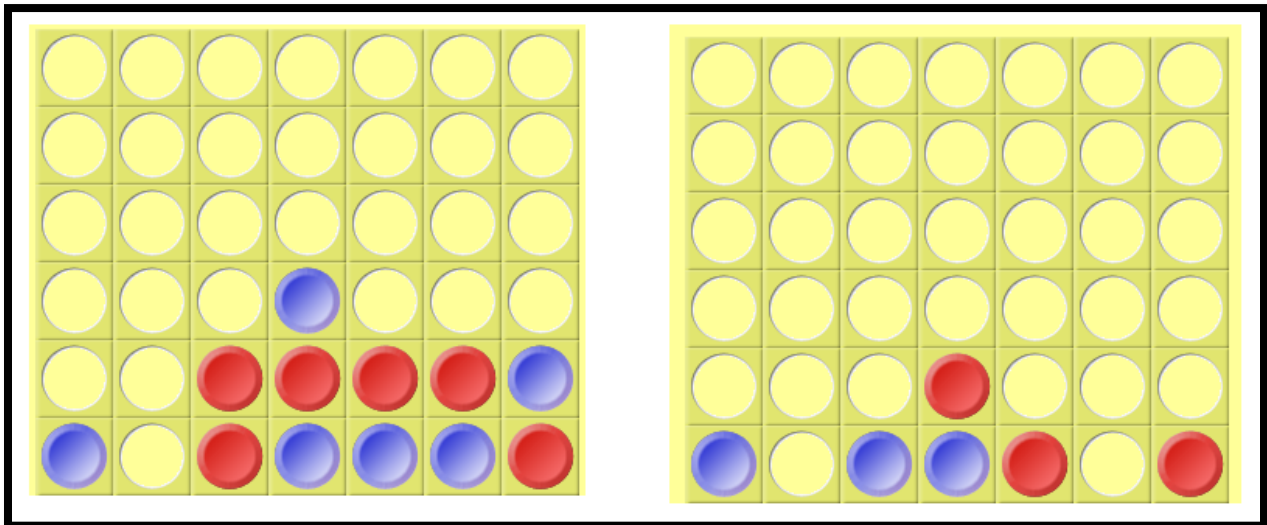
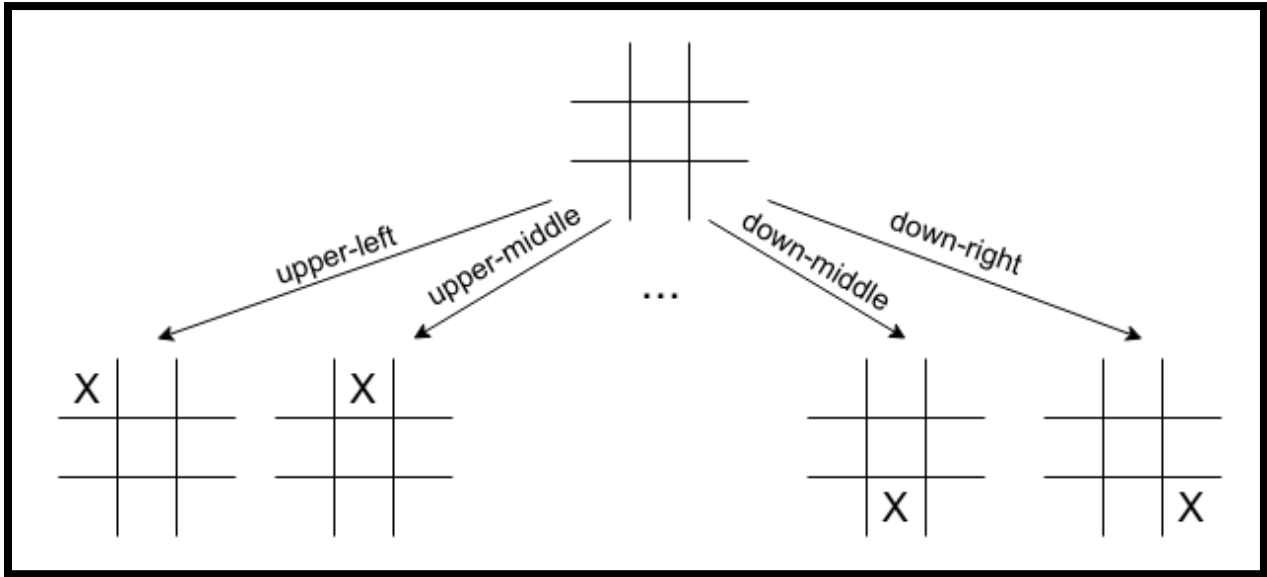




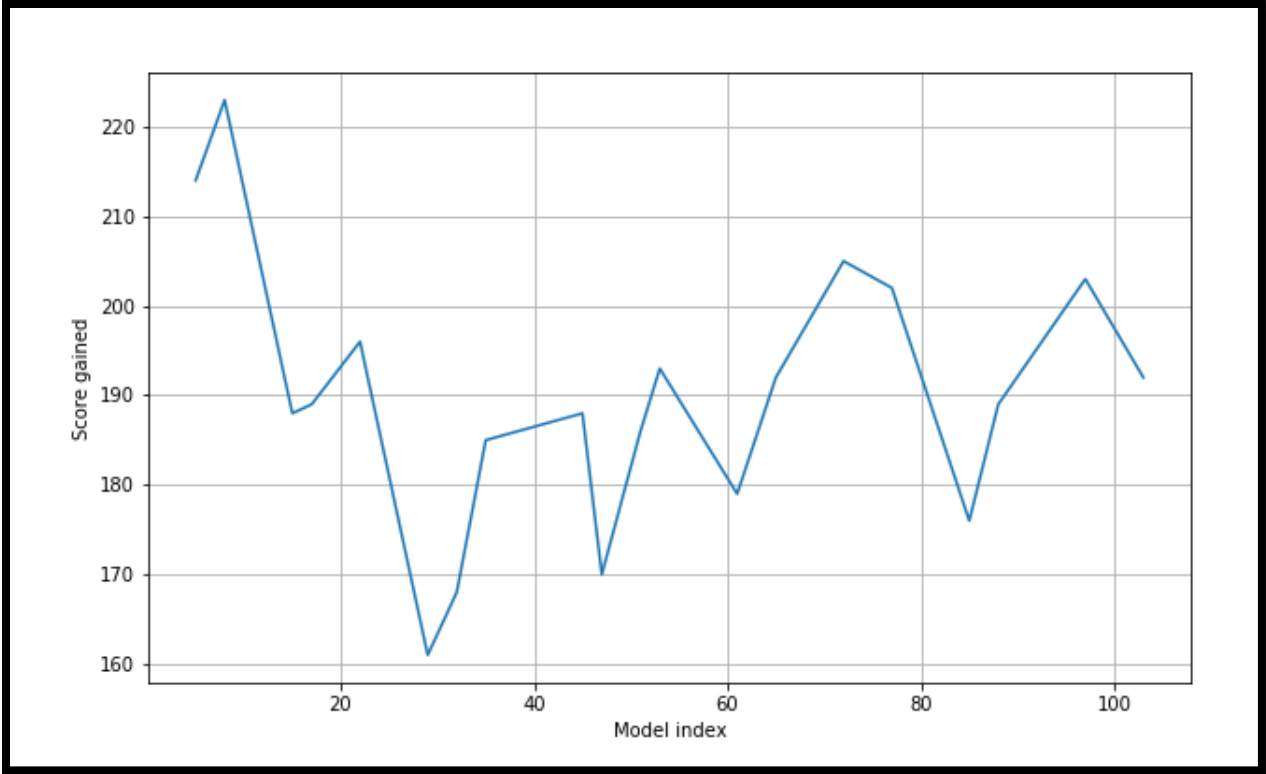




# Chapter 18: AlphaGo Zero









**Shmuma**

/top



**ri\_ch18\_bot**

Leader board

```
Luklusik: won=31, lost=9, draw=0
best_008_02500.dat: won=8, lost=8, draw=0
    Shmuma: won=6, lost=6, draw=0
best_019_11300.dat: won=2, lost=4, draw=0
best_015_09200.dat: won=2, lost=4, draw=0
best_005_01900.dat: won=2, lost=6, draw=0
    mbilan: won=2, lost=3, draw=0
best_007_02400.dat: won=1, lost=0, draw=0
    Roman8i: won=1, lost=3, draw=0
best_010_05900.dat: won=1, lost=0, draw=0
best_001_00800.dat: won=1, lost=2, draw=0
best_014_08600.dat: won=1, lost=0, draw=0
best_004_01600.dat: won=1, lost=0, draw=0
best_009_05700.dat: won=1, lost=2, draw=0
best_002_00900.dat: won=1, lost=5, draw=0
best_022_12200.dat: won=1, lost=5, draw=0
best_061_34100.dat: won=0, lost=2, draw=0
    None: won=0, lost=1, draw=0
best_047_28000.dat: won=0, lost=2, draw=0
```