

# **REINFORCEMENT LEARNING LABORATORY**

## **Atari Freeway-v0**

**Eng. Isac-Daniel CIOBOTĂ**

June 2024

## Table of Contents

1.	Game Description.....	3
2.	Algorithms .....	4
2.1	REINFORCE.....	4
2.2	SARSA.....	4
2.3	Deep Q-Networks (DQN) .....	5
2.4	DQN with target networks .....	5
2.5	Double DQN (DDQN).....	6
2.6	Prioritized Experience Replay (DQN PER) .....	6
2.7	Advantage Actor-Critic (A2C).....	7
2.8	Proximal Policy Optimization (PPO).....	8
2.9	Asynchronous Advantage Actor-Critic (A3C) .....	9
3.	Experimental Results .....	10
3.1	REINFORCE.....	10
3.2	SARSA.....	11
3.3	DQN.....	13
3.4	DQN with target networks .....	13
3.5	DDQN .....	15
3.6	DQN PER .....	17
3.7	A2C n-step returns.....	17
3.8	A2C GAE .....	19
3.9	PPO .....	21
3.10	A3C .....	21
4.	Conclusions.....	24
5.	Reference List .....	25

## 1. Game Description

Atari Freeway is a classic game developed by David Crane and published Activision in 1981. The game was built for the Atari 2600 video game console. The goal is to help the chicken cross the road (a multi-lane highway) without being crushed by a car. The chicken starts from the bottom of the screen and has to cross ten lanes of traffic (five in each direction). The action space is discrete, having three possible actions: UP, DOWN and NOOP (not moving) [1]. If you choose difficulty to be easy, every time you collide with a car, you are knocked back one lane. Otherwise, if you choose difficulty to be hard, every time you collide with a car, you are knocked back to the starting position. Every time you manage to cross all the lanes, you are awarded 1 point. The score is visible at the top of the screen. The game also has a multiplayer mode. Two players can compete against each other. The one who has the most points at the end of the game will be the winner. The total duration of the game is 2 minutes. When you approach the final seconds, the score will start blinking [2]. Screenshots from the game can be observed in Figure 1.

You can select one of the eight available maps (Lake Shore Drive – Chicago – 3 A.M., Interstate 5 – Seattle – 6 A.M., Santa Monica Freeway - Los Angeles - 10 A.M., Bayshore Freeway - San Francisco – Midnight, John Lodge Expressway – Detroit - 9 P.M., The Beltway - Washington D.C. - 6 P.M., LBJ Freeway - Dallas - 5 P.M., Long Island Expressway - New York City - 3 A.M.) [2]. Each map has different traffic intensities, car speed and vehicle types (automobiles and trucks). Maps 5-8 are identical to the first maps, but the speeds of the vehicles are randomly increased or decreased.

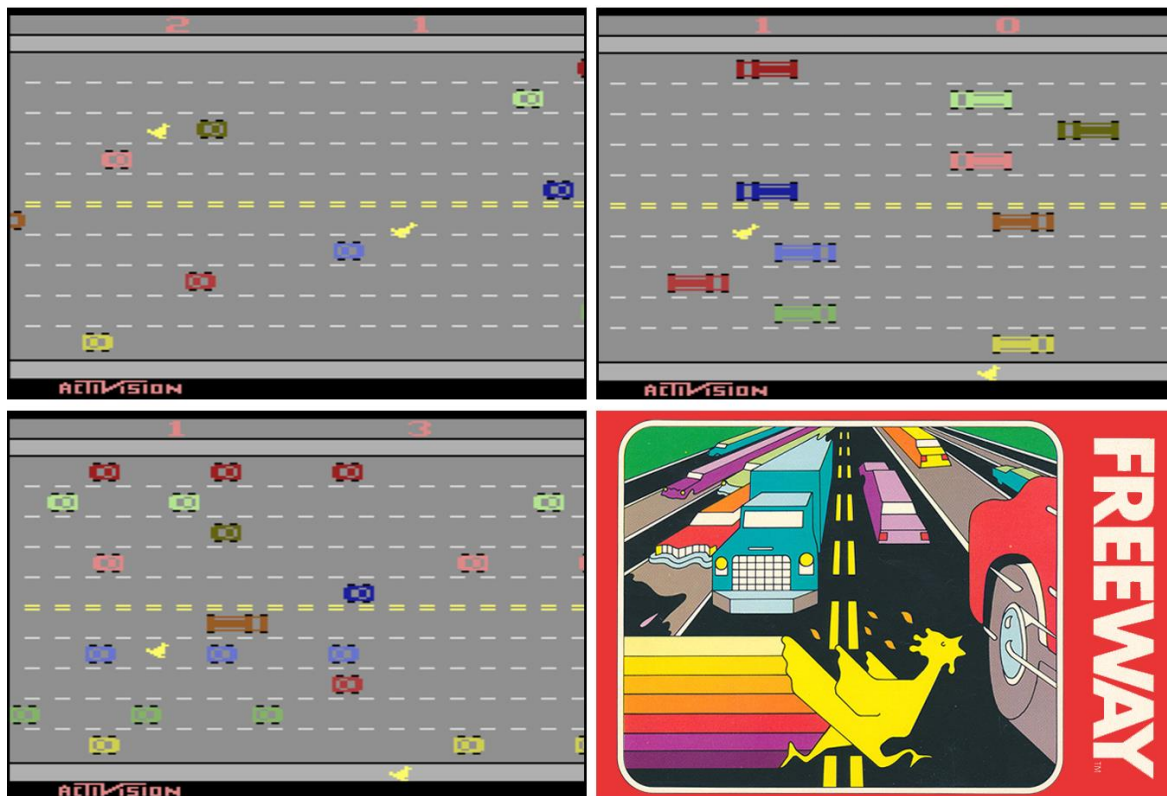


Figure 1. Screenshots from the game (top row and bottom-left). Illustration from the original box of the game (bottom-right). [3]

## 2. Algorithms

### 2.1 REINFORCE

The algorithm was developed by Ronald J. Williams in 1992. It is a basic algorithm. Actions that produced good outcomes will become more probable during learning. On the contrary, actions that produced bad outcomes will become less probable. So, we can say that it „reinforces” the correct actions. The agent will directly learn policies (mapping states to actions). The algorithm can be observed in Figure 2.

```
1: Initialize learning rate  $\alpha$ 
2: Initialize weights  $\theta$  of a policy network  $\pi_\theta$ 
3: for  $episode = 0, \dots, MAX\_EPISODE$  do
4:   Sample a trajectory  $\tau = (s_0, a_0, r_0), \dots, (s_T, a_T, r_T)$ 
5:   Set  $\nabla_\theta J(\pi_\theta) = 0$ 
6:   for  $t = 0, \dots, T$  do
7:      $R_t(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r_{t'}$ 
8:      $\nabla_\theta J(\pi_\theta) = \nabla_\theta J(\pi_\theta) + R_t(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$ 
9:   end for
10:   $\theta \leftarrow \theta + \alpha \nabla_\theta J(\pi_\theta)$ 
11: end for
```

Figure 2. REINFORCE Algorithm [4]

### 2.2 SARSA

It was invented by Rummery and Niranjan in 1994. It is a value-based algorithm used for learning the Q-function and generating actions. The name stands for State Action Reward State Action highlighting the key components involved in updating the Q-values. The algorithm can be observed in Figure 3.

```
1: Initialize learning rate  $\alpha$ 
2: Initialize  $\varepsilon$ 
3: Randomly initialize the network parameters,  $\theta$ 
4: for  $m = 1, \dots, MAX\_STEPS$  do
5:   Gather  $N$  experiences  $(s_i, a_i, r_i, s'_i, a'_i)$  using the current  $\varepsilon$ -greedy policy
6:   for  $i = 1, \dots, N$  do
7:     # Calculate target Q-values for each example
8:      $y_i = r_i + \delta_{s'_i} \gamma Q^{\pi_\theta}(s'_i, a'_i)$  where  $\delta_{s'_i} = 0$  if  $s'_i$  is terminal, 1 otherwise
9:   end for
10:  # Calculate the loss, for example using MSE
11:   $L(\theta) = \frac{1}{N} \sum_i (y_i - Q^{\pi_\theta}(s_i, a_i))^2$ 
12:  # Update the network's parameters
13:   $\theta = \theta - \alpha \nabla_\theta L(\theta)$ 
14:  Decay  $\varepsilon$ 
15: end for
```

Figure 3. SARSA Algorithm with  $\varepsilon$ -greedy policy. [4]

## 2.3 Deep Q-Networks (DQN)

Introduced in 2013 by Mnih et. al. It is similar to SARSA because it is a value-based temporal difference algorithm that learns the Q-function. The difference is that SARSA will approximate the Q-function for the current policy, while DQN will approximate the optimal Q-function, improving the stability and speed of learning. It does this by using not only the action taken in the next state, but all the possible actions from that state, making it an off-policy algorithm. There are several improvements over the “vanilla” version of the algorithm, which will be discussed in the next chapters. The algorithm can be observed in Figure 4.

```

1: Initialize learning rate  $\alpha$ 
2: Initialize  $\tau$ 
3: Initialize number of batches per training step,  $B$ 
4: Initialize number of updates per batch,  $U$ 
5: Initialize batch size,  $N$ 
6: Initialize experience replay memory with max size,  $K$ 
7: Randomly initialize the network parameters,  $\theta$ 
8: for  $m = 1, \dots, MAX\_STEPS$  do
9:   Gather and store  $m$  experiences  $(s_i, a_i, r_i, s'_i, a'_i)$  using the current policy
10:  for  $b = 1, \dots, B$  do
11:    Sample a batch,  $b$ , of experiences from the experience replay memory
12:    for  $u = 1, \dots, U$  do
13:      for  $i = 1, \dots, N$  do
14:        # Calculate target Q-values for each example
15:         $y_i = r_i + \delta_{s'_i} \gamma \max_{a'_i} Q^{\pi_{\theta}}(s'_i, a'_i)$  where  $\delta_{s'_i} = 0$  if  $s'_i$  is terminal, 1 otherwise
16:      end for
17:      # Calculate the loss, for example using MSE
18:       $L(\theta) = \frac{1}{N} \sum_i (y_i - Q^{\pi_{\theta}}(s_i, a_i))^2$ 
19:      # Update the network's parameters
20:       $\theta \leftarrow \theta - \alpha \nabla_{\theta} L(\theta)$ 
21:    end for
22:  end for
23:  Decay  $\tau$ 
24: end for

```

Figure 4. DQN Algorithm with Boltzmann policy. [4]

## 2.4 DQN with target networks

It represents the first improvement of the DQN algorithm. In vanilla DQN, the same Q-network is used for both selecting and evaluating actions, which can lead to instability. To overcome this problem, a separate target network is used. The target network is a copy of the Q-network but is updated less frequently (lagged copy). The target Q-values are then calculated using this target network, which provides more stable targets for the training process. The algorithm is presented in Figure 5.

```

1: Initialize learning rate  $\alpha$ 
2: Initialize  $\tau$ 
3: Initialize number of batches per training step,  $B$ 
4: Initialize number of updates per batch,  $U$ 
5: Initialize batch size,  $N$ 
6: Initialize experience replay memory with max size,  $K$ 
7: Initialize target network update frequency,  $F$ 
8: Randomly initialize the network parameters,  $\theta$ 
9: Initialize the target network parameters  $\phi = \theta$ 
10: for  $m = 1, \dots, MAX\_STEPS$  do
11:   Gather and store  $m$  experiences  $(s_i, a_i, r_i, s'_i, a'_i)$  using the current policy
12:   for  $b = 1, \dots, B$  do
13:     Sample a batch,  $b$  of experiences from the experience replay memory
14:     for  $u = 1, \dots, U$  do
15:       for  $i = 1, \dots, N$  do
16:         # Calculate target Q-values for each example
17:          $y_i = r_i + \delta_{s'_i} \gamma \max_{a'_i} Q^{\pi_\phi}(s'_i, a'_i)$  where  $\delta_{s'_i} = 0$  if  $s'_i$  is terminal, 1 otherwise
18:       end for
19:       # Calculate the loss, for example using MSE
20:        $L(\theta) = \frac{1}{N} \sum_i (y_i - Q^{\pi_\theta}(s_i, a_i))^2$ 
21:       # Update the network's parameters
22:        $\theta \leftarrow \theta - \alpha \nabla_\theta L(\theta)$ 
23:     end for
24:   end for
25:   Decay  $\tau$ 
26:   if  $(m \bmod F) == 0$  then
27:     # Update the target network
28:      $\phi = \theta$ 
29:   end if
30: end for

```

Figure 5. DQN with target networks. [4]

## 2.5 Double DQN (DDQN)

DDQN represents the second improvement of the DQN algorithm. It uses double estimation for computing the Q-function by using a second Q-function trained using different experiences. The algorithm is presented in Figure 6.

$$Q^{\pi}_{tar:DQN}(s, a) = r + \gamma \max_{a'} Q^{\pi_\theta}(s', a')$$

$$Q^{\pi}_{tar:DDQN}(s, a) = r + \gamma \max_{a'} Q^{\pi_\theta}(s', a')$$

Figure 6. DDQN vs DQN Q-function computation. [4]

## 2.6 Prioritized Experience Replay (DQN PER)

It represents the third and last improvement of the DQN algorithm. It was introduced in 2015 by Schaul et. al. It improves upon the standard experience replay mechanism used in DQN. In vanilla DQN, experiences are sampled uniformly from the replay memory. In contrast, PER prioritizes more informative experiences by assigning a priority to each experience. The only difference between DDQN and DQN is the way that the Q-function is computed. The differences can be observed in Figure 7.



```

1: Initialize learning rate  $\alpha$ 
2: Initialize  $\tau$ 
3: Initialize number of batches per training step,  $B$ 
4: Initialize number of updates per batch,  $U$ 
5: Initialize batch size,  $N$ 
6: Initialize experience replay memory with max size,  $K$ 
7: Initialize target network update frequency,  $F$ 
8: Initialize the maximum priority,  $P$ 
9: Initialize  $\varepsilon$ 
10: Initialize the prioritization parameter  $\eta$ 
11: Randomly initialize the network parameters,  $\theta$ 
12: Initialize the target network parameters  $\phi = \theta$ 
13: for  $m = 1, \dots, MAX\_STEPS$  do
14:   Gather and store  $m$  experiences  $(s_i, a_i, r_i, s'_i, p_i)$  where  $p_i = P$ 
15:   for  $b = 1, \dots, B$  do
16:     Sample a prioritized batch,  $b$  of experiences from the experience replay memory
17:     for  $u = 1, \dots, U$  do
18:       for  $i = 1, \dots, N$  do
19:         # Calculate target  $Q$ -values for each example
20:          $y_i = r_i + \delta_{s'_i} \gamma Q^{\pi_\phi}(s'_i, \arg \max_{a'_i} Q^{\pi_\theta}(s'_i, a'_i))$  where  $\delta_{s'_i} = 0$  if  $s'_i$  is terminal, 1 otherwise
21:         # Calculate the absolute TD error for each example
22:          $\omega_i = |y_i - Q^{\pi_\theta}(s_i, a_i)|$ 
23:       end for
24:       # Calculate the loss, for example using MSE
25:        $L(\theta) = \frac{1}{N} \sum_i (y_i - Q^{\pi_\theta}(s_i, a_i))^2$ 
26:       # Update the network's parameters
27:        $\theta \leftarrow \theta - \alpha \nabla_\theta L(\theta)$ 
28:       # Calculate the priorities for each example
29:        $p_i = \frac{(|\omega_i| + \varepsilon)^\eta}{\sum_j (|\omega_j| + \varepsilon)^\eta}$ 
30:       # Update the experience replay memory with the new priorities
31:     end for
32:   end for
33:   Decay  $\tau$ 
34:   if  $(m \bmod F) == 0$  then
35:     # Update the target network
36:      $\phi = \theta$ 
37:   end if
38: end for

```

Figure 7. DQN PER Algorithm. [4]

## 2.7 Advantage Actor-Critic (A2C)

A2C stands for Advantage Actor-Critic. This type of algorithm involves two main components. The first component is an **actor** that learns a policy and tries to output the best action to take in each state. The second component is a **critic** that will learn a value function and will try to provide a baseline for the actor to compare against. The advantage function is used to decide how much better taking a particular action in state is compared to the average action taken in that state. The general A2C algorithm can be observed in Figure 8.

There are two different methods for estimating the advantage function:  $n$ -step returns and Generalized Advantage Estimation (GAE). A2C with  $n$ -step returns will use the returns from  $n$  future steps to update the value estimates and the policy. Instead of waiting until the end of an episode to calculate the total reward A2C updates its estimates after a fixed number of steps. Generalized Advantage Estimation (GAE) represents the second method for estimating the advantage function of an A2C algorithm. It doesn't require choosing a value for  $n$ . Instead, it will use multiple values for  $n$  and will average all the returns.

```

1: Set  $\beta \geq 0$  # entropy regularization weight
2: Set  $\alpha_A \geq 0$  # actor learning rate
3: Set  $\alpha_C \geq 0$  # critic learning rate
4: Randomly initialize the actor and critic parameters  $\theta_A, \theta_C$ 
5: for  $episode = 0, \dots, MAX\_EPISODE$  do
6:   Gather and store data  $(s_t, a_t, r_t, s'_t)$  by acting in the environment using the current policy
7:   for  $t = 0, \dots, T$  do
8:     Calculate predicted  $V$ -value  $\hat{V}^\pi(s_t)$  using the critic network  $\theta_C$ 
9:     Calculate the advantage  $\hat{A}^\pi(s_t, a_t)$  using the critic network  $\theta_C$ 
10:    Calculate  $V_{tar}^\pi(s_t)$  using the critic network  $\theta_C$  and/or trajectory data
11:    Optionally, calculate entropy  $H_t$  of the policy distribution, using the actor network  $\theta_A$ .
    Otherwise set  $\beta = 0$ 
12:  end for
13:  Calculate value loss, for example using MSE:
14:   $L_{val}(\theta_C) = \frac{1}{T} \sum_{t=0}^T (\hat{V}^\pi(s_t) - V_{tar}^\pi(s_t))^2$ 
15:  Calculate policy loss:
16:   $L_{pol}(\theta_A) = \frac{1}{T} \sum_{t=0}^T (-\hat{A}^\pi(s_t, a_t) \log \pi_{\theta_A}(a_t|s_t) - \beta H_t)$ 
17:  Update critic parameters, for example using SGD:
18:   $\theta_C \leftarrow \theta_C - \alpha_C \nabla_{\theta_C} L_{val}(\theta_C)$ 
19:  Update actor parameters, for example using SGD:
20:   $\theta_A \leftarrow \theta_A - \alpha_A \nabla_{\theta_A} L_{pol}(\theta_A)$ 
21: end for

```

Figure 8. A2C Algorithm. [4]

## 2.8 Proximal Policy Optimization (PPO)

PPO was introduced by Schulman et al. in 2017. It improves policy gradient methods by introducing a surrogate objective function, which ensures stable learning by using monotonic policy updates. It also uses off-policy data to be reused (on-policy algorithms cannot reuse data). PPO's balance between simplicity and performance has made it a popular choice for a wide range of reinforcement learning tasks. It has two variants:

1. PPO with adaptive KL penalty: uses an objective function with a dynamically adjusted penalty based on KL divergence to control policy updates.
2. PPO with clipped surrogate objective: uses a clipped objective function to bound the updates directly. Outperformed the KL penalty version in the original paper.

Figure 9 shows the PPO algorithm with clipped surrogate objective.



```

1: Set  $\beta \geq 0$  # entropy regularization weight
2: Set  $\epsilon \geq 0$  # the clipping variable
3: Set  $K$  the number of epochs
4: Set  $N$  the number of actors
5: Set  $T$  the time horizon
6: Set  $M \leq NT$  the minibatch size
7: Set  $\alpha_A \geq 0$  # actor learning rate
8: Set  $\alpha_C \geq 0$  # critic learning rate
9: Randomly initialize the actor and critic parameters  $\theta_A, \theta_C$ 
10: Initialize the "old" actor network  $\theta_{Aold}$ 
11: for  $i = 1, 2, \dots$  do
12:   Set  $\theta_{Aold} = \theta_A$ 
13:   for  $actor = 1, 2, \dots, N$  do
14:     Run policy  $\theta_{Aold}$  in environment for  $T$  time steps and collect the trajectories
15:     Compute advantages  $A_1, \dots, A_T$  using  $\theta_{Aold}$ 
16:     Calculate  $V_{tar,1}^\pi, \dots, V_{tar,T}^\pi$  using the critic network  $\theta$  and/or trajectory data
17:   end for
18:   Let  $batch$  with size  $NT$  consist of the collected trajectories, advantages, and target  $V$ -values
19:   for  $epoch = 1, 2, \dots, K$  do
20:     for minibatch  $m$  in  $batch$  do
21:       The following are computed over the whole minibatch  $m$ :
22:       Calculate  $r_m(\theta_A)$ 
23:       Calculate  $J_m^{CLIP}(\theta_A)$  using the advantages  $A_m$  from the minibatch and  $r_m(\theta_A)$ 
24:       Calculate entropies  $H_m$  using the actor network  $\theta_A$ 
25:       Calculate policy loss:
26:        $L_{pol} = -J_m^{CLIP}(\theta_A) - \beta H_m$ 
27:       Calculate predicted  $V$ -value  $\hat{V}^\pi(s_m)$  using the critic network  $\theta_C$ 
28:       Calculate value loss using the  $V$ -targets from the minibatch:
29:        $L_{val}(\theta_C) = \text{MSE}(\hat{V}^\pi(s_m), V_{tar}^\pi(s_m))$ 
30:       Update actor parameters, for example using SGD:
31:        $\theta_A \leftarrow \theta_A - \alpha_A \nabla_{\theta_A} L_{pol}(\theta_A)$ 
32:       Update critic parameters, for example using SGD:
33:        $\theta_C \leftarrow \theta_C - \alpha_C \nabla_{\theta_C} L_{val}(\theta_C)$ 
34:     end for
35:   end for
36: end for

```

Figure 9. PPO with clipped surrogate objective. [4]

## 2.9 Asynchronous Advantage Actor-Critic (A3C)

It was introduced by Mnih et. al. and uses the A2C algorithm but speeds up the process by training agents in parallel. This leads to faster and more efficient learning. With asynchronous parallelization, the global network will update the parameters whenever it receives data from any worker and each worker will update itself from the global network. The A3C algorithm can be observed in Figure 10.

```

1: Set learning rate  $\alpha$ 
2: Initialize the global network,  $\theta_G$ 
3: Initialize  $N$  worker networks,  $\theta_{W,1}, \theta_{W,2}, \dots, \theta_{W,N}$ 
4: for each worker asynchronously do
5:   Pull parameters from the global network and set  $\theta_{W,i} \leftarrow \theta_G$ 
6:   Collect trajectories
7:   Compute gradient  $\nabla_{\theta_{W,i}}$ 
8:   Push  $\nabla_{\theta_{W,i}}$  to the global network
9: end for
10: On receiving worker gradient, update the global network  $\theta_G \leftarrow \theta_G - \alpha \nabla_{\theta_{W,i}}$ 

```

Figure 10. A3C Algorithm. [4]

### 3. Experimental Results

Every algorithm was trained for 500k steps, using the modular deep reinforcement learning framework in PyTorch called SLM Lab [5]. I will provide the full configuration file for each algorithm, and I will describe each file. Also, there will be graphs of the results, representing the mean returns and the moving average of the mean returns. However, some settings are the same for every file:

- **Network architecture:** Three layers of Convolutional Neural Network, followed by one fully connected layer, with ReLU activation function.
- **Environment:** OpenAI Gym's *Freeway-v0*
- **Training length:** 500,000 time steps

#### 3.1 REINFORCE

The graphs can be observed in Figure 11 and the configuration file can be observed in Figure 12.

- **Algorithm:** The algorithm is Reinforce (line:6). Entropy with a linear decay schedule (line:13).
- **Optimizer:** Adam (line:43), with a learning rate of 0.002 (line:44). The learning rate is constant because the learning rate scheduler is set to null (line:46).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:22). Training frequency is set to 1 (line:19), meaning that the network will be trained every 1 episode.
- **Evaluation:** The agent is evaluated every 10,000 steps (line:65).

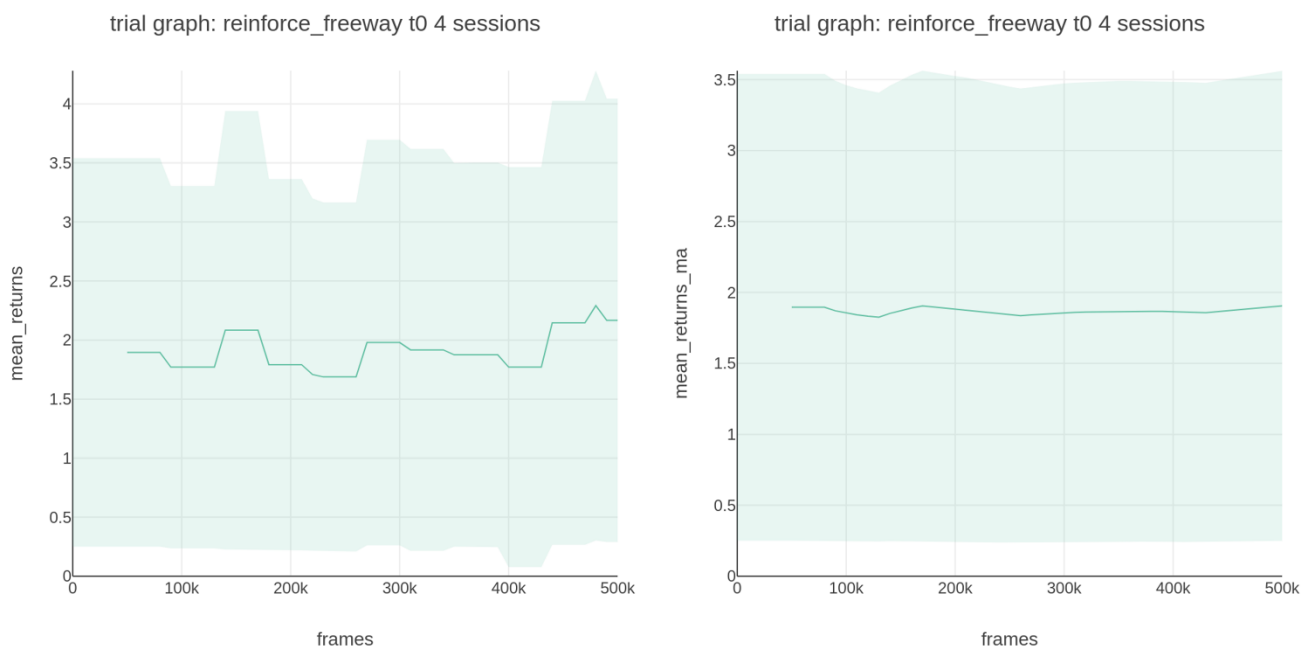


Figure 11. REINFORCE graphs.

```

1  {
2    "reinforce_freeway": {
3      "agent": [{
4        "name": "Reinforce",
5        "algorithm": {
6          "name": "Reinforce",
7          "action_pdtype": "default",
8          "action_policy": "default",
9          "center_return": true,
10         "explore_var_spec": null,
11         "gamma": 0.99,
12         "entropy_coef_spec": {
13           "name": "linear_decay",
14           "start_val": 0.01,
15           "end_val": 0.001,
16           "start_step": 1000,
17           "end_step": 100000
18         },
19         "training_frequency": 1
20       },
21       "memory": {
22         "name": "OnPolicyReplay",
23         "batch_size": 64
24       },
25       "net": {
26         "type": "ConvNet",
27         "shared": true,
28         "conv_hid_layers": [
29           [32, 8, 4, 0, 1],
30           [64, 4, 2, 0, 1],
31           [32, 3, 1, 0, 1]
32         ],
33         "fc_hid_layers": [512],
34         "hid_layers_activation": "relu",
35         "init_fn": "orthogonal_",
36         "normalize": true,
37         "batch_norm": false,
38         "clip_grad_val": 0.5,
39         "loss_spec": {
40           "name": "MSELoss"
41         },
42         "optim_spec": {
43           "name": "Adam",
44           "lr": 0.002
45         },
46         "lr_scheduler_spec": null,
47         "gpu": true
48       }
49     ],
50     "env": [{
51       "name": "Freeway-v0",
52       "frame_op": "concat",
53       "frame_op_len": 4,
54       "reward_scale": "sign",
55       "num_envs": 16,
56       "max_t": null,
57       "max_frame": 500000
58     }],
59     "body": {
60       "product": "outer",
61       "num": 1
62     },
63     "meta": {
64       "distributed": false,
65       "eval_frequency": 10000,
66       "log_frequency": 10000,
67       "max_session": 4,
68       "max_trial": 1
69     }
70   }
71 }

```

Figure 12. REINFORCE config file.

### 3.2 SARSA

The graphs can be observed in Figure 13 and the configuration file can be observed in Figure 14.

- **Algorithm:** The algorithm is SARSA (line:6). Entropy with a linear decay schedule (line:13).
- **Optimizer:** Adam (line:41), with a learning rate of 0.001 (line:42). The learning rate is constant because the learning rate scheduler is set to null (line:44).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:20). Training frequency is set to 1 (line:19), meaning that the network will be trained every 1 episode.
- **Evaluation:** The agent is evaluated every 10,000 steps (line:63).

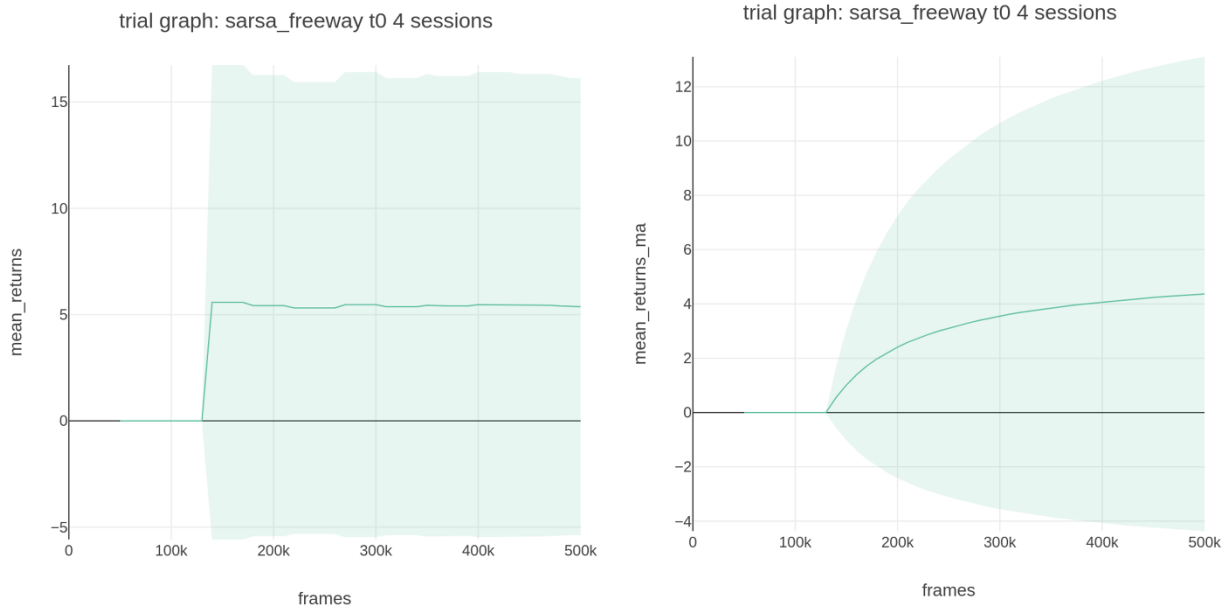


Figure 13. SARSA graphs.

```

1  {
2    "sarsa_freeway": {
3      "agent": [{
4        "name": "SARSA",
5        "algorithm": {
6          "name": "SARSA",
7          "action_pdtype": "Argmax",
8          "action_policy": "epsilon_greedy",
9          "entropy_coef_spec": {
10           "name": "linear_decay",
11           "start_val": 1.0,
12           "end_val": 0.05,
13           "start_step": 0,
14           "end_step": 10000
15         },
16         "gamma": 0.99,
17         "training_frequency": 32
18       }],
19      "memory": {
20        "name": "OnPolicyBatchReplay",
21        "batch_size": 64
22      },
23      "net": {
24        "type": "ConvNet",
25        "shared": true,
26        "conv_hid_layers": [
27          [32, 8, 4, 0, 1],
28          [64, 4, 2, 0, 1],
29          [32, 3, 1, 0, 1]
30        ],
31        "fc_hid_layers": [512],
32        "hid_layers_activation": "relu",
33        "init_fn": "orthogonal_",
34        "normalize": true,
35        "batch_norm": false,
36        "clip_grad_val": 0.5,
37        "loss_spec": {
38          "name": "MSELoss"
39        },
40        "optim_spec": {
41          "name": "Adam",
42          "lr": 0.001
43        },
44        "lr_scheduler_spec": null,
45        "gpu": true
46      }
47    },
48    "env": [{
49      "name": "Freeway-v0",
50      "frame_op": "concat",
51      "frame_op_len": 4,
52      "reward_scale": "sign",
53      "num_envs": 16,
54      "max_t": null,
55      "max_frame": 500000
56    }],
57    "body": {
58      "product": "outer",
59      "num": 1
60    },
61    "meta": {
62      "distributed": false,
63      "eval_frequency": 10000,
64      "log_frequency": 10000,
65      "max_session": 4,
66      "max_trial": 1
67    }
68  }
69

```

Figure 14. SARSA config file.

### 3.3 DQN

The graphs can be observed in Figure 15 and the configuration file can be observed in Figure 16.

- **Algorithm:** The algorithm is VanillaDQN (line:7) (different from DQN, which uses target nets). The action policy is Boltzmann policy (line:9). Entropy with a linear decay schedule (line:11).
- **Optimizer:** Adam (line:47), with a learning rate of 0.001 (line:48). The learning rate is constant because the learning rate scheduler is set to null (line:50).
- **Training frequency:** Training starts after the agent has made 32 steps in the environment (line:21) and occurs every step from then on (line:20). *Replay* memory is used for sampling batches (line:24).
- **Evaluation:** The agent is evaluated every 10,000 steps (line:70).

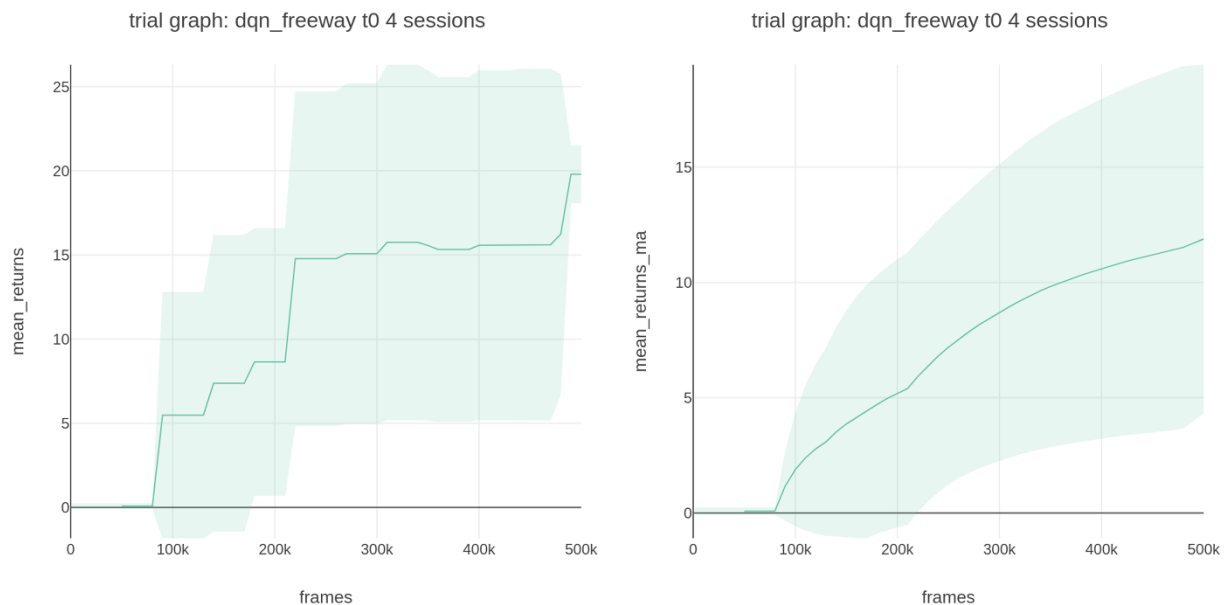


Figure 15. DQN graphs.

### 3.4 DQN with target networks

The graphs can be observed in Figure 17 and the configuration file can be observed in Figure 18.

- **Algorithm:** The algorithm is DQN (line:7). The action policy is Boltzmann policy (line:9). Entropy with a linear decay schedule (line:11).
- **Optimizer:** Adam (line:47), with a learning rate of 0.001 (line:48). The learning rate is constant because the learning rate scheduler is set to null (line:50).
- **Training frequency:** Training starts after the agent has made 32 steps in the environment (line:21) and occurs every step from then on (line:20). *Replay* memory is used for sampling batches (line:24).
- **Evaluation:** The agent is evaluated every 10,000 steps (line:70).



```

1  {
2    "dqn_freeway": {
3      "agent": [
4        {
5          "name": "VanillaDQN",
6          "algorithm": {
7            "name": "VanillaDQN",
8            "action_pdtype": "Argmax",
9            "action_policy": "boltzmann",
10           "explore_var_spec": {
11             "name": "linear_decay",
12             "start_val": 5.0,
13             "end_val": 0.5,
14             "start_step": 0,
15             "end_step": 1000
16           },
17           "gamma": 0.99,
18           "training_batch_iter": 3,
19           "training_iter": 3,
20           "training_frequency": 1,
21           "training_start_step": 32
22         },
23         "memory": {
24           "name": "Replay",
25           "batch_size": 32,
26           "max_size": 10000,
27           "use_cer": false
28         },
29         "net": {
30           "type": "ConvNet",
31           "shared": true,
32           "conv_hid_layers": [
33             [32, 8, 4, 0, 1],
34             [64, 4, 2, 0, 1],
35             [32, 3, 1, 0, 1]
36           ],
37           "fc_hid_layers": [512],
38           "hid_layers_activation": "relu",
39           "init_fn": "orthogonal_",
40           "normalize": true,
41           "batch_norm": false,
42           "clip_grad_val": 0.5,
43           "loss_spec": {
44             "name": "SmoothL1Loss"
45           },
46           "optim_spec": {
47             "name": "Adam",
48             "lr": 0.001
49           },
50           "lr_scheduler_spec": null,
51           "gpu": true
52         }
53       ],
54     },
55     "env": [{
56       "name": "Freeway-v0",
57       "frame_op": "concat",
58       "frame_op_len": 4,
59       "reward_scale": "sign",
60       "num_envs": 16,
61       "max_t": null,
62       "max_frame": 500000
63     }],
64     "body": {
65       "product": "outer",
66       "num": 1
67     },
68     "meta": {
69       "distributed": false,
70       "eval_frequency": 10000,
71       "log_frequency": 10000,
72       "max_session": 4,
73       "max_trial": 1
74     }
75   }
76 }

```

Figure 16. DQN config file.

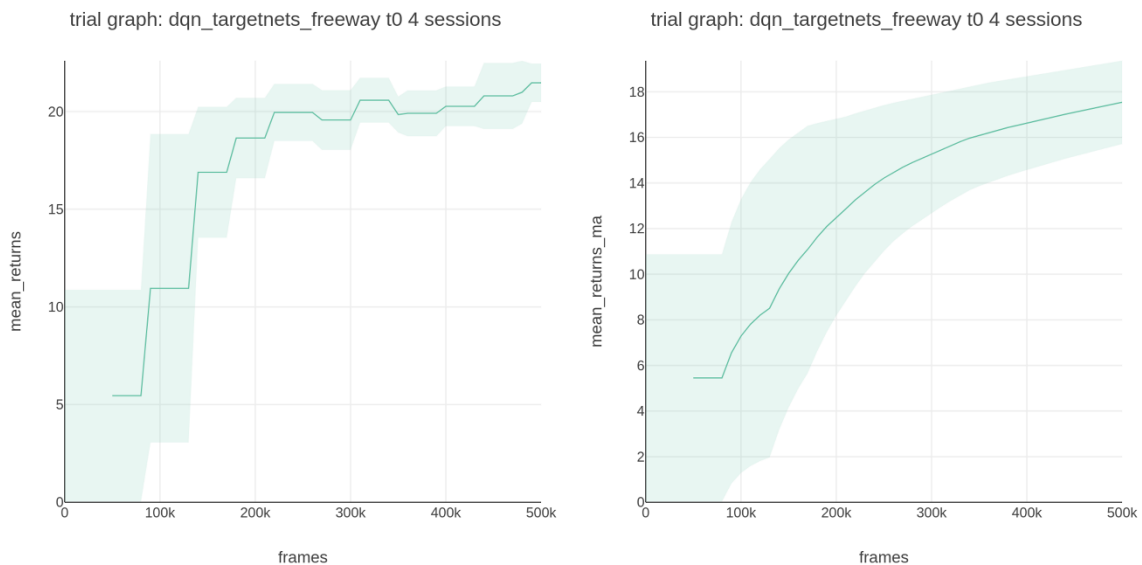


Figure 17. DQN with target nets graphs.



```

1  {
2      "dqn_targetnets_freeway": {
3          "agent": [
4              {
5                  "name": "DQN",
6                  "algorithm": {
7                      "name": "DQN",
8                      "action_pdtype": "Argmax",
9                      "action_policy": "boltzmann",
10                     "explore_var_spec": {
11                         "name": "linear_decay",
12                         "start_val": 5.0,
13                         "end_val": 0.5,
14                         "start_step": 0,
15                         "end_step": 1000
16                     },
17                     "gamma": 0.99,
18                     "training_batch_iter": 3,
19                     "training_iter": 3,
20                     "training_frequency": 1,
21                     "training_start_step": 32
22                 },
23                 "memory": {
24                     "name": "Replay",
25                     "batch_size": 32,
26                     "max_size": 10000,
27                     "use_cer": false
28                 },
29                 "net": {
30                     "type": "ConvNet",
31                     "shared": true,
32                     "conv_hid_layers": [
33                         [32, 8, 4, 0, 1],
34                         [64, 4, 2, 0, 1],
35                         [32, 3, 1, 0, 1]
36                     ],
37                     "fc_hid_layers": [512],
38                     "hid_layers_activation": "relu",
39                     "init_fn": "orthogonal_",
40                     "normalize": true,
41                     "batch_norm": false,
42                     "clip_grad_val": 0.5,
43                     "loss_spec": {
44                         "name": "SmoothL1Loss"
45                     },
46                     "optim_spec": {
47                         "name": "Adam",
48                         "lr": 0.001
49                     },
50                     "lr_scheduler_spec": null,
51                     "gpu": true
52                 }
53             ]
54         },
55         "env": [{
56             "name": "Freeway-v0",
57             "frame_op": "concat",
58             "frame_op_len": 4,
59             "reward_scale": "sign",
60             "num_envs": 16,
61             "max_t": null,
62             "max_frame": 500000
63         }],
64         "body": {
65             "product": "outer",
66             "num": 1
67         },
68         "meta": {
69             "distributed": false,
70             "eval_frequency": 10000,
71             "log_frequency": 10000,
72             "max_session": 4,
73             "max_trial": 1
74         }
75     }
76 }

```

Figure 18. DQN with target nets config file.

### 3.5 DDQN

The graphs can be observed in Figure 19 and the configuration file can be observed in Figure 20.

- **Algorithm:** The algorithm is Double DQN (line:7). The action policy is Boltzmann policy (line:9). Entropy with a linear decay schedule (line:11).
- **Optimizer:** Adam (line:47), with a learning rate of 0.001 (line:48). The learning rate is constant because the learning rate scheduler is set to null (line:50).
- **Training frequency:** Training starts after the agent has made 32 steps in the environment (line:21) and occurs every step from then on (line:20). *Replay* memory is used for sampling batches (line:24).
- **Evaluation:** The agent is evaluated every 10,000 steps (line:70).

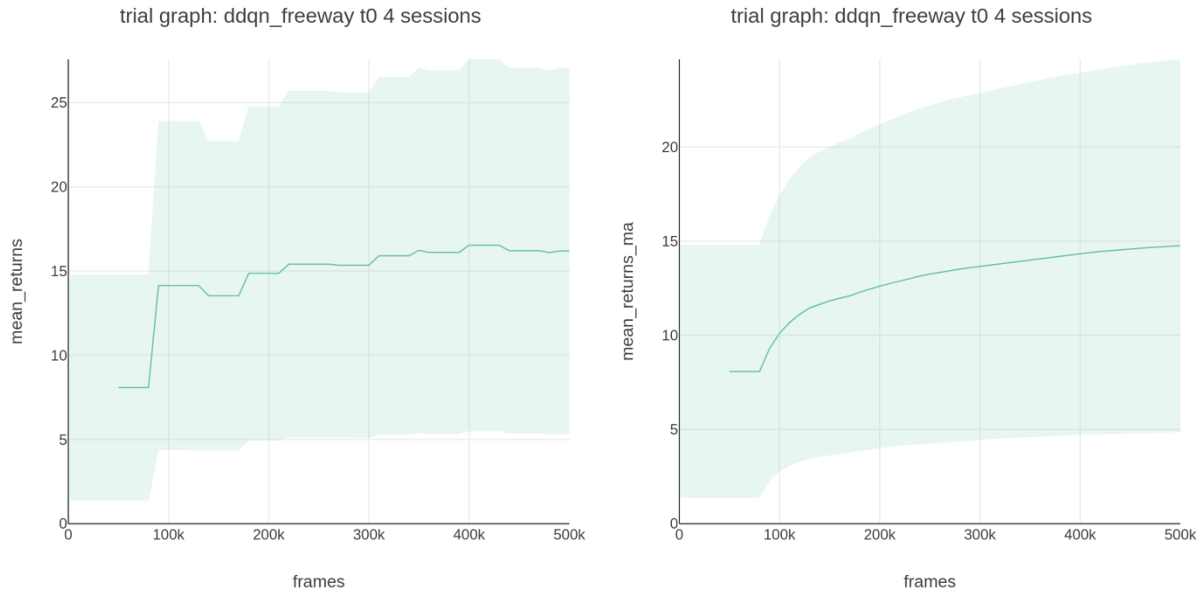


Figure 19. DDQN graphs.

```

1  {
2    "ddqn_freeway": {
3      "agent": [
4        {
5          "name": "DoubleDQN",
6          "algorithm": {
7            "name": "DoubleDQN",
8            "action_pdtype": "Argmax",
9            "action_policy": "boltzmann",
10           "explore_var_spec": {
11             "name": "linear_decay",
12             "start_val": 5.0,
13             "end_val": 0.5,
14             "start_step": 0,
15             "end_step": 1000
16           },
17           "gamma": 0.99,
18           "training_batch_iter": 3,
19           "training_iter": 3,
20           "training_frequency": 1,
21           "training_start_step": 32
22         },
23         "memory": {
24           "name": "Replay",
25           "batch_size": 32,
26           "max_size": 10000,
27           "use_cer": false
28         },
29         "net": {
30           "type": "ConvNet",
31           "shared": true,
32           "conv_hid_layers": [
33             [32, 8, 4, 0, 1],
34             [64, 4, 2, 0, 1],
35             [32, 3, 1, 0, 1]
36           ],
37           "fc_hid_layers": [512],
38           "hid_layers_activation": "relu",
39           "init_fn": "orthogonal_",
40           "normalize": true,
41           "batch_norm": false,
42           "clip_grad_val": 0.5,
43           "loss_spec": {
44             "name": "SmoothL1Loss"
45           },
46           "optim_spec": {
47             "name": "Adam",
48             "lr": 0.001
49           },
50           "lr_scheduler_spec": null,
51           "gpu": true
52         }
53       ],
54     },
55     "env": [{
56       "name": "Freeway-v0",
57       "frame_op": "concat",
58       "frame_op_len": 4,
59       "reward_scale": "sign",
60       "num_envs": 16,
61       "max_t": null,
62       "max_frame": 500000
63     }],
64     "body": {
65       "product": "outer",
66       "num": 1
67     },
68     "meta": {
69       "distributed": false,
70       "eval_frequency": 10000,
71       "log_frequency": 10000,
72       "max_session": 4,
73       "max_trial": 1
74     }
75   }
76 }

```

Figure 20. DDQN config file.

### 3.6 DQN PER

The graphs can be observed in Figure 21 and the configuration file can be observed in Figure 22.

- **Algorithm:** The algorithm is DQN PER (line:7). The action policy is Boltzmann policy (line:9). Entropy with a linear decay schedule (line:11).
- **Optimizer:** Adam (line:49), with a learning rate of 0.001 (line:50). The learning rate is constant because the learning rate scheduler is set to null (line:52).
- **Training frequency:** *PrioritizedExperienceReplay* memory is used for sampling batches (line:24).
- **Evaluation:** The agent is evaluated every 10,000 steps (line:72).

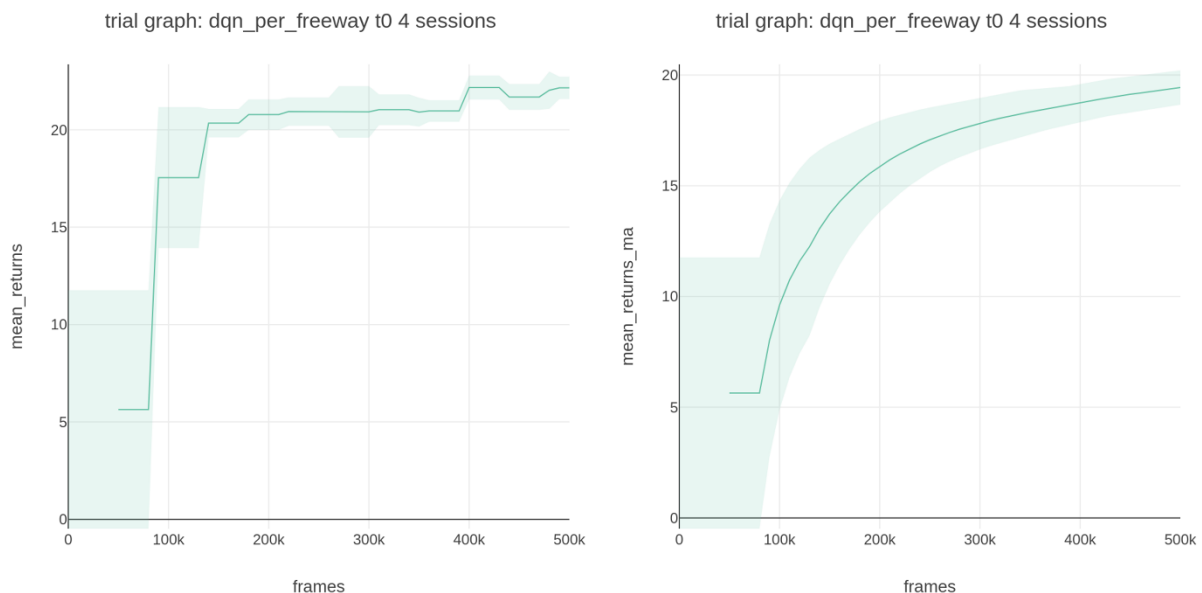


Figure 21. DQN PER graphs.

### 3.7 A2C n-step returns

The graphs can be observed in Figure 23 and the configuration file can be observed in Figure 24.

- **Algorithm:** The algorithm is A2C (line:6) with n-step returns (line:12). Entropy with no decay (line:14).
- **Optimizer:** Same settings for both actor and critic. RMSprop (line:45, line:51), with a learning rate of 0.0007 (line:46, line:52). The learning rate is constant because the learning rate scheduler is set to null (line:56).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:24). Training frequency is set to 5 (line:19), meaning that the network will be trained every 5 episodes.
- **Evaluation:** The agent is evaluated every 10,000 steps (line:75).

**Observation:** I tried various settings in the configuration file but none of them lead to the agent actually learning. The total reward was always 0. Fortunately, I had more luck with A2C GAE.

```

1  {
2    "dqn_per_freeway": {
3      "agent": [
4        {
5          "name": "DQN",
6          "algorithm": {
7            "name": "DQN",
8            "action_ptdtype": "Argmax",
9            "action_policy": "boltzmann",
10           "explore_var_spec": {
11             "name": "linear_decay",
12             "start_val": 5.0,
13             "end_val": 0.5,
14             "start_step": 0,
15             "end_step": 1000
16           },
17           "gamma": 0.99,
18           "training_batch_iter": 3,
19           "training_iter": 3,
20           "training_frequency": 1,
21           "training_start_step": 32
22         },
23         "memory": {
24           "name": "PrioritizedReplay",
25           "alpha": 0.6,
26           "epsilon": 0.0001,
27           "batch_size": 32,
28           "max_size": 10000,
29           "use_cer": false
30         },
31         "net": {
32           "type": "ConvNet",
33           "shared": true,
34           "conv_hid_layers": [
35             [32, 8, 4, 0, 1],
36             [64, 4, 2, 0, 1],
37             [32, 3, 1, 0, 1]
38           ],
39           "fc_hid_layers": [512],
40           "hid_layers_activation": "relu",
41           "init_fn": "orthogonal_",
42           "normalize": true,
43           "batch_norm": false,
44           "clip_grad_val": 0.5,
45           "loss_spec": {
46             "name": "SmoothL1Loss"
47           },
48           "optim_spec": {
49             "name": "Adam",
50             "lr": 0.001
51           },
52           "lr_scheduler_spec": null,
53           "gpu": true
54         }
55       ],
56     },
57     "env": [{
58       "name": "Freeway-v0",
59       "frame_op": "concat",
60       "frame_op_len": 4,
61       "reward_scale": "sign",
62       "num_envs": 16,
63       "max_t": null,
64       "max_frame": 500000
65     }],
66     "body": {
67       "product": "outer",
68       "num": 1
69     },
70     "meta": {
71       "distributed": false,
72       "eval_frequency": 10000,
73       "log_frequency": 10000,
74       "max_session": 4,
75       "max_trial": 1
76     }
77   }
78 }

```

Figure 22. DQN PER config file.

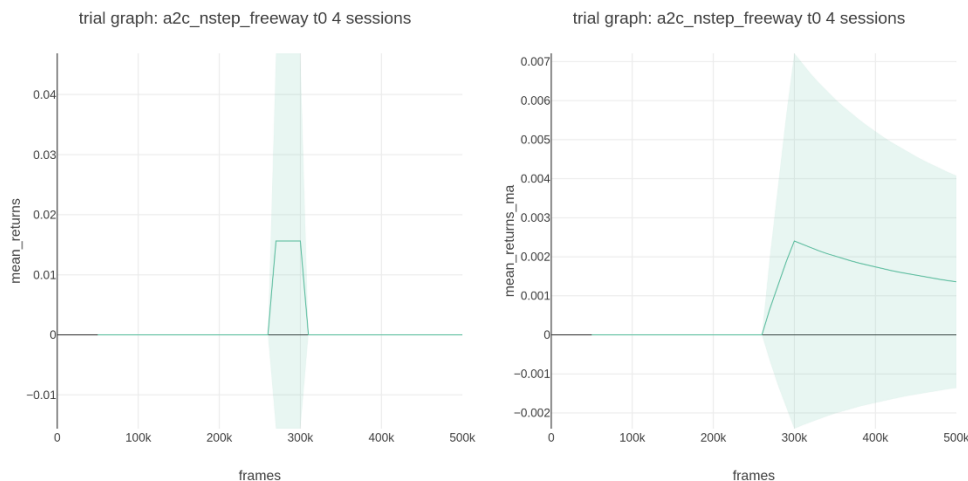


Figure 23. A2C n-step returns graphs.



```

1  {
2      "a2c_nstep_freeway": {
3          "agent": [{
4              "name": "A2C",
5              "algorithm": {
6                  "name": "ActorCritic",
7                  "action_pdtype": "default",
8                  "action_policy": "default",
9                  "explore_var_spec": null,
10                 "gamma": 0.99,
11                 "lam": null,
12                 "num_step_returns": 11,
13                 "entropy_coef_spec": {
14                     "name": "no_decay",
15                     "start_val": 0.01,
16                     "end_val": 0.01,
17                     "start_step": 0,
18                     "end_step": 0
19                 },
20                 "val_loss_coef": 0.5,
21                 "training_frequency": 5
22             },
23             "memory": {
24                 "name": "OnPolicyBatchReplay"
25             },
26             "net": {
27                 "type": "ConvNet",
28                 "shared": true,
29                 "conv_hid_layers": [
30                     [32, 8, 4, 0, 1],
31                     [64, 4, 2, 0, 1],
32                     [32, 3, 1, 0, 1]
33                 ],
34                 "fc_hid_layers": [512],
35                 "hid_layers_activation": "relu",
36                 "init_fn": "orthogonal_",
37                 "normalize": true,
38                 "batch_norm": false,
39                 "clip_grad_val": 0.5,
40                 "use_same_optim": false,
41                 "loss_spec": {
42                     "name": "MSELoss"
43                 },
44                 "actor_optim_spec": {
45                     "name": "RMSprop",
46                     "lr": 7e-4,
47                     "alpha": 0.99,
48                     "eps": 1e-5
49                 },
50                 "critic_optim_spec": {
51                     "name": "RMSprop",
52                     "lr": 7e-4,
53                     "alpha": 0.99,
54                     "eps": 1e-5
55                 },
56                 "lr_scheduler_spec": null,
57                 "gpu": true
58             }
59         ]},
60         "env": [{
61             "name": "Freeway-v0",
62             "frame_op": "concat",
63             "frame_op_len": 4,
64             "reward_scale": "sign",
65             "num_envs": 16,
66             "max_t": null,
67             "max_frame": 500000
68         ]},
69         "body": {
70             "product": "outer",
71             "num": 1
72         },
73         "meta": {
74             "distributed": false,
75             "log_frequency": 10000,
76             "eval_frequency": 10000,
77             "max_session": 4,
78             "max_trial": 1
79         }
80     }
81 }

```

Figure 24. A2C n-step returns config file.

### 3.8 A2C GAE

The graphs can be observed in Figure 25 and the configuration file can be observed in Figure 26.

- **Algorithm:** The algorithm is ActorCritic (line:6) with Generalized Advantage Estimation (line:11). Entropy with no decay (line:14).
- **Optimizer:** Same settings for both actor and critic. Adam (line:45, line:49), with a learning rate of 0.002 (line:46, line:50). The learning rate is constant because the learning rate scheduler is set to null (line:52).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:24). Training frequency is set to 32 (line:21), meaning that the network will be trained every 32 episodes.
- **Evaluation:** The agent is evaluated every 10,000 steps (line:71).

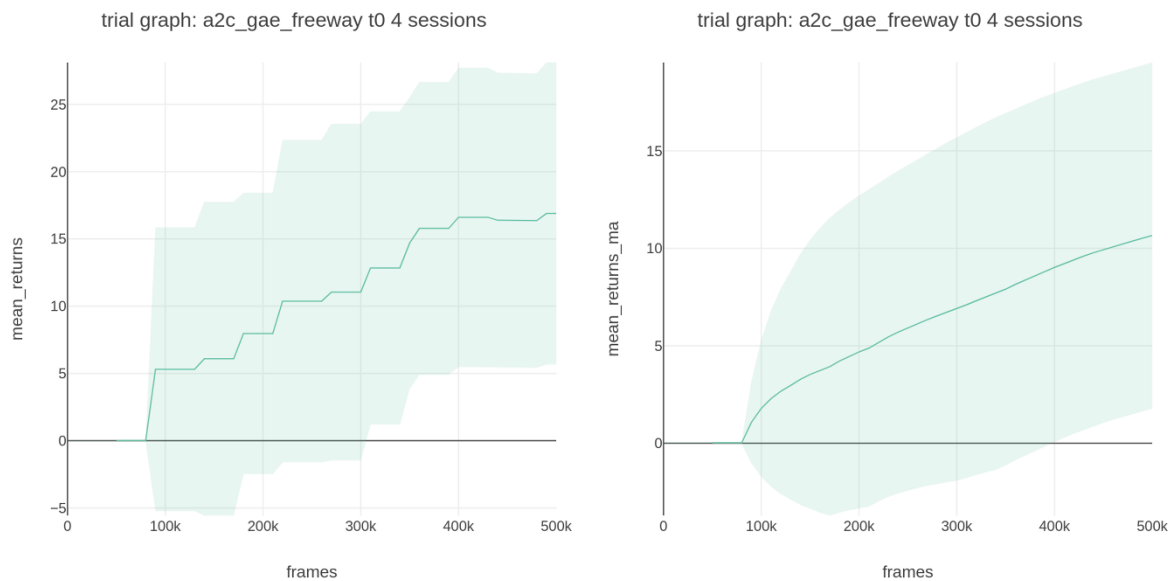


Figure 25. A2C GAE graphs.

```

1  {
2      "a2c_gae_freeway": {
3          "agent": [{
4              "name": "A2C",
5              "algorithm": {
6                  "name": "ActorCritic",
7                  "action_pdtype": "default",
8                  "action_policy": "default",
9                  "explore_var_spec": null,
10                 "gamma": 0.99,
11                 "lam": 0.95,
12                 "num_step_returns": null,
13                 "entropy_coef_spec": {
14                     "name": "no_decay",
15                     "start_val": 0.01,
16                     "end_val": 0.01,
17                     "start_step": 0,
18                     "end_step": 0
19                 },
20                 "val_loss_coef": 0.5,
21                 "training_frequency": 32
22             },
23             "memory": {
24                 "name": "OnPolicyBatchReplay"
25             },
26             "net": {
27                 "type": "ConvNet",
28                 "shared": true,
29                 "conv_hid_layers": [
30                     [32, 8, 4, 0, 1],
31                     [64, 4, 2, 0, 1],
32                     [32, 3, 1, 0, 1]
33                 ],
34                 "fc_hid_layers": [512],
35                 "hid_layers_activation": "relu",
36                 "init_fn": "orthogonal_",
37                 "normalize": true,
38                 "batch_norm": false,
39                 "clip_grad_val": 0.5,
40                 "use_same_optim": false,
41                 "loss_spec": {
42                     "name": "MSELoss"
43                 },
44                 "actor_optim_spec": {
45                     "name": "Adam",
46                     "lr": 0.002
47                 },
48                 "critic_optim_spec": {
49                     "name": "Adam",
50                     "lr": 0.002
51                 },
52                 "lr_scheduler_spec": null,
53                 "gpu": true
54             },
55         }],
56         "env": [{
57             "name": "Freeway-v0",
58             "frame_op": "concat",
59             "frame_op_len": 4,
60             "reward_scale": "sign",
61             "num_envs": 16,
62             "max_t": null,
63             "max_frame": 500000
64         }],
65         "body": {
66             "product": "outer",
67             "num": 1
68         },
69         "meta": {
70             "distributed": false,
71             "log_frequency": 10000,
72             "eval_frequency": 10000,
73             "max_session": 4,
74             "max_trial": 1
75         }
76     }
77 }

```

Figure 26. A2C GAE config file.



### 3.9 PPO

The graphs can be observed in Figure 27 and the configuration file can be observed in Figure 28.

- **Algorithm:** The algorithm is PPO (line:6) and the action policy is the default policy (line:8). GAE is used to estimate advantages (line:11). Entropy with no decay (line:13).
- **Optimizer:** Same settings for both actor and critic. Adam (line:53, line:57), with a learning rate of 0.001 (line:54, line:58). The learning rate is constant because the learning rate scheduler is set to null (line:60).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:24).
- **Evaluation:** The agent is evaluated every 10,000 steps (line:80).

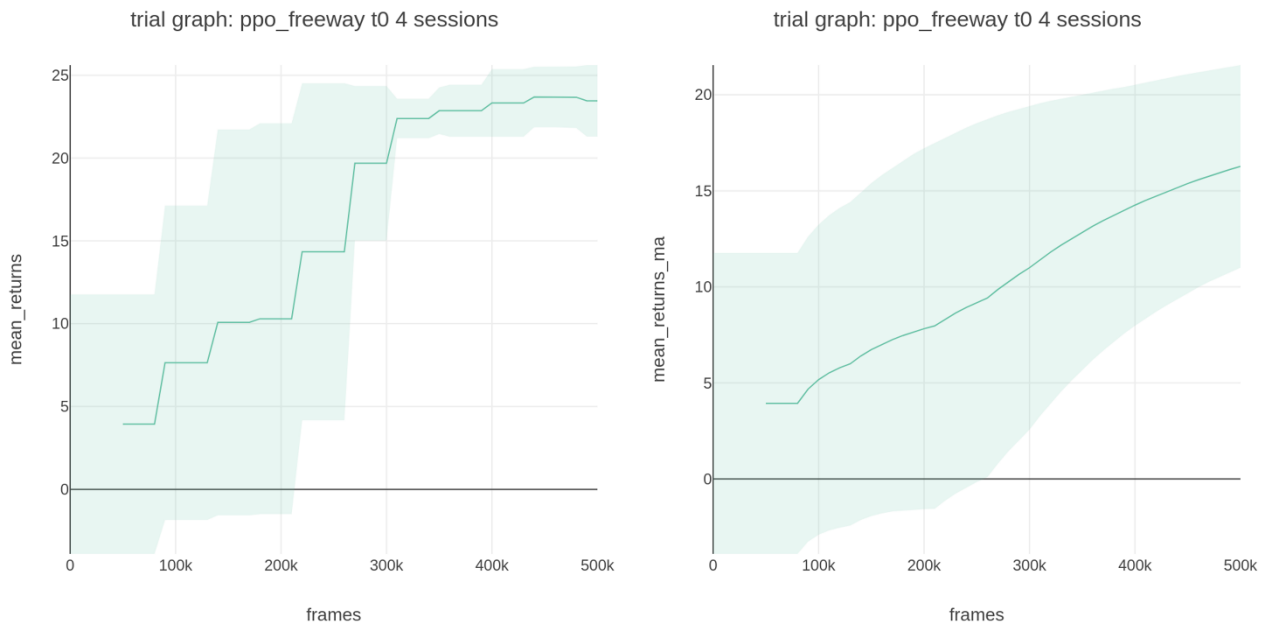


Figure 27. PPO graphs.

### 3.10 A3C

It was trained two times, both with 8 and 16 max sessions (line:73). The graphs for the 8 sessions can be observed in Figure 29 and the graphs for the 16 sessions can be observed in Figure 30. The configuration file can be observed in Figure 31.

- **Algorithm:** The algorithm is ActorCritic (line:6) with Generalized Advantage Estimation (line:11). Entropy with no decay (line:14).
- **Optimizer:** Same settings for both actor and critic. GlobalAdam (line:45, line:49), with a learning rate of 0.007 (line:46, line:50). The learning rate is constant because the learning rate scheduler is set to null (line:52).
- **Training frequency:** Training is episodic, having *OnPolicyReplay* memory (line:24). Training frequency is set to 32 (line:21), meaning that the network will be trained every 32 episodes.
- **Evaluation:** The agent is evaluated every 1,000 steps (line:71).

```

1  {
2    "ppo_freeway": {
3      "agent": [{
4        "name": "PPO",
5        "algorithm": {
6          "name": "PPO",
7          "action_pdtype": "default",
8          "action_policy": "default",
9          "explore_var_spec": null,
10         "gamma": 0.99,
11         "lam": 0.7,
12         "clip_eps_spec": {
13           "name": "no_decay",
14           "start_val": 0.1,
15           "end_val": 0.1,
16           "start_step": 0,
17           "end_step": 0
18         },
19         "entropy_coef_spec": {
20           "name": "no_decay",
21           "start_val": 0.01,
22           "end_val": 0.01,
23           "start_step": 0,
24           "end_step": 0
25         },
26         "val_loss_coef": 0.5,
27         "time_horizon": 128,
28         "minibatch_size": 256,
29         "training_epoch": 4
30       },
31       "memory": {
32         "name": "OnPolicyBatchReplay"
33       },
34       "net": {
35         "type": "ConvNet",
36         "shared": true,
37         "conv_hid_layers": [
38           [32, 8, 4, 0, 1],
39           [64, 4, 2, 0, 1],
40           [32, 3, 1, 0, 1]
41         ],
42         "fc_hid_layers": [512],
43         "hid_layers_activation": "relu",
44         "init_fn": "orthogonal_",
45         "normalize": true,
46         "batch_norm": false,
47         "clip_grad_val": 0.5,
48         "use_same_optim": false,
49         "loss_spec": {
50           "name": "MSELoss"
51         },
52         "actor_optim_spec": {
53           "name": "Adam",
54           "lr": 0.001
55         },
56         "critic_optim_spec": {
57           "name": "Adam",
58           "lr": 0.001
59         },
60         "lr_scheduler_spec": null,
61         "gpu": true
62       },
63     },
64     "env": [{
65       "name": "Freeway-v0",
66       "frame_op": "concat",
67       "frame_op_len": 4,
68       "reward_scale": "sign",
69       "num_envs": 16,
70       "max_t": null,
71       "max_frame": 500000
72     }],
73     "body": {
74       "product": "outer",
75       "num": 1
76     },
77     "meta": {
78       "distributed": false,
79       "log_frequency": 10000,
80       "eval_frequency": 10000,
81       "max_session": 4,
82       "max_trial": 1
83     }
84   }
85 }

```

Figure 28. PPO config file.

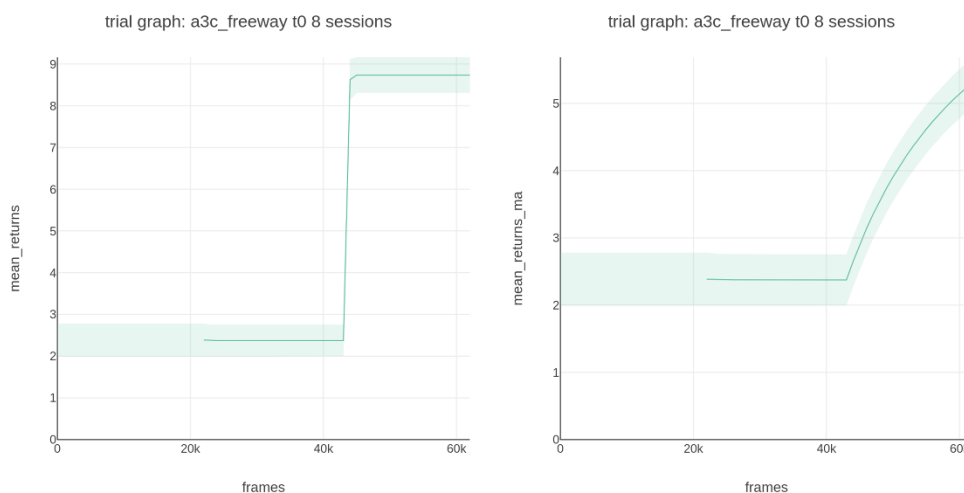


Figure 29. A3C 8 sessions graphs.

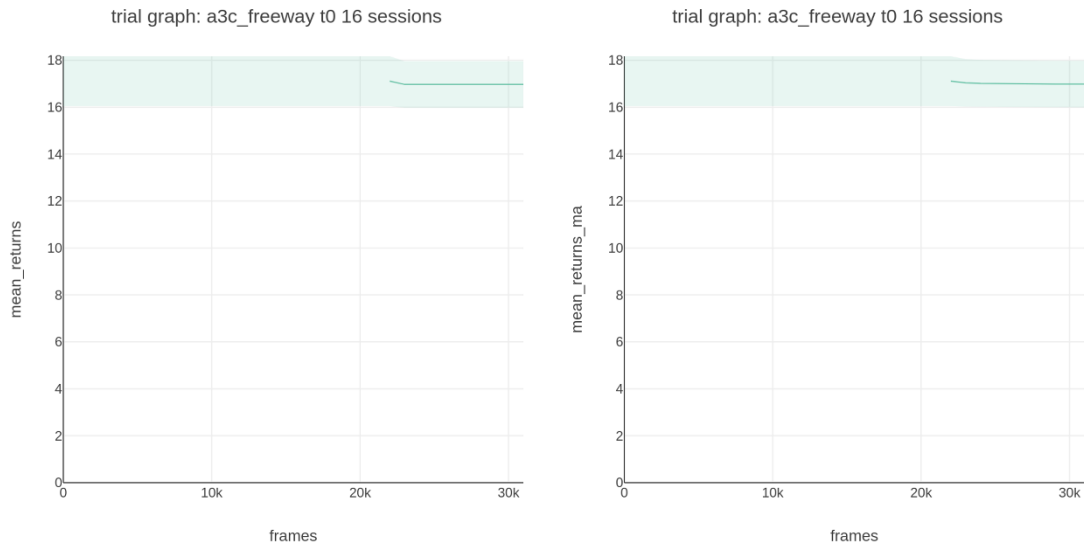


Figure 30. A3C 16 sessions graphs.

```

1  {
2    "a3c_freeway": {
3      "agent": [{
4        "name": "A3C",
5        "algorithm": {
6          "name": "ActorCritic",
7          "action_pdtype": "default",
8          "action_policy": "default",
9          "explore_var_spec": null,
10         "gamma": 0.99,
11         "lam": 0.95,
12         "num_step_returns": null,
13         "entropy_coef_spec": {
14           "name": "no_decay",
15           "start_val": 0.01,
16           "end_val": 0.01,
17           "start_step": 0,
18           "end_step": 0
19         },
20         "val_loss_coef": 0.5,
21         "training_frequency": 32
22       },
23       "memory": {
24         "name": "OnPolicyBatchReplay"
25       },
26       "net": {
27         "type": "ConvNet",
28         "shared": true,
29         "conv_hid_layers": [
30           [32, 8, 4, 0, 1],
31           [64, 4, 2, 0, 1],
32           [32, 3, 1, 0, 1]
33         ],
34         "fc_hid_layers": [512],
35         "hid_layers_activation": "relu",
36         "init_fn": "orthogonal_",
37         "normalize": true,
38         "batch_norm": false,
39         "clip_grad_val": 0.5,
40         "use_same_optim": false,
41         "loss_spec": {
42           "name": "MSELoss"
43         },
44         "actor_optim_spec": {
45           "name": "GlobalAdam",
46           "lr": 0.007
47         },
48         "critic_optim_spec": {
49           "name": "GlobalAdam",
50           "lr": 0.007
51         },
52         "lr_scheduler_spec": null,
53         "gpu": false
54       },
55     },
56     "env": [{
57       "name": "Freeway-v0",
58       "frame_op": "concat",
59       "frame_op_len": 4,
60       "reward_scale": "sign",
61       "num_envs": 8,
62       "max_t": null,
63       "max_frame": 500000
64     }],
65     "body": {
66       "product": "outer",
67       "num": 1
68     },
69     "meta": {
70       "distributed": "synced",
71       "log_frequency": 1000,
72       "eval_frequency": 1000,
73       "max_session": 8,
74       "max_trial": 1
75     }
76   }
77 }

```

Figure 31. A3C config file.

## 4. Conclusions

For the purpose of understanding the results obtained and comparing the different algorithms, I will present three metrics.

**Strength** refers to the average performance of the agent over the entire duration of the trial. This metric is usually calculated as the mean of the returns (cumulative rewards) obtained by the agent across all episodes or evaluation points during the trial. It provides an overall view of the agent's average performance, and it is useful for understanding general learning effectiveness.

**Max Strength** represents the highest performance achieved by the agent at any point during the trial. This metric is determined by identifying the maximum return (cumulative reward) obtained in any single episode or evaluation point. It provides information on the best possible performance the agent achieved, and it is useful for understanding the highest potential of the agent.

**Final Strength** indicates the performance of the agent at the end of the trial. This metric is typically the mean of the returns (cumulative rewards) over the last few episodes or evaluations, providing a measure of the agent's performance after completing its training. It reflects how well the agent performs after it has had the opportunity to learn from the entire training process.

Algorithm	Strength	Max Strength	Final Strength
REINFORCE	1.9053	2.3750	2.1666
SARSA	4.3648	5.5781	5.3750
DQN	11.8797	20.9218	19.7968
DQN w/ target nets	17.5363	21.5468	21.4687
DDQN	14.7540	16.8125	16.1875
DQN PER	19.4317	22.3437	22.1562
A2C n-step returns	0.0013	0.0156	0.0000
A2C GAE	10.6579	17.1875	16.8906
PPO	16.2700	24.2968	23.4531
A3C 8 sessions	5.3196	8.7343	8.7343
A3C 16 sessions	16.9825	17.1618	16.9687

In conclusion, considering the presented metrics, the best overall score during the whole training process was achieved by DQN PER. Also, the best result and the best result after training was achieved by PPO. However, these results do not prove that one algorithm is better than another. The agents were trained on Google Colab and the results might be influenced by the hardware limitations of this platform.

## 5. Reference List

- [1] Gymnasium Documentation. <https://gymnasium.farama.org/environments/atari/freeway/>
- [2] Atari Age. [https://atariage.com/manual\\_html\\_page.php?SoftwareLabelID=192](https://atariage.com/manual_html_page.php?SoftwareLabelID=192)
- [3] Internet Archive. <https://archive.org/details/FreewayAtari2600900DPIScans>
- [4] Reinforcement Learning Course, Prof. Dr. Habil. Eng. Călin-Adrian POPA, 2025
- [5] SLM Lab. <https://slm-lab.gitbook.io/slm-lab>