# Parallelizing Reinforcement Learning

Gabriel S.J. Spadoni

January 2024

## 1  Introduction

Reinforcement Learning is an important and often used Machine Learning algorithm as it can enable an agent to learn about a task or an environment only by knowing what it can do in it. However, the greater the environment, the longer the training can take. In this report we will focus on how we can implement Reinforcement Learning and running parallel training.

First, we will explain what is Reinforcement Learning and how does it work. We will then explain how we managed to parallelize it. Finally, we will look at the results.

## 2  Reinforcement Learning Overview

In this section, we will explain what this Reinforcement Learning and how we implemented it.

So the goal of Reinforcement Learning is to make an agent learn how to reach goal by giving it rewards when the goal is successfully reached. In this way, the agent will try to maximize the rewards it gets with the actions it does and thus reaching the goal. The agent only knows the state and the actions it can do.

We have used as the basis for the code the tutorial on Reinforcement Learning of PyTorch [1]

We are using Deep Q Learning (DQN) which is a type of Reinforcement Learning. The approach used has two models that are learning, we have the policy that is our final model and the target model which is used to keep the old of the policy net.

We will now describe one epoch. So during one epoch the agent will play in the environment until it either reaches the goal or fails. Failure is determined here by reaching a maximum number of actions. For each iteration in the environment, we keep in memory the state, the next state, the chosen action and the reward (see in Fig. 1). Using this, if there is enough data in the memory, we optimize the policy net. That is we optimize it for the state action values using the value of the next state times the Gamma plus the reward. The Gamma is the discount factor, it is a hyper parameter that it used to enforce

Figure 1: One Iteration for the training

good start as early actions will be taken more into account. We are also using the Adam as an optimizer and Huber as Loss function, we are using Huber loss function [3] as it is less sensitive to outliers (See in Fig. 2).



Figure 2: Optimization Function

Following the optimization, we update target network using the state dict of the policy network.

# 3 Parallelization Approach

We will now describe the approach used to parallelize the training of Reinforcement Learning.

We used `pytorch` to parallelize the learning, indeed pytorch allows us to chose on which device a model is. In device, we can specify if we have two gpus which one of them we want to use, for example:

`policy_net = DQN(n_observations, n_actions).to(device)`

We used as a basis for parallelizing the code of the tutorial on dpp of PyTorch [4], which we have adapted to Reinforcement Learning.

For parallelizing, we use the libraries `multiprocessing`, `nn.parallel`, `distributed` and `utils.data.distributed` from pytorch. These different libraries allows us to spawn a number of process groups equal to the number of devices we want to use. For this we use the function `init_process_group` from the `distributed` library. It takes as parameters backend which is the backend used to communicate between groups, here we use ncll which the nvidia backend library for communicated between gpus. There is also rank and world size which are respectively the rank of the current process and the number of processes. So the `init_process_group` will prepare one process at a time.

In order to have multiple processes, we use from the `multiprocessing` library the function `spawn`, that launch a function of a number of processes, it will also pass to the function the index of the process which will the rank for `init_process_group`. So if we have two devices, it will spawn two process and

give them rank 0 and 1, so that one process will use device 0 and the other device 1.

After this, we create the trainer object and we pass it the rank, so that it correctly sends the model and the computation on the wanted device. We also use the function `DistributedDataParallel` from the library `nn.parallel` to synchronize the gradients of each model with all the replicas:

`self.target_net = DDP(target_net, device_ids=[gpu_id])`

and `self.policy_net = DDP(policy_net, device_ids=[gpu_id])`.

This is important such that the learning made in one process is kept for the final model. Here `gpu_id` corresponds to `rank`. After the training is finished for each process, we call the function `destroy_process_group` from the library `distributed`, this way each process is destroyed and we come back into a sequential execution.

## 4 Results

In this section, we will show the results and compare between a standard execution and a parallel execution the time of running and efficacy of learning. We will use the time for a full run, the loss value and also we will check if the agent is able to successfully complete the environment we trained on. In order to compare, we have run both version on the hpc, for a different number of iterations each. For the parallel, we used two gpus. For running it with different number of iterations we created a bash file, that launches the program for each different set up. So we launched with 100, 200, 500 and 1000 iterations.

Training with only one gpu took:

- For 100 iterations: 6.37 seconds

- For 200 iterations: 19.54 seconds

- For 500 iterations: 108.31 seconds
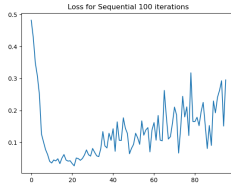
- For 1000 iterations: 735.39 seconds.



Figure 3: Loss for 100 iterations with one GPU

Here we see that loss function for 100 and 200 indeed in 100 we see that the agent is quickly learning and then it stops learning. We have the same for 200 and 500 iterations.
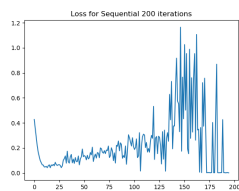
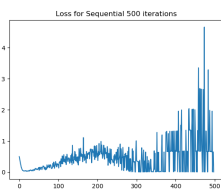Figure 4: Loss for 200 iterations with one GPU



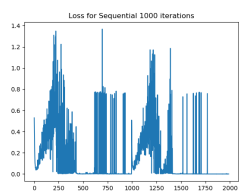Figure 5: Loss for 500 iterations with one GPU



Figure 6: Loss for 1000 iterations with one GPU

However regarding exploitation, the best performing model is the one with 200 iterations, indeed in Fig.7:
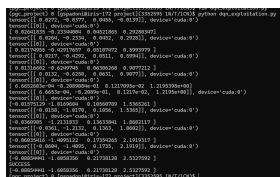


Figure 7: Exploitation for the model with 200 iterations

And the one that performs the worst is the one with 1000 iterations, indeed in Fig. 8:

Here we have the plot for the loss functions of the 100 iterations parallel Reinforcement Learning, see in Fig. 9:

This loss curve is good as we see it is quickly going down. However, the way

4

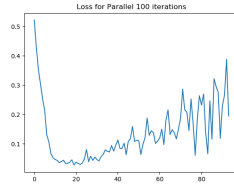Figure 8: Exploitation for the model with 1000 iterations



Figure 9: Loss values for 100 iterations in Parallel

we have set up the learning makes that we have a curve for only one process but the model is for two, which gives the double of total iterations. This explains the good performing model (See in Fig. 10):



Figure 10: Exploitation for 100 parallel iterations

Compared to the exploitation for 100 sequential iterations which has more steps to do before success, see in Fig. 11:

Here we have the curve of the loss for the parallel 200 iterations training (Fig. 13):

Figure 11: Exploitation for 100 iterations



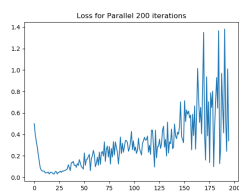Figure 12: Loss values for 200 parallel iterations



Figure 13: Exploitation for 200 parallel iterations

And the result of the exploitation for 200 iterations in parallel (See in Fig. 13):

As we see the number of steps used to atteign a successful state is really low.

We unfortunately don't have the run time for the parallel training.

# 5   Conclusion

As a conclusion, we have implemented parallel training for Reinforcement Learning. It seems to give correct results and help for the training of Reinforcement Learning.

For improving this project, we can apply parallel reinforcement learning on more complex environments as this environment is simple enough that standard RL models will learn it quickly. Indeed, we used a simple environment to help us having quick training times and to speed up the testing. The aim of this project was to get familiar with PyTorch especially for parallelization and to apply to Reinforcement Learning0.

# References

[1]   https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html
      https://github.com/pytorch/tutorials/blob/main/intermediate_source/reinforcement_q_learning.py

[2]  https://en.wikipedia.org/wiki/Q-learning

[3]  https://en.wikipedia.org/wiki/Huber_loss

[4]          https://github.com/pytorch/examples/blob/main/distributed/ddp-tutorial-series/multigpu.py