

# Utilising Hindsight Experience and Attention Mechanisms for On-policy Deep Reinforcement Learning methods in a Continuous Robotic Control environment

Hayden Sampson <sampsoh@edgehill.ac.uk>

2021

Hindsight Experience has been proven to be effective for improving the sample efficiency of Deep Reinforcement Learning (DRL) agents in sparse reward system environments. However, its application has only been utilised with Off-policy algorithms and methods such as Deep Deterministic Policy Gradient (DDPG), which are often in practice much more unstable in training than On-policy methods such as Proximal Policy Optimisation (PPO). Here, we have been able to apply Hindsight Experience in conjunction with PPO methods to provide a more stable and sample efficient application for a continuous robotic control environment. Our results show that the use of hindsight experience with these methods has a significant performance improvement over standard methods.

**Abstract**

## 1 Introduction

The following discussions will be presented as such: Section 2 as a Literature review to explore and discuss the relevant work appropriate to this study, Section 3 contains background describing some of the theory behind the algorithms and methods implemented in this study, Section 4 for a description of our model architecture and environment, followed by Section 5 discussing the setup of our experiments and subsequent results. Finally, Section 6 will draw some conclusions on our work with a discussion for potential future work and improvements.

## 2 Literature Review

There has been extensive research in studies regarding the areas in Deep Reinforcement Learning, including specific focus on the problems of robotic arm

movement and grasping. DeepMind’s Atari [1] first brought into the light the effectiveness of using deep reinforcement learning and ‘Deep Q Networks’ (DQN) for outperforming human performance over a range of Atari games. The novelty of this study showed how these algorithms can make use of graphical pixel images as state inputs to output appropriate actions in the environment. What was also impressive was the breakaway from model-based methods to model-free, where a network could show similar high-scoring performances when trained on a variety of different games.

In terms specific to robotic arm control, Google Brain [2] has explored the subject through a comparative implementation of a number of well-known off-policy learning methods. This presented a good starting place for analysing algorithmic performances for such tasks, and has proven to be highly influential for the motivation behind this proposed research. Where the study lacks is in the comparison and analysis of on-policy methods such as TRPO [3] and PPO [4] learning, which have also been proven through other studies to be effective for robotic control [5]. As well as this, the study has omitted methods that have been proven to improve algorithm performance and sample efficiency, such as Prioritised Experience Replay [Schaul2016] or Hindsight Experience Replay [Andry2018], that also prove to be highly effective for robotic arm environments.

Further research [Blake2020] explores increasing efficiency through attention models, with the aim of reducing the visual field of an agent to only focus on the important aspects of an input image, in order to reduce feature state space and dimensionality.

## 3 Background

### 3.1 Deep Reinforcement Learning

Generally, standard reinforcement learning is concerned with producing an *agent* that learns how to take appropriate *actions* within a given *environment* to collect *rewards* and succeed at a certain task. An indicator that a task has been achieved successfully is the amount of reward an agent will receive, therefore an agent will learn to act in whichever way will allow it to gain the most reward. Agents will make decisions on how to act based on their *observations* of the current *state* of the environment, producing state-action transition pairs. How an agent makes a decision on what action to take is called the agent’s *policy*. Taking an action causes the environment to change and update state as well as giving the agent some reward. The agent can then make a new observation on the updated state and take a new action.

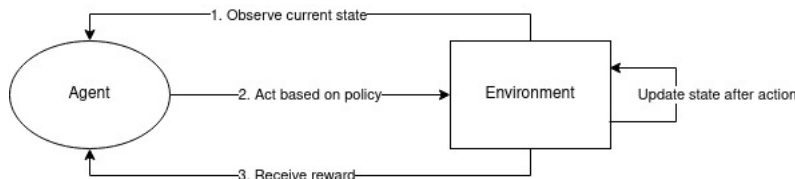


Figure 1: Reinforcement Learning Cycle

This cycle is repeated for  $T$  *timesteps* until the environment reaches its termination state, with the states between and including the initial state and final state defining an *episode*. The learning process therefore involves an agent figuring out which actions to take that will eventually lead to the greatest reward at the end of an episode, also known as its' *cumulative reward*.

The standard reinforcement learning approaches assume a discrete state and action space and hence, suffers from the *curse-of-dimensionality* problem, which simply states that computational complexity increases exponentially with increasing dimensions of state and action space. The problem of discrete state space can be solved by using deep networks that can take continuous values as inputs, including images. In other words, Deep Reinforcement Learning (DRL) methods use deep networks to overcome some of the limitations of standard reinforcement learning. These methods can be broadly divided into two categories: model-based and model-free methods. We are primarily interested in the model-free methods, that do not act based on their knowledge or estimation of the environment, but solely on the observation of the current state of the environment arrived at by using a policy. These model-free methods can be further divided into two sub-categories depending on how the agents learn, namely policy-based and value-based learning. Value-based methods estimate a *Q-value* (cumulative future reward) to every state-action pair by using a Q-function  $Q(s, a)$ , then selects an action which gives a maximum Q-value at a given state. So, in essence, value-based learning is concerned with estimating or approximating the Q-function. On the other hand, policy-based learning aims at estimating a policy directly without estimating a Q function. Policy parameters are updated to maximize or minimize a given objective function. For instance, policy gradient methods update the policy parameters by applying gradient ascent to maximize a given objective function. The advantage of using policy-based learning for this research problem is that it can be much more easily applied to agents with a continuous action space, as an action distribution is being directly output from the network, rather than a Q-value for each possible action available. It is possible to combine these two approaches in the form of an Actor-Critic algorithm [KonTsi2000] that employ a policy-based method to choose an action, with a value-based method to assign values to specific state-action pairs that critique how good that action was. Actor-Critic methods will form the underlying architecture used in this research. Another categorization can be made based on how the agents are trained. The on-policy methods train the agent using data generated using its current policy. This is an online approach

where the agent is learning while doing its job. In contrast, off-policy methods use a set of experience for training which are generated by a different policy than the one currently being used by the agent. Here, an on-policy approach is implemented through the state-of-the-art Proximal Policy Optimisation (PPO) algorithm.

### 3.2 Proximal Policy Optimisation

Proximal Policy Optimisation (PPO) is an on-policy algorithm first introduced by [4], that sought to improve performance of Trust Region Policy Optimisation (TRPO) [3], whilst also simplifying implementation and increasing sample efficiency. The idea behind the algorithm is simple, limit the magnitude in which the policy is updated after each training step to erase the possibility of the gradient update stepping too far, in which doing so controls the stability of the training process. The step control can be implemented through two methods: value clipping or through the use of a KL penalty. For clipping, there is simply a value  $\epsilon$  that acts as the upper and lower boundaries for the step magnitude. Alternatively, a KL Divergence method can be used to describe the similarity between two probability distributions and constrains the objective step function based on this. For this research, PPO-Clipping has been implemented.

PPO training consists of collecting experience data over  $K$  timesteps in the form of  $\langle s, a, r, ns, d \rangle$ , where  $\mathbf{s}$  represents the current state of the environment,  $\mathbf{a}$  the action taken in this state,  $\mathbf{r}$  the reward collected for taking this action in this state,  $\mathbf{ns}$  the next state achieved through taking the action in this state and  $\mathbf{d}$ , the indicator whether the episode has terminated at this state. Thus  $K$  samples in this form are collected, known as the set of *trajectories*. The total *returns* ( $V$ ) and *advantages* ( $A$ ) are then computed on this set using a *Generalised Advantage Estimation* (GAE) method [6], which provides values used in the Actor-Critic step functions.

---

**Algorithm 1** PPO-Clip

---

- 1: Input: initial policy parameters  $\theta_0$ , initial value function parameters  $\phi_0$
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:   Collect set of trajectories  $\mathcal{D}_k = \{\tau_i\}$  by running policy  $\pi_k = \pi(\theta_k)$  in the environment.
- 4:   Compute rewards-to-go  $\hat{R}_t$ .
- 5:   Compute advantage estimates,  $\hat{A}_t$  (using any method of advantage estimation) based on the current value function  $V_{\phi_k}$ .
- 6:   Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left( \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), \quad g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

typically via stochastic gradient ascent with Adam.

- 7:   Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left( V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 8: **end for**
- 

Figure 2: PPO  
OpenAI SpinningUp implementation [7]

### 3.3 Interpolated Policy Gradient

Interpolated Policy Gradient (IPG) is a hybrid algorithm that utilises techniques used in both on- and off-policy methods to *interpolate* and balance the offsets between them, improving stability over off-policy methods whilst controlling bias/variance problems in on-policy methods. IPG achieves this by involving three network models in training, including the usual Actor and Critic, as well as a *Baseline* network that rationalises the difference between old and new policies. For this reason, the baseline network can be implemented simply as the Value function used as the Critic in PPO. Therefore, the IPG critic is then substituted as a DDPG critic, utilising a Replay Buffer for off-policy updates in training. The novelty of IPG is implemented within the Actor policy loss function, where on- and off-policy values are calculated and summed to produce the total loss. The original on-policy loss in the paper is a standard on-policy objective function that can be readily substituted for the PPO policy objective function, which has been implemented here in the PPO-clip form. The off-policy update for the Actor uses Critic Q-values for predicted state-action pairs from a sample of the Replay Buffer to produce an off-policy loss value. These two loss values are then interpolated using an interpolate ratio value  $v$  between 0 and 1 to determine how much effect each loss will have on the total loss and overall update. In IPG, there is also an option to use a *control variate*, to further help the bias-variance balance, however has been omitted in the implementation here.

### 3.4 Hindsight Experience Replay

In regards to improving DRL methods to be more sample efficient, a lot of promise has been presented through the use of Hindsight Experience Replay (HER) [8]. This works more specifically for environments with sparse rewards, such as the one explored in this study. Originally, HER as a method had only been applicable to use with off-policy methods, such as Deep Deterministic Policy Gradient (DDPG) [9]. This is because the hindsight experience used in training is collected from a memory buffer that contained samples from all past policies. In order to make hindsight experience compatible with on-policy methods, Episodic Self Imitation Learning (ESIL) [10] was introduced which combined aspects of Self Imitation Learning (SIL) [11] and HER to allow hindsight experience collected from the current policy only to be used in training. ESIL extends standard Actor-Critic routines by introducing *goals* as states of which the agent is trying to achieve. As well as an initial state, a desired goal is initialised for the agent. The agent will then only receive a reward if the desired goal is reached, which is often not the case in a sparse reward system. For the agent to effectively learn from this, the achieved state of the agent is then considered to be the original desired goal as hindsight experience. This hindsight experience is then stored in memory along with the original trajectory experience. Therefore during training, the agent learns from sampling all of this experience. In the case where the agent fails to reach the original desired goal, it will still be able to learn useful information from the hindsight experience. This is the general concept of goal-based hindsight experience methods.

#### 3.4.1 Applying HER to IPG

### 3.5 Attention Mechanisms

Attention Mechanisms were first introduced in [reference] and were inspired by the way humans perceive and focus on certain aspects of our vision to take in the most important information for whatever task we are undergoing. Two main types of attention were realised, bottom-up saliency attention and top-down task relevant attention (or vice versa, check this). Saliency attention comes from the idea of what grabs our attention that is in our field of vision, in other words what is the most visually interesting thing that demands attention. Task related attention comes from the idea of what we need to pay attention to in order to fulfill a task, for example throwing a ball into a hoop, we would likely focus most of our attention on the hoop. Attention can be implemented through a number of ways [references], though here we opt for a *Self Attention* [reference] approach after the image input has passed through the Convolutional Neural Network process. This allows the Attention mechanism to attend to a much simpler, filtered input sequence, rather than a more complex multi-dimensional image input. In this kind of attention mechanism, there is no need for use of Recurrent Neural Networks (RNNs) or Long Short Term Memory (LSTM) components. Instead, attention mechanisms known as *Transformers* make use of input sequences to produce *key* and *query* vectors that undergo some matrix

manipulation and softmax activation that weight the indices of the sequence that should be attended to. The output of a weighted sequence therefore allows a network to identify which parts of an input have a higher importance in regards to affecting performance. This saves the agent from needing to process and attend the entire input equally, as it learns to be more efficient by only processing the important areas of the input.

## 4 Model Architecture and Environment

### 4.1 Approach

In this research, we primarily make use of actor-critic models as shown in Figure 3 below. In this model, we have two separate networks - one for estimating the policy required for generating the necessary action and the other for approximating the value function that can be used to evaluate the goodness of the action taken. These two networks are updated iteratively to find the optimal policy for a given problem by maximizing the total cumulative reward. This architecture provides flexibility in selecting different architectures and algorithms for each of the actor and critic modules independently. Each of the actor and critic models will use a deep neural network (DNN) to estimate action and value from the input state information. The environment provides the reward and the next state for a given state-action pair. The goal of learning will be to update parameters and so that certain objective cost functions are either maximized and minimized.

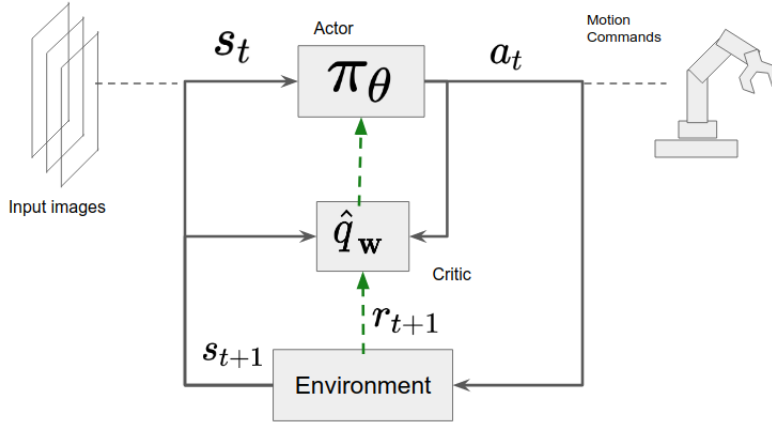


Figure 3: Actor-Critic Architecture

#### 4.1.1 Simulation Environment

Our objective is to develop, test and validate new deep reinforcement learning algorithms for solving robotic problems. In this regard, we have selected Py-

Bullet’s KukaCam [12] as our simulation environment. Some of the screenshots for this environment are shown in Figure 4. In this environment, a 6DoF Kuka robotic arm is expected to pick up different kinds of objects from the tray kept on the table based on the input images seen by an overhead camera. The view of the camera is shown in the top left corner of each image. This image along with the current joint values of the robotic arm (robot pose) forms the current ‘state’ of the environment. The agent is rewarded a value of 1 whenever the robot successfully picks up an item from the tray, otherwise the rewarded value is 0, implying a sparse reward system. Each episode consists of a series of robot movements attempting to pick an item. The objective is to pick the item successfully based on an input image obtained from the overhead camera.

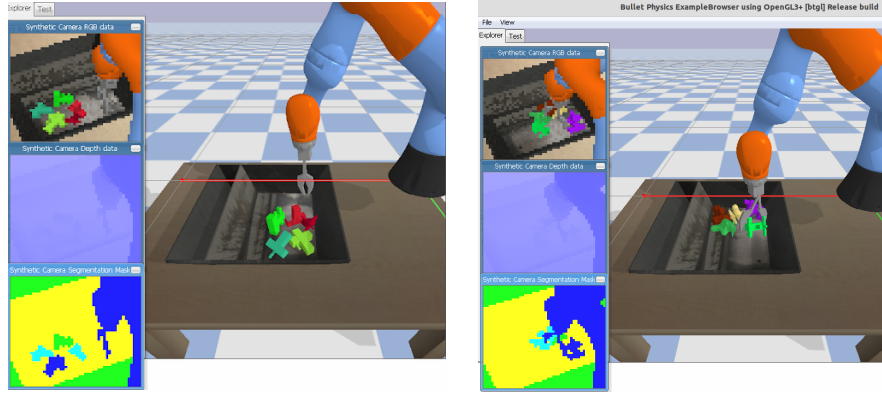


Figure 4: Pybullet KukaCam Diverse Object Environment

#### 4.1.2 Convolutional Neural Network

#### 4.1.3 Actor-Critic Network

## 5 Experiments

### 5.1 Hyperparameters and Setup

All experiments have been performed on an HP Omen GPU laptop for a maximum of 35 seasons. Here, a seasons consists of 1024 timesteps. Therefore, experiments are ran for a maximum total of 35x1024 or 35,840 timesteps. An episode generally consists of 7-8 timesteps with a maximum of 20. If an agent manages to reach a set success value, then it has also been considered to solve the problem and therefore will also end the experiment. For these experiments, when an agent succeeds in picking up an object an average of 70 times over the previous seasons, then it is considered to have solved the problem. Policy training data is collected over a season and then training is completed using batches for a number of epochs. The seasonal data is then cleared so that new data is collected on the updated policy. For IPG, state-action transition data is



stored into a Replay Buffer of maximum size 50,000. A validation procedure is ran after every 10 seasons to test the performance of the current policy. This validation runs 50 episodes of the environment, where the agent’s success rate is calculated. This is to evaluate the performance of the current policy that is separate from the main process of seasonal episodes and timesteps.

Hyper-parameter	Value
Seasons	35
Success value	70
Actor Learning Rate	0.0002
Critic Learning Rate	0.0002
CNN Learning Rate	0.0002
Epochs	10
Training Batch (steps per season)	1024
Batch Size	128
Clipping Value $\epsilon$	0.2
Discount factor $\gamma$	0.993
Lambda $\lambda$	0.7
Interpolate Ratio $v$ (IPG only)	0.2

Table 1: Hyperparameters

## 5.2 Results

# 6 Discussions and Conclusion

## 6.1 Conclusion

## 6.2 Future Work

## References

- [1] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013.
- [2] Deirdre Quillen, Eric Jang, Ofir Nachum, Chelsea Finn, Julian Ibarz, and Sergey Levine. Deep reinforcement learning for vision-based robotic grasping: A simulated comparative evaluation of off-policy methods. *CoRR*, abs/1802.10264, 2018.
- [3] John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR*, abs/1502.05477, 2015.

- [4] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- [5] A. Rupam Mahmood, Dmytro Korenkevych, Gautham Vasan, William Ma, and James Bergstra. Benchmarking reinforcement learning algorithms on real-world robots. *CoRR*, abs/1809.07731, 2018.
- [6] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation, 2018.
- [7] Joshua Achiam. Spinning up in deep reinforcement learning. ., 2018.
- [8] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. *CoRR*, abs/1707.01495, 2017.
- [9] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In Yoshua Bengio and Yann LeCun, editors, *ICLR*, 2016.
- [10] Tianhong Dai, Hengyan Liu, and Anil Anthony Bharath. Episodic self-imitation learning with hindsight, 2020.
- [11] Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. *CoRR*, abs/1806.05635, 2018.
- [12] Y. Bai E. Coumans. pybullet, a python module for physics simulation in robotics, games and machine learning, 2017.