

**General Instructions (summarized, to be removed before submission)**

- Structure — article-like, presenting the findings of our work clearly and objectively, provide well-supported analysis backing up our claims with references or empirical/theoretical evidence where appropriate
- Length — 4 to 8 pages, excluding references and appendices.
- Focus — provide insights valuable for future researchers or practitioners looking to reproduce / build upon the paper
- Evaluation — not solely based on quantitative outcomes, but rather on depth, quality of analysis and clarity of the report

## ATML Report

Jonatan Bella

jonatan.bella@usi.ch

Alessia Berarducci

alessia.berarducci@usi.ch

Jonas Knupp

jonas.knupp@usi.ch

Tobias Erbacher

tobias.erbacher@usi.ch

### Abstract

In this project we aim to analyze the claims of Mahankali et al. (2024) in the paper RANDOM LATENT EXPLORATION FOR DEEP REINFORCEMENT LEARNING which introduces a new technique to explore the state space and thus yield better overall agent scores. Moreover, we intend to reproduce the results obtained by the authors, expand on the findings and verify that the results are not cherry-picked. **<write the main findings of our project>** This report is prepared as part of the course project in *Advanced Topics in Machine Learning* at *Università della Svizzera italiana* in the autumn semester of 2024.

## 1 Introduction

The paper we investigate deals with the problem of motivating an agent in a high-dimensional state space to explore the environment more exhaustively during training and thereby find non-obvious trajectories that can lead to higher long-term rewards for both discrete and continuous action spaces. The paper compares the new *Random Latent Exploration* (RLE) technique to standard *Proximal Policy Optimization* (PPO, see Schulman et al. (2017)), *NoisyNet* (see Fortunato et al. (2018)) and *Random Network Distillation* (RND, see Burda et al. (2019)).

Exploration is one of the major issues of Reinforcement Learning (RL) and its techniques can generally be divided into two categories: Noise-based exploration and bonus-based exploration. There are advantages and disadvantages for both of them but we always have to keep in mind that always choosing the highest short-term (local) reward does not necessarily also yield the highest long-term (global) reward, e.g. picture the environment as shown in figure 1a. Here, the agent starting in the **blue state** will always choose the small reward of 1, i.e. going right, then move back to the initial state, then move right, and so on, dithering between these two states, instead of accepting to collect a reward of 0 by going left in order to be able to collect the much higher reward of 100 in the following step.



Figure 1: An exemplary environment consisting of four state where transitions can occur between neighboring states. The **blue node** is the initial state and the numbers are the rewards.

### 1.1 Random Latent Exploration

The new idea that the authors Mahankali et al. (2024) present is to augment the reward function by adding a randomized term which incentivizes the agent to explore a larger portion of the state space. In particular, we will call this term  $F(s, z)$  where  $s \in \mathcal{S}$  is any state of the state space  $\mathcal{S}$ , and  $\mathbf{z} \in \mathbb{R}^d$  is a  $d$ -dimensional vector sampled from a given distribution  $P_{\mathbf{z}}$  **<elaborate on the distribution>**. So if at time step  $t \in \{0, 1, 2, \dots, T\}$  the agent in state  $s_t$  takes action  $a_t \sim \pi(\cdot | s_t)$  from policy  $\pi$ , then it will obtain the task reward  $r(s_t, a_t)$ . In RLE, we train a so-called latent-conditioned policy network  $\pi(\cdot | s, \mathbf{z})$  and latent-conditioned value network  $V^\pi(s, \mathbf{z})$  to estimate and then maximize the expected sum of rewards from a given state, aware of the random term  $z$ :

$$V^\pi(s, \mathbf{z}) \approx \mathbb{E}_\pi \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + F(s_{t+1}, \mathbf{z})) \right] \quad (1)$$

In equation 1,  $\gamma$  describes the discount factor and  $F(s_{t+1}, \mathbf{z}) = \phi(s) \cdot \mathbf{z}$  where  $\phi(s) : \mathcal{S} \rightarrow \mathbb{R}^d$  is a feature extraction network. Both  $\phi$  and  $\mathbf{z}$  are  $d$ -dimensional vectors. Note that  $z$  is resampled only before a new trajectory  $(s_0, a_0, r_0, s_1, \dots)$  is explored, not before each action that the agent takes. To continue the example from above, in figure 1b the agent would prioritize to go left as this promises the greater reward. It is not hard to imagine that for higher dimensional environments, every time we start a new trajectory and correspondingly sample a new  $z$ , the agent will explore a different region of the state space and thus we can discover non-obvious paths to maximize the rewards. For the detailed pseudocode of this algorithm, see Appendix A **(check link before submission)**.

### 1.2 Proximal Policy Optimization

The algorithm constituting the foundation of RLE is Proximal Policy Iteration (PPO), described in Schulman et al. (2017), which can itself be run to solve RL problems. It is an on-line policy gradient method where after running the policy  $\pi_{\theta_{\text{old}}}$  in parallel with  $N$  actors for  $T \ll E$  timesteps, where  $E$  is the length of an episode, we compute the advantage estimator for each time index  $t \in [0, T]$ :

$$\hat{A}_t = \delta_t + (\gamma\lambda) \delta_{t+1} + \dots + (\gamma\lambda)^{T-t+1} \delta_{T-1} \quad (2)$$

where we have  $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$ . Here,  $\lambda$  is the Generalized Advantage Estimation (GAE) parameter and  $V(s)$  is a learned state-value function. Afterwards, for each  $t$  the surrogate objective is computed and the parameters  $\theta$  are optimized e.g. with minibatch SGD or Adam (note that we optimize by maximizing the objective). This surrogate objective takes the form:

$$L_t^{\text{CLIP}+\text{VF}+\text{S}}(\theta) = \hat{\mathbb{E}}_t [L_t^{\text{CLIP}}(\theta) - c_1 L_t^{\text{VF}}(\theta) + c_2 S[\pi_\theta](s_t)] \quad (3)$$

In equation 3, we have the VF coefficient  $c_1$  and the entropy coefficient  $c_2$ , together with the (optional, to enhance exploration) entropy bonus  $S[\pi_\theta](s_t)$ , which is calculated as the sum of  $-p \log p$  over each of the action-probabilities  $p$  resulting from the policy distribution  $\pi_\theta$  for a state  $s_t$ , and the objectives:

$$L_t^{\text{CLIP}}(\theta) = \hat{\mathbb{E}}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (4)$$

$$L_t^{\text{VF}}(\theta) = \left( V_\theta(s_t) - V_t^{\text{targ}} \right)^2 \quad (5)$$

In equation 4 we denote the probability ratio as  $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ , where  $\pi_{\theta_{\text{old}}}$  is the policy before the update, and the clipping function is:

$$\text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) = \begin{cases} 1 - \epsilon, & r_t(\theta) < 1 - \epsilon \\ r_t(\theta), & 1 - \epsilon \leq r_t(\theta) \leq 1 + \epsilon \\ 1 + \epsilon, & r_t(\theta) > 1 + \epsilon \end{cases} \quad (6)$$

Moreover in equation 5,  $V_t^{\text{targ}}$  represents some target state-value function that is to be obtained. Note that there exist also other objective functions mentioned in Schulman et al. (2017) that can substitute  $L_t^{\text{CLIP}}(\theta)$  in equation 3. In equation 6, we make use of a hyperparameter  $\epsilon \in [0, 1]$  and in combination with the minimizer from equation 4 we can prevent a single policy update to accidentally ruin the policy forever by sending the algorithm to a gradient region with a worse local extremum. The detailed pseudocode of this algorithm can be found in Appendix B ([check link before submission](#)).

In the paper that we investigate, the authors Mahankali et al. (2024) compare the empirical results of the RLE to PPO as the baseline, as well as two other exploration techniques, namely NoisyNet and RND, which will be briefly explained in the following sections.

### 1.3 NoisyNet

The central idea in a NoisyNet architecture brought forward by Fortunato et al. (2018) is to introduce random noise into the weights learned by the network ( $p$  inputs,  $q$  outputs), i.e. if we have a neural network whose output can be parametrized by a vector of weights  $\theta$  as  $\mathbf{y} = f_\theta(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + \mathbf{b}$ , then we can learn the set of parameter vectors  $\zeta = (\boldsymbol{\mu}, \boldsymbol{\Sigma})$  to compute  $\mathbf{w} = (\boldsymbol{\mu}^w + \boldsymbol{\sigma}^w \odot \boldsymbol{\varepsilon}^w)$  and  $\mathbf{b} = (\boldsymbol{\mu}^b + \boldsymbol{\sigma}^b \odot \boldsymbol{\varepsilon}^b)$  with the following dimensionalities:  $\boldsymbol{\mu}^w, \boldsymbol{\sigma}^w, \boldsymbol{\varepsilon}^w \in \mathbb{R}^{q \times p}$  implying  $\mathbf{w} \in \mathbb{R}^{q \times p}$ , and  $\boldsymbol{\mu}^b, \boldsymbol{\sigma}^b, \boldsymbol{\varepsilon}^b \in \mathbb{R}^q$  resulting in  $\mathbf{b} \in \mathbb{R}^q$ . Note that  $\mathbf{x} \in \mathbb{R}^p$ ,  $\mathbf{y} \in \mathbb{R}^q$  and  $\odot$  denotes element-wise multiplication.

We optimize the loss  $\bar{L}(\zeta) = \mathbb{E}_\varepsilon[L(\theta)]$ , i.e. the expectation of the loss  $L(\theta)$  over the noise  $\varepsilon$  which has a mean of zero and fixed statistics, using gradient descent on the parameters  $\zeta$ . The noise can be sampled either independently or in a factorized way from a Gaussian. The former method samples the noise  $\varepsilon_{i,j}^w \in \boldsymbol{\varepsilon}^w$  and  $\varepsilon_j^b \in \boldsymbol{\varepsilon}^b$  for every input to the network leading to a total of  $pq + q$  random samples. Factorized sampling on the other hand greatly reduces computational complexity by factorizing  $\varepsilon_{i,j}^w = f(\varepsilon_i)f(\varepsilon_j)$  and then  $\varepsilon_j^b = f(\varepsilon_j)$  thereby lowering the number of required random samples to  $p + q$ , where  $f: \mathbb{R} \rightarrow \mathbb{R}$ . In Fortunato et al. (2018) the function  $f(x) = \text{sign}(x)\sqrt{|x|}$  is used. After sampling the noise, in theory we compute the gradients  $\nabla \bar{L}(\zeta) = \nabla \mathbb{E}_\varepsilon[L(\theta)] = \mathbb{E}_\varepsilon[\nabla_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} L(\boldsymbol{\mu} + \boldsymbol{\Sigma} \odot \boldsymbol{\varepsilon})]$ , whereas in practice we approximate  $\nabla \bar{L}(\zeta) \approx \nabla_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} L(\boldsymbol{\mu} + \boldsymbol{\Sigma} \odot \boldsymbol{\varepsilon})$  with the Monte Carlo method.

In the code of Mahankali et al. (2024), a NoisyNet is used in conjunction with the A3C-algorithm, in particular the loss is computed from the action-value function is estimated using:

$$\hat{Q}_i = \sum_{j=i}^{k-1} \gamma^{j-i} r_{t+j} + \gamma^{k-i} V(x_{t+k} | \zeta, \varepsilon_i) \quad (7)$$

Moreover, for factorized networks, we initialize  $\mu_{i,j} \sim \mathcal{U} \left[ -\frac{1}{\sqrt{p}}, +\frac{1}{\sqrt{p}} \right]$  and  $\sigma_{i,j} = \frac{\sigma_0}{\sqrt{p}}$  with  $p$  being the number of inputs to the corresponding linear layer,  $\mathcal{U}$  is a uniform distribution over the specified interval, and  $\sigma_0$  is a hyperparameter. The pseudocode for this variant given by Fortunato et al. (2018) can be found in appendix Appendix C [Check link before submission](#).

Are we using independent or factorized sampling in the code? (I'm guessing factorized).

## 1.4 Random Network Distillation

In Random Network Distillation, as described by Burda et al. (2019), we compose the reward which the agent obtains as  $r_t = e_t + i_t$  where  $e_t$  represents a sparse extrinsic (environment's) reward and  $i_t$  is the intrinsic reward for transitioning, i.e. the exploration bonus emanating from the transition at time step  $t$ . The latter measures the novelty of a state and should yield a higher value, the less frequently a state has been visited so far. In case of a finite state space we can use  $i_t = \frac{1}{n_t}$  or  $i_t = \frac{1}{\sqrt{s}}$  where  $n_t(s)$  denotes the number of times state  $s$  has been visited until time step  $t$ . However, there exist alternatives, e.g. we could use state density based approaches to calculate an exploration bonus, which are described in the paper mentioned above. Moreover, they specify that in the paper the intrinsic reward is the prediction error of a randomly generated problem  $i_t = \|\hat{f}(x|\theta) - f(x)\|^2$  involving a target network  $f: \mathcal{O} \rightarrow \mathbb{R}^k$  (fixed and randomly initialized, maps an observation  $\mathcal{O}$  to an embedding  $\mathbb{R}^k$ ) to sample the problem as well as a predictor network  $\hat{f}: \mathcal{O} \rightarrow \mathbb{R}^k$  which was trained on the agent's data by gradient descent to optimize  $i_t$  w.r.t. parameters  $\theta_{\hat{f}}$ .

Since the overall return can be composed as a sum of returns  $R = R_E + R_I$  we can train two heads  $V_E$  and  $V_I$  to estimate the value function  $V = V_E + V_I$  and in extension the advantages. In the end, a policy is trained with standard PPO using the advantage estimators. Note that the intrinsic rewards have to be normalized in order to be useful because otherwise we could not guarantee that they are on a consistent scale, which is done by dividing  $i_t$  by a running estimate of the standard deviations of the intrinsic return. Furthermore, the observation also has to be normalized as  $x \leftarrow \text{clip}(\frac{x-\mu}{\sigma}, -5, 5)$ . The normalization parameters can be initialized by letting a random agent briefly move in the environment before beginning the optimization. The pseudocode given by Burda et al. (2019) can be found in appendix Appendix D [Check link before submission](#).

## 2 Scope of reproducibility

Introduce the specific setting or problem addressed in this work, and list the main claims from the original paper. Think of this as writing out the main contributions of the original paper. Each claim should be relatively concise; some papers may not clearly list their claims, and one must formulate them in terms of the presented experiments.

A claim should be something that can be supported or rejected by your data. An example is, "Finetuning pre-trained BERT on dataset X will have higher accuracy than an LSTM trained with GloVe embeddings." This is concise and is something that can be supported by experiments. An example of a claim that is too vague, and can't be supported by experiments, is "Contextual embedding models have shown strong performance on a number of tasks. We will run experiments evaluating two types of contextual embedding models on datasets X, Y, and Z."

This section roughly tells a reader what to expect in the rest of the report. Clearly itemize the claims you are testing:

- Claim 1
- Claim 2
- Claim 3
- Claim 4

Each experiment in Section 4 will support (at least) one of these claims, so a reader of your report should be able to separately understand the *claims* and the *evidence* that supports them.

### 3 Methodology

Explain your approach – did you use the author’s code, or did you aim to re-implement the approach from the description in the paper? Summarize the resources (code, documentation, GPUs) that you used.

#### 3.1 Model descriptions

Include a description of each model or algorithm used. Be sure to list the type of model, the number of parameters, and other relevant info (e.g. if it’s pre-trained).

#### 3.2 Datasets

Since we are dealing with a reinforcement learning task, we did not use any dataset to train our models but rather four different environments which will explain in section 3.4 ([check link before submission](#)).

#### 3.3 Hyperparameters

Describe how the hyperparameter values were set and what was the source for their value (e.g. paper, code, or your guess). If there was a hyperparameter search done, be sure to include the range of hyperparameters searched over, the method used to search (e.g. manual search, random search, grid search, etc.), and the best hyperparameters found. Include the number of total experiments (e.g. hyperparameter trials). You can also include all results from that search (not just the best-found results).

#### 3.4 Experimental setup and code

Include a description of how the experiments were set up that’s clear enough a reader could replicate the setup. Include a description of the specific measure used to evaluate the experiments (e.g. accuracy, precision@K, BLEU score, etc.). Provide a link to your code or notebook (if available). Add in this section (or in the following) any reference to cloud-based providers you

#### 3.5 Computational requirements

Include a description of the hardware used, such as the GPU or CPU the experiments were run on. For each model, include a measure of the average runtime (e.g. average time to predict labels for a given validation set with a particular batch size). For each experiment, include the total computational requirements (e.g. the total GPU hours spent).

Note: You’ll likely have to record this as you run your experiments, so it’s better to think about it ahead of time.

### 4 Results

Start with a high-level overview of your results. Do your results support the main claims of the original paper? Keep this section as factual and precise as possible, and reserve your judgment and discussion points for the next "Discussion" section.

#### 4.1 Results reproducing original paper

For each experiment, say 1) which claim in Section 2 it supports, and 2) if it successfully reproduced the associated experiment in the original paper. For example,

an experiment training and evaluating a model on a dataset may support a claim that that model outperforms some baseline. Logically group related results into subsections.

#### **4.1.1 Result 1**

#### **4.1.2 Result 2**

### **4.2 Results beyond original paper**

Often papers don't include enough information to fully specify their experiments, so some additional experimentation may be necessary. For example, it might be the case that batch size was not specified, and so different batch sizes need to be evaluated to reproduce the original results. Include the results of any additional experiments here. Note: this won't be necessary for all reproductions.

#### **4.2.1 Additional Result 1**

#### **4.2.2 Additional Result 2**

## **5 Discussion**

Give your judgment on if your experimental results support the claims of the paper. Discuss the strengths and weaknesses of your approach - perhaps you didn't have time to run all the experiments, or perhaps you did additional experiments that further strengthened the claims in the paper.

### **5.1 What was easy**

Give your judgment of what was easy to reproduce. Perhaps the author's code is clearly written and easy to run, so it was easy to verify the majority of original claims. Or, the explanation in the paper was really easy to follow and put into code.

Be careful not to give sweeping generalizations. Something that is easy for you might be difficult for others. Put what was easy in context and explain why it was easy (e.g. code had extensive API documentation and a lot of examples that matched experiments in papers).

### **5.2 What was difficult**

List part of the reproduction study that took more time than you anticipated or you felt was difficult.

Be careful to put your discussion in context. For example, don't say "The maths was difficult to follow", say "The math requires advanced knowledge of calculus to follow".

### **5.3 Communication with original authors**

Document the extent of (or lack of) communication with the original authors. To make sure the reproducibility report is a fair assessment of the original research we recommend getting in touch with the original authors. We advise you to contact them through their email displayed in the paper and by asking specific questions.

## **6 Conclusion**

Try to summarize the achievements of your project and its limits, suggesting (when appropriate) possible extensions and future works.

## Member contributions

Include a section specifying how you organized the work within the group and clearly describe the contributions of each member. Make sure to properly balance the workload among the members.

## References

- Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network distillation. *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=H1lJJnR5Ym>.
- Meire Fortunato, Mohammad G. Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. Noisy networks for exploration. *International Conference on Machine Learning*, 2018. doi: 10.48550/arXiv.1706.10295.
- Srinath Mahankali, Zhang-Wei Hong, Ayush Sekhari, Alexander Rakhlin, and Pulkit Agrawal. Random latent exploration for deep reinforcement learning. *International Conference on Machine Learning*, 2024. doi: 10.48550/arXiv.2407.13755. URL <https://srinathm1359.github.io/random-latent-exploration/>.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. 2017. doi: 10.48550/arXiv.1707.06347.

## A Pseudocode — Random Latent Exploration

For the RLE algorithm, the detailed pseudocode according to Mahankali et al. (2024) is as follows:

RLE:	Detailed Pseudocode
1:	<b>Input:</b> Latent distribution ( $P_{\mathbf{z}}$ ), number of parallel workers ( $N$ ), number of steps per update ( $T$ ), number of steps per sampling ( $S$ ), and feature network update rate ( $\tau$ ).
2:	Randomly initialize a feature network ( $\phi$ ) with the same backbone architecture as the policy and value networks.
3:	Initialize running mean $\mu = \mathbf{0}$ and standard deviation $\sigma = \mathbf{1}$ estimates of $\phi(s)$ over the state space.
4:	Sample an initial latent vector $\mathbf{z} \sim P_{\mathbf{z}}$ for each parallel worker.
5:	<b>Repeat</b> <sup>(1)</sup>
6:	Sample initial state $s_0$ .
7:	<b>For</b> $t = 0, \dots, T$ <b>do</b> <sup>(2)</sup>
8:	Take action $a_t \sim \pi(\cdot   s_t, \mathbf{z})$ and transition to $s_{t+1}$ .
9:	Compute feature $\mathbf{f}(s_{t+1}) = \frac{\phi(s_{t+1}) - \mu}{\sigma}$ .
10:	Compute random reward $F(s_{t+1}, \mathbf{z}) = \frac{\mathbf{f}(s_{t+1})}{\ \mathbf{f}(s_{t+1})\ } \cdot \mathbf{z}$ .
11:	Receive reward $r_t = R(s_t, a_t) + F(s_{t+1}, \mathbf{z})$ .
12:	<b>For</b> $i = 0, 1, \dots, N - 1$ <b>do</b> <sup>(3)</sup>
13:	<b>If</b> worker $i$ terminated <b>or</b> $S$ timesteps passed without resampling, <b>then</b> <sup>(4)</sup>
14:	Resample sample $\mathbf{z} \sim P_{\mathbf{z}}$ for worker $i$ .
15:	<sup>(4)</sup>
16:	<sup>(3)</sup>
17:	<sup>(2)</sup>
18:	Update policy network $\pi$ and value network $V_{\pi}$ with the collected trajectory $(\mathbf{z}, s_0, a_0, r_0, s_1, \dots, s_T)$ .
19:	Update feature network $\phi$ using the value network's parameters $\phi \leftarrow \tau \cdot V^{\pi} + (1 - \tau) \cdot \phi$ .
20:	Update $\mu$ and $\sigma$ using the batch of collected experience.
21:	<sup>(1)</sup> <b>until convergence</b>

## B Pseudocode — Proximal Policy Optimization

For the PPO algorithm, the detailed pseudocode according to Schulman et al. (2017) is as follows:

PPO:	Detailed Pseudocode, Actor-Critic Setup
1:	<b>Input:</b> Number of iterations ( $I$ ), number of actors ( $N$ ), number of timesteps per iteration ( $T$ ), number of epochs ( $K$ ), minibatch size ( $M \leq NT$ ), learned state-value function implicitly in the advantage estimator ( $V(s)$ ), target value ( $V_t^{\text{targ}}$ ) implicitly in the objective $L$ .
2:	<b>For</b> $i = 1, \dots, I$ <b>do</b> <sup>(1)</sup>
3:	<b>For</b> $n = 1, \dots, N$ <b>do</b> <sup>(2)</sup>
4:	Run policy $\pi_{\theta_{\text{old}}}$ in environment for $T$ timesteps.
5:	Compute advantage estimates $\hat{A}_1, \dots, \hat{A}_T$ .
6:	<sup>(2)</sup>
7:	Optimize objective $L$ w.r.t. $\theta$ , for given $K$ and $M$ , using e.g. minibatch SGD or Adam.
8:	Update parameters $\theta_{\text{old}} \leftarrow \theta$ .
9:	<sup>(1)</sup>



## C Pseudocode — NoisyNet

For the NoisyNet + A3C algorithm, the detailed pseudocode according to Fortunato et al. (2018) is as follows:

<b>NoisyNet:</b>	Detailed Pseudocode, for each actor-learner thread
1:	<b>Input:</b> Global shared parameters $(\zeta_\pi, \zeta_V)$ , global shared counter $(T)$ , and maximal time $(T_{\max})$ .
2:	<b>Input (each thread):</b> Thread-specific parameters $(\zeta'_\pi, \zeta'_V)$ , set of random variables $(\varepsilon)$ , thread-specific counter $(t)$ , and roll-out size $(t_{\max})$ .
3:	<b>Output:</b> Policy $\pi(\cdot \zeta_\pi, \varepsilon)$ and value function $V(\cdot \zeta_V, \varepsilon)$ .
4:	$t \leftarrow 1$
5:	<b>Repeat</b> <sup>(1)</sup>
6:	Reset cumulative gradients: $d\zeta_\pi \leftarrow 0, d\zeta_V \leftarrow 0$ .
7:	Synchronize thread-specific parameters: $\zeta'_\pi \leftarrow \zeta_\pi, \zeta'_V \leftarrow \zeta_V$ .
8:	Set counter $c \leftarrow 0$ .
9:	Get state $x_t$ from environment.
10:	Sample noise $\xi \sim \varepsilon$ .
11:	Create lists for rewards $r \leftarrow []$ , actions $a \leftarrow []$ and states $x \leftarrow []$ .
12:	<b>Repeat</b> <sup>(2)</sup>
13:	Choose action $a_t \leftarrow \pi(\cdot x_t, \zeta'_{pi}, \zeta'_V)$ .
14:	$a[-1] \leftarrow a_t$
15:	Obtain reward $r_t$ and new state $x_{t+1}$ .
16:	$r[-1] \leftarrow r_t, x[-1] \leftarrow x_t, t \leftarrow t + 1, T \leftarrow T + 1, c \leftarrow c + 1$
17:	<sup>(2)</sup> <b>until</b> $x_t$ <b>is terminal</b> <b>or</b> $c > t_{\max}$
18:	<b>If</b> $x_t$ <b>is terminal</b> , <b>then</b> <sup>(3)</sup>
19:	$Q \leftarrow 0$
20:	<b>else</b>
21:	$Q \leftarrow V(x_t \zeta'_V, \xi)$
22:	<sup>(3)</sup>
23:	<b>For</b> $i = c - 1, \dots, 0$ <b>do</b> <sup>(4)</sup>
24:	$Q \leftarrow r[i] + \gamma Q$
25:	$d\zeta_\pi \leftarrow d\zeta_\pi \nabla_{\zeta'_\pi} \log(\pi(a[i] x[i], \zeta'_\pi, \xi))[Q - V(x[i] \zeta'_V, \xi)]$
26:	$d\zeta_V \leftarrow d\zeta_V \nabla_{\zeta'_V} [Q - V(x[i] \zeta'_V, \xi)]^2$
27:	<sup>(4)</sup>
28:	Update asynchronously parameter $\zeta_\pi \leftarrow \zeta_\pi + \alpha_\pi d\zeta_\pi$ .
29:	Update asynchronously parameter $\zeta_V \leftarrow \zeta_V - \alpha_V d\zeta_V$ .
30:	<sup>(1)</sup> <b>until</b> $T > T_{\max}$

## D Pseudocode — Random Network Distillation

For the RND algorithm, the detailed pseudocode according to Burda et al. (2019) is as follows:

RND:	Detailed Pseudocode
1:	<b>Input:</b> Number of rollouts ( $N$ ), Number of optimization steps ( $N_{\text{opt}}$ ), and length of initial steps for initializing observation normalization ( $M$ ).
2:	$t \leftarrow 0$
3:	Sample state $s_0 \sim p_0(s_0)$ .
4:	<b>For</b> $m = 1, \dots, M$ <b>do</b> <sup>(1)</sup>
5:	Sample action $a_t \sim \mathcal{U}(a_t)$ . Are we sampling uniformly from all available actions at time $t$ ?
6:	Sample state $s_{t+1} \sim p(s_{t+1} s_t, a_t)$ .
7:	Update observation normalization parameters using $s_{t+1}$ .
8:	$t \leftarrow t + 1$
9:	<sup>(1)</sup>
10:	<b>For</b> $i = 1, \dots, N$ <b>do</b> <sup>(2)</sup>
11:	<b>For</b> $j = 1, \dots, K$ <b>do</b> <sup>(3)</sup>
12:	Sample action $a_t \sim \pi(a_t s_t)$ .
13:	Sample state $s_{t+1}, e_t \sim p(s_{t+1}, e_t s_t, a_t)$ .
14:	Calculate intrinsic reward $i_t = \ \hat{f}(s_{t+1}) - f(s_{t+1})\ ^2$ .
15:	Add $s_t, s_{t+1}, a_t, e_t, i_t$ to optimization batch $B_i$ .
16:	Update running estimate of reward standard deviation using $i_t$ .
17:	$t \leftarrow t + 1$
18:	<sup>(3)</sup>
19:	Normalize the intrinsic rewards contained in $B_i$ .
20:	Calculate returns $R_{I,i}$ and advantages $A_{I,i}$ for intrinsic reward.
21:	Calculate returns $R_{E,i}$ and advantages $A_{E,i}$ for extrinsic reward.
22:	Calculate combined advantages $A_i = A_{I,i} + A_{E,i}$ .
23:	Update observation normalization parameters using $B_i$ .
24:	<b>For</b> $j = 1, \dots, N_{\text{opt}}$ <b>do</b> <sup>(4)</sup>
25:	Optimize $\theta_\pi$ w.r.t. PPO loss on batch $B_i, R_i, A_i$ using Adam.
26:	Optimize $\theta_{\hat{f}}$ w.r.t. distillation loss on $B_i$ using Adam.
27:	<sup>(4)</sup>
28:	<sup>(2)</sup>

## E Citations, figures, and tables

**Citations** When the authors or the publication are included in the sentence, the citation should not be in parenthesis, using `\citet{}` (as in “See (CITATION REMOVED) for more information.”). Otherwise, the citation should be in parenthesis using `\citep{}` (as in “Deep learning shows promise to make progress towards AI (CITATION REMOVED).”).

**Figures** All artwork must be neat, clean, and legible. The figure number and caption always appear after the figure. Place one line space before the figure caption, and one line space after the figure. The figure caption is lowercase (except for the first word and proper nouns). Make sure the figure caption does not get separated from the figure. Leave sufficient space to avoid splitting the figure and figure caption.

**Tables** All tables must be centered, neat, clean and legible. The table number and title always appear before the table. See Table 1. Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lowercase (except for the first word and proper nouns).

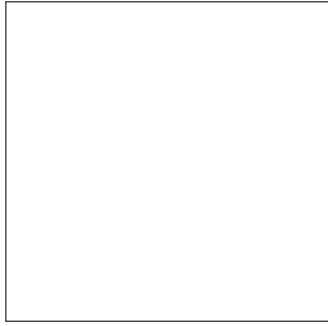


Figure 2: Sample figure caption.

Table 1: Sample table title

PART	DESCRIPTION
Dendrite	Input terminal
Axon	Output terminal
Soma	Cell body (contains cell nucleus)

## F Math Notation

In an attempt to encourage standardized notation, we have included the notation file from the textbook, *Deep Learning* (CITATION REMOVED) available at [https://github.com/goodfeli/dlbook\\_notation/](https://github.com/goodfeli/dlbook_notation/). Use of this style is not required and can be disabled by commenting out `math_commands.tex`.

### Numbers and Arrays

$a$	A scalar (integer or real)
$\mathbf{a}$	A vector
$\mathbf{A}$	A matrix
$\mathbf{A}$	A tensor
$\mathbf{I}_n$	Identity matrix with $n$ rows and $n$ columns
$\mathbf{I}$	Identity matrix with dimensionality implied by context
$\mathbf{e}^{(i)}$	Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at position $i$
$\text{diag}(\mathbf{a})$	A square, diagonal matrix with diagonal entries given by $\mathbf{a}$
$a$	A scalar random variable
$\mathbf{a}$	A vector-valued random variable
$\mathbf{A}$	A matrix-valued random variable

### Sets and Graphs

$\mathbb{A}$	A set
$\mathbb{R}$	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and $n$
$[a, b]$	The real interval including $a$ and $b$
$(a, b]$	The real interval excluding $a$ but including $b$
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of $\mathbb{A}$ that are not in $\mathbb{B}$
$\mathcal{G}$	A graph
$Pa_{\mathcal{G}}(\mathbf{x}_i)$	The parents of $\mathbf{x}_i$ in $\mathcal{G}$

### Indexing

$a_i$	Element $i$ of vector $\mathbf{a}$ , with indexing starting at 1
$\mathbf{a}_{-i}$	All elements of vector $\mathbf{a}$ except for element $i$
$A_{i,j}$	Element $i, j$ of matrix $\mathbf{A}$
$\mathbf{A}_{i,:}$	Row $i$ of matrix $\mathbf{A}$
$\mathbf{A}_{:,i}$	Column $i$ of matrix $\mathbf{A}$
$A_{i,j,k}$	Element $(i, j, k)$ of a 3-D tensor $\mathbf{A}$
$\mathbf{A}_{:,:,i}$	2-D slice of a 3-D tensor
$\mathbf{a}_i$	Element $i$ of the random vector $\mathbf{a}$

### Calculus

$\frac{dy}{dx}$	Derivative of $y$ with respect to $x$
$\frac{\partial y}{\partial x}$	Partial derivative of $y$ with respect to $x$
$\nabla_{\mathbf{x}} y$	Gradient of $y$ with respect to $\mathbf{x}$
$\nabla_{\mathbf{X}} y$	Matrix derivatives of $y$ with respect to $\mathbf{X}$
$\nabla_{\mathbf{x}} y$	Tensor containing derivatives of $y$ with respect to $\mathbf{X}$
$\frac{\partial f}{\partial \mathbf{x}}$	Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or $\mathbf{H}(f)(\mathbf{x})$	The Hessian matrix of $f$ at input point $\mathbf{x}$
$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of $\mathbf{x}$
$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to $\mathbf{x}$ over the set $\mathbb{S}$

### Probability and Information Theory

$P(\mathbf{a})$	A probability distribution over a discrete variable
$p(\mathbf{a})$	A probability distribution over a continuous variable, or over a variable whose type has not been specified
$\mathbf{a} \sim P$	Random variable $\mathbf{a}$ has distribution $P$
$\mathbb{E}_{\mathbf{x} \sim P}[f(\mathbf{x})]$ or $\mathbb{E}f(\mathbf{x})$	Expectation of $f(\mathbf{x})$ with respect to $P(\mathbf{x})$
$\text{Var}(f(\mathbf{x}))$	Variance of $f(\mathbf{x})$ under $P(\mathbf{x})$
$\text{Cov}(f(\mathbf{x}), g(\mathbf{x}))$	Covariance of $f(\mathbf{x})$ and $g(\mathbf{x})$ under $P(\mathbf{x})$
$H(\mathbf{x})$	Shannon entropy of the random variable $\mathbf{x}$
$D_{\text{KL}}(P \  Q)$	Kullback-Leibler divergence of $P$ and $Q$
$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian distribution over $\mathbf{x}$ with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$

### Functions

$f : \mathbb{A} \rightarrow \mathbb{B}$	The function $f$ with domain $\mathbb{A}$ and range $\mathbb{B}$
$f \circ g$	Composition of the functions $f$ and $g$
$f(\mathbf{x}; \boldsymbol{\theta})$	A function of $\mathbf{x}$ parametrized by $\boldsymbol{\theta}$ . (Sometimes we write $f(\mathbf{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation)
$\log x$	Natural logarithm of $x$
$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
$\ \mathbf{x}\ _p$	$L^p$ norm of $\mathbf{x}$
$\ \mathbf{x}\ $	$L^2$ norm of $\mathbf{x}$
$x^+$	Positive part of $x$ , i.e., $\max(0, x)$
$\mathbf{1}_{\text{condition}}$	is 1 if the condition is true, 0 otherwise