

A Multilevel Reinforcement Learning Framework for PDE-based Control

Atish Dixit, Ahmed H. Elsheikh

Herriot-Watt University,
Edinburgh, UK.

October 31, 2022

Abstract

Reinforcement learning (RL) is a promising method to solve control problems (Dixit and ElSheikh, 2022). However, model-free RL algorithms are sample inefficient and require thousands if not millions of samples to learn optimal control policies. A major source of computational cost in RL corresponds to the transition function, which is dictated by the model dynamics. This is especially problematic when model dynamics is represented with coupled PDEs. In such cases, the transition function often involves solving a large-scale discretization of the said PDEs. We propose a multilevel RL framework in order to ease this cost by exploiting sublevel models that correspond to coarser scale discretization (i.e. multilevel models). This is done by formulating an approximate multilevel Monte Carlo estimate (inspired by Giles (2015)) of the objective function of the policy and / or value network instead of Monte Carlo estimates, as done in the classical framework. As a demonstration of this framework, we present a multilevel version of the proximal policy optimization (PPO) algorithm. Here, the level refers to the grid fidelity of the chosen simulation-based environment. We provide two examples of simulation-based environments that employ stochastic PDEs that are solved using finite-volume discretization. For the case studies presented, we observed substantial computational savings using multilevel PPO compared to its classical counterpart.

1 Introduction

Optimal control problem involves finding controls for a dynamical system (often represented by a set of partial differential equations (PDEs)) such that a certain objective function is optimized over a predefined simulation time. In recent years, we have seen a surge in research activities where reinforcement learning (RL) has been demonstrated as an effective method to solve optimal control problems in fields such as energy (Anderlini et al., 2016), fluid dynamics (Rabault et al., 2019), and subsurface flow control (Dixit and ElSheikh, 2022). The reinforcement learning process for optimal control policy often involves a large number of exploration and exploitation attempts of control trajectories. In the context of PDE-based control problems, this corresponds to a large number of simulations of the underlying model dynamics. For large-scale PDE-based problems (i.e., with high-fidelity PDE discretization), this makes RL a computationally expensive process.

Since the introduction of the multilevel Monte Carlo (MLMC) estimate as a computationally cheaper counterpart to classical Monte Carlo estimates, numerous research studies have been conducted in the application of MLMC estimates in uncertainty quantification for stochastic PDEs

(Cliffe et al., 2011; Anderson and Higham, 2012; Giles and Szpruch, 2018). Furthermore, we also see a rise of MLMC estimate applications in certain deep learning research studies. For example, Shi and Cornish (2021) present a framework for MLMC-based unbiased gradient estimation in deep latent variable models. Chada et al. (2022) illustrate how the MLMC method could be applied to Bayesian inference using deep neural networks to compute expectations associated with the posterior distribution where the level corresponds to the sets of neural network parameters under consideration. In this paper, we introduce a novel multilevel framework for reinforcement learning where the learned agent interacts with environments corresponding to simulations of PDEs and the level corresponds to the grid fidelity of the PDE discretization.

We start by presenting the anatomy for classical RL algorithms, which involves estimating the Monte Carlo estimate of the objective function for policy and/or value network. Furthermore, we formulate the approximate MLMC estimation methodology used in the proposed multilevel framework. We then briefly present the mathematical framework that enables synchronized rollouts of task trajectories at different levels of the environment. The data generated through these synchronized rollouts are used to compute the approximate MLMC estimate of the objective function. Using the proposed multilevel framework, we formulate a multilevel variant of the state-of-the-art algorithm titled proximal policy optimization (PPO).

In the experiments presented, we compare the reinforcement learning process for classical and proposed multilevel PPO algorithms. The results are demonstrated for two environments for which the model dynamics is represented by stochastic partial differential equations. Furthermore, we also demonstrate the results of standard MLMC analysis to compare the MLMC and MC estimates for the PPO objective function. These environments were inspired by our research work in Dixit and ElSheikh (2022). In this study, the levels of the environment correspond to the discretization fidelity of the grid of the underlying PDEs.

The following is the outline for the rest of the paper: Section 2 provides the anatomy of the classical RL framework and formally defines the approximate MLMC estimation method. Section 3 introduces the multilevel framework for RL algorithms and further presents the multilevel PPO algorithm along with its analysis methodology. Numerical experiments to demonstrate the proposed multilevel PPO algorithm are detailed in Section 4, and the results of these experiments are delineated in Section 5. Finally, Section 6 concludes with a summary of the research study and an outlook on future research directions.

2 Background

Conventionally, the RL framework consists of the environment \mathcal{E} , which is governed by a Markov decision process described by the tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \mu \rangle$. Here, $\mathcal{S} \subset \mathbb{R}^{n_s}$ is the state-space, $\mathcal{A} \subset \mathbb{R}^{n_a}$ is the action-space, $\mathcal{P}(s'|s, a)$ is a Markov transition probability function between the current state s and the next state s' under action a and $\mathcal{R}(s, a, s')$ is the reward function. The function $\mu(s)$ returns a state from the initial state distribution if s is the terminal state of the episode (e.g., simulation terminal time); otherwise, it returns the same state s . The goal of reinforcement learning is to find the policy $\pi(a|s)$ to take an optimal action a when the state s is observed. In deep reinforcement learning, the policy is denoted $\pi_\theta(a|s)$ and is represented by a neural network with parameters θ either directly (for policy-based algorithms) or indirectly (for value-based algorithms). Learning is initiated with a random policy and then updated by exploring state-action spaces and exploiting the observed rewards in subsequent sampling steps. Each such update is referred to as a policy iteration.

The algorithm 1 outlines a general anatomy of deep reinforcement learning algorithms. Each

Algorithm 1 Anatomy of deep reinforcement learning algorithms

```
1: for policy iteration = 1, 2, ... do
2:   step 1: Generate sequences  $\{s_t, a_t, r_t\}_{t=1}^{t=T}$  using current policy  $\pi_\theta(a|s)$ 
3:   for  $t = 1, 2, \dots, T$  do
4:     generate samples  $s_t, a_t$  and  $r_t$ , where  $s, a, r \sim p_\theta$ 
5:     compute  $\Theta_t$ 
6:   end for
7:   step 2: Compute Monte Carlo estimate of objective function  $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta)]$ :
8:      $\widehat{\mathbb{E}}_{s,a,r \sim p_\theta}^T [J(s, a, r; \theta, \Theta)]$ 
9:   step 3: Update  $\theta$  using the gradient of the estimated objective function
10: end for
```

policy iteration consists of three steps. First, the sequence $\{(s_1, a_1, r_1), \dots, (s_T, a_T, r_T)\}$ is generated by rolling out the current policy $\pi_\theta(a|s)$. In this stage, the RL algorithm utilizes the current policy to interact with the simulated environment by providing actions (aka. controls) and recording the observed rewards. A shorthand notation $s, a, r \sim p_\theta$, is used for the definition of random variables s, a and r . Equation 1 provides a detailed expansion of this shorthand notation.

$$s, a, r \sim p_\theta \begin{cases} s \sim \mu(s) \\ a \sim \pi_\theta(a|s) \\ s' \sim \mathcal{P}(s, a) \\ r = \mathcal{R}(s, a, s') \end{cases} \quad (1)$$

The objective function used to calculate the gradient of the network parameters θ , is of the form $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta)]$, where Θ is a set of other parameters that can vary from one algorithm to another. Appendix A delineates this objective function for various algorithms. The second step consists of computing a Monte Carlo estimate of $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta)]$ which is calculated using the samples generated in the first step. The notation $\widehat{\mathbb{E}}_{x \sim \mathcal{X}}^T[f(x)]$ in algorithm 1 corresponds to the Monte Carlo estimate of $\mathbb{E}_{x \sim \mathcal{X}}[f(x)]$ which is calculated as $T^{-1} \sum_{t=0}^T f(x_t)$, where x_1, \dots, x_T are random samples of the random variable $x \sim \mathcal{X}$. To maintain brevity in the description of Monte Carlo estimate, we use the same notation in the rest of the paper. Finally, in the third step, the policy is updated by updating the network parameters θ using the gradient of the estimated objective function.

2.1 Approximate Multilevel Monte Carlo estimation

Monte Carlo estimate of $\mathbb{E}[f(x^L)]$ for the random variable $x^L \sim \mathcal{X}^L$ is defined as

$$\mathbb{E}_{x^L \sim \mathcal{X}^L}[f(x^L)] \approx \widehat{\mathbb{E}}_{x^L \sim \mathcal{X}^L}^T[f(x^L)],$$

where T denotes the number of samples used in the estimation. Suppose that we have functions φ_L^l that approximate the random variable x^L from level L to l , $\forall l \in \{1, 2, \dots, L\}$ (note that φ_L^L is simply an identity function). Functions φ_L^l are defined so that each decrease in level l corresponds to a proportional decrease in the accuracy and cost of computing the function $f(\varphi_L^l(x^L))$. In PDE-based uncertainty quantification problems, the function f represents the quantity of interest, which implicitly contains the solution for the said PDE. The level refers to the grid discretization used

during the PDE solving; that is, the grid discretization goes from coarsest to finest from level 1 to L . For such a multilevel representation of functions, MLMC estimate of $\mathbb{E}[f(x^L)]$ is defined as

$$\mathbb{E}_{x^L \sim \mathcal{X}^L}[f(x^L)] \approx \sum_{l=1}^L \widehat{\mathbb{E}}_{x^L \sim \mathcal{X}^L}^{T^l}[f(\varphi_L^l(x^L)) - f(\varphi_L^{l-1}(x^L))], \quad (2)$$

where T^l represents the number of samples at each level l , and the value of the function at the zeroth level is predefined at zero (that is, $f(\varphi_L^0(\cdot)) \doteq 0$). The MLMC estimate is introduced by Giles (2015) as a computationally cheaper alternative to the classical Monte Carlo estimate. Readers are referred to the Appendix B where we briefly explain the principle behind the computational savings in the MLMC estimation.

As described in Equation 2, the MLMC estimate is the telescopic sum of Monte Carlo estimates of the difference term $f(\varphi_L^l(x^L)) - f(\varphi_L^{l-1}(x^L)) \forall l \in \{1, 2, \dots, L\}$, for the random variable $x^L \sim \mathcal{X}^L$. We reformulate this MLMC estimate so that we can use approximate samples at each level instead of samples from the finest level L . This is done with the following two approximations: First, we treat $\varphi_L^l(x^L)$ as a random variable $x^l \sim \mathcal{X}^l, \forall l \in \{1, 2, \dots, L\}$. Second, we replace the second difference term from $\varphi_L^{l-1}(x^L)$ to $\varphi_L^{l-1}(x^l)$. In other words, the difference term can now be computed using an approximate random variable x^l as opposed to the random variable on the finest level x^L . Furthermore, this term $\varphi_L^{l-1}(x^l)$, is denoted with \tilde{x}^{l-1} , which represents the synchronized value of x^l at the level $l-1$. We denote this synchronization process by the shorthand notation $\tilde{x}^{l-1} = \mathcal{X}^{l \Rightarrow l-1}$ as a subscript. Taking these approximations into account, we formulate the approximate MLMC estimate as follows.

$$\sum_{l=1}^L \widehat{\mathbb{E}}_{\substack{x^l \sim \mathcal{X}^l \\ \tilde{x}^{l-1} = \mathcal{X}^{l \Rightarrow l-1}}}^{T^l}[f(x^l) - f(\tilde{x}^{l-1})]. \quad (3)$$

Note that with this formulation, we can employ the random variable x^l at each level l . This idea of using approximate samples at each level is at the heart of the proposed multilevel RL framework. In the rest of the paper, we use the Equation 3 notation to formulate the approximate estimate of MLMC.

3 Multilevel RL framework

We introduce a multilevel RL framework formulated as a tuple, $\langle \mathcal{E}, \psi_i^{l'}, \phi_i^{l'} \rangle$ where \mathcal{E} represents a set of multiple environments $\{\mathcal{E}^1, \mathcal{E}^2, \dots, \mathcal{E}^L\}$. An environment \mathcal{E}^L corresponds to the target task described by the tuple $\langle \mathcal{S}^L, \mathcal{A}^L, \mathcal{P}^L, \mathcal{R}^L, \mu^L \rangle$. Its corresponding sublevel tasks are represented as environments $\mathcal{E}^1, \mathcal{E}^2, \dots, \mathcal{E}^L$ such that the computational cost of \mathcal{P}^l and the accuracy of \mathcal{R}^l is lower than \mathcal{P}^{l+1} and \mathcal{R}^{l+1} for all values of $l \in \{1, \dots, L-1\}$. Furthermore, $\psi_i^{l'}(s^l)$ is a mapping function from state on level l (denoted as s^l) to state on level l' (denoted as $s^{l'}$) and similarly $\phi_i^{l'}(a^l)$ is a mapping function from action a^l to $a^{l'}$. The algorithm 2 outlines the anatomy of deep reinforcement learning algorithms with the proposed multilevel framework.

The first step consists of generating the sequence $\{(s_1^l, a_1^l, r_1^l), \dots, (s_{T_l}^l, a_{T_l}^l, r_{T_l}^l)\}$ on level l and its corresponding synchronized sequence $\{(\tilde{s}_1^{l-1}, \tilde{a}_1^{l-1}, \tilde{r}_1^{l-1}), \dots, (\tilde{s}_{T_l}^{l-1}, \tilde{a}_{T_l}^{l-1}, \tilde{r}_{T_l}^{l-1})\}$ on level $l-1$. The shorthand notation $s^l, a^l, r^l \sim p_\theta^l$, for generating rollouts at level l and $\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1} = p_\theta^{l \Rightarrow l-1}$ for

Algorithm 2 Anatomy for multilevel deep reinforcement learning algorithms

```

1: for policy iteration = 1, 2, ... do
2:   step 1: Generate sequences  $\{(s_t^l, a_t^l, r_t^l), (\tilde{s}_t^{l-1}, \tilde{a}_t^{l-1}, \tilde{r}_t^{l-1})\}_{t=1}^{t=T} \}_{l=1}^{l=L}$  with policy  $\pi_\theta(a^L|s^L)$ 
3:   for level  $l = 1, 2, \dots, L$  do
4:      $s_t^l = \psi_{l-1}^l(\tilde{s}_{T_{l-1}}^{l-1})$  ▷ if  $l > 1$ 
5:     for  $t = 1, 2, \dots, T_l$  do
6:       generate samples  $s_t^l, a_t^l, r_t^l$  where  $s^l, a^l, r^l \sim p_\theta^l$ 
7:       compute  $\Theta_t^l$ 
8:       generate synchronised samples  $\tilde{s}_t^{l-1}, \tilde{a}_t^{l-1}, \tilde{r}_t^{l-1}$  ▷ if  $l > 1$ 
9:       where  $\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1} = p_\theta^{l \Rightarrow l-1}$ 
10:      compute  $\tilde{\Theta}_t^{l-1}$  ▷ if  $l > 1$ 
11:    end for
12:  end for
13:  step 2: Compute approximate multilevel Monte Carlo estimate of objective
       $\mathbb{E}_{s,a,r \sim p_\theta} [J(s, a, r; \theta, \Theta)]$ :
14:     $\sum_{l=1}^{l=L} \widehat{\mathbb{E}}^{T_l} \left[ J(s^l, a^l, r^l; \theta, \Theta^l) - J(\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1}; \theta, \tilde{\Theta}^{l-1}) \right]$ ,
      where  $J(\tilde{s}^0, \tilde{a}^0, \tilde{r}^0; \theta, \tilde{\Theta}^0) \doteq 0$ .
15:  step 3: Update  $\theta$  using the gradient of estimated objective function
16:  end for

```

generating its synchronized rollouts at level $l - 1$ are expanded in Equation 4.

$$s^l, a^l, r^l \sim p_\theta^l \left\{ \begin{array}{l} s^l \sim \mu(s^l) \\ s^L = \psi_L^L(s^l) \\ a^L \sim \pi_\theta(a^L|s^L) \\ a^l = \phi_L^l(a^L) \\ s^{l-1} \sim \mathcal{P}^l(s^l, a^l) \\ r^l = \mathcal{R}^l(s^l, a^l, s^l) \end{array} \right. \quad \tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1} = p_\theta^{l \Rightarrow l-1} \left\{ \begin{array}{l} \tilde{s}^{l-1} = \psi_{l-1}^{l-1}(s^l) \\ \tilde{s}^L = \psi_{l-1}^L(\tilde{s}^{l-1}) \\ \tilde{a}^L \sim \pi_\theta(\tilde{a}^L|\tilde{s}^L) \\ \tilde{a}^{l-1} = \phi_{l-1}^L(\tilde{a}^L) \\ \tilde{s}^{l-1} \sim \mathcal{P}^{l-1}(\tilde{s}^{l-1}, \tilde{a}^{l-1}) \\ \tilde{r}^{l-1} = \mathcal{R}^{l-1}(\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{s}^{l-1}) \end{array} \right. \quad (4)$$

Note that since the target task corresponds to the level L , the policy is now represented as $\pi_\theta(a^L|s^L)$. Consequently, during policy rollouts on a certain level l , the state s^l passes through the mapping ψ_l^L and the action obtained a^L passes through the mapping ϕ_L^l . Synchronization from level l to $l - 1$ is obtained by mapping the states: $\tilde{s}^{l-1} = \psi_{l-1}^{l-1}(s^l)$. Figure 1 illustrates the implementations of a policy iteration in the classical and multilevel frameworks. Note that the level $l - 1$ changes to l at the end of steps T_{l-1} (for $l = 2, \dots, L$) and to continue the rollouts at the level l , the state is mapped as $\psi_{l-1}^l(\tilde{s}_{T_{l-1}}^{l-1})$. The generated samples are further used to compute the approximate multilevel Monte Carlo estimate of $\mathbb{E}_{s,a,r \sim p_\theta} [J(s, a, r; \theta, \Theta)]$ which is described as

$$\sum_{l=1}^{l=L} \widehat{\mathbb{E}}^{T_l} \left[J(s^l, a^l, r^l; \theta, \Theta^l) - J(\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1}; \theta, \tilde{\Theta}^{l-1}) \right],$$

where $J(\tilde{s}^0, \tilde{a}^0, \tilde{r}^0; \theta, \tilde{\Theta}^0) \doteq 0$. Since $T_1 > T_2 > \dots > T_L$, most of the computational costs of the rollouts lean towards sublevel environments. As a result, the approximate multilevel Monte Carlo estimate requires an overall lower computational cost than the Monte Carlo estimate in the classical framework. Finally, the network parameters θ are updated using the gradient of the estimated objective function at the end of the policy iteration.

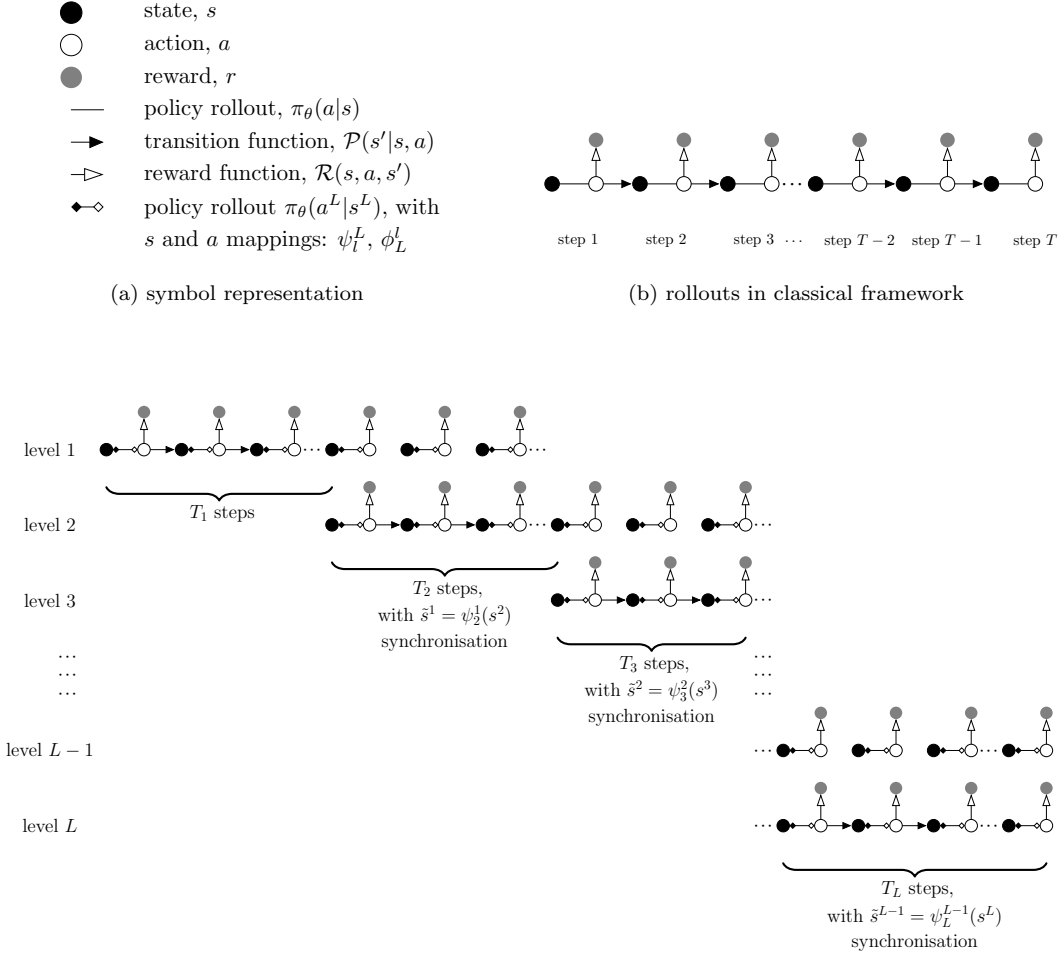


Figure 1: schematics of rollouts for a policy iteration

3.1 Multilevel PPO algorithm

We present the proposed multilevel framework for the state-of-the-art model-free algorithm, proximal policy optimization (PPO) (Schulman et al., 2017). In the context of the multilevel framework, the objective function is defined as

$$\begin{aligned}
J(s, a, r; \theta, \Theta^{ppo}) = & \min(p(\theta)A(s, a), \text{clip}(\mathbf{r}(\theta), 1 - \epsilon, 1 + \epsilon)A(s, a)) \\
& - c_v \left(r + \gamma \max_{s'} V_\theta(s') - V_\theta(s) \right)^2 \\
& + c_e S[\pi_\theta](s).
\end{aligned} \tag{5}$$

The first term of this objective function is called the surrogate policy term, where $p(\theta) = \pi_\theta(a|s)/\pi_{\theta_{old}}(a|s)$ and θ_{old} are the network parameters at the beginning of the policy iteration, $A(s, a)$ is the advantage function, which is estimated using the generalized advantage estimator (Schulman et al., 2015). The second term is referred to as value function error term which correspond to learning value function $V_\theta(s)$, where γ is the discount factor. Finally, the last term $S[\pi_\theta]$, corresponds to the entropy of the learned policy, which is added to ensure sufficient exploration. Parameters Θ^{ppo} refer to the set of the following parameters: $\theta_{old}, A(s, a), \epsilon, c_v, \gamma, V_\theta, s', c_e$ and $S[\pi_\theta](s)$. Readers are referred to the Appendix A for a detailed definition of these parameters. The approximate multilevel Monte Carlo estimate of $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta^{ppo})]$ is defined as

$$\sum_{l=1}^{l=L} \widehat{\mathbb{E}}^{M_l}_{s^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1} = p_\theta^{l \Rightarrow l-1}} \left[J(s^l, a^l, r^l; \theta, \Theta^{ppo^l}) - J(\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1}; \theta, \tilde{\Theta}^{ppo^{l-1}}) \right], \tag{6}$$

where $J(\tilde{s}^0, \tilde{a}^0, \tilde{r}^0; \theta, \tilde{\Theta}^{ppo^0}) \doteq 0$ and M_l is the mini-batch size at level l . The algorithm 3 provides an outline for the multilevel PPO algorithm. The inputs are the same as those of the classical PPO

Algorithm 3 Multilevel Proximal Policy Optimization algorithm

```

1: Input:  $\mathcal{E} = \{\mathcal{E}_1, \dots, \mathcal{E}_L\}, N, \mathbf{T} = \{T_1, \dots, T_L\}, \mathbf{M} = \{M_1, \dots, M_L\}, K$ 
2: for iteration = 1, 2, ... do
3:   for actor = 1, 2, ...,  $N$  do
4:     for level  $l = 1, 2, \dots, L$  do
5:        $s^l = \psi_{l-1}^l(\tilde{s}_{T_{l-1}}^{l-1})$  ▷ if  $l > 1$ 
6:       for  $t = 1, 2, \dots, T_l$  do
7:          $s^l, a^l, r^l \sim p_\theta^l$ 
8:         compute  $\Theta^{ppo^l}$ 
9:          $\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1} = p_\theta^{l \Rightarrow l-1}$  ▷ if  $l > 1$ 
10:        compute  $\tilde{\Theta}^{ppo^{l-1}}$  ▷ if  $l > 1$ 
11:      end for
12:    end for
13:  end for
14:  gather data  $\{(s^l, a^l, r^l), (\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1})\}_{t=1}^{t=NT_l}\}_{l=1}^{l=L}$ , from all actors
15:  optimize equation 6, with  $K$  epochs and minibatch size  $M_l \leq NT_l$ 
16:  update policy network parameters  $\theta$ 
17: end for

```

algorithm, except that multilevel variables are provided as a set of length L : environments at each

level $\mathcal{E} = \{\mathcal{E}^1, \dots, \mathcal{E}^L\}$, number of actors N , number of steps at each level $\mathbf{T} = \{T^1, \dots, T^L\}$, number of batches at each level $\mathbf{M} = \{M^1, \dots, M^L\}$ (such that $NT^l \leq M^l$ and $T^1/M^1 = \dots = T^L/M^L$) and number of epochs K . Note that if the sets \mathcal{E} , \mathbf{T} and \mathbf{M} consist of a single value, this algorithm is the same as the classical PPO algorithm where the objective function is estimated using the Monte Carlo method. We implement this algorithm using a standard RL library, stable baselines 3 (Raffin et al., 2021). The implementation details are delineated in Appendix 5.

3.2 Multilevel PPO analysis methodology

We present an analysis methodology to compare the Monte Carlo estimate and the standard multilevel Monte Carlo estimate of the PPO objective function. The analysis methodology is adopted from Giles (2015), where the strong and weak convergences of the estimates are checked for predefined mean squared error values. For convenience of demonstration, let us consider the following shorthand notation.

$$\begin{aligned} \widehat{\mathbb{E}}^N [Y_l] &\doteq \widehat{\mathbb{E}}^N \left[J(s^l, a^l, r^l; \theta, \Theta^{ppo^l}) - J(\bar{s}^{l-1}, \bar{a}^{l-1}, \bar{r}^{l-1}; \theta, \bar{\Theta}^{ppo^{l-1}}) \right], \\ &\quad \substack{s^l, a^l, r^l \sim p_\theta^l \\ \bar{s}^{l-1}, \bar{a}^{l-1}, \bar{r}^{l-1} = p_\theta^{l \Rightarrow l-1}} \\ \widehat{\mathbb{E}}^N [J_l] &\doteq \widehat{\mathbb{E}}^N \left[J(s^l, a^l, r^l; \theta, \Theta^{ppo^l}) \right]. \end{aligned}$$

As a result, the multilevel Monte Carlo estimate of $\mathbb{E}[J_L]$ is described as

$$Y = \sum_{l=1}^L \widehat{\mathbb{E}}^{M_l} [Y_l]. \quad (7)$$

The mean squared error (*MSE*) for this estimator is defined as

$$\begin{aligned} MSE &= \mathbb{E}[(Y - \mathbb{E}[J_L])^2] \\ &= \mathbb{V}[Y] + (\mathbb{E}[Y] - \mathbb{E}[J_L])^2, \end{aligned}$$

where $\mathbb{V}[Y]$ is the variance of the estimator and $(\mathbb{E}[Y] - \mathbb{E}[J_L])^2$ corresponds to the bias of the estimator. A sufficient condition on $MSE \leq \varepsilon^2$, expands to $\mathbb{V}[Y] = \varepsilon^2/2$ and $(\mathbb{E}[Y] - \mathbb{E}[J_L])^2 \leq \varepsilon^2/2$. Under assumption $\mathbb{V}[Y] = \varepsilon^2/2$, the optimal number of samples at each level M_l and the corresponding total cost of the estimator C_{MLMC} are calculated as

$$M_l = 2\varepsilon^{-2} \left(\sum_{l=1}^L V_l C_l \right) \sqrt{\frac{V_l}{C_l}}, \quad (8)$$

$$C_{\text{MLMC}} = 2\varepsilon^{-2} \left(\sum_{l=1}^L V_l C_l \right) \sqrt{V_l C_l}, \quad (9)$$

where V_l corresponds to the variance estimate $\widehat{\mathbb{V}}^{N_\infty} [Y_l]$ (defined as $\widehat{\mathbb{E}}^{N_\infty} [Y_l^2] - \widehat{\mathbb{E}}^{N_\infty} [Y_l]^2$) for a large number N_∞ , of samples and C_l is the computational cost of each sample of Y_l . The weak convergence test $(\mathbb{E}[Y] - \mathbb{E}[J_L])^2 \leq \varepsilon^2/2$, is ensured by the following inequality:

$$\frac{\max_{l \in \{L-2, L-1, L\}} \widehat{\mathbb{E}}^{N_\infty} [Y_l]}{(2^\alpha - 1)} \leq \frac{\varepsilon}{\sqrt{2}}, \quad (10)$$

where α is assumed to be a positive coefficient that explains the decay in the values of $\widehat{\mathbb{E}}^{N_\infty} [Y_l]$ for the chosen levels in the form $\widehat{\mathbb{E}}^{N_\infty} [Y_l] = c_1 2^{-\alpha l}$. It is estimated using linear regression on $\widehat{\mathbb{E}}^{N_\infty} [Y_l]$ values.

Furthermore, the multilevel estimator Y is compared with the Monte Carlo estimate corresponding to the highest level environment \mathcal{E}^L which is computed as

$$Y_{\text{MC}} = \widehat{\mathbb{E}}^M [J_L]. \quad (11)$$

The number of samples M and the total cost C_{MC} of the Monte Carlo estimate corresponding to the variance of the estimate $\varepsilon^2/2$ are calculated as

$$M = 2\varepsilon^{-2} \frac{V}{C}, \quad (12)$$

$$C_{\text{MC}} = 2\varepsilon^{-2} V \quad (13)$$

where V is the variance estimate $\widehat{V}^{N_\infty} [J_L]$ and C is the computational cost of each sample of J_L .

The multilevel PPO analysis is performed in parallel with learning at certain predefined intervals of policy iterations. An outline of the analysis is presented in algorithm 4. To obtain accurate

Algorithm 4 Analysis of multilevel Proximal Policy Optimization algorithm

- 1: Input: $\mathcal{E} = \{\mathcal{E}_1, \dots, \mathcal{E}_L\}$, $\varepsilon = \{\varepsilon_1, \dots, \varepsilon_n\}$, $\{C_1, \dots, C_L\}$, N_∞
 - 2: generate samples $\{s^L, a^L, r^L\}_{t=1}^{t=N_\infty}$ using the environment \mathcal{E}^L (i.e. $s^L, a^L, r^L \sim p_\theta^L$)
 - 3: generate synchronized samples $\{\{\tilde{s}^l, \tilde{a}^l, \tilde{r}^l\}_{t=1}^{t=N_\infty}\}_{l=L-1}^{l=1}$ on sublevels (where $\tilde{s}^l, \tilde{a}^l, \tilde{r}^l = p_\theta^{L \Rightarrow l}$)
 - 4: compute $\widehat{\mathbb{E}}^{N_\infty} [Y_l]$, $\widehat{\mathbb{E}}^{N_\infty} [J_l]$, $\widehat{V}^{N_\infty} [Y_l]$ and $\widehat{V}^{N_\infty} [J_l]$ using generated data
 - 5: **for** ε in ε **do**
 - 6: compute M_l (equation 8)
 - 7: estimate multilevel Monte Carlo estimate Y (equation 7)
 - 8: compute total cost C_{MLMC} (equation 9)
 - 9: compute M (equation 12)
 - 10: estimate Monte Carlo estimate Y_{MC} (equation 11)
 - 11: compute total cost C_{MC} (equation 13)
 - 12: check weak convergence (equation 10)
 - 13: **end for**
-

estimates of $\widehat{\mathbb{E}}^{N_\infty} [Y_l]$, $\widehat{\mathbb{E}}^{N_\infty} [J_l]$, $\widehat{V}^{N_\infty} [Y_l]$ and $\widehat{V}^{N_\infty} [J_l]$ a high number of samples N_∞ , is chosen. The samples of sequences are rolled out on the finest level L (corresponding to random variable $s^L, a^L, r^L \sim p_\theta^L$) and its synchronized samples are created in parallel on sublevels $l \in \{1, \dots, L-1\}$ (corresponding to random variables $\tilde{s}^l, \tilde{a}^l, \tilde{r}^l = p_\theta^{L \Rightarrow l}$). The notations $s^L, a^L, r^L \sim p_\theta^L$ and $\tilde{s}^l, \tilde{a}^l, \tilde{r}^l = p_\theta^{L \Rightarrow l}$ are delineated in equation 14 as

$$s^L, a^L, r^L \sim p_\theta^L \begin{cases} s^L \sim \mu(s^L) \\ a^L \sim \pi_\theta(a^L | s^L) \\ s'^L \sim \mathcal{P}^L(s^L, a^L) \\ r^L = \mathcal{R}^L(s^L, a^L, s'^L), \end{cases} \quad \tilde{s}^l, \tilde{a}^l, \tilde{r}^l = p_\theta^{L \Rightarrow l} \begin{cases} \tilde{s}^l = \psi_L^l(s^L) \\ \tilde{a}^l \sim \pi_\theta(a^L | s^L) \\ \tilde{a}^l = \phi_L^l(\tilde{a}^L) \\ \tilde{s}^l \sim \mathcal{P}^l(\tilde{s}^l, \tilde{a}^l) \\ \tilde{r}^l = \mathcal{R}^l(\tilde{s}^l, \tilde{a}^l, \tilde{s}^l). \end{cases} \quad (14)$$

The Monte Carlo and multilevel Monte Carlo estimates of the objective function of PPO are computed and compared for a set $\varepsilon = \{\varepsilon_1, \dots, \varepsilon_n\}$, of MSE accuracy values. The computational effectiveness of the multilevel estimator is demonstrated by comparing its total cost C_{MLMC} with the corresponding total cost for the Monte Carlo estimate C_{MC} for each accuracy value.

4 Experiments

We present two case studies of simulation environments in which the transition between states is governed by the solution of two partial differential equations that describe the incompressible flow of a single phase through porous medium. The stochasticity of the environments is attributed to an uncertain field of permeability. The governing equations for a single phase flow c of clean water, through a porous medium with porosity η , consist of the continuity equation coupled with the incompressibility condition that are defined as

$$\eta \frac{dc}{dt} = cq - \nabla \cdot cv; \quad \nabla \cdot v = q \quad \text{in } \Omega \subset \mathbb{R}^2. \quad (15)$$

Flow velocity v and pressure p are related by Darcy's law: $v = -k/\mu \nabla p$, where k is permeability and μ is viscosity. Permeability is treated as a stochastic parameter, and its uncertainty is modeled with a predefined probability distribution. The source and sink are denoted by q , where the source corresponds to the injection rate of uncontaminated fluid (clean water) in the domain Ω , and the sink corresponds to the flow rate of the contaminated fluid at the outlet.

Two environments with distinct parameters and flow scenarios are designed for demonstration of the proposed multilevel PPO algorithm. For both cases, the parameter values emulate those of the benchmark reservoir simulations presented in SPE-10 model 2 (Christie et al., 2001). Environments are denoted ResSim-v1 and ResSim-v2 in the rest of the paper (ResSim is a shorthand term for reservoir simulation).

4.1 ResSim-v1 parameters

Schematics of the domain Ω in ResSim-v1 are illustrated in Figure 2a. Viscosity μ is set to 0.3 cP, while porosity η is set to a constant value of 0.2. According to the convention in geostatistics, the distribution of logarithmic permeability $g = \log(k)$ is assumed to be known. This logarithmic permeability distribution for test case 1 is inspired by the case study conducted by Brouwer et al. (2001). In total, 32 injection locations (illustrated with blue circles) and 32 outlet locations (illustrated with red circles) are placed on the left and right edges of the domain, respectively. The total injection rate is set to a constant value of 2304 ft²/day. As illustrated in Figure 2a, a linear high-permeability channel (shown in gray) passes from the left to the right side of the domain. l_1 and l_2 represent the distance from the top edge of the domain on the left and right sides, while the width of the channel is indicated by w . These parameters follow uniform distributions defined as $w \sim U(120, 360)$, $l_1 \sim U(0, L - w)$ and $l_2 \sim U(0, L - w)$, where L is the domain length. To be specific, the logarithmic permeability g at a location (x, y) is formulated as follows:

$$g(x, y) = \begin{cases} \log(245) & \text{if } \frac{l_2 - l_1}{L}x + l_1 \leq y \leq \frac{l_2 - l_1}{L}x + l_1 + w, \\ \log(0.14) & \text{otherwise,} \end{cases}$$

where x and y are horizontal and vertical distances from the upper left corner of the domain, as illustrated in Figure 2a. The values for permeability in the channel (245 mD) and the rest of the domain (0.14 mD) are inspired from Upperness log-permeability distribution peak values specified in SPE-10 model 2 case.

4.2 ResSim-v2 parameters

Figure 2b shows the reservoir domain for ResSim-v2. It consists of 14 outlets (illustrated with red circles) located symmetrically on the left and right edges (7 on each edge) of the domain and 7

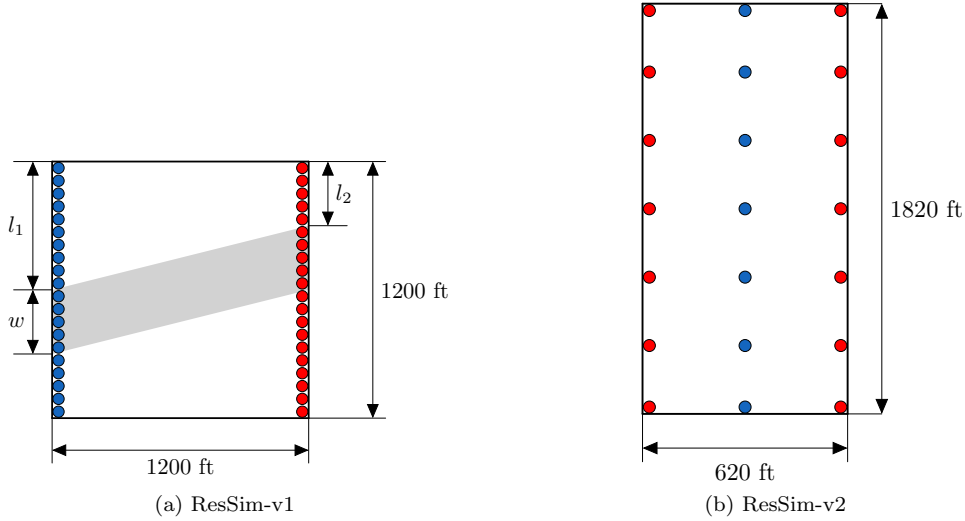


Figure 2: schematic of the spatial domain Ω

injections (illustrated with blue circles) located at the central vertical axis of the domain. The total injection rate is set at a constant value of $9072 \text{ ft}^2/\text{day}$ while viscosity and porosity are set to the same values as in ResSim-v1. The uncertainty distribution of the permeability field is considered to be smoother and spatially correlated, and is modeled as a constrained log-normal distribution. Logarithmic permeability samples are created using ordinary kriging methodology, which are constrained with a constant value of 2.41 logarithmic permeability at injection and outlet locations. The exponential variogram model used for the kriging is defined as

$$\gamma(r) = \sigma^2 \left(1 - \exp \left(-\sqrt{\left(\frac{r_x}{l_x}\right)^2 + \left(\frac{r_y}{l_y}\right)^2} \right) \right)$$

where r_x and r_y are x and y projections of the distance r . The variance of the process σ is set to 5, while the length scales l_x and l_y are set to 620 ft (width of the domain) and 62 ft (10% of domain width), respectively. The samples of permeability fields are further rotated clockwise with the angle $\pi/8$. In this study, the above-mentioned kriging process is performed using the geostatistics library `gstools` (Müller and Schüler, 2019).

4.3 Reinforcement learning task

In the context of reinforcement learning, the state s is represented by a set of variables $\{c, k, \eta, \mu\}$, while the action a is represented by the source/sink term q at each control step. We employ finite-volume discretization of governing equations 15 as detailed in Aarnes et al. (2007) which is treated as a transition function \mathcal{P} between states at time t_m and t_{m+1} . The total time of the simulation is divided into five control steps, which form an episode with finite horizon. As a result, the task is to learn a policy $\pi_\theta(a|s)$ that selects the optimal values of q that maximize the cumulative reward

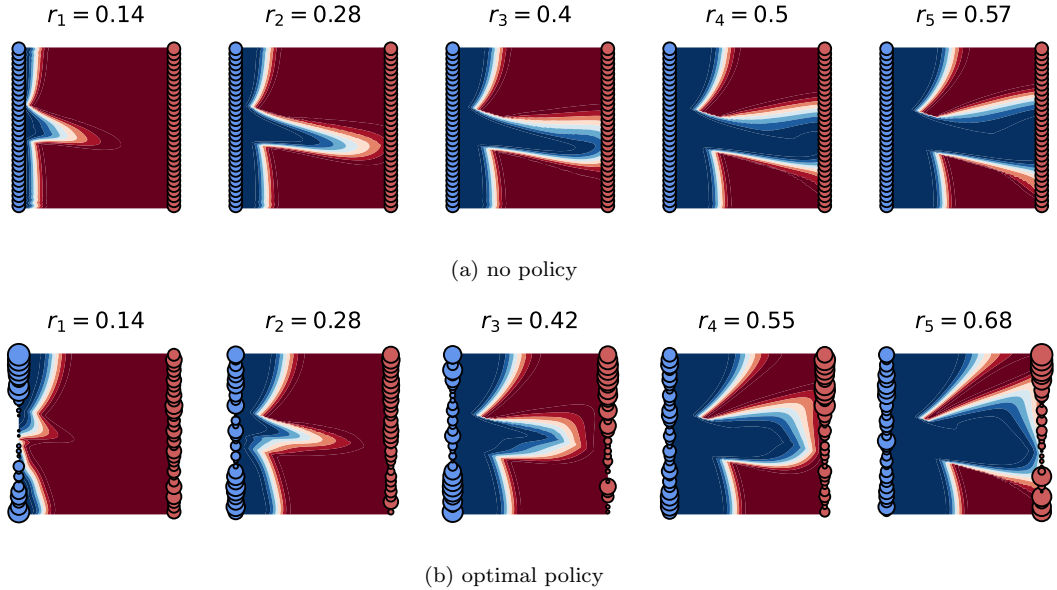


Figure 3: example policy visualization for ResSim-v1

defined as

$$\sum_{m=1}^5 \frac{1}{\phi|\Omega|} \int_{t_m}^{t_{m-1}} \left(\int_{\Omega} \min(q, 0)(1 - c) d\Omega \right) dt, \quad (16)$$

where $|\Omega|$ refers to the area of the domain. This cumulative reward refers to the sweep efficiency of the injected clean water, which ranges from 0 to 1. Furthermore, in the context of temporal difference learning, the reward at time t_m is formulated as the term inside the summation operator of equation 16. To represent the stochasticity of the task, a random sample of permeability is chosen from a finite set of permeabilities for each episode in the learning process. This finite set of permeability samples is achieved with a cluster analysis (please refer to Appendix D for the cluster analysis formulation used in this paper). In order to demonstrate application for a partially observable system, the policy network input is replaced with an observation vector instead of the above-defined state. Here, the observation vector corresponds to values of concentration and fluid pressure at injection and outlet locations. Subsequently, the output of the policy network corresponds to the control vector, which consists of weights (with values ranging between 0.001 to 1) representing flow rates at injection and outlet locations. Note that with such representation of states, the underlying assumption of the Markov property of the transition function is approximated. Such a system is referred to as a partially observable Markov decision process (POMDP). By the definition of POMDP (Spaan, 2012), the policy requires observations and actions from some sort of history or memory of previous control steps to return the action for a certain control step. However, for the case studies presented, observation from only the previous control step is sufficient for policy representation.

Figure 3 illustrates visualization of flow through the domain in ResSim-v1. The injection and outlet locations are indicated with circles in blue and red, respectively, and their radius is proportional to the flow rate. When ResSim-v1 is operated without a policy (that is, constant injection/outlet rate in all locations), most of the concentration flow takes place in the high-permeability channel, causing poor sweep efficiency in the low-permeability region. Figure 3a illustrates the flow scenario

without a policy for a sample of permeability. The concentration flow is highlighted with blue in the domain. Consequently, the reward which refers to the sweep efficiency corresponds to the ratio of domain area highlighted in blue color to the total domain area. In other words, optimal policy refers to the choice of actions that increase the domain area in blue (i.e., swept area of the contaminate where the concentration of the clean water is high). Figure 3b illustrates the optimal policy in which the flow through injection / outlet locations near the high permeability area near the channel is restricted. As indicated by cumulative rewards at each time, we observe an improvement in total reward at the end of the episode (from 0.57 in no policy to 0.68 with optimal policy). Similarly, Figure 4 provides a visualization for the ResSim-v2 environment. The optimal policy in this case is to improve the flow rate at locations near the low-permeability locations while restricting the flow rate in locations near the high-permeability region.

4.4 Multilevel framework formulation

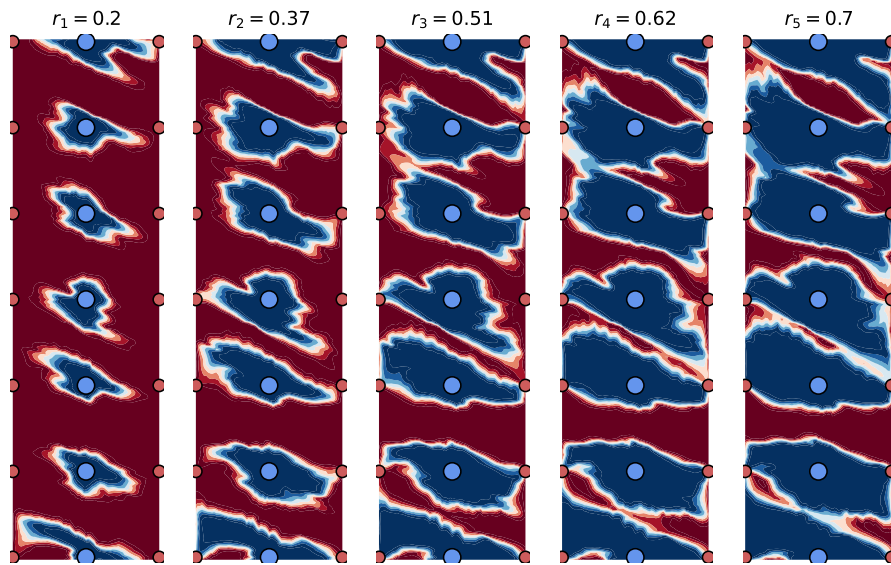
For multilevel formulation, we consider three levels of environment for the ResSim-v1 environment, where the target task is described by the environment on level 3. For ResSim-v2 environment, we consider two levels, where the target task is predefined to be at level 2. The levels for these environments correspond to the grid fidelity of the discretization scheme used to solve the governing equations. Table 1 delineates the grid sizes corresponding to each level in the ResSim-v1 and ResSim-v2 environments. The choice of grid fidelity corresponds to the fact that computational

Table 1: grid size on each level

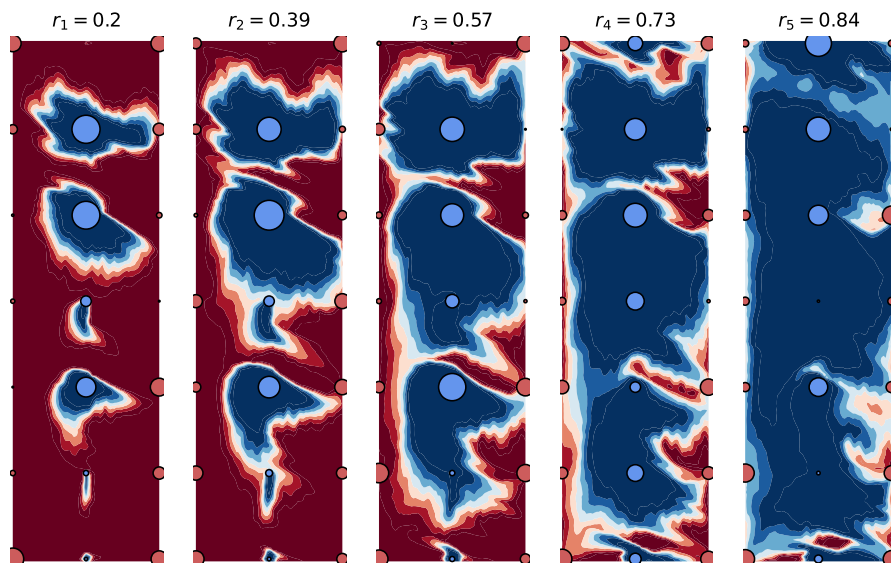
	ResSim-v1	ResSim-v2
level 1	32×32	31×111
level 2	64×64	73×219
level 3	128×128	–

cost and accuracy of model dynamics are proportional to the level. This is due to the fact that the computational cost and accuracy of a PDE are often proportional to the size of the grid. These levels of environment are chosen heuristically for demonstration purpose of the multilevel PPO algorithm. Although there could be a more systematic approach to choosing these levels, we consider this to be outside the scope of this study. The state mapping function $\psi_i^{l'}$ maps the state $\{c^l, k^l, \eta^l, \mu^l\}$ for the environment on level l to the state $\{c^{l'}, k^{l'}, \eta^{l'}, \mu^{l'}\}$ for the environment on level l' . Since porosity η and viscosity μ are set to a constant throughout the domain, we do not need to map them in the function $\psi_i^{l'}$. As a result, $\psi_i^{l'}$ only maps the concentration c and the permeability k between the level l and l' . When l is larger than l' , the mapping occurs from a fine grid to a coarser grid. This is done by super-positioning a fine grid on a coarse grid and creating coarse partitions on the fine grid. The resulting values in each partition are passed through the mean function for concentration values and the harmonic mean for permeability values. On the contrary, when l' is larger than l , the mapping occurs from coarse grid to fine grid. In this case, the coarse value in each partition is simply assigned to fine grid cells in the corresponding partition. When it comes to the action mapping function $\phi_i^{l'}$, note that l' is always larger than l for the proposed multilevel framework. As a result, the action q is always mapped from a coarse grid to a finer one. The coarse to fine mapping is done with the same methodology as $\psi_i^{l'}$, except for the choice of mapping function, which is the sum of the action q . Finally, when l is the same as l' , the mapping functions $\psi_i^{l'}$ and $\phi_i^{l'}$ act as an identity function.

Figure 5a illustrates the comparison between flow through the domain at various levels. Com-



(a) no policy



(b) optimal policy

Figure 4: example of policy visualization for ResSim-v2

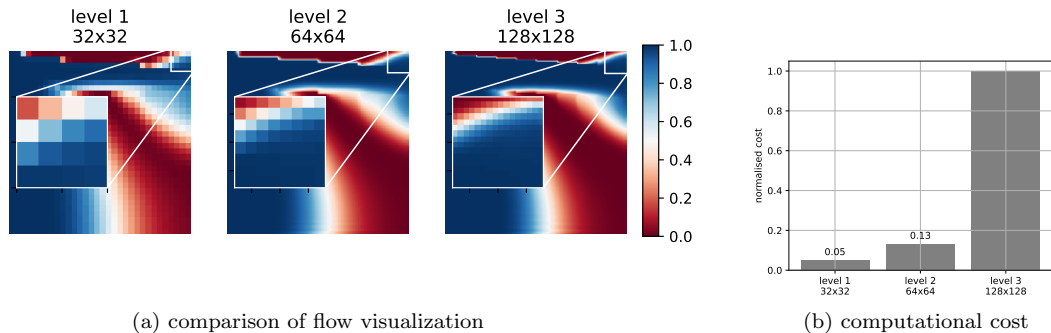


Figure 5: environment levels for ResSim-v1

parison of the computational cost of the transition function for different levels is illustrated with a bar plot in Figure 5b. This computational cost is taken as an average value of 100 simulation trials to account for variability. The computational cost at each level is normalized by dividing it by that corresponding to the target task. Similar plots for the visualization of two levels of ResSim-v2 are shown in Figure 6.

5 Results

We demonstrate the effectiveness of the multilevel PPO algorithm by comparing its results with the results of the classical single-level PPO algorithm for the target task. Table 2 delineates the levels of environments considered in the one-level, two-level, and three-level PPO algorithm (denoted as PPO-1L, PPO-2L, and PPO-3L, respectively). Note that PPO-1L refers to the results of the classical

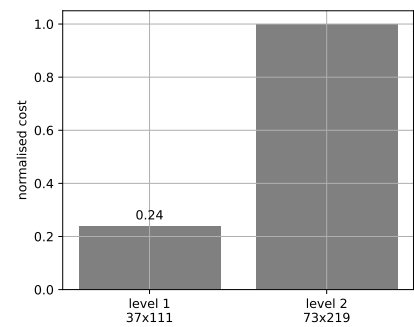
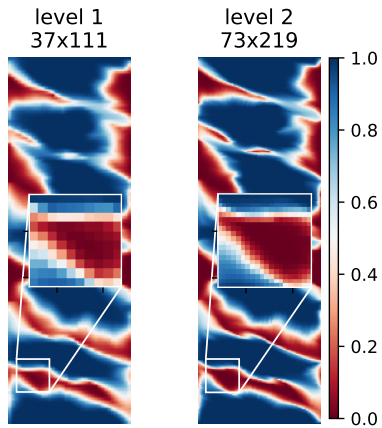
Table 2: levels in each multilevel PPO experiment

	ResSim-v1	ResSim-v2
PPO-1L	{3}	{2}
PPO-2L	{2, 3}	{1, 2}
PPO-3L	{1, 2, 3}	–

single-level PPO algorithm for the target task.

5.1 ResSim-v1 results

First, we present the results for multilevel PPO analysis with PPO-1L, which consists of a total of 300 policy iterations. The analysis is performed every 30 iterations. Figure 7 illustrates the comparison of Monte Carlo and the three-level Monte Carlo estimate of the objective function with a true value that is estimated using 10^5 samples (that is, N_∞ is set to 10^5). The analysis is performed for three values of RMS accuracy: 10^{-2} , 10^{-3} , 10^{-4} . This is done by setting $\varepsilon = \{\sqrt{10^{-2}}, \sqrt{10^{-3}}, \sqrt{10^{-4}}\}$ in the analysis. The cost terms $\{C_1, C_2, C_3\}$, which correspond to the computational cost of each term in the multilevel Monte Carlo estimate, are set to $\{0.1, 0.33, 1.23\}$. These values are chosen from the computational cost on each level, which are illustrated in Figure 5b. To be specific, C_1 refers



(a) comparison of flow visualization

(b) computational cost

Figure 6: environment levels for ResSim-v2

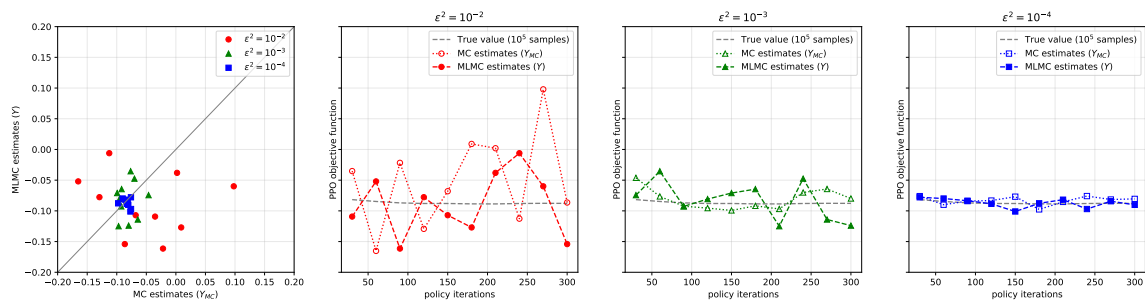


Figure 7: comparison of Monte Carlo and multilevel Monte Carlo estimate of PPO objective function for ResSim-v1

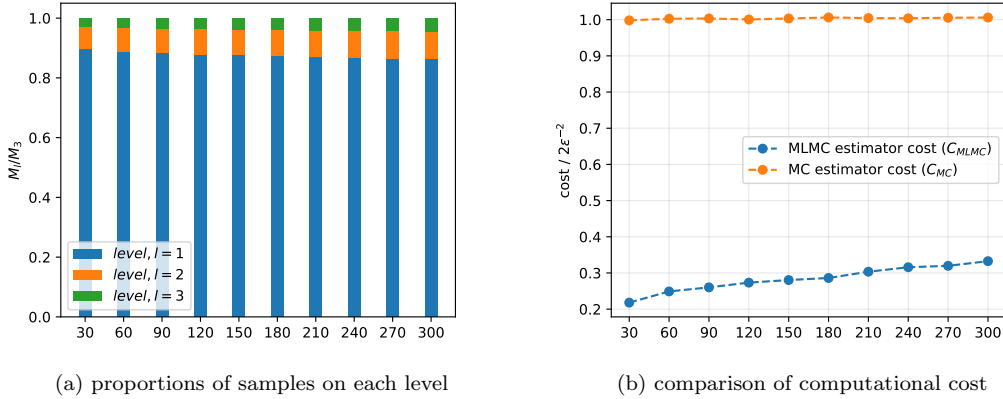


Figure 8: MLMC analysis results for ResSim-v1

to the computational cost on level 1 (that is, 0.1). The term C_2 refers to the computational cost of the difference between synchronized samples at levels 1 and 2, as a result C_2 is set as the sum of the computational cost at levels 1 and 2 (i.e., $0.1+0.23$). Similarly, C_3 is set as the sum of the computational cost at levels 2 and 3 (that is, $0.23 + 1.0$). As can be seen in figure 7, we see a fairly accurate comparison between Monte Carlo and multilevel Monte Carlo estimates. Furthermore, these estimates yield more accurate values as we move towards lower values of ϵ^2 . This is because the number of samples is inversely related to ϵ^2 (as stated in equations 8 and 12). As a result, the number of samples is basically scaled up as we reduce the values of ϵ^2 . Figure 8a illustrates the number of optimal samples at each level of the multilevel estimator. The numbers of samples are normalized to show the proportions of the samples at each level. This is done by dividing M_l (from Equation 8) by M_3 for all $l \in \{1, 2, 3\}$. We see that M_2 is approximately $1/10^{th}$ of M_1 and M_3 is observed to be about half of M_2 throughout the learning process. The comparison between the computational cost of Monte Carlo and the multilevel Monte Carlo estimate is plotted in Figure 8 b. Here, the computational cost terms C_{MC} and C_{MLMC} are divided by $2\epsilon^{-2}$ to obtain the constant cost terms irrespective of RMS accuracy. We observe that the computational cost of the multilevel estimate takes only about 20 to 30% of the Monte Carlo estimate from the analysis.

Figure 9 a shows the superposition learning processes for PPO-1L, PPO-2L, and PPO-3L. The parameters used for the experiments PPO-1L, PPO-2L, and PPO-3L are delineated in the table 3. The learning plots are drawn as the average along with the range of values for three distinct seed

Table 3: parameters of multilevel PPO experiment for ResSim-v1

	T	M	N	K
PPO-1L	{50}	{250}	50	20
PPO-2L	{70, 5}	{350, 25}	50	20
PPO-3L	{80, 10, 5}	{400, 50, 25}	50	20

values. The parameter M , for PPO-1L, PPO-2L, and PPO-3L, is calculated from equation 8 for the RMS value $\epsilon^2 = 7.8 \times 10^{-3}$ where the values of V_l and C_l are taken from the analysis mentioned above. In other words, we compare the results among PPO-1L, PPO-2L, and PPO-3L for a constant RMS

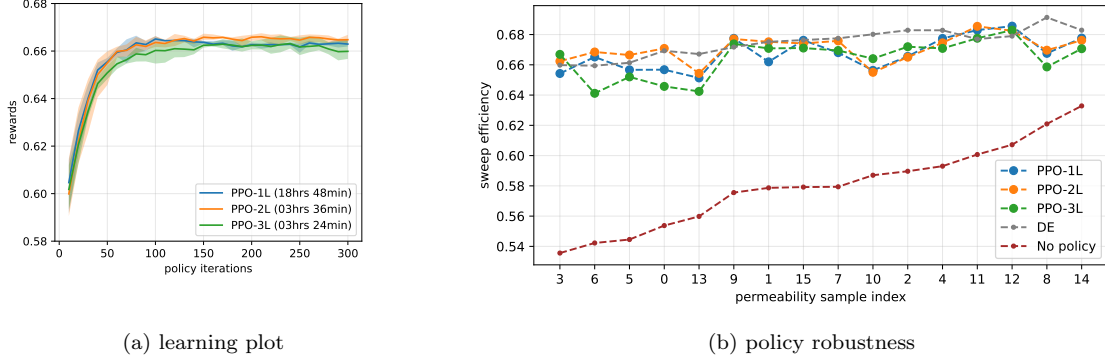


Figure 9: multilevel PPO results for ResSim-v1

accuracy. Note that these choices of values are done only in order to demonstrate a fair comparison among PPO-1L, PPO-2L, and PPO-3L. In practice, it is not required to perform the analysis in order to choose M . Other parameters of the algorithm are tuned to find the convergence for the PPO-3L case first, and these same parameters were used in the PPO-2L and PPO-1L cases. Figure 9 a shows the evaluation of the environment policy corresponding to the target task. This policy evaluation is represented with the average reward corresponding to all the permeability samples used in the learning process. PPO-1L refers to the classical PPO algorithm, which takes around 19 wall clock hours, while PPO-2L and PPO-3L which correspond to the proposed multilevel PPO algorithm achieve the same learning in about three and half hours. In other words, we save around 82% computational costs with the proposed algorithm compared to its classical counterpart. Figure 9b shows the robustness of the learned policies against uncertainty in permeability. This is done by plotting rewards for 16 random permeability samples of the uncertainty distribution that were unseen during the learning process. These results were compared with the optimal solutions obtained using the differential evolution algorithm (implemented using the SciPy library, (Virtanen et al., 2020)), which are denoted DE in Figure 9 b. The algorithm parameters for the PPO and DE algorithms, in this study, are delineated in the Appendix E.

5.2 ResSim-v2 results

Similarly to ResSim-v1, we perform multilevel PPO analysis with PPO-1L, which consists of a total of 1200 policy iterations and is performed every 120 iterations. Figure 10 illustrates the correlation between Monte Carlo and the multilevel Monte Carlo estimate of the objective function. The RMS values in ϵ are set to $\{\sqrt{10^{-2}}, \sqrt{10^{-3}}, \sqrt{10^{-4}}\}$ in the analysis. The cost terms $\{C_1, C_2\}$, which correspond to the computational cost of each term in the multilevel Monte Carlo estimate, are set to $\{0.24, 1.24\}$. As illustrated in Figure 10, we observe a higher correlation between Monte Carlo and multilevel Monte Carlo estimates for higher RMS accuracy values. Figure 11 a illustrates the proportions of the optimal number of samples at both levels of the multilevel estimator. The comparison between the computational cost of Monte Carlo and the multilevel Monte Carlo estimate is plotted in figure 11b. We observe that the computational cost of the multilevel estimate takes only around 25 to 35% of the Monte Carlo estimate from the analysis.

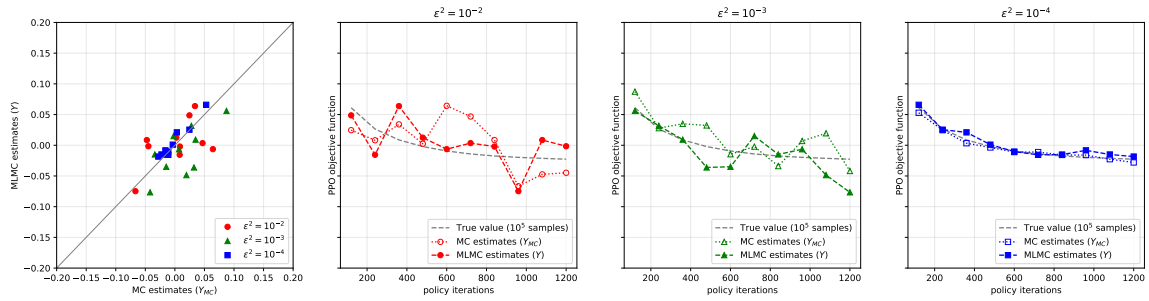


Figure 10: comparison of Monte Carlo and multilevel Monte Carlo estimate of PPO objective function for ResSim-v2

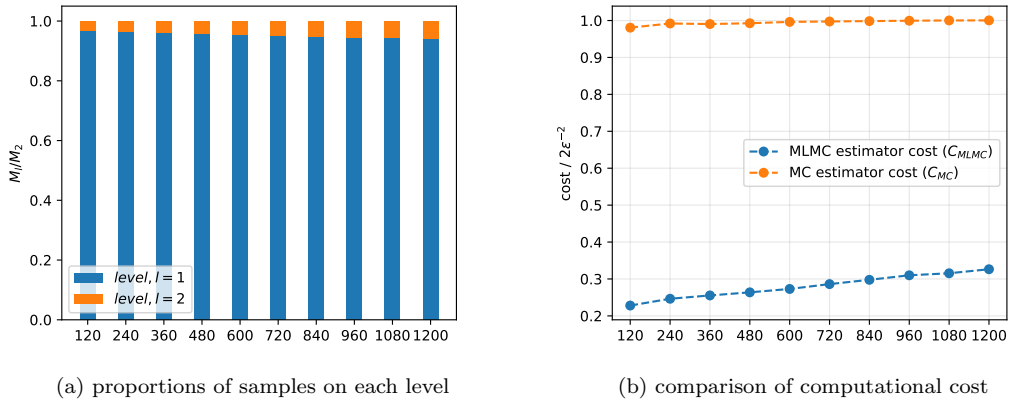


Figure 11: MLMC analysis results for ResSim-v2

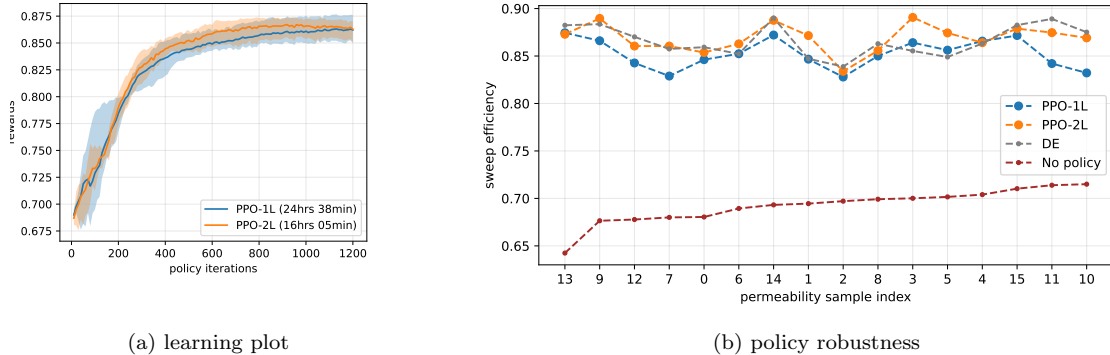


Figure 12: multilevel PPO results for ResSim-v2

Figure 12 a shows the comparison between the learning processes for PPO-1L and PPO-2L. The parameters used for the PPO-1L and PPO-2L experiments are delineated in the table 4. In this case,

Table 4: parameters of multilevel PPO experiment for ResSim-v1

	T	M	N	K
PPO-1L	{100}	{500}	50	20
PPO-2L	{140, 15}	{700, 75}	50	20

the comparison between PPO-1L and PPO-2L is made for the RMS accuracy value $\epsilon^2 = 3.9 \times 10^{-3}$. Similarly to the ResSim-v1 case, the hyperparameters of the algorithm are tuned to find convergence for the PPO-2L case, and the same parameters were used in the PPO-1L case. Figure 9 a shows the average evaluation of the environment policy corresponding to the target task (level 2). PPO-1L refers to the classical PPO algorithm which takes around 24 wall clock hours, while PPO-2L, which corresponds to the proposed multilevel PPO algorithm, achieves the same learning in about 16 hours. In other words, we save around 35% computational costs with the proposed algorithm compared to its classical counterpart. Figure 12b shows the robustness of the learned policies against uncertainty in permeability.

5.3 challenges and further research direction

Albeit successful results in learning and analysis of the proposed framework, we believe that this study deserves deeper mathematical investigation and analysis. In particular, we would like to study and analyze the effect of the approximation introduced in the MLMC estimation. As an introduction to the proposed framework, the experiments presented for multilevel PPO were performed for a specific PDE-based control problem for flow through porous media. Subsequently, we would like to provide a thorough study with a variety of experiments with the proposed multilevel PPO algorithm. This study would be mainly aimed at general benchmark problems in which RL is utilized to achieve superhuman controls, but with excessive computational costs. Furthermore, while tuning the algorithm parameters for the multilevel PPO algorithm, we observed that increasing the number of levels, the learning rate, and the clip range had an adverse effect on the learning convergence.

6 Conclusions

A multilevel framework for deep reinforcement learning is introduced in which the learned agent interacts with multiple levels of PDE-based environments where the level corresponds to the grid fidelity of PDE discretization. We present a mathematical framework that allows the synchronized implementation of task trajectories at multiple environmental levels. The presented approximate MLMC estimate is at the heart of the proposed multilevel framework. We also present a novel multilevel variant of the classical PPO algorithm based on the proposed multilevel framework. The computational efficiency of this multilevel PPO algorithm is illustrated for two environments for which model dynamics is represented with PDEs describing an incompressible single-phase fluid flow through a porous medium. We observe substantial computational savings in the case studies presented (approximately 82% and 35%, respectively).

As a future scope of this study, we aim to analyze the effect of the presented approximation to standard MLMC estimation. For the multilevel PPO algorithm, this can be done by extending the analysis methodology (presented in Section 3.2) for the approximate MLMC estimate. We also aim to provide a future study to benchmark multilevel PPO algorithm performance on a variety of environments.

A Examples of objective functions for different deep RL algorithms

Examples of the objective function $\mathbb{E}_{s,a,r \sim p_\theta} [J(s, a, r; \theta, \Theta)]$ for various deep reinforcement learning algorithms are delineated in table 5. In a value-based algorithm, such as the deep Q network (DQN), the neural network represents a function approximator for the Q function. Q function represents the expected return when the agent takes action a_t in state s_t and is defined as $Q(s, a) = \mathbb{E}_\pi [\sum_m \gamma^m r_{m+t+1} | s_t = s, a_t = a]$ where $\gamma \in [0, 1]$ is the discount factor and $\mathbb{E}_\pi[\dots]$ denotes the expected value given that the agent follows the policy π . The policy refers to taking the action corresponding to the highest Q value. In policy based algorithms like advantage actor-critic (A2C), trust region policy optimization (TRPO) and proximal policy optimization (PPO). The policy is directly modeled as a neural network that maps the state s to the corresponding optimal action a . This network is often integrated with a value network, which maps the state s to its corresponding value $V(s)$. The value function is the expected future return for a particular state s_t and is defined as $V(s) = \mathbb{E}_\pi [\sum_m \gamma^m r_{m+t+1} | s_t = s]$. The policy network objective function corresponds to advantage weighted log-likelihood of chosen actions, where advantage function is defined as the difference between Q-function and value function. Algorithms such as TRPO and PPO employ importance sampling to correct for the estimation of the advantage function according to the old policy $\pi_{\theta_{old}}$ (that is, the policy before it is updated in a given policy iteration). As a result, the policy objective function contains the ratio term $\mathbf{r}(\theta) = \pi_\theta(a|s)/\pi_{\theta_{old}}(a|s)$. Subsequently, the objective function for the integrated network is the sum of policy objective function added and value loss term multiplied by value coefficient c_v . In the TRPO algorithm, the destructive steps of large gradients often encountered in policy gradient algorithms such as A2C are avoided by penalizing the KL-divergence between old and new policies with the factor β . In the PPO algorithm, this is achieved by clipping the ratio $\mathbf{r}(\theta)$ between $1 - \epsilon$ and $1 + \epsilon$ for a small value of $\epsilon \in [0, 1]$. Furthermore, the exploration in policy search is maximized by maximizing the entropy of the learned policy $S[\pi_\theta](s)$ and is added in the objective function with the entropy coefficient c_e .

B Principle behind computational savings of MLMC estimator

Suppose that we estimate the expectation of the quantity $P^L(\omega)$ where ω is a random variable that follows the probability distribution Ω (that is, $\omega \sim \Omega$). The Monte Carlo estimate of this quantity is given by $\widehat{\mathbb{E}}_\Omega^{MC}(P^L(\omega)) = N^{-1} \sum_{i=1}^N P^L(\omega_i)$. If C and V , respectively, correspond to the computational cost and variance of the term $P^L(\omega_i)$, the cost of the estimator $\widehat{\mathbb{E}}_\Omega^{MC}(P^L(\omega))$ is CN while its overall variance is VN^{-1} . That is, to achieve an overall variance of ϵ^2 , we need to choose $N = \epsilon^{-2}V$ (that is, $N \propto V$). Now, if we suppose that we have an approximation of $P^L(\omega)$ defined as $P^l(\omega)$ such that $\mathbb{V}_\Omega[P^l(\omega)] \gg \mathbb{V}_\Omega[P^L(\omega) - P^l(\omega)]$, the two-level Monte Carlo estimator can be written as $\widehat{\mathbb{E}}_\Omega^{2LMC}(P^L(\omega)) = N_l^{-1} \sum_{i=1}^{N_l} P^l(\omega_i) + N_L^{-1} \sum_{i=1}^{N_L} P^L(\omega_i) - P^l(\omega_i)$. If C_L and V_L are the computational cost and variance of the term $P^L(\omega_i) - P^l(\omega_i)$ while C_l and V_l are the computational cost and variance of the term $P^l(\omega_i)$. The total cost of this two-level Monte Carlo estimator can be computed as $N_l C_l + N_L C_L$ where $N_l \propto V_l$ and $N_L \propto V_L$. Since, by definition, $V_l \gg V_L$, we can also conclude that $N_l \gg N_L$. In other words, if $C_l \ll C_L$, computational cost of two-level Monte Carlo estimate $\widehat{\mathbb{E}}_\Omega^{2LMC}(P^L(\omega))$ is much smaller than the Monte Carlo estimate $\widehat{\mathbb{E}}_\Omega^{MC}(P^L(\omega))$. The same concept can be extended to multilevel Monte Carlo instead of two-level Monte Carlo estimate.

Table 5: Objective function $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta)]$ for different deep RL algorithms

Algorithm	Objective function, $\mathbb{E}_{s,a,r \sim p_\theta}[J(s, a, r; \theta, \Theta)]$
DQN (value network)	$\mathbb{E}_{s,a,r \sim p_\theta} \left[(r + \gamma \max_{a'} Q_{\theta_{old}}(s', a') - Q_\theta(s, a))^2 \right]$
A2C (policy + value network)	$\mathbb{E}_{s,a,r \sim p_\theta} \left[\log \pi_\theta(a s) A(s, a) - c_v (r + \gamma \max_{s'} V_{\theta_{old}}(s') - V_\theta(s))^2 \right]$
TRPO (policy + value network)	$\mathbb{E}_{s,a,r \sim p_\theta} \left[\mathbf{r}(\theta) A(s, a) - \beta \text{KL}[\pi_{\theta_{old}}(\cdot a), \pi_\theta(\cdot a)] - c_v (r + \gamma \max_{s'} V_{\theta_{old}}(s') - V_\theta(s))^2 \right]$
PPO (policy + value network)	$\mathbb{E}_{s,a,r \sim p_\theta} \left[\min(\mathbf{r}(\theta) A(s, a), \text{clip}(\mathbf{r}(\theta), 1 - \epsilon, 1 + \epsilon) A(s, a)) - c_v (r + \gamma \max_{s'} V_{\theta_{old}}(s') - V_\theta(s))^2 + c_e S[\pi_\theta](s) \right]$

C Implementation of multilevel PPO in stable baselines 3

The multilevel PPO algorithm is implemented using the Stable Baselines3 (SB3) (Raffin et al., 2021) library, which is a set of reliable implementations of reinforcement learning algorithms in PyTorch. The codes for the multilevel implementation can be found in the fork: <https://github.com/atishdixit16/stable-baselines3>. In the following text, the implementation of the classical PPO in SB3 is explained in detail. Then it is followed by additional implementations corresponding to the multilevel PPO algorithm.

C.1 Classical PPO implementation in stable baselines 3

RL framework consists of the environment \mathcal{E} which is governed by a Markov decision process described by the tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \mu \rangle$. Here, $\mathcal{S} \subset \mathbb{R}^{n_s}$ is the state-space, $\mathcal{A} \subset \mathbb{R}^{n_a}$ is the action-space, $\mathcal{P}(s'|s, a)$ is a Markov transition probability function between the current state s and the next state s' under action a and $\mathcal{R}(s, a, s')$ is the reward function. The function $\mu(s)$ returns a state from the initial state distribution if s is the terminal state of the episode; otherwise, it returns the same state s . The goal of reinforcement learning is to find the policy $\pi_\theta(a|s)$ to take an optimal action a when in the state s , by exploring the state-action space with what are called agent-environment interactions. Figure 13 shows a typical schematic of such agent-environment interaction. The term *agent* refers to the controller that follows the policy $\pi_\theta(a|s)$ while the *environment* consists of the transition function, \mathcal{P} , and the reward function, \mathcal{R} .

The algorithm 5 delimits the simplified implementation of the PPO algorithm in SB3. The algorithm's inputs are: environment E , number of actors N , number of steps in each policy iter-

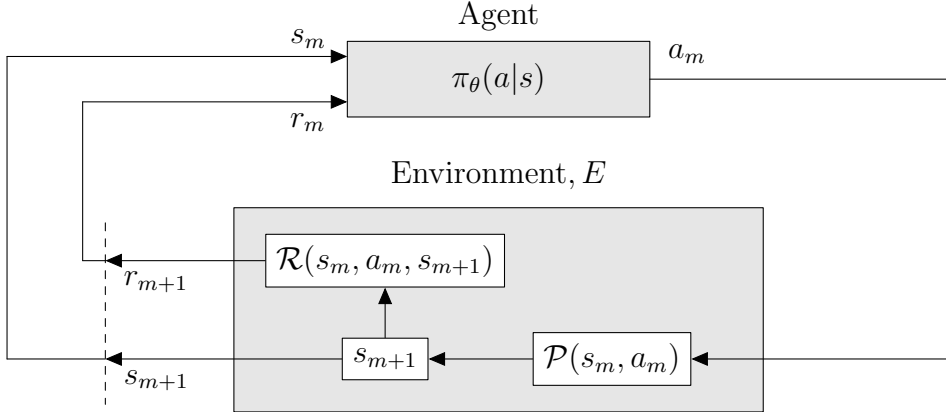


Figure 13: A typical agent-environment interaction for classical framework

ation T , batch size M ($\leq NT$) and number of epochs K . The data obtained through the rollouts of agent-environment interactions is stored in a buffer named RolloutBuffer in the format $[s, a, r, d, V, L_{\text{old}}, R, A]$, where the notation is

- s : state,
- a : action,
- r : reward,
- d : episode terminal boolean (done),
- V : Value function (obtained from policy network rollout),
- L_{old} : log probability value, $\log(\pi_{\theta_{\text{old}}}(a|s))$
- R : Return value (obtained using generalized advantage estimation),
- A : Advantage function (obtained using generalized advantage estimation).

RolloutBuffer accumulates in total $N \times T$ rows of the above data in each iteration. At the beginning of each iteration, the function **CollectRollouts** is used to fill in the data in RolloutBuffer. The total of $N \times T$ data rows is divided into batches of size M , each using the function **GetBatches**. The actor loss term L_a , the value loss term L_v and the entropy loss term L_e (defined in equation 5) are calculated for each such batch using the function **ComputeBatchLosses**. Finally, a Monte Carlo estimate for the loss term is computed as follows.

$$\text{loss}_{\text{MC}} = \text{mean}[L_a + L_v + L_e],$$

which is used to update the policy parameters using automatic differentiation. This is done using the function **UpdatePolicy** and is performed K times for every batch.

The algorithm 6 delineates the steps of the function **CollectRollouts**. For every timestep, the data is obtained using policy rollout, environment transition (using *step* function) and generalized advantage estimation (GAE) computation on all N actors and stored in the RolloutBuffer. Finally,

Algorithm 5 PPO implementation in stable baselines

```
1: Input:  $E, N, T, M, K$ 
2:  $E.reset()$ 
3: Generate empty RolloutBuffer
4: for iteration,  $i = 1, 2, \dots$  do
5:   CollectRollouts( $E, N, T, \text{RolloutBuffer}$ )
6:   for  $epoch = 1, 2, \dots, K$  do
7:     for batch in GetBatches(RolloutBufferArray,  $M$ ): do
8:        $L_a, L_v, L_e = \text{ComputeBatchLosses}$ (batch)
9:        $\text{loss}_{\text{MC}} = \text{mean}[L_a + L_v + L_e]$ 
10:      UpdatePolicy(  $\text{loss}_{\text{MC}}$  )
11:    end for
12:  end for
13: end for
```

Algorithm 6 **CollectRollouts**($E, N, T, \text{RolloutBuffer}$)

```
1: Information: a RolloutBuffer consists of following data:  $[s, a, r, d, V, L_{\text{old}}, R, A]$ 
2: reset RolloutBuffer (i.e. empty the buffer)
3: for  $t$  in  $\text{range}(T)$ : do
4:   rollout current state  $s$ , through policy network to obtain  $a, V, L_{\text{old}}(a)$  on  $N$  actors
5:   if  $s$  is terminal,  $s = E.reset()$ 
6:    $s', r, d, \cdot = E.step(a)$  on  $N$  actors
7:   compute  $R$  and  $A$  using GAE
8:   add  $[s, a, r, d, V, L_{\text{old}}, R, A]$  in the RolloutBuffer
9: end for
```

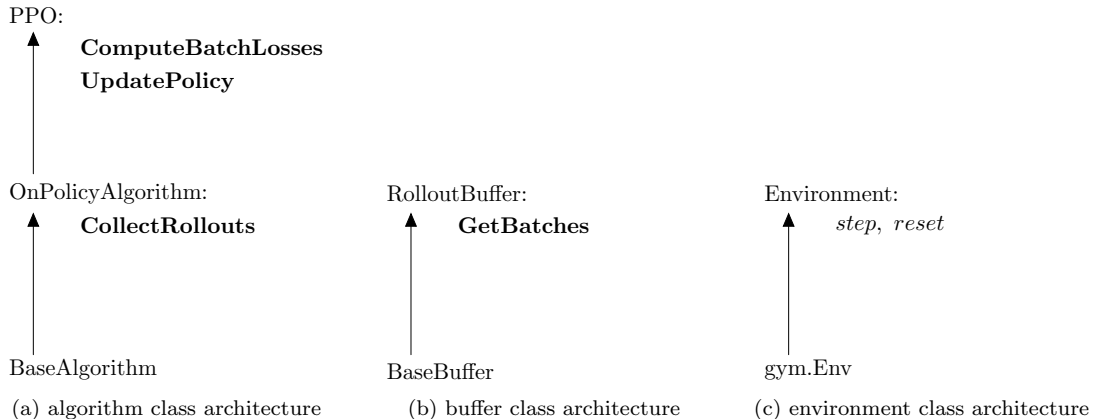


Figure 14: Object-oriented design for the stable baselines implementation of PPO algorithm

ComputeBatchLosses function is illustrated in the algorithm 7. The algorithm lists steps to compute actor loss term L_a , value loss term L_v and entropy loss term L_e for the given batch. Note that the loss terms are the vectors of dimension M , which are added later, and its mean is treated as the final loss term. The mean function in this process indicates the *Monte Carlo* estimator of the PPO loss term.

Algorithm 7 `ComputeBatchLosses(batch)`

- 1: Information: a batch consists of M rows following data: $[s, a, V, L_{\text{old}}, R, A]$
 - 2: compute V_{now} and $L_{\text{now}}(a)$ by rolling out s through policy network
 - 3: compute ratio, $r_t = \exp(L_{\text{now}} - L_{\text{old}})$
 - 4: compute $L_1 = Ar_t$ and $L_2 = A[\text{clip}(r_t, 1 - \epsilon, 1 + \epsilon)]$
 - 5: $L_a = \min(L_1, L_2)$
 - 6: $L_v = C_v |V_{\text{now}} - R|^2$ (C_v is value loss term coefficient)
 - 7: $L_e = -C_e L_{\text{now}}$ (C_e is entropy loss term coefficient)
 - 8: **return** L_a, L_v, L_e
-

The class inheritance schema used in this implementation is shown in Figure 14. The stable baselines use some more classes like Policy, Callbacks etc. but we present only the ones relevant to this discussion. **CollectRollouts** function belongs to OnPolicyAlgorithm which is the child of the BaseAlgorithm class and the parent of the PPO class. The functions **ComputeBatchLosses** and **UpdatePolicy** belong to the PPO class. BaseBuffer is the parent class for the RolloutBuffer class that contains the function **GetBatches**. The Environment class (which is a child of the gym.Env class) contains functions such as *step* and *reset* corresponding to the transition function \mathcal{P} and the initial state function μ , respectively.

C.2 Multilevel PPO implementation in stable baselines 3

Figure 15 illustrates a typical agent-environment interaction in multilevel PPO implementation. Multiple levels of environment are represented with E^1, E^2, \dots, E^{L-1} so that the computational cost of \mathcal{P}^l and the accuracy of \mathcal{R}^l are lower than \mathcal{P}^{l+1} and \mathcal{R}^{l+1} , respectively. The environment

corresponding to the grid fidelity factor l consists of a transition function \mathcal{P}_l , which is achieved by discretizing the dynamical system, and a reward function \mathcal{R}_l . The policy network is designed with states s^L and controls a^L , corresponding to the environment E^L . As a result, state s_{m+1}^l , in the environment, E^l passes through the mapping ψ_l^L which maps the state from level l to level L . Similarly, the action obtained from the policy network is passed through a mapping operator ϕ_L^l , which maps the action from the level L to the level l .

Algorithm 8 illustrates the pseudocode for multilevel implementation of the PPO algorithm in the stable baselines library. The inputs are the same as in classical PPO implementation except multilevel variables are provided as an array of length L : environments at each level $\mathbf{E} = [E^1, E^2, \dots, E^L]$, number of actors N , number of steps in each level $\mathbf{T} = [T^1, T^2, \dots, T^L]$, number of batches in each level $\mathbf{M} = [M^1, M^2, \dots, M^L]$ (such that $NT^l \leq M^l$ and $T^1/M^1 = \dots = T^L/M^L$) and number of epochs K . In multilevel implementation, we formulate the loss term's estimate using multilevel Monte Carlo which is given as

$$\text{loss}_{\text{MLMC}} = \sum_{l=1}^L \text{mean} \left[(L_a^l - \tilde{L}_a^{l-1}) + (L_v^l - \tilde{L}_v^{l-1}) + (L_e^l - \tilde{L}_e^{l-1}) \right],$$

where $\tilde{L}_a^0, \tilde{L}_v^0$ and \tilde{L}_e^0 are set to zero. The outline of a typical agent-environment interaction to obtain synchronized samples of levels l and $l-1$ is illustrated in Figure 15. We use arrays of RolloutBuffers

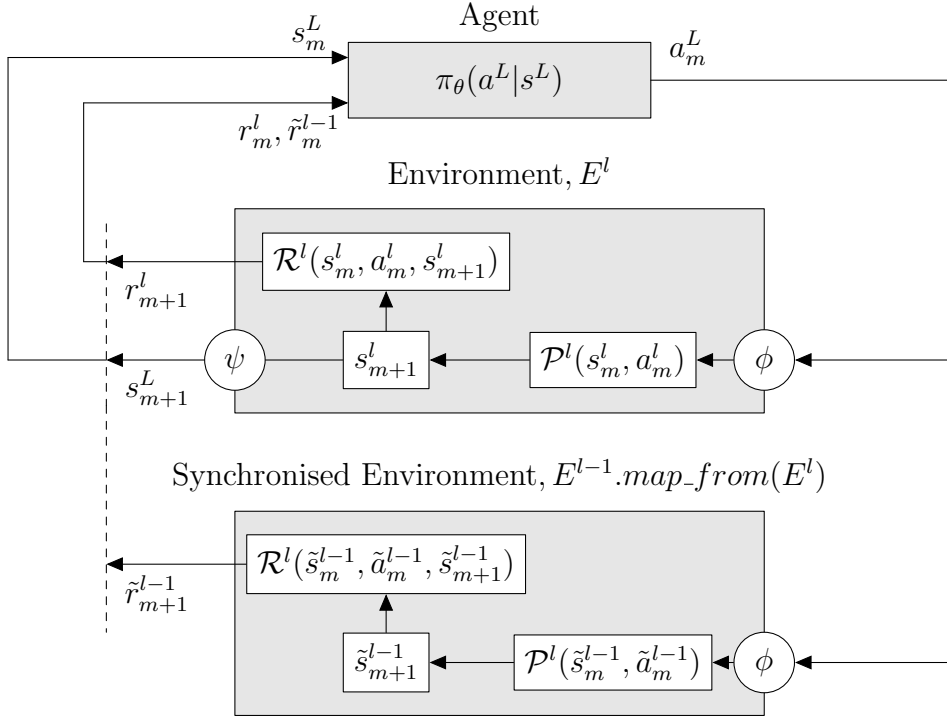


Figure 15: A typical agent-environment interaction for an environment on level l synchronized with environment on level $l-1$

for each level, and each RolloutBuffer ^{l} that collects rollouts at level l has a synchronized buffer

$$\begin{bmatrix} \text{RolloutBuffer}^1 \\ \text{RolloutBuffer}^2 \\ \vdots \\ \text{RolloutBuffer}^L \end{bmatrix} \begin{bmatrix} \text{SyncRolloutBuffer}^1 \\ \text{SyncRolloutBuffer}^2 \\ \vdots \\ \text{SyncRolloutBuffer}^L \end{bmatrix}$$

Figure 16: RolloutBufferArray (on left) and SyncRolloutBufferArray (on right). SyncRolloutBuffer^l consists of synchronized data of RolloutBuffer^l with level l to a level $l - 1$. Each buffer with level l consists of $N \times T_l$ rows of data in $[s, a, r, d, V, L_{\text{old}}, R, A]$ format.

$$\begin{bmatrix} \begin{bmatrix} \text{batch}^1, \text{syncBatch}^0 \\ \text{batch}^2, \text{syncBatch}^1 \\ \vdots \\ \text{batch}^L, \text{syncBatch}^{L-1} \end{bmatrix} & \begin{bmatrix} \text{batch}^1, \text{syncBatch}^0 \\ \text{batch}^2, \text{syncBatch}^1 \\ \vdots \\ \text{batch}^L, \text{syncBatch}^{L-1} \end{bmatrix} & \begin{bmatrix} \dots \\ \dots \\ \dots \\ \dots \end{bmatrix} & \begin{bmatrix} \text{batch}^1, \text{syncBatch}^0 \\ \text{batch}^2, \text{syncBatch}^1 \\ \vdots \\ \text{batch}^L, \text{syncBatch}^{L-1} \end{bmatrix} \end{bmatrix}$$

Figure 17: batch_array which is achieved from **GetBatches** function. It consists of in total NT_l/M_l batches as shown with the columns of the array. Each such batch consists of L batches from RolloutBuffers (denoted by batch^l) and SyncRolloutBuffers (denoted by syncBatch^{l-1}). batch^l and syncBatch^{l-1} consists of M^l rows of data in the format, $[o, a, V, L_{\text{old}}, R, A]$.

SyncRolloutBuffer^l that collects corresponding synchronized data at level $l - 1$. This is achieved using the function **CollectRollouts**. Figure 16 illustrates the RolloutBufferArray and SyncRolloutBufferArray used in this algorithm. Furthermore, the **GetBatches** function is used to generate an array of batches, which is used to compute the multilevel Monte Carlo estimate of the loss term. The batch array consists of in total NT^L/M^L batches, where each batch consists of L batches from RolloutBuffers and L batches from SyncRolloutBuffers. Figure 17 illustrates the batch array used in the algorithm. The $\text{batch}^l, \text{syncBatch}^{l-1}$ from RolloutBuffer^l, SyncRolloutBuffer^l are used to compute the $\text{loss}_{\text{MLMC}}$ terms on the level l . In every batch, these terms are computed at each level and added to obtain $\text{loss}_{\text{MLMC}}$, which is used to update the policy network parameters using the function **UpdatePolicy**.

The algorithm 9 delimits the function **CollectRollouts** used in multilevel implementation. At each level l the RolloutBuffer^l is filled with the data, and the corresponding synchronized data at the level $l - 1$ is filled in the SyncRolloutBuffer^l. Since L_a^0, L_v^0 and L_e^0 are set to zero, the data in SyncRolloutBuffer¹ are filled with None values. The mapping functions $\psi_i^{l'}$ and $\phi_i^{l'}$ are implemented as a set of functions in the definition of the environment E^l . As a result, the mapping of state (ψ_i^l from Equation 4) and action (ϕ_L^l from equation 4) to and from the policy + value network is denoted with shorthand notation ψ and ϕ , respectively. Synchronization of state from level l to l' is indicated by *map_from* function that maps an environment E^l to another environment at level l' , denoted as $E^{l'}$. Algorithm 10 illustrates the pseudocode for the **GetBatches** function, which creates mini-batches (as illustrated in Figure 17) from collected data in RolloutBufferArray and SyncRolloutBufferArray.

The class inheritance schema used in the multilevel implementation is shown in figure 18. **CollectRollouts** function belongs to OnPolicyAlgorithmMultilevel which is the child of BaseAlgorithm

Algorithm 8 Multilevel proximal policy optimization pseudocode

```
1: Input:  $\mathbf{E}, N, \mathbf{T}, \mathbf{M}, K$ 
2:  $E^1.reset()$ 
3: Generate empty RolloutBufferArray, SyncRolloutBufferArray
4: for iteration,  $i = 1, 2, \dots$  do
5:   CollectRollouts( $\mathbf{E}, N, \mathbf{T}$ , RolloutBufferArray, SyncRolloutBufferArray)
6:   for  $epoch = 1, 2, \dots, K$  do
7:     for batch_array in GetBatches(RolloutBufferArray, SyncRolloutBufferArray,  $\mathbf{M}$ ): do
8:       loss_MLMC = 0
9:       for batch $l$ , syncBatch $l-1$  in batch_array do
10:         $L_a^l, L_v^l, L_e^l = \mathbf{ComputeBatchLosses}$ (batch $l$ )
11:        if  $l > 1$  then
12:           $\tilde{L}_a^{l-1}, \tilde{L}_v^{l-1}, \tilde{L}_e^{l-1} = \mathbf{ComputeBatchLosses}$ (syncBatch $l-1$ )
13:        else
14:           $\tilde{L}_a^{l-1}, \tilde{L}_v^{l-1}, \tilde{L}_e^{l-1} = 0$ 
15:        end if
16:         $L^l = \mathit{mean} \left[ (L_a^l - \tilde{L}_a^{l-1}) + (L_v^l - \tilde{L}_v^{l-1}) + (L_e^l - \tilde{L}_e^{l-1}) \right]$ 
17:        loss_MLMC = loss_MLMC +  $L^l$ 
18:      end for
19:      UpdatePolicy( loss_MLMC )
20:    end for
21:  end for
22: end for
```

class and the parent of the PPO_ML class. The functions **ComputeBatchLosses** and **UpdatePolicy** belong to the class PPO_ML. BaseBuffer is the parent class for the RolloutBuffer class that contains the function **GetBatches**. The environment class architecture for multilevel framework is similar to that for classical framework except for the additional mapping functions ψ , ϕ and map_from . The updated definitions of the classes and functions are highlighted in red in figure 18.

D Cluster analysis of permeability uncertainty distribution

A set of permeability samples $\mathbf{k} = \{k_1, \dots, k_l\}$, is chosen to represent the variability in the permeability distribution \mathcal{K} . For the optimal control problem, our main interest is the uncertainty in the dynamical response of permeability, rather than the uncertainty in permeability itself. As a result, the connectivity distance (Park, 2011) is used as a measure of the distance between the permeability field samples. The connectivity distance matrix $\mathbf{D} \in \mathbb{R}^{N \times N}$ among the N samples of \mathcal{K} is formulated as

$$\mathbf{D}(k_i, k_j) = \sum_{x''} \int_{t_0}^T [c(x'', t; k_i) - c(x'', t; k_j)]^2 dt,$$

where N corresponds to a large number of samples of uncertainty distribution, $c(x'', t; k_i)$ is the concentration at the location x'' and at time t , when the permeability is set to k_i and all wells are open equally. The multidimensional scaling of the distance matrix \mathbf{D} is used to produce N two-dimensional coordinates d_1, d_2, \dots, d_N , each representing a permeability sample. The coordinates d_1, d_2, \dots, d_N are obtained such that the distance between d_i and d_j is equivalent to $\mathbf{D}(k_i, k_j)$. In

Algorithm 9 CollectRollouts($\mathbf{E}, N, \mathbf{T}$, RolloutBufferArray, SyncRolloutBufferArray)

```
1: Information: a RolloutBuffer consists of following data:  $[s, a, r, d, V, L_{\text{old}}, R, A]$ 
2: reset RolloutBufferArray, SyncRolloutBufferArray (i.e. empty the buffers)
3: for  $T^l, E^l, \text{RolloutBuffer}^l, \text{SyncRolloutBuffer}^l$  in  $\mathbf{E}, \mathbf{T}$ , RolloutBufferArray, SyncRolloutBuffer-
   Array do
4:   if  $l > 1$  then
5:      $E^l.map\_from(E^{l-1})$ 
6:   end if
7:   for  $t$  in  $\text{range}(T^l)$ : do
8:      $s^l = E^l.reset()$  if  $s^l$  is terminal
9:      $s^L = E^l.\psi(s^l)$ 
10:     $a^L = \pi_\theta(a^L|s^L)$ 
11:     $a^l = \phi(a^L)$ 
12:    compute  $V^l$  and  $L_{\text{old}}(a^L)$ 
13:     $\cdot, r^l, d^l, \cdot = E^l.step(a^l)$  on  $N$  actors
14:    compute  $R^l$  and  $A^l$  using GAE
15:    add  $[s^l, a^l, r^l, d^l, V^l, L_{\text{old}}^l, R^l, A^l]$  in the RolloutBuffer $^l$ 
16:
17:   if  $l > 1$  then
18:      $E^{l-1}.map\_from(E^l)$ 
19:      $\tilde{s}^L = E^{l-1}.\psi(\tilde{s}^{l-1})$ 
20:      $\tilde{a}^{l-1} = a^l$ 
21:      $\tilde{a}^L = \pi_\theta(\tilde{a}^L|\tilde{s}^L)$ 
22:      $\tilde{a}^{l-1} = E^{l-1}.\phi(\tilde{a}^L)$ 
23:     compute  $\tilde{V}^{l-1}$  and  $\tilde{L}_{\text{old}}(a^L)$ 
24:      $\cdot, \tilde{r}^{l-1}, \cdot, \cdot = E^{l-1}.step(\tilde{a}^{l-1})$  on  $N$  actors
25:     compute  $\tilde{R}^{l-1}$  and  $\tilde{A}^{l-1}$  using GAE
26:     add  $[\tilde{s}^{l-1}, \tilde{a}^{l-1}, \tilde{r}^{l-1}, \tilde{d}^l, \tilde{V}^{l-1}, \tilde{L}_{\text{old}}^{l-1}, \tilde{R}^{l-1}, \tilde{A}^{l-1}]$  in the SyncRolloutBuffer $^l$ 
27:   else
28:     add  $[None, \dots, None]$  in the SyncRolloutBuffer $^l$ 
29:   end if
30:   end for
31: end for
```

Algorithm 10 GetBatches(RolloutBufferArray, SyncRolloutBufferArray, \mathbf{M})

```
1: set batch_array to an empty array
2: for RolloutBuffer $^l, \text{SyncRolloutBuffer}^l, M^l$  in RolloutBufferArray, SyncRolloutBufferArray,  $\mathbf{M}$ 
   do
3:   set batches to an empty array
4:   for batch $^l, \text{batch}^{l-1}$  in GetSyncBatches(RolloutBuffer $^l, \text{SyncRolloutBuffer}^l, M^l$ ) do
5:     batches.append([batch $^l, \text{batch}^{l-1}$ ])
6:   end for
7:   batch_array.append(batches)
8: end for
9: return batch_array
```

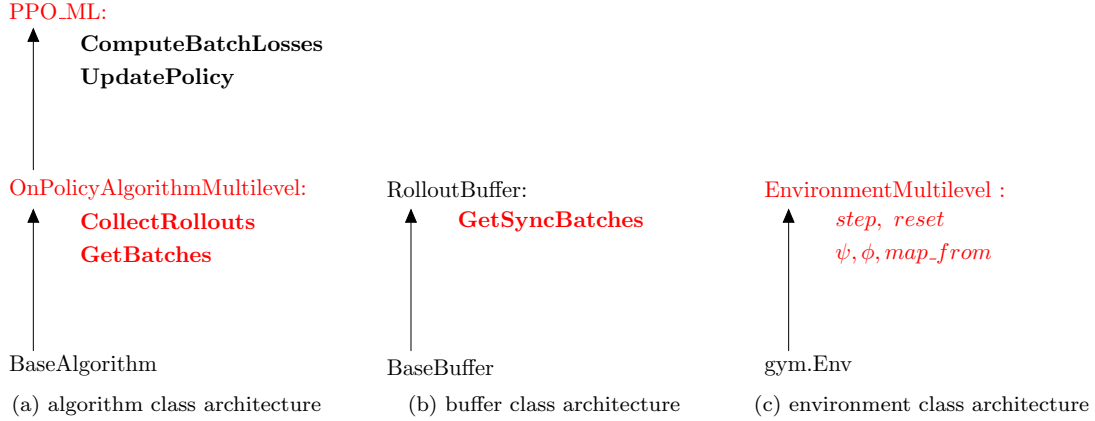


Figure 18: Object-oriented design for the stable baselines implementation of multilevel PPO algorithm. The updated (from classical PPO implementation) definitions of functions and classes are highlighted in red colour.

the k-means clustering process, these coordinates are divided into l sets S_1, S_2, \dots, S_l , obtained by solving the optimization problem:

$$\arg \min_S \sum_i^l \sum_{d_j \in S_i} \|d_j - \mu_{S_i}\|,$$

where μ_{S_i} is the average of all coordinates in the set S_i . The training vector \mathbf{k} is a set of l samples of \mathcal{K} where each of its values k_i corresponds to the closest one to μ_{S_i} . The total number of samples N and clusters l is chosen to be 1000 and 16 for both uncertainty distributions, \mathcal{G}_1 and \mathcal{G}_2 . A training vector \mathbf{k} is obtained with samples k_1, \dots, k_{16} each corresponding to a cluster center. Figures 19a and 19b show cluster plots of permeability samples for ResSim-v1 and ResSim-v2.

E Algorithm Parameters

Parameters used for PPO are tabulated in Table 6 which were tuned using trial and error. For PPO algorithms, the parameters were essentially tuned to find the least variability in the learning plots. The parameters of the DE algorithm are delineated in Table 7. The code repository for both test cases presented in this article can be found at the link: https://github.com/atishdixit16/multilevel_ppo.

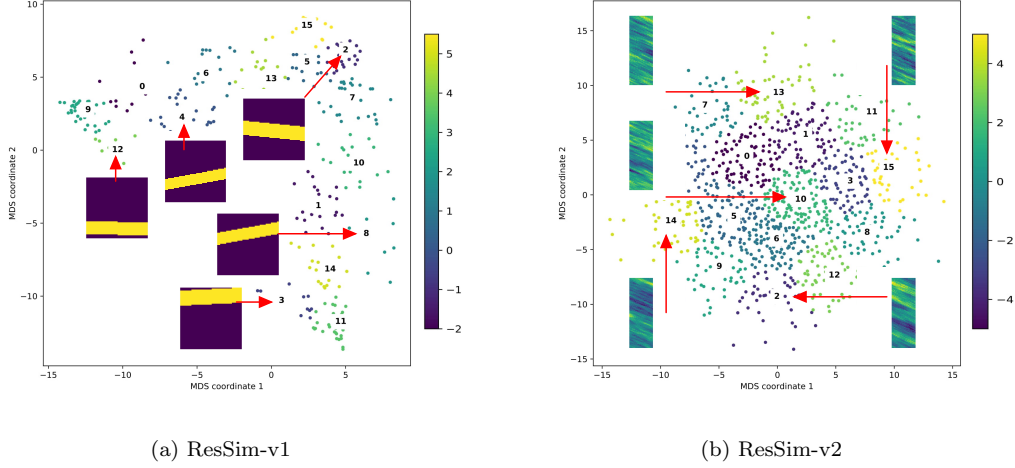


Figure 19: clustering visualization for permeability samples

Table 6: PPO algorithm parameters

	ResSim-v1	ResSim-v2
discount rate, γ	0.99	0.99
clip range, ϵ	0.1	0.15
policy network MLP layers	[93,150,100,80,62]	[35,70,70,50,21]
policy network activation functions	tanh	tanh
policy network optimizers	Adam	Adam
learning rate	3e-6	1e-5

Table 7: DE algorithm parameters

	ResSim-v1	ResSim-v2
number of CPUs	64	64
number of iterations	1024	1024
population size	310	105
recombination factor	0.9	0.9
mutation factor	(0.5,1)	(0.5,1)

References

- Jørg E Aarnes, Tore Gimse, and Knut-Andreas Lie. An introduction to the numerics of flow in porous media using matlab. In *Geometric modelling, numerical simulation, and optimization*, pages 265–306. Springer, 2007.
- Enrico Anderlini, David IM Forehand, Paul Stansell, Qing Xiao, and Mohammad Abusara. Control of a point absorber using reinforcement learning. *IEEE Transactions on Sustainable Energy*, 7(4): 1681–1690, 2016.
- David F Anderson and Desmond J Higham. Multilevel monte carlo for continuous time markov chains, with applications in biochemical kinetics. *Multiscale Modeling and Simulation*, 10(1): 146–179, 2012.
- DR Brouwer, JD Jansen, S Van der Starre, CPJW Van Kruijsdijk, CWJ Berentsen, et al. Recovery increase through water flooding with smart well technology. In *SPE European Formation Damage Conference*. Society of Petroleum Engineers, 2001.
- Neil K Chada, Ajay Jasra, Kody JH Law, and Sumeetpal S Singh. Multilevel bayesian deep neural networks. *arXiv preprint arXiv:2203.12961*, 2022.
- Michael Andrew Christie, MJ Blunt, et al. Tenth SPE comparative solution project: A comparison of upscaling techniques. In *SPE reservoir simulation symposium*. Society of Petroleum Engineers, 2001.
- K Andrew Cliffe, Mike B Giles, Robert Scheichl, and Aretha L Teckentrup. Multilevel monte carlo methods and applications to elliptic pdes with random coefficients. *Computing and Visualization in Science*, 14(1):3–15, 2011.
- Atish Dixit and Ahmed H. ElSheikh. Stochastic optimal well control in subsurface reservoirs using reinforcement learning. *Engineering Applications of Artificial Intelligence*, 114:105106, 2022. ISSN 0952-1976. doi: <https://doi.org/10.1016/j.engappai.2022.105106>. URL <https://www.sciencedirect.com/science/article/pii/S0952197622002469>.
- Michael B Giles. Multilevel monte carlo methods. *Acta numerica*, 24:259–328, 2015.
- Michael B Giles and Lukasz Szpruch. Multilevel monte carlo methods for applications in finance. *High-Performance Computing in Finance*, pages 197–247, 2018.
- Sebastian Müller and Lennart Schüler. Geostat-framework/gstools: Bouncy blue, January 2019. URL <https://doi.org/10.5281/zenodo.2541735>.
- Kwangwon Park. *Modeling uncertainty in metric space*. Stanford University, 2011.
- Jean Rabault, Miroslav Kuchta, Atle Jensen, Ulysse Réglade, and Nicolas Cerardi. Artificial neural networks trained through deep reinforcement learning discover control strategies for active flow control. *Journal of fluid mechanics*, 865:281–302, 2019.
- Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dornmann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. URL <http://jmlr.org/papers/v22/20-1364.html>.

- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Yuyang Shi and Rob Cornish. On multilevel monte carlo unbiased gradient estimation for deep latent variable models. In *International Conference on Artificial Intelligence and Statistics*, pages 3925–3933. PMLR, 2021.
- Matthijs TJ Spaan. Partially observable markov decision processes. In *Reinforcement Learning*, pages 387–414. Springer, 2012.
- Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C J Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020. doi: 10.1038/s41592-019-0686-2.