

Sample-Efficient and Safe Deep Reinforcement Learning via Reset Deep Ensemble Agents

Kaya Unalms, kunalms@princeton.edu.
Lily Gittoes, lgittoes@princeton.edu.

May 18, 2024

Contents

1	Introduction and Motivation	2
1.1	Primacy Bias	2
1.2	Road Map	3
2	Background and Related Work	3
2.1	Overfitting DNN Function Approximators	3
2.2	Agent Resets and Sample Efficiency	4
3	Method	4
3.1	The RDE Algorithm	5
3.2	Implementing RDE with SAC	5
3.3	Implementing RDE with DQN	6
3.4	Experiment Setup	6
4	Results and Discussion	7
4.1	Experimental Results	7
4.2	Evaluation of the Action Selection Coefficient	9
5	Conclusion	10
5.1	Limitations	11
5.2	Future Work	11
A	Algorithms	12
A.1	Pseudocode for RDE with SAC	12
A.1.1	Target Entropy Annealing	13
A.1.2	Change of Variables	13
A.1.3	Suggestions to Improve Target Entropy Annealing	14
A.2	Pseudocode for RDE with SAC	15
B	Hyperparameters	16
C	Hyperparameters	18
D	Compute Times	20
E	Code	21

1 Introduction and Motivation

Many problems in machine learning reduce to, or are usually cast into, an optimization problem to minimize a loss function. Often, it is insufficient to find local minima, and we would like our optimizer to converge to globally interesting minima. This request is demanding, for capturing a global minima in a general optimization space is known to be NP-hard; if we could solve it in polynomial time, any NP-hard problem could be solved efficiently. Therefore, state-of-the-art optimization methods still rely on various versions of gradient descent that are fundamentally local solution finding algorithms. In practice, such gradient descent methods can perform well, so long as they do not get stuck in an undesirable local minima [JNJ17].

This is particularly relevant in cases where the loss function’s gradient is determined in an online manner, or in the more general reinforcement learning (RL) problem, in a manner which is dependent on the trajectories of state-action pairs. In practice, the optimizer moves in the descent direction based on the gradient information of an initial set of samples. If these samples are not representative of the optimization landscape, or even “misleading” in the context of finding a globally interesting minima, the optimizer’s descent will be biased toward an undesirable direction. Depending on the problem, it can be difficult to recover from this bias. Naturally, one attempt to resolve this issue is to proceed with the optimization regardless, collect information about the loss function’s gradient, and then restart the optimization while retaining this information. The hope is that a fresh model with better prior knowledge of the optimization landscape will be more resistant to moving in a descent direction that is biased by early gradient information.

1.1 Primacy Bias

The phenomenon that the model suffers from in the previous section is called *primacy bias*. This concept originates in cognitive science, where results have shown that early impressions have disproportionate impacts on human memories [MW72]. It was extended to deep RL settings by [Nik+22] to describe the overfitting of function approximators to early experiences, which can lead to sub-optimal performance. In such RL settings, function approximators are formulated with deep neural network (DNN) architectures and are used to describe large state-action spaces. They are trained according to samples from a replay buffer that tracks previously observed state transitions. To combat the primacy bias that results in overfitting on early experiences, [Nik+22] proposes a series of periodic agent resets in which their network parameters are reinitialized. As the impacts of primacy bias are compounded when the replay ratio, defined as the number of gradient updates made per environment interaction, is high, the algorithm proposed by [Nik+22] can be applied to improve sample efficiency.

However, immediately following these parameter resets, their model experiences large performance collapses. For certain settings, such as those where the agent must interact with a physical environment (i.e., autonomous cars or robots), these drops in performance raise concerns. With the aim of developing a method that does not suffer from primacy bias nor large safety issues, [Kim+23] proposes the **R**eset-based algorithm by leveraging **D**eep **E**nsemble learning (RDE). By training an ensemble of agents instead of just one, [Kim+23] resets the agents’ network parameters one at a time. This allows [Kim+23] to avoid performance collapses, because their agents do not have to immediately interact with the environment after having their parameters reset. To demonstrate the improvements of RDE over the vanilla reset method by [Nik+22], [Kim+23] implements RDE with two off-policy RL methods: Deep Q-Networks (DQN) and Soft Actor-Critic (SAC) [Mni+15; Haa+19]. They carry out tests over a series of environments, including Atari-100k, MiniGrid, and the DeepMind Control (DMC) Suite [Bel+13; CWP18; Tas+18].

1.2 Road Map

Our project aims to reproduce the results of [Kim+23]. Starting from the algorithm pseudocode they have provided alongside the parameters they tuned, we write code from scratch to employ RDE alongside DQN and SAC. Modifications to the vanilla algorithm that were found necessary for convergence are made. We evaluate our implementation of RDE with SAC using the **Cheetah-Run** setting in the DMC suite and compare our results with those produced by [Kim+23]. After, because of computation and time constraints, rather than performing training in the Atari-100k and MiniGrid environments, which are much harder to learn, we instead evaluate our implementations of RDE with SAC and DQN in the **Mountain Car Continuous** and **Cart Pole Gym** environments [Bro+16].

We lay out the rest of this paper as follows. In Section 2, we describe the different methods employed by others to overcome the challenge of overfitting DNN function approximators and explain how their work relates to ours. Then, in Section 3, we detail the main components of our implementation of RDE with DQN and SAC. Our results are displayed in Section 4, where we also discuss the extent to which our results replicate that of [Kim+23]. Afterwards, we describe the limitations that we ran into, potential paths for future work, and the key conclusions of our paper in Section 5.

2 Background and Related Work

While the concept of primacy bias is new, methods for dealing with overfitting in deep RL settings are not. In this section, we describe methods that have previously been employed by others in relation to the RDE algorithm. For our project, as is used by [Kim+23], we also employ the standard Markov decision process (MDP) formulation. That is, we consider an agent with action space \mathcal{A} that exists in an environment with a state space \mathcal{S} . In a state $s \in \mathcal{S}$, upon choosing an action $a \in \mathcal{A}$, the agent receives a reward $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and transitions to a new state s' according to the probability distribution p . For the discount rate $\gamma \in [0, 1)$, the agent aims to find the optimal policy $\pi : \mathcal{S} \rightarrow \Delta\mathcal{A}$ that maximizes the cumulative discounted rewards $\mathbb{E}[\sum_t \gamma^t r(s_t, a_t)]$.

2.1 Overfitting DNN Function Approximators

DNNs are often used as function approximators to describe the Q-values associated with an MDP’s state-action space. These DNNs are susceptible to overfitting, and a series of different approaches have been suggested for dealing with this, including ensemble learning and prioritized experience replay.

For a policy π and state-action pair (s, a) , the Q-value is defined as $Q^\pi(s, a) = \mathbb{E}^\pi[\sum_t \gamma^t r(s_t, a_t) | s_0 = s]$. We focus on the two off-policy RL algorithms employed by [Kim+23]: DQN and SAC. They both use Temporal Difference (TD) learning, which aims to minimize the difference between $Q^\pi(s, a)$ and $\mathbb{E}_{p(s'|s, a), \pi(a', s')} [r(s, a) + \gamma Q^\pi(s', a')]$. In TD learning settings, gradient descent is implemented according to bootstrapping – namely, regressing towards the model’s own output value function. When this is combined with function approximation in off-policy settings, results are often unstable. In such conditions, it has been found that as more gradient steps are taken, the function approximators lose their expressivity and struggle to fit new functions [Kum+21; LRD22].

The goal of ensemble learning is to integrate a series of agents to act collectively, with the goal of improving the robustness of baseline algorithms. For instance, to improve upon the originally proposed DQN algorithm by [Mni+15], [HGS16] proposes the Double-DQN framework to simultaneously train two sets of Q-networks that update each other. They show that this can alleviate some of the instability problems encountered in the baseline DQN algorithm. Meanwhile, [ABS17] formulates a different ensemble-based algorithm

called Averaged-DQN, in which Q-value estimates are computed as the average across a series of Q-networks. A third method, which is proposed by [Lee+21], suggests training an ensemble of Q-networks that can be used to weight sample transitions according to their level of uncertainty. They implement a weighted Bellman backup, observing that this mitigates some of the error propagation seen in simple Q-learning models.

Meanwhile, others have worked to improve how samples from the replay buffer are selected, with the goal of improving the performance of TD learning methods. Notably, [Sch+16] proposes a sampling method called Prioritized Experience Replay, which assigns weights to samples (s, a, r, s') stored in the replay buffer using the TD difference function $\delta = r + \max_{a'} Q(s', a') - Q(s, a)$. By giving higher weights to samples associated with a larger TD difference, they are able to see significant performance improvements. On the other hand, [Wan+20] suggests a non-uniform sampling method whereby more recent experiences are given higher weights. Although they do not explicitly call it primacy bias, their approach is motivated by preventing overfitting to early experiences.

Inspired by the past successes of ensemble learning in improving the stability of TD learning, [Kim+23] uses this as a key element in their proposed RDE algorithm. In terms of the replay buffer that they maintain, [Kim+23] does not prioritize using certain samples over others when making gradient steps. This is because after performing each agent reset, they want to maintain equal access to all previously seen environment interactions.

2.2 Agent Resets and Sample Efficiency

A key feature of using agent resets to alleviate the impacts of primacy bias is the enhanced sample efficiency it offers. Sample efficiency is improved with higher replay ratios, which allow for a greater number of parameter updates to be made per environment interaction [Fed+20; HHA19]. However, [Nik+22] finds that as the replay ratio increases so does primacy bias. In other words, if too many parameter updates are made early on in the training process, the model could become stuck at a local minima and inflexible to later experiences. To combat the impacts of primacy bias, [Nik+22] introduces the notion of periodically partially or fully resetting an agent’s network parameters. They conjecture that letting an agent start from scratch and make parameter updates according to a preserved replay buffer from prior training steps can prevent overfitting to early environment interactions. By implementing their network reset approach with respect to a series of baseline algorithms, they demonstrate improvements over benchmarks.

Inspired by the work of [Nik+22] on primacy bias, [DOr+23] finds that their agent reset algorithm can be adjusted to allow for improved sample efficiency. Namely, rather than fixing the frequency of agent resets with regards to the number of time steps of environment interactions like [Nik+22] does, [DOr+23] fixes the reset frequency with regards to the number of parameter update steps. By increasing the agent reset frequency alongside the replay ratio, [DOr+23] achieves the high levels of sample efficiency offered by high replay ratios without having to bear the impacts of primacy bias. In their RDE algorithm, [Kim+23] employs an ensemble of agents that are reset sequentially according to the method described by [Nik+22]. Using the results of [DOr+23], their resets occur at a frequency that decreases linearly with the replay ratio. With this, they aim to develop an algorithm that is sample efficient, safe, and avoids primacy bias.

3 Method

In this section, we describe the main components of RDE and provide overviews of the algorithms we implement alongside: SAC and DQN.

3.1 The RDE Algorithm

The RDE algorithm proposed by [Kim+23] that we aim to replicate is divided into the following steps. First, we initialize N ensemble agents, where each agent’s network parameters are denoted by θ_k for $k \in \{1, \dots, N\}$. The N ensemble agents all have the same network architecture.

During model training, the agents $k \in \{1, \dots, N\}$ are reset in a sequential order, starting with $k = 1$. After resetting agent N , we circle back and reset agent 1 again. Resetting an agent consists of resetting all of their network parameters – this includes both the actor and critic networks for SAC and both the Q-network and target Q-network for DQN. Given the input parameter T_{reset} , there are T_{reset}/N training steps between agent resets. Thus, upon completing T_{reset} training steps, all N of the agents will have been reset once. We then continue cycling through the agent resets in this manner for the entirety of the training period.

In terms of interacting with the environment, the ensemble of agents work together as a single unit. According to the agents’ policies $\{\pi_{\theta_1}, \dots, \pi_{\theta_N}\}$, for a given state s , the agents suggest actions $\{a_1, \dots, a_N\}$. The action that is actually taken is selected using the probability distribution $p_{select} = \{p_1, \dots, p_N\}$ over the N possibilities, which is computed as follows. Using the Q-value function \hat{Q} , which is that of the agent that was least recently reset, we calculate Q-values for the suggested actions $\{a_1, \dots, a_N\}$ and the current state s . These are denoted by $\{\hat{Q}(s, a_1), \dots, \hat{Q}(s, a_N)\}$. We then apply the softmax function to these Q-values, using the temperature

$$\alpha = \beta / \max(\hat{Q}(s, a_1), \dots, \hat{Q}(s, a_N)).$$

In this, β is a parameter that is tuned and allows us to decide on the extent to which we assign higher weights to actions with high Q-values and lower weights to those with low Q-values. Thus, the probability distribution p_{select} is given by

$$p_{select} = \{p_1, \dots, p_N\} = \text{softmax}\{\hat{Q}(s, a_1)/\alpha, \dots, \hat{Q}(s, a_N)/\alpha\}. \quad (1)$$

and the probability of selecting action a_k is p_k . We note that before the first reset (which occurs at time step T_{reset}/N), all of the agent’s network parameters are identical. Hence, the probability distribution is uniform: i.e., $p_{select} = \{1/N, \dots, 1/N\}$.

By selecting actions according to the Q-value function of the agent that was least recently reset, we are able to avoid the performance collapses observed in the vanilla reset method proposed by [Nik+22]. As is explained by [Kim+23], this is because choosing an action according to p_{select} from an ensemble of agents allows us to maintain a proper balance between exploration and exploitation. In the vanilla reset method, since there is just one agent, the model is only able to perform exploration immediately after it is reset, leading to drops in performance. Moreover, in the ensemble reset method, choosing actions from a uniform distribution over the N agents is also insufficient to avoid performance collapses. This is because actions selected by agents that were recently reset will be from untrained models, hence they should be taken with a lower probability. By using the probability distribution p_{reset} , these actions are given a lower weight, as actions are selected according to the Q-value function of the least recently reset agent. Thus, agent resets can occur without significant performance drops.

3.2 Implementing RDE with SAC

Following suit with [Kim+23], we start by implementing RDE with SAC. For SAC, we implement the Twin Delayed Deep Deterministic (TD3) algorithm [Haa+19; FHM18]. This algorithm extends upon standard Actor-Critic methods to add encouragement for

policies that achieve a higher entropy, thereby adding incentives for exploration. For each agent $i \in \{1, \dots, N\}$, we initialize two sets of Q-networks alongside two sets of corresponding target Q-networks. In addition, we also initialize a policy network and corresponding target policy network for each agent. By interacting with the environment using the N ensemble agents’ policy networks, experiences are stored in a circular replay buffer \mathcal{B} . In a circular replay buffer, after the number of stored sample transitions reaches the buffer’s capacity, the oldest experiences start to be removed as new ones are added. We note, however, that the buffer’s capacity is set to be very large so that agents that have been reset can still access old samples. Using uniform samples from \mathcal{B} , the agents update their policy networks and Q-networks.

Pseudocode for combining RDE with SAC, including its loss functions, is described in Appendix A. We note that this pseudocode does not account for two characteristics of the SAC algorithm described by [Haa+19]: target entropy annealing and a modified entropy equation. We find that these two techniques are not needed for SAC to be successful in the **Mountain Car Continuous** environment, but they are needed in the **Cheetah-Run** setting, likely because it is more complex. Hence, we only implement them for **Cheetah-Run** and this is further described in Appendix A. Our code for implementing RDE with SAC in these two environments is available in Appendix E.

3.3 Implementing RDE with DQN

Next, we implement RDE with DQN according to the framework proposed by [Mni+15]. In this algorithm, each agent interacts with the environment according to an ϵ -greedy policy and keeps track of their past transitions (s_t, a_t, r_t, s_{t+1}) in a circular replay buffer \mathcal{B} . The agent aims to learn the parameters θ of the DNN that describes their Q-value function $Q(s, a; \theta)$. Pseudocode describing the integration of RDE with DQN is provided in Appendix A. We run simulations of our code in the **Cart Pole** Gym environment. While [Kim+23] runs tests for RDE with DQN in the Atari-100k and MiniGrid environments, because these environments are so complex, our DQN algorithm takes a long time to train. Moreover, [Kim+23] does not provide a complete set of their hyperparameters (in particular, they are missing key details about their DNN architecture used in the MiniGrid environments). Given the long training time for DQN in complex settings with sparse rewards, we find the **Cart Pole** environment to be more suitable for testing the effectiveness of RDE in this project.

3.4 Experiment Setup

In this section, we lay out our setup for the experiments we perform in the **Cheetah-Run**, **Mountain Car Continuous**, and **Cart Pole** settings. For each environment, [Kim+23] runs simulations for a series of nine scenarios. They consider the baseline algorithm (TD3 or DQN) with no ensembles or resets, the algorithm with the vanilla reset method proposed by [Nik+22] (which is identical to the RDE algorithm when $N = 1$), and the algorithm with RDE using $N = 2$ agents. For each, they test replay ratios of 1, 2, and 4. Moreover, [Kim+23] scales the reset interval T_{reset} so that it decreases linearly with the replay ratio. We start by reproducing the results of [Kim+23] in the **Cheetah-Run** setting in the DMC suite [Tas+18]. Then, we repeat these same nine scenarios for the **Mountain Car Continuous** Gym environment [Bro+16].

For the RDE with SAC algorithm in the **Cheetah-Run** setting, we use almost identical hyperparameters to [Kim+23]. However, in order to decrease the computational complexity of our training loop, we make changes such as using fewer hidden units in our DNN layers and sampling smaller batches from the replay buffer, and these are detailed in Appendix B. These changes allow us to closely replicate the results of [Kim+23] with less time and

computational power. Moreover, as the runtime scales linearly with the replay ratio and the number of agents, for the scenario in which we use a replay ratio of 4, we use half the number of training steps as was used when the replay ratio is 1 or 2. Our limitations regarding computational power are discussed in more detail in Section 5.

Then, when implementing RDE with the SAC algorithm in the **Mountain Car Continuous** environment, we use similar parameters to what we used in the **Cheetah-Run** setting, but because the model trains faster in this environment, we examine a fewer number of total timesteps and use an analogously smaller number of timesteps between resets T_{reset} . However, to avoid the extra large computation times in the scenario where we implement RDE with a replay ratio of 4, we use half the number of training steps compared to when the replay ratio is 1 or 2.

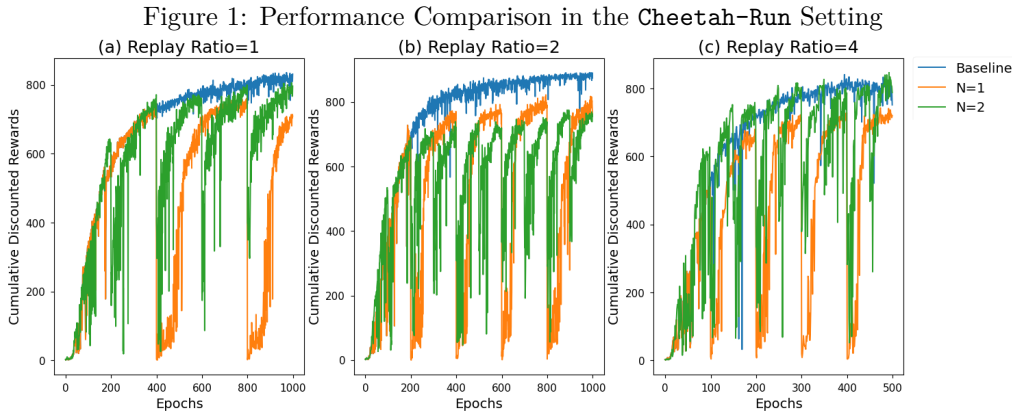
As for when we implement RDE with the DQN algorithm in the **Cart Pole** environment, we use parameters that are similar to those used by [Kim+23] in their MiniGrid experiments. Here, we also halve the number of training steps when applying a replay ratio of 4 because of computational constraints. Moreover, we also use the same parameter T_{reset} for both replay ratios of 2 and 4 rather than decreasing it linearly with the replay ratio as we find that doing so makes the resets so frequent that the agents do not have enough time to recover after. For all three of the environments we consider, more details on our DNN architectures and hyperparameters are provided in Appendix B.

4 Results and Discussion

Using the methods and experimental setups described in Section 3, we run simulations of our code implementations and compare our results to those presented by [Kim+23].

4.1 Experimental Results

In this section, we display the results of our experiments when implementing RDE with SAC and DQN. To begin, we plot a comparison of the rewards achieved in the **Cheetah-Run** setting below in Figure 1. This series of plots aims to replicate plots (d), (e), and (f) in Figure 7 by [Kim+23].

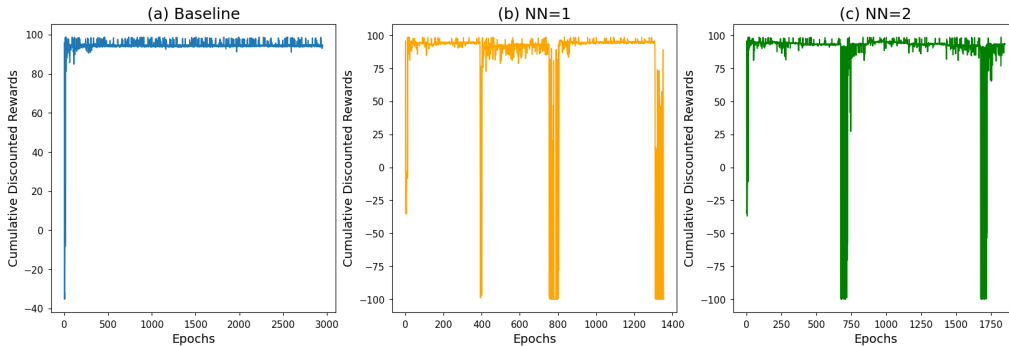


For **Cheetah-Run**, this figure displays the performance of the SAC in the baseline case (in blue), with naive resets with one agent (in orange), and with ensemble resets using RDE with two agents (in green). Results are shown across replay ratios of 1, 2, and 4.

From Figure 1, we can see that when the replay ratio is 1 or 2, the baseline SAC algorithm outperforms both the vanilla reset and RDE methods, though not by a very significant margin. Meanwhile, when the replay ratio is 4, by the end of the training loop, the RDE method does about as well as the baseline SAC algorithm. These results are fairly consistent with [Kim+23], who does not observe significant performance improvements from implementing RDE in the **Cheetah-Run** setting. However, we note that while our implementation of RDE with SAC does not experience performance collapses as significant as those in the vanilla reset method, they are larger than those observed by [Kim+23]. We identify two reasons for this. First, given the long training time that is required, we are only able to run our code once, while [Kim+23] averages their results over 5 seeds, so there is some inherent level of randomness that causes our results to differ. Moreover, given some inconsistencies identified in the discussion that [Kim+23] provides on how to tune the β parameter in computing p_{select} , we believe that this part of our implementation may not exactly match theirs. Further discussion on this is provided in Section 4.2.

Next, we display our results in the **Mountain Car Continuous** setting. We make comparisons across the baseline SAC algorithm without any agent resets, SAC with naive resets, and SAC with RDE for replay ratios of 1, 2, and 4. Our results for a replay ratio of 2 is shown below in Figure 2, while those for replay ratios of 1 and 4 can be found in Appendix C, as we do not find much variation in patterns across the three scenarios.

Figure 2: Performance Comparison in the **Mountain Car Continuous** Setting when $RR=2$



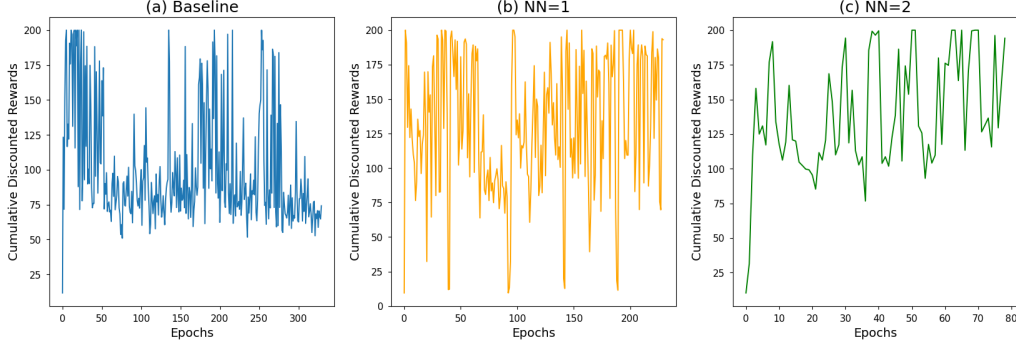
For **Mountain Car Continuous**, we compare the performance of SAC with a replay ratio of two across the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents.

For the **Mountain Car Continuous** setting, we can see that the baseline SAC algorithm is successful on its own at learning the environment. This is likely because of the relatively small state and action space of the environment. Thus, as our algorithm appears to not suffer from primacy bias and overfitting, our results show that agent resets do not bring about significant performance improvements. Nonetheless, across all three tested replay ratios, we see that performance collapses are more frequent in the naive reset method compared to RDE, demonstrating that RDE is successful in improving safety. However, we notice that the extent of the drop in rewards experienced in each performance collapse with RDE when $N = 2$ is very similar to that in the naive reset method. Similar to what we observed in the **Cheetah-Run** environment, we believe that this is the result of our use of the β parameter being different than that described by [Kim+23], and we explain this in more detail in Section 4.2.

Finally, we display results for our simulations in the **Cart Pole** environment. Because

there is not a significant difference in trends across the three replay ratios we test, we focus on results when the replay ratio is 2, and these are shown below in Figure 3 below. Meanwhile, results when the replay ratio is 1 or 4 are shown in Appendix C.

Figure 3: Performance Comparison in the Cart Pole Setting when RR=2



For DQN (using a replay ratio of two) in the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents, we display our results.

We observe that the performance collapses observed when using naive resets are significantly greater than those seen when using RDE. This suggests that RDE indeed does provide a significant improvement over naive resets. However, because DQN is off-policy, uses bootstrapping, and employs function approximators to describe a large state-action space, it suffers from what is known as the deadly triad [Has+18]. When any of these three components are used together, it has been empirically discovered to be more likely to introduce instabilities to standard RL algorithms, impeding convergence of the optimizer to acceptable policies. It is also for this reason that it is particularly difficult for DQN to perform well on many of the more complex MiniGrid environments which have sparse rewards, especially when constrained to a short compute time. The instability of DQN as an algorithm makes it difficult to discern whether some of the performance collapses seen in Figure 3 are from noise or agent resets, especially when using RDE. Nonetheless, because of its added complexity, we note that the RDE algorithm with $N = 2$ agents completes a much smaller number of epochs over the same number of training steps as the baseline DQN or naive reset methods, and it is perhaps a direction for future work to evaluate this algorithm over more episode rollouts.

4.2 Evaluation of the Action Selection Coefficient

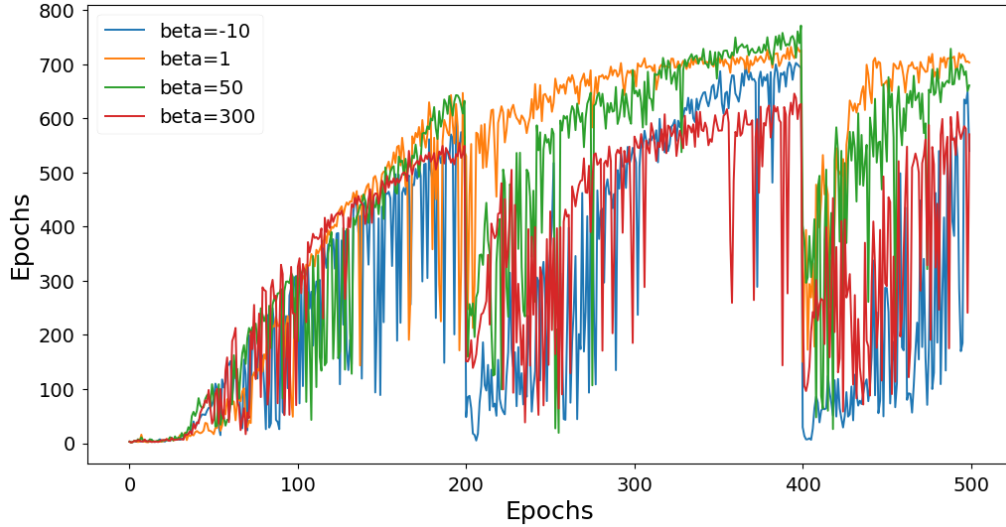
We improve upon the results presented by [Kim+23] by more thoroughly testing the impact of the action selection coefficient β , which is used to compute the softmax temperature α in Equation 1. Here, we notice that there are inconsistencies in their presentation of β . Namely, [Kim+23] claims that high β values yield better performance by adding higher weights to actions chosen by agents that are least recently reset. While we agree with their assessment that negative values of β exacerbate performance collapses, we find that for positive values of $\beta \geq 0$, lower β values place higher weights on less recently reset agents. To see this, consider the definition of the softmax function: for temperature $\alpha \in \mathbb{R}_+$, the softmax function $s_\alpha: \mathbb{R}^c \rightarrow [0, 1]^c$ is defined by

$$s_\alpha: x \mapsto \begin{bmatrix} e^{x(1)/\alpha} \\ \vdots \\ e^{x(c)/\alpha} \end{bmatrix} / \sum_{k=1}^c e^{x(k)/\alpha}.$$

As $\alpha \rightarrow 0$, the sum $\sum_{k=1}^c e^{x(k)/\alpha}$ gets exponentially larger, hence the softmax function places greater weight on larger input values. Since α is defined as being directly proportional to β , the same is true for β . This contradicts the statement by [Kim+23] that higher β values put more weight on actions chosen by less recently reset agents.

To see their claim, [Kim+23] tests values of $\beta \in \{-10, 0, 50, 300\}$. However, considering $\beta = 0$ makes the softmax function in Equation 1 undefined from dividing by 0. For a replay ratio of one, we run experiments of RDE with TD3 using $N = 2$ agents in the **Cheetah-Run** environment (using the same hyperparameters as those presented in Appendix B except that we use 2×10^5 training steps to save time) for varying β values in $\{-10, 1, 50, 300\}$. We plot our results below in Figure 4, where we observe that $\beta = 1$ yields the best performance by placing higher weights on least recently reset agents. Thus, we believe that the discussion [Kim+23] provides on β is inconsistent with their definition of p_{select} .

Figure 4: Performance in the **Cheetah-Run** Setting when RR=1 Across Different β Parameter Values



We display our results of RDE using $N = 2$ agents with SAC in the **Cheetah-Run** setting when the replay ratio is one. We consider β parameters in $\{-10, 1, 50, 300\}$.

5 Conclusion

In conclusion, by reproducing the results of [Kim+23], we find that RDE indeed does allow for less significant performance collapses, and it is an improvement over the naive reset method proposed by [Nik+22]. However, in some cases we have observed that the baseline method with no agent rests performs similarly to the ensemble reset method. In other environments the baseline stochasticity makes it hard to identify whether performance drops and boosts are due to primacy bias or some other influence such as insufficient duration of training, function approximation errors, or instabilities (e.g. such as those that arise from the deadly triad). While primacy bias has a precise cause that we can identify, its future effects are usually difficult to identify with confidence, just like primacy bias in cognitive science. In the sections below, we discuss the limitations of our work

and suggest directions for the future.

5.1 Limitations

The primary limitations in our ability to reproduce the results of [Kim+23] are computational power and time. We list the compute times required for each of our simulations in Appendix D. To this end, we could not fully test the entire suite of environments considered in their paper. Especially for the more complex environments such as Atari-100k and MiniGrid, we did not have the capacity to run a large number of training loops. As [Kim+23] saw varying results across the different environments they considered (including those where RDE significantly outperforms baseline algorithms with no resets), our project is limited in its ability to verify all of their claims.

Moreover, to save time, we also took measures such as using fewer hidden units in our DNNs, taking smaller batches from the replay buffer, taking fewer total training steps, and only running simulations on one seed rather than averaging results over multiple. Meanwhile, as [Kim+23] is inconsistent in their discussion of the β parameter, this could also be causing differences in our experimental results. Thus, we can see that while our results generally align with those found by [Kim+23], differences in implementations do cause some variation in results.

5.2 Future Work

With more time and greater computational resources, there are many directions that can be explored. For instance, varying parameters such as the replay ratio and the reset interval T_{reset} or resetting only some of the layers of the DNNs being trained are all possible future steps. In addition, it would be interesting to run simulations with a larger ensemble. This will allow us to test whether choosing actions based on a larger ensemble of agents will further mitigate primacy bias or if there are diminishing returns after two models. Similarly, experimenting with different methods of selecting an agent could be worthwhile as well, as there could be better ways to formulate the probability distribution p_{reset} than that in Equation 1.

Furthermore, more simulations across a larger number of seeds can help us better understand the extent to which underlying stochasticity is impacting our observed results. To this end, further tests over a wider variety of environments that range in difficulty will help validate the results of [Kim+23]. In addition, considering the impacts of RDE on a wider set of algorithms could also prove to be insightful. As an example, we suggest a possible improvement on temperature annealing in Appendix A.

A Algorithms

A.1 Pseudocode for RDE with SAC

The algorithm below describes the implementation we use of RDE with SAC.

Algorithm 1 Algorithm Pseudocode for RDE+SAC

Require: Learning rate η , Target update frequency f , Target update rate τ , Discount γ , Ensemble size N , Reset frequency T_{reset} , Replay Ratio RR , Coefficient β , Temperature $temp$, Target update weight τ .

- 1: Initialize agent parameters $\{\theta_1, \dots, \theta_N\}$, where each agent θ_i has policy parameters ξ_i , target policy parameters ξ_i^- , parameters for two Q-networks $\phi_{i,1}$ and $\phi_{i,2}$ (which we denote collectively as ϕ_i), and corresponding target Q-network parameters $\phi_{i,1}^-$ and $\phi_{i,2}^-$.
 - 2: Initialize the replay buffer \mathcal{B} .
 - 3: **for** each episode **do**
 - 4: **for** each timestep t **do**
 - 5: **for** $i = 1$ to N **do**
 - 6: Based on state s_t , select an action a_t^i using current policies π_{θ_i} .
 - 7: **end for**
 - 8: Calculate p_{select} using Eq. 1.
 - 9: Sample action a_t from p_{select} and play it.
 - 10: Store observed transition (s_t, a_t, r_t, s_{t+1}) in the replay buffer \mathcal{B} .
 - 11: **for** $j = 1$ to RR **do**
 - 12: Sample a random minibatch from \mathcal{B} .
 - 13: **for** $i = 1$ to N **do**
 - 14: Take gradient steps of size η to update ϕ_i to minimize:
$$L(\phi_i, \mathcal{B}) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(\mathcal{B})} \left[(Q_{target} - Q(s, a, \phi_i))^2 \right]$$

where for each (s, a, r, s') pair, $Q_{target} = r + \gamma(Q(s, a, \phi_i^-) + temp * \mathcal{H}^\pi)$.
 - 15: **if** $t \% f == 0$ **then**
 - 16: Take a gradient steps of size η to update the policy ξ_i to minimize:
$$L(\xi_i, \mathcal{B}) = \mathbb{E}_{s \sim \mathcal{U}(\mathcal{B}), a \sim \pi_{\xi_i}} \left[temp * \log(\pi_{\xi_i}(a|s)) - Q_{\phi_i}(s, a) \right]$$
 - 17: Update $\phi_i^- \leftarrow \tau \phi_i + (1 - \tau) \phi_i^-$
 - 18: Update $\xi_i^- \leftarrow \tau \xi_i + (1 - \tau) \xi_i^-$
 - 19: **end if**
 - 20: **end for**
 - 21: **end for**
 - 22: **if** $t \% (T_{reset}/N) == 0$ **then**
 - 23: Reset θ_k and set $k \leftarrow (k + 1) \% N$
 - 24: **end if**
 - 25: **end for**
 - 26: **end for**
-

In addition to the algorithm above, which we directly apply to the Mountain Car Continuous environment, we also implement target entropy annealing and a change of variables in the Cheetah-Run environment. We find that these features are necessary for TD3 to converge, and we explain them in the following sections.

A.1.1 Target Entropy Annealing

The first technical addition in our TD3 algorithm is target entropy annealing. Without temperature annealing, the only regularization term in our loss function in the form of a constant temperature times the entropy of the distribution from which we select actions at a particular state. However, as the optimizer attempts to minimize the loss on some finite-sized empirical sample of actions and rewards, there are opportunities to overfit to the sample. The resulting prediction for the probability distribution of actions given that state will then have less entropy, which can be undesirable if we are not certain that the sample was representative. We would ideally like to achieve the minimum error as determined by the loss function while retaining maximum entropy, where the latter objective is one that penalizes overfitting, for low entropy is a first order estimate of overfitting. Therefore, we can improve performance by using our prior knowledge to select a target or goal entropy of our distributions. Throughout the optimization, we dynamically change the temperature to encourage the policy’s entropy (according to [Haa+19]) at each state to approach the target entropy via regularization term in loss function. Formally, we append the following term to our loss function

$$R = \alpha(\mathcal{H} - \mathcal{H}_{\text{goal}}) \quad (2)$$

where α specifies the temperature parameter which is tuned by the optimizer so as to minimize loss. We can see that this has the desirable characteristic that when $\mathcal{H} < \mathcal{H}_{\text{goal}}$, α will be increased to minimize R , prioritizing an increase in entropy. Otherwise when $\mathcal{H} > \mathcal{H}_{\text{goal}}$, α will be decreased to minimize R , prioritizing a decrease in the other terms of the loss function, (typically at the expense of reducing entropy). We choose $\mathcal{H}_{\text{goal}}$ to be the negative of the dimension of the action space, in agreement with [Kim+23]. That is, we choose some value that is proportional to the entropy of the action space, so that the optimizer will be able to more easily tune α . Note that negative entropy is not an issue because this is the differential entropy, which is standard practice in reinforcement learning, even though it is not a dimensionless quantity: it has units of length. Therefore, care should be taken to scale the entropy by the measure of the action space. For **Cheetah-Run** environment we consider, the action space consists of unit discs on the real line $[-1, 1]$, so the scaling is up to a small constant already normalized by the measure of the action space (the optimizer can handle the necessary small correction).

A.1.2 Change of Variables

To accomplish the change of variables suggestion by [Haa+19], we draw actions a according to the current state and policy, where $u \sim N(\text{mean}, \text{variance})$ for some Multivariate Gaussian distribution, and $a = \tanh u$. By composing these maps to select an action, we are altering the entropy of the effective distribution which generates the samples, and we needed to account for this for the model to be successful. The modified entropy is

$$\begin{aligned} \mathcal{H}(\pi(a \mid s)) &= \mathbb{E} - \log \pi(a \mid s) \\ &= H(N(\text{mean}, \text{variance})) + \log \|\mathbf{D}_u \tanh u\| \\ &= H(N(\text{mean}, \text{variance})) + \mathbb{E} * \sum_{i=1}^D (1 - \tanh(u)^2) + \log 1 \end{aligned}$$

where the final term is zero since the action space is already normalized to 1. In our code, this is implemented in an automatically differentiable manner via PyTorch’s `distribution.transform(tanh)` and `distribution.log_abs_det_jacobian(a, u)`.

A.1.3 Suggestions to Improve Target Entropy Annealing

For future work, we also suggest how the target entropy annealing strategy described in Appendix A.1.1 can be improved. It is typical in reinforcement learning implementations to choose H_{goal} to be a constant that is reflective of our prior knowledge, in many cases setting it to some linear function of the dimensional of the action space. However, it is a somewhat lazy application of prior knowledge to have a goal entropy that is uniform throughout the state action space. The optimizer will need to work harder to tune the temperature α to overcome this bad approximation. To give a concrete example, suppose the problem at hand is to balance a pole on a cart. When the pole is approximately upright, the choice is not clear which direction to move the cart. For it could have just started falling in the direction of its tilt, or perhaps it experienced an acceleration at an earlier time that is now realized as an angular velocity $\omega > 0$ when the displacement is $\Delta\theta < 0$, so this tilt will soon disappear without external interaction. Therefore, in states with $\Delta\theta \approx 0$ a high target entropy is desirable.¹ In other states, say with large $|\Delta\theta|$, the optimal action is clear and the cart should move in the direction to minimize $|\Delta\theta|$ regardless of the velocity ω as the risk, or in the safety RL context the cost, of not performing corrective action is too large. These states are better suited with a low goal entropy assignment. Making this improvements to the target entropy annealing model would significantly ease the problem given to the optimizer of adaptive selection of the best temperature. However, in general, it can be burdensome to implement this improved target entropy annealing because it can be difficult to determine which states should be assigned low entropy.

¹That is, high entorpy in the direction of movement. The magnitude of the external force applied to balance the pole should still be small, as desired for most states in this problem.

A.2 Pseudocode for RDE with SAC

Next, we display pseudocode for our implementation of RDE with DQN.

Algorithm 2 Algorithm Pseudocode for RDE+DQN

Require: Epsilon ϵ , Learning rate η , Target update frequency f , Discount γ , Ensemble size N , Reset frequency T_{reset} , Replay Ratio RR , Coefficient β .

```

1: Initialize Q-function parameters  $\{\theta_1, \dots, \theta_N\}$  and target Q-function parameters  $\{\theta_1^-, \dots, \theta_N^-\}$ .
2: Initialize the replay buffer  $\mathcal{B}$ .
3: for each episode do
4:   for each timestep  $t$  do
5:     for  $i = 1$  to  $N$  do
6:       Based on state  $s_t$ , select an action  $a_t^i$  according to the  $\epsilon$ -greedy policy.
7:     end for
8:     Calculate  $p_{select}$  using Eq. 1.
9:     Sample action  $a_t$  from  $p_{select}$  and play it.
10:    Store observed transition  $(s_t, a_t, r_t, s_{t+1})$  in the replay buffer  $\mathcal{B}$ .
11:    for  $j = 1$  to  $RR$  do
12:      Sample a random minibatch from  $\mathcal{B}$ .
13:      for  $i = 1$  to  $N$  do
14:        Take gradient step of size  $\eta$  to update  $\theta_i$  to minimize:

$$L(\theta, \mathcal{B}) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(\mathcal{B})} \left[ \left( r + \gamma \max_{a'} Q(s', a'; \theta_i^-) - Q(s, a, \theta_i) \right)^2 \right]$$

15:      end for
16:      if  $f\%t == 0$  then
17:        Copy  $\{\theta_1, \dots, \theta_N\}$  into  $\{\theta_1^-, \dots, \theta_N^-\}$ 
18:      end if
19:    end for
20:    if  $t \% (T_{reset}/N) == 0$  then
21:      Reset  $\theta_k$  and set  $k \leftarrow (k+1)\%N$ 
22:    end if
23:  end for
24: end for

```

B Hyperparameters

In the Table 1 below, we display the hyperparameters we used to train RDE with SAC in the **Cheetah-Run** setting. As we aim to reproduce the results of [Kim+23] exactly, almost all of our network parameters are identical to theirs. However, for the sake of decreasing computational complexity, we adjust a few of their hyperparameters. Namely, we use a minibatch size of 512 instead of 1024 and 256 hidden units in each network layer instead of 1024. Note that the network layers are multi-layer perceptions (MLPs). Moreover, we also implement a delayed network update frequency by setting $f = 2$ unlike [Kim+23] who uses $f = 1$.

Table 1: Hyperparameters for RDE with SAC in the **Cheetah-Run** Setting

Hyperparameter	Value
# of Ensemble Agents N	2
Training steps	$(1 \times 10^6, 5 \times 10^5)$
Discount factor	0.99
Initial collection steps	5000
Minibatch size	512
Optimizer	Adam
Learning rate η	0.0003
Network activation (all)	ReLU
Network layer type (all)	MLP
Number of hidden network layers (all)	2
Number of hidden units per network layer (all)	256
Initial temperature	1
Replay buffer size	1×10^6
Replay ratio	(1, 2, 4)
Target update frequency f	2
Target update weight τ	0.005
Reset interval T_{reset}	$(4 \times 10^5, 2 \times 10^5, 1 \times 10^5)$
Action selection coefficient β	50

In Table 2 below, we display the hyperparameters we used in the **Mountain Car Continuous** environment when using RDE with SAC. Our parameters are very similar to those in Table 1 except because the **Mountain Car Continuous** environment is easier to learn than the **Cheetah-Run** setting, we use fewer training steps and hence also use smaller values for T_{reset} . Lastly, we display the hyperparameters we used to train RDE with DQN in the **Cart Pole** environment below in Table 3.

Table 2: Hyperparameters for RDE with SAC in the **Mountain Car Continuous** Setting

Hyperparameter	Value
# of Ensemble Agents N	2
Training steps	$(2 \times 10^5, 1 \times 10^5)$
Discount factor	0.99
Initial collection steps	5000
Minibatch size	512
Optimizer	Adam
Learning rate η	0.0003
Network activation (all)	ReLU
Network layer type (all)	MLP
Number of hidden network layers (all)	2
Number of hidden units per network layer (all)	256
Temperature	0.01
Replay buffer size	2×10^5
Replay ratio	(1, 2, 4)
Target update frequency f	2
Target update weight τ	0.005
Reset interval T_{reset}	$(8 \times 10^4, 4 \times 10^4, 2 \times 10^4)$
Action selection coefficient β	50

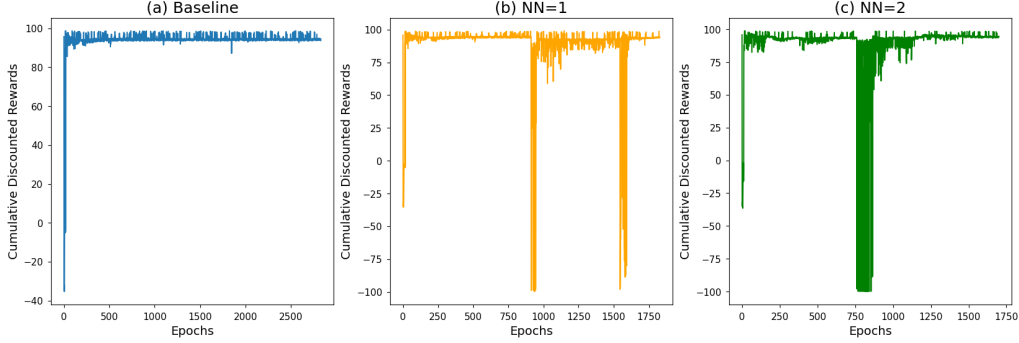
Table 3: Hyperparameters for RDE with DQN in the **Cart Pole** Setting

Hyperparameter	Value
# of Ensemble Agents N	2
Training steps	$(2 \times 10^5, 1 \times 10^5)$
Discount factor	0.99
Initial collection steps	1×10^4
Minibatch size	32
Optimizer	Adam
decay time step	$0.9 \rightarrow 0.05$
Learning rate η	10^5
Network activation (all)	0.0005
Network layer type (all)	ReLU
Number of hidden network layers (all)	MLP
Number of hidden units per network layer (all)	2
Replay buffer size	128
Replay ratio	5×10^4
Target update frequency f	(1, 2, 4)
Network testing frequency	10 episodes
Reset interval T_{reset}	20 episodes
Action selection coefficient β	$(8 \times 10^4, 4 \times 10^4)$
	0.01

C Hyperparameters

In Figure 5, for a replay ratio of 1, we display the results of our performance comparison of RDE with SAC using $N = 2$ agents with the naive reset and baseline SAC algorithms.

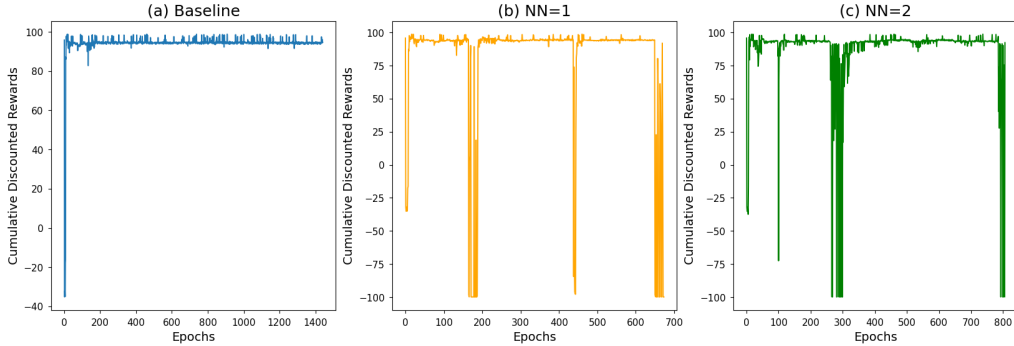
Figure 5: Performance Comparison in the Mountain Car Continuous Setting when $RR=1$



This figure displays our results in the Mountain Car Continuous setting for SAC. We make comparisons across the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents.

Next, in Figure 6, we show the results of our performance comparison when the replay ratio is 4.

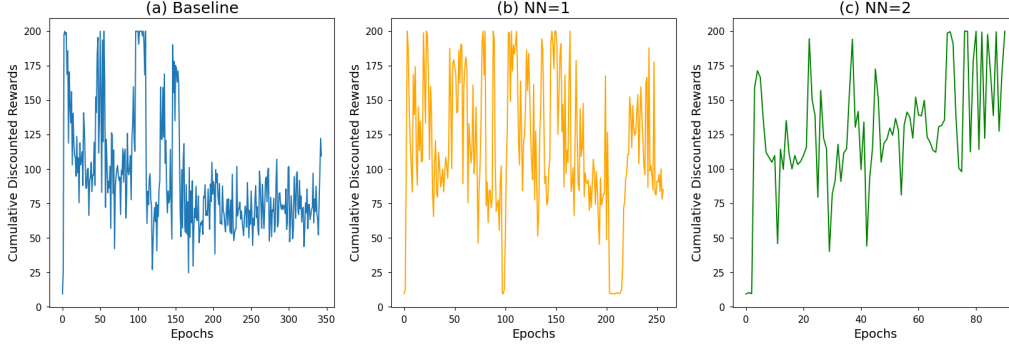
Figure 6: Performance Comparison in the Mountain Car Continuous Setting when $RR=4$



We compare the performance of SAC with a replay ratio of four in the Mountain Car Continuous across the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents.

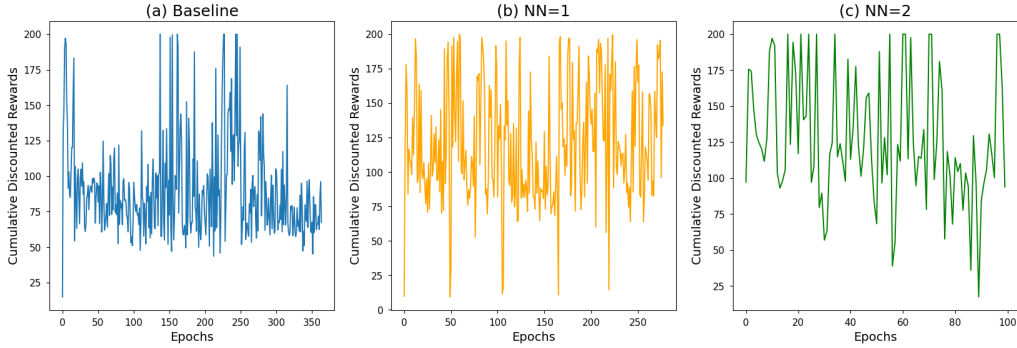
Next, in Figures 7 and 8, we display our results for replay ratios of 1 and 4 in the **Cart Pole** environment.

Figure 7: Performance Comparison in the **Cart Pole** Setting when $RR=1$



For a replay ratio of one, we compare performance in the **Cart Pole** setting for DQN. We consider the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents.

Figure 8: Performance Comparison in the **Cart Pole** Setting when $RR=4$



We display our results in the **Cart Pole** setting for DQN when the replay ratio is four. We examine the baseline case, with naive resets with one agent, and with ensemble agents using RDE with two agents.

D Compute Times

In the tables below, we report the compute times of our experiments in each of the three environments: **Cheetah-Run**, **Mountain Car Continuous**, and **Cart Pole**.

Table 4: Run Times for SAC in the **Cheetah-Run** Setting on GPU

	Baseline	Naive Reset	RDE
$RR = 1$ for 1×10^6 time steps	3 hrs	6 hrs	12 hrs
$RR = 2$ for 1×10^6 time steps	6 hrs	12 hrs	24 hrs
$RR = 4$ for 5×10^5 time steps	6 hrs	12 hrs	24 hrs

Table 5: Run Times for SAC in the **Mountain Car Continuous** Setting on CPU

	Baseline	Naive Reset	RDE
$RR = 1$ for 2×10^5 time steps	1.5 hrs	3 hrs	6 hrs
$RR = 2$ for 2×10^5 time steps	3 hrs	6 hrs	12 hrs
$RR = 4$ for 1×10^5 time steps	3 hrs	6 hrs	12 hrs

Table 6: Run Times for DQN in the **Cart Pole** Setting on CPU

	Baseline	Naive Reset	RDE
$RR = 1$ for 1×10^6 time steps	1 hr	2 hrs	4 hrs
$RR = 2$ for 1×10^6 time steps	2 hrs	4 hrs	8 hrs
$RR = 4$ for 5×10^6 time steps	2 hrs	4 hrs	8 hrs

E Code

In this section, we display the code we implemented for RDE with SAC and DQN in the Cheetah-Run, Mountain Car Continuous, and Cart Pole environments. We note that parts of our code are adapted from our solutions for assignments 4 and 7.

Listing 1: Code for RDE with SAC in the Cheetah-Run Setting

```

1  #!/usr/bin/env python
2  # coding: utf-8

4  # In[1]:

6  # @title Run to install MuJoCo and `dm_control`
7  import distutils.util
8  import os
9  import subprocess

11 if subprocess.run("nvidia-smi").returncode:
12     raise RuntimeError(
13         "Cannot communicate with GPU. "
14         "Make sure you are using a GPU Colab runtime. "
15         "Go to the Runtime menu and select Choose runtime type."
16     )

18 print("Installing dm_control...")
19 get_ipython().system("pip install -q dm_control>=1.0.18")

21 # Configure dm_control to use the EGL rendering backend (requires GPU)
22 get_ipython().run_line_magic("env", "MUJOCO_GL=egl")

24 print("Checking that the dm_control installation succeeded...")
25 try:
26     from dm_control import suite

28     env = suite.load("cartpole", "swingup")
29     pixels = env.physics.render()
30 except Exception as e:
31     raise e from RuntimeError(
32         "Something went wrong during installation. Check the shell output above "
33         "for more information.\n"
34         "If using a hosted Colab runtime, make sure you enable GPU acceleration "
35         "by going to the Runtime menu and selecting \"Choose runtime type\"."
36     )
37 else:
38     del pixels, suite

40 get_ipython().system(
41     'echo Installed dm_control $(pip show dm_control | grep -Po "(?<=Version: ).+")'
42 )

44 # In[2]:

46 # @title All `dm_control` imports required for this tutorial

48 # The basic mujoco wrapper.
49 from dm_control import mujoco

51 # Access to enums and MuJoCo library functions.
52 from dm_control.mujoco.wrapper.mjbindings import enums
53 from dm_control.mujoco.wrapper.mjbindings import mjlib

55 # PyMJCF
56 from dm_control import mjcf

58 # Composer high level imports
59 from dm_control import composer
60 from dm_control.composer.observation import observable
61 from dm_control.composer import variation

63 # Imports for Composer tutorial example
64 from dm_control.composer.variation import distributions
65 from dm_control.composer.variation import noises
66 from dm_control.locomotion.arenas import floors

68 # Control Suite

```

```

69 from dm_control import suite

71 # Run through corridor example
72 from dm_control.locomotion.walkers import cmu_humanoid
73 from dm_control.locomotion.arenas import corridors as corridor_arenas
74 from dm_control.locomotion.tasks import corridors as corridor_tasks

76 # Soccer
77 from dm_control.locomotion import soccer

79 # Manipulation
80 from dm_control import manipulation

82 # In[3]:

84 import random
85 import numpy as np
86 import torch
87 import torch.nn as nn
88 from torch.distributions import Categorical
89 import torch.nn.functional as F
90 import torch.optim as optim
91 import matplotlib.pyplot as plt
92 import sys
93 import copy
94 from typing import Tuple

96 get_ipython().run_line_magic("matplotlib", "inline")
97 import pickle

99 # In[4]:

101 class ReplayBuffer(object):
102     def __init__(self, state_dim, action_dim, max_size=int(1e6)):
103         self.max_size = max_size
104         self.ptr = 0
105         self.size = 0

107         self.state = np.zeros((max_size, state_dim))
108         self.action = np.zeros((max_size, action_dim))
109         self.next_state = np.zeros((max_size, state_dim))
110         self.reward = np.zeros((max_size, 1))
111         self.not_done = np.zeros((max_size, 1))

113         self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

115     def add(self, state, action, next_state, reward, done):
116         self.state[self.ptr] = state
117         self.action[self.ptr] = action
118         self.next_state[self.ptr] = next_state
119         self.reward[self.ptr] = reward
120         self.not_done[self.ptr] = 1.0 - done

122         self.ptr = (self.ptr + 1) % self.max_size
123         self.size = min(self.size + 1, self.max_size)

125     def sample(self, batch_size):
126         ind = np.random.randint(0, self.size, size=batch_size)

128         return (
129             torch.FloatTensor(self.state[ind]).to(self.device),
130             torch.FloatTensor(self.action[ind]).to(self.device),
131             torch.FloatTensor(self.next_state[ind]).to(self.device),
132             torch.FloatTensor(self.reward[ind]).to(self.device),
133             torch.FloatTensor(self.not_done[ind]).to(self.device),
134         )

136 # In[5]:

138 def init_flags():

140     flags = {
141         "env": "cheetah",
142         "env_action": "run",
143         "seed": 0, # random seed
144         "start_timesteps": 5e3, # total steps of free exploration phase
145         "max_timesteps": 6e5, # maximum length of time steps in training

```



```

146         "batch_size": 512,
147         "discount": 0.99,
148         "tau": 0.005, # rate of target update
149         "policy_freq": 2, # delayed policy update frequency in TD3,
150         "N": 2, # number of agents,
151         "RR": 1, # replay ratio,
152         "T": 4e5, # time steps between agent resets ,
153         "beta": 1, # action selection coefficient,
154         "actor_model_file": None, # "actor_model.pt",
155         "critic_model_file": None, # "critic_model.pt",
156         "temp_model_file": None, # "temp_model.pt"
157     }

159     return flags

161 def collect_actions(theta, state):
162     actions = []
163     for theta_i in theta:
164         action = theta_i.select_action(np.array(state))
165         actions.append(torch.from_numpy(action))
166     return actions

168 def get_timestep_info(timestep):
169     ob = timestep.observation
170     state = np.concatenate([*ob.values()]).ravel()
171     return state, timestep.reward, timestep.last()

173 def main(policy_name="DDPG") -> list:
174     """
175     Input:
176     policy_name: str, the method to implement
177     Output:
178     evaluations: list, the reward in every episodes
179     Call DDPG/TD3 trainer and
180     """
181     args = init_flags()
182     env = suite.load(
183         args["env"], args["env_action"], task_kwargs={"random": args["seed"] + 100}
184     )
185     torch.manual_seed(args["seed"])
186     np.random.seed(args["seed"])

188     action_spec = env.action_spec()
189     ob_spec = env.observation_spec()

191     state_dim = 0
192     for _, val in ob_spec.items():
193         state_dim += val.shape[0]
194     state_dim

196     action_dim = action_spec.shape[0]
197     max_action = action_spec.maximum
198     kwargs = {
199         "state_dim": state_dim,
200         "action_dim": action_dim,
201         "max_action": max_action,
202         "discount": args["discount"],
203         "tau": args["tau"],
204         "actor_model_file": args["actor_model_file"],
205         "critic_model_file": args["critic_model_file"],
206         "temp_model_file": args["temp_model_file"],
207     }
208     if policy_name == "TD3":
209         kwargs["policy_freq"] = args["policy_freq"]
210         theta = [TD3(**kwargs) for _ in range(args["N"])]
211     elif policy_name == "DDPG":
212         policy = DDPG(**kwargs)

214     replay_buffer = ReplayBuffer(state_dim, action_dim)
215     evaluations = []
216     timestep = env.reset()
217     state, _, _ = get_timestep_info(timestep)
218     done = False
219     episode_reward = 0
220     episode_timesteps = 0
221     episode_num = 0
222     k = 0

```

```

223     max_episode_steps = 990

225     for t in range(int(args["max_timesteps"])):

227         episode_timesteps += 1

229         if t % 1e5 == 0:
230             theta[0].save_actor_model(filename="actor_model_2.pt")
231             theta[0].save_critic_model(filename="critic_model_2.pt")
232             theta[0].save_temp_model(filename="temp_model_2.pt")
233             np.savetxt("reward_cheetah_run_N2_RR2.csv", evaluations, delimiter=",")

235         # Select action randomly or according to policy
236         entropy = 0
237         mean = 0
238         vari = 0
239         if t < args["start_timesteps"]:
240             action = np.random.uniform(
241                 action_spec.minimum, action_spec.maximum, size=action_spec.shape
242             )
243         else:
244             with torch.no_grad():

246                 actions = collect_actions(theta, state)
247                 # compute Q's, then apply softmax
248                 q_sa = torch.hstack(
249                     [
250                         theta[k].critic.Q1(torch.FloatTensor(state), action)
251                         for action in actions
252                     ]
253                 )
254                 # dim
255                 max_q_sa, _ = torch.max(q_sa, dim=0)
256                 alpha = args["beta"] / max_q_sa
257                 p_select = F.softmax(q_sa / alpha)

259                 if t % 1000 == 0:
260                     print("temperature:", torch.exp(theta[0].log_alpha))

262                 rng = np.random.default_rng()
263                 action = rng.choice(
264                     a=[np.array(tensor) for tensor in actions],
265                     p=p_select.numpy(),
266                     axis=0,
267                 )
268                 action = np.atleast_1d(action)

270             # Perform action
271             next_state, reward, done = get_timestep_info(env.step(action))
272             done_bool = float(done) if episode_timesteps < max_episode_steps else 0

274             # Store data in replay buffer
275             replay_buffer.add(state, action, next_state, reward, done_bool)

277             state = next_state
278             episode_reward += reward

280             # Train agent after collecting sufficient data
281             if t >= args["start_timesteps"]:
282                 for j in range(args["RR"]):
283                     for theta_i in theta:
284                         theta_i.train(replay_buffer, args["batch_size"])

286             if (t % (args["T"] / args["N"])) == 0:
287                 print(k)
288                 # reset just actor or both? ask to confirm
289                 theta[k].actor.reset()
290                 theta[k].critic.reset()
291                 k = (k + 1) % args["N"]

293             if done:
294                 # +1 to account for 0 indexing. +0 on ep_timesteps since it will increment +1 even if done=
True
295                 print(
296                     f"Total T: {t+1} Episode Num: {episode_num+1} Episode T: {episode_timesteps} Reward: {
episode_reward:.3f}"
297                 )

```

```

298         evaluations.append(episode_reward)
299         # Reset environment
300         done = False
301         timestep = env.reset()
302         state, _, _ = get_timestep_info(timestep)
303         episode_reward = 0
304         episode_timesteps = 0
305         episode_num += 1
307     return evaluations
309 # In[6]:
311 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
313 # Construct the actor/critic network for TD3
314 class Actor_TD3(nn.Module):
315     def __init__(self, state_dim: int, action_dim: int, max_action: float):
316
317         super(Actor_TD3, self).__init__()
318         self.l1 = nn.Linear(state_dim, 256)
319         self.l2 = nn.Linear(256, 256)
320         self.l3 = nn.Linear(256, 2 * action_dim)
322
323         if isinstance(max_action, np.ndarray):
324             max_action = torch.from_numpy(max_action)
325         else:
326             assert isinstance(max_action, list)
327             max_action = torch.tensor(max_action)
328             self.max_action = max_action.to(torch.float32)
329             self.action_dim = action_dim
330
331     def forward(self, state: torch.Tensor) -> torch.Tensor:
332         if isinstance(state, np.ndarray):
333             state = torch.from_numpy(state)
334
335         state = state.float()
336         a = F.relu(self.l1(state))
337         a = F.relu(self.l2(a))
338
339         a = self.l3(a)
340
341         mean = a[:, : self.action_dim]
342         cov = nn.functional.softplus(a[:, self.action_dim :]) + 1e-9
343
344         dist = torch.distributions.MultivariateNormal(
345             mean, scale_tril=torch.diag_embed(cov)
346         )
347         return dist
348
349     def reset(self):
350         for layer in self.children():
351             if hasattr(layer, "reset_parameters"):
352                 layer.reset_parameters()
353
354 class Critic_TD3(nn.Module):
355     def __init__(self, state_dim: int, action_dim: int):
356         super(Critic_TD3, self).__init__()
357
358         # Q1 architecture
359         self.l1 = nn.Linear(state_dim + action_dim, 256)
360         self.l2 = nn.Linear(256, 256)
361         self.l3 = nn.Linear(256, 1)
362
363         # Q2 architecture
364         self.l4 = nn.Linear(state_dim + action_dim, 256)
365         self.l5 = nn.Linear(256, 256)
366         self.l6 = nn.Linear(256, 1)
367
368     def forward(
369         self, state: torch.Tensor, action: torch.Tensor
370     ) -> Tuple[torch.Tensor, torch.Tensor]:
371         sa = torch.cat([state, action], 1)
372         q1 = F.relu(self.l1(sa))
373         q1 = F.relu(self.l2(q1))
374         q1 = self.l3(q1)

```

```

375         q2 = F.relu(self.l4(sa))
376         q2 = F.relu(self.l5(q2))
377         q2 = self.l6(q2)
378         return q1, q2

380     def Q1(self, state: torch.Tensor, action: torch.Tensor) -> torch.Tensor:
381         sa = torch.cat([state, action], -1)

383         q1 = F.relu(self.l1(sa))
384         q1 = F.relu(self.l2(q1))
385         q1 = self.l3(q1)
386         return q1

388     def reset(self):
389         for layer in self.children():
390             if hasattr(layer, "reset_parameters"):
391                 layer.reset_parameters()

393 # In[7]:

395 class TD3(object):
396     def __init__(
397         self,
398         state_dim: int,
399         action_dim: int,
400         max_action: float,
401         discount=0.99,
402         tau=0.005,
403         policy_freq=2,
404         init_temperature=1,
405         actor_model_file=None,
406         critic_model_file=None,
407         temp_model_file=None,
408     ):

410         self.actor = Actor_TD3(state_dim, action_dim, max_action).to(device)
411         self.actor_target = copy.deepcopy(self.actor)
412         self.actor_optimizer = torch.optim.Adam(self.actor.parameters(), lr=3e-4)

414         self.critic = Critic_TD3(state_dim, action_dim).to(device)
415         self.critic_target = copy.deepcopy(self.critic)
416         self.critic_optimizer = torch.optim.Adam(self.critic.parameters(), lr=3e-4)

418         if actor_model_file is not None:
419             print("Loading pretrained model from", actor_model_file)
420             self.load_pytorch_actor_model(actor_model_file)

422         if critic_model_file is not None:
423             print("Loading pretrained model from", critic_model_file)
424             self.load_pytorch_critic_model(critic_model_file)

426         self.log_alpha = torch.tensor(np.log(init_temperature)).to(device)
427         self.log_alpha.requires_grad = True
428         self.target_entropy = -1 * action_dim
429         self.log_alpha_optimizer = torch.optim.Adam([self.log_alpha], lr=3e-4)

431         if temp_model_file is not None:
432             print("Loading pretrained model from", temp_model_file)
433             self.load_pytorch_temp_model(temp_model_file)

435         self.max_action = self.actor.max_action
436         self.discount = discount
437         self.tau = tau
438         self.policy_freq = policy_freq

440         self.total_it = 0

442         self.transform = torch.distributions.transforms.ComposeTransform(
443             [
444                 torch.distributions.transforms.TanhTransform(),
445                 torch.distributions.transforms.AffineTransform(
446                     loc=0, scale=self.max_action
447                 ),
448             ]
449         )

451     def select_action(self, state: torch.Tensor) -> torch.Tensor:

```

```

452     state = torch.FloatTensor(state.reshape(1, -1)).to(device)

453
454     actor_dist = self.actor(state)
455     selected_action = self.transform(actor_dist.rsample())
456     return selected_action.data.numpy().flatten()

457
458 def train(self, replay_buffer, batch_size=256):
459     self.total_it += 1
460     temperature = torch.exp(self.log_alpha)

461
462     # Sample replay buffer
463     state, action, next_state, reward, not_done = replay_buffer.sample(batch_size)

464
465     with torch.no_grad():
466         # Select action according to the policy,
467         target_actor_dist = self.actor_target(next_state)
468         raw_next_action = target_actor_dist.rsample()
469         next_action = self.transform(raw_next_action)
470         target_entropy = (
471             target_actor_dist.entropy()
472             + self.transform.log_abs_det_jacobian(raw_next_action, next_action).sum(
473                 -1
474             )
475         ).unsqueeze(-1)
476         # Compute the target Q value
477         target_Q1, target_Q2 = self.critic_target(next_state, next_action)
478         target_Q = torch.min(target_Q1, target_Q2)
479         target_Q = reward + not_done * self.discount * (
480             target_Q + (temperature * target_entropy)
481         )
482         target_Q = target_Q.detach()
483     # Get current Q estimates
484     current_Q1, current_Q2 = self.critic(state, action)

485
486     # Compute critic loss
487     critic_loss = F.mse_loss(current_Q1, target_Q) + F.mse_loss(
488         current_Q2, target_Q
489     )

490
491     # Optimize the critic
492     self.critic_optimizer.zero_grad()
493     critic_loss.backward()
494     self.critic_optimizer.step()

495
496     # Delayed policy updates
497     if self.total_it % self.policy_freq == 0:
498         # don't have to compute critic gradients
499         for param in self.critic.parameters():
500             param.requires_grad = False

501
502         actor_dist = self.actor(state)
503         raw_action = actor_dist.rsample()
504         selected_action = self.transform(raw_action)
505         actor_entropy = (
506             actor_dist.entropy()
507             + self.transform.log_abs_det_jacobian(raw_action, selected_action).sum(
508                 -1
509             )
510         ).unsqueeze(-1)

511
512         alpha_loss = (
513             temperature * (actor_entropy - self.target_entropy).mean().detach()
514         )
515         actor_loss = -(
516             self.critic.Q1(state, selected_action)
517             + (temperature.detach() * actor_entropy)
518         ).mean()

519
520         alpha_actor_loss = alpha_loss + actor_loss

521
522     # Optimize the actor
523     self.actor_optimizer.zero_grad()
524     self.log_alpha_optimizer.zero_grad()
525     alpha_actor_loss.backward()
526     self.actor_optimizer.step()
527     self.log_alpha_optimizer.step()

```

```

529         for param in self.critic.parameters():
530             param.requires_grad = True

532         # Update the frozen target models
533         for param, target_param in zip(
534             self.critic.parameters(), self.critic_target.parameters()
535         ):
536             new_target_params = (
537                 self.tau * param.data + (1 - self.tau) * target_param.data
538             )
539             target_param.data.copy_(new_target_params)

541         for param, target_param in zip(
542             self.actor.parameters(), self.actor_target.parameters()
543         ):
544             new_target_params = (
545                 self.tau * param.data + (1 - self.tau) * target_param.data
546             )
547             target_param.data.copy_(new_target_params)

549     def save_actor_model(self, filename):
550         # pytorch model object like an instance of Actor
551         checkpoint = {
552             "model_state_dict": self.actor.state_dict(),
553             "optimizer_state_dict": self.actor_optimizer.state_dict(),
554         }
555         torch.save(checkpoint, filename)

557     def save_critic_model(self, filename):
558         # pytorch model object like an instance of Actor
559         checkpoint = {
560             "model_state_dict": self.critic.state_dict(),
561             "optimizer_state_dict": self.critic_optimizer.state_dict(),
562         }
563         torch.save(checkpoint, filename)

565     def save_temp_model(self, filename):
566         # pytorch model object like an instance of Actor
567         checkpoint = {
568             "model_state_dict": self.log_alpha,
569             "optimizer_state_dict": self.log_alpha_optimizer.state_dict(),
570         }
571         torch.save(checkpoint, filename)

573     def load_pytorch_actor_model(self, filename):
574         # now to start up a model with the same trained parameters
575         checkpoint = torch.load(filename)
576         self.actor.load_state_dict(checkpoint["model_state_dict"])
577         self.actor_optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
578         # model.eval()

580     def load_pytorch_critic_model(self, filename):
581         # now to start up a model with the same trained parameters
582         checkpoint = torch.load(filename)
583         self.critic.load_state_dict(checkpoint["model_state_dict"])
584         self.critic_optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
585         # model.eval()

587     def load_pytorch_temp_model(self, filename):
588         # now to start up a model with the same trained parameters
589         checkpoint = torch.load(filename)
590         self.log_alpha = checkpoint["model_state_dict"]
591         self.log_alpha_optimizer.load_state_dict(checkpoint["optimizer_state_dict"])
592         # model.eval()

594 # In[ ]:

596 evaluation_td3 = main(policy_name="TD3")

```

Listing 2: Code for RDE with SAC in the Mountain Car Continuous

```

1  #!/usr/bin/env python
2  # coding: utf-8

4  # In[ ]:

6  get_ipython().system("pip3 install swig > /dev/null 2>&1")

```

```

7  get_ipython().system("pip3 uninstall box2d-py -y > /dev/null 2>&1")
8  get_ipython().system("pip3 install box2d-py > /dev/null 2>&1")
9  get_ipython().system("pip3 install box2d box2d-kengz > /dev/null 2>&1")
10 get_ipython().system("apt install xvfb > /dev/null 2>&1")
11 get_ipython().system("pip3 install pyvirtualdisplay > /dev/null 2>&1")
12 get_ipython().system("pip3 install gym==0.25.0 > /dev/null 2>&1")

14 # In[ ]:

16 import gym
17 import random
18 import numpy as np
19 import torch
20 import torch.nn as nn
21 from torch.distributions import Categorical
22 import torch.nn.functional as F
23 import torch.optim as optim
24 import matplotlib.pyplot as plt
25 import sys
26 from pyvirtualdisplay import Display
27 from IPython import display as disp
28 import copy
29 from typing import Tuple

31 get_ipython().run_line_magic("matplotlib", "inline")

33 # In[ ]:

35 # Replay buffer
36 class ReplayBuffer(object):
37     def __init__(self, state_dim, action_dim, max_size=int(2e5)):
38         self.max_size = max_size
39         self.ptr = 0
40         self.size = 0

42         self.state = np.zeros((max_size, state_dim))
43         self.action = np.zeros((max_size, action_dim))
44         self.next_state = np.zeros((max_size, state_dim))
45         self.reward = np.zeros((max_size, 1))
46         self.not_done = np.zeros((max_size, 1))

48         self.device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

50     def add(self, state, action, next_state, reward, done):
51         self.state[self.ptr] = state
52         self.action[self.ptr] = action
53         self.next_state[self.ptr] = next_state
54         self.reward[self.ptr] = reward
55         self.not_done[self.ptr] = 1.0 - done

57         self.ptr = (self.ptr + 1) % self.max_size
58         self.size = min(self.size + 1, self.max_size)

60     def sample(self, batch_size):
61         ind = np.random.randint(0, self.size, size=batch_size)

63         return (
64             torch.FloatTensor(self.state[ind]).to(self.device),
65             torch.FloatTensor(self.action[ind]).to(self.device),
66             torch.FloatTensor(self.next_state[ind]).to(self.device),
67             torch.FloatTensor(self.reward[ind]).to(self.device),
68             torch.FloatTensor(self.not_done[ind]).to(self.device),
69         )

71 # In[ ]:

73 def init_flags():

75     flags = {
76         "env": "MountainCarContinuous",
77         "seed": 0, # random seed
78         "start_timesteps": 5e3, # total steps of free exploration phase
79         "max_timesteps": 2e5, # maximum length of time steps in training
80         "expl_noise": 0.1, # noise strength in exploration
81         "batch_size": 512,
82         "discount": 0.99,
83         "tau": 0.005, # rate of target update

```



```

84     "policy_freq": 2, # delayed policy update frequency in TD3,
85     "N": 1, # number of agents,
86     "RR": 2, # replay ratio,
87     "T": np.inf, # time steps between agent resets (every 8e4 for RR=1) ,
88     "beta": 50, # action selection coefficient
89 }
90
91 return flags
92
93 def collect_actions(theta, state):
94     actions = []
95     for theta_i in theta:
96         action, entropy, mean, vari = theta_i.select_action(np.array(state))
97         actions.append(torch.from_numpy(action))
98     return actions
99
100 def main(policy_name="DDPG") -> list:
101     """
102     Input:
103     policy_name: str, the method to implement
104     Output:
105     evaluations: list, the reward in every episodes
106     Call DDPG/TD3 trainer and
107     """
108     args = init_flags()
109     env = gym.make(args["env"])
110     env.seed(args["seed"] + 100)
111     env.action_space.seed(args["seed"])
112     torch.manual_seed(args["seed"])
113     np.random.seed(args["seed"])
114
115     state_dim = env.observation_space.shape[0]
116     action_dim = env.action_space.shape[0]
117     max_action = float(env.action_space.high[0])
118     kwargs = {
119         "state_dim": state_dim,
120         "action_dim": action_dim,
121         "max_action": max_action,
122         "discount": args["discount"],
123         "tau": args["tau"],
124     }
125     if policy_name == "TD3":
126         kwargs["policy_freq"] = args["policy_freq"]
127         theta = [TD3(**kwargs) for _ in range(args["N"])]
128     elif policy_name == "DDPG":
129         policy = DDPG(**kwargs)
130
131     replay_buffer = ReplayBuffer(state_dim, action_dim)
132     evaluations = []
133     state, done = env.reset(), False
134     episode_reward = 0
135     episode_timesteps = 0
136     episode_num = 0
137     k = 0
138
139     for t in range(int(args["max_timesteps"])):
140
141         episode_timesteps += 1
142
143         # Select action randomly or according to policy
144         entropy = 0
145         mean = 0
146         vari = 0
147         if t < args["start_timesteps"]:
148             action = env.action_space.sample()
149         else:
150             with torch.no_grad():
151
152                 actions = collect_actions(theta, state)
153                 # compute Q's, then apply softmax
154                 q_sa = torch.hstack(
155                     [
156                         theta[k].critic.Q1(torch.FloatTensor(state), action)
157                         for action in actions
158                     ]
159                 )
160                 # dim

```

```

161         max_q_sa, _ = torch.max(q_sa, dim=0)
162         alpha = args["beta"] / max_q_sa
163         p_select = F.softmax(q_sa / alpha)
164
165         if t == args["start_timesteps"]:
166             print(p_select)
167
168         action = np.random.choice(
169             a=torch.hstack(actions).numpy(), p=p_select.numpy()
170         )
171         action = np.atleast_1d(action)
172
173         # Perform action
174         next_state, reward, done, _ = env.step(action)
175         done_bool = float(done) if episode_timesteps < env._max_episode_steps else 0
176
177         # Store data in replay buffer
178         replay_buffer.add(state, action, next_state, reward, done_bool)
179
180         state = next_state
181         episode_reward += reward
182
183         # Train agent after collecting sufficient data
184         if t >= args["start_timesteps"]:
185             for j in range(args["RR"]):
186                 for theta_i in theta:
187                     theta_i.train(replay_buffer, args["batch_size"])
188
189             if (t % (args["T"] / args["N"])) == 0:
190                 print(k)
191                 # reset just actor or both?
192                 theta[k].actor.reset()
193                 theta[k].critic.reset()
194                 k = (k + 1) % args["N"]
195
196         if done:
197             # +1 to account for 0 indexing. +0 on ep_timesteps since it will increment +1 even if done=
True
198             print(
199                 f"Total T: {t+1} Episode Num: {episode_num+1} Episode T: {episode_timesteps} Reward: {
episode_reward:.3f}"
200             )
201             evaluations.append(episode_reward)
202             entropies = []
203             actions_l = []
204             means = []
205             varis = []
206             # Reset environment
207             state, done = env.reset(), False
208             episode_reward = 0
209             episode_timesteps = 0
210             episode_num += 1
211
212         return evaluations
213
214 # In[ ]:
215
216 # Reference Implementation of Twin Delayed Deep Deterministic Policy Gradients (TD3)
217
218 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
219
220 # Construct the actor/critic network for TD3
221 class Actor_TD3(nn.Module):
222     def __init__(self, state_dim: int, action_dim: int, max_action: float):
223         super(Actor_TD3, self).__init__()
224         self.l1 = nn.Linear(state_dim, 256)
225         self.l2 = nn.Linear(256, 256)
226         self.l3 = nn.Linear(256, 2 * action_dim)
227         self.max_action = max_action
228         self.action_dim = action_dim
229
230     def forward(self, state: torch.Tensor) -> torch.Tensor:
231         a = F.relu(self.l1(state))
232         a = F.relu(self.l2(a))
233         a = self.l3(a)
234
235         mean = self.max_action * torch.tanh(a[:, : self.action_dim])

```

```

236         std = nn.functional.softplus(a[:, self.action_dim :]) + 1e-9

238         return mean, std

240     def reset(self):
241         for layer in self.children():
242             if hasattr(layer, "reset_parameters"):
243                 layer.reset_parameters()

245 class Critic_TD3(nn.Module):
246     def __init__(self, state_dim: int, action_dim: int):
247         super(Critic_TD3, self).__init__()

249         # Q1 architecture
250         self.l1 = nn.Linear(state_dim + action_dim, 256)
251         self.l2 = nn.Linear(256, 256)
252         self.l3 = nn.Linear(256, 1)

254         # Q2 architecture
255         self.l4 = nn.Linear(state_dim + action_dim, 256)
256         self.l5 = nn.Linear(256, 256)
257         self.l6 = nn.Linear(256, 1)

259     def forward(
260         self, state: torch.Tensor, action: torch.Tensor
261     ) -> Tuple[torch.Tensor, torch.Tensor]:
262         sa = torch.cat([state, action], 1)
263         q1 = F.relu(self.l1(sa))
264         q1 = F.relu(self.l2(q1))
265         q1 = self.l3(q1)

267         q2 = F.relu(self.l4(sa))
268         q2 = F.relu(self.l5(q2))
269         q2 = self.l6(q2)
270         return q1, q2

272     def Q1(self, state: torch.Tensor, action: torch.Tensor) -> torch.Tensor:
273         sa = torch.cat([state, action], -1)

275         q1 = F.relu(self.l1(sa))
276         q1 = F.relu(self.l2(q1))
277         q1 = self.l3(q1)
278         return q1

280     def reset(self):
281         for layer in self.children():
282             if hasattr(layer, "reset_parameters"):
283                 layer.reset_parameters()

285 # In[ ]:

287 class TD3(object):
288     def __init__(
289         self,
290         state_dim: int,
291         action_dim: int,
292         max_action: float,
293         discount=0.99,
294         tau=0.005,
295         policy_freq=2,
296         temperature=0.01,
297     ):

299         self.actor = Actor_TD3(state_dim, action_dim, max_action).to(device)
300         self.actor_target = copy.deepcopy(self.actor)
301         self.actor_optimizer = torch.optim.Adam(self.actor.parameters(), lr=3e-4)

303         self.critic = Critic_TD3(state_dim, action_dim).to(device)
304         self.critic_target = copy.deepcopy(self.critic)
305         self.critic_optimizer = torch.optim.Adam(self.critic.parameters(), lr=3e-4)

307         self.max_action = max_action
308         self.discount = discount
309         self.tau = tau
310         self.policy_freq = policy_freq

312         self.total_it = 0

```

```

313         self.temperature = temperature

315     def select_action(self, state: torch.Tensor) -> torch.Tensor:
316         state = torch.FloatTensor(state.reshape(1, -1)).to(device)

318         mean, std = self.actor(state)
319         actor_dist = torch.distributions.Normal(mean, std)
320         selected_action = actor_dist.rsample().clamp(-self.max_action, self.max_action)
321         entropy = actor_dist.entropy()
322         vari = torch.square(std)
323         return (
324             selected_action.data.numpy().flatten(),
325             entropy.data.numpy(),
326             mean.data.numpy(),
327             vari.data.numpy(),
328         )

330     def train(self, replay_buffer, batch_size=256):
331         self.total_it += 1

333         # Sample replay buffer
334         state, action, next_state, reward, not_done = replay_buffer.sample(batch_size)

336         with torch.no_grad():
337             # Select action according to the policy,
338             target_mean, target_std = self.actor_target(next_state)
339             target_actor_dist = torch.distributions.Normal(target_mean, target_std)
340             next_action = target_actor_dist.rsample().clamp(
341                 -self.max_action, self.max_action
342             )
343             target_entropy = target_actor_dist.entropy()
344             # Compute the target Q value
345             target_Q1, target_Q2 = self.critic_target(next_state, next_action)

347             target_Q = torch.min(target_Q1, target_Q2)
348             target_Q = reward + not_done * self.discount * (
349                 target_Q + (self.temperature * target_entropy)
350             )
351             target_Q = target_Q.detach()

353         current_Q1, current_Q2 = self.critic(state, action)
354         critic_loss = F.mse_loss(current_Q1, target_Q) + F.mse_loss(
355             current_Q2, target_Q
356         )

358         # Optimize the critic
359         self.critic_optimizer.zero_grad()
360         critic_loss.backward()
361         self.critic_optimizer.step()

363         # Delayed policy updates
364         if self.total_it % self.policy_freq == 0:

366             # Compute actor loss
367             mean, std = self.actor(state)
368             actor_dist = torch.distributions.Normal(mean, std)
369             selected_action = actor_dist.rsample().clamp(
370                 -self.max_action, self.max_action
371             )

373             actor_loss = -(
374                 self.critic.Q1(state, selected_action)
375                 + (self.temperature * actor_dist.entropy())
376             ).mean()

378             # Optimize the actor
379             self.actor_optimizer.zero_grad()
380             actor_loss.backward()
381             self.actor_optimizer.step()

383             # Update the frozen target models
384             for param, target_param in zip(
385                 self.critic.parameters(), self.critic_target.parameters()
386             ):
387                 new_target_params = (
388                     self.tau * param.data + (1 - self.tau) * target_param.data
389                 )

```

```

390         target_param.data.copy_(new_target_params)

392         for param, target_param in zip(
393             self.actor.parameters(), self.actor_target.parameters()
394         ):
395             new_target_params = (
396                 self.tau * param.data + (1 - self.tau) * target_param.data
397             )
398             target_param.data.copy_(new_target_params)

400 # In[ ]:

402 evaluation_td3 = main(policy_name="TD3")

```

Listing 3: Code for RDE with DQN in the Cart Pole Setting

```

1  #!/usr/bin/env python
2  # coding: utf-8

4  # In[ ]:

6  get_ipython().system("apt-get install x11-utils > /dev/null 2>&1")
7  get_ipython().system("pip install pygame > /dev/null 2>&1")
8  get_ipython().system("apt-get install -y xvfb python-opengl > /dev/null 2>&1")

10 get_ipython().system("pip install gym pyvirtualdisplay > /dev/null 2>&1")

12 # In[ ]:

14 import os
15 import pdb
16 import sys
17 import copy
18 import json
19 import argparse
20 from datetime import datetime

22 import gym
23 import torch
24 import torch.nn as nn
25 import torch.nn.functional as F
26 import numpy as np
27 import matplotlib.pyplot as plt
28 from IPython import display as ipythondisplay

30 class DQN(nn.Module):

32     def __init__(self, input, hidden, output):
33         super(DQN, self).__init__()
34         self.fc1 = nn.Linear(input, hidden)
35         self.fc2 = nn.Linear(hidden, hidden)
36         self.fc3 = nn.Linear(hidden, output)

38     def forward(self, x):
39         x = F.relu(self.fc1(x))
40         x = F.relu(self.fc2(x))
41         return self.fc3(x)

43     def reset(self):
44         for layer in self.children():
45             if hasattr(layer, "reset_parameters"):
46                 layer.reset_parameters()

48 class QNetwork:

50     def __init__(self, args, input, output, learning_rate):
51         self.weights_path = "models/%s/%s" % (
52             args["env"],
53             datetime.now().strftime("%Y-%m-%d_%H-%M-%S"),
54         )

56         # Network architecture.
57         self.hidden = 128
58         self.model = DQN(input, self.hidden, output)

60         # Loss and optimizer.
61         self.optim = torch.optim.Adam(self.model.parameters(), lr=learning_rate)

```

```

62     if args["model_file"] is not None:
63         print("Loading pretrained model from", args["model_file"])
64         self.load_model_weights(args["model_file"])

66     def save_model_weights(self, step, i):
67         # Helper function to save your model / weights.
68         if not os.path.exists(self.weights_path):
69             os.makedirs(self.weights_path)
70         torch.save(
71             self.model.state_dict(),
72             os.path.join(self.weights_path, f"model_{step}_{i}.h5"),
73         )

75     def load_model_weights(self, weight_file):
76         # Helper function to load model weights.
77         self.model.load_state_dict(torch.load(weight_file))

79 # In[ ]:

81 class Replay_Memory:
82     def __init__(self, state_dim, action_dim, memory_size=50000, burn_in=10000):
83         # The memory essentially stores transitions recorder from the agent
84         # taking actions in the environment.

86         # Burn in episodes define the number of episodes that are written into the memory from the
87         # randomly initialized agent. Memory size is the maximum size after which old elements in the
88         # memory are replaced.
89         # A simple (if not the most efficient) way to implement the memory is as a list of transitions.
90         self.memory_size = memory_size
91         self.burn_in = burn_in
92         self.states = torch.zeros((self.memory_size, state_dim))
93         self.next_states = torch.zeros((self.memory_size, state_dim))
94         self.actions = torch.zeros((self.memory_size, 1))
95         self.rewards = torch.zeros((self.memory_size, 1))
96         self.dones = torch.zeros((self.memory_size, 1))
97         self.ptr = 0
98         self.burned_in = False
99         self.not_full_yet = True

100     def append(self, states, actions, rewards, next_states, dones):
101         self.states[self.ptr] = states
102         self.actions[self.ptr, 0] = actions
103         self.rewards[self.ptr, 0] = rewards
104         self.next_states[self.ptr] = next_states
105         self.dones[self.ptr, 0] = dones
106         self.ptr += 1

108         if self.ptr > self.burn_in:
109             self.burned_in = True

111         if self.ptr >= self.memory_size:
112             self.ptr = 0
113             self.not_full_yet = False

115     def sample_batch(self, batch_size=32):
116         # This function returns a batch of randomly sampled transitions - i.e. state, action, reward,
117         # next state, terminal flag tuples.
118         # You will feed this to your model to train.
119         if self.not_full_yet:
120             idxs = torch.from_numpy(np.random.choice(self.ptr, batch_size, False))
121         else:
122             idxs = torch.from_numpy(
123                 np.random.choice(self.memory_size, batch_size, False)
124             )

125         states = self.states[idxs]
126         next_states = self.next_states[idxs]
127         actions = self.actions[idxs]
128         rewards = self.rewards[idxs]
129         dones = self.dones[idxs]
130         return states, actions, rewards, next_states, dones

132 # In[ ]:

134 class DQN_Agent:
135     def __init__(self, args, env):
136         # Create an instance of the network itself, as well as the memory.

```

```

137     # Here is also a good place to set environmental parameters,
138     # as well as training parameters - number of episodes / iterations, etc.

140     # Inputs
141     self.args = args
142     self.env = env
143     self.environment_name = args["env"]
144     self.render = self.args["render"]
145     self.epsilon = args["epsilon"]
146     self.network_update_freq = args["network_update_freq"]
147     self.log_freq = args["log_freq"]
148     self.test_freq = args["test_freq"]
149     self.save_freq = args["save_freq"]
150     self.learning_rate = args["learning_rate"]

152     # Other Classes
153     self.q_network = QNetwork(
154         args,
155         self.env.observation_space.shape[0],
156         self.env.action_space.n,
157         self.learning_rate,
158     )
159     self.target_q_network = QNetwork(
160         args,
161         self.env.observation_space.shape[0],
162         self.env.action_space.n,
163         self.learning_rate,
164     )
165     self.batch = list(range(32))

167     # Save hyperparameters
168     self.logdir = "logs/%s/%s" % (
169         self.environment_name,
170         datetime.now().strftime("%Y-%m-%d_%H-%M-%S"),
171     )
172     if not os.path.exists(self.logdir):
173         os.makedirs(self.logdir)
174     with open(self.logdir + "/hyperparameters.json", "w") as outfile:
175         json.dump((self.args), outfile, indent=4)

177     def epsilon_greedy_policy(self, q_values, epsilon):
178         # Creating epsilon greedy probabilities to sample from.

180         # choose random action a fraction epsilon of the time
181         # and a greedy action the rest of the time
182         sample = np.random.rand()
183         if sample < epsilon:
184             return self.env.action_space.sample()
185         else:
186             return torch.argmax(q_values).item()

188     def greedy_policy(self, q_values):
189         return torch.argmax(q_values).item()

191     def td_estimate(self, state, action):
192         # pass through q_network to get Q values
193         Q_values = self.q_network.model.forward(state)
194         action = action.long()
195         return Q_values.gather(1, action)

197     def td_target(self, reward, next_state, done, discount_factor):
198         # pass through target_q_network and take maximum over Q values
199         Q_values = self.target_q_network.model.forward(next_state)
200         max_Q_values, _ = torch.max(Q_values, dim=1, keepdim=True)

202         # compute td_target
203         return reward + discount_factor * (1 - done) * max_Q_values

205     def train_dqn(self, memory, discount_factor):
206         # Sample from the replay buffer.
207         state, action, rewards, next_state, done = memory.sample_batch(batch_size=32)

209         # Optimization step.
210         # For reference, we used F.smooth_l1_loss as our loss function.
211         self.q_network.optim.zero_grad()
212         loss = F.smooth_l1_loss(
213             self.td_estimate(state, action),

```



```

214         self.td_target(rewards, next_state, done, discount_factor),
215     )
216     loss.backward()
217     self.q_network.optim.step()
218
219     return loss
220
221     def hard_update(self):
222         self.target_q_network.model.load_state_dict(self.q_network.model.state_dict())
223
224     @classmethod
225     def plots(cls, reward, td_error):
226         """
227         Plots:
228         1) Avg Cumulative Test Reward over 20 Plots
229         2) TD Error
230         """
231         reward, time = zip(*rewards)
232         plt.figure(figsize=(8, 3))
233         plt.subplot(121)
234         plt.title("Cumulative Reward")
235         plt.plot(time, reward)
236         plt.xlabel("iterations")
237         plt.ylabel("rewards")
238         plt.legend()
239         plt.ylim([0, None])
240
241         loss, time = zip(*td_error)
242         plt.subplot(122)
243         plt.title("Loss")
244         plt.plot(time, loss)
245         plt.xlabel("iterations")
246         plt.ylabel("loss")
247         plt.show()
248
249     def epsilon_decay(self, initial_eps=1.0, final_eps=0.05):
250         if self.epsilon > final_eps:
251             factor = (initial_eps - final_eps) / 10000
252             self.epsilon -= factor
253
254 # In[ ]:
255
256 def init_flags():
257
258     flags = {
259         "env": "CartPole-v0", # Change to "MountainCar-v0" when needed.
260         "render": False,
261         "train": 1,
262         "frameskip": 1,
263         "network_update_freq": 10,
264         "log_freq": 5,
265         "test_freq": 20,
266         "save_freq": 500,
267         "learning_rate": 5e-4,
268         "memory_size": 50000,
269         "epsilon": 0.5,
270         "model_file": None,
271         "N": 2, # number of agents,
272         "RR": 1, # replay ratio,
273         "T": 8e4, # time steps between agent resets,
274         "beta": 0.01, # action selection coefficient
275     }
276
277     return flags
278
279 def get_action(theta, k, epsilon, beta, state, model):
280
281     actions = []
282     for agent in theta:
283         q_value_i = agent.q_network.model.forward((state.reshape(1, -1)))
284         action = agent.epsilon_greedy_policy(q_value_i, epsilon)
285         actions.append(action)
286
287     # now we choose the max q value over our ensemble of agents
288     q_sa = torch.hstack(
289         [
290             theta[k].q_network.model.forward((state.reshape(1, -1)))[:, a]

```

```

291         for a in actions
292     ]
293 )
294 max_q_sa, _ = torch.max(q_sa, dim=0)
295 alpha = beta / max_q_sa
296 p_select = F.softmax(q_sa / alpha)
297
298 action = np.random.choice(a=actions, p=p_select.numpy())
299
300 return action, p_select
301
302 def test_(
303     theta,
304     env,
305     beta,
306     T,
307     k,
308     N,
309     RR,
310     memory,
311     frameskip,
312     discount_factor,
313     t,
314     model_file=None,
315     episodes=100,
316 ):
317     # Evaluate the performance of your agent over 100 episodes, by calculating cumulative rewards for
318     # the 100 episodes.
319     # Here you need to interact with the environment, irrespective of whether you are using a memory.
320     cum_reward = []
321     td_error = []
322     for count in range(episodes):
323         reward, error, _, _ = generate_episode(
324             theta,
325             mode="test",
326             env=env,
327             epsilon=0.05,
328             beta=beta,
329             T=T,
330             k=k,
331             N=N,
332             RR=RR,
333             memory=memory,
334             frameskip=frameskip,
335             discount_factor=discount_factor,
336             t=t,
337         )
338         cum_reward.append(reward)
339         td_error.append(error)
340     cum_reward = torch.tensor(cum_reward)
341     td_error = torch.tensor(td_error)
342     print(
343         "\nTest Rewards: {0} | TD Error: {1:.4f}\n".format(
344             torch.mean(cum_reward), torch.mean(td_error)
345         )
346     )
347     return torch.mean(cum_reward), torch.mean(td_error)
348
349 def burn_in_memory(
350     theta,
351     epsilon,
352     beta,
353     env,
354     T,
355     k,
356     N,
357     RR,
358     memory,
359     discount_factor,
360     t,
361     mode="train",
362     frameskip=1,
363 ):
364     # Initialize your replay memory with a burn_in number of episodes / transitions.
365     while not memory.burned_in:
366         _, _, t, _ = generate_episode(
367             theta,

```

```

367         mode="burn_in",
368         env=env,
369         epsilon=epsilon,
370         beta=beta,
371         T=T,
372         k=k,
373         N=N,
374         RR=RR,
375         memory=memory,
376         frameskip=frameskip,
377         discount_factor=discount_factor,
378         t=t,
379     )
380     print("Burn Complete!")
381
382     return t
383
384 def generate_episode(
385     theta,
386     epsilon,
387     beta,
388     env,
389     T,
390     k,
391     N,
392     RR,
393     memory,
394     discount_factor,
395     t,
396     mode="train",
397     frameskip=1,
398 ):
399     """
400     Collects one rollout from the policy in an environment.
401     """
402     done = False
403     state = torch.from_numpy(env.reset())
404     rewards = 0
405
406     td_error = []
407     while not done:
408         with torch.no_grad():
409             action, p_select = get_action(theta, k, epsilon, beta, state, mode)
410             i = 0
411             while (i < frameskip) and not done:
412                 next_state, reward, done, info = env.step(action)
413                 next_state = torch.from_numpy(next_state)
414                 rewards += reward
415                 i += 1
416
417             if mode == "train":
418                 if t % 1000 == 0:
419                     print(p_select)
420
421             if mode in ["train", "burn_in"]:
422                 memory.append(state, action, reward, next_state, done)
423             if not done:
424                 state = copy.deepcopy(next_state.detach())
425
426     # Train the network.
427     if mode == "train":
428         t += 1
429
430         for j in range(RR):
431             for theta_i in theta:
432                 theta_i.train_dqn(memory, discount_factor)
433
434             if (t % (T / N)) == 0:
435                 print(k)
436                 theta[k].q_network.model.reset()
437                 theta[k].target_q_network.model.reset()
438                 k = (k + 1) % N
439
440     if td_error == []:
441         return rewards, [], t, k
442     return rewards, torch.mean(torch.stack(td_error)), t, k

```

```

444 def main(render=False):
445     args = init_flags()
446     args["render"] = render

448     if args["env"] == "CartPole-v0":
449         env = gym.make(args["env"], render_mode="rgb_array")
450         discount_factor = 0.99
451         num_episodes = 1000
452         max_timesteps = 2e5
453     else:
454         raise Exception("Unknown Environment")

456     memory = Replay_Memory(
457         env.observation_space.shape[0],
458         env.action_space.n,
459         memory_size=args["memory_size"],
460     )

462     theta = [DQN_Agent(args, env) for _ in range(args["N"])]

464     # time step
465     t = 0
466     # number of episodes
467     step = 0
468     # current agent being reset
469     k = 0

471     burn_in_memory(
472         theta,
473         epsilon=args["epsilon"],
474         beta=args["beta"],
475         env=env,
476         T=args["T"],
477         k=k,
478         N=args["N"],
479         RR=args["RR"],
480         memory=memory,
481         mode="train",
482         frameskip=args["frameskip"],
483         discount_factor=discount_factor,
484         t=t,
485     )

487     rewards = []
488     td_error = []

490     while t < max_timesteps:
491         # Generate Episodes using Epsilon Greedy Policy and train the Q network.
492         step += 1
493         _, _, t, k = generate_episode(
494             theta,
495             mode="train",
496             env=env,
497             epsilon=args["epsilon"],
498             beta=args["beta"],
499             T=args["T"],
500             k=k,
501             N=args["N"],
502             RR=args["RR"],
503             memory=memory,
504             frameskip=args["frameskip"],
505             discount_factor=discount_factor,
506             t=t,
507         )

509     # Test the network.
510     if step % args["test_freq"] == 0:
511         print("here")
512         test_reward, test_error = test_(
513             theta,
514             env=env,
515             beta=args["beta"],
516             N=args["N"],
517             T=args["T"],
518             k=k,
519             RR=args["RR"],
520             memory=memory,

```

```
521         discount_factor=discount_factor,
522         t=t,
523         frameskip=args["frameskip"],
524         model_file=None,
525         episodes=20,
526     )
527     rewards.append([test_reward, step])
528     td_error.append([test_error, step])
529
530     # Update the target network.
531     if step % args["network_update_freq"] == 0:
532         for theta_i in theta:
533             theta_i.hard_update()
534
535     # Logging.
536     if step % args["log_freq"] == 0:
537         print("Step: {0:05d}/{1:05d}".format(step, num_episodes))
538
539     # Save the model
540     if step % args["save_freq"] == 0:
541         for i, theta_i in enumerate(theta):
542             theta_i.q_network.save_model_weights(step, i)
543
544     # step += 1
545     for theta_i in theta:
546         theta_i.epsilon_decay()
547
548     return rewards, td_error
549
550 # In[ ]:
551
552 rewards, td_error = main()
553 DQN_Agent.plots(rewards, td_error)
```