

# Sample Efficient Experience Replay in Non-stationary Environments

Anonymous submission

## Abstract

Reinforcement learning often face non-stationary environments in real-world applications, which requires agents to adapt to changing dynamics by correcting outdated experiences during training. Experience replay is a technique used in reinforcement learning to store and reuse past experiences, which can improve the efficiency of agent training. However, existing methods assign sample priorities based on TD-error, which may not effectively support rapid adaptation in non-stationary environments. This is because samples that have not been sufficiently explored tend to exhibit high TD-error, potentially leading to biased prioritization. In this paper, we introduce a new metric called Environmental Difference Error(ED-error), which quantifies the sample value in non-stationary environments by measuring the magnitude of state transition function changes. Based on ED-error, we propose the Environmental Prioritized Experience Replay(EPER), a prioritized experience replay method to address the challenges posed by non-stationary environments. Experimental results show that EPER outperforms several existing experience replays in various non-stationary settings and tasks within the MiniGrid environment.

## Introduction

Reinforcement learning (RL) finds wide-ranging applications in real-world scenarios[?][?][?], with many of these environments being characterized by non-stationarity. In the context of RL, a non-stationary environment refers to a dynamic system where the underlying properties relevant to RL (e.g. environment dynamics, state transition probabilities) change over time[?]. Unlike stationary environments, where these properties remain constant, non-stationary environments introduce additional complexity as the environment’s behavior evolves. The changes in a non-stationary environment is unpredictable, implying the absence of universally applicable patterns to follow. This requires RL agents to continuously adjust their policy to accommodate the evolving environment[?][?].

Efficient sampling, which maximizing learning efficiency by utilizing a limited number of transactions, is particularly important in non-stationary environments for RL[?][?]. The goal of RL is to maximize the cumulative reward by continuously interacting with the environment and learning from experiences to estimate the Q-function’s accumulated reward. In non-stationary environments, the state transition

function changes over time. Inefficient sampling can decrease the learning efficiency of the agent, potentially leading to convergence to suboptimal policies and even reward collapse(Ditzler et al. 2015). Additionally, in real-world applications, low sample efficiency increases the need for a larger number of expensive transactions, resulting in additional time and financial costs. (Yu 2018).

Experience Replay (ER) is a technique used in RL to improve sample efficiency and learning stability. ER stores transactions in a buffer and randomly samples them for training purposes. As an improvement, the use of temporal difference error(TD-error) as a prioritized sampling criterion, is considered a means to further enhance sample efficiency[???]. TD-error measures the discrepancy between the predicted Q-value and the actual observed reward, which reflects the impact of a certain transaction on the Q-function. Therefore, prioritizing transactions with higher TD-error leads to more informative updates of the Q-function and accelerates learning.

However, using TD error as a priority in non-stationary environments fails to accurately reflect the value of transitions for improving sample efficiency. This is because transitions with high TD error in non-stationary environments can be attributed to the following reasons:

- (a) The transition exists in a state-action space that has undergone changes due to the non-stationarity of the environment.
- (b) The transition is in an unexplored state-action space, a low probability event, or is generated by a random policy (e.g.,  $\epsilon$ -greedy policy).

In non-stationary environments, (a) is considered crucial for adapting quickly to changing environments as it reflects changes in the state transition function, while (b) is beneficial for exploration purposes but hinders rapid adaptation. Therefore, assigning high priority to both (a) and (b) impedes the agent’s ability to quickly adapt its policy in non-stationary environments.

In the paper, we introduces the concept of Environmental Difference Error(ED-error), which is a novel metric based on the impact of state transition function changes on the Q-function. ED-error quantifies the difference between the estimated values of the Q-function before and after the state transition function change.

Based on ED-error, we propose the Environmental Prioritized Experience Replay (EPER), an experience replay method for non-stationary environments, which uses ED-error as a priority metric. ED-error assigns higher priority to samples that reflect changes in the state transition function, which are beneficial for the agent to adapt quickly to non-stationary environments. Thus, EPER is only activated when the stationarity affects the agent.

The second challenge is to prevent the agent from overfitting, which causes the agent to forget the learned policy. After a change in the state transfer function is detected, ED-error usually assigns a higher priority to the newly collected samples. However, if there is an excessive reliance on these new samples for training, the high degree of correlation between them can potentially lead to overfitting of the agent. To this end, EPER uses a proportional sampling method. EPER selects a portion of samples from both before and after the change in the state transition function based on the proportion of samples in each group. This process continues until the newly collected samples fill the entire replay buffer.

We evaluated EPER in the MiniGrid environment[???], including several non-stationary settings and tasks. We compare EPER with multiple baselines, including uniform random sampling, PER[???], and CER[???]. Experimental results demonstrate that EPER outperforms these baselines in non-stationary environments.

## Related Work

One of the main techniques of off-policy reinforcement learning is the experience replay method (Sutton and Barto 2018). This method stores the transitions obtained by the agent in a replay buffer and periodically samples a random and uniform batch of transitions from the replay buffer for parameter updating.

Random uniform sampling ignores the importance of each transition, so a series of sampling methods have been proposed to improve sampling efficiency. Schaul et al. proposed the Prioritized Experience Replay (PER), a sampling method that uses TD as a priority (Schaul et al. 2015). PER uses TD to sample proportionally in the replay buffer, achieving better performance than random uniform sampling. Zhang et al. proposed the Combined Experience Replay (CER), which can be seen as an extreme case of PER (Zhang and Sutton 2017). PER makes the latest transition have a higher probability of being sampled, while in CER the latest transition will definitely be sampled. Brittain et al. proposed the Prioritized Sequence Experience Replay (PSER) (Brittain et al. 2019). This sampling method assigns priority to the latest transition while updating the priority of the remaining transitions in the replay buffer.

In other works, other metrics are used as sampling priorities. Sun et al. propose the Attentive Experience Replay (AER) (Sun, Zhou, and Li 2020), a sampling method that prioritizes transitions in the replay buffer based on their similarity to the current state. The Remember and Forget Experience Replay (ReF-ER) classifies transitions as "near-policy" or "far-policy" by the ratio of the current policy to the past policy and samples only the transitions of near-policy (Novati and Koumoutsakos 2019). Fujimoto et

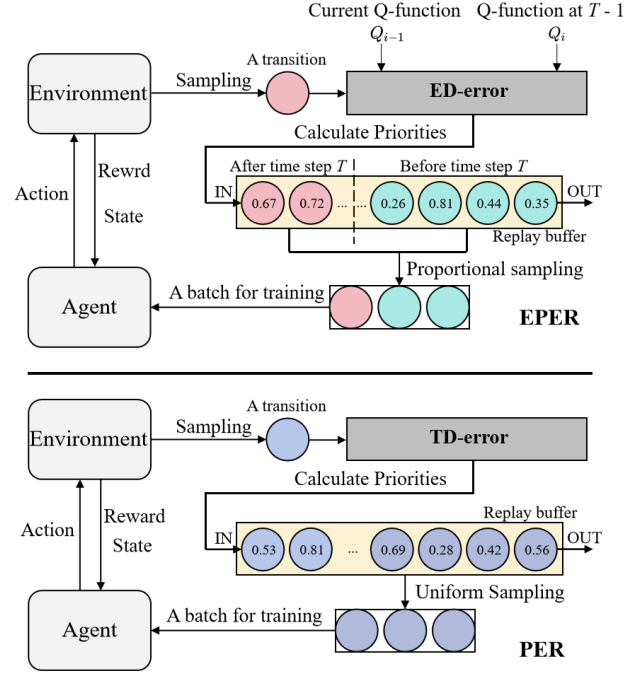


Figure 1: The workflow of EPER

al. equivalently represented the sample priority in PER as a loss function and proposed the Loss-Adjusted Prioritized (LAP) experience replay, a sampling method that optimizes the sample priority through a loss function (Fujimoto, Meger, and Precup 2020).

All of these methods study experience replay in stationary environments and ignore non-stationary environments. The method proposed in the paper is an extension to PER in non-stationary environments. In the rest of the paper, we analyze in detail the problems faced by PER in non-stationary environments and propose our method.

## Methodology

We propose a new experience replay method called EPER, which selects transitions from the replay buffer that reflect environmental non-stationarity to correct the agent's estimates of expected returns. For this, we introduce a novel metric called ED-error in EPER. The ED-error measures the disparity between the Q-function estimation of a transition before the change of the state transfer function and the current Q-function estimation of the same transition. Samples are assigned higher priority if the ED-error indicates a larger difference between the two Q-function estimates. This prioritization mechanism enables EPER to select transitions that reflect significant variations in the environment, ensuring that the agent's estimates of expected returns are appropriately adjusted. Figure 1 shows the workflow of EPER. When the state transition function changes at time step T, EPER stores the Q-function at time step T-1. The stored Q-function is subsequently used to compute the ED-error for newly acquired transitions and transitions in the replay buffer. Then,

EPER selects a proportional subset of transitions from both pre- and post-time step  $T$  to compose a minibatch designated for agent training. Our main contributions are the introduction of ED-error as a priority and the division of the replay buffer into two parts for proportional sampling in EPER, instead of the uniform scanning and sampling of the entire replay buffer in previous work.

## Preliminary

The non-stationary reinforcement learning can be formalized as a Markov decision process(MDP). Specifically, it is defined as a tuple  $\langle \mathcal{S}, \mathcal{A}, P, r, \gamma \rangle$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space,  $P$  is the state transition function,  $r$  is the reward function, and  $\gamma$  is the discount factor. At each time step  $t$ , the agent chooses an action  $a_t \in \mathcal{A}$  to interact with the environment. Then, the environment transits to the next state  $s_{t+1} \in \mathcal{S}$  according to the conditional state transfer function  $P(s_{t+1}, r_t | s_t, a_t)$  and gives a reward  $r_t$  to the agent. The state transition function is defined as:

$$P(s_{t+1}, r_t | s_t, a_t) = \begin{cases} P_0(s_{t+1}, r_t | s_t, a_t), & 0 \leq t < T_0 \\ \dots & \\ P_i(s_{t+1}, r_t | s_t, a_t), & T_{i-1} \leq t < T_i \\ \dots & \\ P_n(s_{t+1}, r_t | s_t, a_t), & T_{n-1} \leq t \leq T \end{cases} \quad (1)$$

Where  $P_i(s_{t+1}, r_t | s_t, a_t)$  denotes the  $i$ -th joint probability distribution of the environment transitioning to state  $s_{t+1}$  and reward  $r_t$  given the state-action pair  $s_t, a_t$ . For brevity, we denote  $P_i(s_{t+1}, r_t | s_t, a_t)$  as  $P_i$ . A non-stationary environment can be regarded as a dynamically changing environment, where the state transition function  $P$  changes with time step  $t$ . When  $n \rightarrow \infty$ , (1) represents scenarios where the state transition function changes at each time step  $t$ , such as fluctuations in natural lighting or temperature.

The goal of the agent is to learn a policy  $\pi$  to maximize the expected discounted return  $\mathbb{E} \left[ \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \right]$ . The value function and Q-function (namely, the action-value function) [??] are proposed to help the agent learn better policies. Determining the state  $s_t$  and the policy  $\pi$ , the value function  $v(s_t)$  is defined as follows:

$$v(s_t) = \mathbb{E}_\pi \left[ \sum_{t'=t}^T \gamma^{t'-t} r_{t'} \right] \quad (2)$$

where  $\mathbb{E}_\pi[\cdot]$  denotes the expectation of a random variable that the agent follows the policy  $\pi$ . Given the state  $s_t$  and action  $a_t$ , the Q-function  $Q(s_t, a_t)$  is defined as the expectation of the value function:

$$Q(s_t, a_t) = \mathbb{E}_\pi [r_t + \gamma v(s_{t+1})] = \sum_{s_{t+1}, r_t} P(s_{t+1}, r_t | s_t, a_t) \left[ r_t + \gamma \mathbb{E} \left[ \sum_{t'=t+1}^T \gamma^{t'-t-1} r_{t'} \right] \right] \quad (3)$$

The optimal policy usually derives from maximizing the expectation of the value function, i.e.  $\pi =$

$\arg \max_{a_t} Q(s_t, a_t)$ . The Q-function  $Q(s_t, a_t)$  is updated by temporal difference learning:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)] \quad (4)$$

where  $\alpha$  is the learning rate.  $r_t + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$  is called the temporal difference error(TD-error) and is denoted as  $\delta^{TD}$  in the paper.

## Sample prioritization based on environmental difference

In previous methods, the priority of a transition is measured by the absolute value of the TD-error. A higher TD-error value implies a greater impact on the agent's Q-function estimation, thus these transitions are considered to have higher priority. During training, transitions with higher priority are selected with a higher probability, which increases the number of training iterations on those important samples and accelerates the convergence of policy. To understand the difference between using TD-error as priority in stationary and non-stationary environments, we first consider a simple and stationary environment, such as the Atari 2600[??] or Grid World[??]. In the environment, given the state  $s$ , action  $a$ , the next state  $s'$  and reward  $r$  are uniquely determined, which is denoted by the state transition function:

$$P(s_{t+1}, r_t | s_t, a_t) = \Pr \{s_{t+1} = s', r_t = r | s_t = s, a_t = a\} = 1 \quad (5)$$

where  $\Pr \{\cdot\}$  denotes the probability. In this case, the environment is stationary, and the agent easily learns each possible pair of next states  $s'$  and rewards  $r$ . Given a transition  $(s_{t-1}, a_{t-1}, r_{t-1}, s_t)$ , the TD-error is denoted as:

$$\delta^{TD} = \sum_{s_{t+1}, r_t} P(s_{t+1}, r_t | s_t, a_t) \left[ \gamma r_t + \gamma^{t'-t+1} \mathbb{E} \left[ \sum_{t'=t+1}^T r_{t'} \right] \right] - \gamma^{t'-t+1} \mathbb{E} \left[ \sum_{t'=t}^T r_{t'} | s_{t-1}, a_{t-1} \right] \quad (6)$$

(6) shows that when transition  $(s_{t-1}, a_{t-1}, r_{t-1}, s_t)$  contains rare events or critical rewards, the transition has high TD-error. The former leads to biased estimation of the state transition function  $P(s_{t+1}, r_t | s_t, a_t)$ , while the latter generates larger disparities in expected returns. Therefore, by assigning higher priority to these rare transitions, the agent can learn crucial experiences that might be overlooked in a uniform random sampling. However, when introducing non-stationary factors into the environment, i.e., the state transition function  $P$  is denoted by (1), the TD-error is denoted as:

$$\begin{aligned} \delta^{TD} &= r_{t-1} + \sum_{s_{t+1}, r_t} P(s_{t+1}, r_t | s_t, a_t) \left[ \gamma r_t + \gamma^{t'-t+1} \right. \\ &\quad \left. \mathbb{E} \left[ \sum_{t'=t+1}^T r_{t'} \right] \right] - \sum_{s_t, r_{t-1}} P(s_t, r_{t-1} | s_{t-1}, a_{t-1}) \cdot \\ &\quad \left[ r_{t-1} + \gamma^{t'-t+1} \mathbb{E} \left[ \sum_{t'=t}^T r_{t'} \right] \right] \end{aligned} \quad (7)$$

(7) shows that the situation is more complex, especially when the state transition function transitions from  $P_{i-1}$  to  $P_i$ . Because the agent's learned experiences are based on the old state transition function  $P_{i-1}$ , it introduces a bias in estimating the expected return, which causes substantial fluctuations in TD-error. Transitions with high TD-error may not represent the most valuable experiences, which can result in incorrect policy updates. Incorrect policy updates cause the agent to take sub-optimal or ineffective actions, which diminish the value of experiences in subsequent transitions. To address the problem, we propose a metric that reflects the difference in expected returns due to a change in the state transfer function. Ideally, if the agent already knows to the expected return of all state-action pairs  $(s, a)$ , then the agent just need to learn the new state transfer function  $P_i$ . We define the metric of the current transition  $(s_{t-1}, a_{t-1}, r_{t-1}, s_t)$  as:

$$\begin{aligned} \delta^{ED} &= \sum_{s_t, r_t} P_{i-1}(s_t, r_{t-1} | s_{t-1}, a_{t-1}) \mathbb{E} \left[ \sum_{t'=t}^T r_{t'} \right] \\ &\quad - \sum_{s_t, r_t} P_i(s_t, r_{t-1} | s_{t-1}, a_{t-1}) \mathbb{E} \left[ \sum_{t'=t}^T r_{t'} \right] \end{aligned} \quad (8)$$

We call the new metric "Environmental Difference Error(ED-error)", denoted as  $\delta^{ED}$ . The intuition behind this is to utilize the estimator  $Q_{i-1}(s_{t-1}, a_{t-1}) - Q_i(s_{t-1}, a_{t-1})$  to approximate  $\delta^{ED}$ , where  $Q_{i-1}$  and  $Q_i$  denote the Q-function under the state transition functions  $P_{i-1}$  and  $P_i$ , respectively. The ED-error quantifies the disparity between  $Q_{i-1}$  and  $Q_i$ , enabling it to guide the agent in gradually transitioning from  $Q_{i-1}$  to  $Q_i$ , thereby adapting to the new state transition function  $P_i$ .

However, it is not feasible to rely on the assumption that the state transition function changes only after the agent has fully explored the entire state-action space in each scenario. Consequently, the agent's estimation of  $\mathbb{E} \left[ \sum_{t'=t}^T r_{t'} \right]$  is also inaccurate. Moreover, the state transition functions  $P_{i-1}$  and  $P_i$  are not explicitly represented in the agent's policy. Thus, we use Q-function approximations at different time steps to estimate the expected returns under each state transition function in ED-error. Specifically, assume that the non-stationary environment is denoted by (1). When time step  $T_{i-1} \leq t \leq T_i$ , we define the ED of the current transition  $(s_{t-1}, a_{t-1}, r_{t-1}, s_t)$  as:

$$\delta^{ED} = r_{t-1} + Q^*(s_t, a_t) - Q(s_{t-1}, a_{t-1}) \quad (9)$$

where  $Q^*$  denotes the Q-function at time step  $T_{i-1} - 1$  and  $Q$  denotes the current Q-function.

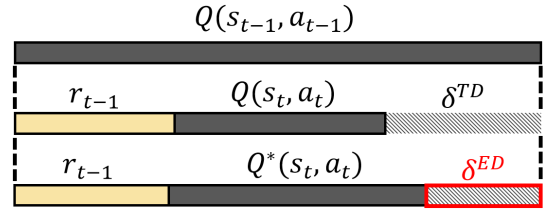


Figure 2: ED-error and TD-error

## Environmental Prioritized Experience Replay

Non-stationary environments is caused by changes in the state transition function. Sudden changes in the state transition function result in inaccurate estimation of expected rewards by the Q-function. The ED-error measures the magnitude of changes in the state transition function of a transition and it indicates the degree to which it surprises the agent's estimate of the expected return. Thus high ED-error transitions can better guide the agent's correction of expected returns. However, using ED-error simply as a priority for sampling has several issues. Firstly, transitions under the current state transfer function usually have a higher ED-error because they directly reflect the degree of instability of the environment. However, there is a high degree of temporal correlation between these transitions, especially when the state transfer function has just changed. This frequent sampling can lead to agent overfitting [????], which in turn reduces the agent's ability to generalize. Secondly, when the state-transfer function changes only locally, the high-ED-error transitions are concentrated in a subspace of the state-action space. This means that the agent will frequently explore in this subspace and ignore the rest.

To address these issues, we propose EPER (see Algorithm 1). EPER does not impose any requirements on the off-policy RL algorithm, which makes EPER easily applicable to many off-policy RL methods (e.g., dqn, ddpg, etc.). EPER saves the Q function for the current time step when it detects a change in the environment and uses ED-error as a priority, and this process continues until the transition after the change in the environment fills the entire replay buffer. Environmental changes are realized through the detection of rewards, which is in fact a change-point detection problem, a problem that has been solved in many works [???]. A simple approach is to infer a change in the state transfer function when a sustained sharp drop in reward is detected. Then, we define the ratio function for the  $j$ -th batch of new samples as:

$$f(j) = \max \left( \epsilon^j, \frac{j \cdot \lceil L/M \rceil}{L} \right) \quad (10)$$

where  $L$  denotes the replay buffer capacity,  $M$  denotes the replay period, and  $\epsilon$  is an exponent that we usually take as 0.9 or 0.95. The probability of sample  $k$  being selected is:

$$P(k) = \frac{p_k^\alpha}{\sum_i p_i^\alpha} \quad (11)$$

where  $\alpha$  is an exponent. Importance-sampling(IS) weights

---

**Algorithm 1:** Environmental prioritized experience replay

---

**Input:** Maximum time step  $T$ , batch size  $N$ , buffer capacity  $L$ , exponents  $\alpha$ ,  $\beta$  and  $\epsilon$ , learning rate  $\eta$ .

```
1: Initialize state  $s_0$  and replay buffer  $\mathcal{D}$  of capacity  $L$ .
2: for  $t = 1$  to  $T$  do
3:   Select and execute  $a_{t-1}$  according to the policy  $\pi$ , observe reward  $r_{t-1}$  and next state  $s_t$ .
4:   Store transition  $(s_{t-1}, a_{t-1}, r_{t-1}, s_t)$  in  $\mathcal{D}$  with maximal priority  $p_t = \max_{i < t} p_i$ 
5:   if environment_change then
6:     Store the Q-function as  $Q^*$ 
7:     Calculate the ED-error for transitions in  $\mathcal{D}$  by Eq.(9) and update the priorities:  $p_k \leftarrow 1/|\delta_k^{ED}|$ 
8:      $j \leftarrow L$ 
9:     if  $j > 0$  then
10:       $j \leftarrow j - 1$ 
11:      Set  $f(j) = \max(\epsilon^j, \frac{j}{L})$ 
12:      Sample a batch of  $N * (1 - f(j))$  and  $N * f(j)$  transitions in  $\mathcal{D}$  by Eq.(11), both before and after the environment change, respectively.
13:      Compute the ED-error for the batch of transitions by Eq.(9) and update the priorities of the transitions after the environment change:  $p_k \leftarrow |\delta_k^{ED}|$ 
14:      for  $k = 1$  to  $N$  do
15:        Compute importance-sampling weight  $w_k = (N * P(k))^{-\beta} / \max_i w_i$ 
16:        Accumulate weight-change  $\Delta \leftarrow \Delta + w_k * \delta_k^{TD} * \nabla_{\theta} Q(s_{k-1}, a_{k-1})$ 
17:        Update weights  $\theta \leftarrow \theta + \eta * \Delta$ , reset  $\Delta = 0$ 
18:      end for
19:    else
20:      Using the TD-error as the priority for experience replay.
21:    end if
22:  end if
23: end for
```

---

are used to correct for bias, and the sampling weight of the final selected sample  $k$  is denoted as:

$$w_k = \left( \frac{1}{N \cdot P(k)} \right)^{\beta} \quad (12)$$

where  $N$  denotes minibatch size and  $\beta$  is an exponent. ED reflects the differences in expected returns due to environmental transition. For a newly acquired sample, ED represents the difference between the post-transition and pre-transition environments reflected by the sample. This difference is exactly the increment that the agent should learn to adapt to the post-transition environment. For samples collected in the pre-transition environment, a larger ED means that the sample responds to more mismatched information. Therefore, the sample priority is defined as  $p_k = |\delta_k^{ED}| + \epsilon$  for samples collected after the environment transition and  $p_k = 1/(|\delta_k^{ED}| + \epsilon)$  for samples collected before the environment transition, where  $\epsilon$  is a small positive constant to prevent  $p_k$  from taking zero or infinity.

## Evaluation

In this section, we examine the sample efficiency of various experience replay methods in reinforcement learning tasks. Our experiments employ a Deep Q-Network (DQN) as the reinforcement learning agent. For all tasks, our DQN model consists of a series of layers: three convolutional layers with 64, 128, and 256 nodes respectively, followed by three fully connected layers with 256, 128, and a number of nodes that

match the number of actions available in the environment. We use the Adam algorithm (Kingma and Ba, 2014) with the learning rate 0.0001. We use a discount factor (gamma) of 0.99 and a batch size of 128.

All experiments are performed on a server with 48 Intel(R) Xeon(R) Silver 4116 CPU @ 2.10GHz and 4 NVIDIA GeForce RTX 2080 Ti 11GB GPU.

## Nonstationary Environment Setups

We evaluated the performance of EPER in the discrete action space using Minigrid(Chevalier-Boisvert, Willems, and Pal 2018), a popular test environment for reinforcement learning research with a collection of 2D grid-world environments with goal-oriented tasks.

To evaluate the effectiveness of the EPER algorithm in non-stationary environments, we introduced specific non-stationary factors at certain time points during the training process. We evaluated the performance of the EPER algorithm and the baseline in non-stationary environments by observing the convergence speed and cumulative rewards of the models after the environment changed. We introduced 2 types of non-stationary factors from the perspectives of environment construction and agent observation during the agent-environment interaction process.

**Observation noise:** Introduce noise to the agent’s environment perception to introduce non-stationarity in the observations. Many works have explored reinforcement learning models in noisy observations(Kilinc and Montana 2018; Shahryari and Doshi 2017; Wang, He, and Tan 2019), as in

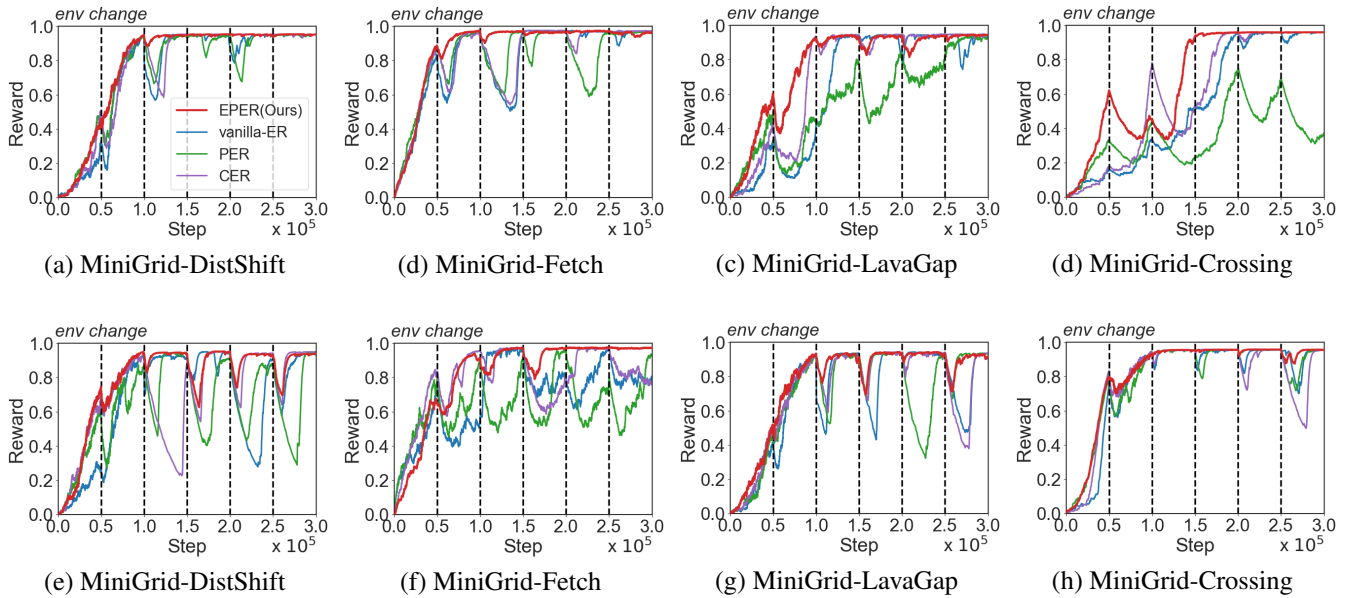


Figure 3: Learning curves of DQN using different sampling methods in 4 non-stationary Minigrid environments. (a)-(d) depict the performance in environments with introduced observation noise. (e)-(h) depict the performance in environments with introduced obstacles. The vertical dashed line indicates the timing of the non-stationarity change in the environment.

many real-world scenarios, observations from sensors such as cameras or microphones may contain inaccuracies due to noise or calibration issues. We follow the approach proposed in (Wang, He, and Tan 2019), where Gaussian noise is generated using a multivariate Gaussian distribution with a given mean and covariance matrix to add noise to the agent’s observations, simulating the nonstationarity commonly encountered in real-world applications.

**Obstacles:** We introduced obstacles into the environment to create nonstationarity in its construction. The number of obstacles introduced depended on the number of empty cells in the map. By incorporating obstacles, we aimed to simulate the nonstationary aspects often encountered in real-world scenarios. When the agent collided with an obstacle, a substantial penalty was subtracted from its reward, and the episode was terminated. This encourages the agent to change its existing strategy to adapt to the environment when it undergoes changes.

## Baseline

We compared our proposed method EPER with several widely used baselines in the field of reinforcement learning that have been shown to improve the performance of deep reinforcement learning algorithms.

**Vanilla-ER (Mnih et al. 2015):** This is the standard experience replay method, where experiences are stored in a replay buffer and randomly sampled during training. The replay buffer is typically limited in size, and new experiences overwrite old ones when the buffer is full.

**Prioritized Experience Replay (PER) (Schaul et al. 2015):** This method assigns a priority to each experience based on its estimated error, with experiences that have higher errors given higher priority. During training, expe-

riences are sampled based on their priority, with higher priority experiences sampled more frequently. This approach has been shown to improve the sample efficiency and overall performance of reinforcement learning algorithms.

**Combined Experience Replay (CER) (Zhang and Sutton 2017):** This method is designed to address the negative influence of a large replay buffer on deep reinforcement learning algorithms. In CER, the latest transition is always added to the replay buffer, while the oldest transition is removed, so the size of the buffer remains constant. The corrected batch is then used to train the agent.

## Comparative Evaluation

Figure X shows the learning curves of different ER methods on four Minigrid non-smooth environments (DistShift, Fetch, LavaGap, Crossing), where the vertical dashed lines indicate the timing of the smoothness change of the environments as we recounted in Section X. The learning curves are shown in Figure X, where (a)-(d) show the change in the environment with the introduction of observation noise, and (e)-(h) show the change in the environment with the introduction of obstacles. In all experiments, the training lasts for  $3e5$  steps, and every  $5e4$  steps a different non-smoothness parameter is introduced, which is reflected in "Observation Noise", as different offsets and noises, and in "Obstacles", as a change of obstacle position. Table X presents the average cumulative rewards obtained during 3,000,000 training steps in the above environments as data. Comparison results show that EPER outperforms the baseline sampling method in 7 out of all 8 test scenarios. It is noteworthy that NERS significantly improves the performance of various nonstrategic RL algorithms in the face of two different noises. These results emphasize that EPER copes well with different types



Environments		Average Episode Reward			
		EPER(Ours)	vanilla-ER	PER	CER
Observation	MiniGrid-DistShift	<b>0.8790</b>	0.8472	0.8435	0.8657
	MiniGrid-Fetch	<b>0.9514</b>	0.9367	0.9217	0.9406
	Noise				
	MiniGrid-LavaGapS6	<b>0.8579</b>	0.8285	0.6591	0.8574
Obstacle	MiniGrid-SimpleCrossing	<b>0.9232</b>	0.6297	0.4189	0.8027
	MiniGrid-DistShift	<b>0.8558</b>	0.7809	0.6378	0.8206
	MiniGrid-Fetch	<b>0.9425</b>	0.8221	0.5256	0.9009
	MiniGrid-LavaGapS6	<b>0.8331</b>	0.7818	0.7613	0.7940
	MiniGrid-SimpleCrossing	<b>0.9277</b>	0.9178	0.9184	0.9130

Table 1: Average episode rewards obtained by different training methods in the non-stationary Minigrid environments, over a total of 3e5 training steps.

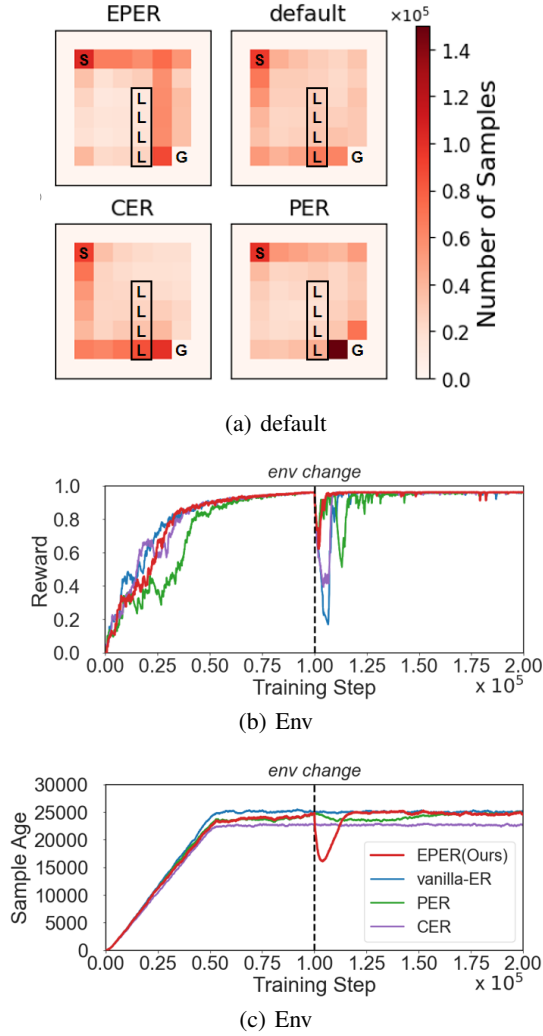


Figure 4: Obstacle

of nonsmoothness.

### Analysis on statistics of samples

We investigated the sampling behavior of EPER during the training process to determine whether EPER can capture informative samples to improve the sample efficiency. In order to determine the sample distribution more intuitively, we chose a most basic scenario Minigrid Empty, as shown in Fig. X(a), in the first stage, the agent needs to reach the end point located at the bottom right corner from the initial position at the top left corner, and in the second stage, we add a 1\*4 magma lattice to introduce non-smoothness. We analyze the changes in sampling behavior by plotting the sampling of the RL agent in the environment using different empirical playback methods.

Figure X(b) shows the position of the agent in the samples selected within 10,000 training steps after the environment change. It can be observed that the samples sampled by the proposed EPER are concentrated near the magma river and reach the end point faster by choosing the path "first to the right, then down". In the vicinity of the magma river, the non-stationarity of the environment introduces an additional challenge, resulting in higher ED errors, so EPER prefers to select these samples for training to better adapt to the new scenario. The results of Vanilla and CER are highly consistent. Vanilla and CER have highly consistent results, as they prefer to sample the "down, then right" route, which is highly rewarded in the first stage, resulting in a higher frequency of samples from this route. However, these samples are no longer relevant in the new scenario, and therefore do not perform well in response to the environmental changes. PER's results show a high concentration of samples near the end of the stage, and a high concentration of samples near the end of the stage.

### Conclusion

This paper proposed EPER, a sample-efficient prioritized experience replay method, which addresses the challenge of applying ER in non-stationary environment: On the one

hand, RL must use as many diverse samples as possible to explore the entire state-action space more comprehensively. On the other hand, RL must focus on a few samples to adapt quickly to the non-stationarity of the environment. We introduce ED, a metric for determining the expected returns between different opponents. EPER detects changes in the environment and using ED error prioritized sampling to balance between new and old transactions. The evaluation results shows that EPER has better sample efficiency than baselines in both discrete action space tasks and continuous action space tasks.

In summary, our proposed EPER method offers an efficient solution to the sample efficiency problem in non-stationary environments. By introducing ED and utilizing prioritized sampling based on it, EPER detects changes in the opponent's behavior and effectively balances between new and old transactions. Our research provides a new approach and method for addressing RL problems in non-stationary environments. Future studies can further explore how to optimize the performance of EPER and apply it to other task, such as multi-agent RL.

## References

- Brittain, M.; Bertram, J.; Yang, X.; and Wei, P. 2019. Prioritized sequence experience replay. *arXiv preprint arXiv:1905.12726*.
- Chevalier-Boisvert, M.; Willems, L.; and Pal, S. 2018. Minimalistic Gridworld Environment for Gymnasium.
- Ditzler, G.; Roveri, M.; Alippi, C.; and Polikar, R. 2015. Learning in nonstationary environments: A survey. *IEEE Computational Intelligence Magazine*, 10(4): 12–25.
- Fujimoto, S.; Meger, D.; and Precup, D. 2020. An equivalence between loss functions and non-uniform sampling in experience replay. *Advances in neural information processing systems*, 33: 14219–14230.
- Kilinc, O.; and Montana, G. 2018. Multi-agent deep reinforcement learning with extremely noisy observations. *arXiv preprint arXiv:1812.00922*.
- Mnih, V.; Kavukcuoglu, K.; Silver, D.; Rusu, A. A.; Veness, J.; Bellemare, M. G.; Graves, A.; Riedmiller, M.; Fidjeland, A. K.; Ostrovski, G.; et al. 2015. Human-level control through deep reinforcement learning. *nature*, 518(7540): 529–533.
- Novati, G.; and Koumoutsakos, P. 2019. Remember and forget for experience replay. In *International Conference on Machine Learning*, 4851–4860. PMLR.
- Schaul, T.; Quan, J.; Antonoglou, I.; and Silver, D. 2015. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*.
- Shahryari, S.; and Doshi, P. 2017. Inverse reinforcement learning under noisy observations. *arXiv preprint arXiv:1710.10116*.
- Sun, P.; Zhou, W.; and Li, H. 2020. Attentive experience replay. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 5900–5907.
- Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement learning: An introduction*. MIT press.
- Wang, Y.; He, H.; and Tan, X. 2019. Robust reinforcement learning in pomdps with incomplete and noisy observations. *arXiv preprint arXiv:1902.05795*.
- Yu, Y. 2018. Towards Sample Efficient Reinforcement Learning. In *IJCAI*, 5739–5743.
- Zhang, S.; and Sutton, R. S. 2017. A deeper look at experience replay. *arXiv preprint arXiv:1712.01275*.