



# Lucid dreaming for experience replay: refreshing past states with the current policy

Yunshu Du<sup>1</sup> · Garrett Warnell<sup>2</sup> · Assefaw Gebremedhin<sup>1</sup> · Peter Stone<sup>3,4</sup> · Matthew E. Taylor<sup>1,5,6</sup>

Received: 16 November 2020 / Accepted: 5 May 2021

© The Author(s), under exclusive licence to Springer-Verlag London Ltd., part of Springer Nature 2021

## Abstract

Experience replay (ER) improves the data efficiency of off-policy reinforcement learning (RL) algorithms by allowing an agent to store and reuse its past experiences in a replay buffer. While many techniques have been proposed to enhance ER by biasing how experiences are sampled from the buffer, thus far they have not considered strategies for *refreshing* experiences inside the buffer. In this work, we introduce *Lucid Dreaming for Experience Replay (LiDER)*, a conceptually new framework that allows replay experiences to be refreshed by leveraging the agent's current policy. LiDER consists of three steps: First, LiDER moves an agent back to a past state. Second, from that state, LiDER then lets the agent execute a sequence of actions by following its current policy—as if the agent were “dreaming” about the past and can try out different behaviors to encounter new experiences in the dream. Third, LiDER stores and reuses the new experience if it turned out better than what the agent previously experienced, i.e., to *refresh* its memories. LiDER is designed to be easily incorporated into off-policy, multi-worker RL algorithms that use ER; we present in this work a case study of applying LiDER to an actor–critic-based algorithm. Results show LiDER consistently improves performance over the baseline in six Atari 2600 games. Our open-source implementation of LiDER and the data used to generate all plots in this work are available at <https://github.com/duyunshu/lucid-dreaming-for-exp-replay>.

**Keywords** Deep reinforcement learning · Experience replay · Self-imitation learning · Behavior cloning

## 1 Introduction

One of the critical components contributing to the recent success of integrating reinforcement learning (RL) with deep learning is the experience replay (ER) mechanism [27]. While deep RL algorithms are often data hungry, ER enhances data efficiency by allowing the agent to store and reuse its past experiences in a replay buffer [24]. Several techniques have been proposed to enhance ER to further reduce data complexity and one of the commonly used techniques is to influence the order of replayed experiences. Instead of replaying experiences uniformly at random [23, 27], studies have found that sampling experiences with different priorities can speed up the learning [7, 31, 37, 40, 47].

Biased experience sampling affects *how* the experiences are replayed. However, it does not consider *what* experience to replay. An experience comprises a state, the action

---

✉ Yunshu Du  
yunshu.du@wsu.edu

Garrett Warnell  
garrett.a.warnell.civ@mail.mil

Assefaw Gebremedhin  
assefaw.gebremedhin@wsu.edu

Peter Stone  
pstone@cs.utexas.edu

Matthew E. Taylor  
matthew.e.taylor@ualberta.ca

<sup>1</sup> Washington State University, Pullman, USA

<sup>2</sup> Army Research Laboratory, Austin, USA

<sup>3</sup> The University of Texas at Austin, Austin, USA

<sup>4</sup> Sony AI, Austin, USA

<sup>5</sup> University of Alberta, Edmonton, Canada

<sup>6</sup> Alberta Machine Intelligence Institute, Edmonton, Canada

taken at that state, and the return<sup>1</sup> obtained by following the agent’s current policy from that state. Existing ER methods usually operate on a fixed set of experiences. That is, once an experience is stored, it remains static inside the buffer until it ages out. An experience from several steps ago may no longer be useful for the current policy to replay because it was generated in the past with a much worse policy. If the agent were given a chance to try again at the same place, its current policy might be able to take different actions that lead to higher returns than what it obtained in the past. *What* the agent should replay is therefore the newer and updated experience, instead of the older one. Given this intuition, we propose in this work *Lucid Dreaming for Experience Replay (LiDER)*, a conceptually new framework that *refreshes* past experiences by leveraging the agent’s current policy, allowing the agent to learn from valuable data generated by its newer self.

LiDER refreshes replay experiences in three steps: First, LiDER moves the agent back to a state it has visited before. Second, LiDER lets the agent follow its current policy to generate a new trajectory from that state. Third, if the new trajectory led to a better outcome than what the agent previously experienced from that state, LiDER stores the new experience into a separate replay buffer and reuses it during training. We refer to this process as “lucid dreaming for experience replay,” because it is as if the agent were “dreaming” about the past and can control the dream to practice again in a past state to achieve better rewards—much like how research in sports science has found that a person’s motor skills can be improved by consciously rehearsing the movements in a lucid dream (e.g., Stumbrys et al. [41]).

One limitation of LiDER is it requires environmental interactions to refresh past states. However, we carefully account for *all* environment interactions, including steps taken to generate new trajectories, and show that LiDER reduces the overall sample complexity of learning compared to methods that do not refresh experiences. LiDER is applicable when a simulator exists for the task—either the task itself is a simulation like a video game or we can build a simulator of the real world—and the simulator is capable of teleporting the agent back to previously visited states and rolling forward in time from there.

The main contributions of this work are as follows:

1. We propose LiDER, a conceptually new framework to *refresh* replay experiences, allowing an agent to revisit

and update past experiences using its current policy in off-policy, multi-worker RL algorithms.

2. LiDER is implemented in an actor–critic-based algorithm as a case study.
3. We experimentally show LiDER outperforms the baseline method (where past experiences were not refreshed) in six Atari 2600 games, including two hard exploration games that are challenging for several RL benchmark algorithms.
4. Analyses and ablation studies help illustrate the functioning of different components of LiDER.
5. Two extensions demonstrate that LiDER is also capable of leveraging policies from external sources, i.e., a policy trained by a different RL algorithm and a behavior cloning policy pre-trained from non-expert human demonstrations.
6. We open-source our implementation of LiDER and the data used to generate all plots in this work for reproducibility at <https://github.com/duyunshu/lucid-dreaming-for-exp-replay>.

## 2 Background

Our algorithm leverages several existing methods, which we briefly review in this section.

### 2.1 Reinforcement learning

We consider an RL problem to be modeled using a Markov decision process, represented by a 5-tuple  $\langle S, A, P, R, \gamma \rangle$ . A state  $s_t \in S$  represents the environment at time  $t$ . An agent learns what action  $a_t \in \mathcal{A}(s)$  to take in  $s_t$  by interacting with the environment. The transition function  $P(s_{t+1}|s_t, a_t)$  denotes the probability of reaching state  $s_{t+1}$  after taking action  $a_t$  at state  $s_t$ . A reward  $r_t \in \mathcal{R} \subset \mathbb{R}$  is given based on  $a_t$  and  $s_{t+1}$ . The goal is to maximize the expected cumulative return  $G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$  from time step  $t$ , where  $\gamma \in [0, 1]$  is a discount factor that determines the relative importance of future and immediate rewards [42].

### 2.2 Asynchronous advantage actor–critic

Policy-based methods such as the asynchronous advantage actor–critic (A3C) algorithm [28] combine a deep neural network with the actor–critic framework. In this work, we leverage the A3C framework to learn both a *policy function*  $\pi(a_t|s_t; \theta)$  (parameterized as  $\theta$ ) and a *value function*  $V(s_t; \theta_v)$  (parameterized as  $\theta_v$ ). The policy function is the *actor* that takes action. The value function is the *critic* that evaluates the quality of the action against a baseline (e.g.,

<sup>1</sup> A one-step reward  $r$  is usually stored instead of the cumulative return (e.g., Mnih et al. [27]). In this work, we follow Oh et al. [32] and store the Monte-Carlo return  $G$ ; we fully describe the buffer structure in Sect. 3.

state value). A3C directly minimizes the policy loss  $L_{\text{policy}}^{a3c}$  as

$$L_{\text{policy}}^{a3c} = \nabla_{\theta} \log(\pi(a_t|s_t; \theta)) (Q^{(n)}(s_t, a_t; \theta, \theta_v) - V(s_t; \theta_v)) + \beta^{a3c} \mathcal{H} \nabla_{\theta} (\pi(s_t; \theta)),$$

where  $Q^{(n)}(s_t, a_t; \theta, \theta_v) = \sum_{k=0}^{n-1} \gamma^k r_{t+k} + \gamma^n V(s_{t+n}; \theta_v)$  is the  $n$ -step bootstrapped value that is bounded by a hyperparameter  $t_{\max}$  ( $n \leq t_{\max}$ ).  $\mathcal{H}$  is an entropy regularizer for policy  $\pi$  (weighted by  $\beta^{a3c}$ ) which helps to prevent premature convergence to sub-optimal policies. The value loss  $L_{\text{value}}^{a3c}$  is

$$L_{\text{value}}^{a3c} = \nabla_{\theta_v} \left( (Q^{(n)}(s_t, a_t; \theta, \theta_v) - V(s_t; \theta_v))^2 \right).$$

The full A3C loss  $L^{a3c}$  given by Mnih et al. [28] is then

$$L^{a3c} = L_{\text{policy}}^{a3c} + \alpha L_{\text{value}}^{a3c}, \tag{1}$$

where  $\alpha$  is a weight for the value loss. A3C’s architecture contains one global policy and  $k$  parallel actor-critic workers. The workers run in parallel and each has its copy of the environment and parameters; each worker updates the global policy asynchronously using the data collected in its own environment. We use the feedforward version of A3C as it runs faster than, but with comparable performance to, the recurrent version [28].

### 2.3 Transformed Bellman operator for A3C

The A3C algorithm uses reward clipping to help stabilize learning. However, Hester et al. [18] showed that clipping rewards to  $[+1, -1]$  results in the agent being unable to distinguish between small and large rewards, thus hurting the performance in the long-term. Pohlen et al. [33] introduced the *transformed Bellman (TB) operator* to overcome this problem in the deep Q-network (DQN) algorithm [27]. The authors consider reducing the scale of the action-value function while keeping the relative differences between rewards which enables DQN to use raw rewards instead of clipping. Pohlen et al. [33] apply a transform function  $h : \mathbb{R} \mapsto \mathbb{R}$  to reduce the scale of  $Q^{(n)}(s_t, a_t; \theta, \theta_v)$  to

$$Q_{TB}^{(n)}(s_t, a_t; \theta, \theta_v) = \sum_{k=0}^{n-1} h(\gamma^k r_{t+k} + \gamma^n h^{-1} V(s_{t+n}; \theta_v)),$$

where

$$h : z \mapsto \text{sign}(z) \left( \sqrt{|z| + 1} - 1 \right) + \varepsilon z$$

and

$$h^{-1} : x \mapsto \text{sign}(x) \left( \left( \frac{\sqrt{1 + 4\varepsilon(|x| + 1 + \varepsilon)} - 1}{2\varepsilon} \right)^2 - 1 \right),$$

and  $\varepsilon$  is a constant that ensures  $h^{-1}$  is Lipschitz continuous with a closed form inverse. Pohlen et al. [33] also prove that the TB operator reduces the variance of the optimization goal while still enabling learning an optimal policy. Given this benefit, our previous work [6] applied the TB operator to A3C, denoted as A3CTB, and showed that A3CTB empirically outperforms A3C.

### 2.4 Self-imitation learning for A3CTB

The *self-imitation learning* (SIL) algorithm [32] is motivated by the intuition that an agent can exploit its own past *good* experiences and thus improve performance. Built upon the actor-critic framework [28], SIL adds a prioritized experience replay buffer  $\mathcal{D} = (S, A, G)$  to store the agent’s past experiences, where  $S$  is a state,  $A$  is the action taken in  $S$ , and  $G$  is the *Monte-Carlo* return from  $S$  (i.e., the return is computed only after a terminal state is reached). In addition to the A3C loss in Equation (1), at each step  $t$ , SIL samples a minibatch from  $\mathcal{D}$  for  $M$  times and optimizes the following off-policy, actor-critic loss  $L_{\text{policy}}^{\text{sil}}$  and  $L_{\text{value}}^{\text{sil}}$ :

$$L_{\text{policy}}^{\text{sil}} = -\log(\pi(a_t|s_t; \theta)) (G_t - V(s_t; \theta_v))_+ \\ L_{\text{value}}^{\text{sil}} = \frac{1}{2} \|(G_t - V(s_t; \theta_v))_+\|^2,$$

where  $G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k} = r_t + \gamma G_{t+1}$  is the discounted cumulative return,  $V$  is the state value. The value of  $(G_t - V(s_t; \theta_v))$  is called the *advantage*, and the max operator  $(\cdot)_+ = \max(\cdot, 0)$  meaning that only experiences with positive advantage values (i.e., good experiences) can contribute to the policy update. The experience buffer is prioritized by  $(G_t - V(s_t; \theta_v))_+$  to increase the chance that a good experience is sampled. The SIL loss  $L^{\text{sil}}$  is then

$$L^{\text{sil}} = L_{\text{policy}}^{\text{sil}} + \beta^{\text{sil}} L_{\text{value}}^{\text{sil}}, \tag{2}$$

where  $\beta^{\text{sil}}$  is a weight for the value loss.

The SIL algorithm minimizes both the A3C loss  $L^{a3c}$  (Equation (1)) and the SIL loss  $L^{\text{sil}}$  (Equation (2)). Minimizing  $L^{a3c}$  lets the agent learn by interacting with the environment and minimizing  $L^{\text{sil}}$  allows the agent to also learn by replaying its past good experiences. In our previous work [6], we leveraged this framework to incorporate SIL into A3CTB, denoted as A3CTBSIL. Specifically, the return  $G_t$  is transformed to the *TB return* using operators  $h$  and  $h^{-1}$  discussed in Sect. 2.3 as:

$$G_t = h(r_t + \gamma h^{-1}(G_{t+1})).$$

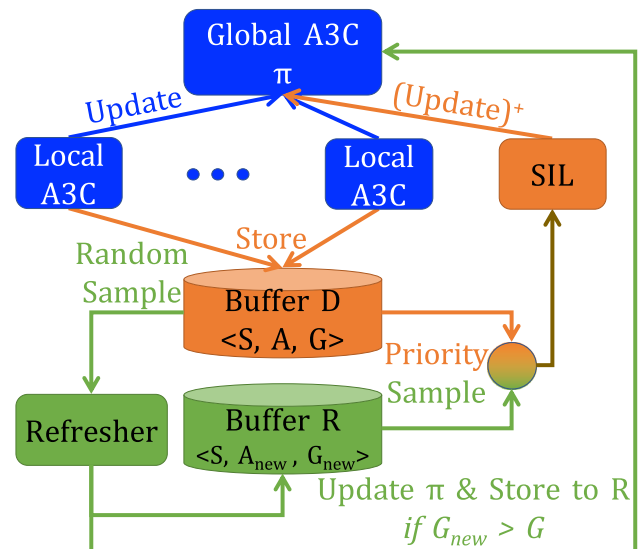
For simplicity, from this point on we will use the word “return” to refer to “TB return.” Our previous work has shown that A3CTBSIL outperformed both the A3C and A3CTB algorithms [6]. This article, therefore, uses an implementation of A3CTBSIL as the baseline.<sup>2</sup>

### 3 Lucid dreaming for experience replay

In this work, we introduce *Lucid Dreaming for Experience Replay (LiDER)*, a conceptually new framework that allows replay experiences to be refreshed by following the agent’s current policy. LiDER consists of three steps: First, LiDER moves an agent back to a past state. Second, from that state, LiDER then lets the agent execute a sequence of actions by following its current policy—as if the agent were “dreaming” about the past and can try out different behaviors to encounter new experiences in the dream. Third, LiDER stores and reuses the new experience if it turned out better than what the agent previously experienced, i.e., to *refresh* its memories. From a high-level perspective, we expect LiDER to help learning by allowing the agent to witness and learn from alternate and advantageous behaviors.

LiDER is designed to be easily incorporated into off-policy, multi-worker RL algorithms that use ER. We implement LiDER in the A3C framework with SIL for two reasons. First, the A3C architecture [28] allows us to conveniently add the “refreshing” component (which we will introduce in the next paragraph) in parallel with A3C and SIL workers, which saves wall-clock time for training. Second, the SIL framework [32] is an off-policy actor-critic algorithm that integrates an experience replay buffer with A3C in a straightforward way, enabling us to directly leverage the return  $G$  of an episode for a policy update—a key component of LiDER.<sup>3</sup>

Figure 1 shows the proposed implementation architecture for LiDER. A3C components are in blue:  $k$  parallel workers interact with their own copies of the environment to update the global policy  $\pi$  [28]. SIL components are in orange: one SIL worker and a prioritized replay buffer  $\mathcal{D}$  are added to A3C [32]. Buffer  $\mathcal{D}$  stores all experiences from the A3C workers in the form of  $\mathcal{D} = \{S, A, G\}$  (as described in Sect. 2). Buffer  $\mathcal{D}$  is prioritized by the advantage value such that good states are more likely to be sampled. The SIL worker runs in parallel with the A3C



**Fig. 1** LiDER architecture. A3C components are in blue and SIL components are in orange. We introduce the novel concept of a refresher worker, in green, to generate new experiences from a randomly sampled past state from buffer  $\mathcal{D}$  by leveraging the agent’s current policy. If the new experiences obtain a higher return than what is currently stored in the replay buffer  $\mathcal{D}$ , they are used to update global policy  $\pi$  and are also stored into replay buffer  $\mathcal{R}$  for reuse

workers but does not interact with the environment; it only samples from buffer  $\mathcal{D}$  and updates  $\pi$  using samples that have positive advantage values (Equation (2)).

We introduce the novel concept of a “refresher” worker in parallel with A3C and SIL to generate new data from past states (shown in green). The refresher has access to the environment and takes randomly sampled states from buffer  $\mathcal{D}$  as input. For each state sampled, the refresher resets the environment to that state and uses the agent’s current policy to perform a rollout until reaching a terminal state (e.g., the agent loses a life). If the Monte-Carlo return of the new trajectory,  $G_{\text{new}}$ , is higher than the previous return,  $G$  (sampled from buffer  $\mathcal{D}$ ), the new trajectory is immediately used to update the global policy  $\pi$ . The update is done in the same way as the A3C workers (Equation (1), replacing  $Q^n$  with  $G_{\text{new}}$ ). The new trajectory is also stored in a prioritized buffer  $\mathcal{R} = \{S, A_{\text{new}}, G_{\text{new}}\}$  (prioritized by advantage, like in buffer  $\mathcal{D}$ ) if  $G_{\text{new}} > G$ . Finally, the SIL worker samples from both buffers as follows. A batch of samples is taken from each of the buffers  $\mathcal{D}$  and  $\mathcal{R}$  (i.e., two batches in total), prioritized by advantage. Samples from both batches are mixed together and put into a temporary buffer, shown in the green-orange circle in Fig. 1; the temporary buffer treats all samples with an equal priority. One batch of samples is then taken (with replacement) from the mixture of the two batches (shown as the brown arrow) and SIL performs updates using the good samples from this

<sup>2</sup> The implementation of A3CTBSIL is open-sourced at <https://github.com/gabrieledcjr/DeepRL>. In de la Cruz Jr et al. [6], we also considered using demonstrations to improve A3CTBSIL, which is not the baseline used in this work.

<sup>3</sup> Note that while the A3C algorithm is on-policy, integrating A3C with SIL makes it an off-policy algorithm (as in Oh et al. [32]).

batch. Note that, although samples in the temporary buffer were initialized with equal priorities, the sampling process is not uniformly at random since we use the implementation of prioritized sampling with stochastic prioritization as describe in Section 3.3 of Schaul et al. [37]. Having this temporary buffer to mix together transitions from buffers  $\mathcal{D}$  and  $\mathcal{R}$  allows the agent to select past and/or refreshed experiences flexibly without needing a fixed sampling strategy. We summarize LiDER’s refresher worker’s procedure in Algorithm 1. Full pseudocode for the A3C and SIL workers is in “Appendix 2.”

## 4 Experiments and analyses

We empirically evaluate LiDER in six Atari 2600 games [3]: Gopher, NameThisGame, Alien, Ms. Pac-Man, Freeway, and Montezuma’s Revenge. We selected these games because they cover a range of properties and difficulties. Based on the Atari game taxonomy defined by Bellemare et al. [2], Gopher and NameThisGame are easy exploration games with dense reward functions; they are relatively easy to learn. Alien and Ms. Pac-Man are hard exploration games with dense reward functions; they are considered to

---

### Algorithm 1 LiDER: Refresher Worker

---

```

1: // Assume shared global policy  $\pi$ , replay buffer  $\mathcal{D}$ , replay buffer  $\mathcal{R}$ 
2: while  $T < T_{max}$  do ▷  $T_{max} = 50$  million
3:   Synchronize refresher’s policy with the global policy:  $\pi_e(\cdot|\theta_e) \leftarrow \pi$ 
4:   Synchronize global step  $T$  from the most recent A3C worker
5:   Initialize  $S \leftarrow \emptyset$ ,  $A_{new} \leftarrow \emptyset$ ,  $R \leftarrow \emptyset$ 
6:   Randomly take a sample  $\{s, a, G\}$  from buffer  $\mathcal{D}$ , reset the environment to  $s$ 
7:   while not terminal do
8:     Execute an action  $s, a, r, s' \sim \pi_e(s|\theta_e)$ 
9:     Store the experience  $S \leftarrow S \cup s$ ,  $A_{new} \leftarrow A_{new} \cup a$ ,  $R \leftarrow R \cup r$ 
10:    Go to next state  $s \leftarrow s'$ 
11:     $T \leftarrow T + 1$ 
12:   end while
13:    $G_{new} = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$ ,  $\forall r \in R$  ▷ Compute the new return
14:   if  $G_{new} > G$  then ▷ Equation (1), replace  $Q^{(n)}$  with  $G_{new}$ 
15:     Update  $\pi$  using  $\{S, A_{new}, G_{new}\}$ 
16:     Store to buffer  $\mathcal{R} \leftarrow \mathcal{R} \cup \{S, A_{new}, G_{new}\}$ 
17:   end if
18: end while

```

---

The main benefit of LiDER is that it allows an agent to leverage its *current* policy to refresh past experiences. However, LiDER does require the refresher to use additional environmental steps (see Algorithm 1 line 11: we account for the refresher steps when measuring the global steps), which can be concerning if acting in the environment is expensive. Despite this shortcoming, we show in our experiments (Sect. 4) that the learning speedup LiDER provides actually reduces the overall number of environment interactions required. It seems that the high *quality* of the refreshed experiences compensates for the additional *quantity* of experiences an agent needs to learn. That is, by leveraging the refresher worker, LiDER can achieve a certain level of performance within a shorter period of time compared to without the refresher—an important benefit as RL algorithms are often data hungry.

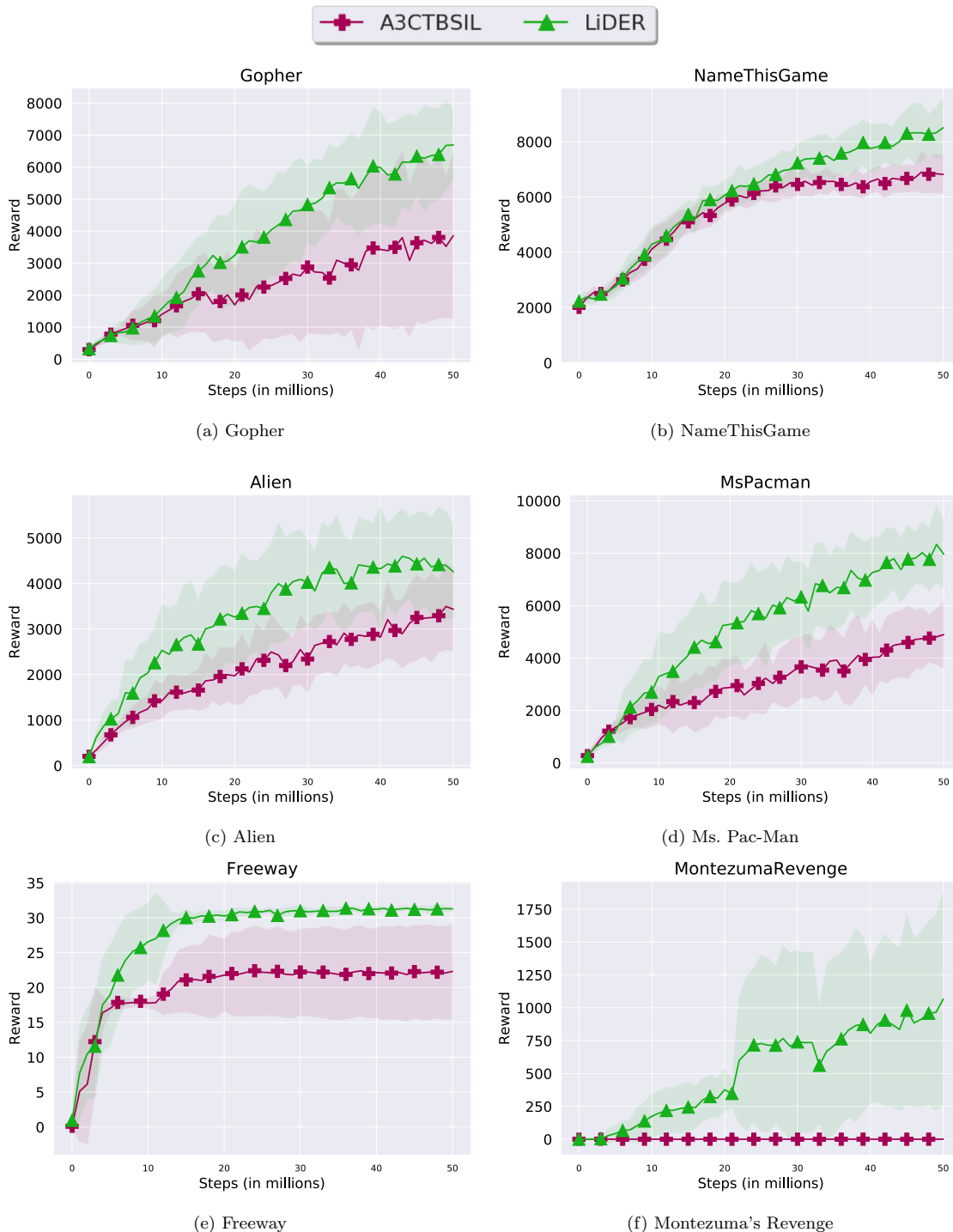
be hard games. Freeway and Montezuma’s Revenge are also hard exploration games but with sparse reward functions; they are considered the hardest games and are challenging for several benchmark RL algorithms (e.g., Bellemare et al. [3], Espeholt et al. [11], and Mnih et al. [28]).

In the next subsection, we compare A3CTBSIL (the baseline method from de la Cruz Jr et al. [6], which uses only the blue and the orange components in Fig. 1) and LiDER (our proposed framework in which the agent’s current policy is used as the refresher) to show LiDER outperforms A3CTBSIL in all games (See “Appendix 1” for implementation details). Section 4.2 then introduces analyses to understand why LiDER helps learning. In Sect. 5, we conduct three ablation studies to validate that our design choices for LiDER were well-founded. Finally, in Sect. 6, we present two extensions and show that LiDER can leverage other policies, rather than its current policy, to refresh past states.

### 4.1 Leveraging the current policy to refresh past states

First, we show that the agent’s current policy can be effectively leveraged to refresh past experiences. Figure 2

shows LiDER outperforms A3CTBSIL in all six games (averaged over eight trials); a one-tailed independent-samples t test confirms statistical significance ( $p \ll 0.001$ , see “Appendix 3” for details of the t tests). We train each trial for 50 million environmental steps. For every 1



**Fig. 2** LiDER performance is compared to A3CTBSIL on six Atari games. The x-axis is the total number of environmental steps: A3CTBSIL counts steps from 16 A3C workers, while LiDER counts

steps from 15 A3C workers plus one refresher worker. The y-axis is the average testing score over eight trials; shaded regions show the standard deviation

million steps, we perform a test of 125,000 steps and report the average testing scores per episode (an episode ends when the agent loses all its lives).

We hypothesize that the performance improvement in the four dense reward games (Gopher, NameThisGame, Alien, and Ms. Pac-Man) was because the likelihood for the refresher to encounter higher-return new trajectories is higher when rewards are dense. In addition, we observe in Ms. Pac-Man that once the return and the action of a state have been refreshed, LiDER prefers to sample and reuse the newer rather than the older state-action-return transition from the same state, which could be another reason for the speed-up in learning—LiDER replays high-rewarding data more frequently. We conduct a detailed analysis of LiDER’s underlying behaviors in the next subsection that supports this hypothesis.

LiDER also learns well in Freeway and Montezuma’s Revenge, the two hard exploration, sparse reward games. In Freeway, the task is difficult because the agent only receives a nonzero reward after successfully crossing the highway. We hypothesize that LiDER is helpful in this case because the refresher can move the agent to an intermediate state (e.g., in the middle of the highway), which shortens the distance between the agent and the rewarding state, and thus allows the agent to learn faster. We can see LiDER’s learning curve in Freeway from Fig. 2e that it consistently finds an optimal path after about 15 million steps of training (the standard deviation becomes negligible) but A3CTBSIL struggles to find a stable solution. The benefit of LiDER is evident particularly in Montezuma’s Revenge. While A3CTBSIL fails to learn anything,<sup>4</sup> LiDER is capable of reaching a reasonable score. Although the absolute performance of our method is not state-of-the-art, we have shown that LiDER is a light-weight addition to a baseline off-policy deep RL algorithm which helps improving performance even in the most difficult Atari games.

## 4.2 Analyses: Why does LiDER help learning?

To understand why LiDER helps improve learning, in this section, we analyze the underlying behavior of LiDER from three perspectives. First, we look at the behaviors of the refresher worker since it is the novel component of LiDER. Inspecting whether the refresher worker can successfully generate higher return trajectories from past states, and how much better the refreshed data is compared to the older data, will give us insight into the quality of the data stored in buffer  $\mathcal{R}$ .

<sup>4</sup> Note the performance in Montezuma’s Revenge differs between A3CTBSIL [6] and the original SIL algorithm [32]—see the discussion in “Appendix 4.”

Second, we examine the SIL worker in LiDER to reveal how SIL makes use of the refreshed data stored in buffer  $\mathcal{R}$ . Not only is it critical for the refresher to be able to generate better data, but the SIL worker must be able to effectively leverage these data to improve learning. The SIL worker should use data from buffer  $\mathcal{R}$  more often for policy updates since the refreshed data is of a higher quality.

Third, we compare the SIL worker between A3CTB-SIL and LiDER. In A3CTBSIL, the SIL worker samples only from one buffer; in LiDER, there are two buffers to sample from. It is thus interesting to investigate whether the SIL worker uses data from the two buffers differently. For example, if samples from buffer  $\mathcal{R}$  have higher returns, we should see more samples with positive advantages in LiDER than in A3CTBSIL. The game of Ms. Pac-Man is used as the running example for all analyses in this section.

### 4.2.1 The refresher worker in LiDER

First, we check two quantities of the refresher worker to get insight into the quality of the data it generated:

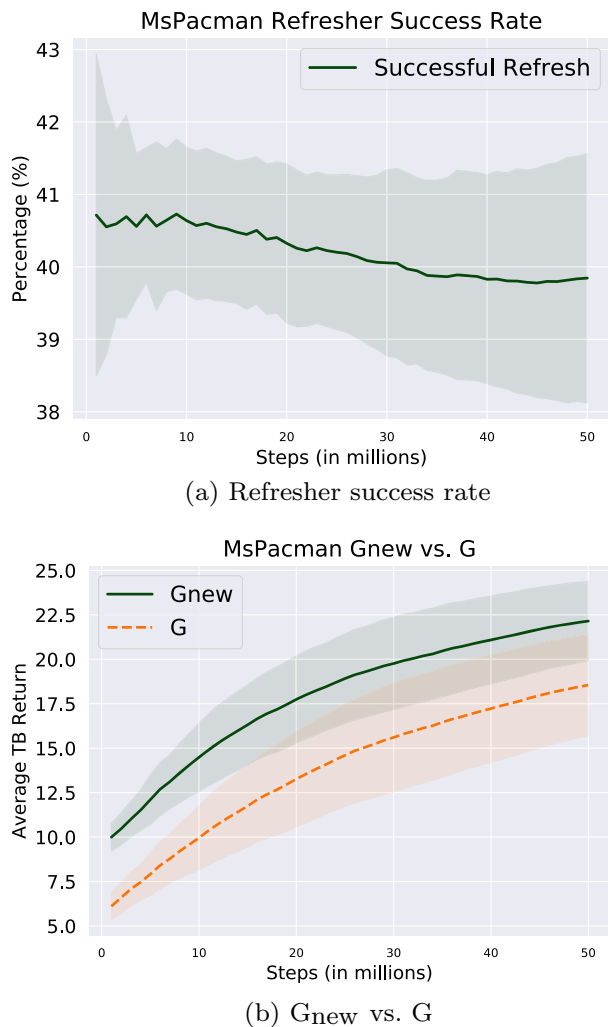
- *Success rate* (Fig. 3a): how often can the refresher worker generate a better trajectory such that  $G_{\text{new}} > G$ .
- $G_{\text{new}}$  vs.  $G$  (Fig. 3b): the average TB return  $G_{\text{new}}$  compared to  $G$  for all successful refresh.

The *success rate* is measured as the percentage of the number of successful rollouts over the total number of rollouts generated. Figure 3a shows that the success rate remains at approximately 40%, indicating that the refresher is able to consistently produce higher return trajectories throughout the training.

The improvement of the refreshed data over the older data can be measured by comparing  $G_{\text{new}}$  to  $G$ .  $G$  is the (TB) return of a state  $S$  sampled from buffer  $\mathcal{D}$ ;  $G_{\text{new}}$  is the refreshed (TB) return of  $S$ . We record the value of  $G_{\text{new}}$  and  $G$  for all successful rollouts then compute their average value. Figure 3b shows that  $G_{\text{new}}$  is indeed higher than  $G$ . Both measures from Fig. 3a, b validate that the refresher worker is able to generate new trajectories with a higher return, and thus data in buffer  $\mathcal{R}$  is expected to be of better quality than data in buffer  $\mathcal{D}$ .

### 4.2.2 The SIL worker in LiDER

We have shown the refresher is able to generate higher return trajectories from past states. Next, we analyze the behavior of the SIL worker in LiDER to check whether it can effectively leverage these data. We inspect the following quantities for buffer  $\mathcal{D}$  and buffer  $\mathcal{R}$ :



**Fig. 3** LiDER’s refresher worker can consistently produce higher return trajectories. The x-axis is the total number of environmental steps. The y-axis value is averaged over eight trials; shaded regions show the standard deviation

- *Old samples used* (Table 1): how many old samples were still used for SIL updates even after the sample has been refreshed to a newer return.
- *Batch sample usage ratio* (Fig. 4a): for one batch of samples, how many samples taken from buffer  $\mathcal{D}$  and buffer  $\mathcal{R}$  were used for SIL updates (i.e., samples with positive advantages), respectively.
- *SIL sample usage ratio* (Fig. 4b): for samples used for SIL updates, how many of them were taken from buffer  $\mathcal{D}$  and buffer  $\mathcal{R}$ , respectively.
- *Return of used samples* (Fig. 4c): the return of samples used for SIL updates.

As mentioned in Sect. 4.1, we hypothesize that once a state’s return has been refreshed, LiDER tends not to reuse the older return. We investigate whether this hypothesis holds by counting how many older samples were used for

SIL updates. Specifically, we assign a *False* Boolean value to each state in buffer  $\mathcal{D}$ . Once a state has been sampled as the input to the refresher worker, we flip the Boolean to *True* for that state. For each SIL update, we count the number of samples with Boolean *True* and compute the ratio of *old samples used* over the total number of samples used for that update.

We show the percentage of the old samples used at 1, 25, and 50 million steps of training in Table 1. It can be seen that less than 0.01% of the older samples were reused throughout training, and the reuse ratio keeps decreasing as training continues. This evidence validates our observation from Sect. 4.1 that LiDER replays the higher-return data in buffer  $\mathcal{R}$  more frequently than the lower-return data in buffer  $\mathcal{D}$ .

Recall that the SIL worker only uses samples with positive advantages (i.e., samples that “pass” the max operator) to update the policy. The percentage of such positive samples in buffer  $\mathcal{D}$  and  $\mathcal{R}$ , respectively, can tell us which buffer has higher quality data. We call this percentage “sample usage ratio” and measure two types of ratios: *batch* and *SIL sample usage ratio*.

*Batch sample usage ratio* measures how many positive samples are from buffer  $\mathcal{D}$  and  $\mathcal{R}$ , respectively, over one batch of samples (the batch size is 32 in our experiments). For example, in one batch of 32 samples, suppose that there were 16 samples with positive advantages that were taken from buffer  $\mathcal{R}$ , and that there were 8 positive samples from buffer  $\mathcal{D}$ . The batch sample usage ratio for buffer  $\mathcal{R}$  is computed as  $\frac{16}{32} = 50\%$ , and for buffer  $\mathcal{D}$  is computed as  $\frac{8}{32} = 25\%$ . Figure 4a shows that, on average, there are more positive samples from buffer  $\mathcal{R}$  than from buffer  $\mathcal{D}$  in one batch of data. This trend indicates that buffer  $\mathcal{R}$ ’s samples are more useful for SIL updates throughout training, and are thus of a higher quality than samples in buffer  $\mathcal{D}$ .

*SIL sample usage ratio* also measures the ratio of positive samples for each buffer, but it is computed over the total number of positive samples instead of the entire batch. For example, suppose that in one batch of 32 samples, there were 24 with positive advantages. 18 out of the 24 samples come from buffer  $\mathcal{R}$  and the other 6 come from buffer  $\mathcal{D}$ . The SIL sample usage ratio for buffer  $\mathcal{R}$  is computed as  $\frac{18}{24} = 75\%$ , and buffer  $\mathcal{D}$ ’s ratio is  $\frac{6}{24} = 25\%$ . Figure 4b shows that buffer  $\mathcal{R}$  always has a higher proportion of positive samples than buffer  $\mathcal{D}$  among all positive samples. Both Fig. 4a, b indicate that the SIL worker is able to effectively leverage the refreshed samples for policy updates.

We can also confirm that the *return of used samples* from buffer  $\mathcal{R}$  is indeed higher than those from buffer  $\mathcal{D}$ . Similar to how  $G_{\text{new}}$  and  $G$  were compared, we compare



**Table 1** Old samples used (%) by the SIL worker in LiDER at 1, 25, and 50 million training steps. LiDER rarely reuses an older state after it has been refreshed. Results were averaged over eight trials

Steps (in millions)	1	25	50
Old samples used (%)	0.0103	0.0053	0.0047
Standard deviation (%)	0.0029	0.0005	0.0003

the average (TB) return of all samples used for SIL updates between buffer  $\mathcal{R}$  and buffer  $\mathcal{D}$ . Figure 4c shows that data in buffer  $\mathcal{R}$  has a higher return than data in buffer  $\mathcal{D}$  during earlier stages of training. The two values then become similar at the end of training—an expected observation as the agent has learned a stable policy.

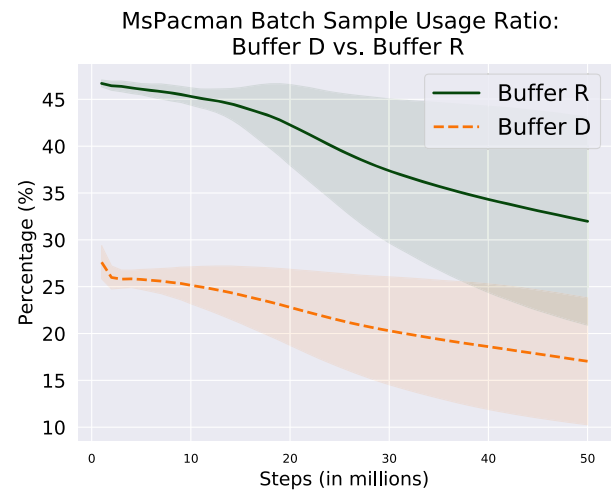
### 4.2.3 The SIL worker in A3CTBSIL versus LiDER

Lastly, we compare the SIL worker between A3CTBSIL and LiDER. We have seen in the previous subsection that the SIL worker in LiDER always prefers to use samples in buffer  $\mathcal{R}$ , which allows more, and higher-quality, data to be leveraged for policy updates. It is thus interesting to inspect what kind of data has been used in A3CTBSIL, and whether the data is better or worse than the data in LiDER. As done in the previous subsection, we examine the *batch sample usage ratio* and *return of used samples* for A3CTBSIL and LiDER.

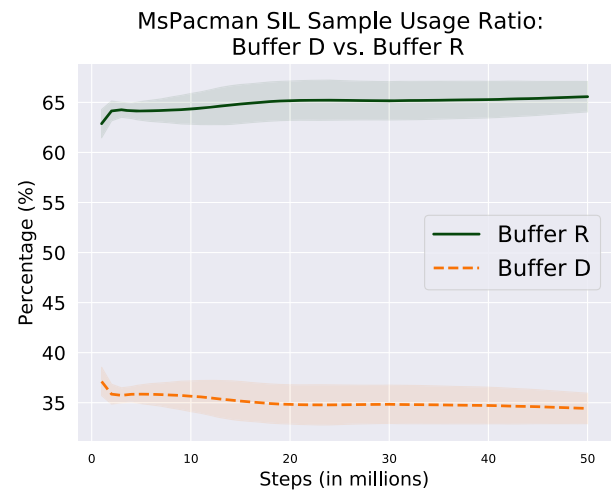
For A3CTBSIL, the batch sample usage ratio and return of used samples are measured in buffer  $\mathcal{D}$  only. For LiDER, we make a small modification that instead of quantifying the two buffers separately, we treat them as one buffer and measure their values together. For example, to compute the batch sample usage ratio, suppose that 6 buffer  $\mathcal{D}$  samples and 18 buffer  $\mathcal{R}$  samples were used for a SIL update, the total batch sample usage ratio for LiDER would be  $\frac{(18+6)}{32} = 75\%$ .

Figure 5a shows that, LiDER has a higher overall batch sample usage ratio than A3CTBSIL. We can also confirm the average (TB) return of used samples is higher in LiDER than in A3CTBSIL (Fig. 5b). This observation indicates that not only can the refresher generate higher return trajectories, but these trajectories are also effectively leveraged by the SIL worker. Both factors contribute to the performance improvement of LiDER over A3CTBSIL.

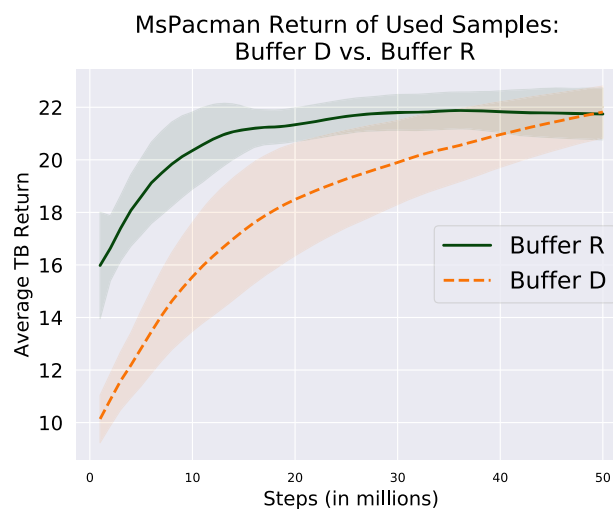
In summary, our analyses show that 1) the refresher can consistently generate good trajectories during training; 2) the SIL worker of LiDER can effectively leverage these good trajectories by sampling from buffer  $\mathcal{R}$ ; and 3) when compared to A3CTBSIL, LiDER performs policy updates with more and higher quality data. All three components



(a) Batch sample usage ratio

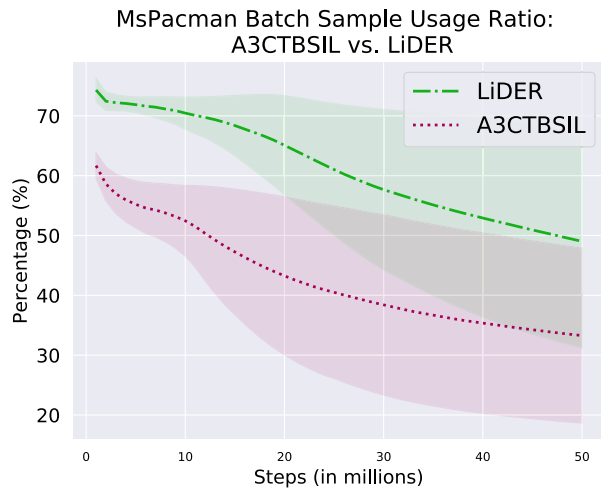


(b) SIL sample usage ratio

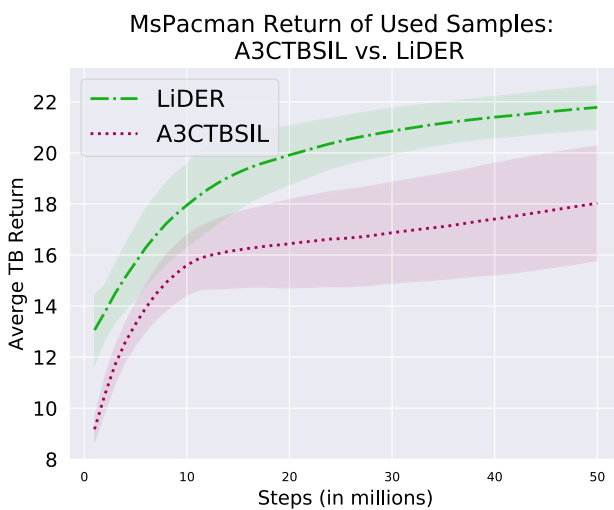


(c) Return of used samples

**Fig. 4** The SIL worker in LiDER leverages more refresher-generated data in buffer  $\mathcal{R}$  than A3C-generated data in buffer  $\mathcal{D}$  because refreshed data has higher returns. The x-axis is the total number of environmental steps. The y-axis value is averaged over eight trials; shaded regions show the standard deviation



(a) Total sample usage ratio



(b) Return of used samples

**Fig. 5** Comparing the SIL worker between A3CTBSIL and LiDER. LiDER leverages more and better quality data for policy updates than A3CTBSIL. The x-axis is the total number of environmental steps. The y-axis value is averaged over eight trials; shaded regions show the standard deviation

contribute to the performance improvement of LiDER over A3CTBSIL.

## 5 Ablation studies

We have shown that LiDER can effectively leverage knowledge from the agent's current policy. In this section, we perform several ablation studies to further validate our design choices.

### 5.1 How does the quality of refresher-generated data affect learning?

As shown in Fig. 3b, LiDER increases the overall data quality because we only use new data when it obtains a higher reward than the old data, i.e.,  $G_{\text{new}} > G$ . We show that it is important to store the refresher-generated experiences and use them to update the global policy only if those experiences are better, i.e., when the new return  $G_{\text{new}}$  computed from the refresher experience is higher than the return  $G$  that the agent previously obtained. This condition ensures that the data in buffer  $\mathcal{R}$  is of a higher quality than that in buffer  $\mathcal{D}$ . Intuitively, LiDER goes back in time to test if its current self can perform better than before and only provide help *where it can*. To validate the importance of this condition, we conduct an experiment in which the refresher adds *all* new experiences to buffer  $\mathcal{R}$ , i.e., without the  $G_{\text{new}} > G$  condition, to check if doing so leads to decreased performance. We denote this experiment as LiDER-AddAll.

### 5.2 How does the buffer architecture affect learning?

The other important design choice of LiDER is the two-buffer architecture: buffer  $\mathcal{D}$  stores A3C-generated data and buffer  $\mathcal{R}$  stores refresher-generated data. One hypothesis could be that LiDER performs better simply because the buffer size is doubled and more experiences can be replayed (e.g., Zhang and Sutton [48] have shown that buffer size can affect learning). We conduct an experiment to show that simply increasing the size of a single buffer does not provide the same performance improvement as LiDER. We modify LiDER to have only buffer  $\mathcal{D}$  and double its size from  $10^5$  to  $2 \times 10^5$ ; both A3C-generated and refresher-generated data are stored in buffer  $\mathcal{D}$ . Prioritized sampling still takes a batch of 32 samples from buffer  $\mathcal{D}$  as the input to the SIL worker, but without using the temporary buffer. We denote this experiment as LiDER-OneBuffer.

### 5.3 How does the sampling ratio affect learning?

LiDER samples from buffer  $\mathcal{D}$  and  $\mathcal{R}$  in a flexible manner as described in Sects. 3, and 4.2.2 we have shown the

samples from buffer  $\mathcal{R}$  are more likely to be used for learning because they have higher returns. The question then arises, “Should the agent always sample from buffer  $\mathcal{R}$  since the refresher-generated data is better?” We conduct an experiment in which the agent only samples from buffer  $\mathcal{R}$ ; a batch of 32 samples are sampled with priority from buffer  $\mathcal{R}$  as the input of the SIL worker, but without using the temporary buffer. Note that although we do not sample from buffer  $\mathcal{D}$ , we still keep it in the architecture since the refresher worker needs to randomly select a past state from buffer  $\mathcal{D}$  to perform the refresh. We denote this experiment as LiDER-SampleR.<sup>5</sup>

## 5.4 Results

Figure 6 shows the results of all ablation studies compared to A3CTBSIL and LiDER. The performance of LiDER-AddAll degraded in four out of six games, except for in Gopher and NameThisGame, where LiDER-AddAll performs comparably to LiDER. This could be because they are easy-exploration games with dense reward functions (as categorized by Bellemare et al. [2]), thus the refresher is more likely to generate better trajectories in these games; adding a few “bad” samples (i.e.,  $G_{\text{new}} \leq G$ ) does not hurt the general performance. In Alien and Montezuma’s Revenge, LiDER-AddAll performs at about the same level as the baseline A3CTBSIL method. Ms. Pac-Man shows the least amount of performance drop for LiDER-AddAll, but it still under-performed LiDER. In Freeway, while LiDER-AddAll eventually reaches the same score as LiDER, it struggled during the early stages of training. These results demonstrate the importance of focusing the exploitation only on places where the refresher can do better than what the agent had previously experienced.

In all games, LiDER-OneBuffer significantly under-performed LiDER ( $p \ll 0.001$ ). Especially in the game of Gopher, NameThisGame, and Ms. Pac-Man where they also performed worse than the baseline A3CTBSIL. These results confirm our analysis in Sect. 4.2.2 that the SIL worker must be able to effectively leverage the high-quality data generated by the refresher to improve learning. When mixing the refresher-generated data and the A3C-generated data into one buffer, it is less likely for the SIL to sample from the good data. Thus, our design of the two-buffer architecture was well chosen.

LiDER-SampleR’s performance was significantly worse than LiDER in five out of six games ( $p \ll 0.001$ ). Except for in Freeway where LiDER-SampleR eventually reaches the same performance as LiDER—but it acts quite unstable (the variance is high). We hypothesize that sampling

only from buffer  $\mathcal{R}$  reduces the amount of state the agent can experience, leading to a lack of exploration which impairs the learning. Therefore, despite that the refresher can generate higher quality data, the agent should learn from both A3C-generated data and refresher-generated data.

In summary, in this section we presented three ablation studies to show the benefits of our design choices of LiDER. Using only experiences where the return is improved, the two-buffer architecture, and the flexible sampling strategy between buffer  $\mathcal{D}$  and buffer  $\mathcal{R}$  indeed improve performance.

## 6 Extensions: leveraging other policies to refresh past states

So far, we have shown in Sect. 4.1 that LiDER outperformed the baseline A3CTBSIL. The analyses in Sect. 4.2 revealed why LiDER helps learning. Section 5 validated the design choices of LiDER through three ablation studies. In this section, we present two extensions to show that LiDER can leverage not only the agent’s current policy, but also policies from external sources to refresh past states.

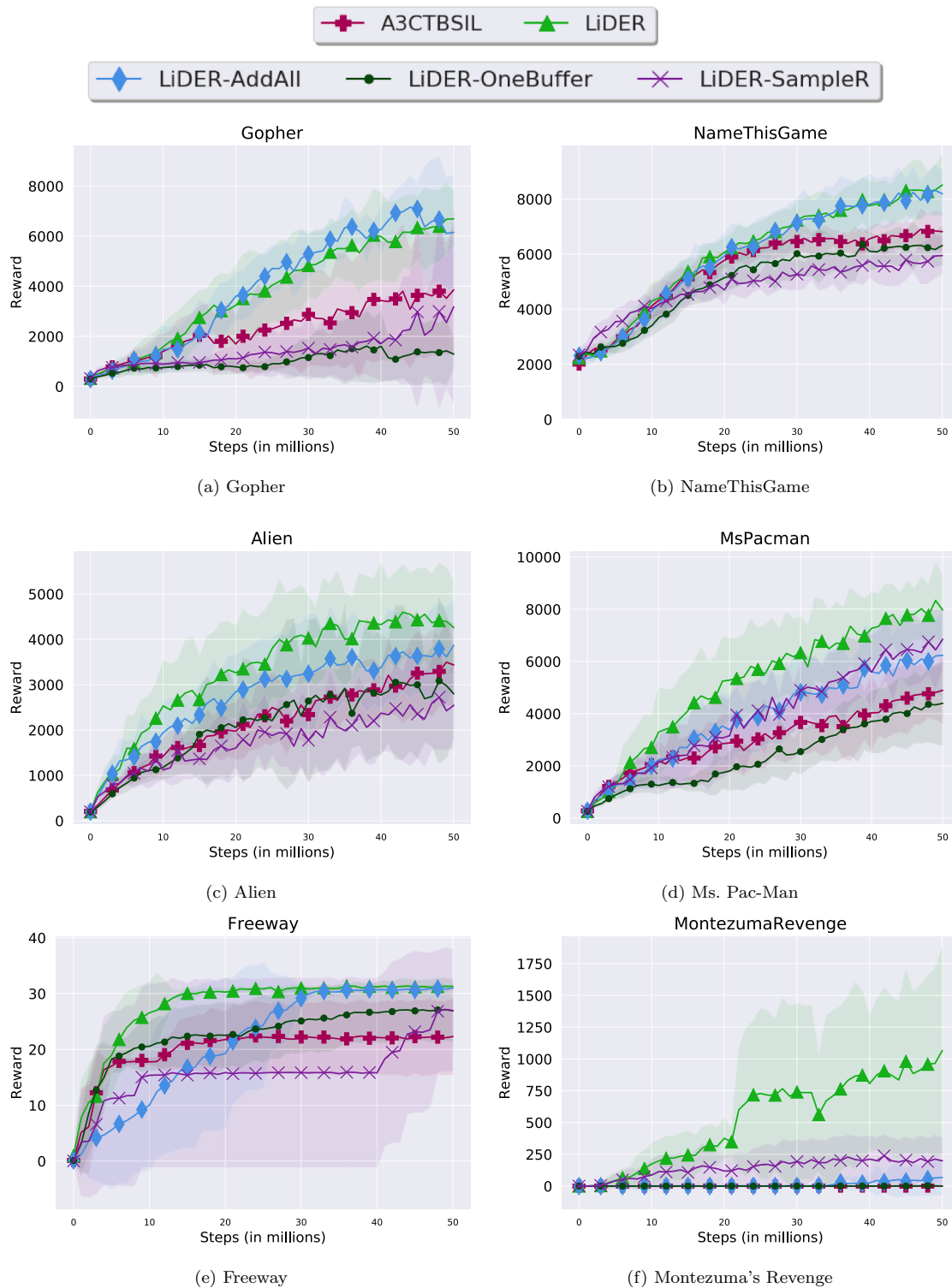
In particular, we consider leveraging a trained agent (TA) and a behavior cloning (BC) model trained from human demonstration data. LiDER-TA uses a trained agent (TA) as the refresher. While the TA could come from any source, we use the best checkpoint from a fully trained LiDER agent from experiments in Sect. 4.1 as the TA. This scenario tests whether LiDER can effectively leverage a high-quality policy.

LiDER-BC uses a behavior cloning (BC) model in the refresher. The BC policy is far from expert and we explore if LiDER can benefit from a sub-optimal policy. The BC model in LiDER-BC is pre-trained with non-expert demonstration data<sup>6</sup>, which was collected in our previous work [6]. Then, we follow the pre-training method introduced in our previous work [6] to jointly pre-train a model with supervised, value, and unsupervised autoencoder losses, which gives us a BC model trained from human demonstration data (see “Appendix 6” for pre-training details). The difference between the TA and BC model used in our setting is that the TA is an RL agent trained with LiDER while the BC model is trained only with human demonstrations.

Figure 7 shows the results of LiDER-TA and LiDER-BC compared with A3CTBSIL and LiDER (averaged over

<sup>5</sup> Note that the baseline A3CTBSIL represents the scenario of SampleD, i.e., always sample from buffer  $\mathcal{D}$ .

<sup>6</sup> The data is publicly available: [https://github.com/gabrieledcjr/atari\\_human\\_demo](https://github.com/gabrieledcjr/atari_human_demo)

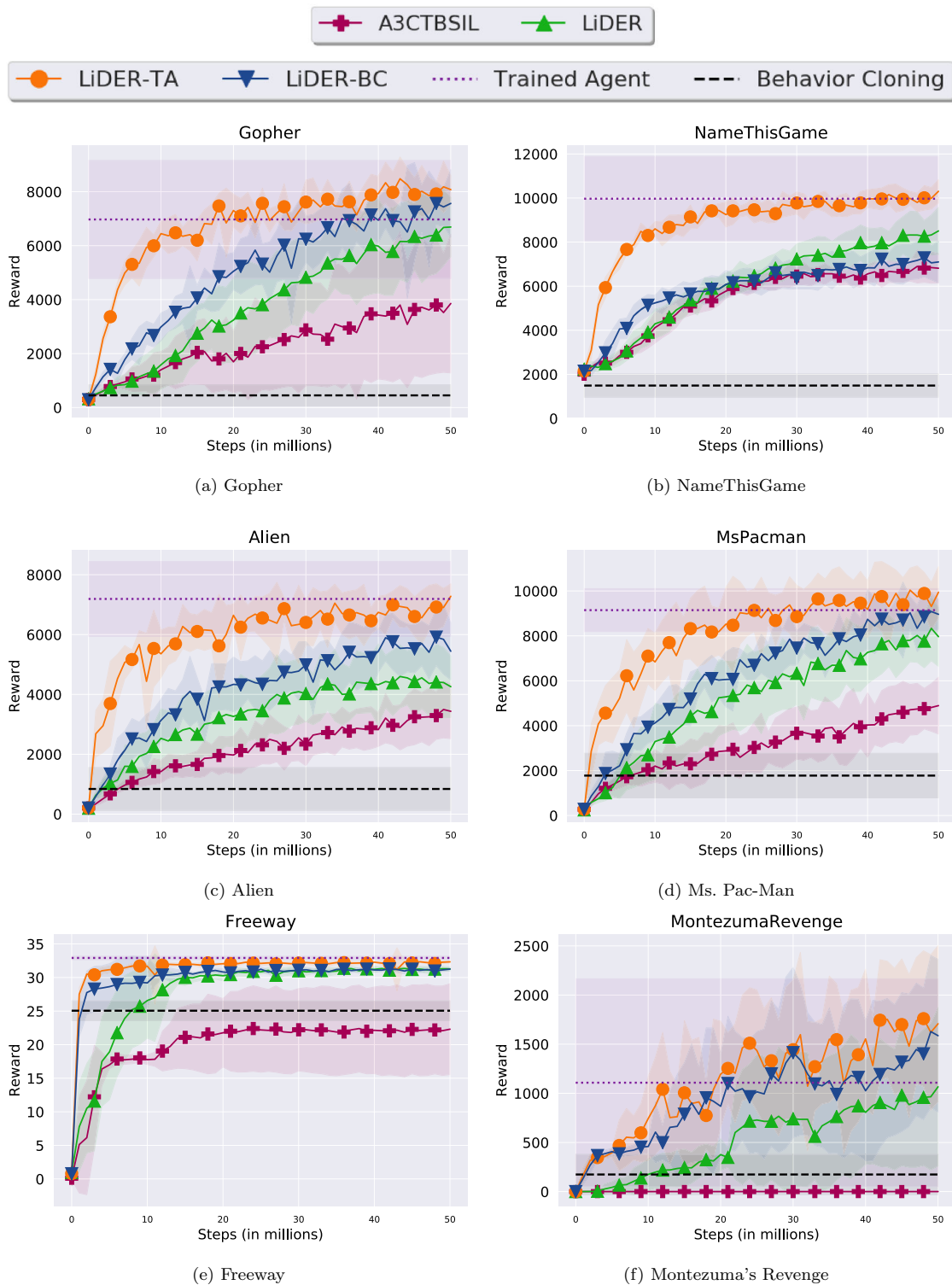


**Fig. 6** Ablation studies on LiDER in six Atari games. Results show that using only experiences where the return is improved, the two-buffer architecture, and the flexible sampling method does indeed improve performance. The x-axis is the total number of

environmental steps: A3CTBSIL counts steps from 16 A3C workers, while LiDER counts steps from 15 A3C workers plus one refresher worker. The y-axis is the average testing score over eight trials; shaded regions show the standard deviation

eight trials). As expected, LiDER-TA performs better than the other three methods, since it uses a trained agent as the

refresher—the learning agent can observe and learn from high-quality data generated by an expert. LiDER-TA even



**Fig. 7** LiDER-TA and LiDER-BC outperform A3CTBSIL and LiDER. The x-axis is the total number of environmental steps: A3CTBSIL counts steps from 16 A3C workers, while LiDER counts

exceeds the TA’s performance in Gopher and Montezuma’s Revenge. The TA’s performance is shown in the purple dotted line (shaded regions show the standard

steps from 15 A3C workers plus one refresher worker. The y-axis is the average testing score over eight trials; shaded regions show the standard deviation

deviation), estimated by executing the TA greedily in the game for 50 episodes. See “Appendix 5” for the score of each TA.

The more interesting result is the performance of LiDER-BC, which demonstrates that LiDER works well even when using a refresher that is far from expert. The black dashed line shows the average performance of the BC model (shaded regions show the standard deviation), estimated by executing the model greedily in the game for 50 episodes (see “Appendix 6” for the score of each BC). LiDER-BC can learn to quickly outperform BC and achieve better results than the baseline. LiDER-BC also slightly outperforms LiDER in five out of the six games, except for NameThisGame in which LiDER-BC outperforms LiDER initially, but later plateaued lower than LiDER. These results suggest that the sub-optimal BC model was able to provide better-than-random data during the early stages of training, which in turn helps the learning in the later stages. LiDER-BC could thus be one method of leveraging imperfect demonstrations to improve RL.

## 7 Related work

LiDER is related to several research directions in the RL literature; we briefly review four of them in this section.

### 7.1 Experience replay and extensions

ER was first introduced to improve the data efficiency of off-policy RL algorithms [24] and has since become an essential component for off-policy deep RL [27]. Many techniques have been proposed to enhance ER for better data efficiency and generally fall into two categories. One category focuses on biasing the sampling strategy such that important experiences are reused more frequently for policy updates [7, 31, 37, 39, 40, 47]. The other category focuses on tuning the replay buffer architecture, such as changing the buffer size [8, 25, 48], combining experiences from multiple workers to generate more data to replay [11, 19, 21], or augmenting the structure or content of replay experiences (e.g., generating additional “goal states” [1] or modifying experiences based on a teacher’s advice [4]).

LiDER does not fall into the first category but is complementary to existing sampling methods. We leverage prioritized experience replay [37] in our experiments: experiences are prioritized by advantages in buffer  $\mathcal{D}$  and buffer  $\mathcal{R}$ , the SIL worker samples from both buffers with priority (although the refresher worker samples randomly from buffer  $\mathcal{D}$ ). LiDER is related to the second category but differs in three ways. First, LiDER uses two replay buffers which double the buffer size, but we have shown that simply extending the size of a single buffer does not achieve the same performance as LiDER. Second, the refresher worker generates additional data, which is similar

to using multiple workers to generate more data, but we kept the total number of workers the same between LiDER and the baseline and accounted for all environmental steps. Third, the refresher-generated data is stored in a separate buffer only when it has a higher return than the old data, which can be viewed as augmenting the quality of the data, but we do not change the data structure when storing them.

Recently, Fedus et al. [12] revisited the fundamentals of ER to study how the architecture and the content of a replay buffer can affect learning. One of their key findings was that performance can be improved by increasing the replay buffer size and decreasing the age of the oldest data stored in the buffer. LiDER’s refreshing mechanism achieves exactly this purpose: the two-buffer architecture doubles the replay size and the refresher worker refreshes older experiences with a newer policy, which reduces the age of the overall policy. Therefore, LiDER can be viewed as a validation of the above finding by Fedus et al. [12].

### 7.2 Experience replay for actor–critic algorithms

The difficulty of combining ER into actor–critic algorithms is caused by the discrepancy between the current policy and the past policy that generated the experience. This problem is usually solved by leveraging various importance sampling techniques, such that the bias from past experiences can be corrected when used for updating the current policy [11, 15, 29, 45, 46]. In this work, we chose to use the SIL algorithm over the other actor–critic with ER algorithms because SIL provides a straightforward way of integrating ER into A3C without importance sampling [32].

As proven theoretically by Oh et al. [32], the SIL objective (Equation 2) updates the policy and the value function directly toward optimal by leveraging the Monte-Carlo return  $G$ , which can be viewed as a form of lower-bound-soft-Q learning. Thus, off-policy correction techniques like importance sampling are not needed even though the SIL worker learns from off-policy data (i.e., from a replay buffer), while the A3C worker learns on-policy. In addition, Oh et al. [32] have shown that the SIL objective is compatible, not conflicting, with off-policy correction algorithms like ACER [45]. LiDER builds upon the SIL objective and thus shares similar properties. Incorporating LiDER into other off-policy RL algorithms is important for future work (as described in Sect. 8).

### 7.3 Learning from past good experiences of oneself

The main idea of LiDER is to allow the agent to learn from past states that have been improved by its current policy. Several existing methods have shown that it is beneficial

for the agent to learn from its past good experiences. For example, the optimality tightening proposed by He et al. [17] constrains the Q function with lower and upper bounds, with the intuition that the Q function should be updated using trajectories that perform better than the current policy. The self-imitation learning (SIL) algorithm was inspired by the lower-bound Q learning from optimality tightening that only trajectories with positive advantages should be used to update the policy [32].

Gangwani et al. [14] and Guo et al. [16] extended the SIL algorithm and found that the performance can be further improved if the past good experiences are also diverse—diversity helps drive exploration. While we did not design LiDER to explicitly leverage exploration techniques, LiDER revisits a past state, and then generates new trajectories using a different policy, which could potentially lead to unseen states and increase the data diversity. This implicit exploration could be one of the reasons that LiDER improves the performance of two hard exploration Atari games.

A generalized form of SIL algorithm was proposed recently by Tang [43]. This new algorithm leverages n-step lower bound Q-learning which improves the original SIL algorithm in two aspects: 1) the agent can now self-imitate partial trajectories while the original SIL algorithm requires learning from a full trajectory, and 2) bootstrapping from learned Q-functions is enabled while the original SIL algorithm does not bootstrap from learned Q-functions. The generalized SIL algorithm can be applied to both deterministic and stochastic RL algorithms and outperforms SIL in a wide range of continuous control tasks. Leveraging the generalized SIL algorithm could be an interesting future work to improve LiDER.

Interestingly, the idea of learning from refreshed past states was also used in the MuZero algorithm [38], a tree-based searching algorithm that combines a learned model. Specifically, MuZero introduced a second variant called MuZero Reanalyze, in which the agent revisits a past time step and performs Monte-Carlo tree search again using the current model parameters. According to Schrittwieser et al. [38], MuZero Reanalyze largely improves the performance of MuZero because the reanalyze process potentially results in better policy than the original search. LiDER's results align with the findings of MuZero Reanalyze: an agent's current policy can be used to generate better quality data from a past state; leveraging these data leads to improved performance.

## 7.4 Relocating the agent to a past state

LiDER assumes there is a simulator for the task where resetting to a previously seen state is possible. The idea of relocating the agent to past states has been explored in the

literature (e.g., Mihalkova and Mooney [26]). Particularly, in the research area of curriculum learning, it is common to assume a simulator is available and the agent can be reset to any arbitrary state at the beginning of training (e.g., Florensa et al. [13]). A similar line of work has found that resetting the agent to a past state, instead of the simulator's default initial state, can benefit the learning. Such a state can be drawn from different distributions over the replay buffer [44] or from human demonstrations [20, 30, 34, 36, 49]. Many simulators are already equipped with the ability to relocate the agent so that they can reset the agent to an initial state when an episode ends. LiDER makes full use of this common feature.

While we can exploit simulators' relocation features if one is available, there are also situations when such a feature does not exist. The recently developed policy-based Go-Explore algorithm learned a goal-conditioned policy to guide the agent to return to a past state, which enables relocating without using the simulator reset feature [10].<sup>7</sup> Concurrently with policy-based Go-Explore, Guo et al. [16] proposed the diverse trajectory-conditioned self-imitation learning (DTSIL) algorithm. It uses similar mechanisms as Ecoffet et al. [10] to train a goal-condition based policy for the relocating process and no simulator reset is needed for DTSIL.<sup>8</sup>

Both the policy-based Go-Explore and DTSIL algorithms share similarities with LiDER in that they first teleport the agent to a past state then explore from there. However, LiDER is distinct from these two algorithms in three perspectives. First, the functionality of the learned policy is different. LiDER learns an actor-critic policy that maximizes the cumulative return; its policy takes a state as input and produces actions that can achieve maximum return. While the policy-based Go-Explore and DTSIL's policies are goal-conditioned and only learn how to return the agent to a past state, not how to maximize return. Their policies take both the agent's current state and the selected state (called a goal state) as input and produce actions that will lead the agent back to the goal state. Second, the state selection strategy is different. LiDER teleports the agent back to a randomly selected state, while the policy-based Go-Explore and DTSIL algorithms select "novel" states to return to (i.e., states that are rarely visited). Lastly, from the relocated state, LiDER then performs a refresh with its current policy. While the policy-based Go-Explore explores with either random actions or actions sampled

<sup>7</sup> The policy-based Go-Explore algorithm is an extension of the Go-Explore without a policy framework, which was presented in an earlier pre-print [9]. Go-Explore without a policy framework also leverages the simulator reset feature.

<sup>8</sup> Ecoffet et al. [10] made a detailed comparison between the policy-based Go-Explore and DTSIL. We refer the interested readers to Ecoffet et al. [10] for further reading.

from the goal-conditioned policy (with equal probability), the DTSIL algorithm only explores randomly from that state.

Besides the three key differences, there are many minor distinctions between LiDER and these two algorithms, such as the main goal and structure of the algorithm, the replay buffer architecture, the state representations, and the hyperparameters. Because of these differences, policy-based Go-Explore and DTSIL are not directly comparable to LiDER. On the other hand, LiDER can be considered as more evidence that supports the core benefits of “agent relocation,” rather than a competing method. Nevertheless, leveraging the relocation mechanism of these two algorithms in LiDER can be an important step toward allowing LiDER to work outside of simulations, as mentioned in the future work discussion in Sect. 8.

## 8 Discussion and future work

In this paper, we proposed *Lucid Dreaming for Experience Replay* (LiDER), a conceptually new framework that allows experiences in the replay buffer to be refreshed by leveraging the agent’s current policy, leading to improved performance compared to the baseline method without refreshing past experiences. We investigated the underlying behavior of the refresher to better understand why LiDER helps learning. We also conducted several ablation studies to validate our design choices of LiDER. Two extensions demonstrated that LiDER is also capable of leveraging knowledge from external policies, such as a trained agent and a behavior cloning model. One potential limitation of LiDER is that it must have access to a simulator that can return to previously visited states before resuming.

This paper opens up several new interesting directions for future work. First, based on the initial positive results reported in this paper, additional computational resources ought to be devoted to evaluating LiDER in a broad variety of domains.

Second, while we have presented in this paper a case study of applying LiDER to a multi-worker, actor–critic-based algorithm, future work could investigate extending LiDER to other types of off-policy RL algorithms that leverage ER. We expect LiDER to be most applicable to the PPO+SIL algorithm Oh et al. [32]. PPO+SIL’s multi-worker, actor–critic architecture allows the refresher worker to be easily added as was done in A3CTBSIL. Similarly, the SIL objective (Equation (2)) of PPO+SIL enables integrating ER, and the agent can learn from both PPO-generated and refresher-generated experiences.

On the other hand, applying LiDER to single-worker, value-based algorithms, such as the deep Q-network

(DQN) algorithm [27], presents more challenges. The first is that it is non-trivial to decide how often one should “pause” the training and perform a refresh. In multi-worker architectures, we do not need to explicitly control this frequency as all workers are running in parallel. With a single worker, the training and the refreshing cannot happen simultaneously. Another challenge is that, in value-based algorithms like DQN, a value function is learned instead of directly learning a policy. The Q value is updated using a one-step TD error instead of the Monte-Carlo return. Since the key concept of LiDER is to leverage higher returns for policy updates, how to integrate returns in value updates should be considered carefully before applying LiDER to DQN. Nevertheless, LiDER has the potential to benefit other off-policy algorithms that use ER, which is a good direction to explore in future work.

Third, the refresher in LiDER-BC uses a fixed policy from behavior cloning. Future work could investigate whether it helps to use different policies during training. For example, one could use the BC policy during the early stages of training, and then once A3C’s current policy outperforms BC, replace it with the A3C policy. Additionally, it is thus natural to consider adding multiple SIL and/or refresher works to enable leveraging multiple policies. Investigating how the proportion among the number of A3C, SIL, and refresher workers affects performance would make for an interesting future study.

Fourth, it would be interesting to allow LiDER to work outside of simulations by returning to a similar, but not identical state, and from there generate new trajectories. For example, in robotics, a robot may be able to return to a position that is close to, but not identical to, a previously experienced state.

## Appendices for “lucid dreaming for experience replay: refreshing past states with the current policy”

We provide further details of our work in the following six appendices:

- “Appendix 1” contains the implementation details of LiDER, including neural network architecture, hyperparameters, and computation resources used for all experiments.
- “Appendix 2” presents the pseudocode for the A3C and SIL workers. Both follow the original work of Mnih et al. [28] and Oh et al. [32], respectively, we add them here for completeness.
- “Appendix 3” provides detailed statistics of the one-tailed independent-samples t tests: 1) A3CTBSIL compared to LiDER, 2) A3CTBSIL compared to the



**Table 2** Hyperparameters for all experiments. We train each game for 50 million steps with a frame skip of 4, i.e., 200 million game frames were consumed for training

Network architecture	Value
Input size	$88 \times 88 \times 4$
Tensorflow Padding method	SAME
Convolutional layer 1	32 filters of size $8 \times 8$ with stride 4
Convolutional layer 2	64 filters of size $4 \times 4$ with stride 2
Convolutional layer 3	64 filters of size $3 \times 3$ with stride 1
Fully connected layer	512
Policy output layer	number of actions
Value output layer	1
<i>Common parameters</i>	
RMSProp initial learning rate	$7 \times 10^{-4}$
RMSProp epsilon	$1 \times 10^{-5}$
RMSProp decay	0.99
RMSProp momentum	0
Maximum gradient norm	0.5
Discount factor $\gamma$	0.99
<i>Parameters for A3CTB</i>	
A3C entropy regularizer weight $\beta^{a3c}$	0.01
A3C maximum bootstrap step $t_{\max}$	20
A3C value loss weight $\alpha$	0.5
$k$ parallel actors	16
Transformed Bellman operator $\varepsilon$	$10^{-2}$
<i>Parameters for SIL</i>	
SIL value loss weight $\beta^{\text{sil}}$	0.5
SIL update per step $M$	4
Replay buffer $\mathcal{D}$ size	$10^5$
Replay buffer priority $\alpha$	0.6
Minibatch size	32
<i>Parameters for LiDER (refresher worker)</i>	
Replay buffer $\mathcal{R}$ size	$10^5$
Minibatch size	32

three ablation studies of LiDER, 3) A3CTBSIL compared to the two extensions of LiDER, 4) LiDER compared to the three ablation studies of LiDER, and 5) LiDER compared to the two extensions of LiDER.

- “Appendix 4” discusses the differences between the A3CTBSIL algorithm in Jr et al. [6] and the original SIL algorithm in Oh et al. [32] (as mentioned in Sect. 4.1).
- “Appendix 5” presents the performance of the trained agents (TA) used in LiDER-TA.
- “Appendix 6” details the pre-training process for obtaining the BC models used in LiDER-BC, including the statistics of the demonstration collected by de la Cruz et al. [5], the network architecture, the hyperparameters used for pre-training, and the performance of the trained BC models.

## Appendix 1: Implementation details

We use the same neural network architecture as in the original A3C algorithm [28] for all A3C, SIL, and refresher workers (the blue, orange, and green components in Fig. 1, respectively). The network consists of three convolutional layers, one fully connected layer, followed by two branches of a fully connected layer: a policy function output layer and a value function output layer. Atari images are converted to grayscale and resized to  $88 \times 88$  with 4 images stacked as the input.

We run each experiment for eight trials due to computation limitations. Each experiment uses one GPU (Tesla K80 or TITAN V), five CPU cores, and 40 GB of memory (each LiDER-OneBuffer experiment uses 64 GB of memory since the buffer size was doubled). The refresher worker runs on GPU to generate data as quickly as possible; the A3C and SIL workers run distributively on CPU

cores. In all games, the wall-clock time is roughly 0.8 to 1 million steps per hour and around 50 to 60 hours to complete one trial of 50 million steps.

The baseline A3CTBSIL is trained with 17 parallel workers; 16 A3C workers and 1 SIL worker. The RMSProp optimizer is used with a learning rate = 0.0007. We use  $t_{\max} = 20$  for  $n$ -step bootstrap  $Q^{(n)}$  ( $n \leq t_{\max}$ ). The SIL worker performs  $M = 4$  SIL policy updates (Equation (2)) per step  $t$  with minibatch size 32 (i.e.,  $32 \times 4 = 128$  total samples per step). Buffer  $\mathcal{D}$  is of size  $10^5$ . The SIL loss weight  $\beta^{\text{sil}} = 0.5$ .

LiDER is also trained with 17 parallel workers: 15 A3C workers, 1 SIL worker, and 1 refresher worker—we keep the total number of workers in A3CTBSIL and LiDER the same to ensure a fair performance comparison. The SIL worker in LiDER also uses a minibatch size of 32; samples are taken from buffer  $\mathcal{D}$  and  $\mathcal{R}$  as described in Sect. 3. All other parameters are identical to that of A3CTBSIL. We summarize the details of the network architecture and experiment parameters in Table 2.

## Appendix 2: Pseudocode for the A3C and SIL workers

---

### Algorithm 2 LiDER: A3C Worker (as in Mnih et al. [28])

---

```

1: // Assume global network parameters  $\theta$  and  $\theta_v$  and global step  $T = 0$ 
2: // Assume replay buffer  $\mathcal{D} \leftarrow \emptyset$ , replay buffer  $\mathcal{R} \leftarrow \emptyset$ 
3: Initialize worker-specified local network parameters,  $\theta', \theta'_v$ 
4: Initialize worker-specified local time step  $t = 0$  and local episode buffer  $\mathcal{E} \leftarrow \emptyset$ 
5: while  $T < T_{\max}$  do ▷  $T_{\max} = 50$  million
6:   Reset gradients:  $d\theta \leftarrow 0, d\theta_v \leftarrow 0$ 
7:   Synchronize local parameters with global parameters  $\theta' \leftarrow \theta$  and  $\theta'_v \leftarrow \theta_v$ 
8:    $t_{\text{start}} \leftarrow t$ 
9:   while  $s_{t+1}$  is not terminal or  $t < t_{\max}$  do ▷  $t_{\max} = 20$ 
10:     Execute an action  $s_t, a_t, r_t, s_{t+1} \sim \pi(a_t | s_t, \theta')$ 
11:     Store transition to local buffer:  $\mathcal{E} \leftarrow \mathcal{E} \cup \{s_t, a_t, r_t, \}$ 
12:      $T \leftarrow T + 1, t \leftarrow t + 1$ 
13:   end while
14:    $G \leftarrow \begin{cases} 0 & \text{if } s_{t+1} \text{ is terminal} \\ V(S_{t+1}; \theta'_v) & \text{otherwise} \end{cases}$  ▷ Perform A3C update [28]
15:   for  $i \in \{t, \dots, t_{\text{start}}\}$  do
16:      $G \leftarrow r_i + \gamma G$ 
17:     Accumulate gradients w.r.t.  $\theta'$ :  $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i | s_i, \theta') (G - V(s_i; \theta'_v))$ 
18:     Accumulate gradients w.r.t.  $\theta'_v$ :  $d\theta_v \leftarrow d\theta_v + \partial (G - V(s_i; \theta'_v))^2 / \partial \theta'_v$ 
19:   end for
20:   if  $s_{t+1}$  is terminal then: ▷ Prepare for SIL worker [32]
21:     compute  $G_t = \sum_k \gamma^{k-t} r_k$  for all  $t$  in  $\mathcal{E}$ 
22:     Store transition to global replay buffer  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, G_t\}$  for all  $t$  in  $\mathcal{E}$ 
23:     Reset local buffer  $\mathcal{E} \leftarrow \emptyset$ 
24:   end if
25:   Asynchronously update global parameters using local parameters
26: end while

```

---

**Algorithm 3** LiDER: SIL Worker (as in Oh et al. [32])

```

1: // Assume global network parameters,  $\theta, \theta_v$ 
2: // Assume (Non-empty) replay buffer  $\mathcal{D}$ , replay buffer  $\mathcal{R}$ 
3: Initialize worker-specific local network parameters,  $\theta', \theta'_v$ 
4: Initialize local buffer  $\mathcal{B} \leftarrow \emptyset$ 
5: while  $T < T_{max}$  do ▷  $T_{max} = 50$  million
6:   Synchronize global step  $T$  from the most recent A3C worker
7:   Synchronize parameters  $\theta' \leftarrow \theta$  and  $\theta'_v \leftarrow \theta_v$ 
8:   for  $m = 1$  to  $M$  do ▷  $M = 4$ 
9:     Sample a minibatch of size 32  $\{s_D, a_D, G_D\}$  from  $\mathcal{D}$ 
10:    Sample a minibatch of size 32  $\{s_R, a_R, G_R\}$  from  $\mathcal{R}$ 
11:    Store both batches into  $\mathcal{B}$ :  $\mathcal{B} \leftarrow \{s_D, a_D, r_D\} \cup \{s_R, a_R, r_R\}$  ▷ Length of  $\mathcal{B}$ =64
12:    Sample a minibatch of 32  $\{s_B, a_B, G_B\}$  from  $\mathcal{B}$  ▷ Perform SIL update [32]
13:    Compute gradients w.r.t.  $\theta'$ :  $d\theta \leftarrow \nabla_{\theta'} \log \pi(a_B | s_B; \theta') (G_B - V(s_B; \theta'_v))_+$ 
14:    Compute gradients w.r.t.  $\theta'_v$ :  $d\theta_v \leftarrow \partial((G_B - V(s_B; \theta'_v))_+)^2 / \partial \theta'_v$ 
15:    Perform asynchronous update of  $\theta$  using  $d\theta$  and  $\theta_v$  using  $d\theta_v$ 
16:    Reset local buffer  $\mathcal{B} \leftarrow \emptyset$ 
17:   end for
18: end while

```

**Appendix 3: One-tailed independent-samples t tests**

We conducted one-tailed independent-samples t tests (equal variances not assumed) in all games to compare the differences in the mean episodic reward among all methods in this paper. For each game, we restored the best model checkpoint from each trial (eight trials per method) and executed the model in the game following a deterministic policy for 100 episodes (an episode ends when the agent loses all its lives) and recorded the reward per episode. This gives us 800 data points for each method in each game. We use a significance level  $\alpha = 0.001$  for all tests.

First, we check the statistical significance of the baseline A3CTBSIL compared to LiDER (Sect. 4.1), the main framework proposed in this paper. We report the detailed statistics in Table 3. Results show that the mean episodic reward of LiDER is significantly higher than A3CTBSIL ( $p \ll 0.001$ ) in all games.

Second, we compare A3CTBSIL to the three ablation studies, LiDER-AddAll, LiDER-OneBuffer, and LiDER-SampleR (Sect. 5). Table 4 shows that all ablations were helpful in Freeway and Montezuma’s Revenge, in which the mean episodic rewards of the ablations are significantly higher than the baseline ( $p \ll 0.001$ ). LiDER-AddAll also performed significantly better than A3CTBSIL in all games ( $p \ll 0.001$ ). LiDER-OneBuffer outperformed A3CTBSIL in Freeway and Montezuma’s Revenge ( $p \ll 0.001$ ), but it performed worse than the other four games ( $p \ll 0.001$ ). LiDER-SampleR outperformed A3CTBSIL in Ms. Pac-Man, Freeway, and Montezuma’s Revenge ( $p \ll 0.001$ ), but under-performed A3CTBSIL in Gopher, Name-ThisGame, and Alien ( $p \ll 0.001$ ).

Third, we compare A3CTBSIL to the two extensions, LiDER-BC and LiDER-TA (Sect. 6). Table 5 shows that

**Table 3** One-tailed independent-samples t test for the differences of the mean episodic reward between A3CTBSIL and LiDER. Equal variances are not assumed

Methods	Mean episodic reward (800 episodes)	SD	One-tailed p-value
<i>Gopher</i>			
A3CTBSIL	4291.20	2913.52	–
LiDER	6618.88	3300.10	$1.24 \times 10^{-47}$
<i>NameThisGame</i>			
A3CTBSIL	6786.75	1275.87	–
LiDER	8332.50	1754.30	$4.09 \times 10^{-80}$
<i>Alien</i>			
A3CTBSIL	3558.58	1596.18	–
LiDER	5065.04	2012.93	$3.77 \times 10^{-57}$
<i>Ms. Pac-Man</i>			
A3CTBSIL	4975.03	1527.05	–
LiDER	8532.34	2477.02	$1.49 \times 10^{-187}$
<i>Freeway</i>			
A3CTBSIL	23.10	5.84	–
LiDER	31.62	0.98	$1.19 \times 10^{-201}$
<i>Montezuma’s Revenge</i>			
A3CTBSIL	0.25	4.99	–
LiDER	987.63	951.69	$3.36 \times 10^{-129}$

the two extensions outperformed the baseline significantly in all games ( $p \ll 0.001$ ).

Fourth, we check the statistical significance of LiDER compared to the three ablation studies, LiDER-AddAll, LiDER-OneBuffer, and LiDER-SampleR (Sect. 5). Results in Table 6 show that most of the ablations significantly under-performed LiDER ( $p \ll 0.001$ ) in terms of the mean

**Table 4** One-tailed independent-samples t test for the differences of the mean episodic reward between A3CTBSIL and LiDER-AddALL, between A3CTBSIL and LiDER-OneBuffer, and between A3CTBSIL and LiDER-SampleR. Equal variances are not assumed

Methods	Mean episodic reward (800 episodes)	SD	One-tailed p-value
<i>Gopher</i>			
A3CTBSIL	4291.20	2913.52	–
(Ablation) LiDER-AddAll	7086.53	3188.04	$2.72 \times 10^{-68}$
(Ablation) LiDER-OneBuffer	1962.05	1872.94	$5.97 \times 10^{-72}$
(Ablation) LiDER-SampleR	3072.40	5146.23	$3.61 \times 10^{-9}$
<i>NameThisGame</i>			
A3CTBSIL	6786.75	1275.87	–
(Ablation) LiDER-AddAll	8200.04	1580.23	$2.86 \times 10^{-77}$
(Ablation) LiDER-OneBuffer	6422.48	1374.87	$2.34 \times 10^{-8}$
(Ablation) LiDER-SampleR	5819.81	1743.05	$3.35 \times 10^{-35}$
<i>Alien</i>			
A3CTBSIL	3558.58	1596.18	–
(Ablation) LiDER-AddAll	4054.28	1837.20	$5.13 \times 10^{-9}$
(Ablation) LiDER-OneBuffer	3204.41	1998.95	$4.74 \times 10^{-5}$
(Ablation) LiDER-SampleR	3104.99	1548.04	$4.86 \times 10^{-9}$
<i>Ms. Pac-Man</i>			
A3CTBSIL	4975.03	1527.05	–
(Ablation) LiDER-AddAll	6828.82	2562.33	$1.81 \times 10^{-62}$
(Ablation) LiDER-OneBuffer	4625.37	1920.67	$2.95 \times 10^{-5}$
(Ablation) LiDER-SampleR	7303.22	1869.98	$4.32 \times 10^{-134}$
<i>Freeway</i>			
A3CTBSIL	23.10	5.84	–
(Ablation) LiDER-AddAll	31.20	0.99	$1.35 \times 10^{-189}$
(Ablation) LiDER-OneBuffer	27.55	5.15	$6.34 \times 10^{-55}$
(Ablation) LiDER-SampleR	27.45	10.46	$4.17 \times 10^{-24}$
<i>Montezuma's Revenge</i>			
A3CTBSIL	0.25	4.99	–
(Ablation) LiDER-AddAll	77.63	144.18	$3.93 \times 10^{-46}$
(Ablation) LiDER-OneBuffer	3.00	24.31	$8.97 \times 10^{-4}$
(Ablation) LiDER-SampleR	265.86	178.74	$8.46 \times 10^{-205}$

episodic reward. Except for Gopher and NameThisGame, in which LiDER-AddAll performs at the same level as LiDER ( $p > 0.001$ ).

Lastly, we compare LiDER to the two extensions, LiDER-TA and LiDER-BC (Sect. 6). Results in Table 7 show that LiDER-TA always outperforms LiDER ( $p \ll 0.001$ ). LiDER-BC outperformed LiDER in Gopher, Alien, Ms. Pac-Man, and Montezuma's Revenge. In Freeway, LiDER-BC performs the same as LiDER ( $p > 0.001$ ), while in NameThisGame LiDER-BC performed worse than LiDER ( $p \ll 0.001$ ).

## Appendix 4: Differences between A3CTBSIL and SIL

There is a performance difference in Montezuma's Revenge between the A3CTBSIL algorithm (our previous work in de la Cruz Jr et al. [6], which is used as the baseline method in this article) and the original SIL algorithm (by Oh et al. [32]). The A3CTBSIL agent fails to achieve any reward while the SIL agent can achieve a score of 1100 (Table 5 in [32]).

We hypothesize that the difference is due to the different number of SIL updates (Equation (2)) that can be performed in A3CTBSIL and SIL; lower numbers of SIL updates would decrease the performance. In particular, Oh et al. [32] proposed to add the “Perform self-imitation

**Table 5** One-tailed independent-samples t test for the differences of the mean episodic reward between A3CTBSIL and LiDER-BC, and between A3CTBSIL and LiDER-TA. Equal variances are not assumed

Methods	Mean episodic reward (800 episodes)	SD	One-tailed p-value
<i>Gopher</i>			
A3CTBSIL	4291.20	2913.52	–
(Extension) LiDER-TA	8133.50	3800.38	$1.97 \times 10^{-98}$
(Extension) LiDER-BC	7775.75	3480.92	$1.11 \times 10^{-91}$
<i>NameThisGame</i>			
A3CTBSIL	6786.75	1275.87	–
(Extension) LiDER-TA	10227.69	2222.20	$7.63 \times 10^{-212}$
(Extension) LiDER-BC	7303.74	1649.01	$1.81 \times 10^{-12}$
<i>Alien</i>			
A3CTBSIL	3558.58	1596.18	–
(Extension) LiDER-TA	7753.54	1681.06	0.000
(Extension) LiDER-BC	6261.79	1865.67	$4.01 \times 10^{-166}$
<i>Ms. Pac-Man</i>			
A3CTBSIL	4975.03	1527.05	–
(Extension) LiDER-TA	10272.18	2035.98	0.000
(Extension) LiDER-BC	9613.89	2875.71	$2.40 \times 10^{-226}$
<i>Freeway</i>			
A3CTBSIL	23.10	5.84	–
(Extension) LiDER-TA	32.42	0.73	$2.81 \times 10^{-223}$
(Extension) LiDER-BC	31.68	0.85	$4.63 \times 10^{-203}$
<i>Montezuma's Revenge</i>			
A3CTBSIL	0.25	4.99	–
(Extension) LiDER-TA	1677.50	1050.33	$2.53 \times 10^{-222}$
(Extension) LiDER-BC	1811.86	994.38	$2.30 \times 10^{-256}$

learning” step in *each* A3C worker (Algorithm 1 of Oh et al. [32]). That is, when running with 16 A3C workers, the SIL agent is actually using 16 SIL workers to update the policy. However, A3CTBSIL only has one SIL worker, which means A3CTBSIL performs strictly fewer SIL updates compared to that of the original SIL algorithm, and thus resulting in lower performance.

We empirically validate the above hypothesis by conducting an experiment in the game of Ms. Pac-Man by modifying the A3CTBSIL algorithm from our previous work [6]. Instead of performing a SIL update whenever the SIL worker can, we force the SIL worker to only perform an update at even global steps; this setting reduces the total number of SIL updates by half. We denote this experiment as A3CTBSIL-ReduceSIL.

Figure 8 shows that A3CTBSIL-ReduceSIL under-performed A3CTBSIL, which provides preliminary evidence that the number of SIL updates is positively correlated to performance. More experiments will be performed in future work to further validate this correlation.

## Appendix 5: The performance of trained agents used in LiDER-TA

Section 6 shows that LiDER can leverage knowledge from a trained agent (TA). While the TA could come from any source, we use the best checkpoint of a fully trained LiDER agent. Table 8 shows the average performance of the TA used in each game. The score is estimated by executing the TA greedily in the game for 50 episodes. An episode ends when the agent loses all its lives.

## Appendix 6: Pre-training the behavior cloning model for LiDER-BC

In Sect. 6, we demonstrated that a BC model can be incorporated into LiDER to improve learning. The BC model is pre-trained using a publicly available human demonstration dataset. Dataset statistics are shown in Table 9.

The BC model uses the same network architecture as the A3C algorithm [28] and pre-training a BC model for A3C

**Table 6** One-tailed independent-samples t test for the differences of the mean episodic reward between LiDER and LiDER-AddAll, between LiDER and LiDER-OneBuffer, and between LiDER and LiDER-SampleR. Equal variances are not assumed. Methods in bold are **not** significant at level  $\alpha = 0.001$

Methods	Mean episodic reward (800 episodes)	SD	One-tailed p-value
<i>Gopher</i>			
LiDER	6618.88	3300.10	–
<b>(Ablation) LiDER-AddAll</b>	7086.53	3188.04	<b>0.002</b>
(Ablation) LiDER-OneBuffer	1962.05	1872.94	$3.65 \times 10^{-186}$
(Ablation) LiDER-SampleR	3072.40	5146.23	$1.38 \times 10^{-55}$
<i>NameThisGame</i>			
LiDER	8332.50	1754.30	–
<b>(Ablation) LiDER-AddAll</b>	8200.04	1580.23	<b>0.056</b>
(Ablation) LiDER-OneBuffer	6422.48	1374.87	$4.25 \times 10^{-110}$
(Ablation) LiDER-SampleR	5819.81	1743.05	$6.54 \times 10^{-147}$
<i>Alien</i>			
LiDER	5065.04	2012.93	–
(Ablation) LiDER-AddAll	4054.28	1837.20	$3.28 \times 10^{-25}$
(Ablation) LiDER-OneBuffer	3204.41	1998.95	$5.92 \times 10^{-70}$
(Ablation) LiDER-SampleR	3104.99	1548.04	$3.55 \times 10^{-92}$
<i>Ms. Pac-Man</i>			
LiDER	8532.34	2477.02	–
(Ablation) LiDER-AddAll	6828.82	2562.33	$9.06 \times 10^{-40}$
(Ablation) LiDER-OneBuffer	4625.37	1920.67	$4.18 \times 10^{-199}$
(Ablation) LiDER-sampleR	7303.22	1869.98	$2.76 \times 10^{-28}$
<i>Freeway</i>			
LiDER	31.62	0.98	–
(Ablation) LiDER-AddAll	31.20	0.99	$1.32 \times 10^{-17}$
(Ablation) LiDER-OneBuffer	27.55	5.15	$1.62 \times 10^{-85}$
(Ablation) LiDER-SampleR	27.45	10.46	$1.22 \times 10^{-27}$
<i>Montezuma's Revenge</i>			
LiDER	987.63	951.69	–
(Ablation) LiDER-AddAll	77.63	144.18	$1.68 \times 10^{-114}$
(Ablation) LiDER-OneBuffer	3.00	24.31	$1.09 \times 10^{-128}$
(Ablation) LiDER-SampleR	265.86	178.74	$5.31 \times 10^{-80}$

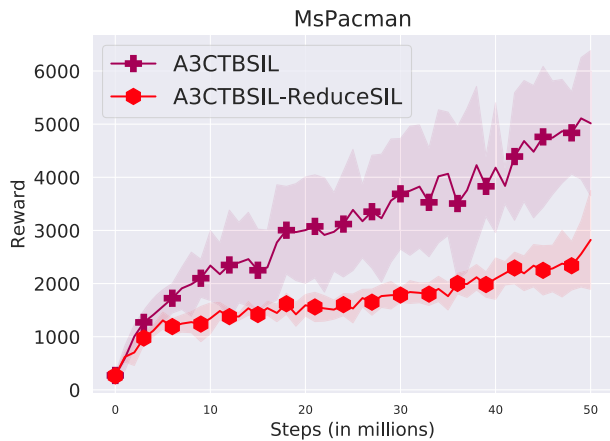
requires a few more steps than just using supervised learning as to how it is normally done in standard imitation learning (e.g., Ross and Bagnell [35]). A3C has two output layers: a policy output layer and a value output layer. The policy output is what we usually train a supervised classifier for. However, the value output layer is usually initialized randomly without being pre-trained. Our previous work [6] observed this inconsistency and leveraged demonstration data to also pre-train the value output layer. In particular, since the demonstration data contains the true return  $G$ , we can obtain a value loss that is almost identical to A3C's value loss  $L_{\text{value}}^{\text{A3C}}$ : instead of using the  $n$ -step bootstrap value  $Q^{(n)}$  to compute the advantage, the true return  $G$  is used.

Inspired by the supervised autoencoder (SAE) framework [22], our previous work [6] also blended in an unsupervised loss for pre-training. In SAE, an image reconstruction loss is incorporated with the supervised loss to help extract better feature representations and achieve better performance. A BC model pre-trained jointly with supervised, value, and unsupervised losses can lead to better performance after fine-tuning with RL, compared to pre-training with the supervised loss only.

We copy this approach by jointly pre-training the BC model for 50,000 steps with a minibatch of size 32. Adam optimizer is used with a learning rate = 0.0005. After training, we perform testing for 50 episodes by executing the model greedily in the game and record the average episodic reward (an episode ends when the agent loses all

**Table 7** One-tailed independent-samples t test for the differences of the mean episodic reward between LiDER and LiDER-TA, and between LiDER and LiDER-BC. Equal variances are not assumed. Methods in bold are not significant at level  $\alpha = 0.001$

Methods	Mean episodic reward (800 episodes)	Standard deviation	One-tailed p-value
<i>Gopher</i>			
LiDER	6618.86	3300.10	–
(Extension) LiDER-TA	8133.50	3800.38	$2.07 \times 10^{-17}$
(Extension) LiDER-BC	7775.75	3480.92	$6.55 \times 10^{-12}$
<i>NameThisGame</i>			
LiDER	8332.50	1754.30	–
(Extension) LiDER-TA	10227.69	2222.20	$3.75 \times 10^{-72}$
(Extension) LiDER-BC	7303.74	1649.01	$1.68 \times 10^{-32}$
<i>Alien</i>			
LiDER	5065.04	2012.93	–
(Extension) LiDER-TA	7753.54	1681.06	$3.53 \times 10^{-148}$
(Extension) LiDER-BC	6261.79	1865.67	$1.05 \times 10^{-33}$
<i>Ms. Pac-Man</i>			
LiDER	8532.34	2477.02	–
(Extension) LiDER-TA	10272.18	2035.98	$8.01 \times 10^{-50}$
(Extension) LiDER-BC	9613.89	2875.71	$7.81 \times 10^{-16}$
<i>Freeway</i>			
LiDER	31.62	0.98	–
(Extension) LiDER-TA	32.42	0.73	$8.54 \times 10^{-69}$
<b>(Extension) LiDER-BC</b>	31.68	0.85	<b>0.104</b>
<i>Montezuma's Revenge</i>			
LiDER	987.63	951.69	–
(Extension) LiDER-TA	1677.50	1050.33	$4.55 \times 10^{-41}$
(Extension) LiDER-BC	1811.88	994.38	$1.53 \times 10^{-59}$



**Fig. 8** A3CTBSIL-ReduceSIL compared to A3CTBSIL in the game of Ms. Pac-Man. The x-axis is the total number of environmental steps. The y-axis is the average testing score over five trials. We ran A3CTBSIL-ReduceSIL for five trials due to limited computing resources; we plot the first five trials out of eight for A3CTBSIL for a fair comparison to the number of trials in A3CTBSIL-ReduceSIL. Shaded regions show the standard deviation

**Table 8** The performance of trained agents used in LiDER-TA, shown as the purple dotted line in Fig. 7. The score is estimated by executing the TA greedily in the game for 50 episodes

Game	Trained TA score	SD
Gopher	6972.4	2190.26
NameThisGame	9969.0	1910.91
Alien	7190.4	1251.27
Ms. Pac-Man	9145.42	955.94
Freeway	32.92	0.27
Montezuma's Revenge	1108.0	1057.14

its lives). For each set of demonstration data, we train five models and use the one with the highest average episodic reward as the BC model in LiDER-BC. The performance of the trained BC models is present in Table 10. All

**Table 9** Demonstration size and quality, collected in de la Cruz et al. [5]. All games are limited to 20 min of demonstration time per episode

Game	Worst score	Best score	# of states	# of episodes
Gopher	1420	5800	16847	8
NameThisGame	2510	4840	17113	4
Alien	3000	8240	12885	5
Ms. Pac-Man	4020	18241	14504	8
Freeway	26	31	24396	12
Montezuma's Revenge	500	10100	18751	9

**Table 10** The performance of behavior cloning models used in LiDER-BC, shown as the black dashed line in Fig. 7. The score is estimated by executing the BC greedily in the game for 50 episodes

Game	Trained BC model score	SD
Gopher	450.8	393.57
NameThisGame	1491.2	530.55
Alien	839.2	718.72
Ms. Pac-Man	1776.6	993.94
Freeway	25.06	1.48
Montezuma's Revenge	174.0	205.72

**Table 11** Hyperparameters for pre-training the behavior cloning (BC) model used in LiDER-BC

Network architecture	Value
Input size	88×88×4
Tensorflow Padding method	SAME
Convolutional layer 1	32 filters of size 8×8 with stride 4
Convolutional layer 2	64 filters of size 4×4 with stride 2
Convolutional layer 3	64 filters of size 3×3 with stride 1
Fully connected layer	512
Classification output layer	number of actions
Value output layer	1
<i>Parameters for pre-training</i>	
Adam learning rate	$5 \times 10^{-4}$
Adam epsilon	$1 \times 10^{-5}$
Adam $\beta_1$	0.9
Adam $\beta_2$	0.999
L2 regularization weight	$1 \times 10^{-5}$
Number of minibatch updates	50,000
Batch size	32

parameters are based on those from our previous work [6] and we summarize them in Table 11.

**Acknowledgements** We thank Gabriel V. de la Cruz Jr. for helpful discussions; his open-source code at <https://github.com/gabrieledcjr/DeepRL> is used for training the behavior cloning models in this work. This research used resources of Kamiak, Washington State

University's high-performance computing cluster. Assefaw Gebremedhin is supported by the NSF award IIS-1553528. Part of this work has taken place in the Intelligent Robot Learning (IRL) Lab at the University of Alberta, which is supported in part by research grants from the Alberta Machine Intelligence Institute (Amii), CIFAR, and NSERC. Part of this work has taken place in the Learning Agents Research Group (LARG) at UT Austin. LARG research is supported in part by NSF (CPS-1739964, IIS-1724157, NRI-1925082), ONR (N00014-18-2243), FLI (RFP2-000), ARL, DARPA, Lockheed Martin, GM, and Bosch. Peter Stone serves as the Executive Director of Sony AI America and receives financial compensation for this work. The terms of this arrangement have been reviewed and approved by the University of Texas at Austin in accordance with its policy on objectivity in research.

## Declarations

**Conflict of interest** The authors declare that they have no conflict of interest.

## References

1. Andrychowicz M, Wolski F, Ray A, Schneider J, Fong R, Welinder P, McGrew B, Tobin J, Pieter Abbeel O, Zaremba W (2017) Hindsight experience replay. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R (eds) Advances in neural information processing systems. Curran Associates, Inc., 30:5048–5058. <https://proceedings.neurips.cc/paper/2017/file/453fadbd8a1a3af50a9df4df899537b5-Paper.pdf>
2. Bellemare M, Srinivasan S, Ostrovski G, Schaul T, Saxton D, Munos R (2016) Unifying count-based exploration and intrinsic motivation. In: Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R (eds) Advances in neural information processing systems. Curran Associates, Inc., vol 29, pp 1471–1479. <https://proceedings.neurips.cc/paper/2016/file/afda332245e2af431fb7b672a68b659d-Paper.pdf>
3. Bellemare MG, Naddaf Y, Veness J, Bowling M (2013) The arcade learning environment: an evaluation platform for general agents. *J Artif Intell Res* 47(1):253–279
4. Chan H, Wu Y, Kiros J, Fidler S, Ba J (2019) ACTRCE: augmenting experience via teacher's advice for multi-goal reinforcement learning. [arXiv:1902.04546](https://arxiv.org/abs/1902.04546)
5. de la Cruz GV, Du Y, Taylor ME (2019) Pre-training with non-expert human demonstration for deep reinforcement learning. *Knowl Eng Rev* 34:e10. <https://doi.org/10.1017/S0269888919000055>
6. de la Cruz Jr GV, Du Y, Taylor ME (2019) Jointly pre-training with supervised, autoencoder, and value losses for deep reinforcement learning. In: Adaptive and learning agents workshop, AAMAS



7. Dao G, Lee M (2019) Relevant experiences in replay buffer. In: 2019 IEEE symposium series on computational intelligence (SSCI), pp 94–101. <https://doi.org/10.1109/SSCI44817.2019.9002745>
8. De Bruin T, Kober J, Tuyls K, Babuška R (2015) The importance of experience replay database composition in deep reinforcement learning. In: Deep reinforcement learning workshop, NIPS
9. Ecoffet A, Huizinga J, Lehman J, Stanley KO, Clune J (2019) Go-explore: a new approach for hard-exploration problems. arXiv preprint [arXiv:1901.10995](https://arxiv.org/abs/1901.10995)
10. Ecoffet A, Huizinga J, Lehman J, Stanley KO, Clune J (2020) First return then explore. arXiv preprint [arXiv:2004.12919](https://arxiv.org/abs/2004.12919)
11. Espeholt L, Soyer H, Munos R, Simonyan K, Mnih V, Ward T, Doron Y, Firoiu V, Harley T, Dunning I, Legg S, Kavukcuoglu K (2018) IMPALA: scalable distributed deep-RL with importance weighted actor-learner architectures. In: Proceedings of Machine learning research 80:1407–1416. <http://proceedings.mlr.press/v80/espeholt18a.html>
12. Fedus W, Ramachandran P, Agarwal R, Bengio Y, Larochelle H, Rowland M, Dabney W (2020) Revisiting fundamentals of experience replay. In: Proceedings of the 37th international conference on machine learning, PMLR. <https://proceedings.icml.cc/paper/2020/hash/5460b9ea1986ec386cb64df22dff37be-Abstract.html>
13. Florensa C, Held D, Wulfmeier M, Zhang M, Abbeel P (2017) Reverse curriculum generation for reinforcement learning. In: Levine S, Vanhoucke V, Goldberg K (eds) Proceedings of machine learning research, PMLR, 78:482–495. <http://proceedings.mlr.press/v78/florensa17a.html>
14. Gangwani T, Liu Q, Peng J (2019) Learning self-imitating diverse policies. In: International conference on learning representations. <https://openreview.net/forum?id=HyxzRsR9Y7>
15. Gruslys A, Dabney W, Azar MG, Piot B, Bellemare M, Munos R (2018) The reactor: a fast and sample-efficient actor-critic agent for reinforcement learning. In: International conference on learning representations. <https://openreview.net/forum?id=rKHVZWZAZ>
16. Guo Y, Choi J, Moczulski M, Feng S, Bengio S, Norouzi M, Lee H (2020) Memory based trajectory-conditioned policies for learning from sparse rewards. In: Advances in neural information processing systems. <https://papers.nips.cc/paper/2020/hash/2df45244f09369e16ea3f9117ca45157-Abstract.html>
17. He FS, Liu Y, Schwing AG, Peng J (2017) Learning to play in a day: faster deep reinforcement learning by optimality tightening. In: International conference on learning representations. <https://openreview.net/forum?id=rJ8Je4clg>
18. Hester T, Vecerik M, Pietquin O, Lanctot M, Schaul T, Piot B, Horgan D, Quan J, Sendonaris A, Osband I, Dulac-Arnold G, Agapiou J, Leibo JZ, Gruslys A (2018) Deep Q-learning from demonstrations. In: Annual meeting of the association for the advancement of artificial intelligence (AAAI), New Orleans (USA)
19. Horgan D, Quan J, Budden D, Barth-Maroon G, Hessel M, van Hasselt H, Silver D (2018) Distributed prioritized experience replay. In: International conference on learning representations. <https://openreview.net/forum?id=H1Dy—0Z>
20. Hosu IA, Rebedea T (2016) Playing atari games with deep reinforcement learning and human checkpoint replay. arXiv preprint [arXiv:1607.05077](https://arxiv.org/abs/1607.05077)
21. Kapturovski S, Ostrovski G, Dabney W, Quan J, Munos R (2019) Recurrent experience replay in distributed reinforcement learning. In: International conference on learning representations. <https://openreview.net/forum?id=r1lyTjAqYX>
22. Le L, Patterson A, White M (2018) Supervised autoencoders: improving generalization performance with unsupervised regularizers. In: Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R (eds) Advances in neural information processing systems, Curran Associates, Inc., 31:107–117. <https://proceedings.neurips.cc/paper/2018/file/2a38a4a9316c49e5a833517c45d31070-Paper.pdf>
23. Lillicrap TP, Hunt JJ, Pritzel A, Heess N, Erez T, Tassa Y, Silver D, Wierstra D (2016) Continuous control with deep reinforcement learning. In: International conference on learning representations. [https://openreview.net/forum?id=tX\\_O8O-8ZI](https://openreview.net/forum?id=tX_O8O-8ZI)
24. Lin LJ (1992) Self-improving reactive agents based on reinforcement learning. Planning and teaching. *Mach Learn* 8(3–4):293–321
25. Liu R, Zou J (2018) The effects of memory replay in reinforcement learning. In: The 56th annual allerton conference on communication, control, and computing, pp 478–485
26. Mihalkova L, Mooney R (2006) Using active relocation to aid reinforcement learning. In: Proceedings of the 19th international FLAIRS conference (FLAIRS-2006), Melbourne Beach, FL, pp 580–585. <http://www.cs.utexas.edu/users/ai-lab/mihalkova/flairs06>
27. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G et al (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540):529
28. Mnih V, Badia AP, Mirza M, Graves A, Lillicrap T, Harley T, Silver D, Kavukcuoglu K (2016) Asynchronous methods for deep reinforcement learning. In: Balcan MF, Weinberger KQ (eds) Proceedings of machine learning research, PMLR, New York, New York, USA, vol 48, pp 1928–1937. <http://proceedings.mlr.press/v48/mniha16.html>
29. Munos R, Stepleton T, Harutyunyan A, Bellemare M (2016) Safe and efficient off-policy reinforcement learning. In: Lee D, Sugiyama M, Luxburg U, Guyon I, Garnett R (eds) Advances in neural information processing systems, Curran Associates, Inc., vol 29, pp 1054–1062. <https://proceedings.neurips.cc/paper/2016/file/c3992e9a68c5ae12bd18488bc579b30d-Paper.pdf>
30. Nair A, McGrew B, Andrychowicz M, Zaremba W, Abbeel P (2018) Overcoming exploration in reinforcement learning with demonstrations. In: 2018 IEEE international conference on robotics and automation (ICRA), pp 6292–6299. <https://doi.org/10.1109/ICRA.2018.8463162>
31. Novati G, Koumoutsakos P (2019) Remember and forget for experience replay. In: Chaudhuri K, Salakhutdinov R (eds) Proceedings of machine learning research, PMLR, Long Beach, California, USA, vol 97, pp 4851–4860. <http://proceedings.mlr.press/v97/novati19a.html>
32. Oh J, Guo Y, Singh S, Lee H (2018) Self-imitation learning. In: Dy J, Krause A (eds) Proceedings of machine learning research, PMLR, Stockholmsmässan, Stockholm Sweden, vol 80, pp 3878–3887. <http://proceedings.mlr.press/v80/oh18b.html>
33. Pohlen T, Piot B, Hester T, Azar MG, Horgan D, Budden D, Barth-Maroon G, van Hasselt H, Quan J, Vecerik M, et al. (2018) Observe and look further: achieving consistent performance on Atari. arXiv preprint [arXiv:1805.11593](https://arxiv.org/abs/1805.11593)
34. Resnick C, Raileanu R, Kapoor S, Peysakhovich A, Cho K, Bruna J (2018) Backplay:” Man muss immer umkehren”. In: Workshop on reinforcement learning in games, AAAI
35. Ross S, Bagnell D (2010) Efficient reductions for imitation learning. In: Teh YW, Titterton M (eds) Proceedings of machine learning research, JMLR workshop and conference proceedings, Chia Laguna Resort, Sardinia, Italy, 9:661–668. <http://proceedings.mlr.press/v9/ross10a.html>
36. Salimans T, Chen R (2018) Learning montezuma’s revenge from a single demonstration. arXiv preprint [arXiv:1812.03381](https://arxiv.org/abs/1812.03381)
37. Schaul T, Quan J, Antonoglou I, Silver D (2016) Prioritized experience replay. In: International conference on learning representations. [arXiv:1511.05952](https://arxiv.org/abs/1511.05952)

38. Schrittwieser J, Antonoglou I, Hubert T, Simonyan K, Sifre L, Schmitt S, Guez A, Lockhart E, Hassabis D, Graepel T, et al. (2019) Mastering Atari, go, chess and shogi by planning with a learned model. arXiv preprint [arXiv:191108265](https://arxiv.org/abs/1911.08265)
39. Sinha S, Song J, Garg A, Ermon S (2020) Experience replay with likelihood-free importance weights. arXiv preprint [arXiv:200613169](https://arxiv.org/abs/2006.13169)
40. Sovrano F (2019) Combining experience replay with exploration by random network distillation. In: 2019 IEEE conference on games (CoG), pp 1–8. <https://doi.org/10.1109/CIG.2019.8848046>
41. Stumbrys T, Erlacher D, Schredl M (2016) Effectiveness of motor practice in lucid dreams: a comparison with physical and mental practice. *J Sports Sci* 34:27–34
42. Sutton RS, Barto AG (2018) Reinforcement learning: an introduction. MIT Press, Cambridge
43. Tang Y (2020) Self-imitation learning via generalized lower bound Q-learning. In: Advances in neural information processing systems, vol 33. <https://papers.nips.cc/paper/2020/file/a0443c8c8c3372d662e9173c18faaa2c-Paper.pdf>
44. Tavakoli A, Levdi V, Islam R, Smith CM, Kormushev P (2018) Exploring restart distributions. [arXiv:181111298](https://arxiv.org/abs/1811.11298)
45. Wang Z, Bapst V, Heess NMO, Mnih V, Munos R, Kavukcuoglu K, de Freitas N (2017) Sample efficient actor-critic with experience replay. In: International conference on learning representations. <https://openreview.net/pdf?id=HyM25Mqel>
46. Wawrzyński P (2009) Real-time reinforcement learning by sequential actor-critics and experience replay. *Neural Netw* 22(10):1484–1497
47. Zha D, Lai KH, Zhou K, Hu X (2019) Experience replay optimization. In: Proceedings of the twenty-eighth international joint conference on artificial intelligence, IJCAI-19, international joint conferences on artificial intelligence organization, pp 4243–4249. <https://doi.org/10.24963/ijcai.2019/589>, <https://doi.org/10.24963/ijcai.2019/589>
48. Zhang S, Sutton RS (2017) A deeper look at experience replay. arXiv preprint [arXiv:171201275](https://arxiv.org/abs/1712.01275)
49. Zhang X, Bharti SK, Ma Y, Singla A, Zhu X (2020) The teaching Dimension of Q-learning. arXiv preprint [arXiv:200609324](https://arxiv.org/abs/2006.09324)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.