

---

# Coreset-Augmented Experience Replay

---

**Anonymous Author(s)**

Affiliation  
Address  
email

## Abstract

1      Experience replay (ER) is a widely used method in reinforcement learning to decor-  
2      relate sequential data and to improve sample efficiency by reusing past transitions  
3      during learning. The most widely used way to select transitions is to keep the  
4      agent’s most recent  $k$  experiences. This limits the agent’s ability to hold onto  
5      important transitions, and can have a large memory footprint. To address this  
6      limitation, we develop a coresnet-based replay strategy designed to keep a compact  
7      yet diverse set of experiences. By employing coresnet-augmented experience replay,  
8      we show that in some environments, it is possible to closely match the performance  
9      of a life-time sized recency buffer with 0.2% of the memory footprint, and storing  
10     transitions from arbitrary points of the agent’s lifetime.

11    **1 Introduction**

12    Experience Replay (ER) [10] has been widely adopted in Reinforcement Learning (RL) to stabilize  
13    and enhance training across various applications [12, 17]. Typically, ER is implemented using a  
14    first-in-first-out buffer that stores recent transitions. Mini-batches are then sampled independently  
15    and uniformly, giving each transition an equal chance of being selected. While this uniform strategy  
16    is simple and effective, several alternatives have been proposed to improve the effectiveness and  
17    efficiency of ER [15, 1]. Among them, Prioritized Experience Replay (PER) [13, 15] has received  
18    particular attention. PER assigns higher sampling probabilities to transitions with larger temporal  
19    difference (TD) errors [18], aiming to focus learning on samples for which the agent has higher  
20    TD-error. However, recent studies [13] suggest that, when combined with function approximation,  
21    PER may lead to overestimation, and poor generalization in early learning. This invites further  
22    investigation into methods that can effectively surpass the performance of recency based Replay  
23    buffers.

24    One approach to increase the sample efficiency of ER is to effectively reduce its size. The literature  
25    presents mixed findings regarding the optimal size of experience replay buffers [4] [23]. Ideally,  
26    we would prefer a solution that mitigates the need to carefully tune buffer size altogether. In this  
27    work, we propose a new perspective on experience replay. Rather than restricting the agent to  
28    learning only from its most recent experiences, we explore retaining experiences from arbitrary  
29    points throughout the agent’s lifetime. While it is theoretically possible to store all past experiences,  
30    doing so is impractical: it may require unbounded memory, which might lead to a degradation in  
31    learning performance [23], and this approach fails to account for the varying importance of different  
32    experiences in shaping the agent’s behavior [8]. What is needed is a method that chooses a subset of  
33    the agent’s experience to store and relearn from, to reduce the memory footprint of the agent.

34    Coreset construction methods are one approach to try to reduce the size of an experience replay buffer.  
35    They have been effectively used to reduce the data needs in machine learning [21, 11]. We introduce  
36    a novel experience replay mechanism called coresnet-augmented replay, which aims to maintain a  
37    limited set of carefully selected experiences from across the agent’s lifetime, according to a selection  
38    criterion. Additionally, to preserve the benefits of recency, we complement the coresnet buffer with a

39 similarly sized recency buffer that captures the most recent experiences, ensuring the agent remains  
 40 responsive to new and relevant information.  
 41 We investigate the effect of using this experience replay method in a variety of benchmark environments.  
 42 To begin with, we test the method in the environment of Open Room, on variations of our  
 43 method to test the usefulness of our method. While promising, our experiments showcase a few of  
 44 the challenges of adapting prototype selection methods to this setting: choosing the right criteria  
 45 such that stored transitions remain useful to the agent, and we showcase this by ejecting the oldest  
 46 samples in the buffer every  $k$  steps. We also show promising results of the algorithm working in well  
 47 known benchmark environments, Mountain Car [14] and Acrobot [19], while using a learned notion  
 48 of similarity/distance between samples. In general, our method showcases a promising direction in  
 49 using more effective sampling selection methods for experience replay.

## 50 2 Background, Problem Setting and Notation

51 We formalize our setting as a finite, discrete-time Markov Decision Process (MDP) [2], described  
 52 by the tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$ . In this framework, at each time step  $t$ , the agent first observes a state  
 53  $S_t \in \mathcal{S}$ . Based on this observation, it selects an action  $A_t \in \mathcal{A}$ , after which the environment  
 54 transitions to a next state  $S_{t+1}$  drawn according to the transition probability distribution  $P(\cdot | S_t, A_t)$ .  
 55 Concurrently, the agent receives a scalar reward signal  $R_{t+1} = R(S_t, A_t)$  as feedback. Ultimately,  
 56 the objective of the agent is to learn a policy  $\pi$  that maximizes the expected cumulative reward,  
 57 defined as  $\mathbb{E}[G_t | S_t = s] = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$ , for each initial state  $S_0$ ,  
 58 where  $\gamma$  is the discount factor  $\gamma \in [0, 1]$ .

59 Q-learning [20] is one way of achieving the above goal, which we focus on in this paper. Q-learning  
 60 estimates the action-value function for the optimal policy, via the off-policy temporal difference  
 61 learning update rule:

$$\hat{q}(S_t, A_t) \leftarrow \hat{q}(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_a \hat{q}(S_{t+1}, a) - \hat{q}(S_t, A_t)] \quad (1)$$

62 for each state-action pair  $(S_t, A_t)$ . The optimal policy  $\pi^*$  is derived from the action-value function  
 63  $\hat{q}(s, \cdot)$ , by choosing the action with the highest value, via  $\arg \max_{a \in \mathcal{A}} \hat{q}(s, a)$ . In environments  
 64 with large or infinite state spaces, explicitly storing action-values  $\hat{q}(s, a)$  for every state-action pair  
 65 becomes infeasible. Therefore, one could use function approximation with neural networks to  
 66 generalize action-values across states.

67 DQN [12] is a reinforcement learning algorithm that incorporates function approximation using  
 68 deep neural networks. To address issues of sample correlation and policy non-stationarity—caused  
 69 by global changes in the network parameters  $\theta$  during stochastic gradient descent—DQN relies on  
 70 experience replay (ER) to stabilize training.

71 Experience replay in its simplest form, operates by storing each new incoming transition in a queue  
 72 of the most  $k$  recent samples, discarding the oldest sample to keep a constant queue size  $k$ . In this  
 73 work, we refer to this type of experience replay buffer as a *recency* replay buffer or *recency* buffer.  
 74 For the purpose of reducing the memory footprint of the agent, we would like to keep a subset of the  
 75 agent’s interactions with the environment, and sample transitions uniformly from that subset. We call  
 76 this subset a *coreset* buffer.

77 In machine learning, a coreset is a small subset of a dataset that approximately preserves key properties  
 78 of the original data, enabling more efficient computation of optimization problems. To construct such  
 79 a subset, especially in settings like reinforcement learning where data arrives continuously and online,  
 80 we must first define a way to measure similarity/distance between samples. In this work, we use a  
 81 *distance representation*  $\phi : (\mathcal{S} \times \mathcal{A} \times \mathcal{R} \times \mathcal{S}) \rightarrow X$ , where  $X$  denotes the image of the representation,  
 82 mapping samples into a space where similarity and diversity can be meaningfully quantified. This  
 83 representation forms the basis for selecting diverse samples during coreset construction.

84 The choice of the distance representation  $\phi$  determines the measure with which we determine how  
 85 diverse or similar samples are to one another. We discuss how this representation can be chosen  
 86 based on domain knowledge or learned automatically in subsequent sections. After the choice of the  
 87 distance representation  $\phi$  is made, The distance between each pair of samples  $x_i, x_j$  is computed  
 88 using their Euclidean distance  $\|x_i - x_j\|_2$ .

89 Here, We specifically focus on online coresets construction [22], where the goal is to incrementally  
 90 select diverse transitions. We adapt the Online Greedy algorithm [16] to be the reinforcement learning  
 91 setting. Online Greedy uses the similarity between the points to construct a matrix of pair wise  
 92 similarities  $K_{ij} = \kappa(x_i, x_j)$ . The method uses the gaussian Radial Basis Function (RBF) Kernel and  
 93 by taking the Euclidean distance of each of the samples' distance representation  $\phi$ , calculates their  
 94 similarities:

$$\kappa(x_i, x_j) = \exp(-\eta \|\phi(x_i) - \phi(x_j)\|_2) \quad (2)$$

95 where  $\eta$  is the shape parameter. The matrix  $K$  is called the Gram matrix [7]. To maximize the diversity  
 96 of samples in the coresets, for each incoming sample, Online Greedy maximizes the logarithm of the  
 97 absolute value of the determinant of the Gram matrix  $K$  (See 3), This objective, captures the diversity  
 98 of the coresets:

$$D(C) = \log |\det(K)| \quad (3)$$

### 99 3 Coreset-Augmented Experience Replay

#### 100 3.1 Online Greedy

101 In order to investigate the effectiveness of diversity-based selection in experience replay, we adapt  
 102 Online Greedy to the reinforcement learning setting. The Online Greedy algorithm (See Algorithm 1)  
 103 maintains a fixed-size coreset by greedily inserting transitions that maximize the logarithm of the  
 104 absolute value of the Gram matrix  $K$ , constructed from the pairwise similarities from current samples  
 105 in the coreset buffer, to ensure that samples stored in the coreset buffer remain as diverse as possible.

---

**Algorithm 1** Online Greedy with coreset buffer

---

**Require:** Similarity function  $\kappa(\cdot, \cdot)$ , threshold for coreset size  $N$ , function  $\log |\det(\cdot)|$ , candidacy rate  $\eta$

1: Initialize coreset  $C \leftarrow \emptyset$   
 2: **while**  $|C| < N$  **do**  
 3:     Receive transition  $t$   
 4:      $C \leftarrow C \cup \{t\}$   
 5: **end while**  
 6: **for all** incoming transition  $t$  **do**  
 7:      $\lambda \sim U(0, 1)$   
 8:     **if**  $\lambda > \eta$  **then**  
 9:         **break**  
 10:     **end if**  
 11:      $G_0 \leftarrow \log |\det(K(C))|$   $\triangleright K(C)$ : Gram matrix of  $C$   
 12:     Initialize  $\text{best\_gain} \leftarrow 0$ ,  $\text{best\_swap\_index} \leftarrow \text{None}$   
 13:     **for**  $i = 1$  to  $N$  **do**  
 14:          $C' \leftarrow C$  with  $C[i]$  replaced by  $t$   
 15:          $G_i \leftarrow \log |\det(K(C'))|$   
 16:         **if**  $G_i - G_0 > \text{best\_gain}$  **then**  
 17:              $\text{best\_gain} \leftarrow G_i - G_0$   
 18:              $\text{best\_swap\_index} \leftarrow i$   
 19:         **end if**  
 20:     **end for**  
 21:     **if**  $\text{best\_gain} > 0$  **then**  
 22:         Replace  $C[\text{best\_swap\_index}]$  with  $t$   
 23:     **end if**  
 24: **end for**

---

106 We adapt the Online Greedy algorithm to our setting as follows. Initially, the buffer is pre-filled  
 107 with the most recent transitions. For each new incoming transition, we consider replacing it with  
 108 each existing transition in the buffer individually. For each potential replacement, we recompute

109 the diversity metric  $D$  over the resulting buffer. The transition is only retained if one of these  
 110 replacements yields an increase in the overall diversity; in that case, we apply the swap that results in  
 111 the greatest increase. If no such improvement is observed, the incoming transition is discarded.

112 In its current form, the algorithm will not be able to react to changes in the environment sufficiently  
 113 quickly. Samples selected for inclusion in the coresset buffer do not necessarily represent the agent's  
 114 most recent interactions with the environment, which may hinder the agent's ability to adapt online  
 115 and learn from its latest experiences. In order to mitigate this, we also include a similarly sized  
 116 recency buffer. We then sample from the coresset and recency buffers in proportion to their respective  
 117 sizes, combine the selected transitions into a mini-batch, and use this batch for training with learning  
 118 algorithms such as DQN.



Figure 1: The composite replay buffer. As the agent interacts with its environment at each time step, new samples arrive and are first put in the recency buffer. Once they exit the recency buffer, they are checked as candidates to be added to the coresset buffer.

### 119 3.1.1 Distance Diversity Criterion

120 It is difficult to scale the size of the coresset in the prior method, as the number of operations scales  
 121 with the order of  $O(n^4)$  where  $n$  is the size of the coresset. To mitigate this, we adapt a dictionary  
 122 sparsification procedure [6] to diversify our coresset buffer. This approach evaluates the similarity  
 123 between the incoming sample and each sample already in the replay buffer, and admits (rejects) the  
 124 new sample only if its similarity scores is smaller (exceeds) a predefined threshold, for every sample  
 125 in the coresset buffer.

$$\kappa(x_i, x_i) - \frac{\kappa(x_i, x_j)^2}{\kappa(x_j, x_j)} \geq \delta^2 \quad (4)$$

126 We simplify this criterion by setting  $\kappa(x_t, x_t) = 1$ , reflecting the assumption that a sample is  
 127 maximally similar to itself. When the buffer reaches capacity, we apply a sparsification procedure  
 128 (see Algorithm 3) to free up space for new samples.

129 The diversity criterion coresset then works by evaluating every incoming sample, and using equation 5  
 130 to decide whether to reject or accept incoming samples.

$$1 - \kappa(x_i, x_j)^2 < \delta^2 \quad (5)$$

131 if the condition holds for all samples  $x_i$  in the coresset, that means that the sample is sufficiently  
 132 dissimilar to all other samples and can be accepted to the coresset (See algorithm 2).

133

## 134 4 Experiments

135 How can these methods be adapted to experience replay to reduce its memory footprint and improve  
 136 sample efficiency? Can they be adapted to a fully specified agent without modification? This paper  
 137 investigates these questions by running experiments with different coresset strategies in standard  
 138 reinforcement learning tasks.

---

**Algorithm 2** Transition Acceptance with Periodic Coreset Sparsification

---

**Require:** Incoming transition  $t$ , similarity function  $\text{sim}(p, t)$ , threshold  $\delta$ , interval  $k$

```
1: Initialize coreset  $C \leftarrow \emptyset$ 
2: Initialize step counter  $i \leftarrow 0$ 
3: function PROCESSTRANSITION( $t$ )
4:    $i \leftarrow i + 1$ 
5:   if  $i \bmod k = 0$  then
6:     Remove the oldest item from  $C$ 
7:   end if
8:   if  $C = \emptyset$  then
9:     Accept  $t$  and add to  $C$ 
10:    return
11:   end if
12:   if coreset is full then
13:     Sparsify  $C$ 
14:   end if
15:   for all  $p \in C$  do
16:     if  $1 - \kappa(p, t)^2 < \delta^2$  then
17:       Reject  $t$ 
18:       return
19:     end if
20:   end for
21:   Accept  $t$  and add to  $C$ 
22: end function
```

---

---

**Algorithm 3** Coreset Sparsification

---

**Require:** Current coreset  $C$ , similarity function  $\text{sim}(\cdot, \cdot)$ , threshold  $\delta$

```
1: Initialize sparsed_coreset  $\leftarrow \emptyset$ 
2:  $\Delta \leftarrow \eta \cdot \delta$                                  $\triangleright$  Increase threshold by factor of  $\eta > 1$ 
3: for all  $p \in \text{reverse}(C)$  do                   $\triangleright$  From most recent to oldest
4:   for all  $t \in \text{sparsed\_coreset}$  do
5:     if  $1 - \kappa(p, t)^2 < \delta^2$  then
6:       continue to next  $p$ 
7:     end if
8:   end for
9:   Add  $p$  to sparsed_coreset
10: end for
11:  $C \leftarrow \text{sparsed\_coreset}$ 
12:  $\delta \leftarrow \Delta$                                       $\triangleright$  Update threshold
```

---

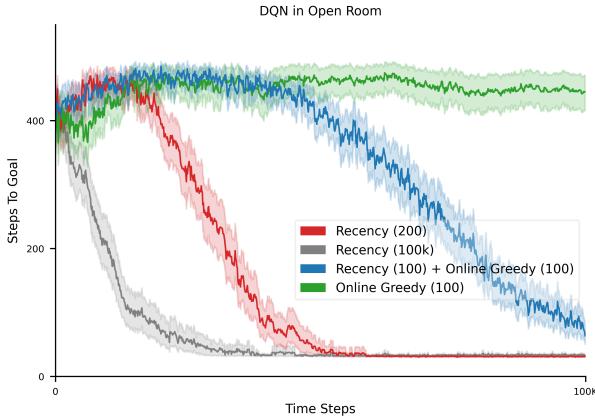
139 **4.1 The Open Room Environment**

- 140 We begin this section with experiments in a simple gridworld environment, called Open Room.
- 141 The environment consists of an empty room, and a square region in the bottom right corner as the goal state. The observation given to the agent is are the  $(x, y)$  coordinates in the room. The agent takes action  $a$  of **left**:  $[-1, 0]$ , **right**:  $[1, 0]$ , **up**:  $[0, 1]$ , and **down**:  $[0, -1]$  directions, and some uniform noise  $\epsilon \sim U(0, 1)$  is added to the action to make the environment stochastic. The episode terminates when the agent reaches the goal region, where the agent receives its only non-zero reward of +1.
- 146 For these experiments, we compare four variants of experience replay: (1) a recency buffer containing all experiences over the agent's lifetime, (2) a recency buffer with the 200 most recent experiences, (3) the Online Greedy coresset construction method with a fixed buffer size of 100, and (4) the composite buffer defined in the methodology section.
- 150 We use DQN as the policy optimization algorithm, with a neural network consisting of two hidden layers, each with 32 units. The target network is updated every 128 time steps. The distance representation  $\phi$  is chosen as the  $(x, y)$  coordinates of the agent in the environment. For the input of



Figure 2: The Open Room Enviornment. The goal of the agent is to reach the goal state.

153 the network, we represent each continuous  $(x, y)$  coordinate as a one-hot vector by discretizing each  
 154 axis into 20 uniform buckets. The resulting feature is a 40-dimensional binary vector with one active  
 155 entry per axis. The  $\eta$  in the RBF kernel is set to 5. To avoid over-computation, we only accept 10%  
 156 of the seen samples by the agent. The intuition is that transitions close to each other are also similar.  
 157 Optimization is performed using Adam [9] with  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . To select the learning rate  
 158 for each variant, we conduct a grid search over values  $[10^{-7}, 5 \times 10^{-7}, 10^{-6}, 5 \times 10^{-6}, 10^{-5}, 5 \times$   
 159  $10^{-5}, 10^{-4}, 5 \times 10^{-4}, 10^{-3}, 5 \times 10^{-3}, 10^{-2}]$  using 30 random seeds. We then evaluate the best  
 160 configuration for each variant over 100 independent seeds to generate the final performance plots.



Using ejection improves the performance the and the algorithm can re-cover the performance of the coresset buffer.

Figure 3: The use of a coresset replay buffer in the Open Room Environment, with DQN. The y-axis represents the number of steps it takes the agent to reach the goal. As one can see, A coresset buffer made using Online Greedy does not learn to reach the goal,

161 The performance of the algorithm is highly degraded in the absence of a recency buffer. This is  
 162 expected since the agent is unable to react to changes or keep track of new experiences, as samples  
 163 can be from arbitrary points of the agents experience.

#### 164 4.1.1 The Effect of Network Size

165 To better understand the impact of buffer composition and network capacity on learning performance,  
 166 we begin by examining a key question: *Why might the network perform poorly during training?*  
 167 One possible explanation is that inaccurate updates—arising from suboptimal samples—propagate

168 through the network and consume a significant portion of its capacity, hindering generalization and  
 169 learning.

170 To investigate this hypothesis, we analyze the effect of increasing the size of the neural network. In  
 171 Figure 4b, subfigure (a) shows performance when using a network with 128 units per layer, while  
 172 subfigure (b) depicts results with only 8 units per layer. In both configurations, agents trained using  
 173 large and small recency buffers (of size  $10^5$  and 200, respectively) are still able to solve the task  
 174 effectively, however, the coresnet buffer method suffers in the case with a smaller network, suggesting  
 175 that the network’s capacity plays a critical role in mitigating the impact of harmful updates.

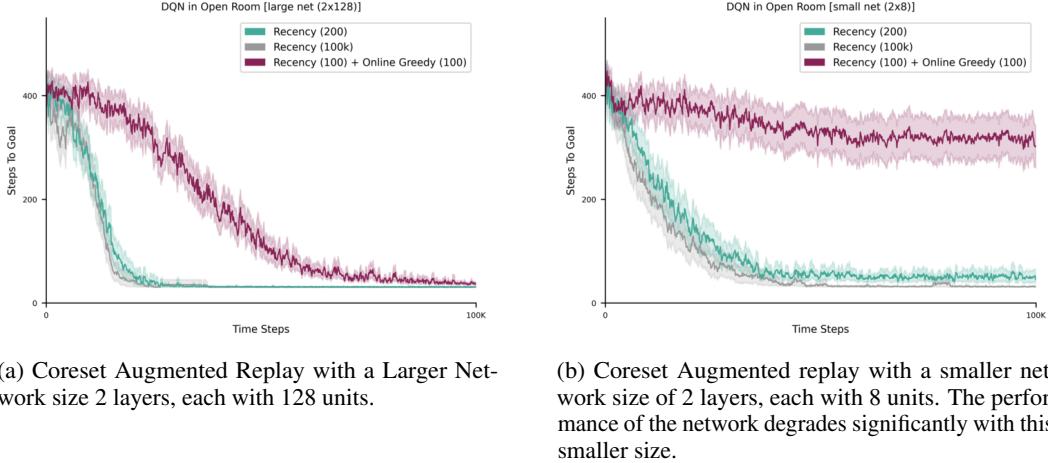


Figure 4: The effect of network size on the performance of coresnet Augmented replay. The performance degrades when using smaller network sizes.

176 Network size affects the performance of a Composite Buffer. Specifically, increasing the size of the  
 177 network (its capacity) seems to help the composite buffer learn better, but learning is still slower than  
 178 both a small and large recency buffers.

#### 179 **4.1.2 Composite coresnet with Ejection**

180 Having established that both small and large networks can learn effectively with recency-based  
 181 buffers, we next ask: *What accounts for the performance degradation observed when using the*  
 182 *coresnet buffer?* By design, the coresnet construction procedure favors diversity, selecting samples  
 183 that are dissimilar from one another. While this promotes broad coverage of the state space, it may  
 184 inadvertently introduce challenges.

185 One plausible explanation is that the resulting buffer becomes highly off-policy, containing older  
 186 samples that no longer reflect the agent’s current behavior, thereby reducing their usefulness for  
 187 learning [4]. Figure 5 supports this hypothesis by illustrating a modification to the composite buffer  
 188 in which the oldest sample in the coresnet are periodically ejected. Specifically, we remove the oldest  
 189 sample every  $k$  steps, evaluating values of  $k = 10, 20, 50$ . We find that more frequent ejection (i.e.,  
 190 lower values of  $k$ ) yields better performance, effectively making the coresnet buffer behave more like  
 191 a recency buffer—an effect clearly reflected in the improved learning outcomes.

#### 192 **4.2 Value-Based Distance Representation**

193 In general, it is not obvious how one can choose a distance representation to adequately address the  
 194 question of sample similarity in an environment. Can we use what the model has learned as a way  
 195 to represent samples? Specifically, can we use the current learned action-value function  $\hat{q}(s, \cdot)$  as a  
 196 distance representation? Does the vector  $\hat{q}(s, \cdot)$  serve as a good distance representation for coresnet  
 197 construction?

198 In order to answer the previous question, we test this in two different control environments, namely  
 199 Mountain Car [14] and Acrobot [19]. For comparison, we use more than one Q-Learning method, by  
 200 using Expected QRC (EQRC) [5] (See Figure 6b and 6d). In order to show the utility of our method

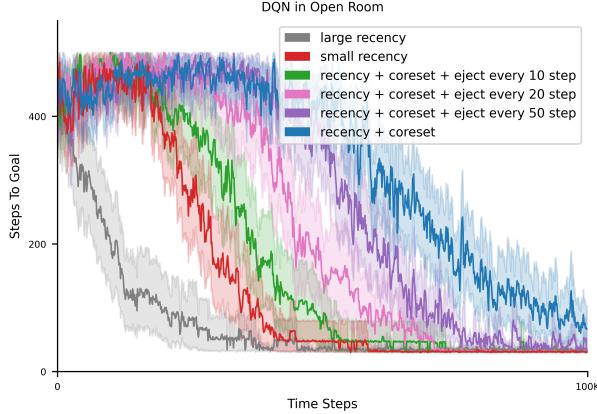


Figure 5: Performance of DQN in the Open Room domain using various configurations of the composite replay buffer. Each line color represents a different buffer variant. The key comparison is between the blue and green curves, which illustrate the performance difference between a composite buffer with recency and one without.

201 for small buffer sizes, we identify the smallest buffer size at which the recency-based method's  
 202 performance begins to deteriorate. In such a setting, we can see that using the composite replay buffer  
 203 using a value-based distance representation  $\phi$  recovers performance. The experiment setup is the  
 204 same as in 3.

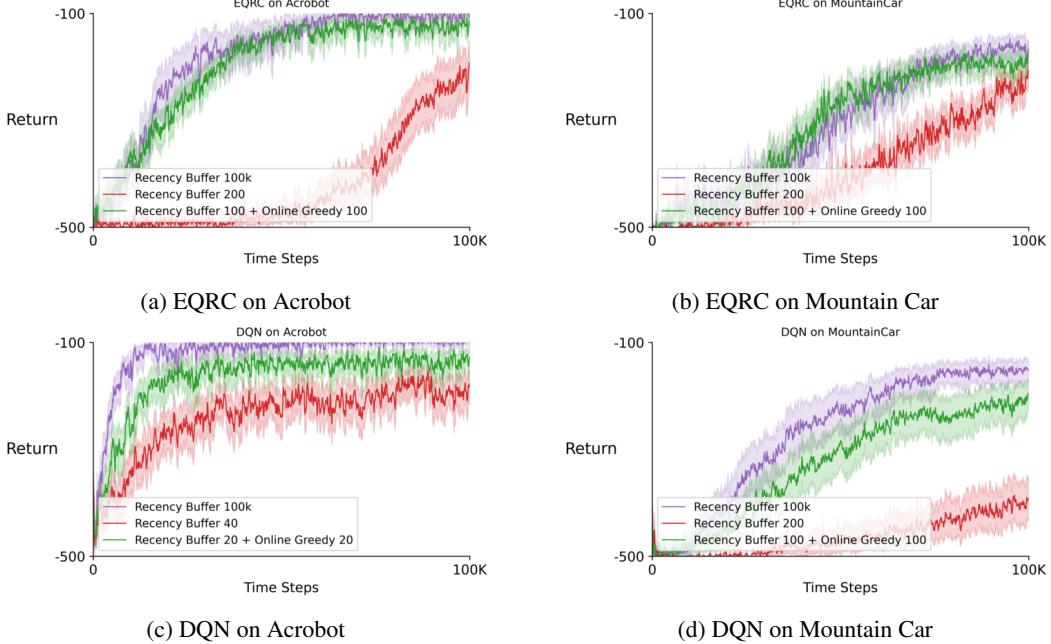


Figure 6: Performance comparison on Acrobot and Mountain Car using a value-based distance representation for coresnet construction. Subfigures (a) and (b) show results with Expected QRC (EQRC), while (c) and (d) show results with DQN. In each case, we compare a small and large recency buffer against a composite buffer using the learned action-value function  $\hat{q}(s, \cdot)$  as a distance representation. The choice of the size of the small and composite buffers are based on when the small buffer fails. The composite buffer improves performance, in this setting where a equally-sized recency buffer alone is not sufficient due to limited capacity.

### 205 4.3 Distance Criterion

206 The computational cost of Online Greedy prohibits the feasibility of experiments with larger coreset  
 207 buffer sizes. To address this, we adapt a coresnet construction method [6], that scales  $O(n^2)$  with the  
 208 size of the coresnet buffer.

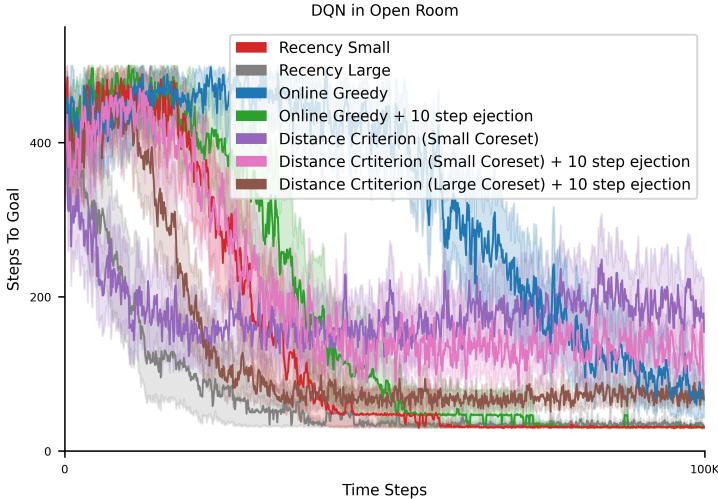


Figure 7: Performance of DQN in the Open Room environment using various replay strategies. The y-axis shows the number of steps required to reach the goal (lower is better), and the x-axis shows the number of training time steps. We compare recency-based buffers (small and large), Online Greedy (with and without 10-step ejection), and distance-based coresets of varying sizes, with and without an ejection mechanism. Ejection improves performance for Online Greedy. Notably, recency-based replay consistently achieves the best performance, while distance-based coresets without ejection struggle to match their effectiveness, especially later in learning.

209 For the distance-based coreset criterion, we use the following hyperparameter values: coreset size  
 210 of 100, an initial diversity threshold  $\delta = 0.05$ , a scaling factor  $\eta = 1.2$ , and an ejection rate  $k = 10$ ,  
 211 which determines how frequently the oldest transition is removed. We observe that, under both  
 212 the ejection and non-ejection variants, the agent fails to learn the optimal policy. To ensure a fair  
 213 comparison, we adopt the same configuration used in Figure 3.

214 The distance-based criterion offers a useful way to improve learning. It helps the agent learn faster  
 215 with a smaller buffer, especially early in training. However, performance tends to drop as training  
 216 continues. This may be because the coreset stops changing once it is full. To test this idea, we added  
 217 a simple ejection mechanism. This change slightly hurts early performance but improves learning  
 218 later on.

## 219 5 Limitations and Future Work

220 While coreset construction shows promise as a direction for improving experience replay mechanisms,  
 221 there are several limitations that warrant further consideration. The assumption that diversity is a  
 222 good criterion for learning remains unproven. There is a need for more theoretical work to determine  
 223 whether diversity is an appropriate and effective objective, or whether alternative criteria could offer  
 224 better performance. The relationship between the chosen criterion for coreset construction and the  
 225 resulting agent performance is not yet fully understood. In particular, the extent to which diversity  
 226 correlates with the utility of selected samples requires deeper investigation. This work assumes the  
 227 ability to freely choose a distance representation for samples, but does not explore or recommend  
 228 which types of representations are most beneficial.

229 There are several directions for future work that could build on this investigation. One promising  
 230 direction is to evaluate the coreset construction approach in additional domains, such as the Atari  
 231 Suite [3]. Understanding the impact of different distance representations on performance, and the  
 232 types of samples prioritized through diversity maximization, requires more investigation. Moreover,  
 233 a broader exploration of alternative diversity criteria and coreset construction methods is warranted.  
 234 The current study does not exhaust the space of possible strategies, and other methods may yield  
 235 better results or improved efficiency. Finally, the computational cost of online coreset construction is

236 non-trivial. Reducing this cost is an important area for future research, particularly for applications  
237 requiring real-time or resource-constrained execution.

## 238 **6 Conclusion**

239 In this paper, we explored the use of coreset-based construction strategies for experience replay in  
240 reinforcement learning. We examined two distinct coreset mechanisms and introduced a value-based  
241 distance representation, which enables compatibility with value-learning methods such as DQN.  
242 Our experiments show that coreset-augmented replay buffers can outperform simple recency-based  
243 baselines in terms of learning performance in some settings. Notably, our learned distance metric  
244 enables the replay buffer to retain diverse transitions, improving learning performance in constrained  
245 memory settings. These findings highlight the potential of coreset-augmented replay as a flexible and  
246 efficient alternative to traditional experience replay strategies.

## 247 **References**

- 248 [1] Marcin Andrychowicz, Filip Wolski, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder,  
249 Bob McGrew, Josh Tobin, OpenAI Pieter Abbeel, and Wojciech Zaremba. Hindsight experience  
250 replay. *Advances in neural information processing systems*, 30, 2017.
- 251 [2] Andrew Barto and Richard S. Sutton. Reinforcement learning: an introduction. 2018.
- 252 [3] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning  
253 environment: An evaluation platform for general agents. *Journal of artificial intelligence  
254 research*, 47:253–279, 2013.
- 255 [4] William Fedus, Prajit Ramachandran, Rishabh Agarwal, Yoshua Bengio, Hugo Larochelle,  
256 Mark Rowland, and Will Dabney. Revisiting fundamentals of experience replay. In Hal Daumé  
257 III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine  
258 Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3061–3071. PMLR,  
259 13–18 Jul 2020.
- 260 [5] Sina Ghassian, Andrew Patterson, Shivam Garg, Dhawal Gupta, Adam White, and Martha  
261 White. Gradient temporal-difference learning with regularized corrections. In Hal Daumé  
262 III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine  
263 Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 3524–3534. PMLR,  
264 13–18 Jul 2020.
- 265 [6] Paul Honeine. Analyzing sparse dictionaries for online learning with kernels. *IEEE Transactions  
266 on Signal Processing*, 63(23):6343–6353, 2015.
- 267 [7] Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2nd  
268 edition, 2013.
- 269 [8] Angelos Katharopoulos and François Fleuret. Not all samples are created equal: Deep learning  
270 with importance sampling. In *International conference on machine learning*, pages 2525–2534.  
271 PMLR, 2018.
- 272 [9] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint  
273 arXiv:1412.6980*, 2014.
- 274 [10] Long-Ji Lin. Self-improving reactive agents based on reinforcement learning, planning and  
275 teaching. *Machine learning*, 8:293–321, 1992.
- 276 [11] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of  
277 machine learning models. In *International Conference on Machine Learning*, pages 6950–6960.  
278 PMLR, 2020.
- 279 [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan  
280 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint  
281 arXiv:1312.5602*, 2013.

- 282 [13] Parham Mohammad Panahi. Investigating the interplay of prioritized replay and generalization.  
283 *Proceedings of the 1st Reinforcement Learning Conference (RLC 2025)*, 2024.
- 284 [14] Andrew William Moore. Efficient memory-based learning for robot control. Technical report,  
285 University of Cambridge, Computer Laboratory, 1990.
- 286 [15] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay.  
287 *arXiv preprint arXiv:1511.05952*, 2015.
- 288 [16] Matthew Schlegel, Yangchen Pan, Jiecao Chen, and Martha White. Adapting kernel rep-  
289 resentations online using submodular maximization. In Doina Precup and Yee Whye Teh,  
290 editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of  
291 *Proceedings of Machine Learning Research*, pages 3037–3046. PMLR, 06–11 Aug 2017.
- 292 [17] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driess-  
293 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mas-  
294 tering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489,  
295 2016.
- 296 [18] Richard S Sutton. Learning to predict by the methods of temporal differences. *Machine learning*,  
297 3:9–44, 1988.
- 298 [19] Richard S Sutton. Generalization in reinforcement learning: Successful examples using sparse  
299 coarse coding. *Advances in neural information processing systems*, 8, 1995.
- 300 [20] Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8:279–292, 1992.
- 301 [21] Yu Yang, Hao Kang, and Baharan Mirzasoleiman. Towards sustainable learning: Coresets  
302 for data-efficient deep learning. In *International Conference on Machine Learning*, pages  
303 39314–39330. PMLR, 2023.
- 304 [22] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. Online coresnet selection  
305 for rehearsal-based continual learning. In *Proceedings of the 10th International Conference on*  
306 *Learning Representations (ICLR)*, 2022.
- 307 [23] Shangtong Zhang and Richard S Sutton. A deeper look at experience replay. *NIPS 2017 Deep*  
308 *Reinforcement Learning Symposium*, 2017.