

# LEARN THE TIME TO LEARN: REPLAY SCHEDULING IN CONTINUAL LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Replay methods are known to be successful at mitigating catastrophic forgetting in continual learning scenarios despite having limited access to historical data. However, storing historical data is cheap in many real-world settings, yet replaying all historical data is often prohibited due to processing time constraints. In such settings, we propose that continual learning systems should learn the time to learn and schedule which tasks to replay at different time steps. We first demonstrate the benefits of our proposal by using Monte Carlo tree search to find a proper replay schedule, and show that the found replay schedules can outperform fixed scheduling policies when combined with multiple replay methods in various continual learning settings. Additionally, we propose a framework for learning replay scheduling policies with reinforcement learning. We show that the learned policies can generalize better in new continual learning scenarios compared to equally replaying all seen tasks, without added computational cost. Our study reveals the importance of learning the time to learn in continual learning, which brings current research closer to real-world needs.

## 1 INTRODUCTION

Many organizations deploying machine learning systems receive large volumes of data daily (Bailis et al., 2017; Hazelwood et al., 2018). Although all historical data are stored in the cloud in practice, retraining machine learning systems on a daily basis is prohibitive both in time and cost. In this setting, the systems often need to continuously adapt to new tasks while retaining the previously learned abilities. Continual learning (CL) methods (Delange et al., 2021; Parisi et al., 2019) address this challenge where, in particular, replay methods (Chaudhry et al., 2019; Hayes et al., 2020) have shown to be effective in achieving great prediction performance. Replay methods mitigate catastrophic forgetting by revisiting a small set of samples, which is feasible to process compared to the size of the historical data. In the traditional CL literature, replay memories are limited due to the assumption that historical data are not available. In real-world settings where historical data are always available, the requirement of small memories remains due to processing time and cost issues.

Recent research on replay-based CL has focused on the quality of memory samples (Aljundi et al., 2019b; Borsos et al., 2020; Chaudhry et al., 2019; Nguyen et al., 2017; Rebuffi et al., 2017; Yoon et al., 2021) or data compression to increase the memory capacity (Hayes et al., 2020; Iscen et al., 2020; Pellegrini et al., 2019). Most previous methods allocate equal memory storage space for samples from old tasks, and replay the whole memory to mitigate catastrophic forgetting. However, in life-long learning settings, this simple strategy would be inefficient as the memory must store a large number of tasks. Furthermore, the commonly used uniform selection policy of samples to revisit ignores the time of which tasks to learn again. This stands in contrast to human learning where education methods focus on scheduling of learning and rehearsal of previous learned knowledge. For example, spaced repetition (Dempster, 1989; Ebbinghaus, 2013; Landauer & Bjork, 1977), where the time interval between rehearsal increases, has been shown to enhance memory retention.

We argue that finding the proper schedule of which tasks to replay in fixed memory settings is critical for CL. To demonstrate our claim, we perform a simple experiment on the Split MNIST (Zenke et al., 2017) dataset where each task consists of learning the digits 0/1, 2/3, etc. arriving in sequence. The replay memory contains 10 samples from task 1 and can only be replayed while learning one of the tasks. Figure 1 shows how the task performances progress over time when the memory is replayed at

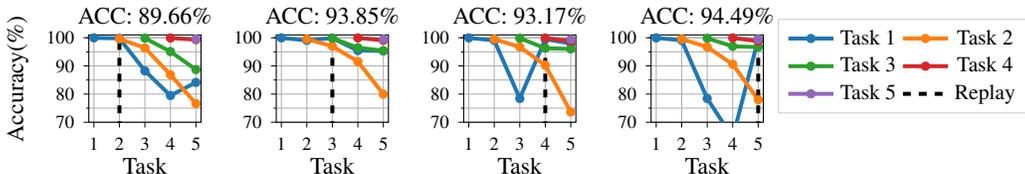


Figure 1: Task accuracies on Split MNIST (Zenke et al., 2017) when replaying only 10 samples of classes 0/1 at a single time step. The black vertical line indicates when replay is used. ACC denotes the average accuracy over all tasks after learning Task 5. Results are averaged over 5 seeds. These results show that the time to replay the previous task is critical for the final performance.

different tasks. In this example, the best average performance is achieved when the memory is used when learning task 5. Note that choosing different time points to replay the same memory leads to noticeably different classification performance. These results indicate that scheduling the time when to apply replay can influence the final performance significantly of a CL system.

In this paper, we propose learning the time to learn for CL systems, in which we learn replay schedules of which tasks to replay at different times inspired from human learning (Dempster, 1989). To demonstrate the benefits of replay scheduling, we perform experiments in an ideal CL environment where multiple trials are allowed to enable searching for the optimal replay schedule. We use Monte Carlo tree search (MCTS) (Coulom, 2006) to find proper replay schedules, which are evaluated by measuring the task performances of a network trained in a CL scenario in where the scheduled replay samples are used for reducing catastrophic forgetting. Furthermore, as using MCTS in real-world large scale CL tasks is infeasible, we propose a framework using reinforcement learning (RL) (Sutton & Barto, 2018) for learning replay scheduling policies. Our goal is to learn a general policy that has implicitly explored task relationships during training, such that the policy can be applied to mitigate catastrophic forgetting in new CL scenarios without additional training at test time. We evaluate the learned policy by comparing its ability to schedule the replay tasks against fixed scheduling policies, such as equally replaying all tasks. In summary, our contributions are:

- We propose a new CL setting where historical data is available while the processing time is limited, in order to adjust current CL research closer to real-world needs (Section 3.1). In this new setting, we introduce replay scheduling where we learn the time of which tasks to replay (Section 3.2).
- To demonstrate the benefits of replay scheduling, we apply MCTS in an ideal CL environment where MCTS searches over a finite set of replay memories at every task (Section 3.2). We show that the found replay schedules efficiently mitigate catastrophic forgetting across multiple benchmarks for various memory selection and replay methods in various CL scenarios (Section 4.1).
- To enable replay scheduling in real-world CL scenarios, propose an RL-based framework for learning policies that can mitigate catastrophic forgetting across different CL environments (Section 3.3). We show that the learned policies can outperform equally replaying all tasks in CL scenarios with new task orders and datasets unseen during training (Section 4.2).

## 2 RELATED WORK

In this section, we give a brief overview of various approaches in CL, especially replay methods. We provide more details on the related work, including spaced repetition in human CL (Dempster, 1989; Hawley et al., 2008; Landauer & Bjork, 1977) and generalization in RL (Igl et al., 2019; Kirk et al., 2021; Zhang et al., 2018a), in Appendix A. Traditional CL can be divided into three main areas, namely regularization-based, architecture-based, and replay-based approaches. Regularization-based methods protect parameters influencing the performance on known tasks from wide changes and use the other parameters for learning new tasks (Adel et al., 2019; Kao et al., 2021; Kirkpatrick et al., 2017; Li & Hoiem, 2017; Nguyen et al., 2017; Schwarz et al., 2018). Architecture-based methods mitigate catastrophic forgetting by maintaining task-specific parameters (Ebrahimi et al., 2020; Mallya & Lazebnik, 2018; Rusu et al., 2016; Serra et al., 2018; Xu & Zhu, 2018; Yoon et al., 2017). Replay methods mix samples from old tasks with the current dataset to mitigate catastrophic forgetting, where the replay samples are stored in an external memory (Aljundi et al., 2019a;b; Borsos et al., 2020; Chaudhry et al., 2019; Chrysakis & Moens, 2020; Hayes et al., 2019, 2020; Jin et al., 2020; Lopez-Paz & Ranzato, 2017; Pellegrini et al., 2019; Rolnick et al., 2018; Verwimp

et al., 2021; Yoon et al., 2021). Selecting the time to replay old tasks has mostly been ignored in the literature, with an exception in Aljundi et al. (2019a) which replays memory samples that would most interfere with a foreseen parameter update. Our replay scheduling approach differs from the above mentioned works since we focus on learning to select which tasks to replay. Nevertheless, our scheduling can be combined with any selection strategy and replay method.

### 3 METHOD

Here, we describe our new problem setting of CL where historical data are available while the processing time is limited when learning new tasks. In Section 3.1 and 3.2, we present the considered problem setting, as well as our idea of learning schedules over which tasks to replay at different time steps to mitigate catastrophic forgetting. Section 3.2 also describes how we use MCTS (Coulom, 2006) to study the benefits of replay scheduling in CL. In Section 3.3, we present an RL-based framework for learning replay scheduling policies that can generalize to different CL scenarios.

#### 3.1 PROBLEM SETTING

We focus on a slightly new setting in CL, where we assume that all historical data is available for mitigating catastrophic forgetting since data storage is cheap. However, as this data volume is typically huge, retraining on all historical data whenever the CL system must adapt to new tasks is impractical. Therefore, we assume there are processing time constraints which limits the system to sample a small replay memory from the historical data only once when adapting to new tasks. The challenge becomes how to select which old tasks to fill the replay memory with, such that the CL system achieves the best possible accuracy and minimize forgetting across all tasks.

The notation of our problem setting resembles the traditional CL setting for image classification. We let the network  $f_\phi$ , parameterized by  $\phi$ , learn  $T$  tasks sequentially from the datasets  $\mathcal{D}_1, \dots, \mathcal{D}_T$  arriving one at a time. The  $t$ -th dataset  $\mathcal{D}_t = \{(\mathbf{x}_t^{(i)}, y_t^{(i)})\}_{i=1}^{N_t}$  consists of  $N_t$  samples where  $\mathbf{x}_t^{(i)}$  and  $y_t^{(i)}$  are the  $i$ -th data point and class label respectively. Furthermore, each dataset is split into a training, validation, and test set, i.e.,  $\mathcal{D}_t = \{\mathcal{D}_t^{(train)}, \mathcal{D}_t^{(val)}, \mathcal{D}_t^{(test)}\}$ . The objective at task  $t$  is to minimize the loss  $\ell(f_\phi(\mathbf{x}_t), y_t)$  where  $\ell(\cdot)$  is the cross-entropy loss in our case.

We assume that historical data from old tasks are accessible at any task  $t$ . However, due to processing time constraints, we can only use a small replay memory  $\mathcal{M}$  of  $M$  historical samples for replay when learning a new task. The challenge then becomes how to select the  $M$  replay samples to efficiently retain knowledge of old tasks. We focus on selecting the samples on task-level by deciding on the task proportion  $(p_1, \dots, p_{t-1})$  of samples to fetch from each task, where  $p_i \geq 0$  is the proportion of  $M$  samples from task  $i$  to place in  $\mathcal{M}$  and  $\sum_{i=1}^{t-1} p_i = 1$ . To simplify the selection of which tasks to replay, we construct a discrete set of possible task proportions that can be used for constructing  $\mathcal{M}$ .

#### 3.2 REPLAY SCHEDULING IN CONTINUAL LEARNING

In this section, we describe our setup for enabling the scheduling for selecting replay memories at different time steps. We define a replay schedule as a sequence  $S = (\mathbf{p}_1, \dots, \mathbf{p}_{T-1})$ , where the task proportions  $\mathbf{p}_i = (p_1, \dots, p_{T-1})$  for  $1 \leq i \leq T-1$  are used for determining how many samples from seen tasks with which to fill the replay memory at task  $i$ . We construct an action space with a discrete number of choices of task proportions that can be selected at each task: At task  $t$ , we have  $t-1$  historical tasks that we can choose samples from. We create  $t-1$  bins  $\mathbf{b}_t = [b_1, \dots, b_{t-1}]$  and sample a task index for each bin  $b_i \in \{1, \dots, t-1\}$ . The bins are treated as interchangeable and we only keep the unique choices. For example, at task 3, we have seen task 1 and 2, so the unique choices of vectors are  $[1, 1], [1, 2], [2, 2]$ , where  $[1, 1]$  indicates that all memory samples are from task 1,  $[1, 2]$  indicates that half memory is from task 1 and the other half are from task etc. We count the number of occurrences of each task index in  $\mathbf{b}_t$  and divide by  $t-1$  to obtain the task proportion, i.e.,  $\mathbf{p}_t = \text{bincount}(\mathbf{b}_t)/(t-1)$ . We round the number of replay samples from task  $i$ , i.e.,  $p_i \cdot M$ , up or down accordingly to keep the memory size  $M$  fixed when filling the memory. From this specification, we can build a tree of different replay schedules to evaluate with the network.

Figure 2 shows an example of a replay schedule tree with Split MNIST where the memory size is  $M = 8$ . Each level corresponds to a CL task, and we show some examples of possible replay memories in the tree that can be evaluated at each task. A replay schedule is represented as a path traversal of different replay memory compositions from task 1 to task 5. At task 1, the memory  $\mathcal{M}_1 = \emptyset$  is empty, while  $\mathcal{M}_2$  is filled with samples from task 1 at task 2. The memory  $\mathcal{M}_3$  can be composed with samples from either task 1 or 2, or equally fill  $\mathcal{M}_3$  with samples from both tasks. All possible paths in the tree are valid replay schedules.

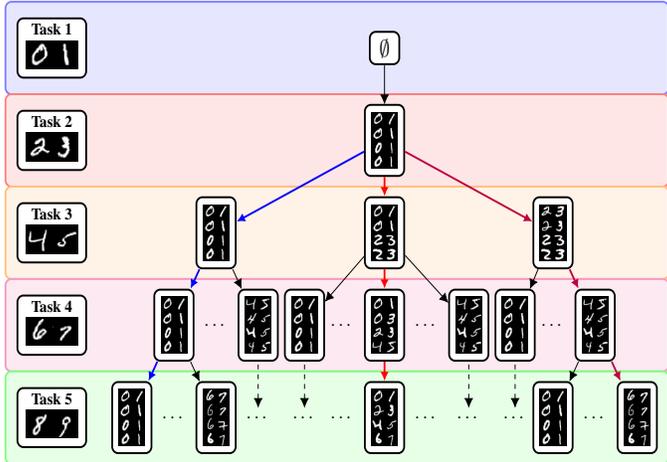


Figure 2: Tree-shaped action space of possible replay memories of size  $M = 8$  at every task for Split MNIST.

We show three examples of possible schedules in Figure 2 for illustration: the blue path represents a replay schedule where only task 1 samples are replayed. The red path represents using memories with equally distributed tasks, and the purple path represents a schedule where the memory is only filled with samples from the most previous task.

**Monte Carlo Tree Search for Replay Schedules.** The tree-shaped action space of task proportions grows fast with the number of tasks. This complicates studying replay scheduling in datasets with long task-horizons, where the action space is too large for using exhaustive searches. We propose to use MCTS since it has been successful in applications with large action spaces (Browne et al., 2012; Chaudhry & Lee, 2018; Silver et al., 2016). We apply MCTS in an ideal setting with multiple rollouts allowed to demonstrate that replay scheduling can be essential for the final CL performance, where MCTS concentrates the search in directions with promising outcomes in the CL environment.

Each replay memory composition in the action space corresponds to a node that MCTS can visit, as can be seen in Figure 2. At level  $t$ , the node  $v_t$  is related to a task proportion  $p_t$  used for retrieving a replay memory composition from the historical data at task  $t$ . One MCTS rollout corresponds to traversing through all tree levels  $1, \dots, T$  to select the replay schedule  $S$  to use during the CL training. Each task proportion  $p_t$  from every visited node is stored in  $S$  during the rollout. When reaching level  $T$ , we start the CL training and use  $S$  for constructing the replay memories at each task. Next, we briefly outline the MCTS steps for performing the search (details in Appendix B.1):

- **Selection.** During a rollout, the current node  $v_t$  either moves randomly to an unvisited child, or selects the next child node  $v_{t+1}$  by evaluating the Upper Confidence Tree (UCT) (Kocsis & Szepesvári, 2006) with the function from Chaudhry & Lee (2018):

$$UCT(v_t, v_{t+1}) = \max(q(v_{t+1})) + C \sqrt{\frac{2 \log(n(v_t))}{n(v_{t+1})}}, \quad (1)$$

where  $q(\cdot)$  is the reward function,  $C$  the exploration constant, and  $n(\cdot)$  the number of node visits.

- **Expansion.** Whenever the current node  $v_t$  has unvisited child nodes, the search tree is expanded with one of the unvisited child nodes  $v_{t+1}$  selected with uniform sampling.
- **Simulation and Reward.** After expansion, the succeeding nodes are selected randomly until reaching a terminal node  $v_T$ . The task proportions from the visited rollout nodes constitutes the replay schedule  $S$ . After training the network using  $S$  for replay, we calculate the reward for the rollout given by  $r = \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(val)}$ , where  $A_{T,i}^{(val)}$  is the validation accuracy of task  $i$  at task  $T$ .
- **Backpropagation.** Reward  $r$  is backpropagated from the expanded node  $v_t$  to the root  $v_1$ , where the reward function  $q(\cdot)$  and number of visits  $n(\cdot)$  are updated at each node.

### 3.3 POLICY LEARNING FRAMEWORK FOR REPLAY SCHEDULING

In this section, we present an RL-based framework for learning replay scheduling policies. We focus on learning a general policy that can be applied in any CL scenario to mitigate catastrophic

forgetting to avoid re-training the policy for every new CL dataset, which would be useful in user personalization applications. Our intuition is that there may exist general patterns regarding replay scheduling, e.g., that tasks that are harder or have been forgotten should be replayed more often. We aim to implicitly explore such task properties by using the task performances of the CL network as states for the policy to select which tasks to replay. Representing the states with task performances also enables transferring the learned policy to reduce forgetting in unseen CL environments. Next, we present our modeling approach for learning the replay scheduling policies using RL.

**CL Environment.** We model the CL environments as Markov Decision Processes (Bellman, 1957) (MDPs) where each MDP is represented as a tuple  $E_i = (\mathcal{S}_i, \mathcal{A}, P_i, R_i, \mu_i, \gamma)$  consisting of the state space  $\mathcal{S}_i$ , action space  $\mathcal{A}$ , state transition probability  $P_i(s'|s, a)$ , reward function  $R_i(s, a)$ , initial state distribution  $\mu_i(s_1)$ , and discount factor  $\gamma$ . Each environment  $E_i$  contains a network  $f_\phi$  and task datasets  $\mathcal{D}_{1:T}$  where the  $t$ -th dataset is learned at time step  $t$ . The state  $s_t$  is defined as the task accuracies  $A_{t,1:t}$  evaluated at task  $t$ , such that  $s_t = [A_{t,1}, \dots, A_{t,t}, 0, \dots, 0]$  where zero-padding is used on future tasks. We obtain the states by evaluating the classifier on the validation datasets to avoid overfitting the policy to the training data. We use the same action space as for MCTS (see Section 3.2), such that  $a_t \in \mathcal{A}$  corresponds to a task proportion  $p_t$  used for sampling the replay memory  $\mathcal{M}_t$ . The state transition distribution  $P_i(s'|s, a)$  represents the dynamics of the environment, which depend on the initialization of  $f_\phi$  and the task order in  $\mathcal{D}_{1:T}$ . We use a dense reward defined as the average validation accuracies at task  $t$ , i.e.,  $r_t = \frac{1}{t} \sum_{i=1}^t A_{t,i}^{(val)}$ , to ease exploration in the action space. The goal for the agent is to maximize the rewards during an episode.

**Policy Training and Evaluation.** The policy interacts with the CL environments by selecting which tasks the network  $f_\phi$  should replay to mitigate catastrophic forgetting. The state  $s_t$  is obtained by evaluating  $f_\phi$  on the validation sets  $\mathcal{D}_{1:t}^{(val)}$  after learning task  $t$ . The action  $a_t$  is selected under the policy  $\pi_\theta(a|s_t)$ , parameterized by  $\theta$ , which is converted into the task proportion  $p_t$  for sampling the replay memory  $\mathcal{M}_t$  from the historical datasets. The network  $f_\phi$  is trained on task  $t+1$  while replaying  $\mathcal{M}_t$ , and we obtain the reward  $r_{t+1}$  and the next state  $s_{t+1}$  by evaluating  $f_\phi$  on the validation sets  $\mathcal{D}_{1:t+1}^{(val)}$ . The collected transitions  $(s_t, a_t, r_{t+1}, s_{t+1})$  are used for updating the policy, and a new episode starts after  $f_\phi$  has learned the final task  $T$ . We let the policy interact with multiple training environments  $\mathcal{E}^{(train)} = \{E_i\}_{i=1}^K$  sampled from a distribution of CL environments, i.e.,  $E_i \sim p(E)$ . To generate diverse CL environments, we let each  $E_i$  have different network initializations of  $f_\phi$  and task orders in the datasets. Our goal is to learn a general replay scheduling policy that can be applied in new CL environments to mitigate catastrophic forgetting. Hence, in Section 4.2, we evaluate the policy in CL environments with new task orders or datasets unseen during training. The policy is applied for only a single CL episode without additional training in the test environment.

## 4 EXPERIMENTS

In this section, we present the experimental results to show the importance of replay scheduling in CL. First, we demonstrate the benefits with replay scheduling by using MCTS for finding replay schedules in Section 4.1. Thereafter, we evaluate our RL-based framework using DQN (Mnih et al., 2013) and A2C (Mnih et al., 2016) for learning policies that generalize to new CL scenarios in Section 4.2. Full details on experimental settings and additional results are in Appendix D and E.

### 4.1 RESULTS ON REPLAY SCHEDULING WITH MONTE CARLO TREE SEARCH

In this section, we show the benefits of replay scheduling in single CL environments using MCTS. We perform extensive evaluation where we apply MCTS with different memory selection and replay methods, varying memory sizes in different CL settings (Van de Ven & Tolias, 2019), and show the potential efficiency of replay scheduling in a tiny memory setting.

**Experimental Setup.** We conduct experiments on several CL benchmark datasets: Split MNIST (LeCun et al., 1998; Zenke et al., 2017), FashionMNIST (Xiao et al., 2017), Split notMNIST (Bulatov, 2011), Permuted MNIST (Goodfellow et al., 2013), and CIFAR-100 (Krizhevsky & Hinton, 2009), and Split miniImagenet (Vinyals et al., 2016). We use a 2-layer MLP with 256 hidden units for Split MNIST, Split FashionMNIST, Split notMNIST, and Permuted MNIST. We apply

Table 1: Performance comparison between MCTS (Ours) and the baselines with various memory selection methods, namely uniform sampling,  $k$ -means, and Mean-of-Features (MoF).

Memory	Schedule	Split MNIST		Split FashionMNIST		Split notMNIST	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Offline	Joint	99.75 ± 0.06	0.01 ± 0.06	99.34 ± 0.08	-0.01 ± 0.14	96.12 ± 0.57	-0.21 ± 0.71
Uniform	Random	94.91 ± 2.52	-6.13 ± 3.16	95.89 ± 2.03	-4.33 ± 2.55	91.84 ± 1.48	-5.37 ± 2.12
	ETS	94.02 ± 4.25	-7.22 ± 5.33	95.81 ± 3.53	-4.45 ± 4.34	91.01 ± 1.39	-6.16 ± 1.82
	Heur-GD	96.02 ± 2.32	-4.64 ± 2.90	97.09 ± 0.62	-2.82 ± 0.84	91.26 ± 3.99	-6.06 ± 4.70
	MCTS	97.93 ± 0.56	-2.27 ± 0.71	98.27 ± 0.17	-1.29 ± 0.20	94.64 ± 0.39	-1.47 ± 0.79
k-means	Random	92.65 ± 1.38	-8.96 ± 1.74	93.11 ± 2.75	-7.76 ± 3.42	93.11 ± 1.01	-3.78 ± 1.43
	ETS	92.89 ± 3.53	-8.66 ± 4.42	96.47 ± 0.85	-3.55 ± 1.07	93.80 ± 0.82	-2.84 ± 0.81
	Heur-GD	96.28 ± 1.68	-4.32 ± 2.11	95.78 ± 1.50	-4.46 ± 1.87	91.75 ± 0.94	-5.60 ± 2.07
	MCTS	98.20 ± 0.16	-1.94 ± 0.22	98.48 ± 0.26	-1.04 ± 0.31	93.61 ± 0.71	-3.11 ± 0.55
MoF	Random	96.96 ± 1.34	-3.57 ± 1.69	96.39 ± 1.69	-3.66 ± 2.17	93.09 ± 1.40	-3.70 ± 1.76
	ETS	97.04 ± 1.23	-3.46 ± 1.50	96.48 ± 1.33	-3.55 ± 1.73	92.64 ± 0.87	-4.57 ± 1.59
	Heur-GD	96.46 ± 2.41	-4.09 ± 3.01	95.84 ± 0.89	-4.39 ± 1.15	93.24 ± 0.77	-3.48 ± 1.37
	MCTS	98.37 ± 0.24	-1.70 ± 0.28	97.84 ± 0.32	-1.81 ± 0.39	94.62 ± 0.42	-1.80 ± 0.56

Memory	Schedule	Permuted MNIST		Split CIFAR-100		Split miniImagenet	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Offline	Joint	95.34 ± 0.13	0.17 ± 0.18	84.73 ± 0.81	-1.06 ± 0.81	74.03 ± 0.83	9.70 ± 0.68
Uniform	Random	72.59 ± 1.52	-25.71 ± 1.76	53.76 ± 1.80	-35.11 ± 1.93	49.89 ± 1.03	-14.79 ± 1.14
	ETS	71.09 ± 2.31	-27.39 ± 2.59	47.70 ± 2.16	-41.69 ± 2.37	46.97 ± 1.24	-18.32 ± 1.34
	Heur-GD	76.68 ± 2.13	-20.82 ± 2.41	57.31 ± 1.21	-30.76 ± 1.45	49.66 ± 1.10	-12.04 ± 0.59
	MCTS	76.34 ± 0.98	-21.21 ± 1.16	56.60 ± 1.13	-31.39 ± 1.11	50.20 ± 0.72	-13.46 ± 1.22
k-means	Random	71.91 ± 1.24	-26.45 ± 1.34	53.20 ± 1.44	-35.77 ± 1.31	49.96 ± 1.46	-14.81 ± 1.18
	ETS	69.40 ± 1.32	-29.23 ± 1.47	47.51 ± 1.14	-41.77 ± 1.30	45.82 ± 0.92	-19.53 ± 1.10
	Heur-GD	75.57 ± 1.18	-22.11 ± 1.22	54.31 ± 3.94	-33.80 ± 4.24	49.25 ± 1.00	-12.92 ± 1.22
	MCTS	77.74 ± 0.80	-19.66 ± 0.95	56.95 ± 0.92	-30.92 ± 0.83	50.47 ± 0.85	-13.31 ± 1.24
MoF	Random	78.80 ± 1.07	-18.79 ± 1.16	62.35 ± 1.24	-26.33 ± 1.25	56.02 ± 1.11	-7.99 ± 1.13
	ETS	77.62 ± 1.12	-20.10 ± 1.26	60.43 ± 1.17	-28.22 ± 1.26	56.12 ± 1.12	-8.93 ± 0.83
	Heur-GD	77.27 ± 1.45	-20.15 ± 1.63	55.60 ± 2.70	-32.57 ± 2.77	52.30 ± 0.59	-9.61 ± 0.67
	MCTS	81.58 ± 0.75	-15.41 ± 0.86	64.22 ± 0.65	-23.48 ± 1.02	57.70 ± 0.51	-5.31 ± 0.55

the ConvNet from Schwarz et al. (2018); Vinyals et al. (2016) for Split CIFAR-100, and the reduced ResNet-18 from Lopez-Paz & Ranzato (2017) for Split miniImagenet. We use multi-head output layers and assume task labels are available at test time unless stated otherwise, except for Permuted MNIST where single-head output layer is used. We measure the CL performances using ACC as the average test accuracy across tasks and BWT for forgetting (Lopez-Paz & Ranzato, 2017), i.e.,

$$\text{ACC} = \frac{1}{T} \sum_{i=1}^T A_{T,i}, \quad \text{BWT} = \frac{1}{T-1} \sum_{i=1}^{T-1} A_{T,i} - A_{i,i}, \quad (2)$$

where  $A_{t,i}$  is the test accuracy for task  $i$  after learning task  $t$ . We compare MCTS to the baselines:

- **Random.** Random policy that randomly selects task proportions from the action space on how to structure the replay memory at every task.
- **Equal Task Schedule (ETS).** Policy that selects equal task proportion such that the replay memory aims to fill the memory with an equal number of samples from every seen task.
- **Heuristic Global Drop (Heur-GD).** Heuristic policy that replays tasks with validation accuracy below a certain threshold proportional to the best achieved validation accuracy on the task.

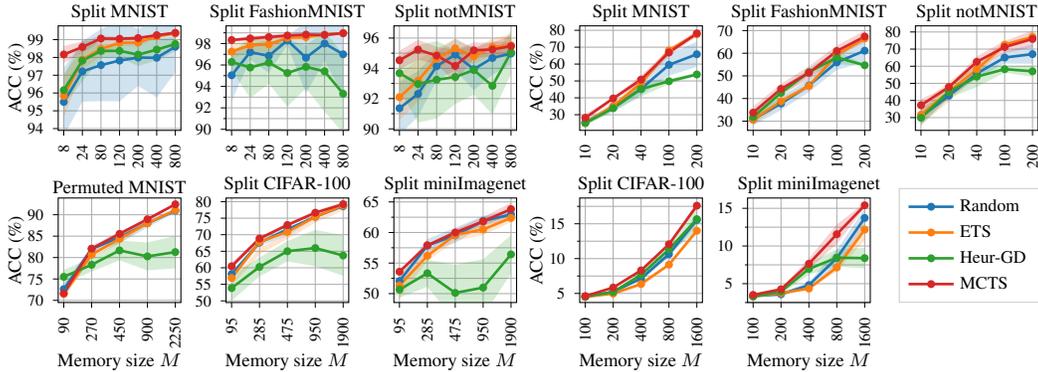
Heur-GD is based on the intuition that forgotten tasks should be replayed. The replay memory is filled with  $M/k$  samples per task, where  $k$  is the number of selected tasks, but skips replay if  $k = 0$ . MCTS and Heur-GD randomly sample 15% of the training data of each task to use for validation. For MCTS, reported results are evaluated on the test set by using replay schedules selected from the validation sets. Memory sizes are set to  $M = 10$  for Split MNIST, Split FashionMNIST, and Split notMNIST, and  $M = 100$  for Permuted MNIST, Split CIFAR-100, and Split miniImagenet, unless stated otherwise. We report means and standard deviations using 5 seeds on all datasets.

**Combine with Different Memory Selection Methods.** We show that our method can be combined with any memory selection method for storing replay samples. In addition to uniform sampling, we apply various memory selection methods commonly used in the CL literature, namely  $k$ -means clustering and Mean-of-Features (MoF) (Rebuffi et al., 2017). Table 1 shows the results across all datasets. We note that using the replay schedule from MCTS outperforms the baselines when using the alternative selection methods, where MoF performs the best on most datasets.

**Applying Scheduling to Recent Replay Methods.** In this experiment, we show that replay scheduling can be combined with any replay method to enhance the CL performance. We combine

Table 2: Performance comparison between scheduling methods MCTS (Ours), Random, ETS, and Heuristic combined with replay methods HAL, MER, and DER.

Method	Schedule	Split MNIST		Split CIFAR-100		Split minilmagenet	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
HAL	Random	96.32 ± 1.77	-3.90 ± 2.28	35.90 ± 2.47	-17.37 ± 3.76	40.86 ± 1.86	-5.12 ± 2.23
	ETS	97.21 ± 1.25	-2.80 ± 1.59	34.90 ± 2.02	-18.92 ± 0.91	38.13 ± 1.18	-8.19 ± 1.73
	Heur-GD	97.69 ± 0.19	-2.22 ± 0.24	35.07 ± 1.29	-24.76 ± 2.41	39.51 ± 1.49	-5.65 ± 0.77
	MCTS	97.96 ± 0.15	-1.85 ± 0.18	40.22 ± 1.57	-12.77 ± 1.30	41.39 ± 1.15	-3.69 ± 1.86
MER	Random	93.00 ± 3.22	-7.96 ± 4.15	42.68 ± 0.86	-35.56 ± 1.39	32.86 ± 0.95	-7.71 ± 0.45
	ETS	92.97 ± 1.73	-8.52 ± 2.15	43.38 ± 1.81	-34.84 ± 1.98	33.58 ± 1.53	-6.80 ± 1.46
	Heur-GD	94.30 ± 2.79	-6.46 ± 3.50	40.90 ± 1.70	-44.10 ± 2.03	34.22 ± 1.93	-7.57 ± 1.63
	MCTS	96.44 ± 0.72	-4.14 ± 0.94	44.29 ± 0.69	-32.73 ± 0.88	32.74 ± 1.29	-5.77 ± 1.04
DER	Random	95.91 ± 2.18	-4.40 ± 2.46	56.17 ± 1.30	-29.03 ± 1.38	35.13 ± 4.11	-10.85 ± 2.92
	ETS	98.17 ± 0.35	-2.00 ± 0.42	52.58 ± 1.49	-32.93 ± 2.04	35.50 ± 2.84	-10.94 ± 2.21
	Heur-GD	94.57 ± 1.71	-6.08 ± 2.09	55.75 ± 1.08	-31.27 ± 1.02	43.62 ± 0.88	-8.18 ± 1.16
	MCTS	99.02 ± 0.10	-0.91 ± 0.13	58.99 ± 0.98	-24.95 ± 0.64	43.46 ± 0.95	-9.32 ± 1.37



(a) Task/Domain-IL

(b) Class-IL

Figure 4: Performance comparison over various memory sizes for the methods, where (a) shows results in the Task- and Domain-Incremental Learning (IL) settings, and (b) in the Class-IL setting.

MCTS with Hindsight Anchor Learning (HAL) (Chaudhry et al., 2021), Meta-Experience Replay (MER) (Riemer et al., 2018), Dark Experience Replay (DER) (Buzzega et al., 2020). Table 2 shows the performance comparison between our the MCTS scheduling against using Random, ETS, and Heuristic schedules for each method. The results confirm that replay scheduling is important for the final performance given the same memory constraints and it can benefit any existing CL framework.

**Replay Schedule Visualization.** We visualize a replay schedule from Split CIFAR-100 with memory size  $M = 100$  to gain insights into the behavior of the scheduling policy from MCTS. Figure 3 shows a bubble plot of the selected task proportions used for filling the replay memory at every task. Each circle color corresponds to a replay task, and its size represents the proportion of replay samples at the current task. The sum of points in all circles at each column is fixed at all current tasks. The task proportions vary dynamically over time in a sophisticated nonlinear way which would be hard to replace by a heuristic method. Moreover, we can observe space repetition-style scheduling on many tasks, e.g., task 1-3 are replayed with similar proportion at the initial tasks but eventually starts varying the time interval between replay. Also, task 4 and 6 need less replay in their early stages, which could potentially be that they are simpler or correlated with other tasks. We provide a similar visualization for Split MNIST in Appendix D.3.

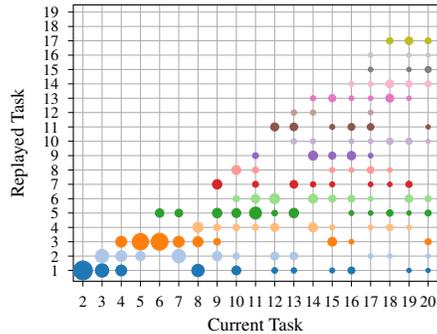


Figure 3: Replay schedule from MCTS on Split CIFAR-100 visualized as bubble plot.

**Varying Memory Size.** We show that our method can improve the CL performance across varying memory sizes in different CL scenarios. Figure 4a shows the results in the Task- and Domain-Incremental Learning (IL) scenarios, where we observe that MCTS generally obtains better task accuracies than ETS, especially for small memory sizes. Both MCTS and ETS perform better

Table 3: Performance comparison in memory setting where only 1 sample/class is available from the historical data for replay. The baselines replay all available samples, while MCTS selects 2 samples for Split MNIST and 50 samples for Permuted MNIST and Split miniImagenet.

Method	Split MNIST		Permuted MNIST		Split miniImagenet	
	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Random	92.56 ± 2.90	-8.97 ± 3.62	70.02 ± 1.76	-28.22 ± 1.92	48.85 ± 1.38	-14.55 ± 1.86
A-GEM	94.97 ± 1.50	-6.03 ± 1.87	64.71 ± 1.78	-34.41 ± 2.05	32.06 ± 1.83	-30.81 ± 1.79
ER-Ring	94.94 ± 1.56	-6.07 ± 1.92	69.73 ± 1.13	-28.87 ± 1.29	49.82 ± 1.69	-14.38 ± 1.57
Uniform	95.77 ± 1.12	-5.02 ± 1.39	69.85 ± 1.01	-28.74 ± 1.17	50.56 ± 1.07	-13.52 ± 1.34
MCTS	96.07 ± 1.60	-4.59 ± 2.01	72.52 ± 0.54	-25.43 ± 0.65	50.70 ± 0.54	-12.60 ± 1.13

than Heur-GD as  $M$  increases, which shows that Heur-GD requires careful tuning of the validation thresholds. We also performed experiments in the Class-IL scenario where task labels are absent. Here, the replay memory is always filled with at least 1 sample/class to avoid fully forgetting non-replayed tasks. Each scheduling method then selects which tasks to replay out of the remaining samples. Figure 4b shows that ETS approaches MCTS when  $M$  increases on the 5-task datasets. However, on the more challenging Split CIFAR-100 and Split miniImagenet, MCTS outperforms ETS clearly as  $M$  increases. These results show that selecting the proper replay schedule is essential in various CL scenarios with both small and large datasets across different backbone choices.

**Efficiency of Replay Scheduling.** We illustrate the efficiency of replay scheduling in a setting where only 1 sample/class is available from the historical data for replay. We consider the scenario where the replay memory size is smaller than the number of classes. The replay memory size for MCTS is set to  $M = 2$  for the 5-task datasets, such that only 2 samples can be selected for replay from the seen tasks. For the larger CL datasets, we set  $M = 50$ . We then compare against the memory efficient CL baselines A-GEM (Chaudhry et al., 2018b) and ER-Ring (Chaudhry et al., 2019), as well as uniform memory selection. Table 3 shows that MCTS, despite using significantly fewer samples for replay, performs mostly on par with the baselines and outperforms them on Permuted MNIST. These results indicate that replay scheduling is an important research direction in CL, since storing every seen class in the memory could be inefficient in settings with large number of tasks.

## 4.2 POLICY GENERALIZATION TO NEW CONTINUAL LEARNING SCENARIOS

In this section, we evaluate how well the learned replay scheduling policies can mitigate catastrophic forgetting in new CL environments. We employ DQN and A2C for policy learning and evaluate their ability to generalize in CL environments with new task orders and datasets unseen during training.

**Experimental Setup.** We conduct experiments on Split MNIST, Split FashionMNIST, Split notMNIST, and Split CIFAR-10 (Krizhevsky & Hinton, 2009). The CL setting is in general the same as in Section 4.1. We evaluate all methods on 10 different test environments, and assess the generalization capability by ranking all methods by comparing their measured ACC per seed in each test environment, since the performance between environments can vary significantly (details in Appendix E.2). The policy is applied for only a single pass over the CL tasks at test time. We add two baselines:

- **Heuristic Local Drop (Heur-LD).** Heuristic policy that replays tasks with validation accuracy below a threshold proportional to the previous achieved validation accuracy on the task.
- **Heuristic Accuracy Threshold (Heur-AT).** Heuristic policy that replays tasks with validation accuracy below a fixed threshold.

Memory sizes are set to  $M = 10$ , and we average the results over 5 seeds.

**Generalization to New Task Orders.** We show that the learned replay scheduling policies can generalize to CL environments with previously unseen task orders. The training and test environments are generated with unique task orders of the CL datasets. The columns **New Task Order** in Table 4 shows the average ranking for the DQN, A2C, and the baselines when being applied in the 10 test environments. Our learned policies obtain the best average ranking across most datasets, where A2C performs better than DQN. To provide further insights, Figure 5 shows the task accuracy progress and the corresponding replay schedule from A2C and ETS from one Split CIFAR-10 test environment. In Figure 5. The replay schedules are visualized with bubble plots showing the selected task proportion to use for composing the replay memories at each task. In Figure 5a, we observe that A2C decides to replay task 2 more than task 1 as the performance on task 2 decreases,

Table 4: Average ranking (lower is better) across methods in the policy generalization experiments. The best and second-best ranks are colored in green and orange respectively.

Method	New Task Order				New Dataset	
	S-MNIST	S-FashionMNIST	S-notMNIST	S-CIFAR-10	S-FashionMNIST	S-notMNIST
Random	3.98	3.44	3.68	4.91	3.99	3.94
ETS	3.82	4.56	4.44	5.38	3.68	4.06
Heur-GD	4.53	4.23	3.44	4.03	2.86	4.61
Heur-LD	4.67	3.63	3.96	3.63	5.08	4.96
Heur-AT	4.38	3.96	5.5	3.43	3.75	4.29
DQN (Ours)	3.46	3.65	3.51	3.83	4.42	3.4
A2C (Ours)	3.16	4.53	3.47	2.79	4.22	2.74

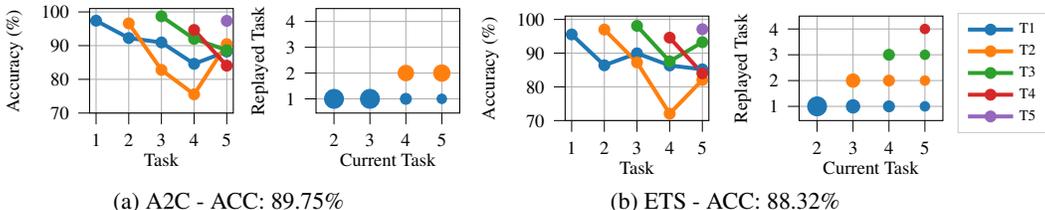


Figure 5: Task accuracies and replay schedules for A2C and ETS for a Split CIFAR-10 environment.

which results in a slightly better ACC metric achieved by A2C than ETS. These results show that the learned policy can flexibly consider replaying forgotten tasks to enhance the CL performance.

**Generalization to New Datasets.** We show that the learned replay scheduling policies are capable of generalizing to CL environments with new datasets unseen in the training environments. We perform two sets of experiments, 1) train with environments generated with Split MNIST and FashionMNIST and test on environments generated with Split notMNIST, and 2) train with environments generated with Split MNIST and notMNIST and test on environments generated with Split FashionMNIST. The columns **New Dataset** in Table 4 shows the average ranking for DQN, A2C, and the baselines when generalization to test environments with the new datasets. We observe that both A2C and DQN successfully generalize to Split notMNIST compared to the baselines. However, the learned RL policies have difficulties generalizing to Split FashionMNIST environments, which could be due to high variations in the state transition dynamics between training and test environments. This shows that learning the replay scheduling policies using RL inherits common challenges with generalization in RL, such as robustness to domain shifts. Potentially, the performance could be improved by generating more training environments for the agent to exhibit more variations of CL scenarios, or by using other advanced RL methods which may generalize better (Igl et al., 2019).

## 5 CONCLUSIONS

We proposed learning the time to learn, i.e., in a real-world CL context, learning schedules of which tasks to replay at different times. To the best of our knowledge, we are the first to consider the time to learn in CL inspired by human learning techniques. We demonstrated the benefits with replay scheduling in CL by showing on several CL benchmarks that replay schedules found with MCTS can outperform replaying all tasks equally or relying on heuristic scheduling rules. Furthermore, we proposed an RL-based framework for learning scheduling policies to enable replay scheduling for real-world CL scenarios. The learned policies are agnostic to the CL dataset, and can be applied to reduce catastrophic forgetting in new CL scenarios without additional training, which would be useful in user personalization applications. Our replay scheduling approach brings current research closer to tackling real-world CL challenges where the number of tasks exceeds the memory size.

**Limitations and Future Work.** Generalization in RL is a challenging research topic by itself. With the current method, large amounts of diverse data and training time is required to enable the learned policy to generalize well. This can be costly since generating the CL environments is expensive as each state transition involves training the network on a CL task. Moreover, we are currently considering a discrete action space which is hard to construct, especially in large-scale CL scenarios. Thus, in future work, we would explore more advanced RL methods which can handle continuous actions and generalize well.

## REFERENCES

- Tameem Adel, Han Zhao, and Richard E Turner. Continual learning with adaptive weights (claw). *arXiv preprint arXiv:1911.09514*, 2019.
- Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. *Advances in neural information processing systems*, 32, 2019a.
- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *arXiv preprint arXiv:1903.08671*, 2019b.
- Hadi Amiri, Timothy Miller, and Guergana Savova. Repeat before forgetting: Spaced repetition for efficient and effective training of neural networks. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 2401–2410, 2017.
- Peter Bailis, Edward Gan, Samuel Madden, Deepak Narayanan, Kexin Rong, and Sahaana Suri. Macrobases: Prioritizing attention in fast data. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pp. 541–556, 2017.
- Philip J Ball, Yingzhen Li, Angus Lamb, and Cheng Zhang. A study on efficiency in continual learning inspired by human learning. *arXiv preprint arXiv:2010.15187*, 2020.
- Philip J Ball, Cong Lu, Jack Parker-Holder, and Stephen Roberts. Augmented world models facilitate zero-shot dynamics generalization from a single offline environment. In *International Conference on Machine Learning*, pp. 619–629. PMLR, 2021.
- Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5): 679–684, 1957.
- Zalán Borsos, Mojmír Mutný, and Andreas Krause. Coresets via bilevel optimization for continual learning and streaming. *arXiv preprint arXiv:2006.03875*, 2020.
- Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfschagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43, 2012.
- Yaroslav Bulatov. The notMNIST dataset. <http://yaroslavvb.com/upload/notMNIST/>, 2011.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *arXiv preprint arXiv:2004.07211*, 2020.
- Yash Chandak, Georgios Theodorou, James Kostas, Scott Jordan, and Philip Thomas. Learning action representations for reinforcement learning. In *International conference on machine learning*, pp. 941–950. PMLR, 2019.
- Yash Chandak, Georgios Theodorou, Chris Nota, and Philip Thomas. Lifelong learning with a changing action set. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3373–3380, 2020.
- Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018a.
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018b.
- Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.

- Arslan Chaudhry, Albert Gordo, Puneet Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6993–7001, 2021.
- Muhammad Umar Chaudhry and Jee-Hyong Lee. Feature selection for high dimensional data using monte carlo tree search. *IEEE Access*, 6:76036–76048, 2018.
- Aristotelis Chrysakis and Marie-Francine Moens. Online continual learning from imbalanced data. In *ICML*, 2020.
- Karl Cobbe, Oleg Klimov, Chris Hesse, Taehoon Kim, and John Schulman. Quantifying generalization in reinforcement learning. In *International Conference on Machine Learning*, pp. 1282–1289. PMLR, 2019.
- Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pp. 72–83. Springer, 2006.
- Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Frank N Dempster. Spacing effects and their implications for theory and practice. *Educational Psychology Review*, 1(4):309–330, 1989.
- Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, Yuhuai Wu, and Peter Zhokhov. Openai baselines. <https://github.com/openai/baselines>, 2017.
- Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer Vision – ECCV 2020*, pp. 86–102. Springer International Publishing, 2020.
- John Dunlosky, Katherine A. Rawson, Elizabeth J. Marsh, Mitchell J. Nathan, and Daniel T. Willingham. Improving students’ learning with effective learning techniques: Promising directions from cognitive and educational psychology. *Psychological Science in the Public Interest*, 14(1): 4–58, 2013. ISSN 15291006. URL <http://www.jstor.org/stable/23484712>.
- Hermann Ebbinghaus. Memory: A contribution to experimental psychology. *Annals of neurosciences*, 20(4):155, 2013.
- Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. *arXiv preprint arXiv:2003.09553*, 2020.
- Jesse Farebrother, Marlos C Machado, and Michael Bowling. Generalization and regularization in dqn. *arXiv preprint arXiv:1810.00123*, 2018.
- Kanyin Feng, Xiao Zhao, Jing Liu, Ying Cai, Zhifang Ye, Chuansheng Chen, and Gui Xue. Spaced learning enhances episodic memory by increasing neural pattern similarity across repetitions. *Journal of Neuroscience*, 39(27):5351–5360, 2019.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pp. 1126–1135. PMLR, 2017.
- Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- Karri S Hawley, Katie E Cherry, Emily O Boudreaux, and Erin M Jackson. A comparison of adjusted spaced retrieval versus a uniform expanded retrieval schedule for learning a name–face association in older adults with probable alzheimer’s disease. *Journal of Clinical and Experimental Neuropsychology*, 30(6):639–649, 2008.
- Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pp. 9769–9776. IEEE, 2019.

- Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. In *European Conference on Computer Vision*, pp. 466–483. Springer, 2020.
- Kim Hazelwood, Sarah Bird, David Brooks, Soumith Chintala, Utku Diril, Dmytro Dzhulgakov, Mohamed Fawzy, Bill Jia, Yangqing Jia, Aditya Kalro, et al. Applied machine learning at facebook: A datacenter infrastructure perspective. In *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 620–629. IEEE, 2018.
- Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- Irina Higgins, Arka Pal, Andrei Rusu, Loic Matthey, Christopher Burgess, Alexander Pritzel, Matthew Botvinick, Charles Blundell, and Alexander Lerchner. Darla: Improving zero-shot transfer in reinforcement learning. In *International Conference on Machine Learning*, pp. 1480–1490. PMLR, 2017.
- Maximilian Igl, Kamil Ciosek, Yingzhen Li, Sebastian Tschjatschek, Cheng Zhang, Sam Devlin, and Katja Hofmann. Generalization in reinforcement learning with selective noise injection and information bottleneck. *Advances in neural information processing systems*, 32, 2019.
- Maximilian Igl, Gregory Farquhar, Jelena Luketina, Wendelin Boehmer, and Shimon Whiteson. Transient non-stationarity and generalisation in deep reinforcement learning. *arXiv preprint arXiv:2006.05826*, 2020.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456. PMLR, 2015.
- Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. In *European Conference on Computer Vision*, pp. 699–715. Springer, 2020.
- David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Ayush Jain, Andrew Szot, and Joseph J Lim. Generalization to new actions in reinforcement learning. *arXiv preprint arXiv:2011.01928*, 2020.
- Xisen Jin, Arka Sadhu, Junyi Du, and Xiang Ren. Gradient based memory editing for task-free continual learning. *arXiv preprint arXiv:2006.15294*, 2020.
- KJ Joseph and Vineeth N Balasubramanian. Meta-consolidation for continual learning. *arXiv preprint arXiv:2010.00352*, 2020.
- Ta-Chu Kao, Kristopher Jensen, Gido van de Ven, Alberto Bernacchia, and Guillaume Hennequin. Natural continual learning: success is a journey, not (just) a destination. *Advances in Neural Information Processing Systems*, 34, 2021.
- Samuel Kessler, Jack Parker-Holder, Philip Ball, Stefan Zohren, and Stephen J Roberts. Same state, different task: Continual reinforcement learning without interference. *arXiv preprint arXiv:2106.02940*, 2021.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A survey of generalisation in deep reinforcement learning. *arXiv preprint arXiv:2111.09794*, 2021.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.

- Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pp. 282–293. Springer, 2006.
- Ilya Kostrikov. Pytorch implementations of reinforcement learning algorithms. <https://github.com/ikostrikov/pytorch-a2c-ppo-acktr-gail>, 2018.
- Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- T. Landauer and Robert Bjork. Optimum rehearsal patterns and name learning. *Practical aspects of memory*, 1, 11 1977.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *arXiv preprint arXiv:1706.08840*, 2017.
- Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7765–7773, 2018.
- Seyed Iman Mirzadeh and Hassan Ghasemzadeh. Cl-gym: Full-featured pytorch library for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, pp. 3621–3627, June 2021.
- Seyed Iman Mirzadeh, Mehrdad Farajtabar, Dilan Gorur, Razvan Pascanu, and Hassan Ghasemzadeh. Linear mode connectivity in multitask and continual learning. *arXiv preprint arXiv:2010.04495*, 2020.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937. PMLR, 2016.
- Anusha Nagabandi, Ignasi Clavera, Simin Liu, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt in dynamic, real-world environments through meta-reinforcement learning. *arXiv preprint arXiv:1803.11347*, 2018.
- Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. *arXiv preprint arXiv:1710.10628*, 2017.
- Pingbo Pan, Siddharth Swaroop, Alexander Immer, Runa Eschenhagen, Richard E Turner, and Mohammad Emtiyaz Khan. Continual deep learning by functional regularisation of memorable past. *arXiv preprint arXiv:2004.14070*, 2020.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- Lorenzo Pellegrini, Gabriele Graffieti, Vincenzo Lomonaco, and Davide Maltoni. Latent replay for real-time continual learning. *arXiv preprint arXiv:1912.01100*, 2019.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.

- Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. Experience replay for continual learning. *arXiv preprint arXiv:1811.11682*, 2018.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pp. 4528–4537. PMLR, 2018.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pp. 4548–4557. PMLR, 2018.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *arXiv preprint arXiv:1705.08690*, 2017.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Paul Smolen, Yili Zhang, and John H Byrne. The right time to learn: mechanisms and optimization of spaced learning. *Nature Reviews Neuroscience*, 17(2):77, 2016.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pp. 23–30. IEEE, 2017.
- Gido M van de Ven and Andreas S Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.
- Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- Gido M van de Ven, Hava T Siegelmann, and Andreas S Tolias. Brain-inspired replay for continual learning with artificial neural networks. *Nature communications*, 11(1):1–14, 2020.
- Eli Verwimp, Matthias De Lange, and Tinne Tuytelaars. Rehearsal revealed: The limits and merits of revisiting samples in continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9385–9394, 2021.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. *arXiv preprint arXiv:1606.04080*, 2016.
- Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.
- Johannes von Oswald, Christian Henning, João Sacramento, and Benjamin F Grewe. Continual learning with hypernetworks. *arXiv preprint arXiv:1906.00695*, 2019.
- Kaixin Wang, Bingyi Kang, Jie Shao, and Jiashi Feng. Improving generalization in reinforcement learning with mixture regularization. *Advances in Neural Information Processing Systems*, 33: 7968–7978, 2020.
- Judy Willis. Review of research: Brain-based teaching strategies for improving students’ memory, learning, and test-taking success. *Childhood Education*, 83(5):310–315, 2007.

- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Ju Xu and Zhanxing Zhu. Reinforced continual learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
- Jaehong Yoon, Saehoon Kim, Eunho Yang, and Sung Ju Hwang. Scalable and order-robust continual learning with additive parameter decomposition. *arXiv preprint arXiv:1902.09432*, 2019.
- Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju Hwang. Online coreset selection for rehearsal-based continual learning. *arXiv preprint arXiv:2106.01085*, 2021.
- Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Relational deep reinforcement learning. *arXiv preprint arXiv:1806.01830*, 2018.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pp. 3987–3995. PMLR, 2017.
- Amy Zhang, Nicolas Ballas, and Joelle Pineau. A dissection of overfitting and generalization in continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*, 2018a.
- Chiyuan Zhang, Oriol Vinyals, Remi Munos, and Samy Bengio. A study on overfitting in deep reinforcement learning. *arXiv preprint arXiv:1804.06893*, 2018b.
- Chenyang Zhao, Olivier Sigaud, Freerk Stulp, and Timothy M Hospedales. Investigating generalisation in continuous deep reinforcement learning. *arXiv preprint arXiv:1902.07015*, 2019.

## APPENDIX

This supplementary material is structured as follows:

- Appendix **A**: Extended related work.
- Appendix **B**: Additional information on the methodology of MCTS for finding replay schedules and our RL-based framework for policy learning.
- Appendix **C**: Additional information on the heuristic scheduling baselines and hyperparameters.
- Appendix **D**: Additional experimental settings and results for Section 4.1.
- Appendix **E**: Additional experimental settings and results for Section 4.2.
- Our code is provided in a zip-file `code.zip` as part of the supplementary material. The code will be made publicly available upon acceptance.

## A EXTENDED RELATED WORK

In this section, we give a brief overview of various approaches in CL, especially replay methods as well as spaced repetition techniques for human CL and generalization in RL.

**Continual Learning.** Traditional CL can be divided into three main areas, namely regularization-based, architecture-based, and replay-based approaches. Regularization-based methods protect parameters influencing the performance on known tasks from wide changes and use the other parameters for learning new tasks (Adel et al., 2019; Kao et al., 2021; Kirkpatrick et al., 2017; Li & Hoiem, 2017; Nguyen et al., 2017; Schwarz et al., 2018; Zenke et al., 2017). Architecture-based methods mitigate catastrophic forgetting by maintaining task-specific parameters (Ebrahimi et al., 2020; Mallya & Lazebnik, 2018; Rusu et al., 2016; Serra et al., 2018; Xu & Zhu, 2018; Yoon et al., 2019; 2017). Replay methods mix samples from old tasks with the current dataset to mitigate catastrophic forgetting, where the replay samples are either stored in an external memory (Aljundi et al., 2019a;b;

Borsos et al., 2020; Chaudhry et al., 2019; Chrysakis & Moens, 2020; Hayes et al., 2019; 2020; Iscen et al., 2020; Isele & Cosgun, 2018; Jin et al., 2020; Lopez-Paz & Ranzato, 2017; Nguyen et al., 2017; Pellegrini et al., 2019; Rebuffi et al., 2017; Rolnick et al., 2018; Verwimp et al., 2021; Yoon et al., 2021) or generated using a generative model (Shin et al., 2017; van de Ven & Tolia, 2018). Regularization-based approaches and dynamic architectures have been combined with replay-based approaches to methods to overcome their limitations (Buzzega et al., 2020; Chaudhry et al., 2018a;b; 2021; Douillard et al., 2020; Ebrahimi et al., 2020; Joseph & Balasubramanian, 2020; Mirzadeh et al., 2020; Nguyen et al., 2017; Pan et al., 2020; Pellegrini et al., 2019; Rolnick et al., 2018; von Oswald et al., 2019). Our work relates most to replay-based methods with external memory which we spend more time on describing in the next paragraph.

**Replay-based Continual Learning.** Much research effort in replay- or memory-based CL has focused on selecting higher quality samples to store in memory (Aljundi et al., 2019b; Borsos et al., 2020; Chaudhry et al., 2019; Chrysakis & Moens, 2020; Hayes et al., 2019; Isele & Cosgun, 2018; Lopez-Paz & Ranzato, 2017; Nguyen et al., 2017; Rebuffi et al., 2017; Yoon et al., 2021). In Chaudhry et al. (2019), several selection strategies are reviewed in scenarios with tiny memory capacity, such as reservoir sampling (Vitter, 1985), first-in first-out buffer (Lopez-Paz & Ranzato, 2017), k-Means, and Mean-of-Features (Rebuffi et al., 2017). However, elaborate selection strategies have been shown to give little benefit over random selection for image classification problems (Chaudhry et al., 2018a; Hayes et al., 2020). More recently, there has been work on compressing raw images to feature representations to increase the number of memory examples for replay (Hayes et al., 2020; Iscen et al., 2020; Pellegrini et al., 2019). Selecting the time to replay old tasks has mostly been ignored in the literature, with an exception in (Aljundi et al., 2019a) which replays memory samples that would most interfere with a foreseen parameter update. Our replay scheduling approach differs from the above mentioned works since we focus on learning to select which tasks to replay. Nevertheless, our scheduling can be combined with any selection strategy and replay-based method.

**Human Continual Learning.** Humans are CL systems in the sense of learning tasks and concepts sequentially. The timing of learning and rehearsal is essential for humans to memorize better (Dempster, 1989; Dunlosky et al., 2013; Willis, 2007). An example technique is spaced repetition where time intervals between rehearsal are gradually increased to improve long-term memory retention (Dempster, 1989; Ebbinghaus, 2013), which has been shown to improve memory retention better uniformly spaced rehearsal times (Hawley et al., 2008; Landauer & Bjork, 1977). Several works in CL with neural networks are inspired by human learning techniques, including spaced repetition (Amiri et al., 2017; Feng et al., 2019; Smolen et al., 2016), sleep mechanisms (Ball et al., 2020; Mallya & Lazebnik, 2018; Schwarz et al., 2018), and memory reactivation (Hayes et al., 2020; van de Ven et al., 2020). Replay scheduling is also inspired by spaced repetition, where we learn schedules of which tasks to replay at different times.

**Generalization in Reinforcement Learning.** Generalization is an active research topic in RL (Kirk et al., 2021) as RL agents tend to overfit to their training environments (Henderson et al., 2018; Zhang et al., 2018a; Zhao et al., 2019; Zhang et al., 2018b). The goal is often to transfer learned policies to environments with new tasks (Finn et al., 2017; Kessler et al., 2021; Higgins et al., 2017) and action spaces (Chandak et al., 2019; Jain et al., 2020; Chandak et al., 2020). Some approaches aim to improve generalization capabilities by generating more diverse training data (Cobbe et al., 2019; Tobin et al., 2017; Wang et al., 2020; Zhang et al., 2018a), using network regularization or inductive biases (Farebrother et al., 2018; Igl et al., 2019; Zambaldi et al., 2018), or learning dynamics models (Ball et al., 2021; Nagabandi et al., 2018). In this paper, we use RL for learning policies for selecting which tasks a CL network should replay. The goal is to learn policies that can be applied in new CL environments for replay scheduling on unseen task orders and datasets without additional computational cost.

## B ADDITIONAL METHODOLOGY

In this section, we provide pseudo-code for MCTS to search for replay schedules in single CL environments in Section B.1 as well as pseudo-code for the RL-based framework for learning the replay scheduling policies in Section B.2.

## B.1 MONTE CARLO TREE SEARCH ALGORITHM FOR REPLAY SCHEDULING

**Algorithm 1** Monte Carlo Tree Search for Replay Scheduling

---

**Require:** Tree nodes  $v_{1:T}$ , Datasets  $\mathcal{D}_{1:T}$ , Learning rate  $\eta$   
**Require:** Replay memory size  $M$

- 1:  $ACC_{best} \leftarrow 0, S_{best} \leftarrow ()$
- 2: **while** within computational budget **do**
- 3:    $S \leftarrow ()$
- 4:    $v_t, S \leftarrow \text{TREEPOLICY}(v_1, S)$
- 5:    $v_T, S \leftarrow \text{DEFAULTPOLICY}(v_t, S)$
- 6:    $ACC \leftarrow \text{EVALUATEREPLAYSCHEDULE}(\mathcal{D}_{1:T}, S, M)$
- 7:    $\text{BACKPROPAGATE}(v_t, ACC)$
- 8:   **if**  $ACC > ACC_{best}$  **then**
- 9:      $ACC_{best} \leftarrow ACC$
- 10:     $S_{best} \leftarrow S$
- 11: **return**  $ACC_{best}, S_{best}$
  
- 12: **function**  $\text{TREEPOLICY}(v_t, S)$
- 13:   **while**  $v_t$  is non-terminal **do**
- 14:     **if**  $v_t$  not fully expanded **then**
- 15:       **return**  $\text{EXPANSION}(v_t, S)$
- 16:     **else**
- 17:        $v_t \leftarrow \text{BESTCHILD}(v_t)$
- 18:        $S.\text{append}(\mathbf{p}_t)$ , where  $\mathbf{p}_t \leftarrow \text{GETTASKPROPORTION}(v_t)$
- 19:   **return**  $v_t, S$
  
- 20: **function**  $\text{EXPANSION}(v_t, S)$
- 21:   Sample  $v_{t+1}$  uniformly among unvisited children of  $v_t$
- 22:    $S.\text{append}(\mathbf{p}_{t+1})$ , where  $\mathbf{p}_{t+1} \leftarrow \text{GETTASKPROPORTION}(v_{t+1})$
- 23:   Add new child  $v_{t+1}$  to node  $v_t$
- 24:   **return**  $v_{t+1}, S$
  
- 25: **function**  $\text{BESTCHILD}(v_t)$
- 26:    $v_{t+1} = \arg \max_{v_{t+1} \in \text{children of } v} \max(Q(v_{t+1})) + C\sqrt{\frac{2\log(N(v_t))}{N(v_{t+1})}}$
- 27:   **return**  $v_{t+1}$
  
- 28: **function**  $\text{DEFAULTPOLICY}(v_t, S)$
- 29:   **while**  $v_t$  is non-terminal **do**
- 30:     Sample  $v_{t+1}$  unFormly among children of  $v_t$
- 31:      $S.\text{append}(\mathbf{p}_{t+1})$ , where  $\mathbf{p}_{t+1} \leftarrow \text{GETTASKPROPORTION}(v_{t+1})$
- 32:     Update  $v_t \leftarrow v_{t+1}$
- 33:   **return**  $v_t, S$
  
- 34: **function**  $\text{EVALUATEREPLAYSCHEDULE}(\mathcal{D}_{1:T}, S, M)$
- 35:   Initialize neural network  $f_\theta$
- 36:   **for**  $t = 1, \dots, T$  **do**
- 37:      $\mathbf{p} \leftarrow S[t-1]$
- 38:      $\mathcal{M} \leftarrow \text{GETREPLAYMEMORY}(\mathcal{D}_{1:t-1}^{(train)}, \mathbf{p}, M)$
- 39:     **for**  $\mathcal{B} \sim \mathcal{D}_t^{(train)}$  **do**
- 40:        $\theta \leftarrow \text{SGD}(\mathcal{B} \cup \mathcal{M}, \theta, \eta)$
- 41:      $A_{1:T}^{(val)} \leftarrow \text{EVALUATEACCURACY}(f_\theta, \mathcal{D}_{1:T}^{(val)})$
- 42:      $ACC \leftarrow \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(val)}$
- 43:   **return**  $ACC$
  
- 44: **function**  $\text{BACKPROPAGATE}(v_t, R)$
- 45:   **while**  $v_t$  is not root **do**
- 46:      $N(v_t) \leftarrow N(v_t) + 1$
- 47:      $Q(v_t) \leftarrow R$
- 48:      $v_t \leftarrow \text{parent of } v_t$

---

We provide pseudo-code in Algorithm 1 outlining the steps for our method using Monte Carlo tree search (MCTS) to find replay schedules described in the main paper (Section 3.2). The MCTS procedure selects actions over which task proportions to fill the replay memory with at every task, where the selected task proportions are stored in the replay schedule  $S$ . The schedule is then passed to `EVALUATEREPLAYSCHEDULE( $\cdot$ )` where the continual learning part executes the training with replay memories filled according to the schedule. The reward for the schedule  $S$  is the average validation accuracy over all tasks after learning task  $T$ , i.e., `ACC`, which is backpropagated through the tree to update the statistics of the selected nodes. The schedule  $S_{best}$  yielding the best `ACC` score is returned to be used for evaluation on the held-out test sets.

The function `GETREPLAYMEMORY( $\cdot$ )` is the policy for retrieving the replay memory  $\mathcal{M}$  from the historical data given the task proportion  $\mathbf{p}$ . The number of samples per task determined by the task proportions are rounded up or down accordingly to fill  $\mathcal{M}$  with  $M$  replay samples in total. The function `GETTASKPROPORTION( $\cdot$ )` simply returns the task proportion that is related to given node.

The following steps are performed during one MCTS rollout (or iteration):

1. **Selection** involves either selecting an unvisited node randomly, or selecting the next node by evaluating the UCT score (see Equation 1) if all children has been visited already. In Algorithm 1, `TREEPOLICY( $\cdot$ )` appends the task proportions  $\mathbf{p}_t$  to the replay schedule  $S$  at every selected node.
2. **Expansion** involves expanding the search tree with one of the unvisited child nodes  $v_{t+1}$  selected with uniform sampling. `EXPANSION( $\cdot$ )` in Algorithm 1 appends the task proportions  $\mathbf{p}_t$  to the replay schedule  $S$  of the expanded node.
3. **Simulation** involves selecting the next nodes randomly until a terminal node  $v_T$  is reached. In Algorithm 1, `DEFAULTPOLICY( $\cdot$ )` appends the task proportions  $\mathbf{p}_t$  to the replay schedule  $S$  at every randomly selected node until reaching the terminal node.
4. **Reward** The reward for the rollout is given by the `ACC` of the validation sets for each task. In Algorithm 1, `EVALUATEREPLAYSCHEDULE( $\cdot$ )` involves learning the tasks  $t = 1, \dots, T$  sequentially and using the replay schedule to sample the replay memories to use for mitigating catastrophic forgetting when learning a new task. The reward  $r$  for the rollout is calculated after task  $T$  has been learnt.
5. **Backpropagation** involves updating the reward function  $q(\cdot)$  and number of visits  $n(\cdot)$  from the expansion node up to the root node. See `BACKRPROPAGATE( $\cdot$ )` in Algorithm 1.

## B.2 RL FRAMEWORK ALGORITHM

We provide pseudo-code for the RL-based framework for learning the replay scheduling policy with either DQN (Mnih et al., 2013) or A2C (Mnih et al., 2016) in Algorithm 2. The procedure collects experience from all training environments in  $\mathcal{E}^{(train)}$  at every time step  $t$ . The datasets and classifiers are specific for each environment  $E_i \in \mathcal{E}^{(train)}$ . At  $t = 1$ , we obtain the initial state  $s_1^{(i)}$  by evaluating the classifier on the validation set  $\mathcal{D}_1^{(val)}$  after training the classifier on the task 1. Next, we get the replay memory for mitigating catastrophic forgetting when learning the next task  $t + 1$  by 1) taking action  $a_t^{(i)}$  under policy  $\pi_\theta$ , 2) converting action  $a_t^{(i)}$  into the task proportion  $\mathbf{p}_t$ , and 3) sampling the replay memory  $\mathcal{M}_t$  from the historical datasets given the selected proportion. We then obtain the reward  $r_t$  and the next state  $s_{t+1}$  by evaluating the classifier on the validation sets  $\mathcal{D}_{1:t+1}^{(val)}$  after learning task  $t + 1$ . The collected experience from each time step is stored in the experience buffer  $\mathcal{B}$  for both DQN and A2C. In `UPDATEPOLICY( $\cdot$ )`, we outline the steps for updating the policy parameters  $\theta$  with either DQN or A2C.

**Algorithm 2** RL Framework for Learning Replay Scheduling Policy

---

**Require:**  $\mathcal{E}^{(train)}$ : Training environments,  $\theta$ : Policy parameters,  $\gamma$ : Discount factor  
**Require:**  $\eta$ : Learning rate,  $n_{episodes}$ : Number of episodes,  $M$ : Replay memory size  
**Require:**  $n_{steps}$ : Number of steps for A2C

- 1:  $\mathcal{B} = \{\}$  ▷ Initialize experience buffer
- 2: **for**  $i = 1, \dots, n_{episodes}$  **do**
- 3:     **for**  $t = 1, \dots, T - 1$  **do**
- 4:         **for**  $E_i \in \mathcal{E}^{(train)}$  **do**
- 5:              $\mathcal{D}_{1:t+1} = \text{GETDATASETS}(E_i, t)$  ▷ Get datasets from environment  $E_i$
- 6:              $f_{\phi}^{(i)} = \text{GETCLASSIFIER}(E_i)$  ▷ Get classifier from environment  $E_i$
- 7:             **if**  $t == 1$  **then**
- 8:                  $\text{TRAIN}(f_{\phi}^{(i)}, \mathcal{D}_t^{(train)})$  ▷ Train classifier  $f_{\phi}^{(i)}$  on task 1
- 9:                  $A_{1:t}^{(val)} = \text{EVAL}(f_{\phi}^{(i)}, \mathcal{D}_{1:t}^{(val)})$  ▷ Evaluate classifier  $f_{\phi}^{(i)}$  on task 1
- 10:                  $s_t^{(i)} = A_{1:t}^{(val)} = [A_{1,1}^{(val)}, 0, \dots, 0]$  ▷ Get initial state
- 11:                  $a_t^{(i)} \sim \pi_{\theta}(a, s_t^{(i)})$  ▷ Take action under policy  $\pi_{\theta}$
- 12:                  $\mathbf{p}_t = \text{GETTASKPROPORTION}(a_t^{(i)})$
- 13:                  $\mathcal{M}_t \sim \text{GETREPLAYMEMORY}(\mathcal{D}_{1:t}^{(train)}, \mathbf{p}_t, M)$
- 14:                  $\text{TRAIN}(f_{\phi}^{(i)}, \mathcal{D}_{t+1}^{(train)} \cup \mathcal{M}_t)$  ▷ Train classifier  $f_{\phi}^{(i)}$
- 15:                  $A_{1:t+1}^{(val)} = \text{EVAL}(f_{\phi}^{(i)}, \mathcal{D}_{1:t+1}^{(val)})$  ▷ Evaluate classifier  $f_{\phi}^{(i)}$
- 16:                  $s_{t+1}^{(i)} = A_{1:t+1}^{(val)} = [A_{t+1,1}^{(val)}, \dots, A_{t+1,t+1}^{(val)}, 0, \dots, 0]$  ▷ Get next state
- 17:                  $r_t^{(i)} = \frac{1}{t+1} \sum_{j=1}^{t+1} A_{1:j}^{(val)}$  ▷ Compute reward
- 18:                  $\mathcal{B} = \mathcal{B} \cup \{(s_t^{(i)}, a_t^{(i)}, r_t^{(i)}, s_{t+1}^{(i)})\}$  ▷ Store transition in buffer
- 19:                 **if** time to update policy **then**
- 20:                      $\theta, \mathcal{B} = \text{UPDATEPOLICY}(\theta, \mathcal{B}, \gamma, \eta, n_{steps})$  ▷ Update policy with experience
- 21:     **return**  $\theta$  ▷ Return policy
- 22: **function**  $\text{UPDATEPOLICY}(\theta, \mathcal{B}, \gamma, \eta, n_{steps})$
- 23:     **if** DQN **then**
- 24:          $(s_j, a_j, r_j, s'_j) \sim \mathcal{B}$  ▷ Sample mini-batch from buffer
- 25:          $y_j = \begin{cases} r_j & \text{if } s'_j \text{ is terminal} \\ r_j + \gamma \max_a Q_{\theta^-}(s'_j, a) & \text{else} \end{cases}$  ▷ Compute  $y_j$  with target net  $\theta^-$
- 26:          $\theta = \theta - \eta \nabla_{\theta} (y_j - Q_{\theta}(s_j, a_j))^2$  ▷ Update  $Q$ -function
- 27:     **else if** A2C **then**
- 28:          $s_t = \mathcal{B}[n_{steps}]$  ▷ Get last state in buffer
- 29:          $R = \begin{cases} 0 & \text{if } s_t \text{ is terminal} \\ V_{\theta_v}(s_t) & \text{else} \end{cases}$  ▷ Bootstrap from last state
- 30:         **for**  $j = n_{steps} - 1, \dots, 0$  **do**
- 31:              $s_j, a_j, r_j = \mathcal{B}[j]$  ▷ Get state, action, and reward at step  $j$
- 32:              $R = r_j + \gamma R$
- 33:              $\theta = \theta - \eta \nabla_{\theta} \log \pi_{\theta}(a_j, s_j)(R - V_{\theta_v}(s_j))$  ▷ Update policy
- 34:              $\theta_v = \theta_v - \eta \nabla_{\theta_v} (R - V_{\theta_v}(s_j))^2$  ▷ Update value function
- 35:          $\mathcal{B} = \{\}$  ▷ Reset experience buffer
- 36:     **return**  $\theta, \mathcal{B}$

---

## C HEURISTIC SCHEDULING BASELINES

We implemented three heuristic scheduling baselines to compare against our proposed methods. These heuristics are based on the intuition of re-learning tasks when they have been forgotten. We keep a validation set for each task to determine whether any task should be replayed by comparing the validation accuracy against a hand-tuned threshold. If the validation accuracy is below the threshold, then the corresponding task is replayed. Let  $A_{t,i}^{(val)}$  be the validation accuracy for task  $t$  evaluated at time step  $i$ . The threshold is set differently in each of the baselines:

- **Heuristic Global Drop (Heur-GD)**. Heuristic policy that replays tasks with validation accuracy below a certain threshold proportional to the best achieved validation accuracy on the task. The best achieved validation accuracy for task  $i$  is given by  $A_{t,i}^{(best)} = \max\{A_{1,i}^{(val)}, \dots, A_{t,i}^{(val)}\}$ . Task  $i$  is replayed if  $A_{t,i}^{(val)} < \tau A_{t,i}^{(best)}$  where  $\tau \in [0, 1]$  is a ratio representing the degree of how much the validation accuracy of a task is allowed to drop. Note that Heur-GD (denoted as Heuristic) is the only one used in the experiments with MCTS in single CL environments in Section 4.1.
- **Heuristic Local Drop (Heur-LD)**. Heuristic policy that replays tasks with validation accuracy below a threshold proportional to the previous achieved validation accuracy on the task. Task  $i$  is replayed if  $A_{t,i}^{(val)} < \tau A_{t-1,i}^{(val)}$  where  $\tau$  again represents the degree of how much the validation accuracy of a task is allowed to drop.
- **Heuristic Accuracy Threshold (Heur-AT)**. Heuristic policy that replays tasks with validation accuracy below a fixed threshold. Task  $i$  is replayed if  $A_{t,i}^{(val)} < \tau$  where  $\tau \in [0, 1]$  represents the least tolerated accuracy before we need to replay the task.

The replay memory is filled with  $M/k$  samples from each selected task, where  $k$  is the number of tasks that need to be replayed according to their decrease in validation accuracy. We skip replaying any tasks if no tasks are selected for replay, i.e.,  $k = 0$ .

**Grid search for  $\tau$  in Single CL Environments.** We performed a coarse-to-fine grid search for the parameter  $\tau$  on each dataset with Heur-GD to compare against the MCTS replay schedules. The best value for  $\tau$  is selected according to the highest mean accuracy on the validation set averaged over 5 seeds. The validation set consists of 15% of the training data and is the same for MCTS. We use the same experimental settings as described in Appendix D. The memory sizes are set to  $M = 10$  and  $M = 100$  for the 5-task datasets and the 10/20-task datasets respectively, and we apply uniform sampling as the memory selection method. We provide the ranges for  $\tau$  that was used on each dataset and put the best value in **bold**:

- Split MNIST:  $\tau = \{0.9, 0.93, 0.95, \mathbf{0.96}, 0.97, 0.98, 0.99\}$
- Split FashionMNIST:  $\tau = \{0.9, 0.93, 0.95, 0.96, \mathbf{0.97}, 0.98, 0.99\}$
- Split notMNIST:  $\tau = \{0.9, 0.93, 0.95, 0.96, 0.97, \mathbf{0.98}, 0.99\}$
- Permuted MNIST:  $\tau = \{0.5, 0.55, 0.6, 0.65, 0.7, \mathbf{0.75}, 0.8, 0.9, 0.95, 0.97, 0.99\}$
- Split CIFAR-100:  $\tau = \{0.3, 0.4, 0.45, \mathbf{0.5}, 0.55, 0.6, 0.65, 0.7, 0.8, 0.9, 0.95, 0.97, 0.99\}$
- Split miniImagenet:  $\tau = \{0.5, 0.6, 0.65, 0.7, \mathbf{0.75}, 0.8, 0.85, 0.9, 0.95, 0.97, 0.99\}$

Note that we use these values for  $\tau$  on all experiments with Heur-GD for the corresponding datasets. The performance for the heuristics highly depends on careful tuning for the ratio  $\tau$  when the memory size or memory selection method changes, as can be seen in Figure 4 and Table 6. We also provide the ranges for  $\tau$  that was used on each dataset in the Class Incremental Learning setting and put the best value in **bold**:

- Split MNIST:  $\tau = \{0.2, 0.3, 0.5, \mathbf{0.75}, 0.9\}$
- Split FashionMNIST:  $\tau = \{0.2, 0.3, \mathbf{0.5}, 0.75, 0.9\}$
- Split notMNIST:  $\tau = \{0.2, 0.3, \mathbf{0.5}, 0.75, 0.9\}$
- Split CIFAR-100:  $\tau = \{\mathbf{0.01}, 0.025, 0.05, 0.1, 0.25, 0.5\}$
- Split miniImagenet:  $\tau = \{\mathbf{0.01}, 0.025, 0.05, 0.1, 0.25, 0.5\}$

Table 5: The threshold parameter  $\tau$  used in the heuristic scheduling baselines Heuristic Global Drop (Heur-GD), Heuristic Local Drop (Heur-LD), and Heuristic Accuracy Threshold (Heur-AT). The search range is  $\tau \in \{0.90, 0.95, 0.999\}$  for all methods and we display the number of environments used for selecting the parameter used at test time.

Method	New Task Order						New Dataset			
	S-MNIST		S-FashionMNIST		S-CIFAR-10		S-notMNIST		S-FashionMNIST	
	$\tau$	#Envs	$\tau$	#Envs	$\tau$	#Envs	$\tau$	#Envs	$\tau$	#Envs
Heur-GD	0.9	10	0.95	20	0.9	10	0.9	10	0.9	10
Heur-LD	0.9	10	0.999	20	0.999	10	0.95	10	0.999	10
Heur-AT	0.9	10	0.999	20	0.9	10	0.9	10	0.95	10

**Grid search for  $\tau$  in Multiple CL Environments.** We performed a grid search for the parameter  $\tau$  for the three heuristic scheduling baselines for each experiment to compare against the learned replay scheduling policies. We select the parameter based on ACC scores achieved in the same number of training environments used by either DQN or A2C. The search range we use is  $\tau \in \{0.90, 0.95, 0.999\}$ . In Table 5, we show the selected parameter value of  $\tau$  and the number of environments used for selecting the value for each method and experiment in Section 4.2. The same parameters are used to generate the results on the heuristics in Table 4.

## D ADDITIONAL EXPERIMENTAL SETTINGS AND RESULTS FOR REPLAY SCHEDULING USING MCTS

This section is structured as follows:

- Appendix D.1: Full details on the experimental settings.
- Appendix D.2: Performance progress of MCTS as sanity check.
- Appendix D.3: Visualization of replay schedule from MCTS on Split MNIST.
- Appendix D.4: Additional results on Memory Selection Methods experiment.
- Appendix D.5: Additional results on Applying Replay Scheduling to Recent Replay Methods experiment.
- Appendix D.6: Additional results on Efficiency of Replay Scheduling experiment.
- Appendix D.7: Additional results on Varying Memory Size experiment.

The additional results appendices provide Welch’s t-tests for statistical significance between MCTS and the baselines.

### D.1 EXPERIMENTAL SETTINGS FOR MCTS IN SINGLE CL ENVIRONMENTS

Here, we provide details on the experimental settings for the experiments with MCTS in single CL environments.

**Datasets.** We conduct experiments on six datasets commonly used in the CL literature. Split MNIST (Zenke et al., 2017) is a variant of the MNIST (LeCun et al., 1998) dataset where the classes have been divided into 5 tasks incoming in the order 0/1, 2/3, 4/5, 6/7, and 8/9. Split FashionMNIST (Xiao et al., 2017) is of similar size to MNIST and consists of grayscale images of different clothes, where the classes have been divided into the 5 tasks T-shirt/Trouser, Pullover/Dress, Coat/Sandals, Shirt/Sneaker, and Bag/Ankle boots. Similar to MNIST, Split notMNIST (Bulatov, 2011) consists of 10 classes of the letters A-J with various fonts, where the classes are divided into the 5 tasks A/B, C/D, E/F, G/H, and I/J. We use training/test split provided by Ebrahimi et al. (2020) for Split notMNIST. Permuted MNIST (Goodfellow et al., 2013) dataset consists of applying a unique random permutation of the pixels of the images in original MNIST to create each task, except for the first task that is to learn the original MNIST dataset. We reduce the original MNIST dataset to 10k samples and create 9 unique random permutations to get a 10-task version of Permuted MNIST. In Split CIFAR-100 (Krizhevsky & Hinton, 2009), the 100 classes are divided into

20 tasks with 5 classes for each task (Lopez-Paz & Ranzato, 2017; Rebuffi et al., 2017). Similarly, Split miniImagenet (Vinyals et al., 2016) consists of 100 classes randomly chosen from the original Imagenet dataset where the 100 classes are divided into 20 tasks with 5 classes per task.

**CL Network Architectures.** We use a 2-layer MLP with 256 hidden units and ReLU activation for Split MNIST, Split FashionMNIST, Split notMNIST, and Permuted MNIST. We use a multi-head output layer for each dataset except Permuted MNIST where the network uses single-head output layer. For Split CIFAR-100, we use a multi-head CNN architecture built according to the CNN in Adel et al. (2019); Schwarz et al. (2018); Vinyals et al. (2016), which consists of four 3x3 convolutional blocks, i.e. convolutional layer followed by batch normalization (Ioffe & Szegedy, 2015), with 64 filters, ReLU activations, and 2x2 Max-pooling. For Split miniImagenet, we use the reduced ResNet-18 from Lopez-Paz & Ranzato (2017) with multi-head output layer.

**CL Hyperparameters.** We train all networks with the Adam optimizer (Kingma & Ba, 2014) with learning rate  $\eta = 0.001$  and hyperparameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Note that the learning rate for Adam is not reset before training on a new task. Next, we give details on number of training epochs and batch sizes specific for each dataset:

- Split MNIST: 10 epochs/task, batch size 128.
- Split FashionMNIST: 30 epochs/task, batch size 128.
- Split notMNIST: 50 epochs/task, batch size 128.
- Permuted MNIST: 20 epochs/task, batch size 128.
- Split CIFAR-100: 25 epochs/task, batch size 256.
- Split miniImagenet: 1 epoch/task (task 1 trained for 5 epochs as warm up), batch size 32.

**Monte Carlo Tree Search.** We run RS-MCTS for 100 iterations in all experiments. The replay schedules used in the reported results on the held-out test sets are from the replay schedule that gave the highest reward on the validation sets. The exploration constant for UCT in Equation 1 is set to  $C = 0.1$  in all experiments (Chaudhry & Lee, 2018).

**Computational Cost.** All experiments were performed on one NVIDIA GeForce RTX 2080Ti on an internal GPU cluster. The wall clock time for ETS on Split MNIST was around 1.5 minutes, and RS-MCTS and BFS takes 40 seconds on average to run one iteration, where BFS runs 1050 iterations in total for Split MNIST.

**Implementations.** We adapted the implementation released by Borsos et al. (2020) for the memory selection strategies Uniform sampling,  $k$ -means clustering,  $k$ -center clustering (Nguyen et al., 2017), and Mean-of-Features (Rebuffi et al., 2017). For HAL (Chaudhry et al., 2021), MER (Riemer et al., 2018), DER (Buzzega et al., 2020), and DER++, we follow the implementations released by Buzzega et al. (2020) for each method to apply them to our replay scheduling methods. Furthermore, we follow the implementations released by Chaudhry et al. (2019) and Mirzadeh & Ghasemzadeh (2021) for A-GEM (Chaudhry et al., 2018b) and ER-Ring (Chaudhry et al., 2019). For MCTS, we adapted the implementation from <https://github.com/int8/monte-carlo-tree-search> to search for replay schedules.

**Experimental Settings for Single Task Replay Memory Experiment.** We motivated the need for replay scheduling in CL with Figure 1 in Section 1. This simple experiment was performed on Split MNIST where the replay memory only contains samples from the first task, i.e., learning the classes 0/1. Furthermore, the memory can only be replayed at one point in time and we show the performance on each task when the memory is replayed at different time steps. We set the memory size to  $M = 10$  samples such that the memory holds 5 samples from both classes. We use the same network architecture and hyperparameters as described above for Split MNIST. The ACC metric above each subfigure corresponds to the ACC for training a network with the single task memory replay at different tasks. We observe that choosing different time points to replay the same memory leads to noticeably different results in the final performance, and in this example, the best final performance is achieved when the memory is used when learning task 5. Therefore, we argue that finding the proper schedule of what tasks to replay at what time in the fixed memory situation can be critical for CL.

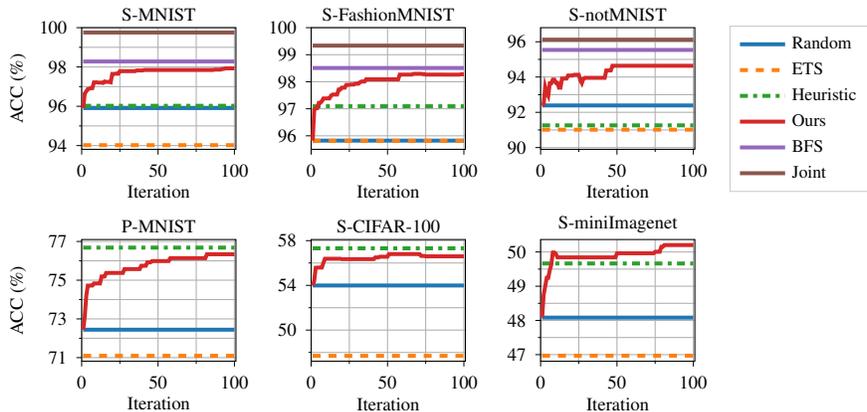


Figure 6: Average test accuracies over tasks after learning the final task (ACC) over the MCTS simulations for all datasets, where ‘S’ and ‘P’ are used as short for ‘Split’ and ‘Permuted’. We compare performance for MCTS (Ours) against random replay schedules (Random), Equal Task Schedule (ETS), and Heuristic Global Drop (Heuristic) baselines. For the first three datasets, we show the best ACC found from a breadth-first search (BFS) as well as the ACC achieved by training on all seen datasets jointly at every task (Joint). All results have been averaged over 5 seeds. These results show that replay scheduling can improve over ETS and outperform or perform on par with Heuristic across different datasets and network architectures.

## D.2 PERFORMANCE PROGRESS OF MCTS

In the first experiments, we show that the replay schedules from MCTS yield better performance than replaying an equal amount of samples per task. The replay memory size is fixed to  $M = 10$  for Split MNIST, FashionMNIST, and notMNIST, and  $M = 100$  for Permuted MNIST, Split CIFAR-100, and Split miniImagenet. Uniform sampling is used as the memory selection method for all methods in this experiment. For the 5-task datasets, we provide the optimal replay schedule found from a breadth-first search (BFS) over all 1050 possible replay schedules in our action space (which corresponds to a tree with depth of 4) as an upper bound for MCTS. As the search space grows fast with the number of tasks, BFS becomes computationally infeasible when we have 10 or more tasks.

Figure 6 shows the progress of ACC over iterations by MCTS for all datasets. We also show the best ACC metrics for Random, ETS, Heuristic, and BFS (where appropriate) as straight lines. Furthermore, we include the ACC achieved by training on all seen datasets jointly at every task (Joint) for the 5-task datasets. We observe that MCTS outperforms ETS successively with more iterations. Furthermore, MCTS approaches the upper limit of BFS on the 5-task datasets. For Permuted MNIST and Split CIFAR-100, the Heuristic baseline and MCTS perform on par after 50 iterations. This shows that Heuristic with careful tuning of the validation accuracy threshold can be a strong baseline when comparing replay scheduling methods. The top row of Table 1 shows the ACC for each method for this experiment. We note that MCTS outperforms ETS significantly on most datasets and performs on par with Heuristic.

## D.3 REPLAY SCHEDULE VISUALIZATION FOR SPLIT MNIST

In Figure 7, we show the progress in test classification performance for each task when using ETS and MCTS with memory size  $M = 10$  on Split MNIST. For comparison, we also show the performance from a network that is fine-tuning on the current task without using replay. Both ETS and MCTS overcome catastrophic forgetting to a large degree compared to the fine-tuning network. Our method MCTS further improves the performance compared to ETS with the same memory, which indicates that learning the time to learn can be more efficient against catastrophic forgetting. Especially, Task 1 and 2 seems to be the most difficult task to remember since it has the lowest final performance using the fine-tuning network. Both ETS and MCTS manage to retain their performance on Task 1 using replay, however, MCTS remembers Task 2 better than ETS by around 5%.

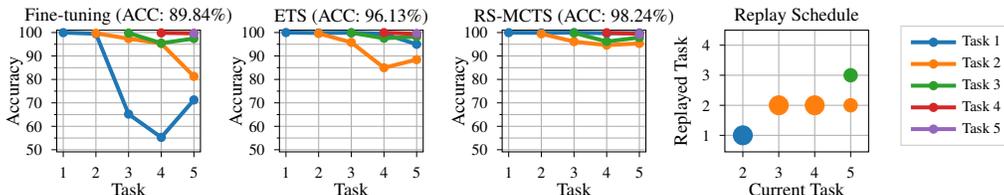


Figure 7: Comparison of test classification accuracies for Task 1-5 on Split MNIST from a network trained without replay (Fine-tuning), ETS, and MCTS. The ACC metric for each method is shown on top of each figure. We also visualize the replay schedule found by MCTS as a bubble plot to the right. The memory size is set to  $M = 10$  with uniform memory selection for ETS and MCTS. Results are shown for 1 seed.

To bring more insights to this behavior, we have visualized the task proportions of the replay examples using a bubble plot showing the corresponding replay schedule from MCTS in Figure 7(right). At Task 3 and 4, we see that the schedule fills the memory with data from Task 2 and discards replaying Task 1. This helps the network to retain knowledge about Task 2 better than ETS at the cost of forgetting Task 3 slightly when learning Task 4. This shows that the learned policy has considered the difficulty level of different tasks. At the next task, the MCTS schedule has decided to rehearse Task 3 and reduces replaying Task 2 when learning Task 5. This behavior is similar to spaced repetition, where increasing the time interval between rehearsals helps memory retention. We emphasize that even on datasets with few tasks, using learned replay schedules can overcome catastrophic forgetting better than standard ETS approaches.

#### D.4 ALTERNATIVE MEMORY SELECTION METHODS

We show that our method can be combined with any memory selection method for storing replay samples. In addition to uniform sampling, we apply various memory selection methods commonly used in the CL literature, namely  $k$ -means clustering,  $k$ -center clustering (Nguyen et al., 2017), and Mean-of-Features (MoF) (Rebuffi et al., 2017). The replay memory sizes are  $M = 10$  for the 5-task datasets and  $M = 100$  for the 10- and 20-task datasets. Table 6 shows the results across all datasets, and present the statistical significance tests between MCTS and the baselines in Table 7. We note that our method in general achieves significantly higher ACC comparing to the baselines showing that learning the time to learn is important.

#### D.5 APPLYING SCHEDULING TO RECENT REPLAY METHODS

In Section 4.1, we showed that MCTS can be applied to any replay method. We combined MCTS together with four recent replay methods, namely Hindsight Anchor Learning (HAL) (Chaudhry et al., 2021), Meta Experience Replay (MER) (Riemer et al., 2018), and Dark Experience Replay (DER) (Buzzega et al., 2020). We present the hyperparameters used for each method in Table 8. The hyperparameters for each method are denoted as

- **HAL.**  $\eta$ : learning rate,  $\lambda$ : regularization,  $\gamma$ : mean embedding strength,  $\beta$ : decay rate,  $k$ : gradient steps on anchors
- **MER.**  $\gamma$ : across batch meta-learning rate,  $\beta$ : within batch meta-learning rate
- **DER.**  $\alpha$ : loss coefficient for memory logits
- **DER++.**  $\alpha$ : loss coefficient for memory logits,  $\beta$ : loss coefficient for memory labels

For the experiments, we used the same architectures and hyperparameters as described in Appendix D.1 for all datasets if not mentioned otherwise. We used the Adam optimizer with learning rate  $\eta = 0.001$  for MER, DER, and DER++. For HAL, we used the SGD optimizer since using Adam made the model diverge in our experiments. Table 9 shows the ACC and BWT for all methods combined with the scheduling from Random, ETS, Heuristic, and MCTS. We observe that MCTS can further improve the performance for each of the replay methods across the different datasets. Table 10 shows statistical significance tests between MCTS and the baselines for every considered replay method, where we note that our method in general achieves significantly higher ACC comparing to the baselines.

Table 6: Performance comparison with ACC and BWT metrics for all datasets between MCTS (Ours) and the baselines with various memory selection methods. We provide the metrics for training on all seen task datasets jointly (Joint) as an upper bound. Furthermore, we include the results from a breadth-first search (BFS) with Uniform memory selection for the 5-task datasets. The memory size is set to  $M = 10$  and  $M = 100$  for the 5-task and 10/20-task datasets respectively. We report the means and standard deviations averaged over 5 seeds.

		Split MNIST		Split FashionMNIST		Split notMNIST	
Memory	Schedule	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Offline	Joint	99.75 ± 0.06	0.01 ± 0.06	99.34 ± 0.08	-0.01 ± 0.14	96.12 ± 0.57	-0.21 ± 0.71
Uniform	BFS	98.28 ± 0.49	-1.84 ± 0.63	98.51 ± 0.23	-1.03 ± 0.28	95.54 ± 0.67	-1.04 ± 0.87
Uniform	Random	94.91 ± 2.52	-6.13 ± 3.16	95.89 ± 2.03	-4.33 ± 2.55	91.84 ± 1.48	-5.37 ± 2.12
	ETS	94.02 ± 4.25	-7.22 ± 5.33	95.81 ± 3.53	-4.45 ± 4.34	91.01 ± 1.39	-6.16 ± 1.82
	Heur-GD	96.02 ± 2.32	-4.64 ± 2.90	97.09 ± 0.62	-2.82 ± 0.84	91.26 ± 3.99	-6.06 ± 4.70
	MCTS	97.93 ± 0.56	-2.27 ± 0.71	98.27 ± 0.17	-1.29 ± 0.20	94.64 ± 0.39	-1.47 ± 0.79
k-means	Random	92.65 ± 1.38	-8.96 ± 1.74	93.11 ± 2.75	-7.76 ± 3.42	93.11 ± 1.01	-3.78 ± 1.43
	ETS	92.89 ± 3.53	-8.66 ± 4.42	96.47 ± 0.85	-3.55 ± 1.07	93.80 ± 0.82	-2.84 ± 0.81
	Heur-GD	96.28 ± 1.68	-4.32 ± 2.11	95.78 ± 1.50	-4.46 ± 1.87	91.75 ± 0.94	-5.60 ± 2.07
	MCTS	98.20 ± 0.16	-1.94 ± 0.22	98.48 ± 0.26	-1.04 ± 0.31	93.61 ± 0.71	-3.11 ± 0.55
k-center	Random	95.48 ± 0.82	-5.40 ± 1.05	93.24 ± 2.84	-7.64 ± 3.51	91.70 ± 1.94	-5.33 ± 2.80
	ETS	94.84 ± 1.40	-6.20 ± 1.77	97.28 ± 0.50	-2.58 ± 0.66	91.08 ± 2.48	-6.39 ± 3.46
	Heur-GD	94.55 ± 2.79	-6.47 ± 3.50	94.08 ± 3.72	-6.59 ± 4.57	92.06 ± 1.20	-4.70 ± 2.09
	MCTS	98.24 ± 0.36	-1.93 ± 0.44	98.06 ± 0.35	-1.59 ± 0.45	94.26 ± 0.37	-1.97 ± 1.02
MoF	Random	96.96 ± 1.34	-3.57 ± 1.69	96.39 ± 1.69	-3.66 ± 2.17	93.09 ± 1.40	-3.70 ± 1.76
	ETS	97.04 ± 1.23	-3.46 ± 1.50	96.48 ± 1.33	-3.55 ± 1.73	92.64 ± 0.87	-4.57 ± 1.59
	Heur-GD	96.46 ± 2.41	-4.09 ± 3.01	95.84 ± 0.89	-4.39 ± 1.15	93.24 ± 0.77	-3.48 ± 1.37
	MCTS	98.37 ± 0.24	-1.70 ± 0.28	97.84 ± 0.32	-1.81 ± 0.39	94.62 ± 0.42	-1.80 ± 0.56
		Permuted MNIST		Split CIFAR-100		Split miniImagenet	
Memory	Schedule	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Offline	Joint	95.34 ± 0.13	0.17 ± 0.18	84.73 ± 0.81	-1.06 ± 0.81	74.03 ± 0.83	9.70 ± 0.68
Uniform	Random	72.59 ± 1.52	-25.71 ± 1.76	53.76 ± 1.80	-35.11 ± 1.93	49.89 ± 1.03	-14.79 ± 1.14
	ETS	71.09 ± 2.31	-27.39 ± 2.59	47.70 ± 2.16	-41.69 ± 2.37	46.97 ± 1.24	-18.32 ± 1.34
	Heur-GD	76.68 ± 2.13	-20.82 ± 2.41	57.31 ± 1.21	-30.76 ± 1.45	49.66 ± 1.10	-12.04 ± 0.59
	MCTS	76.34 ± 0.98	-21.21 ± 1.16	56.60 ± 1.13	-31.39 ± 1.11	50.20 ± 0.72	-13.46 ± 1.22
k-means	Random	71.91 ± 1.24	-26.45 ± 1.34	53.20 ± 1.44	-35.77 ± 1.31	49.96 ± 1.46	-14.81 ± 1.18
	ETS	69.40 ± 1.32	-29.23 ± 1.47	47.51 ± 1.14	-41.77 ± 1.30	45.82 ± 0.92	-19.53 ± 1.10
	Heur-GD	75.57 ± 1.18	-22.11 ± 1.22	54.31 ± 3.94	-33.80 ± 4.24	49.25 ± 1.00	-12.92 ± 1.22
	MCTS	77.74 ± 0.80	-19.66 ± 0.95	56.95 ± 0.92	-30.92 ± 0.83	50.47 ± 0.85	-13.31 ± 1.24
k-center	Random	71.39 ± 1.87	-27.04 ± 2.05	48.29 ± 2.11	-40.88 ± 2.28	44.40 ± 1.35	-20.03 ± 1.31
	ETS	69.11 ± 1.69	-29.58 ± 1.81	44.13 ± 1.06	-45.28 ± 1.04	41.35 ± 1.23	-23.71 ± 1.45
	Heur-GD	74.33 ± 2.00	-23.45 ± 2.27	50.32 ± 1.97	-37.99 ± 2.14	44.13 ± 0.95	-18.26 ± 1.05
	MCTS	76.55 ± 1.16	-21.06 ± 1.32	51.37 ± 1.63	-37.01 ± 1.62	46.76 ± 0.96	-16.56 ± 0.90
MoF	Random	78.80 ± 1.07	-18.79 ± 1.16	62.35 ± 1.24	-26.33 ± 1.25	56.02 ± 1.11	-7.99 ± 1.13
	ETS	77.62 ± 1.12	-20.10 ± 1.26	60.43 ± 1.17	-28.22 ± 1.26	56.12 ± 1.12	-8.93 ± 0.83
	Heur-GD	77.27 ± 1.45	-20.15 ± 1.63	55.60 ± 2.70	-32.57 ± 2.77	52.30 ± 0.59	-9.61 ± 0.67
	MCTS	81.58 ± 0.75	-15.41 ± 0.86	64.22 ± 0.65	-23.48 ± 1.02	57.70 ± 0.51	-5.31 ± 0.55

Table 7: Two-tailed Welch’s  $t$ -test results for the various memory selection methods presented in Table 6.

Memory	Schedule	Split MNIST		Split FashionMNIST		Split notMNIST	
		$t$	$p$	$t$	$p$	$t$	$p$
Uniform	MCTS vs Random	2.34	0.074	2.34	0.078	3.66	<b>0.017</b>
	MCTS vs ETS	1.82	0.140	1.39	0.236	5.04	<b>0.005</b>
	MCTS vs Heur-GD	1.60	0.178	3.64	<b>0.017</b>	1.69	0.166
k-means	MCTS vs Random	7.97	<b>0.001</b>	3.90	<b>0.017</b>	0.81	0.445
	MCTS vs ETS	3.00	<b>0.040</b>	4.54	<b>0.007</b>	-0.36	0.732
	MCTS vs Heur-GD	2.27	0.085	3.55	<b>0.022</b>	3.15	<b>0.015</b>
k-center	MCTS vs Random	6.15	<b>0.001</b>	3.37	<b>0.027</b>	2.59	0.057
	MCTS vs ETS	4.71	<b>0.007</b>	2.56	<b>0.037</b>	2.53	0.062
	MCTS vs Heur-GD	2.62	0.057	2.13	0.099	3.51	<b>0.019</b>
MoF	MCTS vs Random	2.07	0.103	1.68	0.163	2.10	0.093
	MCTS vs ETS	2.13	0.095	1.99	0.110	4.11	<b>0.007</b>
	MCTS vs Heur-GD	1.58	0.188	4.21	<b>0.008</b>	3.15	<b>0.019</b>

Memory	Schedule	Permuted MNIST		Split CIFAR-100		Split miniImagenet	
		$t$	$p$	$t$	$p$	$t$	$p$
Uniform	MCTS vs Random	4.14	<b>0.005</b>	2.67	<b>0.033</b>	0.49	0.636
	MCTS vs ETS	4.18	<b>0.007</b>	7.30	<b>0.000</b>	4.52	<b>0.003</b>
	MCTS vs Heur-GD	-0.29	0.780	-0.87	0.412	0.82	0.441
k-means	MCTS vs Random	7.91	<b>0.000</b>	4.39	<b>0.003</b>	0.61	0.565
	MCTS vs ETS	10.83	<b>0.000</b>	12.93	<b>0.000</b>	7.42	<b>0.000</b>
	MCTS vs Heur-GD	3.05	<b>0.019</b>	1.31	0.255	1.87	0.099
k-center	MCTS vs Random	4.70	<b>0.003</b>	2.32	0.051	2.85	<b>0.024</b>
	MCTS vs ETS	7.25	<b>0.000</b>	7.46	<b>0.000</b>	6.94	<b>0.000</b>
	MCTS vs Heur-GD	1.92	0.100	0.82	0.437	3.89	<b>0.005</b>
MoF	MCTS vs Random	4.26	<b>0.004</b>	2.67	<b>0.037</b>	2.75	<b>0.036</b>
	MCTS vs ETS	5.86	<b>0.001</b>	5.69	<b>0.001</b>	2.57	<b>0.045</b>
	MCTS vs Heur-GD	5.26	<b>0.002</b>	6.22	<b>0.002</b>	13.90	<b>0.000</b>

Table 8: Hyperparameters for replay-based methods HAL, MER, DER and DER++ used in experiments on applying MCTS to recent replay-based methods in Section 4.1.

Method	Hyperparam.	5-task Datasets			10- and 20-task Datasets		
		S-MNIST	S-FashionMNIST	S-notMNIST	P-MNIST	S-CIFAR-100	S-miniImagenet
HAL	$\eta$	0.1	0.1	0.1	0.1	0.03	0.03
	$\lambda$	0.1	0.1	0.1	0.1	1.0	0.03
	$\gamma$	0.5	0.1	0.1	0.1	0.1	0.1
	$\beta$	0.7	0.5	0.5	0.5	0.5	0.5
	$k$	100	100	100	100	100	100
MER	$\gamma$	1.0	1.0	1.0	1.0	1.0	1.0
	$\beta$	1.0	0.01	1.0	1.0	0.1	0.1
DER	$\alpha$	0.2	0.2	0.1	1.0	1.0	0.1
DER++	$\alpha$	0.2	0.2	0.1	1.0	1.0	0.1
	$\beta$	1.0	1.0	1.0	1.0	1.0	1.0

Table 9: Performance comparison with ACC and BWT metrics between scheduling methods MCTS (Ours), Random, ETS, and Heuristic when combining them with replay-based methods Hindsight Anchor Learning (HAL), Meta Experience Replay (MER), Dark Experience Replay (DER), and DER++. Replay memory sizes are  $M = 10$  and  $M = 100$  for the 5-task and 10/20-task datasets respectively. We report the mean and standard deviation averaged over 5 seeds. Results on Heuristic where some seed did not converge is denoted by \*. Applying MCTS to each method can enhance the performance compared to using the baseline schedules.

		Split MNIST		Split FashionMNIST		Split notMNIST	
Method	Schedule	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
HAL	Random	96.32 ± 1.77	-3.90 ± 2.28	90.42 ± 4.26	-10.75 ± 5.45	93.50 ± 1.10	-3.14 ± 1.56
	ETS	97.21 ± 1.25	-2.80 ± 1.59	96.75 ± 0.50	-2.84 ± 0.75	92.16 ± 1.82	-5.04 ± 2.24
	Heur-GD	97.69 ± 0.19	-2.22 ± 0.24	*74.16 ± 11.19	*-31.26 ± 14.00	93.64 ± 0.93	-2.80 ± 1.20
	MCTS	97.96 ± 0.15	-1.85 ± 0.18	97.56 ± 0.51	-2.02 ± 0.63	94.47 ± 0.82	-1.67 ± 0.64
MER	Random	93.00 ± 3.22	-7.96 ± 4.15	96.20 ± 2.10	-2.31 ± 2.59	89.10 ± 2.57	-8.82 ± 3.26
	ETS	92.97 ± 1.73	-8.52 ± 2.15	84.88 ± 3.85	-3.34 ± 5.59	90.56 ± 0.83	-6.11 ± 1.06
	Heur-GD	94.30 ± 2.79	-6.46 ± 3.50	96.91 ± 0.62	-1.34 ± 0.76	90.90 ± 1.30	-6.24 ± 1.96
	MCTS	96.44 ± 0.72	-4.14 ± 0.94	86.67 ± 4.09	0.85 ± 3.85	92.44 ± 0.77	-3.63 ± 1.06
DER	Random	95.91 ± 2.18	-4.40 ± 2.46	50.00 ± 0.00	-12.20 ± 0.07	78.76 ± 12.73	-11.91 ± 4.45
	ETS	98.17 ± 0.35	-2.00 ± 0.42	97.69 ± 0.58	-2.05 ± 0.71	94.74 ± 1.05	-1.94 ± 1.17
	Heur-GD	94.57 ± 1.71	-6.08 ± 2.09	*72.49 ± 19.32	*-20.88 ± 11.46	*77.88 ± 12.58	*-12.66 ± 4.17
	MCTS	99.02 ± 0.10	-0.91 ± 0.13	98.33 ± 0.51	-1.26 ± 0.63	95.02 ± 0.33	-0.97 ± 0.81
DER++	Random	90.09 ± 10.02	-11.73 ± 12.38	*50.00 ± 0.00	*-12.20 ± 0.07	61.83 ± 9.84	-14.40 ± 10.67
	ETS	97.98 ± 0.52	-2.24 ± 0.66	98.12 ± 0.40	-1.59 ± 0.52	94.53 ± 1.02	-1.82 ± 1.02
	Heur-GD	92.35 ± 2.42	-8.83 ± 2.99	*67.31 ± 21.20	*-24.86 ± 16.34	93.88 ± 1.33	-2.86 ± 1.49
	MCTS	98.84 ± 0.21	-1.14 ± 0.26	98.38 ± 0.43	-1.17 ± 0.51	94.73 ± 0.20	-1.21 ± 1.12
		Permuted MNIST		Split CIFAR-100		Split miniImagenet	
Method	Schedule	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
HAL	Random	88.93 ± 0.53	-6.77 ± 0.64	35.90 ± 2.47	-17.37 ± 3.76	40.86 ± 1.86	-5.12 ± 2.23
	ETS	88.46 ± 0.86	-7.26 ± 0.90	34.90 ± 2.02	-18.92 ± 0.91	38.13 ± 1.18	-8.19 ± 1.73
	Heur-GD	*66.63 ± 28.50	*-29.68 ± 27.90	35.07 ± 1.29	-24.76 ± 2.41	39.51 ± 1.49	-5.65 ± 0.77
	MCTS	89.14 ± 0.74	-6.29 ± 0.74	40.22 ± 1.57	-12.77 ± 1.30	41.39 ± 1.15	-3.69 ± 1.86
MER	Random	87.25 ± 0.47	-8.77 ± 0.59	42.68 ± 0.86	-35.56 ± 1.39	32.86 ± 0.95	-7.71 ± 0.45
	ETS	73.01 ± 0.96	-25.19 ± 1.10	43.38 ± 1.81	-34.84 ± 1.98	33.58 ± 1.53	-6.80 ± 1.46
	Heur-GD	83.86 ± 3.19	-12.48 ± 3.60	40.90 ± 1.70	-44.10 ± 2.03	34.22 ± 1.93	-7.57 ± 1.63
	MCTS	79.72 ± 0.71	-17.42 ± 0.78	44.29 ± 0.69	-32.73 ± 0.88	32.74 ± 1.29	-5.77 ± 1.04
DER	Random	90.67 ± 0.31	-5.20 ± 0.30	56.17 ± 1.30	-29.03 ± 1.38	35.13 ± 4.11	-10.85 ± 2.92
	ETS	85.71 ± 0.75	-11.15 ± 0.87	52.58 ± 1.49	-32.93 ± 2.04	35.50 ± 2.84	-10.94 ± 2.21
	Heur-GD	81.56 ± 2.28	-15.06 ± 2.51	55.75 ± 1.08	-31.27 ± 1.02	43.62 ± 0.88	-8.18 ± 1.16
	MCTS	90.11 ± 0.18	-5.89 ± 0.23	58.99 ± 0.98	-24.95 ± 0.64	43.46 ± 0.95	-9.32 ± 1.37
DER++	Random	89.83 ± 0.92	-6.03 ± 0.98	60.90 ± 0.89	-23.45 ± 1.34	46.78 ± 1.96	3.28 ± 1.35
	ETS	85.25 ± 0.88	-11.60 ± 1.03	52.54 ± 1.06	-33.22 ± 1.51	41.36 ± 2.90	-4.07 ± 2.28
	Heur-GD	79.17 ± 2.44	-17.68 ± 2.68	56.70 ± 1.27	-30.33 ± 1.41	45.73 ± 0.84	-6.09 ± 1.24
	MCTS	89.84 ± 0.22	-6.13 ± 0.29	59.23 ± 0.83	-24.61 ± 0.91	49.45 ± 0.68	-3.12 ± 0.89

Table 10: Two-tailed Welch’s  $t$ -test results for the alternative replay methods results presented in Table 9.

Memory	Schedule	Split MNIST		Split FashionMNIST		Split notMNIST	
		$t$	$p$	$t$	$p$	$t$	$p$
HAL	MCTS vs Random	1.84	0.139	3.33	<b>0.028</b>	1.41	0.198
	MCTS vs ETS	1.20	0.295	2.27	0.053	2.31	0.064
	MCTS vs Heur-GD	2.26	0.056	4.18	<b>0.014</b>	1.34	0.218
MER	MCTS vs Random	2.08	0.099	-4.15	<b>0.006</b>	2.48	0.059
	MCTS vs ETS	3.71	<b>0.012</b>	0.64	0.542	3.32	<b>0.011</b>
	MCTS vs Heur-GD	1.48	0.204	-4.96	<b>0.007</b>	2.04	0.084
DER	MCTS vs Random	2.85	<b>0.046</b>	190.21	<b>0.000</b>	2.56	0.063
	MCTS vs ETS	4.64	<b>0.007</b>	1.64	0.139	0.51	0.635
	MCTS vs Heur-GD	5.21	<b>0.006</b>	2.67	0.055	2.73	0.053
DER++	MCTS vs Random	1.75	0.156	227.62	<b>0.000</b>	6.68	<b>0.003</b>
	MCTS vs ETS	3.09	<b>0.026</b>	0.89	0.397	0.39	0.712
	MCTS vs Heur-GD	5.36	<b>0.006</b>	2.93	<b>0.043</b>	1.26	0.272

Memory	Schedule	Permuted MNIST		Split CIFAR-100		Split miniImagenet	
		$t$	$p$	$t$	$p$	$t$	$p$
HAL	MCTS vs Random	0.46	0.657	2.96	0.022	0.49	0.640
	MCTS vs ETS	1.20	0.266	4.16	<b>0.004</b>	3.97	<b>0.004</b>
	MCTS vs Heur-GD	1.58	0.189	5.07	<b>0.001</b>	2.01	0.082
MER	MCTS vs Random	-17.61	<b>0.000</b>	2.92	<b>0.020</b>	-0.14	0.889
	MCTS vs ETS	11.24	<b>0.000</b>	0.95	0.386	-0.84	0.425
	MCTS vs Heur-GD	-2.54	0.059	3.70	<b>0.013</b>	-1.27	0.244
DER	MCTS vs Random	-3.12	<b>0.019</b>	3.46	<b>0.010</b>	3.96	<b>0.014</b>
	MCTS vs ETS	11.36	<b>0.000</b>	7.18	<b>0.000</b>	5.31	<b>0.003</b>
	MCTS vs Heur-GD	7.47	<b>0.002</b>	4.44	<b>0.002</b>	-0.25	0.809
DER++	MCTS vs Random	0.03	0.981	-2.75	<b>0.025</b>	2.57	0.051
	MCTS vs ETS	10.17	<b>0.000</b>	9.94	<b>0.000</b>	5.43	<b>0.004</b>
	MCTS vs Heur-GD	8.72	<b>0.001</b>	3.34	<b>0.013</b>	6.89	<b>0.000</b>

## D.6 EFFICIENCY OF REPLAY SCHEDULING

We illustrate the efficiency of replay scheduling in a setting where only 1 sample/class is available from the historical data for replay. Table 11 shows that MCTS, despite using significantly fewer samples for replay, performs mostly on par with the baselines and outperforms them on Permuted MNIST. Table 12 shows statistical significance tests between MCTS and the baselines for the corresponding results. Our method mostly performs on par with the baselines, but is significantly better than the baselines on Permuted MNIST.

We visualize the memory usage in the experiment on efficiency of replay scheduling in Section 4.1. For the 5-task datasets, the replay memory size for MCTS is set to  $M = 2$ , such that only 2 samples can be selected for replay at all times. Similarly, we set  $M = 50$  for the 10- and 20-task datasets which have 100 classes to learn in total. The baselines A-GEM (Chaudhry et al., 2018b), ER-Ring (Chaudhry et al., 2019), and Uniform use an incremental memory in order to replay 1 sample/class at all tasks. We visualize the memory usage for our method and the baselines for the 5-task datasets in Figure 8. Here, the memory capacity is reached at task 2, while the baselines must increment their memory size. Figure 9 shows the memory usage for Permuted MNIST and the 20-task datasets Split CIFAR-100 and Split miniImagenet.

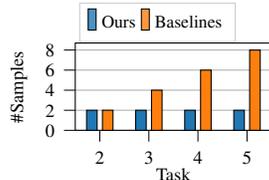


Figure 8: Number of replayed samples per task for the 5-task datasets in the tiny memory setting.

Table 11: Performance comparison with ACC and BWT metrics for all datasets between MCTS and the baselines in the setting where only 1 sample per class can be replayed. The memory sizes are set to  $M = 10$  and  $M = 100$  for the 5-task and 10/20-task datasets respectively. MCTS (Ours) and Random uses  $M = 2$  and  $M = 50$  for the 5-task and 10/20-task datasets respectively. We report the means and standard deviations averaged over 5 seeds. MCTS performs on par with the best baselines for both metrics on all datasets, except on Permuted MNIST where MCTS outperforms the baselines.

Method	Split MNIST		Split FashionMNIST		Split notMNIST	
	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Random	92.56 +/- 2.90	-8.97 +/- 3.62	92.70 +/- 3.78	-8.24 +/- 4.75	89.53 +/- 3.96	-8.13 +/- 5.02
A-GEM	94.97 +/- 1.50	-6.03 +/- 1.87	94.81 +/- 0.86	-5.65 +/- 1.06	92.27 +/- 1.16	-4.17 +/- 1.39
ER-Ring	94.94 +/- 1.56	-6.07 +/- 1.92	95.83 +/- 2.15	-4.38 +/- 2.59	91.10 +/- 1.89	-6.27 +/- 2.35
ER-Uniform	95.77 +/- 1.12	-5.02 +/- 1.39	97.12 +/- 1.57	-2.79 +/- 1.98	92.14 +/- 1.45	-4.90 +/- 1.41
MCTS	96.07 +/- 1.60	-4.59 +/- 2.01	97.17 +/- 0.78	-2.64 +/- 0.99	93.41 +/- 1.11	-3.36 +/- 1.56

Method	Permuted MNIST		Split CIFAR-100		Split minilmagenet	
	ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
Random	70.02 +/- 1.76	-28.22 +/- 1.92	48.62 +/- 1.02	-39.95 +/- 1.10	48.85 +/- 1.38	-14.55 +/- 1.86
A-GEM	64.71 +/- 1.78	-34.41 +/- 2.05	42.22 +/- 2.13	-46.90 +/- 2.21	32.06 +/- 1.83	-30.81 +/- 1.79
ER-Ring	69.73 +/- 1.13	-28.87 +/- 1.29	53.93 +/- 1.13	-34.91 +/- 1.18	49.82 +/- 1.69	-14.38 +/- 1.57
ER-Uniform	69.85 +/- 1.01	-28.74 +/- 1.17	52.63 +/- 1.62	-36.43 +/- 1.81	50.56 +/- 1.07	-13.52 +/- 1.34
MCTS	72.52 +/- 0.54	-25.43 +/- 0.65	51.50 +/- 1.19	-37.01 +/- 1.08	50.70 +/- 0.54	-12.60 +/- 1.13

Table 12: Two-tailed Welch’s  $t$ -test results for efficiency of replay scheduling presented in Table 11.

Methods	Split MNIST		Split FashionMNIST		Split notMNIST	
	$t$	$p$	$t$	$p$	$t$	$p$
MCTS vs Random	2.12	0.077	2.32	0.076	1.89	0.122
MCTS vs A-GEM	1.00	0.347	4.08	<b>0.004</b>	1.41	0.195
MCTS vs ER-Ring	1.01	0.342	1.18	0.292	2.11	0.076
MCTS vs Uniform	0.30	0.770	0.07	0.950	1.39	0.203

Methods	Permuted MNIST		Split CIFAR-100		Split minilmagenet	
	$t$	$p$	$t$	$p$	$t$	$p$
MCTS vs Random	2.72	<b>0.044</b>	3.67	<b>0.007</b>	2.48	0.054
MCTS vs A-GEM	8.40	<b>0.001</b>	7.60	<b>0.000</b>	19.52	<b>0.000</b>
MCTS vs ER-Ring	4.46	<b>0.005</b>	-2.96	<b>0.018</b>	0.99	0.369
MCTS vs Uniform	4.68	<b>0.003</b>	-1.13	<b>0.296</b>	0.23	0.824

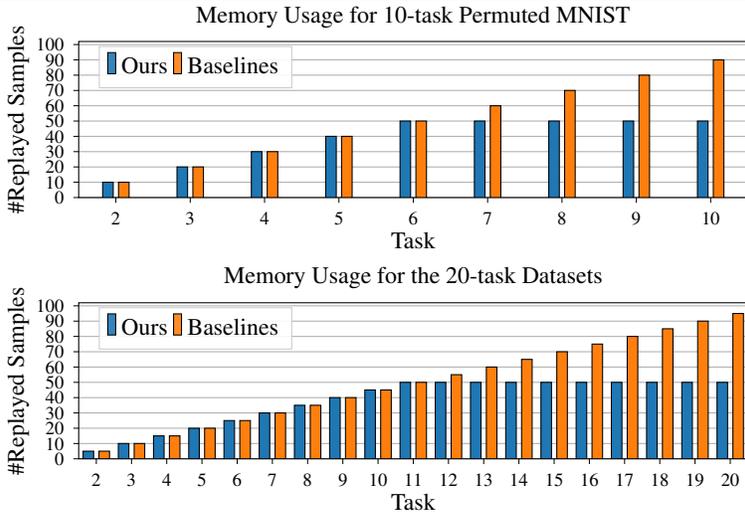


Figure 9: Number of replayed samples per task for 10-task Permuted MNIST (top) and the 20-task datasets in the experiment in Section 4.1. The fixed memory size  $M = 50$  for our method is reached after learning task 6 and task 11 on the Permuted MNIST and the 20-task datasets respectively, while the baselines continue incrementing their number of replay samples per task.

## D.7 VARYING MEMORY SIZE IN DIFFERENT CONTINUAL LEARNING SETTING

Here, we provide the ACC and BWT metrics from the varying memory size experiments from figure 4 in Section 4.1. We also provide the p-values from Welch’s t-tests to show the statistical significance between the methods.

- **Domain-IL, 10-task Permuted MNIST:** Metrics in Table 13, t-tests in Table 14.
- **Task-IL, 5-task Split MNIST/FashionMNIST/notMNIST:** Metrics in Table 15, t-tests in Table 16.
- **Task-IL, 20-task Split CIFAR-100/miniImagenet:** Metrics in Table 17, t-tests in Table 18.
- **Class-IL, 5-task Split MNIST/FashionMNIST/notMNIST:** Metrics in Table 19, t-tests in Table 20.
- **Class-IL, 20-task Split CIFAR-100/miniImagenet:** Metrics in Table 21, t-tests in Table 22.

We note that MCTS mostly performs significantly better than the baselines on the 10-task Permuted MNIST and, interestingly, on the 20-task Split CIFAR-100 and Split miniImagenet in both Task-IL and Clas-IL settings, which shows the importance of replay scheduling for CL datasets with long task horizons.

Table 13: Performance comparison in the Domain Incremental Learning setting over various memory sizes for the methods on Permuted MNIST.

Memory Size	Method	Permuted MNIST	
		ACC (%)	BWT (%)
M=90	Random	72.63 ± 1.06	-25.62 ± 1.20
	ETS	71.49 ± 1.15	-26.92 ± 1.26
	Heur-GD	75.50 ± 1.64	-22.14 ± 1.93
	MCTS	71.66 ± 0.67	-28.70 ± 0.81
M=270	Random	82.01 ± 0.79	-15.24 ± 0.86
	ETS	80.68 ± 0.66	-16.72 ± 0.77
	Heur-GD	78.32 ± 1.58	-19.01 ± 1.83
	MCTS	82.08 ± 0.35	-17.18 ± 0.38
M=450	Random	84.72 ± 1.39	-12.16 ± 1.53
	ETS	84.38 ± 1.12	-12.54 ± 1.20
	Heur-GD	81.66 ± 2.30	-15.27 ± 2.48
	MCTS	85.53 ± 0.42	-13.34 ± 0.49
M=900	Random	88.02 ± 1.03	-8.51 ± 1.09
	ETS	88.02 ± 0.46	-8.52 ± 0.47
	Heur-GD	80.27 ± 3.26	-16.77 ± 3.76
	MCTS	88.94 ± 0.43	-9.55 ± 0.47
M=2250	Random	90.94 ± 0.46	-5.26 ± 0.43
	ETS	91.07 ± 0.23	-5.11 ± 0.21
	Heur-GD	81.28 ± 3.79	-15.68 ± 4.32
	MCTS	92.45 ± 0.34	-5.63 ± 0.41

Table 14: Two-tailed Welch’s t-test results for the varying memory size experiments presented in Table 15 in the Domain Incremental Learning setting on Permuted MNIST.

Memory size	Methods	Permuted MNIST	
		t	p
M=90	MCTS vs Random	-1.55	0.166
	MCTS vs ETS	0.26	0.806
	MCTS vs Heur-GD	-4.32	<b>0.007</b>
M=270	MCTS vs Random	0.17	0.872
	MCTS vs ETS	3.75	<b>0.009</b>
	MCTS vs Heur-GD	4.66	<b>0.008</b>
M=450	MCTS vs Random	1.13	0.313
	MCTS vs ETS	1.93	0.111
	MCTS vs Heur-GD	3.31	<b>0.027</b>
M=900	MCTS vs Random	1.66	0.153
	MCTS vs ETS	2.92	<b>0.019</b>
	MCTS vs Heur-GD	5.28	<b>0.006</b>
M=2250	MCTS vs Random	5.31	<b>0.001</b>
	MCTS vs ETS	6.74	<b>0.000</b>
	MCTS vs Heur-GD	5.88	<b>0.004</b>

Table 15: Performance comparison in the Task Incremental Learning setting over various memory sizes for the methods on the 5-task datasets.

Memory Size	Method	Split MNIST		Split FashionMNIST		Split notMNIST	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
M=8	Random	95.50 ± 1.92	-5.38 ± 2.38	95.04 ± 2.43	-5.42 ± 3.06	91.36 ± 1.76	-6.15 ± 1.98
	ETS	95.83 ± 1.23	-4.96 ± 1.53	97.25 ± 0.68	-2.66 ± 0.82	92.10 ± 1.53	-5.10 ± 1.97
	Heur-GD	96.16 ± 0.83	-4.45 ± 1.05	96.29 ± 0.79	-3.83 ± 1.05	93.68 ± 1.11	-3.06 ± 1.54
	MCTS	98.17 ± 0.48	-1.99 ± 0.61	98.33 ± 0.11	-1.27 ± 0.11	94.53 ± 0.62	-1.67 ± 1.03
M=24	Random	97.21 ± 1.70	-3.24 ± 2.14	97.20 ± 0.71	-2.72 ± 0.92	92.31 ± 1.50	-4.88 ± 1.75
	ETS	97.87 ± 0.58	-2.41 ± 0.74	97.84 ± 0.55	-1.91 ± 0.64	93.19 ± 1.41	-3.84 ± 1.48
	Heur-GD	97.82 ± 1.17	-2.40 ± 1.46	95.76 ± 2.84	-4.48 ± 3.52	92.97 ± 2.55	-3.54 ± 3.07
	MCTS	98.58 ± 0.32	-1.47 ± 0.41	98.48 ± 0.24	-1.06 ± 0.32	95.23 ± 0.66	-1.08 ± 1.72
M=80	Random	97.57 ± 2.03	-2.78 ± 2.53	96.85 ± 1.58	-3.10 ± 2.05	94.14 ± 0.95	-2.50 ± 1.59
	ETS	98.47 ± 0.40	-1.65 ± 0.49	97.93 ± 0.85	-1.76 ± 0.98	94.63 ± 0.67	-1.82 ± 0.74
	Heur-GD	98.36 ± 0.40	-1.72 ± 0.50	96.22 ± 1.96	-3.92 ± 2.45	93.24 ± 2.56	-3.73 ± 2.71
	MCTS	99.06 ± 0.16	-0.89 ± 0.21	98.60 ± 0.24	-0.87 ± 0.34	94.84 ± 0.58	-0.97 ± 1.28
M=120	Random	97.82 ± 2.29	-2.47 ± 2.88	98.29 ± 0.44	-1.26 ± 0.56	94.91 ± 1.03	-1.65 ± 1.04
	ETS	98.82 ± 0.18	-1.22 ± 0.26	98.53 ± 0.28	-0.96 ± 0.36	95.32 ± 0.66	-1.21 ± 1.30
	Heur-GD	98.37 ± 0.38	-1.71 ± 0.48	95.26 ± 3.14	-5.12 ± 3.98	93.42 ± 1.78	-3.47 ± 2.22
	MCTS	99.05 ± 0.11	-0.88 ± 0.15	98.75 ± 0.19	-0.77 ± 0.25	94.16 ± 1.08	-1.99 ± 1.69
M=200	Random	97.99 ± 1.59	-2.25 ± 2.00	96.68 ± 3.33	-3.38 ± 4.19	93.94 ± 1.40	-2.12 ± 1.70
	ETS	98.83 ± 0.23	-1.19 ± 0.28	98.60 ± 0.25	-0.99 ± 0.29	94.79 ± 0.50	-1.58 ± 1.03
	Heur-GD	98.15 ± 0.64	-1.97 ± 0.81	95.83 ± 2.00	-4.40 ± 2.52	93.88 ± 1.63	-2.71 ± 1.69
	MCTS	99.09 ± 0.08	-0.83 ± 0.11	98.83 ± 0.11	-0.65 ± 0.15	95.19 ± 0.53	-0.49 ± 0.47
M=400	Random	97.98 ± 2.23	-2.28 ± 2.80	98.00 ± 1.59	-1.74 ± 2.00	94.68 ± 1.11	-1.77 ± 1.09
	ETS	99.18 ± 0.10	-0.78 ± 0.13	98.83 ± 0.09	-0.71 ± 0.12	95.41 ± 0.56	-0.20 ± 1.07
	Heur-GD	98.44 ± 0.60	-1.64 ± 0.77	95.41 ± 3.77	-4.93 ± 4.69	92.84 ± 1.88	-3.95 ± 2.05
	MCTS	99.25 ± 0.05	-0.63 ± 0.09	98.80 ± 0.11	-0.62 ± 0.17	95.25 ± 0.44	-0.97 ± 1.20
M=800	Random	98.60 ± 1.42	-1.51 ± 1.77	97.00 ± 3.93	-2.98 ± 4.94	94.97 ± 0.73	-1.86 ± 0.77
	ETS	99.34 ± 0.06	-0.57 ± 0.04	98.97 ± 0.06	-0.51 ± 0.10	95.52 ± 0.74	-0.61 ± 1.21
	Heur-GD	98.76 ± 0.41	-1.23 ± 0.55	93.32 ± 3.67	-7.54 ± 4.65	95.06 ± 1.44	-1.17 ± 1.30
	MCTS	99.38 ± 0.08	-0.45 ± 0.10	98.96 ± 0.12	-0.45 ± 0.18	95.46 ± 0.70	-0.48 ± 0.81

Table 16: Two-tailed Welch’s  $t$ -test results for the varying memory size experiments presented in Table 15 in the Task Incremental Learning setting on the 5-task datasets.

Memory size	Methods	Split MNIST		Split FashionMNIST		Split notMNIST	
		$t$	$p$	$t$	$p$	$t$	$p$
M=8	MCTS vs Random	2.70	<b>0.048</b>	2.70	0.054	3.41	<b>0.019</b>
	MCTS vs ETS	3.53	<b>0.016</b>	3.13	<b>0.033</b>	2.94	<b>0.030</b>
	MCTS vs Heur-GD	4.18	<b>0.005</b>	5.12	<b>0.006</b>	1.34	0.227
M=24	MCTS vs Random	1.59	0.183	3.39	<b>0.020</b>	3.55	<b>0.014</b>
	MCTS vs ETS	2.15	0.074	2.14	<b>0.080</b>	2.63	<b>0.041</b>
	MCTS vs Heur-GD	1.26	0.268	1.90	0.128	1.72	0.153
M=80	MCTS vs Random	1.47	0.215	2.18	0.091	1.26	0.251
	MCTS vs ETS	2.76	<b>0.038</b>	1.53	0.191	0.48	0.645
	MCTS vs Heur-GD	3.26	<b>0.021</b>	2.41	0.072	1.22	0.285
M=120	MCTS vs Random	1.07	0.344	1.95	0.105	-1.01	0.343
	MCTS vs ETS	2.12	0.073	1.31	0.231	-1.83	0.112
	MCTS vs Heur-GD	3.45	<b>0.020</b>	2.22	0.090	0.70	0.507
M=200	MCTS vs Random	1.38	0.240	1.29	0.266	1.67	0.154
	MCTS vs ETS	2.06	0.094	1.69	0.146	1.11	0.301
	MCTS vs Heur-GD	2.91	<b>0.042</b>	3.00	<b>0.040</b>	1.53	0.190
M=400	MCTS vs Random	1.14	0.318	1.00	0.373	0.94	0.389
	MCTS vs ETS	1.25	0.258	-0.42	0.684	-0.47	0.655
	MCTS vs Heur-GD	2.69	0.054	1.80	0.146	2.49	0.061
M=800	MCTS vs Random	1.10	0.332	1.00	0.373	0.96	0.364
	MCTS vs ETS	0.82	0.435	-0.06	0.955	-0.11	0.912
	MCTS vs Heur-GD	3.05	<b>0.035</b>	3.08	<b>0.037</b>	0.50	0.636

Table 17: Performance comparison in the Task Incremental Learning setting over various memory sizes for the methods on the 20-task datasets.

Memory Size	Method	Split CIFAR-100		Split miniImagenet	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)
M=95	Random	56.92 ± 0.72	-31.98 ± 0.69	52.06 ± 0.94	-12.49 ± 1.41
	ETS	55.19 ± 0.77	-33.70 ± 0.82	51.33 ± 2.01	-13.32 ± 2.15
	Heur-GD	53.86 ± 3.89	-34.51 ± 4.25	50.66 ± 1.36	-10.99 ± 1.39
	MCTS	60.46 ± 1.05	-27.67 ± 1.23	53.59 ± 0.24	-9.49 ± 0.51
M=285	Random	67.79 ± 1.19	-20.31 ± 1.12	57.85 ± 1.09	-6.33 ± 1.30
	ETS	65.02 ± 0.98	-23.36 ± 1.10	56.18 ± 0.83	-8.43 ± 1.37
	Heur-GD	60.26 ± 2.88	-27.44 ± 2.87	53.33 ± 2.21	-7.96 ± 2.18
	MCTS	68.85 ± 0.56	-18.41 ± 0.63	57.91 ± 0.09	-5.18 ± 0.54
M=475	Random	70.86 ± 1.24	-17.19 ± 1.26	59.54 ± 1.25	-4.22 ± 1.96
	ETS	69.40 ± 0.73	-18.68 ± 0.82	59.60 ± 0.98	-4.75 ± 1.31
	Heur-GD	65.00 ± 3.08	-22.58 ± 3.20	50.12 ± 4.60	-11.79 ± 5.08
	MCTS	72.93 ± 0.54	-14.15 ± 0.78	60.00 ± 0.48	-2.85 ± 0.34
M=950	Random	75.39 ± 0.46	-12.24 ± 0.46	61.87 ± 0.85	-2.32 ± 1.06
	ETS	74.24 ± 0.61	-13.49 ± 0.44	60.51 ± 0.95	-4.01 ± 1.28
	Heur-GD	65.97 ± 5.51	-21.40 ± 5.58	50.99 ± 4.64	-11.21 ± 4.62
	MCTS	76.65 ± 0.62	-10.31 ± 0.84	61.82 ± 0.69	-1.38 ± 0.55
M=1900	Random	78.86 ± 0.38	-8.64 ± 0.37	62.95 ± 0.69	-1.56 ± 0.99
	ETS	77.66 ± 0.38	-9.72 ± 0.36	62.38 ± 0.68	-1.55 ± 0.90
	Heur-GD	63.70 ± 6.12	-24.07 ± 6.60	56.42 ± 3.15	-5.43 ± 3.63
	MCTS	79.24 ± 0.59	-7.43 ± 0.64	63.83 ± 1.04	0.41 ± 1.10

Table 18: Two-tailed Welch’s  $t$ -test results for the varying memory size experiments presented in Table 17 in the Task Incremental Learning setting on the 20-task datasets.

Memory size	Methods	Split CIFAR-100		Split miniImagenet	
		$t$	$p$	$t$	$p$
M=95	MCTS vs Random	5.56	<b>0.001</b>	3.15	<b>0.029</b>
	MCTS vs ETS	8.11	<b>0.000</b>	2.23	0.088
	MCTS vs Heur-GD	3.27	<b>0.025</b>	4.24	<b>0.012</b>
M=285	MCTS vs Random	1.62	0.159	0.11	0.918
	MCTS vs ETS	6.76	<b>0.000</b>	4.14	<b>0.014</b>
	MCTS vs Heur-GD	5.87	<b>0.003</b>	4.13	<b>0.014</b>
M=475	MCTS vs Random	3.05	<b>0.025</b>	0.68	0.525
	MCTS vs ETS	7.78	<b>0.000</b>	0.73	0.496
	MCTS vs Heur-GD	5.07	<b>0.006</b>	4.28	<b>0.012</b>
M=950	MCTS vs Random	3.25	<b>0.013</b>	-0.09	0.929
	MCTS vs ETS	5.51	<b>0.001</b>	2.25	0.058
	MCTS vs Heur-GD	3.85	<b>0.017</b>	4.62	<b>0.009</b>
M=1900	MCTS vs Random	1.08	0.317	1.41	0.203
	MCTS vs ETS	4.53	<b>0.003</b>	2.34	0.052
	MCTS vs Heur-GD	5.05	<b>0.007</b>	4.47	<b>0.007</b>

Table 19: Performance comparison in the Class Incremental Learning setting over various memory sizes for the methods on the 5-task datasets.

Memory Size	Method	Split MNIST		Split FashionMNIST		Split notMNIST	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)	ACC (%)	BWT (%)
M=10	Random	25.53 ± 1.57	-92.74 ± 1.99	30.78 ± 2.74	-85.70 ± 3.47	30.48 ± 5.16	-81.80 ± 6.18
	ETS	25.10 ± 1.04	-93.27 ± 1.32	30.67 ± 2.36	-85.83 ± 3.03	31.67 ± 6.03	-80.22 ± 7.36
	Heur-GD	25.05 ± 1.51	-93.27 ± 1.85	31.81 ± 3.17	-84.33 ± 3.96	29.86 ± 3.49	-82.07 ± 4.26
M=20	MCTS	28.19 ± 2.18	-89.34 ± 2.71	33.77 ± 2.71	-81.91 ± 3.45	37.23 ± 3.53	-73.93 ± 3.89
	Random	34.13 ± 2.04	-81.98 ± 2.57	37.80 ± 2.92	-76.83 ± 3.62	42.78 ± 1.24	-65.75 ± 1.92
	ETS	35.17 ± 2.06	-80.65 ± 2.54	38.76 ± 1.58	-75.63 ± 1.93	45.21 ± 5.92	-63.10 ± 7.33
M=40	Heur-GD	33.87 ± 2.15	-82.20 ± 2.65	42.52 ± 1.98	-70.94 ± 2.47	45.13 ± 3.68	-63.32 ± 4.31
	MCTS	39.62 ± 0.73	-75.05 ± 0.92	44.27 ± 2.19	-68.71 ± 2.74	47.83 ± 0.82	-59.91 ± 1.11
	Random	45.34 ± 3.71	-67.96 ± 4.64	45.78 ± 2.96	-66.84 ± 3.72	55.70 ± 2.96	-49.58 ± 3.36
M=100	ETS	48.57 ± 1.90	-63.93 ± 2.37	45.54 ± 1.21	-67.15 ± 1.48	58.07 ± 3.88	-46.85 ± 4.92
	Heur-GD	45.22 ± 4.54	-68.05 ± 5.65	51.67 ± 1.59	-59.43 ± 2.02	53.90 ± 6.10	-52.09 ± 7.73
	MCTS	50.79 ± 2.45	-61.07 ± 3.08	51.33 ± 2.96	-59.87 ± 3.70	62.55 ± 3.23	-41.20 ± 4.81
M=200	Random	59.53 ± 6.37	-50.12 ± 7.96	55.96 ± 3.31	-53.98 ± 4.17	65.06 ± 3.48	-37.86 ± 4.30
	ETS	68.09 ± 1.06	-39.49 ± 1.33	59.26 ± 0.95	-49.92 ± 1.14	72.55 ± 1.65	-28.88 ± 2.31
	Heur-GD	49.81 ± 1.58	-62.17 ± 1.95	58.14 ± 2.96	-51.24 ± 3.71	58.28 ± 2.02	-46.70 ± 2.63
M=200	MCTS	66.86 ± 2.21	-40.96 ± 2.75	60.97 ± 2.43	-47.63 ± 2.99	71.18 ± 3.58	-30.27 ± 4.15
	Random	65.82 ± 7.50	-42.27 ± 9.36	61.11 ± 5.19	-47.26 ± 6.48	67.07 ± 5.65	-35.36 ± 6.96
	ETS	77.73 ± 1.31	-27.42 ± 1.63	66.51 ± 0.75	-40.60 ± 0.99	77.15 ± 0.24	-22.61 ± 0.89
M=200	Heur-GD	53.85 ± 0.88	-57.05 ± 1.11	54.69 ± 0.32	-55.42 ± 0.43	57.05 ± 3.28	-47.66 ± 4.09
	MCTS	78.14 ± 1.67	-26.81 ± 2.00	67.38 ± 2.60	-39.41 ± 3.26	75.91 ± 4.88	-24.15 ± 5.98

Table 20: Two-tailed Welch’s *t*-test results for the varying memory size experiments presented in Table 19 in the Class Incremental Learning setting on the 5-task datasets.

Memory size	Methods	Split MNIST		Split FashionMNIST		Split notMNIST	
		<i>t</i>	<i>p</i>	<i>t</i>	<i>p</i>	<i>t</i>	<i>p</i>
M=10	MCTS vs Random	1.98	0.086	1.56	0.158	2.16	0.067
	MCTS vs ETS	2.56	<b>0.045</b>	1.73	0.123	1.59	0.159
	MCTS vs Heur-GD	2.37	<b>0.049</b>	0.94	0.374	2.97	<b>0.018</b>
M=20	MCTS vs Random	5.06	<b>0.004</b>	3.54	<b>0.009</b>	6.80	<b>0.000</b>
	MCTS vs ETS	4.07	<b>0.010</b>	4.07	<b>0.004</b>	0.88	0.429
	MCTS vs Heur-GD	5.08	<b>0.004</b>	1.18	0.271	1.43	0.219
M=40	MCTS vs Random	2.45	<b>0.044</b>	2.65	<b>0.029</b>	3.13	<b>0.014</b>
	MCTS vs ETS	1.43	0.192	3.62	<b>0.014</b>	1.78	0.115
	MCTS vs Heur-GD	2.16	0.073	-0.21	0.843	2.51	<b>0.046</b>
M=100	MCTS vs Random	2.18	0.082	2.44	<b>0.043</b>	2.45	<b>0.040</b>
	MCTS vs ETS	-1.00	0.356	1.31	0.246	-0.69	0.515
	MCTS vs Heur-GD	12.55	<b>0.000</b>	1.48	0.180	6.29	<b>0.001</b>
M=200	MCTS vs Random	3.21	<b>0.029</b>	2.16	0.075	2.37	<b>0.046</b>
	MCTS vs ETS	0.39	0.707	0.65	0.548	-0.51	0.639
	MCTS vs Heur-GD	25.79	<b>0.000</b>	9.69	<b>0.001</b>	6.41	<b>0.000</b>

Table 21: Performance comparison in the Class Incremental Learning setting over various memory sizes for the methods on the 20-task datasets.

Memory Size	Method	Split CIFAR-100		Split miniImagenet	
		ACC (%)	BWT (%)	ACC (%)	BWT (%)
M=100	Random	4.49 ± 0.04	-85.73 ± 0.34	3.36 ± 0.23	-60.11 ± 0.83
	ETS	4.51 ± 0.13	-85.62 ± 0.34	3.34 ± 0.21	-60.25 ± 0.47
	Heur-GD	4.55 ± 0.13	-84.62 ± 0.44	3.29 ± 0.14	-57.01 ± 0.45
	MCTS	4.62 ± 0.10	-84.73 ± 0.23	3.51 ± 0.19	-57.57 ± 1.13
M=200	Random	5.24 ± 0.12	-84.53 ± 0.28	3.60 ± 0.24	-58.83 ± 0.81
	ETS	4.98 ± 0.14	-85.05 ± 0.19	3.75 ± 0.14	-58.92 ± 1.03
	Heur-GD	5.20 ± 0.16	-83.33 ± 0.29	3.95 ± 0.36	-54.81 ± 0.58
	MCTS	5.84 ± 0.16	-82.89 ± 0.28	4.26 ± 0.19	-55.19 ± 0.73
M=400	Random	7.13 ± 0.25	-81.93 ± 0.29	4.81 ± 0.40	-56.36 ± 0.44
	ETS	6.35 ± 0.25	-82.92 ± 0.23	4.35 ± 0.15	-57.28 ± 0.83
	Heur-GD	7.66 ± 0.60	-79.82 ± 0.60	6.97 ± 0.78	-49.44 ± 1.04
	MCTS	8.31 ± 0.21	-79.58 ± 0.51	7.64 ± 0.70	-48.73 ± 1.71
M=800	Random	10.64 ± 0.51	-77.27 ± 0.62	8.40 ± 1.07	-50.95 ± 1.85
	ETS	9.13 ± 0.20	-79.54 ± 0.43	7.15 ± 0.51	-53.95 ± 0.74
	Heur-GD	11.33 ± 0.47	-74.75 ± 0.73	8.44 ± 1.15	-48.89 ± 1.43
	MCTS	12.06 ± 0.24	-74.70 ± 0.29	11.56 ± 1.30	-44.30 ± 1.63
M=1600	Random	15.54 ± 0.69	-70.42 ± 0.55	13.72 ± 1.30	-44.19 ± 1.73
	ETS	13.99 ± 0.32	-72.42 ± 0.24	12.17 ± 1.40	-47.55 ± 1.39
	Heur-GD	15.66 ± 1.84	-67.24 ± 1.75	8.39 ± 1.35	-49.35 ± 1.77
	MCTS	17.59 ± 0.42	-66.59 ± 0.66	15.40 ± 0.34	-42.27 ± 0.73

Table 22: Two-tailed Welch’s  $t$ -test results for the varying memory size experiments presented in Table 21 in the Class Incremental Learning setting on the 20-task datasets.

Memory size	Methods	Split CIFAR-100		Split miniImagenet	
		$t$	$p$	$t$	$p$
M=100	MCTS vs Random	2.35	0.065	1.03	0.335
	MCTS vs ETS	1.38	0.208	1.20	0.266
	MCTS vs Heur-GD	0.87	0.410	1.86	0.103
M=200	MCTS vs Random	5.95	<b>0.000</b>	4.34	<b>0.003</b>
	MCTS vs ETS	8.07	<b>0.000</b>	4.35	<b>0.003</b>
	MCTS vs Heur-GD	5.72	<b>0.000</b>	1.51	0.182
M=400	MCTS vs Random	7.15	<b>0.000</b>	7.06	<b>0.000</b>
	MCTS vs ETS	11.98	<b>0.000</b>	9.22	<b>0.000</b>
	MCTS vs Heur-GD	2.01	0.101	1.29	0.232
M=800	MCTS vs Random	5.05	<b>0.003</b>	3.75	<b>0.006</b>
	MCTS vs ETS	18.79	<b>0.000</b>	6.32	<b>0.001</b>
	MCTS vs Heur-GD	2.77	<b>0.033</b>	3.60	<b>0.007</b>
M=1600	MCTS vs Random	5.08	<b>0.002</b>	2.50	0.059
	MCTS vs ETS	13.60	<b>0.000</b>	4.49	<b>0.008</b>
	MCTS vs Heur-GD	2.04	0.104	10.04	<b>0.000</b>

## E ADDITIONAL EXPERIMENTAL SETTINGS AND RESULTS FOR REPLAY SCHEDULING POLICY EXPERIMENTS

This section is structured as follows:

- Appendix E.1: Full details on the experimental settings.
- Appendix E.2: Details of the ranking method to assess the generalization abilities of the scheduling policies.
- Appendix E.3: Additional experimental results for the Replay Scheduling Policy Generalization experiments with Welch’s t-tests for statistical significance between the RL algorithms (DQN and A2C) and the baselines.
- Appendix E.4: Task splits in the continual learning environments used for testing.

### E.1 EXPERIMENTAL SETTINGS FOR RL-BASED FRAMEWORK

Here, we provide details on the experimental settings for the experiments with our RL-based framework where we use multiple CL environments for learning replay scheduling policies that generalize.

**Datasets.** We conduct experiments on CL environments with four datasets common CL benchmarks, namely, Split MNIST (Zenke et al., 2017), Split Fashion-MNIST (Xiao et al., 2017), Split notMNIST (Bulatov, 2011), and Split CIFAR-10 (Krizhevsky & Hinton, 2009). All datasets consists of 5 tasks with 2 classes/task.

**CL Network Architectures.** We use a 2-layer MLP with 256 hidden units and ReLU activation for Split MNIST, Split FashionMNIST, and Split notMNIST. For Split CIFAR-10, we use the same ConvNet architecture as used for Split CIFAR-100 in Appendix D.1. We use a multi-head output layer for each dataset and assume task labels are available at test time for selecting the correct output head related to the task.

**CL Hyperparameters.** We train all networks with the Adam optimizer (Kingma & Ba, 2014) with learning rate  $\eta = 0.001$  and hyperparameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . Note that the learning rate for Adam is not reset before training on a new task. Next, we give details on number of training epochs and batch sizes specific for each dataset:

- Split MNIST: 10 epochs/task, batch size 128.
- Split FashionMNIST: 10 epochs/task, batch size 128.
- Split notMNIST: 20 epochs/task, batch size 128.
- Split CIFAR-10: 20 epochs/task, batch size 256.

**Generating CL Environments.** We generate multiple CL environments with pre-set random seeds for initializing the network parameters  $\phi$  and shuffling the task order. The pre-set random seeds are in the range 0 – 49, such that we have 50 environments for each dataset. We shuffle the task order by permuting the class order and then split the classes into 5 pairs (tasks) with 2 classes/pair. For environments with seed 0, we keep the original task order in the dataset. Taking a step at task  $t$  in the CL environments involves training the CL network on the  $t$ -th dataset with a replay memory  $\mathcal{M}_t$  from the discrete action space described in Section 3.2. Therefore, to speed up the experiments with the RL algorithms, we run a breadth-first search (BFS) through the discrete action space and save the classification results for re-use during policy learning. Note that the action space has 1050 possible paths of replay schedules for the datasets with  $T = 5$  tasks, which makes the environment generation time-consuming. Hence, we only generate environments where the replay memory size  $M = 10$  have been used, and leave analysis of different memory sizes as future work. The CL environments we used are provided in the code submission.

**DQN and A2C Architectures.** The input layer has size  $T - 1$  where each unit is inputting the task performances since the states are represented by the validation accuracies  $s_t = [A_{t,1}^{(val)}, \dots, A_{t,t}^{(val)}, 0, \dots, 0]$ . The current task can therefore be determined by the number of non-zero state inputs. The output layer has 35 units representing the possible actions at  $T = 5$  with the discrete action space we have constructed in Section 3.2. We use action masking on the output units to prevent the network from selection invalid actions for constructing the replay memory at the current task. The DQN is a 2-layer MLP with 512 hidden units and ReLU activations. For A2C, we use

Table 23: DQN hyperparameters for the experiments on **New Task Orders** in Section 4.2.

Hyperparameters	Split MNIST	Split FashionMNIST	Split CIFAR-10
Training Environments	30	20	10
Learning Rate	0.0001	0.0003	0.0003
Optimizer	Adam	Adam	Adam
Buffer Size	10k	10k	10k
Target Update per step	500	500	500
Batch Size	32	32	32
Discount Factor $\gamma$	1.0	1.0	1.0
Exploration Start $\epsilon_{start}$	1.0	1.0	1.0
Exploration Final $\epsilon_{final}$	0.02	0.02	0.02
Exploration Annealing (episodes)	2.5k	2.5k	2.5k
Training Episodes	10k	10k	10k

Table 24: A2C hyperparameters for the experiments on **New Task Orders** in Section 4.2.

Hyperparameters	Split MNIST	Split FashionMNIST	Split CIFAR-10
Training Environments	10	10	10
Learning Rate	0.0001	0.0003	0.00003
Optimizer	RMSProp	RMSProp	RMSProp
Gradient Clipping	0.5	0.5	0.5
GAE parameter $\lambda$	0.95	0.95	0.95
VF coefficient	0.5	0.5	0.5
Entropy coefficient	0.01	0.01	0.01
Number of steps $n_{steps}$	5	5	5
Discount Factor $\gamma$	1.0	1.0	1.0
Training Episodes	100k	100k	100k

separate networks for parameterizing the policy and the value function, where both networks are 2-layer MLPs with 64 hidden units of Tanh activations.

**DQN and A2C Hyperparameters.** We provide the hyperparameters for the both DQN and A2C in Table 23-27. Table 23 and 24 includes the hyperparameters on the New Task Order experiment for DQN and A2C respectively, while Table 26 and 27 includes the hyperparameters on the New Dataset experiment for DQN and A2C respectively. Regarding the training environments in Table 26 and 27, we use two different datasets in the training environments to increase the diversity. When Split notMNIST is for testing, half the amount of training environments are using Split MNIST and the other half uses Split FashionMNIST. For example, in Table 27, A2C uses 10 training environments which means that there are 5 Split MNIST environments and 5 Split FashionMNIST environments. Similarly, half the amount of training environments are using Split MNIST and the other half uses Split notMNIST when the testing environments uses Split FashionMNIST.

**Computational Cost.** All experiments were performed on one NVIDIA GeForce RTX 2080Ti on an internal GPU cluster. Generating a CL environment for one seed with Split MNIST took on around 9.5 hours averaged over 10 runs of BFS. Similarly for Split CIFAR-10, generating one CL environment took on average 16.1 hours. Table 25 shows a time-cost ablation experiment w/ or w/o a DQN for selecting which tasks to replay in Split MNIST. We measured the wall clock time for training and evaluating the CL model on the 5 Split MNIST tasks w/ and w/o the DQN, and show the wall clock time averaged over 10 different DQN seeds. The time difference when w/ DQN is only 3.2 seconds, since selecting which tasks to replay is only a forward pass with the RL policy.

Table 25: Time-cost ablation experiment w/ or w/o a DQN for replay scheduling on Split MNIST.

Time Cost	With DQN	Without DQN	Difference
Avg. Time (in sec)	84.6	81.4	3.2

**Implementations.** The code for DQN was adapted from OpenAI baselines (Dhariwal et al., 2017) and the PyTorch (Paszke et al., 2019) tutorial on DQN [https://pytorch.org/tutorials/intermediate/reinforcement\\_q\\_learning.html](https://pytorch.org/tutorials/intermediate/reinforcement_q_learning.html). For A2C, we followed the implementations released by Kostrikov (2018) and Igl et al. (2020).

Table 26: DQN hyperparameters for the experiments on **New Dataset** in Section 4.2. Split notM-NIST and Split FashionMNIST indicate the dataset used in the test environments.

<b>Hyperparameters</b>	<b>Split notMNIST</b>	<b>Split FashionMNIST</b>
Training Environments	30	30
Learning Rate	0.0001	0.0001
Optimizer	Adam	Adam
Buffer Size	10k	10k
Target Update per step	500	500
Batch Size	32	32
Discount Factor $\gamma$	1.0	1.0
Exploration Start $\epsilon_{start}$	1.0	1.0
Exploration Final $\epsilon_{final}$	0.02	0.02
Exploration Annealing (episodes)	2.5k	2.5k
Training Episodes	10k	10k

Table 27: A2C hyperparameters for the experiments on **New Dataset** in Section 4.2. Split notM-NIST and Split FashionMNIST indicate the dataset used in the test environments.

<b>Hyperparameters</b>	<b>Split notMNIST</b>	<b>Split FashionMNIST</b>
Training Environments	10	10
Learning Rate	0.0001	0.0003
Optimizer	RMSProp	RMSProp
Gradient Clipping	0.5	0.5
GAE parameter $\lambda$	0.95	0.95
VF coefficient	0.5	0.5
Entropy coefficient	0.01	0.01
Number of steps $n_{steps}$	5	5
Discount Factor $\gamma$	1.0	1.0
Training Episodes	100k	100k

## E.2 ASSESSING GENERALIZATION WITH RANKING METHOD

We use a ranking method based on the CL performance in every test environment for performance comparison between the methods in Section 4.2. We use rankings because the performances can vary greatly between environments with different task orders and datasets. To measure the CL performance in the environments, we use the average test accuracy over all tasks after learning the final task, i.e.,

$$\text{ACC} = \frac{1}{T} \sum_{i=1}^T A_{T,i}^{(test)},$$

where  $A_{t,i}^{(test)}$  is the test accuracy of task  $i$  after learning task  $t$ . Each method are ranked in descending order based on the ACC achieved in an environment. For example, assume that we want to compare the CL performance from using learned replay scheduling policies with DQN and A2C against a Random scheduling policy in one environment. The CL performances achieved for each method are given by

$$[\text{ACC}_{\text{Random}}, \text{ACC}_{\text{DQN}}, \text{ACC}_{\text{A2C}}] = [90\%, 99\%, 95\%].$$

We get the following ranking order between the methods based on their corresponding ACC:

$$\text{ranking}([\text{ACC}_{\text{Random}}, \text{ACC}_{\text{DQN}}, \text{ACC}_{\text{A2C}}]) = [3, 1, 2],$$

where DQN is ranked in 1st place, A2C in 2nd, and Random in 3rd. When there are multiple environments for evaluation, we compute the average ranking across the ranking positions in every environment for each method to compare.

The average ranking for DQN and A2C are computed over the seed for initializing the network parameters as well as the seed of the environment. Similarly, the Random baseline is affected by the seed setting the random selection of actions and the environment seed. However, the performance of the ETS and Heuristic baselines are affected by the seed of the environment as these policies are fixed. We use copied values of the performance in environments for the ETS and Heuristic baselines when we need to compare across different random seeds for Random, DQN, and A2C. We show an example of such ranking calculation for ETS, a Heuristic baseline, DQN, and A2C. Consider the following performances for one environment:

$$\begin{bmatrix} \text{ACC}_{\text{ETS}}^1 & \text{ACC}_{\text{Heur}}^1 & \text{ACC}_{\text{DQN}}^1 & \text{ACC}_{\text{A2C}}^1 \\ \text{ACC}_{\text{ETS}}^2 & \text{ACC}_{\text{Heur}}^2 & \text{ACC}_{\text{DQN}}^2 & \text{ACC}_{\text{A2C}}^2 \end{bmatrix} = \begin{bmatrix} 90\% & 95\% & 95\% & 99\% \\ * & * & 97\% & 98\% \end{bmatrix},$$

where \* denotes a copy of the ACC value in the first row. The subscript on ACC denotes the method and the superscript the seed used for initializing the policy network  $\theta$ . Therefore, we copy the values for ETS and Heur such that the  $\text{ACC}_{\text{DQN}}^2$  for seed 2 can be compared against ETS and Heur. Note that there is a tie between  $\text{ACC}_{\text{Heur}}^1$  and  $\text{ACC}_{\text{DQN}}^1$  as they have ACC 95%. We handle ties by assigning tied methods the average of their ranks, such that the ranks for both seeds will be

$$\begin{aligned} & \text{ranking} \left( \begin{bmatrix} \text{ACC}_{\text{ETS}}^1 & \text{ACC}_{\text{Heur}}^1 & \text{ACC}_{\text{DQN}}^1 & \text{ACC}_{\text{A2C}}^1 \\ \text{ACC}_{\text{ETS}}^2 & \text{ACC}_{\text{Heur}}^2 & \text{ACC}_{\text{DQN}}^2 & \text{ACC}_{\text{A2C}}^2 \end{bmatrix}, \text{axis}=-1, \text{keepdim}=\text{True} \right) \\ & = \text{ranking} \left( \begin{bmatrix} 90\% & 95\% & 95\% & 99\% \\ 90\% & 95\% & 97\% & 98\% \end{bmatrix}, \text{axis}=-1, \text{keepdim}=\text{True} \right) \\ & = \begin{bmatrix} 4 & 2.5 & 2.5 & 1 \\ 4 & 3 & 2 & 1 \end{bmatrix}, \end{aligned}$$

where we inserted the copied values, such that  $\text{ACC}_{\text{ETS}}^1 = \text{ACC}_{\text{ETS}}^2 = 90\%$  and  $\text{ACC}_{\text{Heur}}^1 = \text{ACC}_{\text{Heur}}^2 = 95\%$ . The mean ranking across the seeds thus becomes

$$\text{mean} \left( \begin{bmatrix} 4 & 2.5 & 2.5 & 1 \\ 4 & 3 & 2 & 1 \end{bmatrix}, \text{axis}=0 \right) = [4 \quad 2.75 \quad 2.25 \quad 1]$$

where A2C comes in 1st place, DQN in 2nd, Heur. in 3rd, and ETS on 4th place. We average across seeds and environments to obtain the final ranking score for each method for comparison.

### E.3 ADDITIONAL RESULTS FOR REPLAY SCHEDULING POLICY GENERALIZATION EXPERIMENTS

Here, we display the ACC and BWT metrics for each method averaged across 5 seeds and the average rank in every test environment. Note that the ACC and BWT from ETS and the heuristic scheduling baselines have standard deviation zero since these policies are fixed. Averaging the ranks over all test environments yields the corresponding average rank in Table 4. Furthermore, we provide the p-values from Welch’s t-test to show whether the statistical significance of the results.

- **New Task Order, Split MNIST:** Metrics in Table 28, and Welch’s t-test in Table 29.
- **New Task Order, Split FashionMNIST:** Metrics in Table 31, and Welch’s t-test in Table 30.
- **New Task Order, Split notMNIST:** Metrics in Table 32, and Welch’s t-test in Table 33.
- **New Task Order, Split CIFAR-10:** Metrics in Table 35, and Welch’s t-test in Table 34.
- **New Dataset, Split FashionMNIST:** Metrics in Table 36, and Welch’s t-test in Table 37.
- **New Dataset, Split notMNIST:** Metrics in Table 39, and Welch’s t-test in Table 38.

The improvement with the learning replay scheduling policies is less significant than in the MCTS experiments, however, such behaviour is common when RL is used for generalizing to new environments.

Table 28: Performance comparison in every test environment with seed (10-19) with with **Split MNIST** for **New Task Order** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 10			Test Env. Seed 11		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.95 ± 2.68	-6.95 ± 3.33	4.2	92.13 ± 1.34	-9.46 ± 1.68	4.6
ETS	89.89 ± 0.00	-12.01 ± 0.00	6.8	93.77 ± 0.00	-7.41 ± 0.00	2.8
Heur-GD	94.64 ± 0.00	-6.06 ± 0.00	4.5	91.34 ± 0.00	-10.47 ± 0.00	5.6
Heur-LD	95.63 ± 0.00	-4.80 ± 0.00	1.8	91.34 ± 0.00	-10.47 ± 0.00	5.6
Heur-AT	94.64 ± 0.00	-6.06 ± 0.00	4.5	91.34 ± 0.00	-10.47 ± 0.00	5.6
DQN	94.68 ± 1.16	-6.00 ± 1.44	3.2	94.08 ± 1.75	-7.01 ± 2.19	2.8
A2C	95.31 ± 0.00	-5.21 ± 0.00	3	96.41 ± 0.18	-4.09 ± 0.23	1
Method	Test Env. Seed 12			Test Env. Seed 13		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	92.22 ± 2.35	-9.16 ± 2.93	4.4	93.65 ± 2.06	-7.47 ± 2.58	4.8
ETS	91.69 ± 0.00	-9.82 ± 0.00	6.4	94.95 ± 0.00	-5.88 ± 0.00	3
Heur-GD	94.82 ± 0.00	-5.89 ± 0.00	1.7	94.66 ± 0.00	-6.14 ± 0.00	4.7
Heur-LD	91.93 ± 0.00	-9.50 ± 0.00	5.4	93.17 ± 0.00	-8.00 ± 0.00	6.8
Heur-AT	94.82 ± 0.00	-5.89 ± 0.00	1.7	94.66 ± 0.00	-6.14 ± 0.00	4.7
DQN	93.05 ± 1.37	-8.05 ± 1.71	4.6	95.59 ± 1.22	-4.94 ± 1.53	2
A2C	93.62 ± 0.00	-7.34 ± 0.00	3.8	95.56 ± 0.00	-5.01 ± 0.00	2
Method	Test Env. Seed 14			Test Env. Seed 15		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	85.74 ± 3.47	-17.19 ± 4.32	3.8	94.54 ± 1.49	-6.20 ± 1.87	4.8
ETS	87.29 ± 0.00	-15.23 ± 0.00	2.4	95.32 ± 0.00	-5.23 ± 0.00	4.6
Heur-GD	81.20 ± 0.00	-23.00 ± 0.00	6.3	95.92 ± 0.00	-4.49 ± 0.00	2.7
Heur-LD	81.20 ± 0.00	-23.00 ± 0.00	6.3	96.05 ± 0.00	-4.30 ± 0.00	1
Heur-AT	82.36 ± 0.00	-21.52 ± 0.00	4.8	95.92 ± 0.00	-4.49 ± 0.00	2.7
DQN	91.22 ± 3.23	-10.45 ± 4.02	1.6	94.37 ± 0.74	-6.45 ± 0.90	6.4
A2C	88.16 ± 5.81	-14.27 ± 7.25	2.8	94.82 ± 0.00	-5.94 ± 0.00	5.8
Method	Test Env. Seed 16			Test Env. Seed 17		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	83.05 ± 3.07	-20.78 ± 3.84	6.2	95.86 ± 0.48	-4.52 ± 0.59	2.4
ETS	79.38 ± 0.00	-25.36 ± 0.00	6.8	95.69 ± 0.00	-4.77 ± 0.00	2.8
Heur-GD	91.16 ± 0.00	-10.57 ± 0.00	3.2	93.48 ± 0.00	-7.39 ± 0.00	5.8
Heur-LD	91.16 ± 0.00	-10.57 ± 0.00	3.2	93.48 ± 0.00	-7.39 ± 0.00	5.8
Heur-AT	91.16 ± 0.00	-10.57 ± 0.00	3.2	93.48 ± 0.00	-7.39 ± 0.00	5.8
DQN	92.93 ± 1.19	-8.41 ± 1.48	1.8	94.67 ± 2.13	-5.91 ± 2.65	4
A2C	91.11 ± 0.98	-10.69 ± 1.23	3.6	96.28 ± 0.21	-3.89 ± 0.27	1.4
Method	Test Env. Seed 18			Test Env. Seed 19		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	91.30 ± 2.91	-10.28 ± 3.63	2.8	95.85 ± 2.20	-4.71 ± 2.76	1.8
ETS	92.89 ± 0.00	-8.28 ± 0.00	1.4	97.40 ± 0.00	-2.78 ± 0.00	1.2
Heur-GD	87.82 ± 0.00	-14.53 ± 0.00	5.8	88.34 ± 0.00	-14.21 ± 0.00	5
Heur-LD	87.82 ± 0.00	-14.53 ± 0.00	5.8	88.34 ± 0.00	-14.21 ± 0.00	5
Heur-AT	87.82 ± 0.00	-14.53 ± 0.00	5.8	88.34 ± 0.00	-14.21 ± 0.00	5
DQN	90.05 ± 0.88	-11.74 ± 1.12	3.8	89.04 ± 1.82	-13.28 ± 2.28	4.4
A2C	91.64 ± 0.00	-9.75 ± 0.00	2.6	87.64 ± 1.88	-15.03 ± 2.33	5.6

Table 29: Two-tailed Welch’s  $t$ -test results for **Split MNIST** in **New Task Order** experiment.

Methods	Test Env. Seed 10		Test Env. Seed 11		Test Env. Seed 12		Test Env. Seed 13		Test Env. Seed 14	
	$t$	$p$								
DQN vs Random	0.50	0.636	1.77	0.117	0.62	0.560	1.62	0.152	2.32	0.049
DQN vs ETS	8.30	<b>0.001</b>	0.36	0.740	1.99	0.117	1.05	0.352	2.44	0.071
DQN vs Heur-GD	0.08	0.942	3.13	<b>0.035</b>	-2.57	0.062	1.53	0.201	6.21	<b>0.003</b>
DQN vs Heur-LD:	-1.64	0.175	3.13	<b>0.035</b>	1.63	0.177	3.97	<b>0.017</b>	6.21	<b>0.003</b>
DQN vs Heur-AT	0.08	0.942	3.13	<b>0.035</b>	-2.57	0.062	1.53	0.201	5.49	<b>0.005</b>
DQN vs A2C	-1.08	0.341	-2.65	0.056	-0.83	0.455	0.06	0.957	0.92	0.391
A2C vs Random	1.01	0.369	6.32	<b>0.003</b>	1.20	0.298	1.85	0.138	0.72	0.499
A2C vs ETS	inf	<b>0.000</b>	28.80	<b>0.000</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	0.30	0.778
A2C vs Heur-GD	inf	<b>0.000</b>	55.32	<b>0.000</b>	-inf	<b>0.000</b>	inf	<b>0.000</b>	2.40	0.075
A2C vs Heur-LD	-inf	<b>0.000</b>	55.32	<b>0.000</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	2.40	0.075
A2C vs Heur-AT	inf	<b>0.000</b>	55.32	<b>0.000</b>	-inf	<b>0.000</b>	inf	<b>0.000</b>	2.00	0.117
A2C vs DQN	1.08	0.341	2.65	0.056	0.83	0.455	-0.06	0.957	-0.92	0.391
	Test Env. Seed 15		Test Env. Seed 16		Test Env. Seed 17		Test Env. Seed 18		Test Env. Seed 19	
Methods	$t$	$p$								
DQN vs Random	-0.21	0.844	6.01	<b>0.002</b>	-1.09	0.332	-0.82	0.452	-4.76	<b>0.002</b>
DQN vs ETS	-2.60	0.060	22.75	<b>0.000</b>	-0.95	0.396	-6.48	<b>0.003</b>	-9.17	<b>0.001</b>
DQN vs Heur-GD	-4.22	<b>0.013</b>	2.98	<b>0.041</b>	1.12	0.324	5.10	<b>0.007</b>	0.77	0.484
DQN vs Heur-LD:	-4.57	<b>0.010</b>	2.98	<b>0.041</b>	1.12	0.324	5.10	<b>0.007</b>	0.77	0.484
DQN vs Heur-AT	-4.22	<b>0.013</b>	2.98	<b>0.041</b>	1.12	0.324	5.10	<b>0.007</b>	0.77	0.484
DQN vs A2C	-1.24	0.283	2.36	<b>0.047</b>	-1.50	0.205	-3.63	<b>0.022</b>	1.07	0.315
A2C vs Random	0.38	0.722	5.00	0.005	1.60	0.164	0.24	0.824	-5.67	<b>0.001</b>
A2C vs ETS	-inf	<b>0.000</b>	23.87	<b>0.000</b>	5.56	<b>0.005</b>	-inf	<b>0.000</b>	-10.39	<b>0.000</b>
A2C vs Heur-GD	-inf	<b>0.000</b>	-0.10	0.925	26.08	<b>0.000</b>	inf	<b>0.000</b>	-0.75	0.498
A2C vs Heur-LD	-inf	<b>0.000</b>	-0.10	0.925	26.08	<b>0.000</b>	inf	<b>0.000</b>	-0.75	0.498
A2C vs Heur-AT	-inf	<b>0.000</b>	-0.10	0.925	26.08	<b>0.000</b>	inf	<b>0.000</b>	-0.75	0.498
A2C vs DQN	1.24	0.283	-2.36	<b>0.047</b>	1.50	0.205	3.63	0.022	-1.07	0.315

Table 30: Two-tailed Welch’s  $t$ -test results for **Split FashionMNIST** in **New Task Order** experiment.

Methods	Test Env. Seed 10		Test Env. Seed 11		Test Env. Seed 12		Test Env. Seed 13		Test Env. Seed 14	
	$t$	$p$								
DQN vs Random	-0.05	0.962	-0.17	0.871	0.18	0.864	3.28	0.024	0.24	0.819
DQN vs ETS	1.78	0.150	0.76	0.490	2.06	0.109	12.55	<b>0.000</b>	4.19	<b>0.014</b>
DQN vs Heur-GD	-3.61	<b>0.023</b>	-1.29	0.265	-0.33	0.761	4.03	<b>0.016</b>	-0.85	0.445
DQN vs Heur-LD:	-3.61	<b>0.023</b>	-1.99	0.118	-1.07	0.343	1.51	0.205	-0.56	0.604
DQN vs Heur-AT	13.80	<b>0.000</b>	-1.45	0.221	-0.88	0.427	-2.62	0.059	11.80	<b>0.000</b>
DQN vs A2C	-0.61	0.574	2.26	0.084	4.35	<b>0.012</b>	-1.51	0.180	-0.13	0.901
A2C vs Random	0.09	0.934	-5.96	<b>0.002</b>	-2.74	0.052	3.87	<b>0.016</b>	0.64	0.556
A2C vs ETS	22.94	<b>0.000</b>	-10.13	<b>0.001</b>	-inf	<b>0.000</b>	24.85	<b>0.000</b>	81.68	<b>0.000</b>
A2C vs Heur-GD	-28.72	<b>0.000</b>	-23.78	<b>0.000</b>	-inf	<b>0.000</b>	10.04	<b>0.001</b>	-13.50	<b>0.000</b>
A2C vs Heur-LD	-28.72	<b>0.000</b>	-28.38	<b>0.000</b>	-inf	<b>0.000</b>	5.65	<b>0.005</b>	-8.14	<b>0.001</b>
A2C vs Heur-AT	138.22	<b>0.000</b>	-24.81	<b>0.000</b>	-inf	<b>0.000</b>	-1.54	0.199	225.61	<b>0.000</b>
A2C vs DQN	0.61	0.574	-2.26	0.084	-4.35	0.012	1.51	0.180	0.13	0.901
	Test Env. Seed 15		Test Env. Seed 16		Test Env. Seed 17		Test Env. Seed 18		Test Env. Seed 19	
Methods	$t$	$p$								
DQN vs Random	-7.48	<b>0.000</b>	-1.05	0.327	-1.11	0.300	3.23	<b>0.028</b>	-2.26	0.067
DQN vs ETS	-13.28	<b>0.000</b>	-4.06	<b>0.015</b>	4.28	<b>0.013</b>	5.49	<b>0.005</b>	-2.31	0.082
DQN vs Heur-GD	34.42	<b>0.000</b>	12.59	<b>0.000</b>	-3.97	<b>0.017</b>	7.21	<b>0.002</b>	-0.34	0.752
DQN vs Heur-LD:	-10.25	<b>0.001</b>	7.27	<b>0.002</b>	-3.64	<b>0.022</b>	9.52	<b>0.001</b>	2.41	0.073
DQN vs Heur-AT	3.44	<b>0.026</b>	0.12	0.912	0.16	0.883	3.72	<b>0.021</b>	-1.24	0.281
DQN vs A2C	27.78	<b>0.000</b>	2.23	0.056	-0.64	0.540	3.54	<b>0.010</b>	-1.79	0.143
A2C vs Random	-26.79	<b>0.000</b>	-3.58	<b>0.010</b>	-0.50	0.634	2.39	0.073	-1.34	0.235
A2C vs ETS	-109.10	<b>0.000</b>	-6.83	<b>0.002</b>	5.87	<b>0.004</b>	2.18	0.094	-2.52	0.066
A2C vs Heur-GD	11.50	<b>0.000</b>	8.66	<b>0.001</b>	-3.56	<b>0.024</b>	4.99	<b>0.008</b>	7.67	<b>0.002</b>
A2C vs Heur-LD	-101.45	<b>0.000</b>	3.71	<b>0.021</b>	-3.19	<b>0.033</b>	8.76	<b>0.001</b>	21.88	<b>0.000</b>
A2C vs Heur-AT	-66.83	<b>0.000</b>	-2.94	<b>0.042</b>	1.15	0.313	-0.71	0.517	3.00	<b>0.040</b>
A2C vs DQN	-27.78	<b>0.000</b>	-2.23	0.056	0.64	0.540	-3.54	0.010	1.79	0.143

Table 31: Performance comparison in every test environment with seed (10-19) with with **Split FashionMNIST** for **New Task Order** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 10			Test Env. Seed 11		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	96.79 ± 3.02	-3.11 ± 3.76	2.2	90.84 ± 1.64	-8.02 ± 2.04	4.4
ETS	96.10 ± 0.00	-3.98 ± 0.00	5.6	88.84 ± 0.00	-10.55 ± 0.00	5.8
Heur-GD	97.96 ± 0.00	-1.93 ± 0.00	2.3	93.21 ± 0.00	-4.96 ± 0.00	3.4
Heur-LD	97.96 ± 0.00	-1.93 ± 0.00	2.3	94.68 ± 0.00	-3.12 ± 0.00	1
Heur-AT	91.95 ± 0.00	-9.30 ± 0.00	6.8	93.54 ± 0.00	-4.70 ± 0.00	2
DQN	96.71 ± 0.69	-3.47 ± 0.86	4.5	90.46 ± 4.25	-8.60 ± 5.32	4.6
A2C	96.93 ± 0.07	-3.20 ± 0.09	4.3	85.60 ± 0.64	-14.65 ± 0.79	6.8
Method	Test Env. Seed 12			Test Env. Seed 13		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.97 ± 4.60	-6.60 ± 5.74	2.8	91.66 ± 3.08	-9.38 ± 3.85	6.4
ETS	91.25 ± 0.00	-10.08 ± 0.00	5.6	91.24 ± 0.00	-9.85 ± 0.00	6.6
Heur-GD	94.97 ± 0.00	-5.26 ± 0.00	4.2	95.09 ± 0.00	-5.10 ± 0.00	5
Heur-LD	96.14 ± 0.00	-3.89 ± 0.00	1.8	96.23 ± 0.00	-3.69 ± 0.00	3.8
Heur-AT	95.84 ± 0.00	-4.21 ± 0.00	3	98.10 ± 0.00	-1.38 ± 0.00	1.2
DQN	94.46 ± 3.12	-6.00 ± 3.95	3.8	96.91 ± 0.90	-2.83 ± 1.12	2.6
A2C	87.67 ± 0.00	-14.49 ± 0.00	6.8	97.70 ± 0.52	-1.87 ± 0.67	2.4
Method	Test Env. Seed 14			Test Env. Seed 15		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	94.17 ± 1.37	-3.79 ± 1.73	3.6	93.74 ± 0.95	-4.53 ± 1.18	1.4
ETS	90.04 ± 0.00	-8.92 ± 0.00	6	93.51 ± 0.00	-4.81 ± 0.00	1.6
Heur-GD	95.37 ± 0.00	-2.26 ± 0.00	1.6	79.33 ± 0.00	-22.05 ± 0.00	7
Heur-LD	95.07 ± 0.00	-2.65 ± 0.00	3.2	92.61 ± 0.00	-5.44 ± 0.00	3
Heur-AT	81.98 ± 0.00	-18.88 ± 0.00	7	88.54 ± 0.00	-10.50 ± 0.00	4.8
DQN	94.47 ± 2.12	-3.12 ± 2.66	2.4	89.56 ± 0.59	-9.25 ± 0.74	4.2
A2C	94.61 ± 0.11	-2.96 ± 0.14	4.2	80.68 ± 0.24	-20.34 ± 0.29	6
Method	Test Env. Seed 16			Test Env. Seed 17		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	90.96 ± 1.68	-7.18 ± 2.06	2.4	98.99 ± 0.26	-0.66 ± 0.34	3.8
ETS	94.41 ± 0.00	-2.96 ± 0.00	1	98.11 ± 0.00	-1.77 ± 0.00	7
Heur-GD	73.82 ± 0.00	-28.91 ± 0.00	7	99.37 ± 0.00	-0.24 ± 0.00	1
Heur-LD	80.40 ± 0.00	-20.66 ± 0.00	6	99.32 ± 0.00	-0.30 ± 0.00	2
Heur-AT	89.24 ± 0.00	-9.68 ± 0.00	3.6	98.74 ± 0.00	-0.99 ± 0.00	4.8
DQN	89.39 ± 2.47	-9.49 ± 3.09	3.4	98.76 ± 0.31	-0.97 ± 0.38	5
A2C	85.33 ± 2.66	-14.62 ± 3.36	4.6	98.89 ± 0.27	-0.80 ± 0.33	4.4
Method	Test Env. Seed 18			Test Env. Seed 19		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	89.92 ± 3.44	-11.26 ± 4.30	5.6	97.64 ± 0.79	-1.58 ± 1.06	1.8
ETS	93.56 ± 0.00	-6.74 ± 0.00	4.4	97.49 ± 0.00	-1.81 ± 0.00	2
Heur-GD	92.92 ± 0.00	-7.31 ± 0.00	5.4	95.79 ± 0.00	-3.94 ± 0.00	5.4
Heur-LD	92.06 ± 0.00	-8.40 ± 0.00	6.4	93.42 ± 0.00	-6.85 ± 0.00	6.8
Heur-AT	94.22 ± 0.00	-5.66 ± 0.00	2.2	96.57 ± 0.00	-3.01 ± 0.00	4.2
DQN	95.60 ± 0.74	-3.92 ± 0.93	1	95.50 ± 1.72	-4.42 ± 2.16	5
A2C	94.06 ± 0.46	-5.88 ± 0.58	3	97.07 ± 0.33	-2.41 ± 0.40	2.8

Table 32: Performance comparison in every test environment with seed (10-19) with **Split notM-NIST** for **New Task Order** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 10			Test Env. Seed 11		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	89.19 ± 3.60	-9.93 ± 4.54	5	91.90 ± 1.05	-2.01 ± 1.35	4
ETS	93.24 ± 0.00	-4.74 ± 0.00	1.6	91.79 ± 0.00	-2.28 ± 0.00	3.4
Heur-GD	91.58 ± 0.00	-5.60 ± 0.00	4	90.92 ± 0.00	-4.21 ± 0.00	6.8
Heur-LD	90.88 ± 0.00	-6.51 ± 0.00	6.4	91.60 ± 0.00	-3.31 ± 0.00	5.6
Heur-AT	91.36 ± 0.00	-5.48 ± 0.00	5	91.70 ± 0.00	-3.44 ± 0.00	4.6
DQN	93.70 ± 1.02	-3.37 ± 0.72	1.8	92.73 ± 0.77	-2.12 ± 0.65	1.8
A2C	91.78 ± 0.60	-5.63 ± 0.49	4.2	92.66 ± 0.36	-2.31 ± 0.33	1.8
Method	Test Env. Seed 12			Test Env. Seed 13		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	90.59 ± 2.85	-7.94 ± 3.69	4	91.63 ± 2.40	-3.73 ± 2.70	3.8
ETS	83.06 ± 0.00	-17.48 ± 0.00	7	86.28 ± 0.00	-10.76 ± 0.00	7
Heur-GD	92.89 ± 0.00	-4.45 ± 0.00	2.2	94.06 ± 0.00	-1.22 ± 0.00	1.2
Heur-LD	92.96 ± 0.00	-3.57 ± 0.00	1.2	92.21 ± 0.00	-3.53 ± 0.00	4
Heur-AT	91.20 ± 0.00	-5.92 ± 0.00	4.6	88.77 ± 0.00	-7.46 ± 0.00	5.8
DQN	91.59 ± 0.87	-5.46 ± 0.95	3.8	92.30 ± 1.48	-3.42 ± 1.86	3.5
A2C	91.02 ± 0.26	-5.97 ± 0.19	5.2	93.52 ± 0.60	-2.05 ± 0.44	2.7
Method	Test Env. Seed 14			Test Env. Seed 15		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	92.42 ± 0.80	-3.95 ± 0.94	3	89.42 ± 1.66	-6.65 ± 1.97	3.6
ETS	92.67 ± 0.00	-3.89 ± 0.00	2.4	89.32 ± 0.00	-6.41 ± 0.00	4.2
Heur-GD	93.63 ± 0.00	-1.46 ± 0.00	1	92.08 ± 0.00	-4.12 ± 0.00	1
Heur-LD	89.18 ± 0.00	-7.02 ± 0.00	5.8	89.58 ± 0.00	-7.25 ± 0.00	3
Heur-AT	88.14 ± 0.00	-8.58 ± 0.00	6.8	84.47 ± 0.00	-14.12 ± 0.00	6.8
DQN	90.69 ± 0.79	-5.25 ± 0.95	4.8	86.99 ± 1.50	-10.77 ± 1.87	5.6
A2C	91.47 ± 1.79	-4.64 ± 2.34	4.2	89.04 ± 2.55	-8.28 ± 3.07	3.8
Method	Test Env. Seed 16			Test Env. Seed 17		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.66 ± 0.54	-0.94 ± 0.77	1	90.30 ± 2.35	-7.20 ± 2.85	6
ETS	91.91 ± 0.00	-2.65 ± 0.00	2.8	92.56 ± 0.00	-3.99 ± 0.00	3.4
Heur-GD	89.22 ± 0.00	-5.95 ± 0.00	4.7	91.30 ± 0.00	-5.69 ± 0.00	5.9
Heur-LD	89.22 ± 0.00	-5.95 ± 0.00	4.7	91.30 ± 0.00	-5.69 ± 0.00	5.9
Heur-AT	86.29 ± 0.00	-9.36 ± 0.00	7	93.87 ± 0.00	-1.30 ± 0.00	2
DQN	88.94 ± 1.41	-5.84 ± 1.93	5.2	94.53 ± 0.90	-1.41 ± 1.11	1.2
A2C	91.85 ± 1.45	-1.88 ± 1.43	2.6	92.58 ± 0.67	-3.78 ± 0.99	3.6
Method	Test Env. Seed 18			Test Env. Seed 19		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	94.78 ± 2.97	-2.78 ± 3.76	2	92.83 ± 2.19	-2.62 ± 2.73	4.4
ETS	92.43 ± 0.00	-5.57 ± 0.00	5.8	90.89 ± 0.00	-5.17 ± 0.00	6.8
Heur-GD	92.79 ± 0.00	-5.69 ± 0.00	4.4	94.16 ± 0.00	-2.21 ± 0.00	3.2
Heur-LD	94.65 ± 0.00	-3.71 ± 0.00	1.8	94.94 ± 0.00	-1.01 ± 0.00	1.2
Heur-AT	88.15 ± 0.00	-11.30 ± 0.00	7	93.46 ± 0.00	-3.99 ± 0.00	5.4
DQN	93.32 ± 0.69	-4.77 ± 1.02	3.8	94.08 ± 0.96	-3.10 ± 1.02	3.6
A2C	93.64 ± 0.22	-4.58 ± 0.48	3.2	93.66 ± 0.80	-3.82 ± 0.94	3.4

Table 33: Two-tailed Welch’s  $t$ -test results for **Split notMNIST** in **New Task Order** experiment.

Methods	Test Env. Seed 10		Test Env. Seed 11		Test Env. Seed 12		Test Env. Seed 13		Test Env. Seed 14	
	$t$	$p$								
DQN vs Random	-1.32	0.226	2.81	<b>0.023</b>	2.64	<b>0.033</b>	2.22	0.083	-0.41	0.692
DQN vs ETS	-0.36	<b>0.739</b>	3.08	<b>0.037</b>	0.68	0.535	8.68	<b>0.001</b>	0.20	0.849
DQN vs Heur-GD	-1.20	0.296	6.06	<b>0.004</b>	9.21	<b>0.001</b>	-7.65	<b>0.002</b>	-0.72	0.512
DQN vs Heur-LD:	-1.20	0.296	6.06	<b>0.004</b>	9.21	<b>0.001</b>	-7.65	<b>0.002</b>	-1.56	0.194
DQN vs Heur-AT	-1.20	0.296	6.06	<b>0.004</b>	9.21	<b>0.001</b>	-7.65	<b>0.002</b>	-0.72	0.512
DQN vs A2C	-0.63	0.554	0.78	0.458	1.29	0.245	-3.88	<b>0.005</b>	-1.64	0.172
A2C vs Random	-1.04	0.333	1.74	0.122	2.10	0.091	3.57	<b>0.019</b>	1.68	0.159
A2C vs ETS	0.67	0.537	1.57	0.191	-1.56	0.194	12.84	<b>0.000</b>	9.59	<b>0.001</b>
A2C vs Heur-GD	-0.96	0.390	4.08	<b>0.015</b>	15.93	<b>0.000</b>	-1.58	0.189	4.86	<b>0.008</b>
A2C vs Heur-LD	-0.96	0.390	4.08	<b>0.015</b>	15.93	<b>0.000</b>	-1.58	0.189	0.56	0.606
A2C vs Heur-AT	-0.96	0.390	4.08	<b>0.015</b>	15.93	<b>0.000</b>	-1.58	0.189	4.86	0.008
A2C vs DQN	0.63	0.554	-0.78	0.458	-1.29	0.245	3.88	<b>0.005</b>	1.64	0.172
Methods	Test Env. Seed 5		Test Env. Seed 6		Test Env. Seed 7		Test Env. Seed 8		Test Env. Seed 9	
	$t$	$p$								
DQN vs Random	1.82	0.108	1.05	0.332	-1.16	0.283	0.85	0.438	-3.30	<b>0.023</b>
DQN vs ETS	3.61	<b>0.023</b>	2.01	0.115	-3.07	<b>0.037</b>	7.25	<b>0.002</b>	-4.29	<b>0.013</b>
DQN vs Heur-GD	1.67	0.169	2.66	0.056	15.47	<b>0.000</b>	-0.78	0.479	2.69	0.055
DQN vs Heur-LD:	1.67	0.169	9.49	<b>0.001</b>	8.95	<b>0.001</b>	8.72	<b>0.001</b>	2.69	0.055
DQN vs Heur-AT	1.67	0.169	2.66	0.056	0.32	0.767	-0.78	0.479	2.69	0.055
DQN vs A2C	-0.38	0.722	1.63	0.142	-0.39	0.710	-1.43	0.226	-2.71	0.053
A2C vs Random	2.56	0.060	-0.07	0.944	-0.83	0.430	8.71	<b>0.001</b>	-2.39	0.073
A2C vs ETS	18.34	<b>0.000</b>	-0.40	0.707	-2.50	0.067	inf	<b>0.000</b>	-39.30	<b>0.000</b>
A2C vs Heur-GD	9.47	<b>0.001</b>	0.18	0.867	15.87	<b>0.000</b>	inf	<b>0.000</b>	135.29	<b>0.000</b>
A2C vs Heur-LD	9.47	<b>0.001</b>	6.23	<b>0.003</b>	9.41	<b>0.001</b>	inf	<b>0.000</b>	135.29	<b>0.000</b>
A2C vs Heur-AT	9.47	<b>0.001</b>	0.18	0.867	0.86	0.439	inf	<b>0.000</b>	135.29	<b>0.000</b>
A2C vs DQN	0.38	0.722	-1.63	0.142	0.39	0.710	1.43	0.226	2.71	0.053

Table 34: Two-tailed Welch’s  $t$ -test results for **Split CIFAR-10** in **New Task Order** experiment.

Methods	Test Env. Seed 10		Test Env. Seed 11		Test Env. Seed 12		Test Env. Seed 13		Test Env. Seed 14	
	$t$	$p$								
DQN vs Random	0.70	0.504	8.12	<b>0.000</b>	-3.77	<b>0.006</b>	0.48	0.653	2.03	0.091
DQN vs ETS	-0.07	0.948	12.90	<b>0.000</b>	-0.79	0.473	12.33	<b>0.000</b>	12.12	<b>0.000</b>
DQN vs Heur-GD	1.18	0.303	-0.09	0.932	-10.41	<b>0.000</b>	5.18	<b>0.007</b>	-3.77	<b>0.020</b>
DQN vs Heur-LD:	0.61	0.574	-6.72	<b>0.003</b>	-5.31	<b>0.006</b>	-1.29	0.265	1.24	0.283
DQN vs Heur-AT	1.77	0.152	-5.15	<b>0.007</b>	-9.97	<b>0.001</b>	0.84	0.448	-3.26	<b>0.031</b>
DQN vs A2C	3.97	<b>0.017</b>	-10.54	<b>0.000</b>	-2.38	0.051	-2.70	<b>0.049</b>	-1.86	0.127
A2C vs Random	-2.29	0.084	13.83	<b>0.000</b>	-2.31	0.061	1.63	0.178	3.15	<b>0.033</b>
A2C vs ETS	inf	<b>0.000</b>	inf	<b>0.000</b>	3.19	<b>0.033</b>	68.53	<b>0.000</b>	50.37	<b>0.000</b>
A2C vs Heur-GD	inf	<b>0.000</b>	inf	<b>0.000</b>	-12.02	<b>0.000</b>	36.03	<b>0.000</b>	-6.59	<b>0.003</b>
A2C vs Heur-LD	inf	<b>0.000</b>	inf	<b>0.000</b>	-3.96	<b>0.017</b>	6.66	<b>0.003</b>	11.37	<b>0.000</b>
A2C vs Heur-AT	inf	<b>0.000</b>	inf	<b>0.000</b>	-11.34	<b>0.000</b>	16.34	<b>0.000</b>	-4.77	<b>0.009</b>
A2C vs DQN	-3.97	<b>0.017</b>	10.54	<b>0.000</b>	2.38	0.051	2.70	0.049	1.86	0.127
Methods	Test Env. Seed 15		Test Env. Seed 16		Test Env. Seed 17		Test Env. Seed 18		Test Env. Seed 19	
	$t$	$p$								
DQN vs Random	-2.67	<b>0.045</b>	0.07	0.945	1.35	0.222	1.53	0.165	4.03	<b>0.010</b>
DQN vs ETS	-1.66	0.172	-1.08	0.342	2.69	0.055	0.20	0.850	5.14	<b>0.007</b>
DQN vs Heur-GD	-1.38	0.239	7.09	<b>0.002</b>	0.56	0.608	-0.38	0.725	4.68	<b>0.009</b>
DQN vs Heur-LD:	-1.98	0.119	9.15	<b>0.001</b>	-0.54	0.621	-3.08	<b>0.037</b>	5.53	<b>0.005</b>
DQN vs Heur-AT	-1.98	0.119	7.09	<b>0.002</b>	-0.74	0.502	-0.38	0.725	4.68	<b>0.009</b>
DQN vs A2C	-2.65	0.054	-2.90	<b>0.044</b>	-0.33	0.756	-1.86	0.119	-1.31	0.259
A2C vs Random	-0.35	0.740	0.50	0.643	3.17	<b>0.034</b>	4.00	<b>0.009</b>	4.80	<b>0.009</b>
A2C vs ETS	6.26	<b>0.003</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	5.35	<b>0.006</b>	inf	<b>0.000</b>
A2C vs Heur-GD	7.95	<b>0.001</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	3.95	<b>0.017</b>	inf	<b>0.000</b>
A2C vs Heur-LD	4.31	<b>0.013</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	-2.58	0.061	inf	<b>0.000</b>
A2C vs Heur-AT	4.31	<b>0.013</b>	inf	<b>0.000</b>	inf	<b>0.000</b>	3.95	<b>0.017</b>	inf	<b>0.000</b>
A2C vs DQN	2.65	0.054	2.90	<b>0.044</b>	0.33	0.756	1.86	0.119	1.31	0.259

Table 35: Performance comparison in every test environment with seed (10-19) with with **Split CIFAR-10** for **New Task Order** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 10			Test Env. Seed 11		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	86.35 ± 1.57	-13.27 ± 1.98	4.2	73.05 ± 1.42	-27.40 ± 1.91	6.8
ETS	87.10 ± 0.00	-12.31 ± 0.00	1.8	75.16 ± 0.00	-25.14 ± 0.00	6.2
Heur-GD	86.31 ± 0.00	-13.19 ± 0.00	4	79.45 ± 0.00	-19.79 ± 0.00	4.6
Heur-LD	86.67 ± 0.00	-12.74 ± 0.00	3	81.64 ± 0.00	-17.05 ± 0.00	2
Heur-AT	85.94 ± 0.00	-13.69 ± 0.00	5.2	81.12 ± 0.00	-17.74 ± 0.00	3
DQN	87.06 ± 1.26	-12.42 ± 1.57	2.8	79.42 ± 0.66	-19.64 ± 0.83	4.4
A2C	84.55 ± 0.00	-15.59 ± 0.00	7	82.90 ± 0.00	-15.26 ± 0.00	1
Method	Test Env. Seed 12			Test Env. Seed 13		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	87.86 ± 1.12	-10.36 ± 1.44	3.4	80.07 ± 1.43	-15.40 ± 1.78	3.3
ETS	85.53 ± 0.00	-13.33 ± 0.00	6.6	76.85 ± 0.00	-19.34 ± 0.00	7
Heur-GD	89.76 ± 0.00	-7.69 ± 0.00	1	78.93 ± 0.00	-16.48 ± 0.00	5.8
Heur-LD	87.52 ± 0.00	-10.40 ± 0.00	3.8	80.81 ± 0.00	-14.20 ± 0.00	2.8
Heur-AT	89.57 ± 0.00	-8.15 ± 0.00	2	80.19 ± 0.00	-14.77 ± 0.00	4.3
DQN	85.18 ± 0.88	-13.31 ± 1.12	6.4	80.43 ± 0.58	-14.47 ± 0.75	3.4
A2C	86.42 ± 0.56	-11.90 ± 0.75	4.8	81.24 ± 0.13	-13.35 ± 0.15	1.4
Method	Test Env. Seed 14			Test Env. Seed 15		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	82.04 ± 2.37	-17.84 ± 2.95	6	84.86 ± 0.62	-14.77 ± 0.79	2
ETS	77.79 ± 0.00	-23.26 ± 0.00	6.8	83.78 ± 0.00	-16.15 ± 0.00	5.2
Heur-GD	86.86 ± 0.00	-11.70 ± 0.00	1	83.52 ± 0.00	-16.29 ± 0.00	6.2
Heur-LD	84.00 ± 0.00	-15.24 ± 0.00	4.8	84.08 ± 0.00	-15.58 ± 0.00	3.5
Heur-AT	86.57 ± 0.00	-12.05 ± 0.00	2	84.08 ± 0.00	-15.58 ± 0.00	3.5
DQN	84.71 ± 1.14	-14.62 ± 1.36	4.2	82.22 ± 1.88	-17.83 ± 2.38	5.8
A2C	85.81 ± 0.32	-13.35 ± 0.40	3.2	84.74 ± 0.31	-14.65 ± 0.39	1.8
Method	Test Env. Seed 16			Test Env. Seed 17		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	87.60 ± 1.72	-10.51 ± 2.15	4	72.33 ± 2.14	-26.27 ± 2.73	6
ETS	87.80 ± 0.00	-10.20 ± 0.00	2.8	70.35 ± 0.00	-28.79 ± 0.00	6.6
Heur-GD	86.77 ± 0.00	-12.23 ± 0.00	5.1	74.14 ± 0.00	-23.30 ± 0.00	4.8
Heur-LD	86.51 ± 0.00	-12.63 ± 0.00	6.6	76.08 ± 0.00	-20.88 ± 0.00	2.6
Heur-AT	86.77 ± 0.00	-12.23 ± 0.00	5.1	76.44 ± 0.00	-20.49 ± 0.00	1.4
DQN	87.66 ± 0.25	-11.02 ± 0.33	2.8	75.13 ± 3.56	-21.96 ± 4.45	3
A2C	88.03 ± 0.00	-10.59 ± 0.00	1.6	75.72 ± 0.00	-21.18 ± 0.00	3.6
Method	Test Env. Seed 18			Test Env. Seed 19		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	87.43 ± 0.88	-10.75 ± 1.03	6.4	74.89 ± 2.33	-25.80 ± 2.89	7
ETS	88.32 ± 0.00	-10.19 ± 0.00	5.8	77.67 ± 0.00	-22.16 ± 0.00	5
Heur-GD	88.59 ± 0.00	-10.34 ± 0.00	4.3	77.87 ± 0.00	-22.98 ± 0.00	3.5
Heur-LD	89.85 ± 0.00	-8.95 ± 0.00	1.2	77.50 ± 0.00	-23.50 ± 0.00	6
Heur-AT	88.59 ± 0.00	-10.34 ± 0.00	4.3	77.87 ± 0.00	-22.98 ± 0.00	3.5
DQN	88.41 ± 0.93	-10.73 ± 1.15	3.8	79.91 ± 0.87	-20.52 ± 1.07	1.7
A2C	89.35 ± 0.39	-9.47 ± 0.48	2.2	80.48 ± 0.00	-19.78 ± 0.00	1.3

Table 36: Performance comparison in every test environment with seed (0-9) with with **Split FashionMNIST** for **New Dataset** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 0			Test Env. Seed 1		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	94.33 ± 3.03	-6.15 ± 3.78	5.4	92.26 ± 3.84	-5.96 ± 4.78	4.2
ETS	97.10 ± 0.00	-2.70 ± 0.00	5.8	94.10 ± 0.00	-3.59 ± 0.00	3.2
Heur-GD	97.90 ± 0.00	-1.59 ± 0.00	1	95.01 ± 0.00	-2.69 ± 0.00	1.4
Heur-LD	97.59 ± 0.00	-1.99 ± 0.00	3	90.03 ± 0.00	-8.90 ± 0.00	6.8
Heur-AT	97.41 ± 0.00	-2.21 ± 0.00	4.6	94.09 ± 0.00	-3.78 ± 0.00	4.2
DQN	96.82 ± 1.50	-3.06 ± 1.87	3.8	93.74 ± 1.86	-4.24 ± 2.33	3.2
A2C	95.74 ± 3.33	-4.39 ± 4.15	4.4	92.80 ± 1.56	-5.27 ± 1.95	5
Method	Test Env. Seed 2			Test Env. Seed 3		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.74 ± 2.97	-6.36 ± 3.71	5.8	94.12 ± 4.14	-6.99 ± 5.16	3.9
ETS	86.72 ± 0.00	-15.17 ± 0.00	7	89.44 ± 0.00	-12.86 ± 0.00	6.4
Heur-GD	97.41 ± 0.00	-1.87 ± 0.00	2.8	96.69 ± 0.00	-3.68 ± 0.00	3
Heur-LD	97.30 ± 0.00	-1.91 ± 0.00	3.8	90.61 ± 0.00	-11.26 ± 0.00	4.8
Heur-AT	97.65 ± 0.00	-1.58 ± 0.00	1	99.40 ± 0.00	-0.26 ± 0.00	1.1
DQN	96.13 ± 0.97	-3.47 ± 1.20	5	94.66 ± 5.00	-6.24 ± 6.26	3.6
A2C	97.31 ± 0.55	-1.99 ± 0.71	2.6	92.12 ± 3.75	-9.40 ± 4.68	5.2
Method	Test Env. Seed 4			Test Env. Seed 5		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	83.64 ± 5.22	-16.86 ± 6.48	5.4	89.76 ± 2.46	-7.91 ± 3.06	4
ETS	87.07 ± 0.00	-12.61 ± 0.00	4.2	91.53 ± 0.00	-5.61 ± 0.00	1.2
Heur-GD	91.29 ± 0.00	-7.30 ± 0.00	2.4	90.33 ± 0.00	-7.25 ± 0.00	3.4
Heur-LD	86.75 ± 0.00	-12.87 ± 0.00	5.2	88.01 ± 0.00	-10.02 ± 0.00	4.8
Heur-AT	83.28 ± 0.00	-17.02 ± 0.00	6.4	90.53 ± 0.00	-6.81 ± 0.00	2.4
DQN	88.76 ± 4.73	-10.37 ± 5.95	3.2	87.31 ± 1.09	-10.89 ± 1.38	6.2
A2C	91.98 ± 0.17	-6.36 ± 0.18	1.2	87.93 ± 0.00	-10.07 ± 0.00	6
Method	Test Env. Seed 6			Test Env. Seed 7		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	95.34 ± 0.74	-2.57 ± 1.01	1.4	94.40 ± 0.77	-3.53 ± 0.94	4
ETS	95.48 ± 0.00	-2.50 ± 0.00	1.6	95.31 ± 0.00	-2.21 ± 0.00	2.6
Heur-GD	90.01 ± 0.00	-9.09 ± 0.00	4.2	91.76 ± 0.00	-6.84 ± 0.00	5.6
Heur-LD	85.44 ± 0.00	-14.96 ± 0.00	5.2	95.14 ± 0.00	-2.64 ± 0.00	3.6
Heur-AT	76.02 ± 0.00	-26.58 ± 0.00	7	96.78 ± 0.00	-0.67 ± 0.00	1.2
DQN	86.60 ± 5.92	-13.56 ± 7.39	4.8	89.22 ± 2.34	-10.06 ± 2.88	6.8
A2C	89.14 ± 3.66	-10.35 ± 4.57	3.8	93.49 ± 2.58	-4.77 ± 3.18	4.2
Method	Test Env. Seed 8			Test Env. Seed 9		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.15 ± 4.08	-7.22 ± 5.09	3	94.81 ± 3.66	-5.14 ± 4.51	2.8
ETS	94.80 ± 0.00	-5.18 ± 0.00	2.8	96.71 ± 0.00	-2.76 ± 0.00	2
Heur-GD	95.00 ± 0.00	-4.83 ± 0.00	1.6	95.27 ± 0.00	-4.75 ± 0.00	3.2
Heur-LD	87.86 ± 0.00	-13.66 ± 0.00	6.8	90.02 ± 0.00	-11.16 ± 0.00	6.8
Heur-AT	93.93 ± 0.00	-6.07 ± 0.00	3.8	90.29 ± 0.00	-10.82 ± 0.00	5.8
DQN	91.67 ± 2.85	-8.94 ± 3.61	5	95.32 ± 1.76	-4.60 ± 2.19	2.6
A2C	93.57 ± 0.00	-6.61 ± 0.00	5	91.31 ± 1.08	-9.62 ± 1.37	4.8

Table 37: Two-tailed Welch’s  $t$ -test results for **Split FashionMNIST** in **New Dataset** experiment.

Methods	Test Env. Seed 0		Test Env. Seed 1		Test Env. Seed 2		Test Env. Seed 3		Test Env. Seed 4	
	$t$	$p$								
DQN vs Random	1.47	0.193	0.69	0.515	1.53	0.187	0.17	0.872	1.46	0.184
DQN vs ETS	-0.38	0.725	-0.38	0.720	19.36	<b>0.000</b>	2.09	0.105	0.72	0.513
DQN vs Heur-GD	-1.45	0.222	-1.36	0.245	-2.63	0.058	-0.81	0.463	-1.07	0.345
DQN vs Heur-LD:	-1.03	0.360	3.98	0.016	-2.40	0.074	1.62	0.181	0.85	0.442
DQN vs Heur-AT	-0.79	0.473	-0.37	0.728	-3.12	<b>0.035</b>	-1.90	0.131	2.32	0.081
DQN vs A2C	0.59	0.579	0.77	0.462	-2.11	0.076	0.81	0.442	-1.36	0.245
A2C vs Random	0.63	0.547	0.26	0.804	2.37	0.073	-0.72	0.495	3.19	<b>0.033</b>
A2C vs ETS	-0.82	0.461	-1.67	0.171	38.38	<b>0.000</b>	1.43	0.226	57.27	<b>0.000</b>
A2C vs Heur-GD	-1.30	0.265	-2.84	<b>0.047</b>	-0.35	0.746	-2.44	0.072	8.05	<b>0.001</b>
A2C vs Heur-LD	-1.11	0.329	3.56	<b>0.024</b>	0.05	0.962	0.81	0.465	61.00	<b>0.000</b>
A2C vs Heur-AT	-1.00	0.373	-1.65	0.173	-1.22	0.290	-3.88	0.018	101.48	<b>0.000</b>
A2C vs DQN	-0.59	0.579	-0.77	0.462	2.11	0.076	-0.81	0.442	1.36	0.245
Methods	Test Env. Seed 5		Test Env. Seed 6		Test Env. Seed 7		Test Env. Seed 8		Test Env. Seed 9	
	$t$	$p$								
DQN vs Random	-1.82	0.122	-2.93	<b>0.041</b>	-4.21	<b>0.009</b>	-0.59	0.570	0.25	0.812
DQN vs ETS	-7.73	<b>0.002</b>	-3.00	<b>0.040</b>	-5.20	<b>0.007</b>	-2.20	0.093	-1.58	0.189
DQN vs Heur-GD	-5.53	<b>0.005</b>	-1.15	0.313	-2.17	0.096	-2.34	0.079	0.05	0.961
DQN vs Heur-LD:	-1.28	0.269	0.39	0.714	-5.06	<b>0.007</b>	2.67	0.056	6.00	<b>0.004</b>
DQN vs Heur-AT	-5.90	<b>0.004</b>	3.58	<b>0.023</b>	-6.46	<b>0.003</b>	-1.59	0.187	5.70	<b>0.005</b>
DQN vs A2C	-1.14	0.320	-0.73	0.490	-2.45	<b>0.040</b>	-1.34	0.253	3.87	<b>0.007</b>
A2C vs Random	-1.49	0.210	-3.32	<b>0.026</b>	-0.68	0.529	0.21	0.846	-1.84	0.129
A2C vs ETS	-inf	<b>0.000</b>	-3.47	<b>0.026</b>	-1.41	0.231	-inf	<b>0.000</b>	-10.02	<b>0.001</b>
A2C vs Heur-GD	-inf	<b>0.000</b>	-0.47	0.660	1.34	0.252	-inf	<b>0.000</b>	-7.35	<b>0.002</b>
A2C vs Heur-LD	-inf	<b>0.000</b>	2.03	0.113	-1.28	0.270	inf	<b>0.000</b>	2.39	0.075
A2C vs Heur-AT	-inf	<b>0.000</b>	7.18	<b>0.002</b>	-2.55	0.063	-inf	<b>0.000</b>	1.89	0.131
A2C vs DQN	1.14	0.320	0.73	0.490	2.45	<b>0.040</b>	1.34	0.253	-3.87	<b>0.007</b>

Table 38: Two-tailed Welch’s  $t$ -test results for **Split notMNIST** in **New Dataset** experiment.

Methods	Test Env. Seed 0		Test Env. Seed 1		Test Env. Seed 2		Test Env. Seed 3		Test Env. Seed 4	
	$t$	$p$								
DQN vs Random	2.41	<b>0.065</b>	1.27	0.242	0.67	0.534	0.47	0.651	-3.09	<b>0.015</b>
DQN vs ETS	0.89	0.423	2.45	0.071	19.64	<b>0.000</b>	8.12	<b>0.001</b>	-5.03	<b>0.007</b>
DQN vs Heur-GD	4.15	<b>0.014</b>	4.72	<b>0.009</b>	-3.00	<b>0.040</b>	-2.37	0.077	-7.47	<b>0.002</b>
DQN vs Heur-LD:	5.52	<b>0.005</b>	2.96	<b>0.041</b>	-3.16	<b>0.034</b>	0.13	0.906	3.84	<b>0.018</b>
DQN vs Heur-AT	4.58	<b>0.010</b>	2.70	0.054	0.88	0.428	4.76	<b>0.009</b>	6.47	<b>0.003</b>
DQN vs A2C	3.23	<b>0.016</b>	0.17	0.870	1.25	0.271	-1.52	0.186	-0.80	0.457
A2C vs Random	1.42	0.226	1.36	0.234	0.30	0.777	1.52	0.194	-0.97	0.372
A2C vs ETS	-4.85	<b>0.008</b>	4.73	<b>0.009</b>	60.22	<b>0.000</b>	24.32	<b>0.000</b>	-1.34	0.251
A2C vs Heur-GD	0.67	0.537	9.51	<b>0.001</b>	-14.14	<b>0.000</b>	-1.82	0.142	-2.42	0.073
A2C vs Heur-LD	2.99	<b>0.040</b>	5.82	<b>0.004</b>	-14.66	<b>0.000</b>	4.40	<b>0.012</b>	2.57	0.062
A2C vs Heur-AT	1.40	0.233	5.27	<b>0.006</b>	-1.39	0.238	15.96	<b>0.000</b>	3.73	<b>0.020</b>
A2C vs DQN	-3.23	<b>0.016</b>	-0.17	0.870	-1.25	0.271	1.52	0.186	0.80	0.457
Methods	Test Env. Seed 15		Test Env. Seed 16		Test Env. Seed 17		Test Env. Seed 18		Test Env. Seed 19	
	$t$	$p$								
DQN vs Random	-2.17	0.062	-6.27	<b>0.001</b>	3.36	<b>0.019</b>	-0.96	0.387	1.04	0.340
DQN vs ETS	-3.11	<b>0.036</b>	-4.23	<b>0.013</b>	4.40	<b>0.012</b>	2.60	0.060	6.64	<b>0.003</b>
DQN vs Heur-GD	-6.80	<b>0.002</b>	-0.40	0.708	7.21	<b>0.002</b>	1.56	0.194	-0.16	0.881
DQN vs Heur-LD:	-3.45	<b>0.026</b>	-0.40	0.708	7.21	<b>0.002</b>	-3.85	<b>0.018</b>	-1.79	0.149
DQN vs Heur-AT	3.37	<b>0.028</b>	3.77	<b>0.020</b>	1.47	0.214	15.03	<b>0.000</b>	1.29	0.267
DQN vs A2C	-1.39	0.211	-2.89	<b>0.020</b>	3.50	<b>0.009</b>	-0.87	0.426	0.68	0.514
A2C vs Random	-0.25	0.809	-2.32	0.067	1.86	0.126	-0.77	0.483	0.70	0.512
A2C vs ETS	-0.22	0.837	-0.07	0.945	0.06	0.958	11.20	<b>0.000</b>	6.87	<b>0.002</b>
A2C vs Heur-GD	-2.39	0.075	3.62	<b>0.022</b>	3.84	<b>0.018</b>	7.88	<b>0.001</b>	-1.26	0.277
A2C vs Heur-LD	-0.42	0.694	3.62	<b>0.022</b>	3.84	<b>0.018</b>	-9.37	<b>0.001</b>	-3.20	<b>0.033</b>
A2C vs Heur-AT	3.59	<b>0.023</b>	7.66	<b>0.002</b>	-3.89	<b>0.018</b>	50.85	<b>0.000</b>	0.48	0.659
A2C vs DQN	1.39	0.211	2.89	<b>0.020</b>	-3.50	<b>0.009</b>	0.87	0.426	-0.68	0.514

Table 39: Performance comparison in every test environment with seed (0-9) with with **Split notM-NIST** for **New Dataset** experiment. Under each column named 'Test Env. Seed X', we show the mean and stddev. of ACC and BWT, and the Rank averaged over the RL seeds for the corresponding method.

Method	Test Env. Seed 0			Test Env. Seed 1		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	93.09 ± 2.31	-3.66 ± 2.82	2.4	92.13 ± 0.80	-3.91 ± 0.93	3.8
ETS	91.08 ± 0.00	-6.61 ± 0.00	5.8	92.46 ± 0.00	-3.75 ± 0.00	2.8
Heur-GD	92.41 ± 0.00	-3.55 ± 0.00	3.8	91.18 ± 0.00	-6.26 ± 0.00	5.8
Heur-LD	92.41 ± 0.00	-3.55 ± 0.00	3.8	91.18 ± 0.00	-6.26 ± 0.00	5.8
Heur-AT	92.41 ± 0.00	-3.55 ± 0.00	3.8	91.18 ± 0.00	-6.26 ± 0.00	5.8
DQN	90.52 ± 3.15	-6.92 ± 3.57	4.6	93.79 ± 0.86	-1.65 ± 1.03	1.6
A2C	91.62 ± 1.62	-5.47 ± 2.49	3.8	93.26 ± 1.02	-2.12 ± 1.41	2.4
Method	Test Env. Seed 2			Test Env. Seed 3		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	87.13 ± 3.01	-10.47 ± 3.60	4.4	89.04 ± 3.01	-8.60 ± 3.86	6.4
ETS	91.23 ± 0.00	-5.23 ± 0.00	2	89.33 ± 0.00	-8.28 ± 0.00	6.4
Heur-GD	82.58 ± 0.00	-15.83 ± 0.00	5.8	95.26 ± 0.00	-0.99 ± 0.00	2
Heur-LD	82.58 ± 0.00	-15.83 ± 0.00	5.8	95.26 ± 0.00	-0.99 ± 0.00	2
Heur-AT	82.58 ± 0.00	-15.83 ± 0.00	5.8	95.26 ± 0.00	-0.99 ± 0.00	2
DQN	91.92 ± 2.03	-4.77 ± 2.31	1.6	92.49 ± 0.73	-4.69 ± 0.54	5.2
A2C	90.46 ± 0.99	-6.03 ± 1.24	2.6	94.61 ± 0.82	-2.06 ± 0.46	4
Method	Test Env. Seed 4			Test Env. Seed 5		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	91.37 ± 1.63	-5.04 ± 2.02	4.6	91.06 ± 1.45	-5.87 ± 1.82	5.6
ETS	90.52 ± 0.00	-5.81 ± 0.00	6.2	90.73 ± 0.00	-6.00 ± 0.00	6.4
Heur-GD	91.64 ± 0.00	-6.55 ± 0.00	4.3	91.80 ± 0.00	-3.52 ± 0.00	4
Heur-LD	92.66 ± 0.00	-5.28 ± 0.00	2.4	91.80 ± 0.00	-3.52 ± 0.00	4
Heur-AT	91.64 ± 0.00	-6.55 ± 0.00	4.3	91.80 ± 0.00	-3.52 ± 0.00	4
DQN	90.76 ± 2.44	-6.15 ± 3.21	4.4	92.72 ± 1.11	-3.49 ± 1.10	2.2
A2C	92.80 ± 0.48	-4.09 ± 0.48	1.8	92.94 ± 0.24	-3.00 ± 0.29	1.8
Method	Test Env. Seed 6			Test Env. Seed 7		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	91.80 ± 2.56	-4.63 ± 3.25	3	93.72 ± 0.78	-3.20 ± 0.94	2.4
ETS	91.99 ± 0.00	-5.12 ± 0.00	2.6	94.03 ± 0.00	-1.98 ± 0.00	1.4
Heur-GD	91.55 ± 0.00	-4.47 ± 0.00	4.7	88.81 ± 0.00	-8.97 ± 0.00	7
Heur-LD	87.03 ± 0.00	-10.13 ± 0.00	7	90.65 ± 0.00	-5.98 ± 0.00	6
Heur-AT	91.55 ± 0.00	-4.47 ± 0.00	4.7	93.08 ± 0.00	-2.95 ± 0.00	3.8
DQN	93.32 ± 1.33	-2.59 ± 1.27	2.4	93.17 ± 0.56	-2.79 ± 0.55	4
A2C	91.69 ± 1.50	-4.35 ± 1.87	3.6	93.32 ± 0.57	-2.46 ± 0.59	3.4
Method	Test Env. Seed 8			Test Env. Seed 9		
	ACC (%)	BWT (%)	Rank	ACC (%)	BWT (%)	Rank
Random	92.62 ± 0.30	-4.15 ± 0.61	4.6	92.48 ± 1.34	-5.36 ± 1.50	2.2
ETS	89.04 ± 0.00	-8.36 ± 0.00	6	94.29 ± 0.00	-3.55 ± 0.00	1
Heur-GD	93.56 ± 0.00	-2.86 ± 0.00	2.9	79.09 ± 0.00	-22.95 ± 0.00	5.8
Heur-LD	88.21 ± 0.00	-9.54 ± 0.00	7	79.09 ± 0.00	-22.95 ± 0.00	5.8
Heur-AT	93.56 ± 0.00	-2.86 ± 0.00	2.9	79.09 ± 0.00	-22.95 ± 0.00	5.8
DQN	93.12 ± 1.13	-3.20 ± 1.38	3.4	84.95 ± 4.36	-15.42 ± 5.12	4.6
A2C	93.92 ± 0.00	-2.53 ± 0.00	1.2	90.87 ± 0.17	-8.58 ± 0.22	2.8

## E.4 TASK SPLITS IN TEST ENVIRONMENTS IN POLICY GENERALIZATION EXPERIMENTS

Here, we provide the task splits of the test environments used in the policy generalization experiments in Section 4.2. We evaluated all methods using 10 test environments in all experiments. The test environments in the New Task Order experiments were generated with seeds 10-19. We show the task splits for the Split MNIST, Split FashionMNIST, and Split CIFAR-10 environments in Table 40, 41, and 42 respectively. The test environments in the New Dataset experiments were generated with seeds 0-9. We show the task splits for the Split notMNIST and Split FashionMNIST environments in Table 43 and 44 respectively.

Table 40: Task splits with their corresponding seed for test environments of Split MNIST datasets in the **New Task Orders** experiments in Section 4.2.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
10	8, 2	5, 6	3, 1	0, 7	4, 9
11	7, 8	2, 6	4, 5	1, 3	0, 9
12	5, 8	7, 0	4, 9	3, 2	1, 6
13	3, 5	6, 1	4, 7	8, 9	0, 2
14	3, 9	0, 5	4, 2	1, 7	6, 8
15	2, 6	1, 3	7, 0	9, 4	5, 8
16	6, 2	0, 7	8, 4	3, 1	5, 9
17	7, 2	5, 3	4, 0	9, 8	6, 1
18	7, 9	0, 4	2, 1	6, 5	8, 3
19	1, 7	9, 6	8, 4	3, 0	2, 5

Table 41: Task splits with their corresponding seed for test environments of Split FashionMNIST datasets in the **New Task Orders** experiments in Section 4.2.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
10	Bag, Pullover	Sandal, Shirt	Dress, Trouser	T-shirt/top, Sneaker	Coat, Ankle boot
11	Sneaker, Bag	Pullover, Shirt	Coat, Sandal	Trouser, Dress	T-shirt/top, Ankle boot
12	Sandal, Bag	Sneaker, T-shirt/top	Coat, Ankle boot	Dress, Pullover	Trouser, Shirt
13	Dress, Sandal	Shirt, Trouser	Coat, Sneaker	Bag, Ankle boot	T-shirt/top, Pullover
14	Dress, Ankle boot	T-shirt/top, Sandal	Coat, Pullover	Trouser, Sneaker	Shirt, Bag
15	Pullover, Shirt	Trouser, Dress	Sneaker, T-shirt/top	Ankle boot, Coat	Sandal, Bag
16	Shirt, Pullover	T-shirt/top, Sneaker	Bag, Coat	Dress, Trouser	Sandal, Ankle boot
17	Sneaker, Pullover	Sandal, Dress	Coat, T-shirt/top	Ankle boot, Bag	Shirt, Trouser
18	Sneaker, Ankle boot	T-shirt/top, Coat	Pullover, Trouser	Shirt, Sandal	Bag, Dress
19	Trouser, Sneaker	Ankle boot, Shirt	Bag, Coat	Dress, T-shirt/top	Pullover, Sandal

Table 42: Task splits with their corresponding seed for test environments of Split CIFAR-10 datasets in the **New Task Orders** experiments in Section 4.2.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
10	Ship, Bird	Dog, Frog	Cat, Automobile	Airplane, Horse	Deer, Truck
11	Horse, Ship	Bird, Frog	Deer, Dog	Automobile, Cat	Airplane, Truck
12	Dog, Ship	Horse, Airplane	Deer, Truck	Cat, Bird	Automobile, Frog
13	Cat, Dog	Frog, Automobile	Deer, Horse	Ship, Truck	Airplane, Bird
14	Cat, Truck	Airplane, Dog	Deer, Bird	Automobile, Horse	Frog, Ship
15	Bird, Frog	Automobile, Cat	Horse, Airplane	Truck, Deer	Dog, Ship
16	Frog, Bird	Airplane, Horse	Ship, Deer	Cat, Automobile	Dog, Truck
17	Horse, Bird	Dog, Cat	Deer, Airplane	Truck, Ship	Frog, Automobile
18	Horse, Truck	Airplane, Deer	Bird, Automobile	Frog, Dog	Ship, Cat
19	Automobile, Horse	Truck, Frog	Ship, Deer	Cat, Airplane	Bird, Dog

Table 43: Task splits with their corresponding seed for test environments of Split notMNIST datasets in the **New Dataset** experiments in Section 4.2.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
0	A, B	C, D	E, F	G, H	I, J
1	C, J	G, E	A, D	B, H	I, F
2	E, B	F, A	H, C	D, G	J, I
3	F, E	B, C	J, G	H, A	D, I
4	D, I	E, J	C, G	A, B	F, H
5	J, F	C, E	H, B	A, I	G, D
6	I, B	H, A	G, F	C, E	D, J
7	I, F	A, C	B, J	H, D	G, E
8	I, G	J, A	C, F	H, B	E, D
9	I, E	H, C	B, J	D, A	G, F

Table 44: Task splits with their corresponding seed for test environments of Split FashionMNIST datasets in the **New Dataset** experiments in Section 4.2.

Seed	Task 1	Task 2	Task 3	Task 4	Task 5
0	T-shirt/top, Trouser	Pullover, Dress	Coat, Sandal	Shirt, Sneaker	Bag, Ankle boot
1	Pullover, Ankle boot	Shirt, Coat	T-shirt/top, Dress	Trouser, Sneaker	Bag, Sandal
2	Coat, Trouser	Sandal, T-shirt/top	Sneaker, Pullover	Dress, Shirt	Ankle boot, Bag
3	Sandal, Coat	Trouser, Pullover	Ankle boot, Shirt	Sneaker, T-shirt/top	Dress, Bag
4	Dress, Bag	Coat, Ankle boot	Pullover, Shirt	T-shirt/top, Trouser	Sandal, Sneaker
5	Ankle boot, Sandal	Pullover, Coat	Sneaker, Trouser	T-shirt/top, Bag	Shirt, Dress
6	Bag, Trouser	Sneaker, T-shirt/top	Shirt, Sandal	Pullover, Coat	Dress, Ankle boot
7	Bag, Sandal	T-shirt/top, Pullover	Trouser, Ankle boot	Sneaker, Dress	Shirt, Coat
8	Bag, Shirt	Ankle boot, T-shirt/top	Pullover, Sandal	Sneaker, Trouser	Coat, Dress
9	Bag, Coat	Sneaker, Pullover	Trouser, Ankle boot	Dress, T-shirt/top	Shirt, Sandal