



The Hebrew University of Jerusalem  
The Rachel and Selim Benin School of Computer Science and Engineering

# **Tiny Memory Experience Replay in Class Incremental Learning**

**Shahar Shaul-Ariel**

Thesis submitted in partial fulfillment of the requirements  
for the Master of Sciences degree  
in Computer Science

Under the supervision of **Prof. Daphna Weinshall**

**September 2024**



האוניברסיטה העברית בירושלים  
בית הספר להנדסה ולמדעי המחשב על שם רחל וסלים בנין  
החוג למדעי המחשב

## Tiny Memory Experience Replay in Class Incremental Learning

זכרון זעיר בלמידת מחלקות מתמשכת

מוגש על ידי  
שחר שאול-אריאל

עבודת גמר לתואר מוסמך במדעי המחשב

עבודה זו הונחתה על ידי  
פרופ' דפנה וינשל

אב תשפ"ד

# תקציר

למידה מתמשכת היא אתגר בלתי פתור, שהרלוונטיות שלו עולה כאשר בוחנים יישומים מודרניים. שלא כמו המוח האנושי, רשתות עצביות עמוקות מאומנות סובלות מתופעה הנקראת שכחה קטסטרופלית, שבה הן מאבדות בהדרגה ידע שנרכש בעבר עם לימוד משימות חדשות. כדי להקל על בעיה זו, פותחו שיטות רבות, רבות מהן מסתמכות על שידור חוזר של דוגמאות מהעבר במהלך אימון משימות חדשות. עם זאת, ככל שהזיכרון שהוקצה לשידור החוזר פוחת, היעילות של גישות אלו פוחתת. מצד שני, שמירה על זיכרון גדול לצורך שידור חוזר היא לא יעילה ולרוב לא מעשית. כאן אנו מציגים את TEAL, גישה חדשנית לאכלוס הזיכרון בדוגמאות, שניתן לשלב בשיטות שונות של שידור חוזר ולשפר משמעותית את הביצועים שלהן על מאגרי זיכרון קטנים. אנו מראים ש-TEAL משפר את הדיוק הממוצע של שיטות 'למידת מחלקות מתמשכת' קיימות, מגיע לביצועים טובים יותר ביחס לאסטריטגיות בחירה אחרות, ומשיג ביצועים מתקדמים אפילו עם מאגרי זיכרון קטנים של 1-3 דוגמאות למחלקה במשימה הסופית. זה מאשש את ההשערה הראשונית שלנו שכאשר הזיכרון דל, עדיף לתעדף את הנתונים האופייניים ביותר.

# Abstract

Continual Learning is an unresolved challenge, whose relevance increases when considering modern applications. Unlike the human brain, trained deep neural networks suffer from a phenomenon called *catastrophic forgetting*, wherein they progressively lose previously acquired knowledge upon learning new tasks. To mitigate this problem, numerous methods have been developed, many relying on the replay of past exemplars during new task training. However, as the memory allocated for replay decreases, the effectiveness of these approaches diminishes. On the other hand, maintaining a large memory for the purpose of replay is inefficient and often impractical. Here we introduce *TEAL*, a novel approach to populate the memory with exemplars, that can be integrated with various experience-replay methods and significantly enhance their performance with small memory buffers. We show that *TEAL* enhances the average accuracy of existing class-incremental methods and outperforms other selection strategies, achieving state-of-the-art performance even with small memory buffers of 1-3 exemplars per class in the final task. This confirms our initial hypothesis that when memory is scarce, it is best to prioritize the most typical data.

# Acknowledgements

I would like to thank my supervisor, Prof. Daphna Weinshall, for her valuable guidance and support throughout my research. Her expertise and insights have been instrumental in shaping this work. I also wish to thank my lab members for their helpful contributions and collaboration.

To Eviatar, my husband — thank you for being my anchor, always. Your love, patience, and support helped me through it all.

I also want to thank my parents, Yedidya and Pnina Ariel, who have encouraged me since childhood to expand my knowledge and pursue an academic path, always believing in my abilities.

Finally, I want to mention my late grandfather, Prof. Michael Schieber, who has always been my inspiration. He survived the Holocaust as a teenager and went on to earn a PhD in physics, despite everything he had endured. I have always admired his motivation for learning and researching, and I hope he is proud of me.

# Contents

|          |  |           |
|----------|--|-----------|
| <b>1</b> | <b>Introduction</b>                                      | <b>1</b>  |
| <b>2</b> | <b>Previous Work</b>                                     | <b>4</b>  |
| 2.1      | Incremental learning . . . . .                           | 4         |
| 2.2      | Selection strategies . . . . .                           | 4         |
| 2.3      | ER-IL methods . . . . .                                  | 5         |
| <b>3</b> | <b>Methods</b>   | <b>6</b>  |
| 3.1      | Problem formulation . . . . .                            | 6         |
| 3.2      | Our method: <i>TEAL</i> . . . . .                        | 6         |
| <b>4</b> | <b>Empirical Evaluation</b>                              | <b>10</b> |
| 4.1      | Methodology . . . . .                                    | 10        |
| 4.1.1    | ER-IL baselines . . . . .                                | 11        |
| 4.1.2    | Selection Strategies baselines . . . . .                 | 12        |
| 4.1.3    | Metrics . . . . .  | 12        |
| 4.1.4    | Datasets . . . . .                                       | 13        |
| 4.1.5    | Architectures . . . . .                                  | 13        |
| 4.2      | Main results . . . . .                                   | 13        |
| 4.2.1    | <i>TEAL</i> integrated with SOTA ER-IL methods . . . . . | 13        |
| 4.2.2    | Comparison of selection strategies . . . . .             | 16        |
| 4.2.3    | Results of Task-IL . . . . .                             | 17        |
| 4.3      | Ablation study . . . . .                                 | 18        |
| 4.3.1    | Iterative process of <i>TEAL</i> . . . . .               | 18        |
| 4.3.2    | Different architecture . . . . .                         | 19        |
| 4.4      | Exploratory analysis and insights . . . . .              | 19        |

|          |   |           |
|----------|---|-----------|
| 4.4.1    | Importance weighting . . . . .                      | 20        |
| 4.4.2    | Integrating <i>TEAL</i> with <i>iCarl</i> . . . . . | 21        |
| <b>5</b> | <b>Discussion and Future Work</b>                   | <b>23</b> |
| 5.1      | <i>Herding</i> vs. <i>TEAL</i> . . . . .            | 23        |
| 5.2      | Strengths . . . . .                                 | 23        |
| 5.3      | Limitations . . . . .                               | 23        |
| <b>A</b> | <b>Additional Empirical Results</b>                 | <b>29</b> |
| <b>B</b> | <b>Implementation details</b>                       | <b>32</b> |
| B.1      | <i>TEAL</i> implementation . . . . .                | 32        |
| B.2      | Stand-alone setting . . . . .                       | 32        |
| B.3      | Integrated setting . . . . .                        | 33        |
| B.4      | Baselines setting . . . . .                         | 33        |
| B.5      | Compute resources . . . . .                         | 33        |

# List of Figures

|     |   |    |
|-----|---|----|
| 1.1 | Illustration of CIL with Experience Replay. . . . .   | 1  |
| 3.1 | Illustration of <i>TEAL</i> 's iterative class selection process. . . . .   | 7  |
| 4.1 | Baseline ER-IL methods comparison on Split CIFAR-100 with different fixed buffer sizes. . . . .                                 | 12 |
| 4.2 | The performance of ER-IL methods with and without <i>TEAL</i> . . . . .   | 14 |
| 4.3 | Performance improvement of <i>TEAL</i> when integrated with <i>XDER</i> over various buffer sizes. . . . .                      | 16 |
| 4.4 | The difference between the improvement obtained by <i>TEAL</i> and the one obtained by <i>Herdling</i> . . . . .                | 17 |
| 4.5 | Split CIFAR-100: <i>TEAL</i> performance gain compared to 4 baseline selection strategies. . . . .                              | 17 |
| 4.6 | Importance weighting results . . . . .  | 20 |
| 4.7 | <i>iCarl</i> comparison results . . . . .   | 22 |
| A.1 | Performance improvement of <i>TEAL</i> when integrated with <i>ER</i> and with <i>ER-ACE</i> over various buffer sizes. . . . . | 29 |
| A.2 | Split CIFAR-100 with 20 tasks instead of 10 . . . . .   | 29 |
| A.3 | The performance of ER-IL methods with and without <i>TEAL</i> . . . . .   | 30 |
| A.4 | Split CIFAR-100: <i>TEAL</i> performance gain compared to 4 baseline selection strategies. . . . .                              | 31 |

# List of Tables

|     |  |    |
|-----|--|----|
| 4.1 | Final average accuracy for method with and without <i>TEAL</i> . . . . .                                 | 15 |
| 4.2 | Task-IL: accuracy with and without <i>TEAL</i> . . . . .   | 18 |
| 4.3 | The improvement of 2 variants of <i>TEAL</i> when integrated with different ER-IL methods . . . . .      | 18 |
| 4.4 | Final average accuracy for method with and without <i>TEAL</i> using the ArchCraft architecture. . . . . | 19 |

# 1 Introduction

With the recent advances in deep neural networks, there has been a growing research interest in incremental learning. The need to integrate new task knowledge into an already trained network has become increasingly important, especially considering the time-consuming nature of training on large datasets. Retraining the network from scratch on both the original and new task data is often impractical, and access to the original training data may be limited or unavailable. In this context, *catastrophic forgetting*, as described by McCloskey and Cohen [25], can be particularly severe.

To address this challenge, various methods have been developed across different frameworks. Van de Ven et al. [31] categorizes these methods into three types: task-incremental, domain-incremental, and class-incremental learning. The fundamental idea behind incremental learning is that a model must sequentially learn tasks, one after the other. Among these approaches, Class-Incremental Learning (CIL) is recognized as the most challenging. Here, each task introduces new classes, and the model must accurately identify the class of each input without access to the corresponding task ID, see Fig. 1.1.

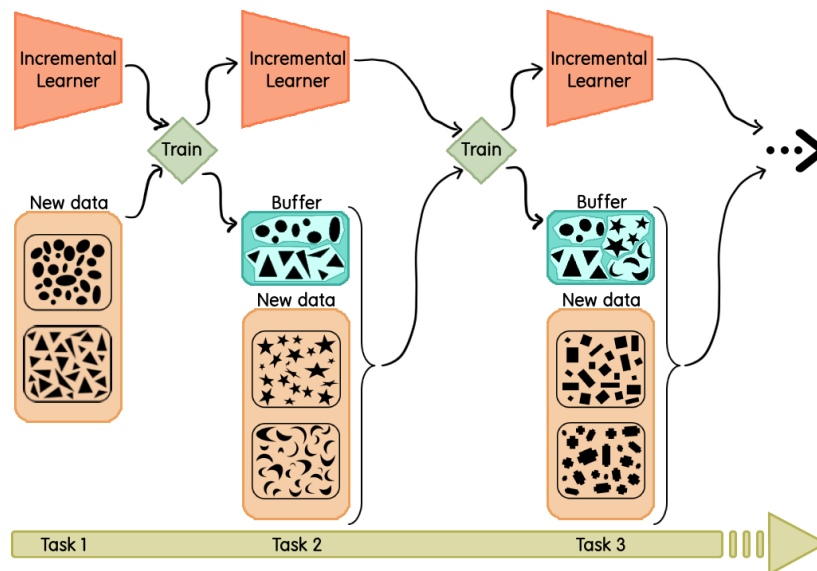


Figure 1.1: Illustration of CIL with Experience Replay.

Incremental learning can be approached in various ways, each with its own assumptions and configurations. In this paper, we follow the constrained setup outlined by De Lange et al. [10], which does not depend on task boundaries during either train-

ing or testing. This approach maintains a fixed memory size throughout incremental training, ensuring it stays within a predefined limit set from the beginning. At each incremental step, new exemplars can be added to the memory buffer only after sufficient space has been vacated.

We focus our attention on a prevalent and rather successful framework called Experience Replay (ER), which involves storing a set of exemplars in memory and reusing them for rehearsal purposes while training on new tasks. Within this framework, several strategies exist for selecting which exemplars to retain in memory. Not surprisingly, the smaller the memory buffer is, the less effective the strategy is at mitigating catastrophic forgetting. This leads us to the following question: Are these strategies necessarily optimal for all memory sizes? In other words, can different strategies be found suitable for different sizes of the memory buffer?

In active learning, it has been shown (empirically and formally) that when the number of labeled examples is small, it is best to choose the most typical examples for training [13]. In the context of CIL and when the memory buffer is too small to truly represent the distribution of each class, we propose to adopt this strategy for the selection of points that are intended to populate the replay memory buffer. In other words, when the buffer is considerably small, it should contain representative exemplars and thus retain a more significant fraction of the previously acquired knowledge.

Our proposed method *TEAL*, Typicality Election Approach to continual Learning, is primarily targeted at scenarios with small buffers. While the use of small buffers may seem too strict, there are situations where maintaining a large memory for replay is simply not feasible. For example, applications running on mobile devices face severe memory constraints and could greatly benefit from methods specifically designed for small buffers.

Accordingly, *TEAL* aims to identify a set of representative exemplars that also exhibit diversity. An exemplar is deemed representative if its likelihood, when considered within the distribution of all points, is high. To ensure diversity, data clustering is leveraged. Central to the success of our approach is the ability to derive an appropriate latent space, where data clustering and likelihood estimation can be reliably obtained. Ideally, any mechanism for memory population can be combined with any compet-

itive Experience-Replay-Incremental-Learning (ER-IL) method in which this mechanism is a separate module. In accordance, we evaluate our approach in Section 4 by considering alternative ER-IL methods, where we replace their native mechanism for buffer population by *TEAL*. The method, methodology and its experimental evaluation are described in the rest of this paper, with emphasis on the enhancements *TEAL* offers to various existing class-incremental methods.

We present *TEAL*, a novel strategy designed to select effective exemplars for small memory buffers. *TEAL* is seamlessly integrated with a number of SOTA replay-based class-incremental learning method, and significantly enhances their performance.

## 2 Previous Work

### 2.1 Incremental learning

Several approaches exist for incremental learning [see 10], including Experience Replay (ER), Generative Replay (GR) [29], Parameter Isolation, and Regularization. Similar to ER, GR replays data from previous tasks during new task training, but uses a generative model to create new samples instead of retaining exemplars seen by the model [9, 11, 12]. Parameter Isolation assigns distinct parameters to each task to reduce forgetting [23] by fixing parameters assigned to previous tasks, while Regularization-based methods [18] incorporate additional terms into the loss function to retain prior knowledge while learning from new data.

These methods differ in the ways they utilize memory: ER stores exemplars, GR stores a generative model, and Parameter Isolation stores task-specific parameters. Regularization-based methods do not rely on memory. This diversity makes it challenging to compare methods directly. In particular, we note that since generative models tend to be very large, GR is hardly suitable to the domain of small memory buffers addressed here. A similar concern can be raised concerning Parameter Isolation methods. Hence, our focus in this paper is on ER.

Note that while few-shot incremental learning [30] may appear similar to our work on CIL with a small memory buffer, there is a significant difference. Specifically, we assume that the data for new tasks is initially sufficient for effective learning, whereas few-shot incremental learning inherently deals with a scarcity of labeled samples.

### 2.2 Selection strategies

When comparing strategies for the population of the memory buffer, Masana et al. [24] demonstrate that the most successful strategies are either random-sampling or *Herding* as defined by Welling [33]. The latter strategy involves retaining a set of exemplars whose mean is closest to the class mean. Another successful strategy presented by Bang et al. [2] leverages classification uncertainty and data augmentation to enhance

the diversity of data instances (*Uncertainty*). Other selection strategies, such as GSS [1] and selecting the exemplars with the highest entropy of the softmax outputs [5], have been shown in previous studies to be inferior to random sampling and *Herding* [24, 26]. Therefore, we do not include them in our comparative empirical evaluation.

Recently, Hacoheh and Tuytelaars [14] showed that models with smaller buffer sizes tend to remember quickly learned examples. To address this issue, they propose a selection strategy called Goldilocks, which retains exemplars learned at an intermediate pace, accommodating various buffer sizes. Since the results presented in this paper are primarily within a task-incremental framework (which yields higher average accuracy) and the paper does not provide code, we do not include this method in our comparisons.

## 2.3 ER-IL methods

There are many ER-IL methods, each utilizing experience replay in a unique way. Some of the key methods we discuss in our work are *XDER* [3], which updates the memory buffer by integrating current information with past memories and preparing the model for new tasks; *ER-ACE* [4], or ER with Asymmetric Cross-Entropy, which applies separate loss functions for new and past tasks; and *BiC* [34], which introduces a bias correction layer to address class imbalance in incremental learning. We also discuss *iCaRL* [28], which combines incremental classifier and representation learning with a nearest-mean-of-exemplars strategy, using *Herding* for exemplar selection. *Herding* is crucial for *iCaRL* as it ensures that the features of the exemplars in the buffer approximate the original class means. Additionally, *GEM* [21] employs Gradient Episodic Memory to prevent gradient interference between new and past tasks through constrained optimization, while *GDumb* [27] uses a greedy sampling strategy for memory storage and retraining the neural network from scratch during inference. Finally, *ER* [7] explores different selection strategies within a simple ER framework.

# 3 Methods

## 3.1 Problem formulation

A CIL problem  $\mathcal{T}$  consists of a sequence of  $T$  tasks. Each task  $t \in T$  contains a set of classes  $C^t = (c^1, \dots, c^{n_t})$  and labeled samples from these classes  $X_t = (X^1, \dots, X^{n_t})$ , where  $X^i = \{(x_1, i), \dots, (x_{m_i}, i)\}$ . The tasks do not share classes, i.e.,  $C^{t_i} \cap C^{t_j} = \emptyset \forall i \neq j$ . We denote by  $N^t$  the total number of classes in all tasks up to task  $t$ :  $N^t = \sum_{i=1}^t n_i$ .

We consider the scenario where there is access to a memory buffer  $\mathcal{M}$ , which has a fixed size throughout the learning of  $\mathcal{T}$ . An incremental learner is a learning model (we consider only deep neural networks) that is trained sequentially on the tasks. Accordingly, the training on task  $t$  is performed on data  $X_t \cup \mathcal{M}$ , where  $\mathcal{M}$  contains stored exemplars from tasks  $1, \dots, t-1$ . At test time, a CIL method is required to classify a given example into its predicted class while considering all previously seen classes.

Our proposed method *TEAL* is described in Section 3.2. It aims to address only one component of the general Incremental Learning (IL) problem, namely, how to populate the memory buffer at the end of each IL iteration. Accordingly, after task  $t$ , *TEAL* should select for each seen class  $n = \frac{|\mathcal{M}|}{N_t}$  exemplars to populate memory buffer  $\mathcal{M}$ .

## 3.2 Our method: *TEAL*

We begin with a definition of point *typicality*. For each exemplar  $x$ , let

$$\text{Typicality}(x) = \left( \frac{1}{K} \sum_{x_i \in K\text{-NN}(x)} \|x - x_i\|_2 \right)^{-1}$$

where  $K$  is a fixed number of nearest neighbors<sup>1</sup>. The essence of *TEAL* is to populate the memory buffer with points that are both *typical* (as captured by the definition above) and *diverse*.

After training on task  $t$ , the incremental learner has to update the memory buffer  $\mathcal{M}$

---

<sup>1</sup>We use  $K = 20$ , but other options yield similar results.

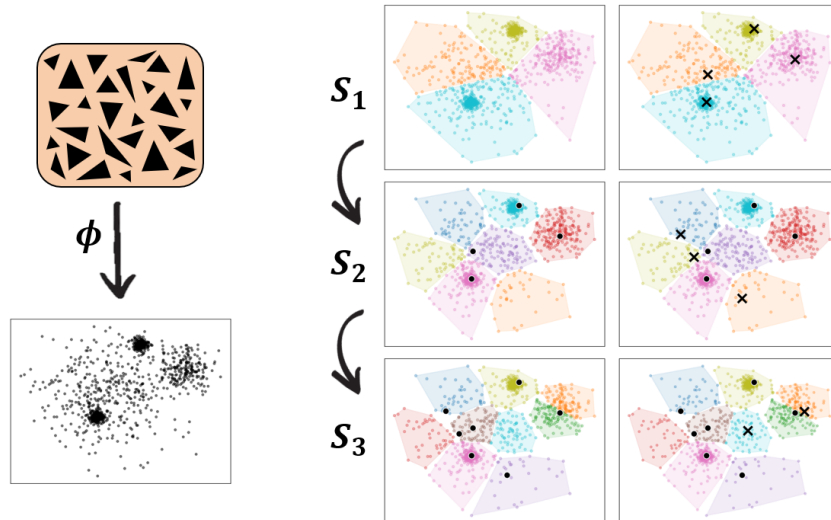


Figure 3.1: Illustration of *TEAL*'s iterative class selection process, which establishes a priority order for the selected set. Initially, an embedding space is generated separately for each class (shown on the left). Samples are then selected iteratively with  $s_1 = 4$ ,  $s_2 = 7$ , and  $s_3 = 9$  (see text for details). Each row on the right panel represents one iteration: the left image displays the  $s_i$  clusters (obtained using K-means) with the previously selected points  $S_{i-1}$  marked with 'o', while the right image shows the updated set  $S_i$ , with newly selected samples marked with 'x'. In the last iteration 3 clusters remain uncovered, but only 2 samples are selected, leaving the red cluster uncovered by  $S_3$ .

with data from the new classes  $\{X^{N^{t-1}+1}, \dots, X^{N^t}\}$ . Given the fixed size of  $\mathcal{M}$ , the learner must first reduce the number of exemplars already in the buffer to accommodate new ones from the new classes. For each class, this involves repeatedly removing examples from the buffer after each IL step.

To address this challenge, *TEAL* maintains a priority list of selected exemplars from each class, which reflects the order in which they should be removed from  $\mathcal{M}$ . This approach allows the learner to remove the least typical points from the buffer when new classes emerge. The list is structured so that the most typical and diverse exemplars, which should be retained as long as possible, are positioned at the top, while those slated for earlier removal are positioned at the bottom.

More specifically, consider a fixed class, and let  $n$  denote the number of exemplars assigned to that class. *TEAL* repeatedly selects a small fraction of  $n$ , generating a sequence of subsets  $S_1 \subseteq S_2 \subseteq \dots \subseteq S_k$ , where  $s_i = |S_i|$  and  $s_1 \leq \dots \leq s_k = n$ . The sizes  $s_i$  define the selection pace, indicating the rate at which exemplars are accumulated.

The inclusion relation induces the priority order: points in  $S_k \setminus S_{k-1}$  are removed first, while points in  $S_1$  are removed last. This process is illustrated in Fig. 3.1.

To initialize the selection process, we need a suitable embedding space for the class exemplars. To this end we train a deep model on all the available training data, and use activations in its penultimate layer as a representation for the new classes. Subsequently, the selection process in the  $i^{\text{th}}$  iteration involves 2 steps (see pseudocode in Alg. 1):

1. **Step 1: Clustering.** In order to ensure diversity, we seek typical exemplars from different regions of the embedding space. When constructing set  $S_i$ , we achieve this by dividing the set of labeled points into  $s_i$  clusters using K-Means [19].
2. **Step 2: Typicality.** We then select the most typical point from each of the  $s_i - s_{i-1}$  largest uncovered clusters, where an uncovered cluster is a cluster from which no point has already been selected.

---

**Algorithm 1** *ConstructExemplarSet*


---

**Input:** a set of exemplars from class  $c$   $X^c$ , number of examples to choose  $n$

**Require:** current feature function  $\phi : \mathcal{X} \rightarrow \mathbb{R}^d$ , iterations pace  $s_1, \dots, s_k$

**Output:** a set of  $n$  exemplars  $\mathcal{M}_c \subseteq X^c$

$X_{emb}^c \leftarrow \phi(X^c)$

$\mathcal{M}_c \leftarrow \emptyset$

**for all**  $i = 1, \dots, k$  **do**

$C_1, \dots, C_{s_i} \leftarrow \text{clustering\_algorithm}(X_{emb}^c, s_i)$  #  $|C_1| \geq \dots \geq |C_{s_i}|$

**for all**  $j = 1, \dots, s_i - s_{i-1}$  **do**

**if**  $C_j$  is uncovered **then**

add  $\text{argmax}_{x \in C_j} \{ \text{Typicality}(x) \}$  to  $\mathcal{M}_c$

**end if**

**end for**

**end for**

**return**  $\mathcal{M}_c$

---

The integration of *TEAL* into a general class-incremental algorithm is described in Alg. 2.

---

**Algorithm 2** *IncrementalTraining*


---

**Input:** training examples  $X_t$

**Require:** current exemplar sets  $\mathcal{M} = (\mathcal{M}_1, \dots, \mathcal{M}_{N^{t-1}})$ , current model parameters  $\theta^{t-1}$

$\theta^t \leftarrow \text{training\_model}(X^t, \mathcal{M}; \theta^{t-1})$

$n \leftarrow \frac{|\mathcal{M}|}{N^t}$

**for**  $c \in \{1, \dots, N^{t-1}\}$  **do**

$\mathcal{M}_c \leftarrow$  first  $n$  exemplars in  $\mathcal{M}_c$

**end for**

**for**  $c \in \{N^{t-1} + 1, \dots, N^t\}$  **do**

$\mathcal{M}_c \leftarrow \text{ConstructExemplarSet}(X^c, n, \theta^t)$

**end for**

$\mathcal{M} \leftarrow (\mathcal{M}_1, \dots, \mathcal{M}_{N^t})$

---

# 4 Empirical Evaluation

We report two settings:

1. *Integrated*: we evaluate the beneficial contribution of *TEAL* to competitive ER-IL methods, by replacing their native selection strategy (*vanilla* version) with *TEAL* (Section 4.2.1).
2. *Alternative selection strategies*: we evaluate the beneficial contribution of *TEAL* as compared to other selection strategies (Section 4.2.2).

In both cases, we incrementally train a deep model with a fixed-size replay buffer and monitor the average accuracy defined below upon completion of each task  $t$ .

## 4.1 Methodology

The majority of the experiments are conducted using the open-source Continual Learning library *Avalanche* [20], with the exception of those involving *XDER* [3], which is not integrated into *Avalanche*. For *XDER*, we utilize the code provided by the authors. In order to guarantee a fair comparison and in all experiments, we employ identical network architectures and maintain consistent experimental conditions.

We make certain in both settings that the buffer remains class-balanced, namely, each update maintains an equal representation of exemplars across all classes<sup>1</sup>. To maintain balance, we follow this procedure: Prior to adding new exemplars we calculate  $n = \frac{|\mathcal{M}|}{N^t}$ , where  $N^t$  is the total number of existing and new classes. Subsequently, we adjust the number of exemplars from existing classes to accommodate  $n$  by removing redundant points, followed by the addition of  $n$  exemplars from each new class.

In the second setting, we employ a simple baseline ER-IL model. This model updates a fixed-size buffer of exemplars after training each task  $t$  and replays it while training on task  $t + 1$ . Additionally, we utilize a weighted data loader to ensure a balanced mix of data from both the new classes and the exemplars stored in the buffer in each batch.

---

<sup>1</sup>When the buffer size is not evenly divisible by the number of classes, there may be classes with an additional exemplar.

Across experiments, the only variation lies in the selection strategy employed: some experiments utilize one of the selection strategies baselines described below, and the remaining experiments employ *TEAL*.

Other than this change in the method’s buffer population mechanism, everything else remains the same.

### 4.1.1 ER-IL baselines

The following ER-IL methods are used: *XDER* [3], *ER-ACE* [4], *ER* [7], *BiC* [34], *iCaRL* [28], *GEM* [21], and *GDumb* [27].

To assess the suitability of these methods under the CIL conditions studied here, we used Split CIFAR-100 and set the buffer size to 300, 500, and 2000. Figure 4.1 presents the results. As clearly illustrated in Fig. 4.1, *GEM* and *GDumb* perform poorly in this scenario, likely due to their unsuitability for a very small memory buffer. Similarly, *BiC* remains competitive only when the buffer size is 2000. While *iCarl* remains competitive across all buffer sizes, it is inherently non-modular, relying heavily on its native *Herding* selection strategy and a K-Means classifier. Hence we cannot integrate it with *TEAL*.

We are therefore left with the methods *XDER*, *ER-ACE*, and *ER*, which are both suitable and competitive for further examination with and without *TEAL* as the selection mechanism. For *ER-ACE*, which requires populating the buffer with exemplars at the start of the task before training on new classes, we cannot rely on the representation provided by a trained model. To address this, we adjust the integration of *TEAL* by filling the buffer in two stages. First, at the beginning of task  $t$ , we populate the buffer with exemplars from the new classes using a random-sampling selection strategy. Second, after completing the training for task  $t$ , we replace the exemplars from the new classes in the buffer with those selected using *TEAL*. It’s important to note that throughout this process, the buffer maintains a fixed size, ensuring compliance with the class-incremental framework.

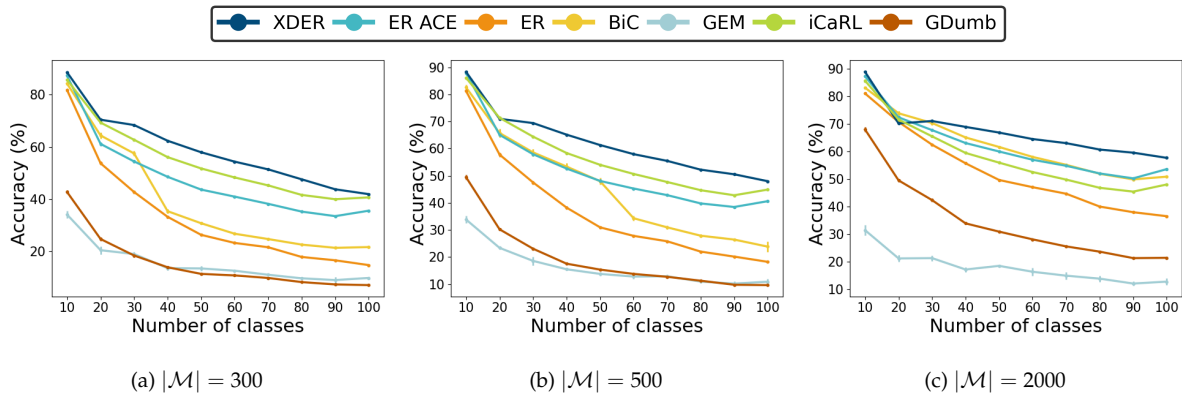


Figure 4.1: Baseline ER-IL methods comparison on Split CIFAR-100 with different fixed buffer sizes. 7 baseline methods are shown, reporting the average accuracy  $A_t$  as new tasks are learned.

### 4.1.2 Selection Strategies baselines

We compare with the following selection strategies (see Sec. 2.2): Random sampling, *Herding*, *Uncertainty*, and *Centered* (closest-to-center) - another common strategy that selects the exemplars closest to the center of all elements in feature space. *Centered* focuses on individual exemplars, unlike *Herding*, which keeps the mean of the whole constructed set close to the center.

### 4.1.3 Metrics

As customary, we employ a score that reflects the accuracy during each stage of the incremental training. Let  $T$  denote the total number of tasks, and  $a_{t,i}$  denote the accuracy of task  $i$  at the end of task  $t$  ( $i \leq t \leq T$ ). For every task  $t = 1, \dots, T$ , define its average accuracy  $A_t = \frac{1}{t} \sum_{i=1}^t a_{t,i}$ . This metric provides a single value for each incremental step, enabling us to directly compare different methods at each step. Since by construction each task has the same number of classes, there is no need for additional weighting terms. We do not include the metric of *forgetting*, which estimates how much the model has forgotten about previous tasks, since it is less suitable for the class-incremental setting: with the addition of new classes, performance inevitably drops across all classes [24].

#### 4.1.4 Datasets

We use several well known Continual Learning benchmarks. (i) **Split CIFAR-100** [8, 28], a dataset created by splitting CIFAR-100 [16] into 10 tasks, each containing 10 different classes of 32X32 images, with 500 images per class for training and 100 for testing. We also use a variation of this dataset by splitting CIFAR-100 into 20 tasks instead of 10, as reported in Fig. A.2. (ii) **Split tinyImageNet**, a dataset created by splitting tinyImageNet [17] into 10 tasks, each containing 20 different classes of 64X64 images, with 500 images per class for training and 50 for testing. (iii) **Split CUB-200** [6], a dataset created by splitting into 20 tasks the CUB-200 [32] high-resolution image classification dataset, consisting of 200 categories of birds, with around 30 images per class for training and 30 for testing.

#### 4.1.5 Architectures

In the first setting, except for one experiment, we employ a ResNet-18 model in all our experiments. The exception is the experiment involving the training of *XDER* on the Split CUB-200 dataset. Due to computational limitations, we used instead a pre-trained ResNet-50 model [15], and this was maintained under all the relevant conditions. Additionally, we run some experiments with another architecture called ArchCraft, a ResNet variant designed for improved performance in CL [22]. In the second setting, we use a smaller version of ResNet-18 [15] as a simple baseline ER-IL model [see 21].

## 4.2 Main results

### 4.2.1 *TEAL* integrated with SOTA ER-IL methods

As mentioned above, we investigate the baseline methods *XDER*, *ER-ACE*, and *ER*, comparing their performance using either their native selection strategies or *TEAL*. We conducted experiments on Split CIFAR-100 with buffer sizes ranging from 100 to 4000, Split tinyImageNet with buffer sizes from 200 to 6000, and Split CUB-200 with

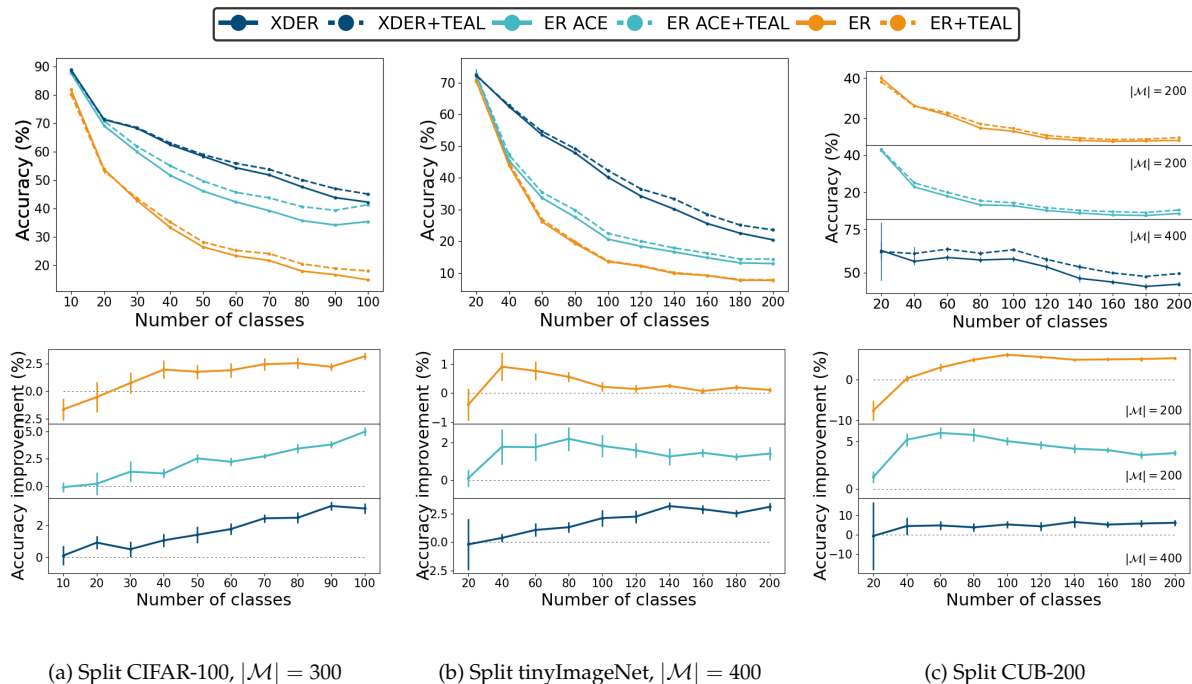


Figure 4.2: The performance of ER-IL methods with and without *TEAL*. First row displays the average accuracy after training incrementally on a different number of classes. Each color corresponds to a different ER-IL method, where the continuous line represents the vanilla method, while the dashed line represents the method with *TEAL* as its selection strategy. The error bars correspond to standard error based on 4-10 repetitions. Second row depicts the difference in accuracy between *TEAL* and another method (*XDER*, *ER-ACE*, and *ER*) across all tasks.

buffer sizes of  $200^2$  and 400. The results for final average accuracy  $A_t$  and the differences between ‘improved’ and ‘vanilla’ are presented in Table 4.1. First row of Fig. 4.2 shows the average accuracy  $A_t$  after each task  $t$ , while the second row illustrates the corresponding improvement from *TEAL* in selected experiments. More figures can be found in Suppl A.

As demonstrated in Fig. 4.2, in most cases, the improvement achieved by *TEAL* increases as incremental training progresses, resulting in a more significant improvement by the final task. Indeed, we do not expect significant differences immediately after the first task, as catastrophic forgetting hasn’t occurred yet. Reassuringly, for the majority of the experiments, the improvement becomes apparent starting from the second task.

<sup>2</sup>Due to computational limitations of the relevant package, we did not run *XDER* with a buffer size of 200.

Table 4.1: Final average accuracy for method with and without *TEAL*.

| Dataset               | $ \mathcal{M} $ | XDER       |               |             | ER-ACE     |               |             | ER         |               |              |
|-----------------------|-----------------|------------|---------------|-------------|------------|---------------|-------------|------------|---------------|--------------|
|                       |                 | Vanilla    | + <i>TEAL</i> | Improvement | Vanilla    | + <i>TEAL</i> | Improvement | Vanilla    | + <i>TEAL</i> | Improvement  |
| Split<br>CIFAR-100    | 100             | 27.98±0.18 | 31.08±0.29    | <b>3.1</b>  | 24.37±0.42 | 31.06±0.27    | <b>6.9</b>  | 10.31±0.03 | 10.24±0.03    | <b>-0.07</b> |
|                       | 300             | 41.97±0.25 | 45.05±0.24    | <b>3.08</b> | 35.37±0.29 | 41.28±0.2     | <b>5.92</b> | 14.82±0.21 | 17.97±0.25    | <b>3.15</b>  |
|                       | 500             | 47.97±0.22 | 50.29±0.2     | <b>2.32</b> | 40.7±0.25  | 45.99±0.33    | <b>5.29</b> | 18.21±0.29 | 22.39±0.31    | <b>4.18</b>  |
|                       | 1000            | 53.69±0.32 | 55.02±0.34    | <b>1.33</b> | 48.16±0.14 | 51.85±0.22    | <b>3.69</b> | 26.66±0.81 | 28.58±0.65    | <b>1.92</b>  |
|                       | 2000            | 57.69±0.21 | 58.79±0.17    | <b>1.1</b>  | 55.99±0.1  | 58.7±0.17     | <b>2.71</b> | 36.54±0.31 | 38.17±0.37    | <b>1.62</b>  |
|                       | 3000            | 58.96±0.25 | 60.09±0.25    | <b>1.13</b> | 59.75±0.25 | 62.3±0.19     | <b>2.55</b> | 43.72±1.63 | 46.03±1.49    | <b>2.32</b>  |
|                       | 4000            | 59.89±0.16 | 60.42±0.28    | <b>0.54</b> | 62.6±0.22  | 64.81±0.2     | <b>2.22</b> | 49.03±0.89 | 52.86±0.62    | <b>3.82</b>  |
| Split<br>tinyImageNet | 200             | 14.6±0.16  | 16.76±0.3     | <b>2.16</b> | 11.3±0.08  | 12.47±0.22    | <b>1.17</b> | 7.95±0.02  | 8.01±0.03     | <b>0.07</b>  |
|                       | 400             | 20.42±0.12 | 23.58±0.18    | <b>3.16</b> | 13.02±0.2  | 14.43±0.29    | <b>1.41</b> | 7.75±0.06  | 7.86±0.05     | <b>0.11</b>  |
|                       | 600             | 25.74±0.09 | 28.71±0.23    | <b>2.97</b> | 14.48±0.13 | 15.7±0.14     | <b>1.22</b> | 7.54±0.02  | 7.8±0.07      | <b>0.26</b>  |
|                       | 1000            | 32.68±0.24 | 34.28±0.24    | <b>1.59</b> | 17.37±0.34 | 18.88±0.21    | <b>1.51</b> | 7.7±0.06   | 7.99±0.06     | <b>0.29</b>  |
|                       | 2000            | 39.76±0.31 | 40.76±0.23    | <b>1.0</b>  | 21.71±0.17 | 24.12±0.24    | <b>2.4</b>  | 8.27±0.08  | 8.82±0.09     | <b>0.56</b>  |
|                       | 4000            | 43.91±0.14 | 44.42±0.26    | <b>0.51</b> | 27.15±0.13 | 29.62±0.18    | <b>2.47</b> | 11.61±0.1  | 13.59±0.21    | <b>1.99</b>  |
|                       | 6000            | 44.71±0.15 | 45.44±0.39    | <b>0.73</b> | 30.3±0.23  | 33.3±0.23     | <b>3</b>    | 16.53±0.15 | 19.07±0.18    | <b>2.54</b>  |
| Split<br>CUB-200      | 200             | –          | –             | –           | 8.56±0.15  | 10.41±0.16    | <b>1.85</b> | 8.98±0.18  | 10.37±0.23    | <b>1.39</b>  |
|                       | 400             | 43.63±1.78 | 49.55±0.47    | <b>5.96</b> | 11.25±0.99 | 12.33±0.26    | <b>1.08</b> | 12.01±0.52 | 14.13±0.3     | <b>2.12</b>  |

Table 4.1 shows that the difference (indicated in the ‘improvement’ row) is consistently positive, demonstrating that *TEAL* always enhances performance. However, it is important to note that when catastrophic forgetting in the vanilla version is too severe, the integration of *TEAL* can no longer mend the damage. Consequently, the enhancements to *ER* and *ER-ACE* on the Split tinyImageNet dataset are mostly notable with larger buffer sizes.

Note that the improvement provided by *TEAL* tends to increase as the buffer size decreases. To further explore the relationship between buffer size and relative improvement, we present a graph showing the final average accuracy improvement across various buffer sizes. The results for Split CIFAR-100 and Split tinyImageNet using *XDER* are illustrated in Fig. 4.3 (see results with *ER-ACE* and *ER* in Fig. A.1). Clearly, while improvements do occur across all buffer sizes, larger improvements are seen for smaller buffers.

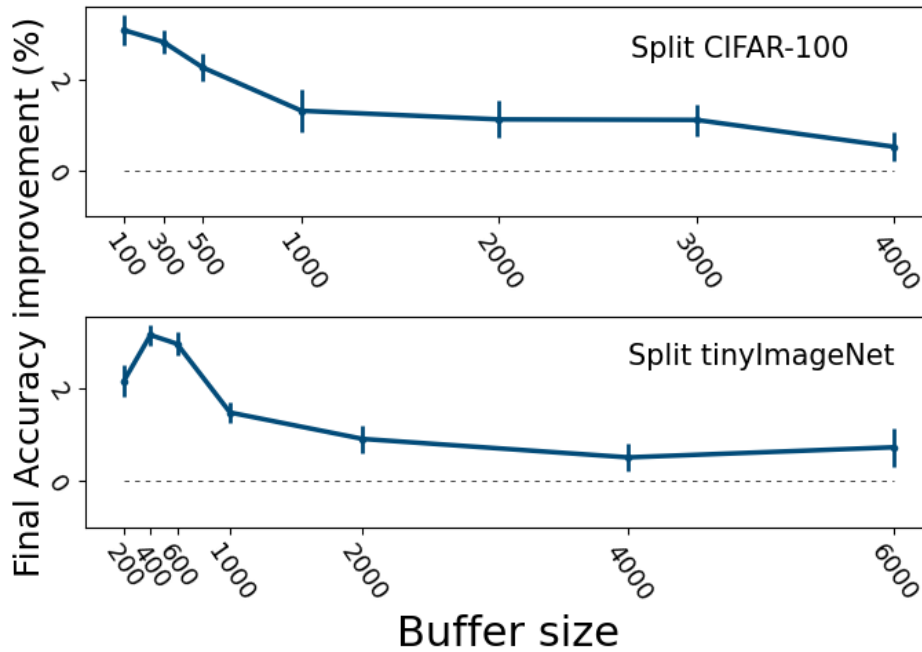


Figure 4.3: Performance improvement of *TEAL* when integrated with *XDER* over various buffer sizes.

## 4.2.2 Comparison of selection strategies

We begin with a simple ER-IL model on Split CIFAR-100<sup>3</sup> with various buffer sizes, and 5 different selection strategies: random, *Herding*, *Uncertainty*, *Centered* and *TEAL*. As depicted in Fig. 4.5, our method enhances performance as compared to *Uncertainty* by 3-4%, and random selection by almost 3% for smaller buffers (300, 400, 500) and about 2.5% for larger buffers (1000, 2000). When compared to *Herding*, our method is more effective with smaller buffer sizes, enhancing performance by up to 1.2% for a buffer size of 300, but less effective with large buffers. It is possible that the mechanism of retaining a set of the nearest neighbors around the average sample in each class becomes less effective when the buffer is too small. Additionally, *Herding* lacks a component for selecting a diverse set, which may lead to the retention of similar exemplars, an issue that becomes more pronounced with smaller buffer sizes. Interestingly, despite its similarity to *Herding*, *Centered* performs worse on larger buffer sizes, resulting in a 1.5-2.5% improvement by *TEAL*. Similar results, using different class ordering, are shown in Fig. A.4.

<sup>3</sup>The order of the classes and the partition to tasks is randomly selected and fixed in all conditions.

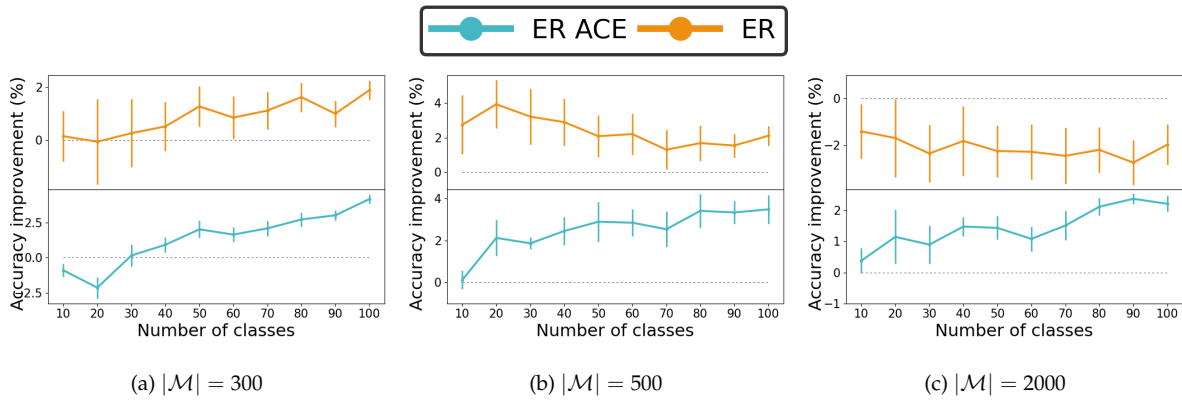


Figure 4.4: The difference between the improvement obtained by *TEAL* and the one obtained by *Herding*, showing 2 ER-IL methods and 3 buffer sizes while training on the Split CIFAR-100 dataset.

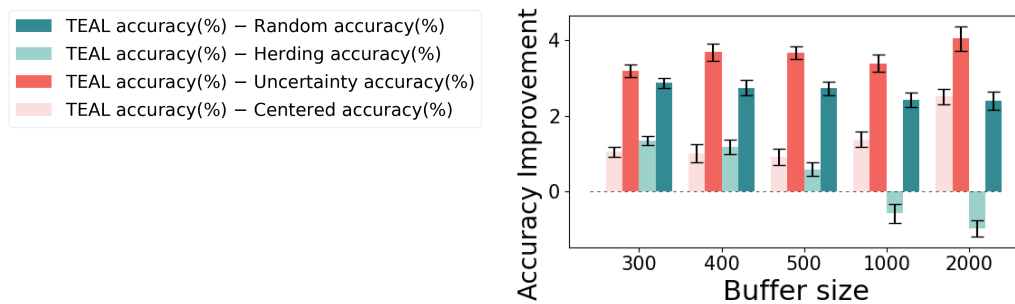


Figure 4.5: Split CIFAR-100: *TEAL* performance gain compared to 4 baselines: random-sampling, *Herding*, *Centered*, and *Uncertainty*, each bar corresponds to a fixed buffer size. The error bars correspond to standard error based on 10-20 repetitions.

Given the results shown in Fig. 4.5, we repeat this comparison while integrating both methods - *TEAL* and *Herding* - into competitive ER-IL methods. Fig. 4.4 shows the difference between the improvement obtained by *TEAL* and the one obtained by *Herding*. Since in almost all cases this difference is positive, we may conclude that *TEAL* is significantly more beneficial than *Herding*.

### 4.2.3 Results of Task-IL

The incorporation of *TEAL* into existing ER-IL method is also beneficial in the task incremental scenario, although the advantage is less pronounced, likely because the scenario is easier. Results are shown in Table 4.2.

Table 4.2: Task-IL: accuracy with and without *TEAL*

| Dataset               | $ \mathcal{M} $ | XDER             | XDER+ <i>TEAL</i> | Improvement |
|-----------------------|-----------------|------------------|-------------------|-------------|
| Split<br>CIFAR-100    | 100             | 73.75 $\pm$ 0.59 | 76.56 $\pm$ 0.27  | <b>2.81</b> |
|                       | 300             | 83.5 $\pm$ 0.06  | 84.46 $\pm$ 0.26  | <b>0.96</b> |
|                       | 500             | 85.96 $\pm$ 0.14 | 86.1 $\pm$ 0.2    | <b>0.14</b> |
| Split<br>tinyImageNet | 200             | 42.48 $\pm$ 0.38 | 46.55 $\pm$ 0.37  | <b>4.07</b> |
|                       | 400             | 55.03 $\pm$ 0.15 | 59.2 $\pm$ 0.47   | <b>4.17</b> |
|                       | 600             | 63.69 $\pm$ 0.21 | 65.26 $\pm$ 0.28  | <b>1.57</b> |
|                       | 1000            | 69.72 $\pm$ 0.31 | 70.61 $\pm$ 0.23  | <b>0.89</b> |

## 4.3 Ablation study

### 4.3.1 Iterative process of *TEAL*

We explore an alternative variant of *TEAL*, which selects the set of exemplars in a single pass instead of by iterations. We call this variant *TEAL.OneTime*. While still selecting exemplars for each class separately, this variant first partitions the exemplars of each class into  $n = \frac{|\mathcal{M}|}{N^t}$  clusters, and then selects the most typical exemplar from each cluster in descending order of cluster size, preserving the order of selections.

Table 4.3: The improvement of 2 variants of *TEAL* when integrated with different ER-IL methods: *XDER*, *ER-ACE*, and *ER*. In almost all cases, the original *TEAL* matches or outperforms the *OneTime* variant.

| $ \mathcal{M} $ | Selection Strategy  | XDER                             | ER-ACE                          | ER                               |
|-----------------|---------------------|----------------------------------|---------------------------------|----------------------------------|
| 300             | <i>TEAL.OneTime</i> | 44.13 $\pm$ 0.26                 | 35.81 $\pm$ 0.42                | 17.64 $\pm$ 0.12                 |
|                 | <i>TEAL</i>         | <b>45.05<math>\pm</math>0.24</b> | <b>36.4<math>\pm</math>0.12</b> | <b>18.1<math>\pm</math>0.1</b>   |
| 500             | <i>TEAL.OneTime</i> | 49.87 $\pm$ 0.22                 | 40.16 $\pm$ 0.3                 | 22.01 $\pm$ 0.18                 |
|                 | <i>TEAL</i>         | <b>50.29<math>\pm</math>0.2</b>  | 40.26 $\pm$ 0.23                | 22.07 $\pm$ 0.13                 |
| 2000            | <i>TEAL.OneTime</i> | 59.05 $\pm$ 0.14                 | 51.52 $\pm$ 0.24                | <b>40.65<math>\pm</math>0.32</b> |
|                 | <i>TEAL</i>         | 58.79 $\pm$ 0.17                 | 51.68 $\pm$ 0.35                | 39.96 $\pm$ 0.23                 |

The experiments on the two variants are conducted using Split CIFAR-100 with the

Table 4.4: Final average accuracy for method with and without *TEAL* using the ArchCraft architecture.

| Dataset               | $ \mathcal{M} $ | ER-ACE           |                  |             | ER               |                  |             |
|-----------------------|-----------------|------------------|------------------|-------------|------------------|------------------|-------------|
|                       |                 | Vanilla          | + <i>TEAL</i>    | Imp.        | Vanilla          | + <i>TEAL</i>    | Imp.        |
| Split<br>CIFAR-100    | 300             | 36.75 $\pm$ 0.28 | 42.47 $\pm$ 0.08 | <b>5.72</b> | 14.59 $\pm$ 0.08 | 18.01 $\pm$ 0.35 | <b>3.42</b> |
|                       | 500             | 41.98 $\pm$ 0.24 | 47.53 $\pm$ 0.35 | <b>5.56</b> | 18.54 $\pm$ 0.4  | 22.3 $\pm$ 0.4   | <b>3.76</b> |
|                       | 2000            | 57.33 $\pm$ 0.18 | 60.18 $\pm$ 0.17 | <b>2.85</b> | 41.03 $\pm$ 0.37 | 43.44 $\pm$ 0.28 | <b>2.41</b> |
| Split<br>tinyImageNet | 400             | 12.24 $\pm$ 0.11 | 13.2 $\pm$ 0.38  | <b>0.96</b> | 7.67 $\pm$ 0.13  | 7.84 $\pm$ 0.19  | <b>0.17</b> |
|                       | 1000            | 15.64 $\pm$ 0.37 | 16.7 $\pm$ 0.2   | <b>1.07</b> | 8.12 $\pm$ 0.05  | 8.36 $\pm$ 0.18  | <b>0.24</b> |
|                       | 2000            | 19.08 $\pm$ 0.36 | 20.43 $\pm$ 0.37 | <b>1.35</b> | 9.65 $\pm$ 0.11  | 10.66 $\pm$ 0.03 | <b>1.01</b> |

smaller version of ResNet-18 mentioned in Section 4.1.5 and buffer sizes of 300, 500 and 2000. We investigate the performance enhancement of each variant when integrated with *XDER*, *ER-ACE*, and *ER*. The results are shown in Table 4.3. Clearly, whenever there is a significant difference in performance, the iterative variant outperforms the other one in 4 out of 5 cases.

### 4.3.2 Different architecture

Some experiments were replicated using ArchCraft instead of the ResNet-18 architecture, obtaining similar results (see Table 4.4).

## 4.4 Exploratory analysis and insights

During the research, we conducted experiments that did not produce results that advanced *TEAL*. In this section, we will discuss the approaches we took and the lessons learned from these outcomes.

### 4.4.1 Importance weighting

*TEAL* focuses on sampling exemplars that are the most typical and diverse, utilizing *Importance Sampling*. However, after sampling, the weights of these exemplars are only used to determine which exemplars to remove from the buffer after each task, not during training; each exemplar then has an equal probability of appearing in the next training session. The only weighting applied is between the new data and the buffer data, as the new data typically outweighs the buffer in quantity. No weighting is applied among the buffer data itself during training.

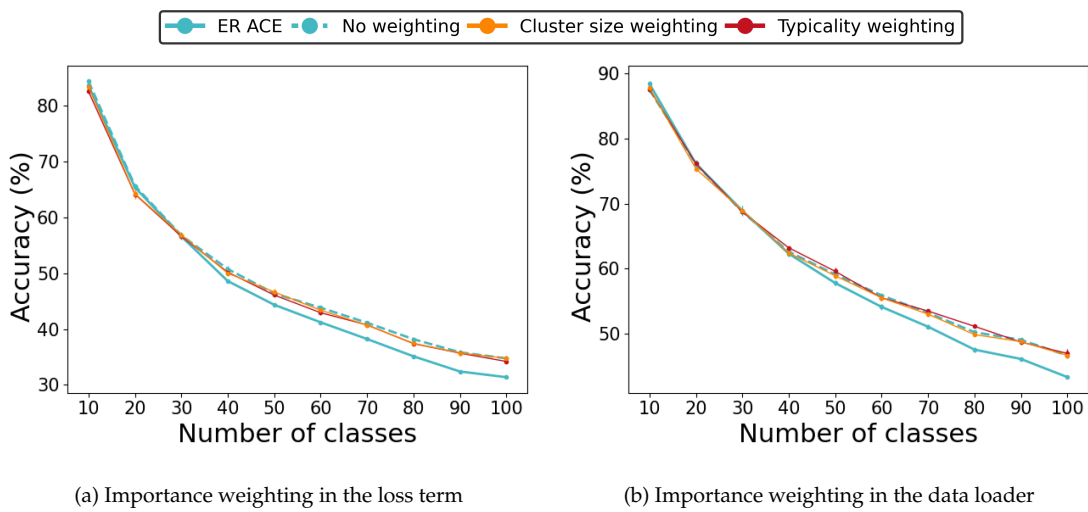


Figure 4.6: Importance weighting results: The experiments are conducted on Split CIFAR-100 with a buffer size of 1000. The solid turquoise line represents the vanilla *ER-ACE*. The other three lines represent *ER-ACE* with *TEAL* as the selection strategy: ‘No weighting’ indicates that no importance weighting is applied (as described in Sec. 4.2.1); ‘Cluster size weighting’ means the weights are based on the size of the cluster to which the exemplars belong at the last iteration of *TEAL*; ‘Typicality weighting’ means the weights are the typicality scores of the exemplars.

In theory, *TEAL* could benefit from incorporating *Importance Weighting* during training on the buffer data, as the exemplars in the buffer are selected with a specific order, making some more typical and diverse than others. To test that we use the ER-IL method *ER-ACE* as a baseline, and we combine *TEAL* with importance weighting in two ways: (i) in the loss term of *ER-ACE*; (ii) in the buffer’s data loader. While *ER-ACE* trains on very small batches of size 10, we conduct the experiments in (i) with a batch size of 512. This larger batch size is necessary because weights in the loss term would

not have a significant impact with such a small batch size. In both (i) and (ii), we set the weights to be either the typicality score of each exemplar or the size of the cluster to which the exemplar belongs at the last iteration of *TEAL*. The weights are normalized per class in the buffer, as *TEAL* selects exemplars from each class separately.

As can be seen in Fig. 4.6, in all the tested scenarios, there is almost no difference between integrating *TEAL* with and without importance weighting. A possible explanation for this is the small size of the buffer, and consequently, the limited number of exemplars per class. The weights are intended to prioritize exemplars within each class in the buffer, but with so few exemplars, their impact is minimal.

#### 4.4.2 Integrating *TEAL* with *iCarl*

To find suitable ER-IL baselines for integrating *TEAL*, we conduct experiments using *iCarl*. As previously mentioned, *iCarl* employs *Herding* as its selection strategy. Surprisingly, despite Fig. 4.5 clearly showing that *TEAL* enhances the accuracy of *Herding* by over 1% in a simple ER-IL model with a buffer of 300 exemplars, switching *iCarl*'s selection strategy to *TEAL* with the same buffer size results in a decrease in accuracy.

*iCarl* utilizes a nearest class mean (NCM) classifier, which represents each class with a vector that is the mean of the class's features. Classification is then performed by identifying the class with the closest vector. *Herding* selects exemplars by constructing a set that preserves the original mean of all features within a specific class, ensuring that the mean of the selected exemplars' features closely matches the original class feature distribution.

Our assumption is that *Herding* is crucial for *iCarl* because the classifier relies on the mean of each class, which *Herding* preserves, while *TEAL* does not. To test this, we also run a variation of *iCarl* that uses a linear classifier instead of an NCM classifier. The results for both variations of *iCarl* with and without *TEAL* as their selection strategy are shown in Fig. 4.7. As can be seen, using a linear classifier with *iCarl* results in lower accuracy, indicating that the NCM classifier is a key component of the method. While substituting *Herding* with *TEAL* is not beneficial in either variation, the accuracy degradation is less pronounced in the final few tasks when using a linear classifier.

Based on these results and our understanding of *iCarl*'s algorithm, we believe that *iCarl* is highly fine-tuned to work specifically with *Herding*, making it unsuitable for integration with other selection strategies like *TEAL*.

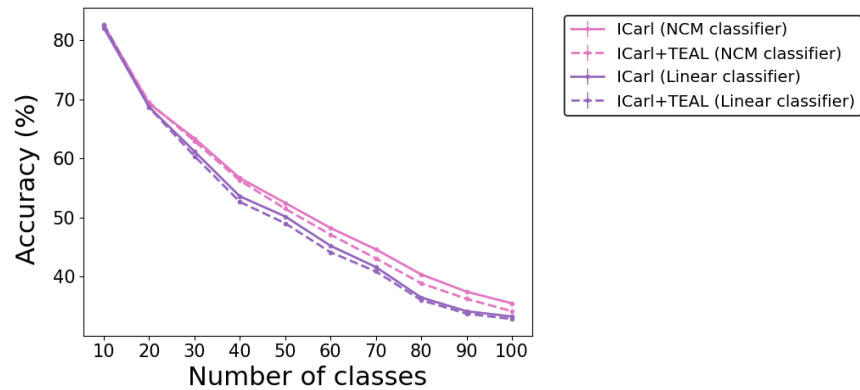


Figure 4.7: *iCarl* comparison results: The experiments are conducted on Split CIFAR-100 with a buffer size of 300. The pink lines represent *iCarl* with its original nearest class mean (NCM) classifier, while the purple lines represent *iCarl* with a linear classifier. Solid lines indicate *iCarl* using *Herding* as the selection strategy, while dashed lines indicate the use of *TEAL* as the selection strategy.

# 5 Discussion and Future Work

## 5.1 *Herding* vs. *TEAL*

As shown in Figs. 4.5, A.4a and A.4b, *TEAL* begins to outperform *Herding* at a buffer size somewhere between 500 and 1000. This trend is also observed with the vanilla *ER* method in Fig. 4.4. However, it remains unclear why *Herding* consistently outperforms *TEAL* across three different class orderings after reaching a certain buffer size, while the other selection strategies do not. Most notably, *Centered*, which selects exemplars in a manner very similar to *Herding*, does not exhibit the same behavior. This phenomenon is left for future work to investigate.

## 5.2 Strengths

We proposed a new mechanism to select exemplars for the memory buffer in replay-based CIL methods. This method is based on the principles of diversity and representativeness. When the memory buffer is relatively small, our method *TEAL* is shown to outperform both the native mechanism of each ER-IL method (usually random selection) and alternative selection mechanisms, including *Herding*, *Uncertainty* and *Centered*. Even when the buffer is large, our method is beneficial in almost all cases.

*TEAL* can be easily integrated with various ER-IL methods by simply replacing the existing exemplar selection mechanism with *TEAL*'s approach.

## 5.3 Limitations

*TEAL* heavily relies on a meaningful representation, as both the clustering and typicality calculations are performed in an embedding space. If this space does not accurately reflect the data distribution, the results may not be useful. Consequently, *TEAL* is not recommended when there is insufficient data to learn an effective representation. Additionally, even with adequate data, if there are limitations on the number of times each data point is introduced during training (as in online learning), it can

be challenging to develop a strong representation, potentially reducing the effectiveness of *TEAL*. There may be a solution to this issue, perhaps by employing a different approach to achieve the representation, but this has not been explored in the current research.

Another limitation is related to the ER-IL method with which *TEAL* is integrated. If the underlying method does not manage catastrophic forgetting effectively and experiences significant forgetting during incremental training, the integration of *TEAL* may not substantially improve performance. For *TEAL* to contribute effectively, there must be some retained data in the process to ensure that the exemplars selected by *TEAL* have a chance to aid in retention and enhance the method's ability to remember more. This limitation is exemplified in Table 4.1, specifically the results for the *ER* method on Split CIFAR-100 with a buffer size of 100 and on Split tinyImageNet with buffer sizes ranging from 200 to 2000. The performance of the vanilla method in these experiments is notably low, as indicated by the average accuracy at the final task, which suggests that the incremental learner nearly randomizes its guesses for previous classes, reflecting significant forgetting. Consequently, the performance with *TEAL* does not show substantial improvement compared to the baseline.

# Bibliography

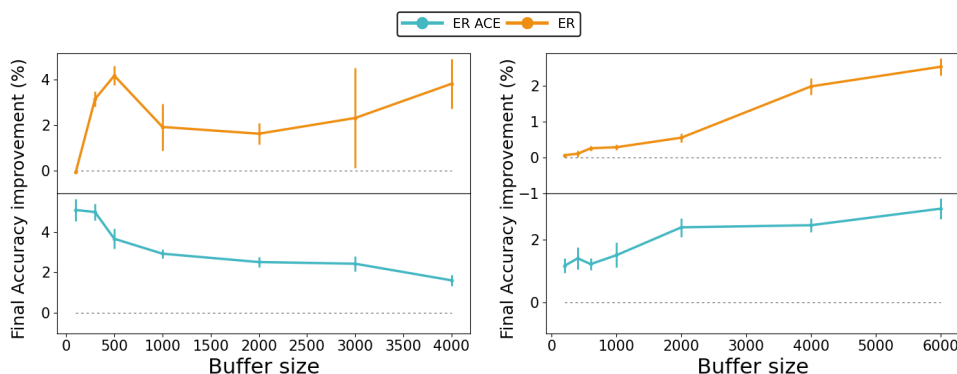
- [1] Aljundi, R., Lin, M., Goujaud, B., and Bengio, Y. (2019). Gradient based sample selection for online continual learning. *Advances in neural information processing systems*, 32.
- [2] Bang, J., Kim, H., Yoo, Y., Ha, J.-W., and Choi, J. (2021). Rainbow memory: Continual learning with a memory of diverse samples. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8218–8227.
- [3] Boschini, M., Bonicelli, L., Buzzega, P., Porrello, A., and Calderara, S. (2022). Class-incremental continual learning into the extended der-verse. *IEEE transactions on pattern analysis and machine intelligence*, 45(5):5497–5512.
- [4] Caccia, L., Aljundi, R., Asadi, N., Tuytelaars, T., Pineau, J., and Belilovsky, E. (2021). New insights on reducing abrupt representation change in online continual learning. *arXiv preprint arXiv:2104.05025*.
- [5] Chaudhry, A., Dokania, P. K., Ajanthan, T., and Torr, P. H. (2018a). Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European conference on computer vision (ECCV)*, pages 532–547.
- [6] Chaudhry, A., Ranzato, M., Rohrbach, M., and Elhoseiny, M. (2018b). Efficient lifelong learning with A-GEM. *CoRR*, abs/1812.00420.
- [7] Chaudhry, A., Rohrbach, M., Elhoseiny, M., Ajanthan, T., Dokania, P., Torr, P., and Ranzato, M. (2019a). Continual learning with tiny episodic memories. In *Workshop on Multi-Task and Lifelong Reinforcement Learning*.
- [8] Chaudhry, A., Rohrbach, M., Elhoseiny, M., Ajanthan, T., Dokania, P. K., Torr, P. H. S., and Ranzato, M. (2019b). Continual learning with tiny episodic memories. *CoRR*, abs/1902.10486.
- [9] Choi, Y., El-Khamy, M., and Lee, J. (2021). Dual-teacher class-incremental learning with data-free generative replay. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3543–3552.

- [10] De Lange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., and Tuytelaars, T. (2021). A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385.
- [11] Gao, R. and Liu, W. (2023). Ddgr: Continual learning with deep diffusion-based generative replay. In *International Conference on Machine Learning*, pages 10744–10763. PMLR.
- [12] Gautam, C., Parameswaran, S., Mishra, A., and Sundaram, S. (2024). Generative replay-based continual zero-shot learning. In *Towards Human Brain Inspired Lifelong Learning*, pages 73–100. World Scientific.
- [13] Hacothen, G., Dekel, A., and Weinshall, D. (2022). Active learning on a budget: Opposite strategies suit high and low budgets. In *International Conference on Machine Learning*. PMLR.
- [14] Hacothen, G. and Tuytelaars, T. (2024). Forgetting order of continual learning: Examples that are learned first are forgotten last. *arXiv preprint arXiv:2406.09935*.
- [15] He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society.
- [16] Krizhevsky, A., Hinton, G., et al. (2009). Learning multiple layers of features from tiny images. *Online*.
- [17] Le, Y. and Yang, X. S. (2015). Tiny imagenet visual recognition challenge.
- [18] Li, Z. and Hoiem, D. (2017). Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947.
- [19] Lloyd, S. (1982). Least squares quantization in pcm. *IEEE transactions on information theory*, 28(2):129–137.
- [20] Lomonaco, V., Pellegrini, L., Cossu, A., Carta, A., Graffieti, G., Hayes, T. L., De Lange, M., Masana, M., Pomponi, J., Van de Ven, G. M., et al. (2021). Avalanche:

- an end-to-end library for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3600–3610.
- [21] Lopez-Paz, D. and Ranzato, M. (2017). Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30.
- [22] Lu, A., Feng, T., Yuan, H., Song, X., and Sun, Y. (2024). Revisiting neural networks for continual learning: An architectural perspective. *arXiv preprint arXiv:2404.14829*.
- [23] Mallya, A. and Lazechnik, S. (2018). Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7765–7773.
- [24] Masana, M., Liu, X., Twardowski, B., Menta, M., Bagdanov, A. D., and Van De Weijer, J. (2022). Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533.
- [25] McCloskey, M. and Cohen, N. J. (1989). Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier.
- [26] Prabhu, A., Al Kader Hammoud, H. A., Dokania, P. K., Torr, P. H., Lim, S.-N., Ghanem, B., and Bibi, A. (2023). Computationally budgeted continual learning: What does matter? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3698–3707.
- [27] Prabhu, A., Torr, P. H., and Dokania, P. K. (2020). Gdumb: A simple approach that questions our progress in continual learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 524–540. Springer.
- [28] Rebuffi, S.-A., Kolesnikov, A., Sperl, G., and Lampert, C. H. (2017). icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010.
- [29] Shin, H., Lee, J. K., Kim, J., and Kim, J. (2017). Continual learning with deep generative replay. *Advances in neural information processing systems*, 30.

- [30] Tian, S., Li, L., Li, W., Ran, H., Ning, X., and Tiwari, P. (2024). A survey on few-shot class-incremental learning. *Neural Networks*, 169:307–324.
- [31] Van de Ven, G. M., Tuytelaars, T., and Tolias, A. S. (2022). Three types of incremental learning. *Nature Machine Intelligence*, 4(12):1185–1197.
- [32] Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. (2011). The caltech-ucsd birds-200-2011 dataset.
- [33] Welling, M. (2009). Herding dynamical weights to learn. In *Proceedings of the 26th annual international conference on machine learning*, pages 1121–1128.
- [34] Wu, Y., Chen, Y., Wang, L., Ye, Y., Liu, Z., Guo, Y., and Fu, Y. (2019). Large scale incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 374–382.

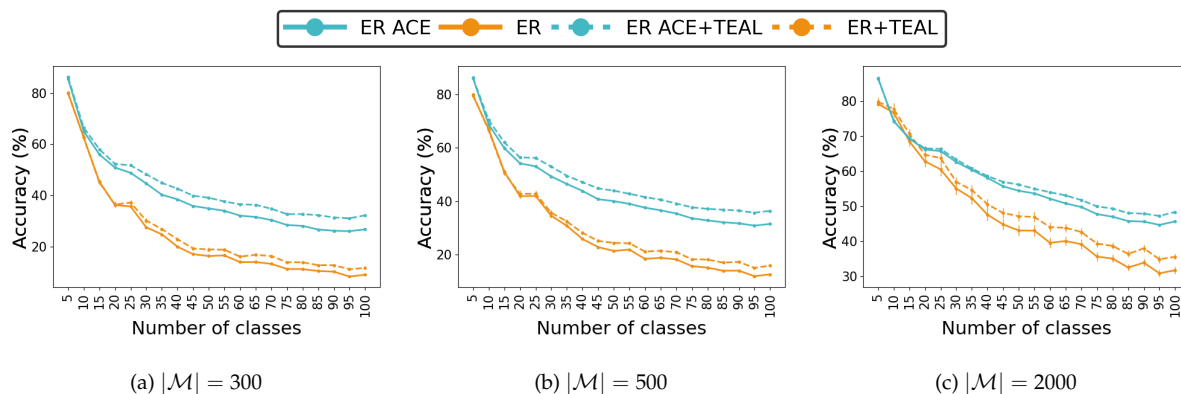
# A Additional Empirical Results



(a) Split CIFAR-100

(b) Split tinyImageNet

Figure A.1: Performance improvement of *TEAL* when integrated with *ER* and with *ER-ACE* over various buffer sizes.

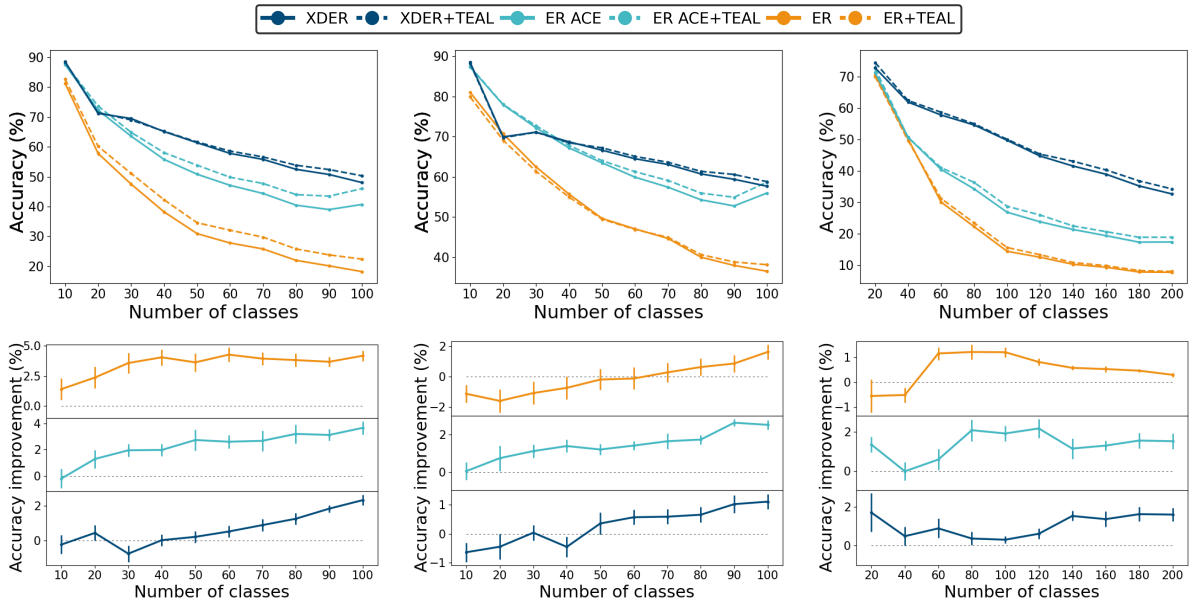


(a)  $|\mathcal{M}| = 300$

(b)  $|\mathcal{M}| = 500$

(c)  $|\mathcal{M}| = 2000$

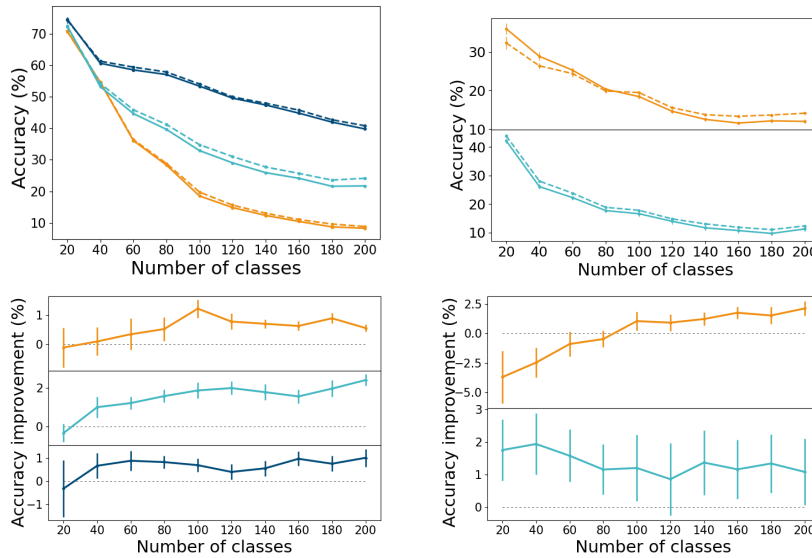
Figure A.2: Split CIFAR-100: The performance of ER-IL methods with and without *TEAL* with splitting CIFAR-100 into 20 tasks instead of 10. Each color corresponds to a different ER-IL method, where the continuous line represents the vanilla method, while the dashed line represents the method with *TEAL* as its selection strategy. The error bars correspond to standard error based on 6 repetitions.



(a) Split CIFAR-100,  $|\mathcal{M}| = 500$

(b) Split CIFAR-100,  $|\mathcal{M}| = 2000$

(c) Split tinyImageNet,  $|\mathcal{M}| = 1000$



(d) Split tinyImageNet,  $|\mathcal{M}| = 2000$

(e) Split CUB-200,  $|\mathcal{M}| = 400$

Figure A.3: The performance of ER-IL methods with and without *TEAL*. First row displays the average accuracy after training incrementally on a different number of classes. Each color corresponds to a different ER-IL method, where the continuous line represents the vanilla method, while the dashed line represents the method with *TEAL* as its selection strategy. The error bars correspond to standard error based on 4-10 repetitions. Second row depicts the difference in accuracy between *TEAL* and another method (*XDER*, *ER-ACE*, and *ER*) across all tasks.

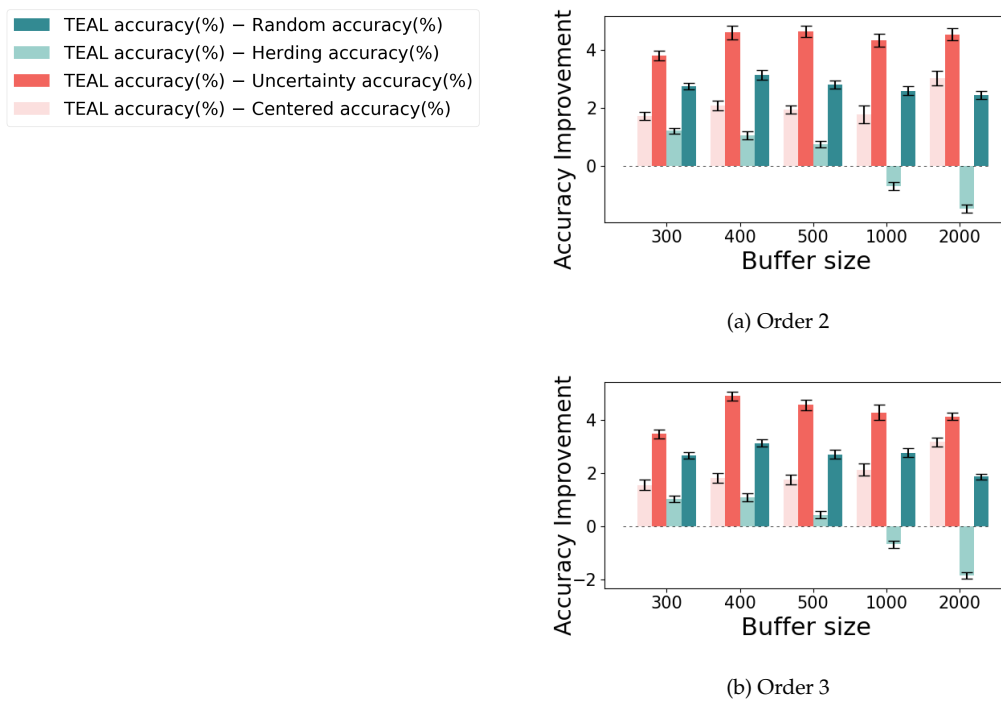


Figure A.4: Split CIFAR-100: *TEAL* Enhanced accuracy of *TEAL* compared to 4 baselines: random-sampling, *Herding*, *Centered*, and *Uncertainty*, in two different classes order.

# B Implementation details

## B.1 TEAL implementation

**Clustering algorithm** We used scikit-learn KMeans implementation.

**Iterations pace** *TEAL*, as described in Alg. 1, requires an iterations pace  $s_1 \leq \dots \leq s_k = n$  which indicates the pace of selecting exemplars from new class data. In both settings we use a logarithmic pace. We set a base  $b = 1.4$ , and define  $s_1 = \lfloor b^4 \rfloor, s_2 = \lfloor b^5 \rfloor, \dots, s_k = \lfloor b^{\log_b n} \rfloor$ .

## B.2 Stand-alone setting

For all selection strategies, we use a smaller ResNet-18, as mentioned above, trained for 200 epochs. Optimization is performed with an SGD optimizer using Nesterov momentum of 0.9, a weight decay of 0.0002, and a learning rate that starts at 0.1 and decays by a factor of 0.3 every 66 epochs. Training is conducted with a batch size of 128 examples, and data augmentation is applied through random cropping and horizontal flips.

We run this setting on Split CIFAR-100 with three different random orders of classes:

1. The order we display on Fig. 4.5 is: [44, 19, 93, 90, 71, 69, 37, 95, 53, 91, 81, 42, 80, 85, 74, 56, 76, 63, 82, 40, 26, 92, 57, 10, 16, 66, 89, 41, 97, 8, 31, 24, 35, 30, 65, 7, 98, 23, 20, 29, 78, 61, 94, 15, 4, 52, 59, 5, 54, 46, 3, 28, 2, 70, 6, 60, 49, 68, 55, 72, 79, 77, 45, 1, 32, 34, 11, 0, 22, 12, 87, 50, 25, 47, 36, 96, 9, 83, 62, 84, 18, 17, 75, 67, 13, 48, 39, 21, 64, 88, 38, 27, 14, 73, 33, 58, 86, 43, 99, 51]
2. The order in Fig. A.4a is: [45, 15, 90, 32, 35, 63, 17, 72, 79, 96, 48, 36, 16, 11, 23, 80, 22, 58, 3, 62, 50, 33, 66, 99, 43, 76, 7, 57, 81, 82, 6, 10, 24, 52, 95, 73, 91, 21, 38, 31, 85, 59, 13, 69, 75, 70, 64, 8, 77, 34, 46, 39, 92, 0, 44, 98, 49, 9, 4, 61, 12, 83, 28, 78, 40, 88, 54, 5, 26, 41, 89, 20, 84, 2, 1, 55, 19, 74, 25, 37, 42, 14, 30, 18, 67, 71, 68, 27, 60, 51, 29, 56, 93, 47, 97, 94, 86, 87, 65, 53]

3. The order in Fig. A.4b is: [48, 97, 1, 81, 90, 49, 10, 8, 7, 20, 70, 73, 75, 14, 91, 38, 47, 21, 74, 52, 80, 98, 59, 12, 71, 85, 6, 34, 55, 82, 95, 63, 78, 15, 94, 60, 99, 76, 25, 40, 88, 0, 62, 96, 87, 51, 16, 18, 9, 19, 29, 45, 86, 53, 56, 31, 28, 61, 30, 33, 4, 67, 64, 58, 50, 54, 3, 13, 37, 27, 66, 77, 84, 69, 2, 41, 22, 92, 42, 44, 11, 36, 46, 79, 65, 72, 23, 17, 39, 5, 89, 35, 24, 83, 43, 57, 93, 32, 68, 26]

### B.3 Integrated setting

Here we train a ResNet-18 model for 100 epochs using the same optimizer as in B.2, with the same batch size and data augmentations, with some exceptions. *ER-ACE* starts with a learning rate of 0.01 and train on batch size of 10 examples as in the original paper, and all experiments on Split CUB-200 are conducted for 30 epochs with a batch size of 16 due to the dataset size and the resolution of its images. In experiments using the ArchCraft model, we used the ResAC-A model as implemented in the official paper’s code.

### B.4 Baselines setting

The setting for the experiments of the baseline methods is the same as in B.3 with some exceptions. For *BiC*, the number of training epochs is 250, and the learning rate scheduler decays the learning rate by a factor of 0.1 on epochs 100, 150 and 200. For *GEM*, the batch size is 32 and the learning rate starts from 0.03. For *GDumb*, the batch size is 32. For *iCaRL* the learning rate starts from 2, the weight decay is 0.00001 and the learning rate scheduler decays the learning rate by a factor of 0.2 on epochs 49 and 63.

### B.5 Compute resources

All experiments involved training deep learning models, necessitating the use of GPUs. For the Split CIFAR-100 experiments, 10 GB of GPU memory was used. The Split tiny-ImageNet experiments required 22 GB of GPU memory, while the Split CUB-200 experiments utilized 45 GB of GPU memory. Any other experiments conducted for the

full research and not reported in the paper required the same compute resources.