

# LEVERAGING SELF-SUPERVISED AND SUPERVISED EMBEDDINGS FOR MEMORY-EFFICIENT EXPERIENCE-REPLAY CONTINUAL LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Catastrophic forgetting remains a key challenge in Continual Learning (CL). In replay-based CL with severe memory constraints, performance critically depends on the sample selection strategy - that is, which examples are stored for replay. Most existing approaches construct memory buffers using embeddings learned under supervised objectives. However, class-agnostic, self-supervised representations often encode rich, class-relevant semantics that are overlooked. We propose a new method, *MERS- Multiple Embedding Replay Selection*, which replaces the buffer selection module with a graph-based approach that integrates both supervised and self-supervised embeddings. Empirical results show consistent improvements over state-of-the-art selection strategies across a range of continual learning algorithms, with particularly strong gains in low-memory regimes. On CIFAR-100 and TinyImageNet, *MERS* outperforms single-embedding baselines without adding model parameters or increasing replay volume, making it a practical, drop-in enhancement for replay-based continual learning.

## 1 INTRODUCTION

*Continual Learning* (CL) focuses on acquiring knowledge from a continuous stream of data, where the distribution of information can change over time. Unlike traditional machine learning, which relies on static datasets and assumes that the data distribution remains fixed, many real-world scenarios involve environments that evolve over time, such as autonomous driving.

A core challenge in CL is *catastrophic forgetting* [18], the tendency of neural networks to lose previously acquired knowledge when trained on new tasks. Without mechanisms to preserve past information, models quickly forget earlier concepts, leading to degraded performance over time. This issue is especially pronounced in the *class-incremental learning* (CIL) setting, a particularly challenging variant of continual learning. In CIL, each task introduces entirely new classes, and the model must learn to classify all new and previously seen classes jointly.

Owing to this difficulty, *replay-based methods* have emerged as an effective approach to mitigate catastrophic forgetting in CIL. These methods maintain a small buffer of selected past examples and replay them during training on new tasks. However, when the buffer is small, as is often the case, the strategy for selecting the examples to retain becomes critical to overall performance. Most existing selection strategies rely on representations learned from supervised models, which may fail to capture the full diversity or structure of the data.

We propose a novel algorithm that unifies multiple representation spaces, each capturing distinct aspects of the data, within a graph-based selection framework (see Fig. 1). By leveraging nonparametric density estimation and localized coverage strategies, our method selects examples that offer better coverage and diversity across all embedding spaces. This multi-embedding approach can be integrated into existing selection algorithms (see review of related work below). It adapts their hyperparameters in a data-driven manner to the geometry of each embedding space.

We demonstrate that our method consistently outperforms single-embedding baselines across several continual learning scenarios and datasets. Notably, our improvements are most pronounced in low-buffer regimes, where efficient use of memory is crucial.

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

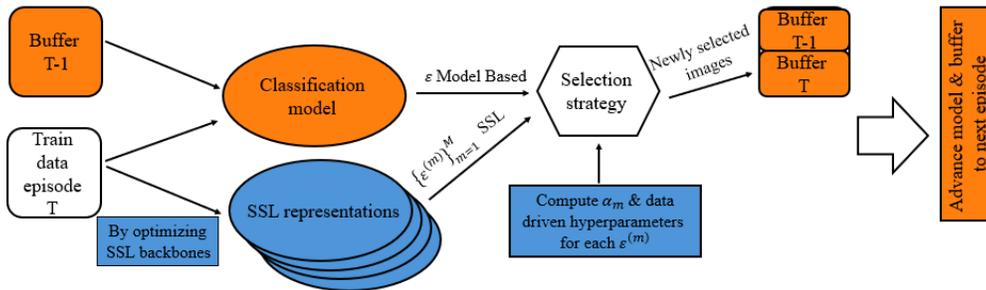


Figure 1: Illustration of our *MERS* in the class-incremental learning (CIL) setup: Training data is split into disjoint episodes, where at each episode only the current data is accessible and becomes unavailable afterward. Orange components denote the standard CIL process: the classification model and buffer are retained and advanced to the next episode. Blue components highlight our contribution, *MERS*, which augments CIL by training a self-supervised model from scratch at each episode, extracting its embeddings, and using them (together with the supervised model outputs) to guide buffer updates through a data-driven selection strategy.

## 2 RELATED WORK

**Continual learning paradigms** CL approaches are often grouped into (i) regularization-based methods that constrain parameter updates to preserve prior knowledge (e.g., EWC [15], LwF [17]), (ii) architecture-based methods that expand capacity across tasks (e.g. HAT [22], DAN [28]) and (iii) Replay-based methods that maintain a small memory of exemplars for replay (e.g., ER [21], ER-ACE [7]). In CIL, rehearsal is particularly competitive under tight memory budgets because it is able to preserve decision boundaries as the label set grows [14].

**Selection strategies** A central problem for Replay-based methods is *exemplar selection*. The iCaRL method employed a *Herding* algorithm to approximate each class centroid within a fixed feature space [14]. More recent strategies fall into two complementary families: (i) Gradient or conflict-oriented approaches, such as GSS [2], that prioritize samples whose loss gradients would be most altered by the impending parameter update, thereby directly mitigating catastrophic interference. (ii) Representativeness-oriented approaches, such as TEAL [23], that keep only the samples with the highest typicality, the inverse of their mean distance to the  $K$  nearest neighbors, ensuring that the buffer stores the most representative points.

**Coverage-based selection and its guarantees** *ProbCover* formalizes tiny-buffer selection as covering a  $k$ -NN graph and offers a  $(1 - 1/e)$  greedy guarantee [27]. *MaxHerding* smooths this objective with kernels, retains submodularity, and is robust to hyperparameters [3].

Prior Continual Learning heuristics, such as iCaRL’s Herding-based exemplar selection [20], and Rainbow Memory’s diversity through uncertainty strategy [5], compute coverage in a single embedding at a fixed scale, rendering them brittle to task heterogeneity. Our method introduces the notion of coverage to this line of work, expands coverage to multiple embeddings, and *adapts* locality per embedding using nonparametric statistics, which we find crucial in tiny-buffer regimes.

**Self-supervised representations for CL** Self-supervised learning (SSL) captures class-agnostic invariances that naturally complement supervised features [24]. Contrastive methods such as SimCLR [11] and redundancy-reduction objectives like VICReg [6] yield rich embeddings without label supervision, while teacher-student paradigms like DINO [8] tighten view consistency.

These SSL representations have already demonstrated effective transfer to object detection, semantic segmentation, depth estimation, robotics manipulation, and few-shot recognition, often rivaling or surpassing supervised pretraining [24]. Yet most rehearsal-based CIL methods still choose exemplars solely in the *supervised* feature space of the current classifier, with only a handful operating purely in an SSL space, as mention in the Selection strategies part.

We instead *jointly exploit supervised and SSL embeddings*, ensuring that the memory preserves both class-discriminative and class-agnostic geometry. This dual-space strategy harnesses complementary signals and, as our experiments confirm, delivers consistent gains in the tiny-buffer continual-learning regime.

**Multi-view learning** This is an ML paradigm where data is represented through multiple distinct feature sets or "views" (e.g., text and image) [29]. Common approaches include co-training and multi-view representation learning [30]. The central idea is to leverage the complementary information in these views to improve performance, often by enforcing consistency or agreement across them. In contrast, our approach aims to exploit variability among representation in order to achieve a more representative set of examples, rather than achieving a single coherent view of the data.

### 3 OUR METHOD: *MERS*

The proposed method (see Fig. 1), called *Multi-Embedding Replay Selection (MERS)*, is intended to enhance any replay-based approach within the CIL (Class Incremental Learning) framework. It involves 2 main steps: (i) replace the buffer selection method with a graph-based method; (ii) expand the graph-based method to integrate a supervised and self-supervised embeddings. The method's primary advantage is expected to emerge in low memory buffer scenarios.

The optimization problem, which lies at the heart of the new method, can be shown to be a known variant of the  $k$ -coverage problem, which is defined as follows:

**Definition 1** (Element-Weighted Maximum  $k$ -Coverage with two groups). *Let  $U$  be a universe of elements, partitioned into two disjoint subsets  $U^1$  and  $U^2$  such that  $U = U^1 \cup U^2$  and  $U^1 \cap U^2 = \emptyset$ . Each element  $e \in U^i$  is associated with a nonnegative weight  $\alpha_i(e) \in \mathbb{R}_{\geq 0}$ , where the weight functions  $\alpha_1, \alpha_2$  may differ between the two groups.*

*Let  $\mathcal{S} = \{S_1, S_2, \dots, S_l\}$  be a family of subsets of  $U$ , and let  $k \in \mathbb{N}$  be a budget parameter. For a subcollection  $\mathcal{A} \subseteq \mathcal{S}$ , define the coverage weight as*

$$\text{Coverage}(\mathcal{A}) = \sum_{e \in \bigcup_{S \in \mathcal{A}} S \cap U^1} \alpha_1(e) + \sum_{e \in \bigcup_{S \in \mathcal{A}} S \cap U^2} \alpha_2(e). \quad (1)$$

*The goal is to select a subcollection  $\mathcal{A} \subseteq \mathcal{S}$  of size at most  $k$  that maximizes  $\text{Coverage}(\mathcal{A})$ .*

The problem can be extended to multiple groups, corresponding to multiple embeddings. A natural greedy algorithm, which iteratively selects the set with the largest marginal gain in coverage, achieves  $(1 - 1/e)$ -approximation for this problem [25].

#### NOTATIONS

Let  $X = \{x_i\}_{i=1}^n$  represent a set of  $N$  data points, where  $x_i \in \mathcal{X}$ . For this dataset define the graph  $G = (V, E)$ , with vertices  $V = \{v_i\}_{i=1}^n$  where  $v_i \leftrightarrow x_i$ , and edges  $e_{i,j} = D(x_i, x_j)$  for some distance metric  $D: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_{\geq 0}$ .

With multiple embeddings, each dataset can now be represented by a collection of graphs  $\{V, E^{(m)}\}$ , where  $m \in [M]$  indexes the embeddings,  $v_i \leftrightarrow x_i$  and  $e_{i,j}^{(m)} = D(z_i^{(m)}, z_j^{(m)})$  for an embedding  $f^{(m)}: \mathcal{X} \rightarrow \mathcal{Z}^{(m)}$ .

**Definition 2** (Cover ball). *Fix  $\delta > 0$ , and consider an embedding  $f^{(m)}: \mathcal{X} \rightarrow \mathcal{Z}^{(m)}$  where  $z_x^{(m)} = f^{(m)}(x)$ . Define*

$$B_\delta^{(m)}(x) = \{x' \in \mathcal{X} \mid D(z_{x'}^{(m)}, z_x^{(m)}) \leq \delta\}$$

*$B_\delta^{(m)}(x)$  denotes the set of points whose embedding lies inside the ball of radius  $\delta$  centered at  $x$  in embedding  $m$ <sup>1</sup>.*

#### 3.1 COVERAGE-BASED SELECTION

Our method is designed to enhance any coverage-based selection method within the framework of *replay-based CIL*. These methods aim to select a small, representative subset of the original dataset; their core rationale is to reduce the selection problem to one of maximizing graph coverage. Thus, the selection of subset of  $b$  elements is reduced to max probability coverage as follows:

<sup>1</sup>Superscript  $(m)$  can be omitted with a single embedding.

**Definition 3** (Max Probability Cover). Fix  $\delta > 0$ , and obtain a subset  $\mathcal{M} \subset X$  with  $|\mathcal{M}| = b$  that maximizes the probability of the covered area:

$$\mathcal{M} = \arg \max_{L \subseteq X, |L|=b} P \left( \bigcup_{x \in L} B_\delta(x) \right) \quad (2)$$

In *ProbCover* [27], the probability of the covered area is estimated by way of the empirical likelihood -  $P(\bigcup_{x \in L} B_\delta(x)) = |\bigcup_{x \in L} B_\delta(x)|$ . In other words, we seek a subset  $\mathcal{M}$  for which the number of points in the original dataset that lie within distance  $\delta$  of the points in  $\mathcal{M}$  is as large as possible. Finally, in *MaxHerding* [3], the probability of the covered area is estimated with an RBF kernel, centered at each of the selected points in set  $\mathcal{L}$ .

### 3.2 BUFFER SELECTION, MULTIPLE EMBEDDINGS

We begin by generalizing (2) to multiple embeddings:

**Definition 4** (Buffer selection, weighted coverage). Obtain a subset of elements  $\mathcal{M} \subset \mathbb{X}$ ,  $|\mathcal{M}| = b$ , that maximizes

$$\mathcal{M} = \arg \max_{L \subseteq X, |L|=b} \sum_{m=1}^M \alpha_m \left| \bigcup_{x \in L} B_\delta^{(m)}(x) \right| \quad (3)$$

In (3), we seek a subset  $\mathcal{M}$  that maximizes the coverage probability in all the embeddings. This is achieved by defining a notion of *weighted coverage*, which sums over the coverage in each embedding with weight  $\alpha_m$ . The weights reflect the relative importance (or effectiveness) of each embedding, and need to be determined by the algorithm.

The optimization problem in (3) is an instance of the *weighted maximum k-coverage with m groups* (Def. 1), under the following correspondence: (i)  $U^m$  corresponds to the set  $X$  in embedding  $m$ , i.e., for every  $x_i \in X$  there exists a corresponding  $u_i^m \in U^m$ . (ii) For each datapoint  $x_i$  we associate the subset  $S_i = \bigcup_m B_{\delta_m}^{(m)}(u_i^m)$ ,  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ .

### 3.3 EMBEDDING ALIGNMENT

In each coverage-based selection method described above, there is at least one parameter that captures the range of similarities in the data, and how well the data is partitioned. Each algorithm is sensitive in different ways to this parameter, and therefore its automatic evaluation from the data is crucial to the effectiveness of the method. This problem is exacerbated in our work, as we integrate multiple embedding spaces  $\{\mathcal{E}^{(m)}\}_{m=1}^M$  that originate from distinct backbones and therefore may exhibit markedly different geometric characteristics. Therefore, a fixed one-size-fits-all solution is not likely to be effective.

For the purpose of embedding alignment, we use *k*-Nearest Neighbor (*k*-NN) density estimation, which is a non-parametric technique widely used in machine learning to estimate an unknown sample distribution without making any parametric assumptions. This model is used in [27; 3] for the generalization analysis of both *ProbCover* and *MaxHerding* because it depends exclusively on distances from a set of training examples and does not involve any additional inductive bias.

Specifically, when integrating either *ProbCover* or *MaxHerding* (defined in Section 3.1) into *MERS*, the relevant parameters are  $\delta$  from Def. 2, and the bandwidth of the RBF kernel  $\sigma$  respectively. In order to *adapt* each hyper-parameter to the statistics *within* the relevant embedding, we propose to use the median of *k*-NN distances for  $\delta$  where *k* is determined by the memory-aware ratio, and the median heuristic for  $\sigma$ . This allows the method to align the covering sets and kernel similarities with the true geometry and sparsity of  $\mathcal{E}^{(m)}$ .

This alignment is critical: it guarantees that *MaxHerding* selects representatives at the right granularity, and that *ProbCover* strikes an optimal balance between coverage and diversity, *regardless of which embedding space is active at any point in the stream*, as demonstrated in the ablation study (Section 6).

### 3.3.1 SELECTION OF $\delta$ IN *ProbCover*

We estimate the radius  $\delta$  of each *Cover ball* (see Def. 2) by computing the median of the  $k$ -nearest neighbor ( $k$ -NN) distances for every feature vector in the current stream:

More specifically, let  $\mathcal{M}_c = \{x \in X \mid y(x) = c\}$ . For each  $\mathbf{x}_i \in \mathcal{M}_c$ , let  $\mathcal{N}_K(\mathbf{x}_i)$  be the set of its  $k$  nearest neighbors in  $\mathcal{M}_c \setminus \mathbf{x}_i$ . Compute the median distance from  $\mathbf{x}_i$  to its  $k$  nearest neighbors:  $r_i := \text{median}_{\mathbf{x}_j \in \mathcal{N}_K(\mathbf{x}_i)} \|\mathbf{x}_i - \mathbf{x}_j\|$ . Finally, take the median over set  $\{r_i\}$  and fix  $\delta = \text{median}_{i \in [N]} r_i$  to be the radius of  $\mathcal{B}_\delta$  in *ProbCover*.

The neighborhood size  $k$  is also a data-driven hyper-parameter that is being adapted to both the stream statistics and the class-specific memory budget:  $K = \frac{|\mathcal{D}_c|}{\mathcal{M}_c}$ , where  $|\mathcal{D}_c|$  denotes the number of samples of class  $c$  observed in the current episode, and  $\mathcal{M}_c$  is the buffer capacity allocated to that class. This ratio partitions the feature space into as many covering regions as the buffer can hold: a larger buffer yields a finer subdivision (larger  $k$ ), whereas a smaller buffer enforces coarser regions. For the Model based embedding, we fix  $K = 1$  to enforce localized neighborhood selection. Because this representation space is denser, a smaller  $\delta$  is required to prevent neighborhoods from overlapping excessively, thereby maintaining adequate sample diversity in the buffer.

### 3.3.2 BANDWIDTH SELECTION FOR THE RBF KERNEL IN *MaxHerding*

When *MaxHerding* is used for coverage-based selection, it employs the radial basis function (RBF) kernel  $\kappa(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / (2\sigma^2))$ . Following the widely adopted *median heuristic* [13], we set the bandwidth  $\sigma$  to the median cosine distances among all exemplars in the current episode.

### 3.3.3 WEIGHTING EACH EMBEDDING

Finally, we discuss the estimation of the vector of weights  $\{\alpha_m\}$  defined in (3).

First, we recall the definition of the  $k$ -NN density estimation. Once again, let  $\mathcal{M}_c = \{x \in X \mid y(x) = c\}$ . For any  $x \in \mathcal{M}_c$ , let  $\mathcal{N}_K^{(m)}(\mathbf{x})$  denote the set of its  $k$  nearest neighbors in  $\mathcal{M}_c \setminus \mathbf{x}$  in embedding  $\mathcal{E}^{(m)}$ . Let  $\rho_k^{(m)}(x)$  denote the mean distance from  $x$  to set  $\mathcal{N}_K(\mathbf{x})$ . In embedding  $m$ , the  $k$ NN density estimate at  $x$  is defined as follows:

$$\hat{f}_k^{(m)}(x) = \frac{k}{\rho_k^{(m)}(x)} \quad (4)$$

For embedding  $m$ , we now defined its weight as follows:

$$\alpha_m = \frac{\text{median}(\hat{f}_k^{(m)}(x))}{\text{median}(\hat{f}_1^{(m)}(x))} \quad (5)$$

The reasoning behind this definition is as follows: if two point clouds differ only by a scale factor, the distribution of  $\alpha$  remains unchanged, resulting in  $\alpha_1 = \alpha_2$ . In practice, however, the supervised embedding  $\mathcal{E}_{\text{Model Based}}$  tends to exhibit *micro-clusters* - tightly grouped, nearly identical samples within a class - more so than the self-supervised embedding  $\mathcal{E}_{\text{self-supervised}}$ . These geometric irregularities disrupt scale invariance:  $\rho_1/\rho_k$  in  $\mathcal{E}_{\text{sup}}$  and deflate it in  $\mathcal{E}_{\text{self}}$ , so that

$$\beta = \frac{\alpha_{\text{Model Based}}}{\alpha_{\text{self-supervised}}} > 1. \quad (6)$$

Our greedy algorithm maximizes the *weighted coverage score* defined in (3). Because the algorithm also enforces *diversity* through disjoint  $k$ -NN balls, dense supervised balls contain far fewer candidate edges than large self-supervised balls. Multiplying the supervised edge count by  $\beta$  therefore equalizes the **effective area** (i.e. edge mass) that each selected point can cover, ensuring that the sampler does not over-represent the sparse self-supervised space and achieves a balanced, diverse subset across both embeddings.

### 3.4 PSEUDO-CODE

We propose two variants of our method *MERS*, that are distinguished by the coverage-based method they incorporate - *ProbCover* or *MaxHerding*. Pseudo-code for these 2 variants is provided in Algorithm 1 and Algorithm 2 respectively.

---

#### Algorithm 1 *MERS ProbCover*

---

**Input:** Set  $(X_m)$  of exemplars from  $m$  embeddings, Memory buffer  $\mathcal{M}$ , Ball-size  $\delta$

**Output:**  $\mathcal{M}$

```

1:  $B_\delta^{(m)}(x) = \{x' \in \mathcal{X} \mid D(z_{x'}^{(m)}, z_x^{(m)}) \leq \delta\}$ 
2:  $G = \{V, E^{(m)} = \{(x, x') : x' \in B_\delta^{(m)}(x)\}\}$ 
3: for  $i \in \{1, \dots, b\}$  do
4:    $\mathcal{M} \leftarrow \operatorname{argmax} \left( \sum_{m=1}^M \alpha_m \left| \bigcup_{x \in X_m} B_\delta^{(m)}(x) \right| \right)$ 
5:   Remove all the incoming edges to covered vertices,
    $E \leftarrow E \setminus \{(x, x') : x' \in B_\delta^{(m)}(x)\}$ 
6: end for
7: return  $\mathcal{M}$ 

```

---



---

#### Algorithm 2 *MERS MaxHerding*

---

**Input:** Set  $(X_m)$  of exemplars from  $m$  embeddings, Memory buffer  $\mathcal{M}$

**Output:**  $\mathcal{M}$

```

1: Compute integrated kernel:
    $k(x, x') = \sum_{m=1}^M \alpha_m k_m(x^{(m)}, x'^{(m)})$ 
2:  $\mathcal{B} \leftarrow \emptyset$ ;  $\mathbf{k} \in \mathbb{R}^{|\mathcal{C}|}$  with  $k_i = 0$ 
3: for  $b \in \{1, \dots, B\}$  do
4:   Select
    $x_b^* = \operatorname{argmax}_{\tilde{x} \in \mathcal{C}} \frac{1}{|\mathcal{C}|} \sum_{n=1}^{|\mathcal{C}|} \max(k(x_n, \tilde{x}) - k_n, 0)$ 
5:   Update  $k_i \leftarrow \max(k(x_i, x_b^*), k_i)$ 
6:    $\mathcal{B} \leftarrow \mathcal{B} \cup \{x_b^*\}_{b=1}^B$ 
7:    $\mathcal{C} \leftarrow \mathcal{C} \setminus \{x_b^*\}_{b=1}^B$ 
8: end for

```

---

## 4 METHODOLOGY

In our empirical evaluation, we separately evaluate each of the two aforementioned variants of our method. We report 3 scenarios: (i) **SimCLR *MERS*** uses a single unsupervised embedding - the SimCLR method recalled below; (ii) **Model Based *MERS*** uses a single supervised embedding, obtained from the learned classifier; (iii) **integrated *MERS*** as defined in Algorithms 1-2, which integrate the two embeddings.

Our method is evaluated while enhancing 3 distinct experience replay continual learning algorithms, detailed in Section 4.1. It is compared to the vanilla version of each method, and to another method that can be used to enhance the buffer selection step, which is described in Section 4.2. The two coverage-based methods, employed by the variants of *MERS*, are expanded on in Section 4.3. Section 4.4 describes the two datasets used in our evaluation, following customary practice in the evaluation of CIL methods. Common evaluation metrics are described in Section 4.5.

**SimCLR [11]** is a self-supervised contrastive learning method that learns visual representations by leveraging data augmentation and a contrastive loss, without requiring any labeled data. The core idea is to train a neural network to recognize that different augmented views of the same image should have similar representations, while views of different images should be distinct. This is achieved by applying random augmentations to each image in the dataset, creating pairs of correlated views that serve as positive examples. All other images in the batch serve as implicit negatives, encouraging the model to distinguish between different inputs based solely on appearance.

**VICReg [6]** is a self-supervised learning method that combines three complementary regularization terms: (i) The invariance term encourages representations of different augmented views of the same image to be close, (ii) the variance term prevents representation collapse by ensuring each dimension has non-trivial variance across a batch, (iii) and the covariance term reduces redundancy by decorrelating feature dimensions. Together, these constraints produce stable and diverse embeddings that transfer well to downstream tasks, while simplifying training compared to contrastive approaches.

**DINO [9]** is a self-supervised distillation method where a student matches soft targets from a teacher on augmented views of the same image. The teacher is an exponential moving average (EMA) of the student, providing stable targets without labels. DINOv2 [19] extends this with Vision Transformers, a 142M-image dataset, and training refinements, producing versatile representations transferable to tasks like classification, retrieval, and clustering.

#### 4.1 CONTINUAL LEARNING ALGORITHMS

We evaluated *MERS* with three rehearsal-based continual learning algorithms. First, **ER** [21] stores a subset of past examples in a memory buffer and replays them during training to reduce forgetting. Then, **ER-ACE** [7] extends ER by separating the loss contributions of new data and replayed samples. Finally, **MIR** [1] selects from the buffer the samples whose loss increases most after a gradient step on the current batch, focusing rehearsal on knowledge most at risk of forgetting.

#### 4.2 ALTERNATIVE SELECTION STRATEGIES

Herding [26; 20] is one of the earliest exemplar selection strategies. It constructs a representative memory by sequentially selecting samples whose inclusion best approximates the class mean in feature space, ensuring that the chosen exemplars collectively act as a centroid for their class. Rainbow Memory [4] takes a complementary approach by explicitly balancing multiple selection criteria, such as diversity, uncertainty, and class balance, when building the memory buffer. TEAL [23] is a low-budget exemplar selection strategy for CIL. It clusters class samples in feature space and selects the most typical point from each cluster, ensuring both diversity and representativeness.

#### 4.3 DIFFERENT COVERAGE STRATEGIES

**ProbCover** [27] An active-learning algorithm that maximizes coverage under a small budget by building an  $r$ -neighborhood graph with radius  $\delta$  and iteratively selecting the node with maximal uncovered degree. For continual learning, we adapt it by treating the memory buffer as the pool and the exemplar set as labeled data (see Section 3.4).

**MaxHerding** [3] An active learning method that generalizes *ProbCover* by replacing hard  $\delta$ -ball coverage with a continuous kernel-based similarity measure. In its greedy variant, *MaxHerding* evaluates, at each step, a gain function that favors points in dense regions while penalizing redundancy with the current exemplar set, and then selects the sample with the highest gain.

#### 4.4 DATASETS

Two datasets, that are commonly used to evaluate CIL methods, are used here: **(i) Split CIFAR-100** [10; 20], created by splitting CIFAR-100 and is divided into 10 episodes, each containing 10 different classes with 500 train images, and 100 test images. **(ii) Split TinyImageNet**, [16] created by splitting TinyImageNet and is divided into 10 episodes, each contain 20 different classes with 500 train images and 50 test images.

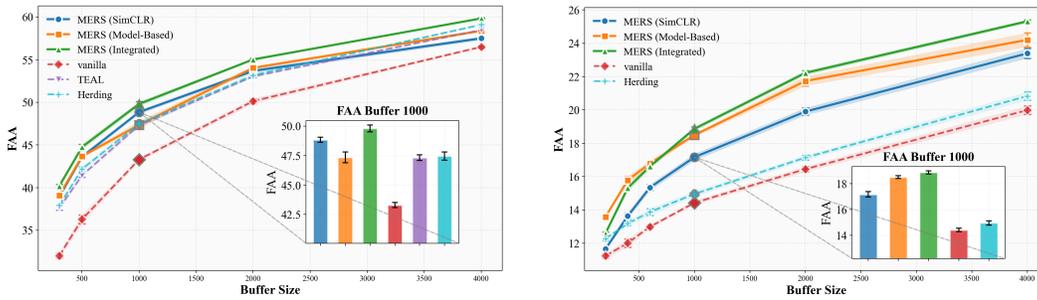
#### 4.5 EVALUATION METRICS IN CIL

The **Average Accuracy** ( $AA_t$ ) is the mean accuracy over all tasks up to task  $t$ . The **Final Average Accuracy** (**FAA**) is the average accuracy after training on the last task  $T$ , i.e.,  $FAA = AA_{t=T}$ , which measures overall performance across all tasks. The **Any-time Average Accuracy** (**AAA**) is the mean of  $AA_t$  across all  $T$  tasks:  $AAA = \frac{1}{T} \sum_{t=1}^T AA_t$ .

## 5 EMPIRICAL RESULTS

### 5.1 MAIN RESULTS

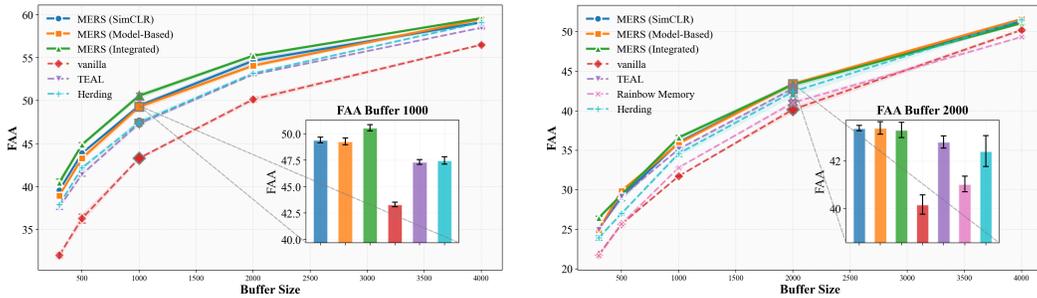
In our empirical evaluation, we assess two variants of our *MERS* that rely on two different coverage based method, denoted *MERS ProbCover* and *MERS MaxHerding*. Each method is evaluated in 3 conditions: (i) constrained to use a single unsupervised SimCLR embedding; (ii) constrained to use only the supervised embedding derived from the classifier; (iii) allowed to exploit both embeddings as defined in Algorithms 1–2. These variants are used to enhance three existing experience replay continual learning algorithms, and are compared against both the vanilla versions of those algorithms and 3 alternative buffer selection methods - TEAL, Herding and Rainbow Memory - described in Section 4.2. To assess robustness to memory constraints, we varied the replay-buffer capacity across multiple sizes, from 100 to 1000 on the Split CIFAR-100 benchmark, and 200 to 6000 on the Split



(a) CIFAR-100

(b) TinyImageNet

Figure 2: **MERS ProbCover**: FAA of ER-ACE as a function of  $|\mathcal{M}|$ , on CIFAR-100 (left) and TinyImageNet (right). Three variants of **MERS ProbCover** are evaluated and compared against alternative selection strategies.



(a) ER-ACE

(b) MIR

Figure 3: **MERS MaxHerding**: FAA of ER-ACE (left) and MIR (right) on *CIFAR-100*, as a function of  $|\mathcal{M}|$ . Results with **MERS MaxHerding** are compared against other selection strategies.

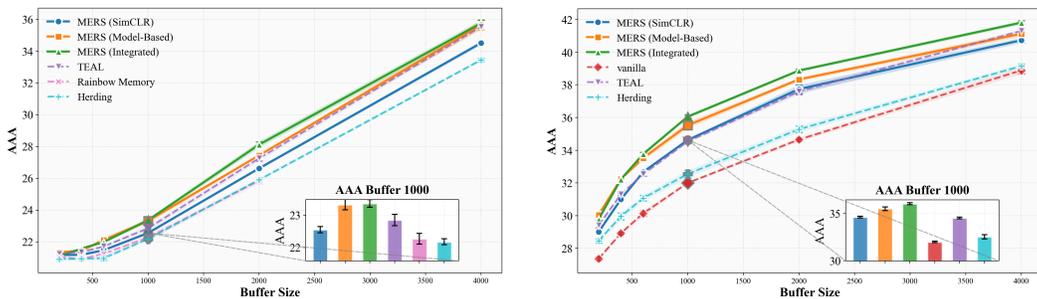
TinyImageNet benchmark. Final Average Accuracy (FAA) is reported in Figs. 2–4, while the complete results, including FAA and AAA, are provided in Appendix A (see Tables 1–10).

## 5.2 IMPACT OF PRETRAINED VS. EPISODIC EMBEDDINGS ON MERS

Following the main results protocol (Section 5.1), we evaluate **MERS** with ER-ACE on Split CIFAR-100 across different buffer sizes using embeddings beyond SimCLR. (i) **VICReg** is trained from scratch at each episode using only the current episode’s training data, identical to the SimCLR protocol. (ii) **DINOv2** embeddings are extracted from a foundational model (Section 4). Results are presented in Fig. 5, with complete FAA and AAA tables reported in Appendix A.

## 5.3 DISCUSSION

Across every buffer size, experience replay base-method and dataset, **MERS ProbCover** achieves the most competitive results. Its integrated variant matches or exceeds its constrained variants,



(a) ER

(b) ER-ACE

Figure 4: **MERS ProbCover**: AAA of ER (left) and ER-ACE (right) on *TinyImageNet*. Three variants of **MERS ProbCover** are evaluated, compared to other selection strategies

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

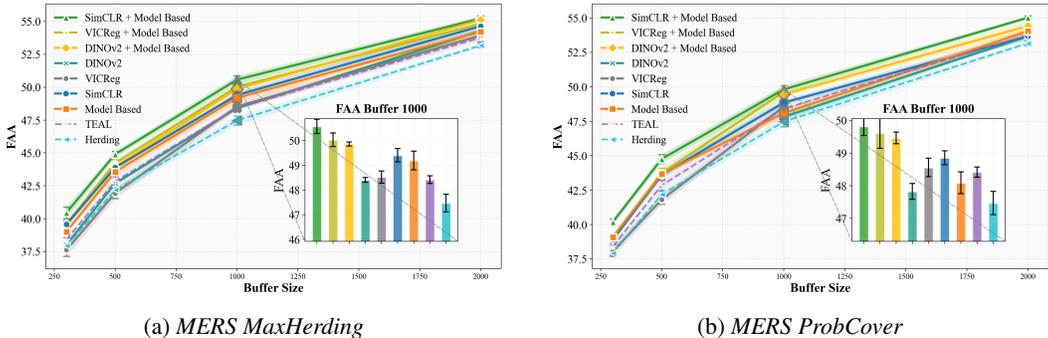


Figure 5: FAA of *MERS MaxHerding* (left) and *MERS ProbCover* (right) with ER-ACE on *CIFAR-100* with different embeddings: SimCLR, VICReg and DINOv2

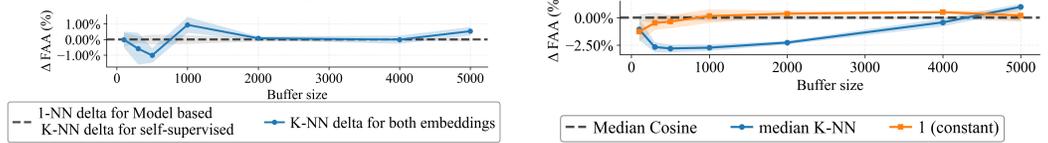


Figure 6: Improvements in FAA on *CIFAR-100* as a function of  $|\mathcal{M}|$  while varying *ProbCover*'s radius  $\delta$ .

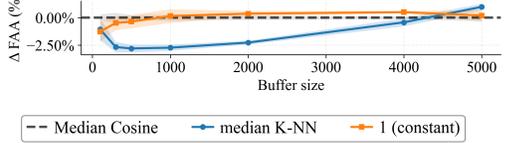


Figure 7: Improvements in FAA on *CIFAR-100* as a function of  $|\mathcal{M}|$  while varying the RBF bandwidth  $\sigma$  in *MaxHerding*.

an advantage that is most pronounced in the low-budget regime (up to 1000 exemplars), where it opens a clear gap over the 2 constrained variants. As  $|\mathcal{M}|$  increases the gap indeed narrows, yet the integrated *MERS* still *retains first place*, sharing the highest score with one of the constrained variants. *MERS MaxHerding* shows a split pattern, staying ahead of the constrained variants across all buffer sizes only on *CIFAR-100*. Overall, *MERS* significantly outperforms either embedding in isolation, with the integrated design providing a distinct advantage under limited memory. In particular, episodic embeddings drive stronger task adaptation when memory is constrained.

## 6 ABLATION STUDY

We conducted targeted ablations to identify which design choices of our *MERS* are most critical:

**Adaptive *ProbCover* radius  $\delta$ :** Fig. 6 shows how varying  $\delta$  affects *MERS ProbCover*. When used for active learning, the algorithm suffers from high sensitivity to  $\delta$ , as argued in [3]. We note that the optimal value of  $\delta$  differs between Model based and self-supervised embeddings, reflecting the distinct statistical properties of supervised versus contrastive representations.

**RBF bandwidth  $\sigma$  in *MaxHerding*** We tested three settings for  $\sigma$ : (i) median cosine distances, (ii)  $\sigma = 1$ , and (iii) median  $k$ -NN distances. On *CIFAR-100*, (i) and (iii) coincide, while the constant value reduces FAA by  $\approx 1\%$  in the small-buffer regime. As (i) is dataset-agnostic and robust across budgets, we adopt it as the default.

We conducted an ablation study on the embedding weight  $\alpha$  using different density estimators. The results show a slight improvement when using the  $\alpha$  defined in (5), as reported in Appendix A.

## 7 SUMMARY

We present Multi-Embedding Replay Selection (*MERS*), a plug-and-play sampler for replay-based continual learning that merges supervised and self-supervised feature spaces in a complementary manner. By building  $k$ -NN coverage graphs in each space, re-scaling them with density-aware weights, and greedily selecting exemplars that maximize a combined coverage score, *MERS* fills both class-discriminative and invariant regions of the data manifold. Across Split *CIFAR-100* and Split *TinyImageNet*, it boosts final-average accuracy over single-embedding baselines when memory is tight, all without increasing the buffer size or changing model parameters. The method is plug-and-play, incurs only double selection-time overhead and self-supervised training. The approach opens avenues for dynamic, task-aware embedding integration in future work.

## REFERENCES

- 486  
487  
488 [1] Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min  
489 Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval.  
490 *Advances in neural information processing systems*, 32, 2019.
- 491 [2] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample se-  
492 lection for online continual learning. *Advances in neural information processing systems*, 32,  
493 2019.
- 494 [3] Wonho Bae, Junhyug Noh, and Danica J Sutherland. Generalized coverage for more ro-  
495 bust low-budget active learning. In *European Conference on Computer Vision*, pp. 318–334.  
496 Springer, 2024.
- 497 [4] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow mem-  
498 ory: Continual learning with a memory of diverse samples. In *Proceedings of the IEEE/CVF*  
499 *conference on computer vision and pattern recognition*, pp. 8218–8227, 2021.
- 500 [5] Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow  
501 memory: Continual learning with a memory of diverse samples, 2021. URL [https://](https://arxiv.org/abs/2103.17230)  
502 [arxiv.org/abs/2103.17230](https://arxiv.org/abs/2103.17230).
- 503 [6] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regular-  
504 ization for self-supervised learning. *arXiv preprint arXiv:2105.04906*, 2021.
- 505 [7] Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene  
506 Belilovsky. New insights on reducing abrupt representation change in online continual learn-  
507 ing. *arXiv preprint arXiv:2104.05025*, 2021.
- 508 [8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
509 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceed-*  
510 *ings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- 511 [9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski,  
512 and Armand Joulin. Emerging properties in self-supervised vision transformers. In *2021*  
513 *IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 9630–9640, 2021. doi:  
514 10.1109/ICCV48922.2021.00951.
- 515 [10] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, P Doka-  
516 nia, P Torr, and M Ranzato. Continual learning with tiny episodic memories. In *Workshop on*  
517 *Multi-Task and Lifelong Reinforcement Learning*, 2019.
- 518 [11] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework  
519 for contrastive learning of visual representations. In *International conference on machine*  
520 *learning*, pp. 1597–1607. PmLR, 2020.
- 521 [12] Victor Guilherme Turrise da Costa, Enrico Fini, Moin Nabi, Nicu Sebe, and Elisa Ricci. solo-  
522 learn: A library of self-supervised methods for visual representation learning. *Journal of*  
523 *Machine Learning Research*, 23(56):1–6, 2022. URL [http://jmlr.org/papers/v23/](http://jmlr.org/papers/v23/21-1155.html)  
524 [21-1155.html](http://jmlr.org/papers/v23/21-1155.html).
- 525 [13] Damien Garreau, Wittawat Jitkrittum, and Motonobu Kanagawa. Large sample analysis of the  
526 median heuristic, 2018. URL <https://arxiv.org/abs/1707.07269>.
- 527 [14] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified  
528 classifier incrementally via rebalancing. In *Proceedings of the IEEE/CVF conference on com-*  
529 *puter vision and pattern recognition*, pp. 831–839, 2019.
- 530 [15] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, An-  
531 dree A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al.  
532 Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy*  
533 *of sciences*, 114(13):3521–3526, 2017.
- 534 [16] Yann Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.

- 540 [17] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern*  
541 *analysis and machine intelligence*, 40(12):2935–2947, 2017.  
542
- 543 [18] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks:  
544 The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp.  
545 109–165. Elsevier, 1989.  
546
- 547 [19] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khali-  
548 dov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud As-  
549 sran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan  
550 Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou, Julien Mairal,  
551 Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual fea-  
552 tures without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 553 [20] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:  
554 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on*  
555 *Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.  
556
- 557 [21] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Ex-  
558 perience replay for continual learning. *Advances in neural information processing systems*, 32,  
559 2019.  
560
- 561 [22] Joan Serrà, Dídac Surís, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic  
562 forgetting with hard attention to the task, 2018. URL [https://arxiv.org/abs/1801.](https://arxiv.org/abs/1801.01423)  
563 [01423](https://arxiv.org/abs/1801.01423).
- 564 [23] Shahar Shaul-Ariel and Daphna Weinshall. Teal: New selection strategy for small buffers in  
565 experience replay class incremental learning. *arXiv preprint arXiv:2407.00673*, 2024.  
566
- 567 [24] Tobias Uelwer, Jan Robine, Stefan Sylvius Wagner, Marc Höftmann, Eric Upschulte, Sebastian  
568 Konietzny, Maike Behrendt, and Stefan Harmeling. A survey on self-supervised methods for  
569 visual representation learning. *Machine Learning*, 114(4):1–56, 2025.  
570
- 571 [25] Vijay V Vazirani. *Approximation algorithms*, volume 1. Springer, 2001.  
572
- 573 [26] Max Welling. Herding dynamical weights to learn. In *Proceedings of the 26th annual interna-*  
574 *tional conference on machine learning*, pp. 1121–1128, 2009.
- 575 [27] Ofer Yehuda, Avihu Dekel, Guy Hacohen, and Daphna Weinshall. Active learning through a  
576 covering lens. *Advances in Neural Information Processing Systems*, 35:22354–22367, 2022.  
577
- 578 [28] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dy-  
579 namically expandable networks, 2018. URL <https://arxiv.org/abs/1708.01547>.  
580
- 581 [29] Zhiwen Yu, Ziyang Dong, Chenchen Yu, Kaixiang Yang, Ziwei Fan, and CL Philip Chen. A  
582 review on multi-view learning. *Frontiers of Computer Science*, 19(7):197334, 2025.  
583
- 584 [30] Qinghai Zheng, Jihua Zhu, Zhongyu Li, Zhiqiang Tian, and Chen Li. Comprehensive multi-  
585 view representation learning. *Information Fusion*, 89:198–209, 2023.  
586

## 587 A APPENDIX

### 588 A.1 USE OF LLMs

589  
590 A large language model (ChatGPT, GPT-5 by OpenAI) was employed solely for minor editorial  
591 assistance. All methodological design, experimental results, and scientific conclusions are entirely  
592 the authors’ own.  
593

## 594 A.2 TIME AND SPACE COMPLEXITY OF *MERS*

595 We analyse the computational cost under the standard setting in which the selection strategy is  
 596 invoked *once per training episode*. Let  $n$  be the number of examples from the current episode  
 597 that belong to class  $c$ ,  $M$  the number of distinct embedding spaces,  $d$  the dimensionality of each  
 598 embedding, and  $b$  the class-wise memory-buffer budget ( the number of items that  $|\mathcal{M}|$  may store  
 600 for class  $c$ ).

601 **Self-supervised stage.** During every episode, *MERS* is called exactly once. Running SimCLR for  
 602  $E_{\text{ssl}}$  epochs on  $A = 2$  views of the  $n$  episode images costs

$$603 T_{\text{SimCLR}} = O(E_{\text{ssl}} A n P)$$

604 with  $P$  trainable parameters. Self-supervised training consumes

$$605 S_{\text{SimCLR}} = O(P + s f)$$

606 space model parameters  $P$  plus the current batch’s  $s$  activations of size  $f$ , and the batch size  $s$ . The  
 607 SimCLR weights are discarded after each episode, persistent memory is dominated by the replay  
 608 images.

### 609 A.2.1 *MERS ProbCover*

610 The algorithm consists of two stages:

611 **(i) Ball-graph construction.** For every embedding  $m \in \{1, \dots, M\}$  we compute all pair-  
 612 wise cosine distances in  $\mathbb{R}^d$  to obtain the  $\delta$ -neighbourhoods  $B_\delta^{(m)}(x)$ . This step costs  $T_{\text{graph}} =$   
 613  $O(M n^2 \max\{d, b\})$  and stores  $S_{\text{graph}} = O(M n^2)$  adjacency edges.

614 **(ii) Greedy covering.** Across  $b$  iterations we repeatedly pick the vertex that covers the largest  
 615 number of still-uncovered neighbours. The work per iteration yields  $T_{\text{cover}} = O(|E| + b n) \subseteq$   
 616  $O(M n^2 + b n)$ .

617 **Overall complexity.**

$$618 T_{\text{MERS-ProbCover}} = O(M n^2 \max d, b),$$

$$619 S_{\text{MERS-ProbCover}} = O(M n^2),$$

620 The original ProbCover analysis [27] reports a running time of  $O(n^2 \max\{d, b\})$ . Our derivation  
 621 shows that the multi-embedding extension, *MERS-ProbCover*, retains the same quadratic depen-  
 622 dence on  $n$  and on  $\max d, b$ , differing only by the multiplicative factor  $M$  (which equals 2 in all of  
 623 our experiments).

624 **(ii) Greedy *MaxHerding* selection.**

625 **(i) Integrated-kernel construction.** We assemble the Gram matrix

$$626 K_{ij} = k(x_i, x_j) = \sum_{m=1}^M \alpha_m k_m(x_i^{(m)}, x_j^{(m)}).$$

627 Forming its  $\frac{1}{2}n(n-1)$  entries costs

$$628 T_{\text{kernel}} = O(m n^2 d), \quad S_{\text{kernel}} = O(n^2).$$

629 **(ii) Greedy selection.** Each of the  $b$  iterations scans all candidates ( $\leq n$ ) and exploits the pre-  
 630 computed kernel:

$$631 T_{\text{MaxHerding}} = O(b n^2), \quad S_{\text{MaxHerding}} = O(n).$$

632 **Overall complexity.**

$$633 T_{\text{MERS-MaxHerding}} = O(m n^2 (d + b)),$$

$$634 S_{\text{MERS-MaxHerding}} = O(n^2 + n d).$$

648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701

### A.3 HYPERPARAMETERS

#### A.3.1 CLASSIFICATION MODEL

we employ a ResNet-18 backbone trained for 100 epochs with a batch size of 10. The **ER-ACE** configuration begins with a learning rate of 0.01. The **ER** and **MIR** configuration begins with a learning rate of 0.1, for all configurations, SGD optimization includes Nesterov momentum of 0.9 and weight decay 0.0002. The learning rate is decayed by a factor of 0.3 every 66 epochs. All experiments were run with five random seeds (0-4).

#### A.3.2 CLASS ORDER

We follow the canonical class order for each benchmark: Split CIFAR-100 uses classes [1 . . . 100], and Split TinyImageNet uses classes [1 . . . 200].

#### A.3.3 SELF-SUPERVISED TRAINING

Our SimCLR implementation is adapted from solo-learn[12], and is available in the source code. The self-supervised model is trained on the images observed in the current episode only, never on the full dataset.

#### A.3.4 FEATURE NORMALIZATION

Each feature vector is divided by its  $\ell_2$  norm, yielding unit-norm representations. Similarities are therefore computed with the cosine distance.

### A.4 COMPUTE RESOURCES

Each experiment trained deep-learning models on GPUs, consuming up to 22 GB of GPU memory and no more than 20 GB of system RAM.

### A.5 SOURCE CODE

The complete source code is provided in the supplementary ZIP file and will be publicly released on GitHub upon acceptance. The source code includes a README that lists the commands required to reproduce all of the experiments described in this paper.

### A.6 MAIN RESULTS TABLES

The tables 1- 10 presents the complete tables underlying Figs. 2- 4, evaluated with both the FAA and AAA metrics.

### A.7 TRANSFER LEARNING TABLES

The tables 11- 12 presents the complete tables underlying Fig. 5, evaluated with both the FAA and AAA metrics.

Table 1: **Final Averaged Accuracy (FAA) On CIFAR-100** averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best accuracy is in bold; results within the standard error of the best are also bolded.

(a) CIFAR-100 with the ER-ACE framework.

$ \mathcal{M} $	ER-ACE (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	21.80 $\pm$ 0.34	29.88 $\pm$ 0.24	28.61 $\pm$ 0.21	29.61 $\pm$ 0.43	29.79 $\pm$ 0.25	30.00 $\pm$ 0.32	<b>30.99</b> $\pm$ 0.29
300	32.01 $\pm$ 0.30	39.09 $\pm$ 0.22	38.92 $\pm$ 0.21	<b>40.19</b> $\pm$ 0.22	38.93 $\pm$ 0.08	39.56 $\pm$ 0.28	<b>40.45</b> $\pm$ 0.44
500	36.29 $\pm$ 0.52	43.68 $\pm$ 0.17	43.61 $\pm$ 0.28	<b>44.77</b> $\pm$ 0.32	43.31 $\pm$ 0.34	43.89 $\pm$ 0.11	<b>44.90</b> $\pm$ 0.23
1000	43.30 $\pm$ 0.21	47.36 $\pm$ 0.46	48.85 $\pm$ 0.21	49.82 $\pm$ 0.28	49.28 $\pm$ 0.33	49.40 $\pm$ 0.27	<b>50.57</b> $\pm$ 0.29
2000	50.14 $\pm$ 0.30	54.05 $\pm$ 0.18	53.69 $\pm$ 0.27	<b>55.03</b> $\pm$ 0.13	54.07 $\pm$ 0.22	54.64 $\pm$ 0.33	<b>55.22</b> $\pm$ 0.20
4000	56.50 $\pm$ 0.13	58.38 $\pm$ 0.17	57.51 $\pm$ 0.22	<b>59.84</b> $\pm$ 0.06	59.48 $\pm$ 0.17	59.11 $\pm$ 0.16	59.59 $\pm$ 0.12
5000	58.28 $\pm$ 0.26	60.08 $\pm$ 0.08	58.61 $\pm$ 0.17	<b>61.07</b> $\pm$ 0.12	60.33 $\pm$ 0.12	60.28 $\pm$ 0.17	60.66 $\pm$ 0.35

(b) CIFAR-100 with the ER framework.

$ \mathcal{M} $	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	10.07 $\pm$ 0.13	<b>11.71</b> $\pm$ 0.06	10.81 $\pm$ 0.04	11.30 $\pm$ 0.10	11.28 $\pm$ 0.08	10.91 $\pm$ 0.07	<b>11.76</b> $\pm$ 0.16
300	13.25 $\pm$ 0.10	<b>18.29</b> $\pm$ 0.28	16.83 $\pm$ 0.17	<b>18.45</b> $\pm$ 0.32	17.56 $\pm$ 0.14	17.65 $\pm$ 0.35	<b>18.36</b> $\pm$ 0.24
500	17.69 $\pm$ 0.30	<b>23.61</b> $\pm$ 0.26	21.80 $\pm$ 0.20	<b>23.64</b> $\pm$ 0.36	23.23 $\pm$ 0.29	22.59 $\pm$ 0.17	<b>23.45</b> $\pm$ 0.35
1000	26.04 $\pm$ 0.24	31.99 $\pm$ 0.28	31.29 $\pm$ 0.36	33.13 $\pm$ 0.35	32.46 $\pm$ 0.26	32.11 $\pm$ 0.19	<b>33.31</b> $\pm$ 0.08
2000	38.30 $\pm$ 0.23	42.55 $\pm$ 0.52	42.91 $\pm$ 0.26	<b>43.68</b> $\pm$ 0.22	42.37 $\pm$ 0.29	41.83 $\pm$ 0.96	<b>43.85</b> $\pm$ 0.25
4000	50.63 $\pm$ 0.10	53.12 $\pm$ 0.08	53.39 $\pm$ 0.19	<b>53.97</b> $\pm$ 0.19	53.34 $\pm$ 0.22	52.73 $\pm$ 0.20	53.21 $\pm$ 0.29
5000	53.86 $\pm$ 0.37	56.07 $\pm$ 0.26	55.87 $\pm$ 0.49	<b>56.27</b> $\pm$ 0.14	55.91 $\pm$ 0.31	55.66 $\pm$ 0.39	56.01 $\pm$ 0.17

(c) CIFAR-100 with the MIR framework.

$ \mathcal{M} $	MIR (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	17.80 $\pm$ 0.32	19.80 $\pm$ 0.28	19.48 $\pm$ 0.25	<b>20.32</b> $\pm$ 0.14	19.89 $\pm$ 0.15	19.71 $\pm$ 0.46	19.78 $\pm$ 0.23
300	21.78 $\pm$ 0.17	26.03 $\pm$ 0.25	24.63 $\pm$ 0.31	25.70 $\pm$ 0.20	24.93 $\pm$ 0.36	25.13 $\pm$ 0.20	<b>26.46</b> $\pm$ 0.22
500	25.68 $\pm$ 0.24	29.40 $\pm$ 0.12	28.88 $\pm$ 0.15	<b>29.66</b> $\pm$ 0.20	<b>29.88</b> $\pm$ 0.26	29.10 $\pm$ 0.44	29.41 $\pm$ 0.22
1000	31.74 $\pm$ 0.06	35.36 $\pm$ 0.42	34.85 $\pm$ 0.28	35.89 $\pm$ 0.31	35.88 $\pm$ 0.24	36.01 $\pm$ 0.23	<b>36.62</b> $\pm$ 0.29
2000	40.17 $\pm$ 0.41	<b>43.27</b> $\pm$ 0.16	42.85 $\pm$ 0.27	<b>43.50</b> $\pm$ 0.24	<b>43.40</b> $\pm$ 0.26	<b>43.40</b> $\pm$ 0.11	<b>43.32</b> $\pm$ 0.33
4000	50.23 $\pm$ 0.20	51.29 $\pm$ 0.29	51.31 $\pm$ 0.31	<b>51.74</b> $\pm$ 0.35	<b>51.61</b> $\pm$ 0.21	51.35 $\pm$ 0.25	51.09 $\pm$ 0.16
5000	52.56 $\pm$ 0.34	53.53 $\pm$ 0.29	53.62 $\pm$ 0.22	<b>54.17</b> $\pm$ 0.40	<b>53.88</b> $\pm$ 0.27	53.74 $\pm$ 0.28	53.70 $\pm$ 0.39

## A.8 ABLATION STUDY

Fig. 8 presents an ablation study on the effect of the embedding weight parameter  $\alpha$  when using the median K-NN density defined in Eq. 4, applied to *MERS ProbCover* on CIFAR-100 under the ER-ACE setting. The results indicate a slight but consistent improvement when using the formulation of  $\alpha$  given in Eq. 5.

We further investigate the influence of the hyperparameter  $\delta$  in *MERS ProbCover*. Fig. 9 reports results when using a single embedding (either SimCLR or model-based). In both cases, the values of  $\delta$  defined in Subsection 3.3.1 yield the best performance, confirming their suitability.

Table 2: **Final Averaged Accuracy (FAA)**, on **TinyImageNet** averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best accuracy is in bold; results within the standard error of the best are also bolded.

(a) **TinyImageNet** with the **ER-ACE** framework.

$\mathcal{M}$	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	11.24 $\pm$ 0.16	<b>13.59</b> $\pm$ 0.04	11.66 $\pm$ 0.12	12.59 $\pm$ 0.10	12.86 $\pm$ 0.09	12.57 $\pm$ 0.20	12.67 $\pm$ 0.31
400	12.01 $\pm$ 0.25	<b>15.79</b> $\pm$ 0.23	13.63 $\pm$ 0.19	15.32 $\pm$ 0.21	14.62 $\pm$ 0.22	14.76 $\pm$ 0.22	15.02 $\pm$ 0.17
600	12.99 $\pm$ 0.12	<b>16.79</b> $\pm$ 0.17	15.34 $\pm$ 0.18	16.61 $\pm$ 0.14	15.75 $\pm$ 0.20	15.88 $\pm$ 0.10	16.29 $\pm$ 0.13
1000	14.42 $\pm$ 0.13	18.49 $\pm$ 0.11	17.16 $\pm$ 0.21	<b>18.86</b> $\pm$ 0.13	18.06 $\pm$ 0.20	17.33 $\pm$ 0.17	17.73 $\pm$ 0.21
2000	16.45 $\pm$ 0.17	21.73 $\pm$ 0.29	19.90 $\pm$ 0.22	<b>22.22</b> $\pm$ 0.14	20.70 $\pm$ 0.23	19.96 $\pm$ 0.18	20.36 $\pm$ 0.19
4000	19.99 $\pm$ 0.26	24.20 $\pm$ 0.42	23.39 $\pm$ 0.30	<b>25.32</b> $\pm$ 0.10	24.27 $\pm$ 0.24	22.99 $\pm$ 0.11	23.78 $\pm$ 0.21
6000	23.14 $\pm$ 0.28	26.36 $\pm$ 0.38	26.75 $\pm$ 0.15	<b>27.88</b> $\pm$ 0.13	26.83 $\pm$ 0.24	25.65 $\pm$ 0.30	26.54 $\pm$ 0.23

(b) **TinyImageNet** with the **ER** framework.

$\mathcal{M}$	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	6.65 $\pm$ 0.04	6.69 $\pm$ 0.05	6.72 $\pm$ 0.10	6.71 $\pm$ 0.08	<b>6.84</b> $\pm$ 0.06	6.68 $\pm$ 0.04	6.70 $\pm$ 0.08
400	6.34 $\pm$ 0.07	<b>6.56</b> $\pm$ 0.09	6.44 $\pm$ 0.04	<b>6.56</b> $\pm$ 0.03	6.52 $\pm$ 0.03	6.45 $\pm$ 0.04	6.51 $\pm$ 0.02
600	6.18 $\pm$ 0.09	<b>6.53</b> $\pm$ 0.03	6.37 $\pm$ 0.06	<b>6.58</b> $\pm$ 0.06	<b>6.54</b> $\pm$ 0.04	6.43 $\pm$ 0.07	6.45 $\pm$ 0.05
1000	6.23 $\pm$ 0.03	<b>7.05</b> $\pm$ 0.06	6.53 $\pm$ 0.10	6.80 $\pm$ 0.07	6.81 $\pm$ 0.05	6.75 $\pm$ 0.06	6.63 $\pm$ 0.08
2000	7.41 $\pm$ 0.12	9.11 $\pm$ 0.10	8.51 $\pm$ 0.13	<b>9.26</b> $\pm$ 0.13	8.77 $\pm$ 0.16	8.38 $\pm$ 0.12	8.56 $\pm$ 0.13
4000	11.87 $\pm$ 0.15	15.68 $\pm$ 0.29	14.53 $\pm$ 0.26	<b>15.94</b> $\pm$ 0.20	15.05 $\pm$ 0.35	13.75 $\pm$ 0.15	14.72 $\pm$ 0.30
6000	18.70 $\pm$ 0.39	21.71 $\pm$ 0.33	20.36 $\pm$ 0.32	<b>22.24</b> $\pm$ 0.35	20.87 $\pm$ 0.42	20.08 $\pm$ 0.13	20.47 $\pm$ 0.37

$\mathcal{M}$	Herding	TEAL	Rainbow Memory
	Model Based	Model Based	Model Based
100	19.38 $\pm$ 0.04	<b>19.97</b> $\pm$ 0.30	17.54 $\pm$ 0.32
300	23.95 $\pm$ 0.30	<b>24.97</b> $\pm$ 0.24	21.68 $\pm$ 0.29
500	27.02 $\pm$ 0.16	<b>29.15</b> $\pm$ 0.23	25.69 $\pm$ 0.33
1000	34.60 $\pm$ 0.37	<b>35.16</b> $\pm$ 0.15	32.83 $\pm$ 0.06
2000	42.42 $\pm$ 0.66	<b>42.81</b> $\pm$ 0.25	41.03 $\pm$ 0.33

Table 3: FAA on **CIFAR-100** with the **MIR** framework. Comparison of Herding, TEAL, and Rainbow Memory across different buffer sizes

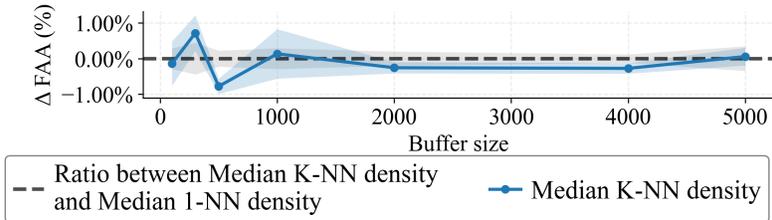


Figure 8: **MERS MaxHerding**. Ablation of the embedding weight  $\alpha$  using K-NN density estimators on CIFAR-100 with ER-ACE. The baseline corresponds to Eq. (5), and a minor but consistent improvement is observed with this weighting.

Table 4: FAA on **CIFAR-100** with the **ER-ACE** framework. Comparison of Herding and TEAL Memory across different buffer sizes

	Herding	TEAL
$ \mathcal{M} $	Model Based	Model Based
100	29.93 $\pm$ 0.31	29.67 $\pm$ 0.13
300	37.92 $\pm$ 0.13	38.40 $\pm$ 0.16
500	42.18 $\pm$ 0.27	42.84 $\pm$ 0.38
1000	47.47 $\pm$ 0.35	48.42 $\pm$ 0.16
2000	53.17 $\pm$ 0.17	53.76 $\pm$ 0.29

Table 5: FAA on **CIFAR-100** with the **ER** framework. Comparison of Herding, TEAL, and Rainbow Memory across different buffer sizes

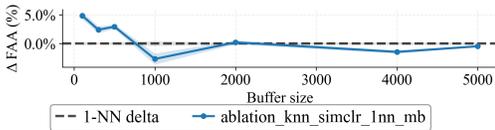
	Herding	TEAL	Rainbow Memory
$ \mathcal{M} $	Model Based	Model Based	Model Based
100	10.86	11.84 $\pm$ 0.12	10.15 $\pm$ 0.10
300	16.13	17.06 $\pm$ 0.13	13.46 $\pm$ 0.10
500	20.21	22.49 $\pm$ 0.20	16.98 $\pm$ 0.60
1000	30.01	31.92 $\pm$ 0.43	26.72 $\pm$ 0.17
2000	41.62	42.22 $\pm$ 0.51	38.40 $\pm$ 0.22

Table 6: FAA on **TinyImageNet** with the **ER-ACE** framework. Comparison of Herding, TEAL, and Rainbow Memory across different buffer sizes

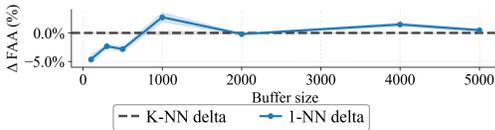
	Herding	TEAL
$ \mathcal{M} $	Model Based	Model Based
200	12.27 $\pm$ 0.12	<b>13.32</b> $\pm$ 0.29
400	13.18 $\pm$ 0.12	<b>14.91</b> $\pm$ 0.15
600	13.88 $\pm$ 0.21	<b>15.66</b> $\pm$ 0.01
1000	14.95 $\pm$ 0.17	<b>17.35</b> $\pm$ 0.17
2000	17.14 $\pm$ 0.14	<b>20.11</b> $\pm$ 0.38

Table 7: FAA on **TinyImageNet** with the **ER** framework. Comparison of Herding, TEAL, and Rainbow Memory across different buffer sizes

	Herding	TEAL	RM
$ \mathcal{M} $	Model Based	Model Based	Model Based
200	6.60 $\pm$ 0.08	6.68 $\pm$ 0.03	<b>6.78</b> $\pm$ 0.03
400	6.33 $\pm$ 0.03	6.49 $\pm$ 0.03	6.31 $\pm$ 0.11
600	6.09 $\pm$ 0.05	6.50 $\pm$ 0.06	6.38 $\pm$ 0.03
1000	6.18 $\pm$ 0.08	6.68 $\pm$ 0.07	6.43 $\pm$ 0.10
2000	7.51 $\pm$ 0.11	8.66 $\pm$ 0.10	7.79 $\pm$ 0.41



(a) *MERS ProbCover* Model Based



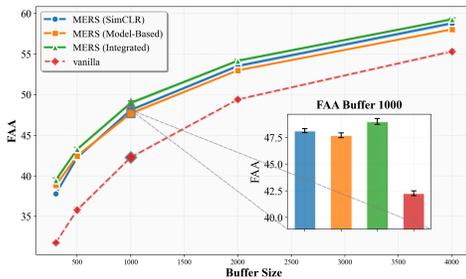
(b) *MERS ProbCover* SimCLR

Figure 9: *MERS ProbCover*. Ablation study of the hyperparameter  $\delta$  on CIFAR-100 with ER-ACE, using a single embedding. Results are shown separately for (a) model-based embeddings and (b) SimCLR embeddings.

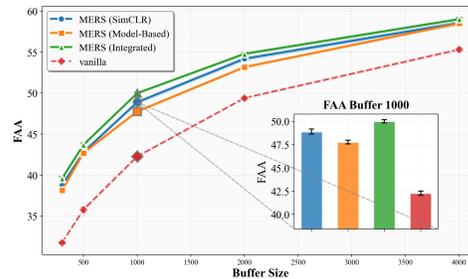
Table 8: **Average Accumulated Accuracy (AAA)**, averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best aaa is in bold; results within the standard error of the best are also bolded.

(a) **TinyImageNet** with the **ER-ACE** framework (AAA).

$ \mathcal{M} $	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	27.37 $\pm$ 0.12	<b>30.05</b> $\pm$ 0.18	29.01 $\pm$ 0.14	29.78 $\pm$ 0.08	29.51 $\pm$ 0.13	29.48 $\pm$ 0.16	29.64 $\pm$ 0.18
400	28.92 $\pm$ 0.08	<b>32.25</b> $\pm$ 0.08	31.00 $\pm$ 0.11	<b>32.22</b> $\pm$ 0.09	31.58 $\pm$ 0.13	31.91 $\pm$ 0.19	32.03 $\pm$ 0.05
600	30.14 $\pm$ 0.17	<b>33.52</b> $\pm$ 0.15	32.67 $\pm$ 0.11	<b>33.77</b> $\pm$ 0.11	32.87 $\pm$ 0.12	32.92 $\pm$ 0.14	33.43 $\pm$ 0.15
1000	32.00 $\pm$ 0.07	<b>35.52</b> $\pm$ 0.19	34.62 $\pm$ 0.08	<b>36.05</b> $\pm$ 0.10	35.25 $\pm$ 0.16	34.76 $\pm$ 0.13	34.81 $\pm$ 0.06
2000	34.65 $\pm$ 0.12	<b>38.33</b> $\pm$ 0.11	37.74 $\pm$ 0.25	<b>38.87</b> $\pm$ 0.09	38.16 $\pm$ 0.12	37.12 $\pm$ 0.11	37.74 $\pm$ 0.17
4000	38.88 $\pm$ 0.21	<b>41.12</b> $\pm$ 0.17	40.72 $\pm$ 0.16	<b>41.80</b> $\pm$ 0.11	40.98 $\pm$ 0.19	40.37 $\pm$ 0.09	40.91 $\pm$ 0.12
6000	41.24 $\pm$ 0.19	<b>43.08</b> $\pm$ 0.13	43.26 $\pm$ 0.21	<b>43.75</b> $\pm$ 0.10	42.94 $\pm$ 0.10	42.09 $\pm$ 0.32	42.49 $\pm$ 0.13



(a) *MERS ProbCover*

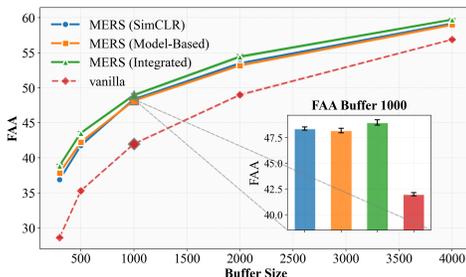


(b) *MERS MaxHerding*

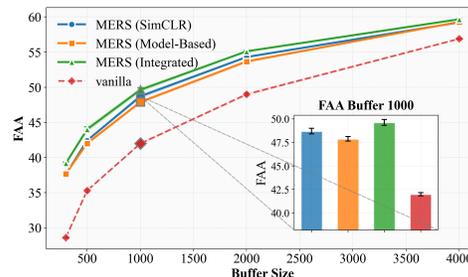
Figure 11: FAA of ER-ACE on with *MERS ProbCover* (left) and *MERS MaxHerding* (right) is shown, as a function of the replay-buffer size ( $|\mathcal{M}|$ ). The class order is generated with **seed 42**, as described in A.9. Three variants of *MERS ProbCover* are evaluated (see text), and are being compared to the original ER-ACE (vanilla).

#### A.9 ROBUSTNESS TO EPISODE CLASS ORDER IN CONTINUAL LEARNING

As in the experiments presented in Tables 1–2, we repeated them using different episode Class orders, defined by seeds 42 and 35. Below are the Final Averaged Accuracy and the Anytime Averaged Accuracy for seed 42: Tables 13, 15 and for seed 35: Tables 14, 16. and the results are analyzed in Fig. 10- 11



(a) *MERS ProbCover*



(b) *MERS MaxHerding*

Figure 10: FAA of ER-ACE on with *MERS ProbCover* (left) and *MERS MaxHerding* (right) is shown, as a function of the replay-buffer size ( $|\mathcal{M}|$ ). The class order is generated with **seed 35**, as described in A.9. Three variants of *MERS ProbCover* are evaluated (see text), and are being compared to the original ER-ACE (vanilla).

Table 9: **Average Accumulated Accuracy (AAA)**, averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best aaa is in bold; results within the standard error of the best are also bolded.

(a) CIFAR-100 with the ER-ACE framework (AAA).

$ \mathcal{M} $	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	41.31 $\pm$ 0.30	46.91 $\pm$ 0.14	46.83 $\pm$ 0.17	47.30 $\pm$ 0.27	47.01 $\pm$ 0.15	47.80 $\pm$ 0.12	<b>48.13</b> $\pm$ 0.21
300	49.90 $\pm$ 0.28	54.44 $\pm$ 0.17	54.80 $\pm$ 0.29	55.17 $\pm$ 0.17	54.59 $\pm$ 0.13	55.19 $\pm$ 0.08	<b>56.10</b> $\pm$ 0.17
500	53.72 $\pm$ 0.20	58.15 $\pm$ 0.22	58.23 $\pm$ 0.18	58.64 $\pm$ 0.31	57.99 $\pm$ 0.26	<b>58.82</b> $\pm$ 0.35	<b>59.02</b> $\pm$ 0.20
1000	58.88 $\pm$ 0.09	61.87 $\pm$ 0.24	62.18 $\pm$ 0.34	63.00 $\pm$ 0.28	62.57 $\pm$ 0.17	62.85 $\pm$ 0.35	<b>63.44</b> $\pm$ 0.26
2000	64.21 $\pm$ 0.35	66.44 $\pm$ 0.21	65.97 $\pm$ 0.21	67.06 $\pm$ 0.17	66.41 $\pm$ 0.15	66.87 $\pm$ 0.15	<b>67.59</b> $\pm$ 0.21
4000	68.98 $\pm$ 0.11	70.27 $\pm$ 0.07	68.84 $\pm$ 0.24	<b>70.64</b> $\pm$ 0.13	70.43 $\pm$ 0.21	69.46 $\pm$ 0.32	70.07 $\pm$ 0.09
5000	70.77 $\pm$ 0.33	71.47 $\pm$ 0.10	68.86 $\pm$ 0.20	<b>71.93</b> $\pm$ 0.17	71.00 $\pm$ 0.18	71.18 $\pm$ 0.34	71.59 $\pm$ 0.29

(b) CIFAR-100 with the ER framework (AAA).

$ \mathcal{M} $	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	28.68 $\pm$ 0.15	30.93 $\pm$ 0.23	30.10 $\pm$ 0.10	31.05 $\pm$ 0.13	30.97 $\pm$ 0.11	30.77 $\pm$ 0.10	<b>31.49</b> $\pm$ 0.06
300	34.58 $\pm$ 0.24	39.22 $\pm$ 0.16	38.58 $\pm$ 0.24	39.81 $\pm$ 0.17	38.86 $\pm$ 0.14	38.72 $\pm$ 0.43	<b>40.59</b> $\pm$ 0.17
500	40.41 $\pm$ 0.30	44.97 $\pm$ 0.18	44.08 $\pm$ 0.33	45.30 $\pm$ 0.44	44.71 $\pm$ 0.29	45.12 $\pm$ 0.28	<b>45.89</b> $\pm$ 0.18
1000	49.84 $\pm$ 0.45	52.99 $\pm$ 0.37	53.07 $\pm$ 0.28	53.91 $\pm$ 0.24	53.90 $\pm$ 0.10	53.71 $\pm$ 0.40	<b>54.40</b> $\pm$ 0.13
2000	59.89 $\pm$ 0.15	61.63 $\pm$ 0.48	<b>62.59</b> $\pm$ 0.15	62.03 $\pm$ 0.27	61.49 $\pm$ 0.30	61.60 $\pm$ 0.49	<b>62.59</b> $\pm$ 0.25
4000	68.41 $\pm$ 0.09	69.13 $\pm$ 0.09	69.53 $\pm$ 0.13	<b>70.18</b> $\pm$ 0.16	69.13 $\pm$ 0.24	68.83 $\pm$ 0.18	69.05 $\pm$ 0.32
5000	70.04 $\pm$ 0.13	71.19 $\pm$ 0.19	70.71 $\pm$ 0.32	<b>71.45</b> $\pm$ 0.06	70.64 $\pm$ 0.29	71.00 $\pm$ 0.21	70.77 $\pm$ 0.16

(c) CIFAR-100 with the ER-ACE framework (AAA).

$ \mathcal{M} $	Herding	TEAL	RM
	Model Based	Model Based	Model Based
100	46.88 $\pm$ 0.17	42.91 $\pm$ 0.13	10.73 $\pm$ 0.03
300	53.75 $\pm$ 0.21	53.13 $\pm$ 0.19	10.77 $\pm$ 0.04
500	57.09 $\pm$ 0.23	56.80 $\pm$ 0.14	10.69 $\pm$ 0.06
1000	61.67 $\pm$ 0.29	61.10 $\pm$ 0.28	10.75 $\pm$ 0.01
2000	66.47 $\pm$ 0.13	65.66 $\pm$ 0.24	10.79 $\pm$ 0.02
4000	70.69 $\pm$ 0.09	70.08 $\pm$ 0.26	10.66 $\pm$ 0.05

(d) CIFAR-100 with the MIR framework (AAA).

$ \mathcal{M} $	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	38.09 $\pm$ 0.20	40.41 $\pm$ 0.22	39.81 $\pm$ 0.16	<b>40.67</b> $\pm$ 0.14	<b>40.79</b> $\pm$ 0.12	40.42 $\pm$ 0.36	40.09 $\pm$ 0.29
300	44.25 $\pm$ 0.22	<b>46.80</b> $\pm$ 0.13	46.11 $\pm$ 0.36	<b>46.78</b> $\pm$ 0.28	45.41 $\pm$ 0.37	46.48 $\pm$ 0.19	<b>46.88</b> $\pm$ 0.11
500	47.89 $\pm$ 0.32	49.95 $\pm$ 0.20	49.98 $\pm$ 0.07	<b>50.15</b> $\pm$ 0.21	49.78 $\pm$ 0.50	<b>50.37</b> $\pm$ 0.28	<b>50.13</b> $\pm$ 0.26
1000	53.83 $\pm$ 0.14	55.37 $\pm$ 0.27	55.55 $\pm$ 0.21	55.81 $\pm$ 0.29	55.44 $\pm$ 0.14	55.91 $\pm$ 0.19	<b>56.32</b> $\pm$ 0.08
2000	60.52 $\pm$ 0.17	61.55 $\pm$ 0.18	<b>61.64</b> $\pm$ 0.26	<b>61.57</b> $\pm$ 0.23	<b>61.70</b> $\pm$ 0.14	<b>61.60</b> $\pm$ 0.07	<b>61.67</b> $\pm$ 0.06
4000	66.84 $\pm$ 0.21	<b>67.60</b> $\pm$ 0.12	67.44 $\pm$ 0.24	<b>67.68</b> $\pm$ 0.21	<b>67.57</b> $\pm$ 0.15	<b>67.50</b> $\pm$ 0.11	67.12 $\pm$ 0.19
5000	68.41 $\pm$ 0.18	<b>69.09</b> $\pm$ 0.21	<b>68.93</b> $\pm$ 0.17	<b>68.96</b> $\pm$ 0.15	68.76 $\pm$ 0.20	68.81 $\pm$ 0.11	68.86 $\pm$ 0.11

Table 10: **Average Accumulated Accuracy (AAA)**, averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best aaa is in bold; results within the standard error of the best are also bolded.

(a) **TinyImageNet** with the **ER** framework (AAA).

$\mathcal{M}$	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	20.97 $\pm$ 0.09	<b>21.34</b> $\pm$ 0.07	21.16 $\pm$ 0.04	21.10 $\pm$ 0.04	21.13 $\pm$ 0.12	21.12 $\pm$ 0.04	21.14 $\pm$ 0.07
400	20.91 $\pm$ 0.13	21.44 $\pm$ 0.08	21.18 $\pm$ 0.09	<b>21.57</b> $\pm$ 0.04	<b>21.55</b> $\pm$ 0.05	21.32 $\pm$ 0.08	<b>21.28</b> $\pm$ 0.09
600	21.03 $\pm$ 0.07	<b>22.11</b> $\pm$ 0.11	21.47 $\pm$ 0.06	21.98 $\pm$ 0.08	21.96 $\pm$ 0.08	21.62 $\pm$ 0.05	21.78 $\pm$ 0.08
1000	22.17 $\pm$ 0.10	<b>23.32</b> $\pm$ 0.16	22.55 $\pm$ 0.10	<b>23.36</b> $\pm$ 0.12	<b>23.24</b> $\pm$ 0.05	22.99 $\pm$ 0.10	22.91 $\pm$ 0.08
2000	25.38 $\pm$ 0.09	27.44 $\pm$ 0.06	26.62 $\pm$ 0.09	<b>28.12</b> $\pm$ 0.19	27.59 $\pm$ 0.07	26.36 $\pm$ 0.05	26.78 $\pm$ 0.15
4000	33.17 $\pm$ 0.29	<b>35.66</b> $\pm$ 0.29	34.52 $\pm$ 0.12	<b>35.77</b> $\pm$ 0.20	35.22 $\pm$ 0.20	33.62 $\pm$ 0.24	34.17 $\pm$ 0.24
6000	40.14 $\pm$ 0.22	<b>41.61</b> $\pm$ 0.20	40.53 $\pm$ 0.26	41.41 $\pm$ 0.07	40.23 $\pm$ 0.23	39.82 $\pm$ 0.21	39.97 $\pm$ 0.23

(b) **TinyImageNet** with the **ER-ACE** framework (AAA).

$\mathcal{M}$	Herding	TEAL
	Model Based	Model Based
200	28.46 $\pm$ 0.12	29.41 $\pm$ 0.16
400	29.96 $\pm$ 0.22	31.30 $\pm$ 0.08
600	31.07 $\pm$ 0.18	32.53 $\pm$ 0.11
1000	32.54 $\pm$ 0.23	34.53 $\pm$ 0.08
2000	35.28 $\pm$ 0.20	37.59 $\pm$ 0.19
4000	39.15 $\pm$ 0.12	41.29 $\pm$ 0.13
6000	42.14 $\pm$ 0.19	43.38 $\pm$ 0.24

(c) **CIFAR-100** with the **MIR** framework (AAA).

$\mathcal{M}$	Herding	RM
	Model Based	Model Based
100	39.23 $\pm$ 0.47	37.69 $\pm$ 0.23
300	45.39 $\pm$ 0.16	43.88 $\pm$ 0.23
500	48.74 $\pm$ 0.19	48.04 $\pm$ 0.21
1000	55.64 $\pm$ 0.28	54.31 $\pm$ 0.14
2000	61.74 $\pm$ 0.21	60.66 $\pm$ 0.28
4000	67.94 $\pm$ 0.23	66.17 $\pm$ 0.08
5000	69.56 $\pm$ 0.21	68.22 $\pm$ 0.06

Table 11: **Final Averaged Accuracy (FAA)**, averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best AAA is in bold; results within the standard error of the best are also bolded.

(a) CIFAR-100 with the ER-ACE framework.

$ \mathcal{M} $	<i>MERS MaxHerding</i>			
	VICReg	VICReg + Model Based	DINOv2	DINOv2 + Model Based
100	25.33 $\pm$ 0.39	28.96 $\pm$ 0.31	27.96 $\pm$ 0.22	<b>29.80</b> $\pm$ 0.19
300	37.62 $\pm$ 0.49	<b>39.61</b> $\pm$ 0.21	38.01 $\pm$ 0.20	<b>39.52</b> $\pm$ 0.15
500	41.93 $\pm$ 0.40	<b>44.23</b> $\pm$ 0.27	42.68 $\pm$ 0.13	<b>43.95</b> $\pm$ 0.33
1000	48.52 $\pm$ 0.25	<b>50.03</b> $\pm$ 0.27	48.40 $\pm$ 0.10	<b>49.86</b> $\pm$ 0.08
2000	53.93 $\pm$ 0.18	54.82 $\pm$ 0.17	54.27 $\pm$ 0.18	<b>55.13</b> $\pm$ 0.09
4000	<b>59.77</b> $\pm$ 0.31	<b>59.82</b> $\pm$ 0.18	59.44 $\pm$ 0.12	59.44 $\pm$ 0.19
5000	<b>60.82</b> $\pm$ 0.25	<b>60.84</b> $\pm$ 0.19	60.30 $\pm$ 0.27	<b>60.73</b> $\pm$ 0.19

(b) CIFAR-100 with the ER-ACE framework (FAA).

$ \mathcal{M} $	<i>MERS MaxHerding</i>			
	VICReg	VICReg + Model Based	DINOv2	DINOv2 + Model Based
100	25.33 $\pm$ 0.39	28.96 $\pm$ 0.31	27.96 $\pm$ 0.22	<b>29.80</b> $\pm$ 0.19
300	37.62 $\pm$ 0.49	<b>39.61</b> $\pm$ 0.21	38.01 $\pm$ 0.20	<b>39.52</b> $\pm$ 0.15
500	41.93 $\pm$ 0.40	<b>44.23</b> $\pm$ 0.27	42.68 $\pm$ 0.13	<b>43.95</b> $\pm$ 0.33
1000	48.52 $\pm$ 0.25	<b>50.03</b> $\pm$ 0.27	48.40 $\pm$ 0.10	<b>49.86</b> $\pm$ 0.08
2000	53.93 $\pm$ 0.18	54.82 $\pm$ 0.17	54.27 $\pm$ 0.18	<b>55.13</b> $\pm$ 0.09
4000	<b>59.77</b> $\pm$ 0.31	<b>59.82</b> $\pm$ 0.18	59.44 $\pm$ 0.12	59.44 $\pm$ 0.19
5000	<b>60.82</b> $\pm$ 0.25	<b>60.84</b> $\pm$ 0.19	60.30 $\pm$ 0.27	<b>60.73</b> $\pm$ 0.19

Table 12: **Average Accumulated Accuracy (AAA)**, averaged over 5 independent runs (mean  $\pm$  standard error). For each buffer size, the best AAA is in bold; results within the standard error of the best are also bolded.

(a) CIFAR-100 with the ER-ACE framework (AAA).

$ \mathcal{M} $	<i>MERS MaxHerding</i>			
	VICReg	VICReg + Model Based	DINOv2	DINOv2 + Model Based
100	45.07 $\pm$ 0.27	<b>47.02</b> $\pm$ 0.15	45.75 $\pm$ 0.17	46.54 $\pm$ 0.11
300	53.70 $\pm$ 0.29	<b>55.25</b> $\pm$ 0.13	54.58 $\pm$ 0.16	55.01 $\pm$ 0.24
500	57.84 $\pm$ 0.36	<b>58.54</b> $\pm$ 0.32	58.13 $\pm$ 0.11	<b>58.54</b> $\pm$ 0.18
1000	62.50 $\pm$ 0.16	<b>63.15</b> $\pm$ 0.33	62.33 $\pm$ 0.31	<b>62.86</b> $\pm$ 0.22
2000	66.41 $\pm$ 0.19	<b>66.99</b> $\pm$ 0.10	66.55 $\pm$ 0.25	<b>67.13</b> $\pm$ 0.23
4000	<b>70.62</b> $\pm$ 0.26	<b>70.47</b> $\pm$ 0.14	<b>70.52</b> $\pm$ 0.07	69.94 $\pm$ 0.17
5000	<b>71.56</b> $\pm$ 0.37	<b>71.22</b> $\pm$ 0.09	71.11 $\pm$ 0.20	<b>71.42</b> $\pm$ 0.20

(b) CIFAR-100 with the ER-ACE framework (AAA).

$ \mathcal{M} $	<i>MERS ProbCover</i>			
	VICReg	VICReg + Model Based	DINOv2	DINOv2 + Model Based
100	44.55 $\pm$ 0.17	<b>45.95</b> $\pm$ 0.25	44.51 $\pm$ 0.10	<b>45.74</b> $\pm$ 0.22
300	53.71 $\pm$ 0.48	<b>54.72</b> $\pm$ 0.22	53.89 $\pm$ 0.18	<b>54.64</b> $\pm$ 0.19
500	57.22 $\pm$ 0.21	<b>58.62</b> $\pm$ 0.40	57.67 $\pm$ 0.19	<b>58.51</b> $\pm$ 0.15
1000	62.58 $\pm$ 0.20	<b>63.35</b> $\pm$ 0.39	61.61 $\pm$ 0.14	<b>62.97</b> $\pm$ 0.15
4000	70.37 $\pm$ 0.17	<b>71.10</b> $\pm$ 0.05	70.63 $\pm$ 0.22	70.73 $\pm$ 0.19
5000	<b>71.90</b> $\pm$ 0.29	<b>71.83</b> $\pm$ 0.36	71.61 $\pm$ 0.15	71.60 $\pm$ 0.20

Table 13: **FAA** for the class order defined by **seed 42**, averaged over 5 independent runs (mean  $\pm$  standard error). Several sample-selection strategies and embedding spaces are compared across multiple replay-buffer sizes ( $|\mathcal{M}|$ ).

(a) FAA on CIFAR-100 ER ACE.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	20.79 $\pm$ 0.27	29.81 $\pm$ 0.20	27.82 $\pm$ 0.25	29.35 $\pm$ 0.24	29.42 $\pm$ 0.11	29.20 $\pm$ 0.33	29.89 $\pm$ 0.22
300	31.76 $\pm$ 0.07	38.84 $\pm$ 0.19	37.78 $\pm$ 0.14	39.47 $\pm$ 0.28	38.16 $\pm$ 0.35	38.73 $\pm$ 0.36	39.60 $\pm$ 0.25
500	35.80 $\pm$ 0.32	42.46 $\pm$ 0.23	42.25 $\pm$ 0.19	43.28 $\pm$ 0.23	42.72 $\pm$ 0.23	42.82 $\pm$ 0.21	43.71 $\pm$ 0.21
1000	42.27 $\pm$ 0.22	47.69 $\pm$ 0.23	48.11 $\pm$ 0.19	48.98 $\pm$ 0.27	47.77 $\pm$ 0.21	48.89 $\pm$ 0.27	50.00 $\pm$ 0.16
2000	49.41 $\pm$ 0.18	52.99 $\pm$ 0.07	53.53 $\pm$ 0.24	54.17 $\pm$ 0.19	53.18 $\pm$ 0.20	54.20 $\pm$ 0.28	54.80 $\pm$ 0.19
4000	55.32 $\pm$ 0.24	58.03 $\pm$ 0.11	58.79 $\pm$ 0.24	59.28 $\pm$ 0.18	58.52 $\pm$ 0.20	58.56 $\pm$ 0.34	59.03 $\pm$ 0.19
5000	57.96 $\pm$ 0.21	60.10 $\pm$ 0.11	60.79 $\pm$ 0.19	60.84 $\pm$ 0.27	59.85 $\pm$ 0.18	59.90 $\pm$ 0.11	60.07 $\pm$ 0.13

(b) FAA on CIFAR-100 ER.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	10.50 $\pm$ 0.14	13.02 $\pm$ 0.11	11.44 $\pm$ 0.13	12.18 $\pm$ 0.08	12.53 $\pm$ 0.18	12.24 $\pm$ 0.07	12.79 $\pm$ 0.12
300	14.67 $\pm$ 0.24	20.32 $\pm$ 0.26	19.01 $\pm$ 0.34	20.33 $\pm$ 0.21	19.62 $\pm$ 0.17	18.83 $\pm$ 0.56	19.52 $\pm$ 0.33
500	19.86 $\pm$ 0.31	25.37 $\pm$ 0.18	23.68 $\pm$ 0.35	25.01 $\pm$ 0.44	25.73 $\pm$ 0.22	24.25 $\pm$ 0.22	25.74 $\pm$ 0.55
1000	28.48 $\pm$ 0.22	34.37 $\pm$ 0.33	33.62 $\pm$ 0.29	35.18 $\pm$ 0.21	34.54 $\pm$ 0.34	34.43 $\pm$ 0.35	35.40 $\pm$ 0.20
2000	40.45 $\pm$ 0.23	43.84 $\pm$ 0.32	44.34 $\pm$ 0.31	45.38 $\pm$ 0.28	44.62 $\pm$ 0.19	44.76 $\pm$ 0.37	45.58 $\pm$ 0.40
4000	51.23 $\pm$ 0.22	53.91 $\pm$ 0.20	54.37 $\pm$ 0.20	54.97 $\pm$ 0.27	54.69 $\pm$ 0.32	54.01 $\pm$ 0.16	54.81 $\pm$ 0.16
5000	55.03 $\pm$ 0.20	56.75 $\pm$ 0.13	57.16 $\pm$ 0.25	57.79 $\pm$ 0.22	56.67 $\pm$ 0.23	56.49 $\pm$ 0.21	57.06 $\pm$ 0.23

(c) FAA on TinyImageNet ER ACE.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	11.89 $\pm$ 0.13	13.95 $\pm$ 0.17	12.58 $\pm$ 0.05	13.33 $\pm$ 0.09	13.54 $\pm$ 0.07	13.50 $\pm$ 0.14	13.91 $\pm$ 0.25
400	13.27 $\pm$ 0.12	15.69 $\pm$ 0.12	14.33 $\pm$ 0.12	15.16 $\pm$ 0.11	14.72 $\pm$ 0.21	15.00 $\pm$ 0.10	15.35 $\pm$ 0.19
600	13.47 $\pm$ 0.08	16.44 $\pm$ 0.06	15.42 $\pm$ 0.19	16.64 $\pm$ 0.24	15.71 $\pm$ 0.18	16.07 $\pm$ 0.10	16.37 $\pm$ 0.31
1000	14.50 $\pm$ 0.16	18.16 $\pm$ 0.19	16.99 $\pm$ 0.19	18.44 $\pm$ 0.11	17.51 $\pm$ 0.16	17.26 $\pm$ 0.15	17.89 $\pm$ 0.12
2000	16.59 $\pm$ 0.15	20.50 $\pm$ 0.21	19.71 $\pm$ 0.19	21.03 $\pm$ 0.09	20.08 $\pm$ 0.23	19.50 $\pm$ 0.20	20.26 $\pm$ 0.27
4000	19.11 $\pm$ 0.13	23.09 $\pm$ 0.15	22.94 $\pm$ 0.17	24.45 $\pm$ 0.20	23.18 $\pm$ 0.21	22.43 $\pm$ 0.33	22.94 $\pm$ 0.20
6000	22.57 $\pm$ 0.06	25.41 $\pm$ 0.20	25.42 $\pm$ 0.30	26.40 $\pm$ 0.26	25.66 $\pm$ 0.19	24.66 $\pm$ 0.18	25.02 $\pm$ 0.15

Table 14: **FAA** for the class order defined by **seed 35**, averaged over 5 independent runs (mean  $\pm$  standard error). Several sample-selection strategies and embedding spaces are compared across multiple replay-buffer sizes ( $|\mathcal{M}|$ ).

(a) FAA on CIFAR-100 **ER ACE**.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	20.59 $\pm$ 0.23	27.64 $\pm$ 0.44	25.67 $\pm$ 0.45	27.42 $\pm$ 0.32	28.10 $\pm$ 0.30	28.30 $\pm$ 0.44	29.35 $\pm$ 0.25
300	28.61 $\pm$ 0.05	37.84 $\pm$ 0.14	36.90 $\pm$ 0.31	38.88 $\pm$ 0.21	37.70 $\pm$ 0.34	37.63 $\pm$ 0.31	39.19 $\pm$ 0.19
500	35.30 $\pm$ 0.19	42.23 $\pm$ 0.19	41.75 $\pm$ 0.18	43.55 $\pm$ 0.23	42.02 $\pm$ 0.25	42.42 $\pm$ 0.13	44.02 $\pm$ 0.38
1000	41.99 $\pm$ 0.16	48.18 $\pm$ 0.22	48.36 $\pm$ 0.15	48.96 $\pm$ 0.25	47.89 $\pm$ 0.23	48.71 $\pm$ 0.30	49.63 $\pm$ 0.30
2000	49.00 $\pm$ 0.23	53.20 $\pm$ 0.22	53.51 $\pm$ 0.08	54.44 $\pm$ 0.28	53.67 $\pm$ 0.22	54.30 $\pm$ 0.14	55.11 $\pm$ 0.10
4000	56.89 $\pm$ 0.11	59.01 $\pm$ 0.27	59.18 $\pm$ 0.12	59.73 $\pm$ 0.15	59.30 $\pm$ 0.05	59.23 $\pm$ 0.14	59.67 $\pm$ 0.10
5000	58.75 $\pm$ 0.22	60.45 $\pm$ 0.18	60.65 $\pm$ 0.11	61.40 $\pm$ 0.22	60.17 $\pm$ 0.11	60.37 $\pm$ 0.09	60.94 $\pm$ 0.07

(b) FAA on CIFAR-100 **ER**.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	9.95 $\pm$ 0.07	11.45 $\pm$ 0.11	10.32 $\pm$ 0.06	11.17 $\pm$ 0.17	11.19 $\pm$ 0.05	11.05 $\pm$ 0.18	11.51 $\pm$ 0.01
300	13.71 $\pm$ 0.09	18.87 $\pm$ 0.11	17.16 $\pm$ 0.14	18.71 $\pm$ 0.25	18.10 $\pm$ 0.34	18.01 $\pm$ 0.22	18.71 $\pm$ 0.25
500	17.41 $\pm$ 0.38	23.75 $\pm$ 0.39	22.37 $\pm$ 0.34	24.66 $\pm$ 0.14	24.11 $\pm$ 0.15	23.40 $\pm$ 0.15	24.66 $\pm$ 0.29
1000	27.44 $\pm$ 0.48	33.50 $\pm$ 0.17	32.51 $\pm$ 0.48	33.99 $\pm$ 0.27	33.70 $\pm$ 0.20	33.27 $\pm$ 0.16	34.64 $\pm$ 0.31
2000	39.78 $\pm$ 0.30	43.73 $\pm$ 0.01	43.74 $\pm$ 0.30	44.02 $\pm$ 0.20	44.06 $\pm$ 0.30	44.01 $\pm$ 0.20	45.22 $\pm$ 0.16

(c) FAA on TinyImagenet **ER ACE**.

Buffer	ER (vanilla)	<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	11.39 $\pm$ 0.10	13.23 $\pm$ 0.12	12.22 $\pm$ 0.14	12.75 $\pm$ 0.13	13.11 $\pm$ 0.08	12.71 $\pm$ 0.11	12.95 $\pm$ 0.09
400	11.98 $\pm$ 0.24	15.09 $\pm$ 0.21	13.70 $\pm$ 0.16	14.71 $\pm$ 0.24	13.84 $\pm$ 0.15	14.12 $\pm$ 0.21	14.51 $\pm$ 0.16
600	12.90 $\pm$ 0.13	16.18 $\pm$ 0.12	14.68 $\pm$ 0.08	15.78 $\pm$ 0.18	14.97 $\pm$ 0.19	15.48 $\pm$ 0.13	15.14 $\pm$ 0.06
1000	14.14 $\pm$ 0.09	17.67 $\pm$ 0.26	16.21 $\pm$ 0.22	17.47 $\pm$ 0.15	16.61 $\pm$ 0.11	16.32 $\pm$ 0.18	16.77 $\pm$ 0.10
2000	15.94 $\pm$ 0.16	19.88 $\pm$ 0.24	18.60 $\pm$ 0.23	20.42 $\pm$ 0.29	19.70 $\pm$ 0.34	19.01 $\pm$ 0.14	19.21 $\pm$ 0.22
4000	19.42 $\pm$ 0.22	22.86 $\pm$ 0.13	23.05 $\pm$ 0.35	24.08 $\pm$ 0.07	22.80 $\pm$ 0.12	21.84 $\pm$ 0.28	21.84 $\pm$ 0.23
6000	22.05 $\pm$ 0.25	25.98 $\pm$ 0.34	25.63 $\pm$ 0.30	26.53 $\pm$ 0.13	25.14 $\pm$ 0.23	24.43 $\pm$ 0.28	25.23 $\pm$ 0.25

Table 15: AAA for the class order defined by **seed 42**, averaged over 5 independent runs (mean  $\pm$  standard error). Several sample-selection strategies and embedding spaces are compared across multiple replay-buffer sizes ( $|\mathcal{M}|$ ).

(a) AAA on CIFAR-100 ER ACE.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based		Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	39.94 $\pm$ 0.09		46.09 $\pm$ 0.08	45.60 $\pm$ 0.09	46.90 $\pm$ 0.18	46.21 $\pm$ 0.18	46.17 $\pm$ 0.11	46.78 $\pm$ 0.27
300	49.19 $\pm$ 0.10		53.33 $\pm$ 0.07	53.62 $\pm$ 0.09	54.30 $\pm$ 0.13	53.46 $\pm$ 0.30	53.92 $\pm$ 0.19	54.68 $\pm$ 0.14
500	52.85 $\pm$ 0.08		56.55 $\pm$ 0.15	56.76 $\pm$ 0.12	57.07 $\pm$ 0.26	56.52 $\pm$ 0.12	57.34 $\pm$ 0.10	57.77 $\pm$ 0.07
1000	57.60 $\pm$ 0.13		60.65 $\pm$ 0.09	60.90 $\pm$ 0.10	61.46 $\pm$ 0.06	60.60 $\pm$ 0.23	61.25 $\pm$ 0.06	61.97 $\pm$ 0.17
2000	62.35 $\pm$ 0.14		64.36 $\pm$ 0.20	64.85 $\pm$ 0.11	64.91 $\pm$ 0.12	64.29 $\pm$ 0.13	65.01 $\pm$ 0.08	65.26 $\pm$ 0.15
4000	66.73 $\pm$ 0.17		68.22 $\pm$ 0.09	68.39 $\pm$ 0.17	68.67 $\pm$ 0.14	68.20 $\pm$ 0.14	67.97 $\pm$ 0.08	68.16 $\pm$ 0.07
5000	68.32 $\pm$ 0.16		69.36 $\pm$ 0.08	69.92 $\pm$ 0.15	69.80 $\pm$ 0.17	69.40 $\pm$ 0.14	69.07 $\pm$ 0.04	69.20 $\pm$ 0.10

(b) AAA on CIFAR-100 ER.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based		Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
100	28.19 $\pm$ 0.11		30.89 $\pm$ 0.07	29.72 $\pm$ 0.24	30.51 $\pm$ 0.16	30.31 $\pm$ 0.11	30.37 $\pm$ 0.15	30.49 $\pm$ 0.15
300	34.42 $\pm$ 0.42		38.55 $\pm$ 0.37	38.12 $\pm$ 0.33	39.30 $\pm$ 0.13	38.44 $\pm$ 0.25	37.92 $\pm$ 0.46	38.37 $\pm$ 0.42
500	40.31 $\pm$ 0.21		43.90 $\pm$ 0.27	42.60 $\pm$ 0.39	43.55 $\pm$ 0.34	43.74 $\pm$ 0.09	43.58 $\pm$ 0.37	44.68 $\pm$ 0.31
1000	48.62 $\pm$ 0.24		51.78 $\pm$ 0.20	51.86 $\pm$ 0.22	52.58 $\pm$ 0.35	51.67 $\pm$ 0.37	52.27 $\pm$ 0.31	52.71 $\pm$ 0.25
2000	58.66 $\pm$ 0.23		59.70 $\pm$ 0.47	60.57 $\pm$ 0.20	61.48 $\pm$ 0.18	60.68 $\pm$ 0.13	60.73 $\pm$ 0.37	60.49 $\pm$ 0.28
4000	66.49 $\pm$ 0.18		67.67 $\pm$ 0.23	67.41 $\pm$ 0.17	68.30 $\pm$ 0.32	68.10 $\pm$ 0.12	67.35 $\pm$ 0.27	68.12 $\pm$ 0.11
5000	68.89 $\pm$ 0.17		69.58 $\pm$ 0.18	69.53 $\pm$ 0.22	70.13 $\pm$ 0.21	69.02 $\pm$ 0.21	69.21 $\pm$ 0.32	69.43 $\pm$ 0.05

(c) AAA on TinyImageNet ER ACE.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based		Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>
200	25.92 $\pm$ 0.07		28.32 $\pm$ 0.06	27.43 $\pm$ 0.10	28.17 $\pm$ 0.11	27.91 $\pm$ 0.09	27.93 $\pm$ 0.09	28.20 $\pm$ 0.08
400	27.78 $\pm$ 0.16		30.60 $\pm$ 0.13	29.61 $\pm$ 0.07	30.50 $\pm$ 0.09	29.73 $\pm$ 0.05	29.94 $\pm$ 0.10	30.10 $\pm$ 0.10
600	28.96 $\pm$ 0.07		31.60 $\pm$ 0.13	30.94 $\pm$ 0.09	31.82 $\pm$ 0.08	31.18 $\pm$ 0.09	31.40 $\pm$ 0.15	31.56 $\pm$ 0.13
1000	30.49 $\pm$ 0.05		33.60 $\pm$ 0.15	33.05 $\pm$ 0.13	34.08 $\pm$ 0.10	33.43 $\pm$ 0.12	33.12 $\pm$ 0.20	33.22 $\pm$ 0.11
2000	33.23 $\pm$ 0.13		36.09 $\pm$ 0.13	35.84 $\pm$ 0.11	36.86 $\pm$ 0.05	36.08 $\pm$ 0.14	35.51 $\pm$ 0.14	36.04 $\pm$ 0.20
4000	36.95 $\pm$ 0.13		39.32 $\pm$ 0.13	39.06 $\pm$ 0.10	39.87 $\pm$ 0.12	39.10 $\pm$ 0.12	38.47 $\pm$ 0.09	38.57 $\pm$ 0.12
6000	39.66 $\pm$ 0.12		40.90 $\pm$ 0.12	41.19 $\pm$ 0.10	41.67 $\pm$ 0.08	40.94 $\pm$ 0.15	40.08 $\pm$ 0.10	40.26 $\pm$ 0.16

Table 16: **AAA** for the class order defined by **seed 35**, averaged over 5 independent runs (mean  $\pm$  standard error). Several sample-selection strategies and embedding spaces are compared across multiple replay-buffer sizes ( $|\mathcal{M}|$ ).

(a) AAA on CIFAR-100 **ER ACE**.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>	
100	41.79 $\pm$ 0.14	47.76 $\pm$ 0.18	47.34 $\pm$ 0.14	48.44 $\pm$ 0.08	48.51 $\pm$ 0.17	48.63 $\pm$ 0.15	49.18 $\pm$ 0.11	
300	51.26 $\pm$ 0.11	56.10 $\pm$ 0.12	56.54 $\pm$ 0.18	57.50 $\pm$ 0.13	56.33 $\pm$ 0.17	57.15 $\pm$ 0.14	57.78 $\pm$ 0.13	
500	56.12 $\pm$ 0.28	59.79 $\pm$ 0.14	60.32 $\pm$ 0.15	61.35 $\pm$ 0.08	60.28 $\pm$ 0.12	60.63 $\pm$ 0.12	61.45 $\pm$ 0.11	
1000	61.84 $\pm$ 0.20	64.49 $\pm$ 0.20	65.41 $\pm$ 0.12	65.54 $\pm$ 0.14	64.56 $\pm$ 0.12	65.04 $\pm$ 0.17	65.79 $\pm$ 0.22	
2000	66.46 $\pm$ 0.05	68.70 $\pm$ 0.17	68.92 $\pm$ 0.18	69.02 $\pm$ 0.15	68.87 $\pm$ 0.12	69.22 $\pm$ 0.09	69.29 $\pm$ 0.17	
4000	71.57 $\pm$ 0.10	72.34 $\pm$ 0.17	72.52 $\pm$ 0.11	73.00 $\pm$ 0.03	72.53 $\pm$ 0.06	72.30 $\pm$ 0.21	72.44 $\pm$ 0.13	
5000	72.95 $\pm$ 0.09	73.61 $\pm$ 0.10	73.48 $\pm$ 0.09	74.22 $\pm$ 0.12	73.14 $\pm$ 0.09	73.36 $\pm$ 0.14	73.66 $\pm$ 0.03	

(b) AAA on CIFAR-100 **ER**.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>	
100	29.82 $\pm$ 0.13	32.42 $\pm$ 0.04	31.62 $\pm$ 0.12	32.58 $\pm$ 0.20	32.28 $\pm$ 0.08	32.09 $\pm$ 0.16	32.58 $\pm$ 0.17	
300	37.89 $\pm$ 0.08	42.18 $\pm$ 0.13	41.20 $\pm$ 0.14	42.32 $\pm$ 0.15	41.33 $\pm$ 0.18	41.74 $\pm$ 0.14	42.39 $\pm$ 0.04	
500	43.07 $\pm$ 0.17	47.15 $\pm$ 0.13	47.20 $\pm$ 0.09	48.48 $\pm$ 0.09	47.70 $\pm$ 0.09	47.64 $\pm$ 0.12	48.52 $\pm$ 0.19	
1000	52.56 $\pm$ 0.13	55.93 $\pm$ 0.05	55.92 $\pm$ 0.29	56.60 $\pm$ 0.31	56.30 $\pm$ 0.08	56.35 $\pm$ 0.20	57.20 $\pm$ 0.11	
2000	62.59 $\pm$ 0.15	64.42 $\pm$ 0.21	64.67 $\pm$ 0.11	64.56 $\pm$ 0.04	64.45 $\pm$ 0.13	64.52 $\pm$ 0.07	64.96 $\pm$ 0.10	

(c) AAA on TinyImageNet **ER-ACE**.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>	
200	26.65 $\pm$ 0.04	28.49 $\pm$ 0.11	27.57 $\pm$ 0.06	28.24 $\pm$ 0.09	28.16 $\pm$ 0.19	28.19 $\pm$ 0.05	28.31 $\pm$ 0.12	
400	28.01 $\pm$ 0.07	30.93 $\pm$ 0.23	29.82 $\pm$ 0.13	30.79 $\pm$ 0.15	30.06 $\pm$ 0.02	30.28 $\pm$ 0.12	30.49 $\pm$ 0.12	
600	29.02 $\pm$ 0.12	32.01 $\pm$ 0.09	31.05 $\pm$ 0.16	32.21 $\pm$ 0.14	31.76 $\pm$ 0.18	31.45 $\pm$ 0.09	31.66 $\pm$ 0.08	
1000	31.03 $\pm$ 0.15	33.92 $\pm$ 0.13	32.97 $\pm$ 0.07	34.43 $\pm$ 0.16	33.52 $\pm$ 0.14	33.15 $\pm$ 0.16	33.11 $\pm$ 0.12	
2000	34.01 $\pm$ 0.16	36.25 $\pm$ 0.22	35.97 $\pm$ 0.17	36.96 $\pm$ 0.10	36.65 $\pm$ 0.11	36.11 $\pm$ 0.18	36.11 $\pm$ 0.19	
4000	37.83 $\pm$ 0.13	39.51 $\pm$ 0.11	39.78 $\pm$ 0.21	40.28 $\pm$ 0.12	39.37 $\pm$ 0.13	38.70 $\pm$ 0.10	39.05 $\pm$ 0.11	
6000	40.15 $\pm$ 0.24	42.05 $\pm$ 0.26	41.66 $\pm$ 0.15	42.57 $\pm$ 0.15	41.37 $\pm$ 0.14	40.76 $\pm$ 0.24	41.43 $\pm$ 0.04	

(d) AAA on TinyImageNet **ER**.

Buffer	ER (vanilla)		<i>MERS ProbCover</i>			<i>MERS MaxHerding</i>		
	Model Based	Model Based	SimCLR	<i>MERS</i>	Model Based	SimCLR	<i>MERS</i>	
200	21.09 $\pm$ 0.10	21.06 $\pm$ 0.04	21.03 $\pm$ 0.02	21.00 $\pm$ 0.09	21.12 $\pm$ 0.13	21.22 $\pm$ 0.02	21.16 $\pm$ 0.12	
400	20.94 $\pm$ 0.07	21.48 $\pm$ 0.11	21.06 $\pm$ 0.04	21.59 $\pm$ 0.06	21.56 $\pm$ 0.05	21.34 $\pm$ 0.05	21.33 $\pm$ 0.09	
600	21.16 $\pm$ 0.10	22.15 $\pm$ 0.09	21.59 $\pm$ 0.09	21.91 $\pm$ 0.11	22.17 $\pm$ 0.10	21.75 $\pm$ 0.08	21.78 $\pm$ 0.06	
1000	21.91 $\pm$ 0.15	23.30 $\pm$ 0.05	22.91 $\pm$ 0.15	23.64 $\pm$ 0.12	23.30 $\pm$ 0.10	22.82 $\pm$ 0.09	22.83 $\pm$ 0.08	
2000	25.57 $\pm$ 0.14	27.72 $\pm$ 0.14	26.72 $\pm$ 0.10	27.64 $\pm$ 0.09	27.25 $\pm$ 0.15	26.46 $\pm$ 0.16	27.03 $\pm$ 0.10	
4000	33.03 $\pm$ 0.11	35.37 $\pm$ 0.30	34.41 $\pm$ 0.18	36.29 $\pm$ 0.15	35.24 $\pm$ 0.17	34.08 $\pm$ 0.08	34.45 $\pm$ 0.23	
6000	39.88 $\pm$ 0.15	41.56 $\pm$ 0.17	41.02 $\pm$ 0.17	41.70 $\pm$ 0.14	40.87 $\pm$ 0.24	39.76 $\pm$ 0.18	40.65 $\pm$ 0.07	