



EFFICIENT SINGLE-STEP FRAMEWORK FOR INCREMENTAL CLASS LEARNING IN NEURAL NETWORKS

A PREPRINT

 **Alejandro Dopico-Castro***

 Oscar Fontenla-Romero

 Bertha Guijarro-Berdiñas

 Amparo Alonso-Betanzos

Universidade da Coruña, CITIC

Facultade de Informática, Campus de Elviña s/n, A Coruña, Spain

{alejandro.dopico2, oscar.fontenla, berta.guijarro, amparo.alonso.betanzos}@udc.es

ABSTRACT

Incremental learning remains a critical challenge in machine learning, as models often struggle with catastrophic forgetting—the tendency to lose previously acquired knowledge when learning new information. These challenges are even more pronounced in resource-limited settings. Many existing Class Incremental Learning (CIL) methods achieve high accuracy by continually adapting their feature representations; however, they often require substantial computational resources and complex, iterative training procedures. This work introduces CIFNet (Class Incremental and Frugal Network), a novel CIL approach that addresses these limitations by offering a highly efficient and sustainable solution.

CIFNet’s key innovation lies in its novel integration of several existing, yet separately explored, components: a pre-trained and frozen feature extractor, a compressed data buffer, and an efficient non-iterative one-layer neural network for classification. A pre-trained and frozen feature extractor eliminates computationally expensive fine-tuning of the backbone. This, combined with a compressed buffer for efficient memory use, enables CIFNet to perform efficient class-incremental learning through a single-step optimization process on fixed features, minimizing computational overhead and training time without requiring multiple weight updates.

Experiments on benchmark datasets confirm that CIFNet effectively mitigates catastrophic forgetting at the classifier level, achieving high accuracy comparable to that of existing state-of-the-art methods, while substantially improving training efficiency and sustainability. CIFNet represents a significant advancement in making class-incremental learning more accessible and pragmatic in environments with limited resources, especially when strong pre-trained feature extractors are available.

Keywords Class Incremental Learning · Continual Learning · Lifelong Learning · Catastrophic Forgetting · Frugal AI · Green AI · Sustainable AI

1 Introduction

In the real world, knowledge acquisition is not a static or complete process from its inception. Just as humans continuously encounter and learn to recognise new categories of objects throughout their lives, artificial intelligence systems must adapt to evolving environments where new data classes emerge over time. This dynamic nature of learning presents a fundamental challenge in machine learning: how can we design systems that efficiently incorporate new categories of information without compromising their existing knowledge?.

Class incremental learning (CIL) [1, 2] addresses this challenge by allowing machine learning models to progressively learn new classes incrementally while maintaining performance on previously learned ones. This is critical

*Corresponding author: alejandro.dopico2@udc.es

in applications such as robotics, particularly in autonomous navigation systems, where new object categories (e.g., previously unencountered obstacles or traffic signs) may emerge, or in industrial automation, where robots must adapt to new product variants or environmental changes. Retraining models from scratch as new classes emerge is impractical due to memory, privacy, and computational constraints, especially in resource-constrained environments. CIL offers a more practical solution by allowing models to adapt incrementally, avoiding catastrophic forgetting [3] of previously learned information. This approach both reduces computational overhead and aligns with the way knowledge naturally accumulates in real-world scenarios. Beyond computational efficiency, CIL is critical for creating adaptive and sustainable Artificial Intelligence (AI) systems that evolve in response to human needs.

As said, incremental learning aims to build models that can adapt over time, incorporating new information while retaining previous knowledge. A widely adopted strategy to support this is the use of backbone architectures—deep neural networks that extract transferable features that can be reused across tasks. Despite significant progress, most state-of-the-art [4, 5, 6] methods rely on continually adapting the deep feature extractor (backbone). While this paradigm is effective for handling significant domain shifts, it often results in complex training procedures, high computational costs, and substantial energy consumption. These factors can be prohibitively expensive for many real-world applications, especially those deployed on edge devices or in environments with strict energy budgets. There is a critical need for CIL solutions that prioritize computational and energy efficiency without severely compromising accuracy. This work focuses on practical scenarios in which powerful, pre-trained feature extractors are readily available, providing a strong foundation for learning. This setup is common in many computer vision applications, and while we acknowledge the need for adaptability in more distinct domains, we leave such considerations for future work.

This work introduces CIFNet (**C**lass **I**ncremental and **F**rugal **N**etwork), a CIL approach that revisits the trade-off between feature adaptation and efficiency. Unlike conventional methods that continually update the entire neural network backbone, CIFNet relies on a pre-trained feature extractor that remains fixed throughout incremental learning. This choice is motivated by the established effectiveness of pre-trained models, particularly when new tasks are related to the pre-training domain [7]. By leveraging such stable representations, CIFNet shifts the computational focus from backbone adaptation to classifier learning, thereby improving efficiency and stability.

We evaluate CIFNet on the widely used CIFAR-100 [8] and ImageNet-100 [9] benchmarks. Results demonstrate that the method achieves competitive accuracy while significantly reducing computational and memory costs compared to state-of-the-art approaches that fine-tune the backbone.

The main contributions of this work are as follows:

- We introduce CIFNet, a Class-Incremental Learning framework designed for computational and energy efficiency by employing a frozen, pre-trained feature extractor.
- We propose a non-interactive classification layer specifically tailored for incremental learning on fixed features, enabling efficient single-step optimization and extending the applicability of such approaches to complex visual data.
- We justify the use of a compressed embedding buffer for replay and demonstrate its superior memory efficiency compared to raw image storage.
- We show that CIFNet achieves competitive performance on standard benchmarks while providing substantial improvements in training time, energy consumption, and memory usage.

The paper is organized as follows: Section 2 discusses related work, including catastrophic forgetting and state-of-the-art solutions. Section 3 introduces key concepts like incremental learning and replay strategies. Section 4 details our approach, and Section 5 presents a comparative analysis, including sustainability metrics like energy consumption and training time.

2 Related Work

Continual Learning (CL) in Machine Learning addresses the challenge of learning from a continuous stream of data while retaining knowledge of previous tasks [10]. A task² refers to a learning unit characterized by a set of data representing some specific classes. These tasks are presented sequentially, each time learning new classes, emulating human sequential knowledge acquisition. While vital for real-world applicability, CL presents significant challenges, primarily catastrophic forgetting, where models rapidly lose previously acquired knowledge when trained on new tasks [2]. This necessitates strategies that balance plasticity (acquiring new knowledge) with stability (retaining old knowledge). Continual learning includes different approaches to learning over time. In this context, CIL focuses on

²Alternatively referred to as ‘stage’ or ‘phase’.

scenarios where new classes are introduced progressively, without explicit task boundaries, where the model must recognise both old and new classes. Strategies to address this have been proposed by researchers, including memory-based methods, which store and replay past samples [2, 11], and dynamic architectures, which adapt their structure to accommodate new knowledge [12].

2.1 Replay Methods

Experience replay is a core and highly effective strategy in CIL to mitigate catastrophic forgetting. It involves storing a subset of past data (exemplars) and replaying these samples alongside current data during the training of new tasks [2, 13, 14]. This joint training on both old and new data helps the model explicitly revisit and retain its knowledge of previously learned classes, while simultaneously acquiring expertise in the new ones.

However, a limitation of traditional replay methods, especially in scenarios with large image datasets, is the considerable memory overhead associated with storing raw data. This can become prohibitive in resource-constrained environments or for long learning sequences. There are existing solutions to this approach, like storing extracted features [11] or employing gradient-based sample selection [15]. Generative Replay [16, 17, 18] offers an alternative by synthesizing samples, avoiding explicit data storage but introducing challenges related to generative model training and sample fidelity.

2.2 Dynamic Architectures

Dynamic expansion strategies involve adding neurons or layers to accommodate new tasks, effectively mitigating forgetting by dedicating new capacity for new knowledge. DEN [12] and Progressive Networks [19] expand the network while aiming to maintain prior knowledge through careful connections and regularization.

However, the most significant drawback of dynamic architectures is their uncontrolled model growth. As new tasks arrive, the model size continuously increases, leading to prohibitively large and computationally expensive models over time. This limits their applicability in scenarios with many tasks, resource-constrained environments (e.g., edge devices), or where inference speed is critical due to the ever-increasing model complexity. While methods like PackNet [20] and CPG [21] attempt to optimize capacity allocation within the feature extractor to balance adaptation and efficiency, they still involve modifying and managing the backbone’s architecture.

While these methods effectively adapt the backbone for new tasks, their reliance on continuous updates and architectural expansion can lead to substantial computational overhead and scalability issues, particularly in long-term or resource-constrained scenarios. An alternative direction explored in this work seeks to leverage frozen [22] representations for enhanced efficiency, inherently avoiding these limitations.

3 Preliminaries

In this section, we provide the background for understanding our approach. We first introduce the incremental learning framework and the challenge of catastrophic forgetting. Next, we review the strategy for mitigating catastrophic forgetting. Finally, we describe ROLANN, the classifier layer used in CIFNet.

3.1 Problem Formulation and Challenges in Class Incremental Learning

Class Incremental Learning (CIL) involves acquiring knowledge through a sequence of tasks, each one introducing a unique set of classes (see Figure 1). At each incremental step, the model receives data from a new task and no longer has access to the previous set of tasks. This makes the problem significantly challenging, as the model must classify all seen classes without task-specific clues.

Formally, consider a sequence of tasks T_1, T_2, \dots, T_k , where each T_i represents a disjoint set of classes, such that $T_i \cap T_j = \emptyset$ for $i \neq j$. Each task T_i has a dataset $D_i = (\mathcal{X}_i, \mathcal{Y}_i)$, where (x, y) pairs consist of an input sample x (e.g., an image) and its class label y . The model learns these tasks sequentially, and once a task is completed, its D is no longer available. The primary challenge is to enable the model to classify all classes seen so far—from the current task and all previous ones—without forgetting past knowledge.

Formally, at each incremental step k , a model f_θ parametrized by θ is trained on the dataset D_k . The objective is to learn a function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ that can classify samples from the union of all observed classes:

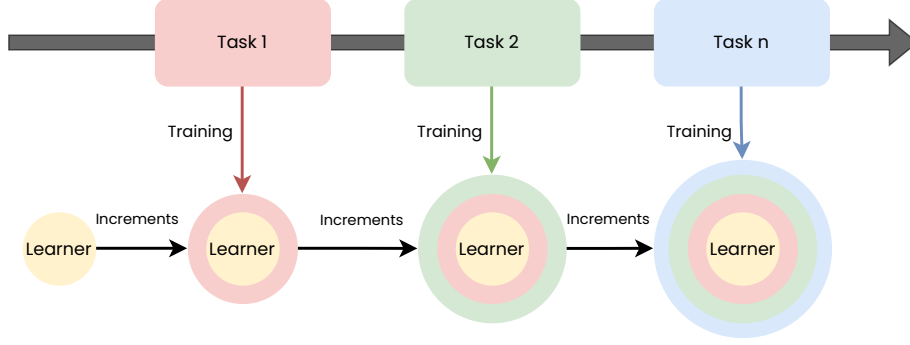


Figure 1: Class incremental learning setup: tasks arrive sequentially, each containing non-overlapping classes. The model learns these classes incrementally, expanding its knowledge over time.

$$\mathcal{Y} = \bigcup_{i=1}^k \mathcal{Y}_i. \quad (1)$$

The model is optimized by minimizing the expected loss, typically cross-entropy, on the current task’s data. A major challenge in this process is catastrophic forgetting (CF), where a neural network loses the knowledge of previous tasks when trained on new data. This occurs because standard training methods, like Stochastic Gradient Descent (SGD), update network parameters globally, focusing on the current task and overwriting those optimized for past knowledge.

To mitigate CF, replay-based methods use a small memory buffer B_{k-1} to store a subset of examples from past tasks. At each step k , the model is trained on a combination of the current task’s data (D_k) and the replayed samples from the buffer. This dual-training approach helps maintain the model’s stability (retaining old knowledge) while allowing for plasticity (acquiring new knowledge).

3.2 Regularized One-Layer Neural Network

In our proposed CIFNet architecture, the final classification layer is based on Regularized One-Layer Neural Network (ROLANN) [23], a non-iterative alternative to conventional linear layers trained via gradient descent. ROLANN optimizes the mean square error (MSE) before applying the activation function of the output neurons and incorporates L2 regularization, offering computational efficiency and suitability for large-scale or distributed learning tasks. This inherent characteristic is the key enabler for CIFNet’s single-step optimization process, drastically reducing the training time associated with incremental learning. The ability to directly compute weights without iterative updates makes ROLANN uniquely suited for highly efficient and scalable CIL applications.

Given an input $X \in \mathbb{R}^{m \times n}$, where m is the number of input variables and n is the number of samples, ROLANN computes the weights W using the following steps.

First, the method computes the product of the input matrix and a diagonal matrix F of activation function derivatives:

$$XF, \quad \text{where} \quad F = \text{diag}(f'(\bar{d}_1), f'(\bar{d}_2), \dots, f'(\bar{d}_n)). \quad (2)$$

where $f'(\bar{d})$ denotes the derivative of the activation function evaluated at the pre-activation outputs \bar{d} .

Next, Singular Value Decomposition (SVD) is applied to this product:

$$[U, S, \sim] = \text{SVD}(XF) \quad (3)$$

U and S containing the left singular vectors and singular values, respectively.

The next step involves computing the moment matrix M :

$$M = X \cdot (f' \odot f' \odot \bar{d}), \quad (4)$$

where \odot denotes the element-wise (Hadamard) product.

Finally, the weight matrix W is obtained by:

$$W = U \cdot (S \cdot S + \lambda I)^{-1} \cdot U^T \cdot M, \quad (5)$$

where λ is a regularization parameter that controls the weight decay, and I is the identity matrix. This closed-form solution avoids iterative optimization.

ROLANN naturally supports incremental and distributed learning by enabling the integration of new data partitions without retraining from scratch. Let (U_k, S_k, M_k) denote the knowledge extracted from previous partition k and (U_p, S_p, M_p) denote the knowledge extracted from a new partition p . New matrices $U_{k|p}, S_{k|p}$ containing Knowledge from both partitions are obtained by concatenating the existing and the new singular value decompositions as:

$$[U_{k|p}, S_{k|p}, \sim] = \text{SVD}([U_k S_k \| U_p S_p]), \quad (6)$$

where $\|$ denotes the horizontal concatenation of the matrices. The moment matrices are combined by:

$$M_{k|p} = M_k + M_p \quad (7)$$

Then, the updated weights are computed analogously:

$$W = U_{k|p} \cdot (S_{k|p} \cdot S_{k|p} + \lambda I)^{-1} \cdot U_{k|p}^T \cdot M_{k|p}. \quad (8)$$

This design allows ROLANN to handle streaming data and distributed datasets efficiently, making it a practical choice for federated or incremental scenarios.

Algorithm 2, provided in the supplementary material, summarizes the training procedure, including the steps for incremental and distributed learning. Note that the weights associated with each of the network outputs can be calculated independently.

4 Proposed Method

Our proposed method, CIFNet, aims to facilitate incremental learning of new classes while robustly preserving previously acquired knowledge, thus mitigating the phenomenon of catastrophic forgetting. In addition, the method is designed to achieve efficient training times and rapid inference. This approach comprises three main components: a frozen pre-trained visual backbone (F_θ) for feature extraction, a ROLANN classifier layer (G_ϕ), and an expansion buffer (B), as illustrated in Figure 2.

F_θ acts as a high-quality feature extractor, transforming input images $x \in \mathbb{R}^{h \times w \times c}$ into lower-dimensional, discriminative feature vectors $\mathcal{E} = F_\theta(x)$ where $\mathcal{E} \in \mathbb{R}^d$, $d \ll h \times w \times c$, and where h is the height, w the width, and c the number of channels of the input image. This dimensionality reduction serves two purposes: first, it captures essential visual patterns while removing redundant information, and second, it ensures computational efficiency during both training and inference phases. Through this transformation, the backbone provides a compact yet discriminative representation that serves as input to the subsequent classifier layer.

The ROLANN classifier layer G_ϕ operates on these extracted features, introducing an efficient approach to incremental learning that fundamentally differs from traditional gradient descent-based methods. Instead of requiring multiple iterative weight updates across numerous training epochs, ROLANN achieves network adaptation through a single direct weight update step, considerably reducing training time.

While the ROLANN was originally designed for incremental learning in data streams (i.e, new samples arriving sequentially), its inherent property of independent output neuron training makes it uniquely suitable for extension to class-incremental scenarios.

In this work, we propose a novel adaptation and extension of the ROLANN layer to address CIL. Unlike its initial formulation, our adapted layer dynamically expands its architecture by adding new output neurons specifically for each new class as they emerge. This allows the model to incrementally learn new categories while operating on the rich representations provided by a frozen backbone. Therefore, the ROLANN layer will incrementally expand its architecture by adding $|C_k|$ output neurons in each task k , where $|C_k|$ corresponds to the number of new classes introduced. At the

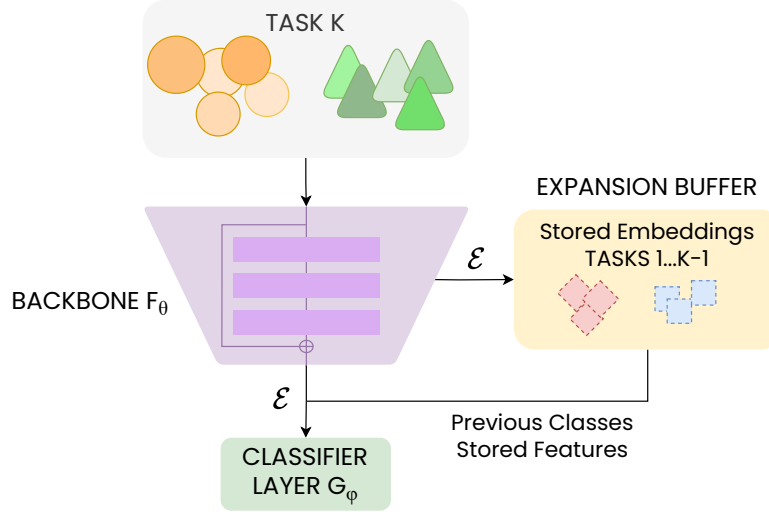


Figure 2: Architecture of the proposed model. Images of the task k are fed to the frozen feature extractor F_θ , generating embeddings \mathcal{E} that update the expansion buffer B_k . Concurrently, the ROLANN classifier layer G_ϕ learns to classify new task features as well as the sample embeddings of the previous classes.

beginning of each task, the weights w and matrices M, U and S of these new neurons are empty, which means that they do not contain any prior knowledge, including information from classes seen previously. Subsequently, these matrices and the corresponding weights will be updated with the data received on the new classes in the training set.

A key challenge with independently trained sigmoid neurons is their tendency to over-activate for recently introduced classes, leading to biased predictions. Our framework mitigates this by using the expansion buffer B to expose new neurons to feature representations of past classes, ensuring a more calibrated response across all tasks.

The expansion buffer B is a key component in our proposed pipeline, designed to facilitate the integration of new classes while maintaining overall classification performance. Unlike conventional replay buffers in CIL, which mainly prevent catastrophic forgetting by re-exposing past data, our ROLANN classifier already preserves prior knowledge due to its dynamic weight computation from incrementally updated matrices (M , U , and S). Instead, the expansion buffer specifically addresses the challenge posed by ROLANN’s sigmoid-activated, independently trained neurons: without calibration, new neurons can dominate the output distribution. To counter this, the buffer maintains a randomly selected subset of previously seen examples, which serve as negative examples for training new neurons. By exposing them to the diversity of earlier data, the buffer ensures more balanced decision boundaries relative and prevents excessive activation for unknown classes.

Formally, the output probability for class i defined as:

$$y_i = \frac{1}{1 + e^{-(\mathbf{w}_i^T \mathbf{x} + b_i)}} \quad (9)$$

where w_i and b_i are the weight vector and bias, respectively, for class i , the final prediction is obtained as $\text{pred} = \text{argmax}(y_i)$. As discussed, without proper calibration, neurons introduced in later tasks may produce disproportionately high activations for instances of earlier classes, overshadowing correct predictions. The expansion buffer alleviates this by enforcing more uniform activations across negative classes, resulting in a more balanced and reliable classifier response. Figure 3 clearly illustrates the impact of the expansion buffer.

We adopt random sampling for the expansion buffer rather than employing instance selection techniques such as herding or heuristic-based criteria. While these approaches attempt to preserve the most “representative” samples of the current task, they risk biasing the buffer toward features that are optimal for the present distribution but less useful for distinguishing future classes or maintaining balanced knowledge across past tasks. Random sampling, in contrast, provides a simple yet effective mechanism to retain a diverse subset of examples, reducing the likelihood of overfitting the buffer to specific task characteristics. This strategy increases the probability that the buffer remains

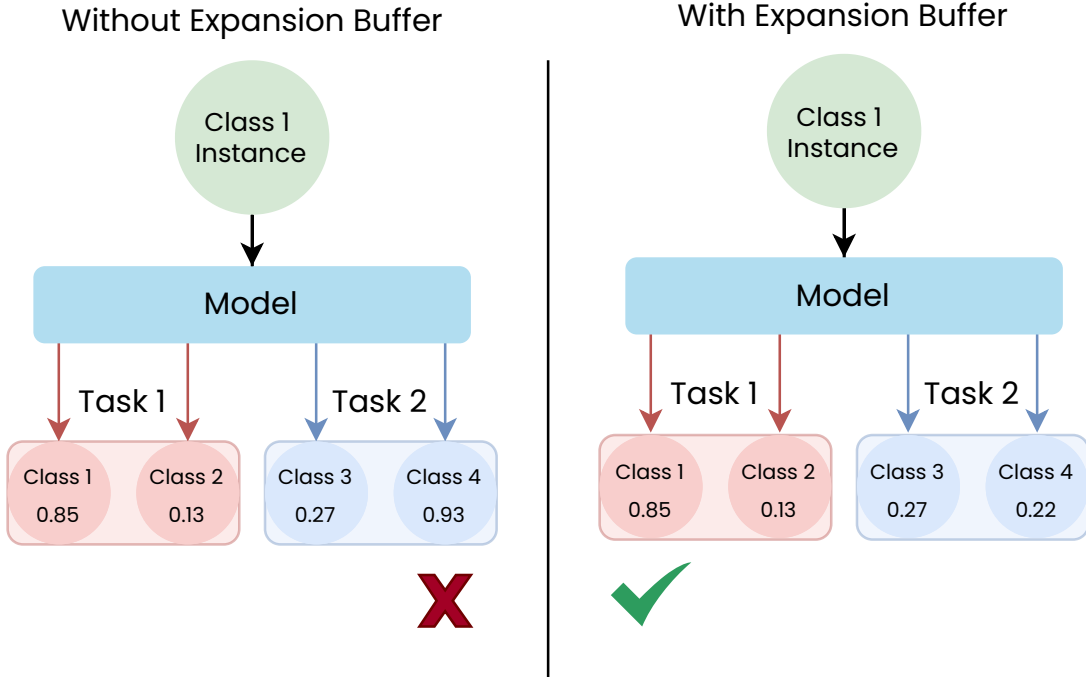


Figure 3: An instance of class 1 (corresponding to task 1) enters the network, showing how the expansion buffer mitigates the dominance of neuron activation: Without the buffer, neurons 3 and 4 from new task 2 show high activation for unknown classes, while with the buffer these neurons show lower, more calibrated responses due to exposure to small representations of previous classes, resulting in a more balanced probability distribution across negative classes.

broadly representative over time, thereby supporting more stable knowledge retention and effective calibration of new ROLANN neurons.

A primary feature of the buffer is that, since the backbone F_θ is frozen, it is unnecessary to store the original images. Instead, the extracted features $\mathcal{E} \in \mathbb{R}^d$, representing a compressed embedding of the images, are stored. This greatly reduces memory requirements, as the buffer B_k stores compact feature representations rather than high-dimensional raw data. For example, storing an image of size $3 \times 224 \times 224$ (ImageNet format) in integer format requires ~ 150 KB, while its extracted feature representation (e.g., a 512-dim float32 vector such as the output of ResNet-18 [24]) occupies only 2 KB, reducing memory usage by a factor of 75. This design improves computational and energy efficiency, making the method more scalable and suitable for resource-constrained environments.

It is critical to note that combining the buffer with data from the current task creates a strong imbalance: new classes contribute many fresh samples, while past classes are represented by only a few stored samples. This bias drives the optimization process to prioritize current classes, leading to a degradation in accuracy for earlier ones.

To counter this, we introduce a temporal oversampling strategy applied to the expansion buffer before training. For each class c previously stored in the buffer B_{k-1} , let n'_c denote its number of stored samples, and let $n_{\max} = \max(\{n_c\}_{c \in \mathcal{C}_k})$ be the size of the largest class in the current task \mathcal{D}_k . Each buffered class c is then replicated $r_c = \lfloor n_{\max}/n'_c \rfloor$ times, producing a balanced dataset in which past classes appear in proportion to the dominant current class.

This oversampling ensures that all classes –old and new– contribute more equally to training, mitigating bias toward the most recent task. As a result, past knowledge is reinforced more effectively, stabilizing the learning process and improving robustness against forgetting.

The training procedure of CIFNet is finally detailed in Algorithm 1.

Algorithm 1 Training Phase**Input:** $\mathcal{D}_k = \{(x_i, y_i)\}_{i=1}^{N_k}$, training data for task k , where $x_i \in \mathbb{R}^{h \times w \times c}$ and $y_i \in \mathcal{C}_k$ F_θ , frozen feature extractor $G_{\phi_{k-1}}$, classifier layer $B_{k-1} = \{(x_i, y_i)\}_{i=1}^{N_{k-1}}$, expansion buffer m , buffer Size per class**Output:** Updated classifier layer G_{ϕ_k} , Updated expansion buffer B_k

```

1: function INCREMENTALTRAIN( $\mathcal{D}_k, F_\theta, G_{\phi_{k-1}}, B_{k-1}, m$ )
2:   Expand  $G_{\phi_{k-1}}$  by adding  $\mathcal{C}_k$  output neurons for the new classes in  $\mathcal{C}_k$ 
3:    $\tilde{B} \leftarrow \text{OVERSAMPLEBUFFER}(\mathcal{D}_k, B_{k-1})$ 
4:    $\mathcal{E}_{D_k} \leftarrow \{(F_\theta(x_i), y_i) \mid (x_i, y_i) \in \mathcal{D}_k\}$ 
5:    $X_k, Y_k \leftarrow \mathcal{E}_{D_k} \cup \tilde{B}$ 
6:    $G_{\phi_k} \leftarrow \text{TRAINROLANN}(X_k, Y_k, G_{\phi_{k-1}})$ 
7:    $B_k \leftarrow \text{UPDATEBUFFER}(\mathcal{E}_{D_k}, B_{k-1}, m)$ 
8:   return  $G_{\phi_k}, B_k$ 
9: end function

```

5 Experimental Evaluation

In this section, a detailed experimental evaluation of the proposed method is presented, with a focus on a comparative analysis of its performance against several state-of-the-art methods in the context of CIL. The evaluation encompasses not only accuracy but also sustainability metrics, such as training duration, emissions, and energy consumption. This enables a comprehensive evaluation of the efficiency of the proposed method. In addition to performance comparisons, we conduct an ablation study to assess the impact of different design choices on the overall performance, justifying the selected approach.

5.1 Experimental Setup

Before presenting the results, we provide an overview of the benchmark datasets used, the hyperparameters and backbone employed, including buffer size configurations and other relevant experiment details.

Benchmark Datasets We validate our method on two benchmark datasets widely used in CIL: CIFAR-100 [8] and ImageNet-100 [9]. CIFAR-100 contains 50,000 training images and 10,000 test images (32×32 pixels) evenly distributed over 100 classes. ImageNet-100 is a subset of ImageNet, which includes 130,000 training and 5,000 validation images (224×224 pixels) across 100 classes. These classes are selected as the first 100 classes after a random shuffle³. Both datasets were divided into sequential groups of equal size to simulate incremental learning scenarios.

Selected Methods

5.1.1 Selected Methods

We compared CIFNet with state-of-the-art CIL methods from the survey by Zhou et al. [25], grouped as follows:

- **Data Replay methods:** Replay [26], RMM [27], GEM [28], DER [6]. Storage and periodic reuse of samples from previous tasks to reduce forgetting.
- **Data Regularization methods:** GEM [28], DER [6]. Update constraints that preserve performance on previous tasks.
- **Dynamic Networks methods:** RMM-FOSTER [29], MEMO [30], DER [31]. Network expansion or structural adaptation when new tasks arrive.
- **Knowledge Distillation methods:** iCaRL [2], PODnet [32], Coil [33], DER [6], RMM-FOSTER [29]. Transfer of knowledge from previous models through distillation losses.

³The division of these 100 classes is provided in the code presented in [25], and all experiments associated with this dataset use the same subset of classes.

- **Model Rectifying methods:** WA [34], BiC [35]. Post-hoc correction of decision boundaries or bias after learning new tasks.

We also included the *Finetune* baseline, which retrain the entire network on new classes without retaining prior knowledge, leading to catastrophic forgetting. All the selected methods are based on CNNs. We excluded Visual Transformers due to their high computational complexity.

CIFNet was implemented in PyTorch, the code is available on our GitHub⁴. All the other methods were chosen from the repository of Zhou et al.⁵. The experiments were carried out on an NVIDIA GeForce RTX 3090 Ti GPU and an Intel Core i7-12700K CPU.

Evaluation Metrics Model performance in CIL is primarily assessed using accuracy-based metrics. Following Zhou et al. [25], A_k denotes the top-1 accuracy after task k , i.e., the proportion of correctly classified instances at that stage. As CIL models progressively incorporate new tasks and classes, a decline in accuracy is commonly observed due to the challenge of retaining knowledge of previously learned classes while adapting to new ones. The final accuracy A_K reflects the performance after the last task, while the average accuracy (\bar{A}) (Eq. 10) captures learning dynamics across all stages.

$$\bar{A} = \frac{1}{K} \sum_{k=1}^K A_k \quad (10)$$

To account for efficiency, we additionally report green metrics—training time (minutes), carbon emissions (CO₂-eq), and energy consumption (kWh)—using the CodeCarbon tool⁶. We argue that an incremental model should achieve accuracy improvements with minimal environmental costs.

Buffer Configuration and Hyperparameters All buffer-based methods were constrained to a maximum of 2,000 raw images to ensure consistent evaluation and to avoid unfair advantages from larger memory allowances. For CIFNet, the buffer stores embeddings rather than raw images, resulting in a markedly lower memory footprint. As shown in B, storing 2,000 ImageNet-100 samples as raw float32 images requires 1206.40 MB, while the corresponding embeddings are only 4.10 MB.

With regard to hyperparameters, our method is hyperparameter-free, and for the other methods, these were adopted from Zhou et al. [25].

Backbone Selection and Training ResNet-32 was used for CIFAR-100 and ResNet-18 for ImageNet-100, reflecting the relative complexity of the datasets and the need to balance capacity with efficiency. In CIFNet, the backbone F_θ is pre-trained on ImageNet and kept frozen during incremental learning. The decision is based on the well-established effectiveness of pre-trained models in providing expressive and transferable representations for computer vision tasks. Since these representations already capture pertinent information about ImageNet images, utilising them is a rational strategy to enhance incremental learning stability and efficiency. By focusing learning exclusively on the classifier layer, our approach reduces computational costs and ensures greater training stability compared to methods that continually adjust backbone weights, which introduces higher computational demands and risks disrupting previously learned knowledge. While CIFNet differs from backbone-updating approaches, our goal is not to evaluate representation learning, but to demonstrate that a fixed-feature approach can achieve competitive and robust performance efficiently, as will be shown in our results.

5.2 Comparison on CIFAR-100 Dataset

Our experiments on the CIFAR-100 are designed to rigorously evaluate performance in settings that emulate real-world conditions, focusing on scenarios with 5, 10, 20, and 50 new classes to be learned per increment task.

Tables 1, 2, 3, and 4 summarize the performance of the model, including the number of parameters (#P, in millions), the accuracy, training time, and the sustainability measures, with arrows indicating the improvement direction for each metric and best results marked in bold. The results demonstrate the advantages of CIFNet, consistently balancing accuracy and computational efficiency.

⁴<https://github.com/AlejandroDopico2/CIFNet>

⁵https://github.com/zhoudw-zdw/CIL_Survey/

⁶<https://github.com/mlco2/codecarbon>

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	0.46	14.17	0.012	0.068	17.58	5.25
BiC	0.46	113.28	0.089	0.512	53.65	26.91
Coil	0.46	309.88	0.180	1.031	59.14	34.61
FOSTER	0.46	75.96	0.075	0.430	63.54	48.43
GEM	0.46	136.72	0.145	0.834	20.87	6.99
iCaRL	0.46	44.94	0.041	0.237	54.58	34.87
PODnet	0.46	62.98	0.056	0.322	49.22	28.41
Replay	0.46	16.33	0.015	0.084	54.49	34.05
RMM-FOSTER	0.46	76.36	0.075	0.429	67.03	51.10
WA	0.46	45.61	0.042	0.239	59.30	42.75
MEMO	7.14	66.80	0.064	0.367	68.49	54.34
DER	9.27	102.44	0.116	0.667	71.34	57.34
CIFNet (Ours)	0.46	15.13	0.011	0.063	73.56	59.26

Table 1: Comparison of energy consumption, environmental impact, and accuracy of various methods on the CIFAR-100 dataset with a 5 classes per task setting. The best results are marked in bold.

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	0.46	11.22	0.010	0.059	28.18	9.02
BiC	0.46	64.83	0.055	0.314	63.52	40.33
Coil	0.46	187.89	0.111	0.637	62.12	39.54
FOSTER	0.46	51.06	0.052	0.301	68.93	55.19
GEM	0.46	117.31	0.133	0.765	29.20	10.10
iCaRL	0.46	30.53	0.030	0.170	61.51	40.76
PODnet	0.46	40.69	0.037	0.215	57.52	37.86
Replay	0.46	12.09	0.011	0.065	61.04	41.53
RMM-FOSTER	0.46	50.70	0.052	0.298	70.86	57.12
WA	0.46	31.07	0.030	0.172	66.97	51.90
MEMO	7.14	36.51	0.035	0.202	71.99	58.46
DER	9.27	48.46	0.053	0.302	70.80	59.03
CIFNet (Ours)	0.46	14.35	0.010	0.060	72.93	59.57

Table 2: Comparison of energy consumption, environmental impact, and accuracy of various methods on the CIFAR-100 dataset with a 10 classes per task setting. The best results are marked in bold.

In the most realistic small-increment settings (5 and 10 classes per task), CIFNet establishes its dominance. In the 5-class setting, it achieves the highest final accuracy (A_K) of 59.26%, outperforming all other methods. The closest competitor, DER, trails at 57.34%. More critically, this superior accuracy is achieved in just 15.13 minutes, whereas DER requires 102.44 minutes, a 7x increase in training time, and an 11x increase in energy consumption for a worse result. This finding directly challenges the prevailing notion that complex, computationally-heavy methods are required to attain state-of-the-art accuracy.

When we analyse the 10-class setting, the results are similar. CIFNet again achieves a final top-tier accuracy of 59.57%. Although the simple Replay method is marginally faster (12.09 vs 14.35 minutes), it incurs a catastrophic accuracy penalty, scoring only 41.45%. This demonstrates that Replay’s speed is achieved by sacrificing a primary goal of learning. Conversely, methods that are slightly below CIFNet’s accuracy, like MEMO (58.46%) or DER (59.03%), demand 2.5 to 3.5 times the computational cost.

This pattern of superior efficiency continues in the larger-increment settings. In the 20-class scenario, RMM-FOSTER achieves a marginally higher accuracy (63.01% vs. 59.98%). However, this minor 3% accuracy gain comes at the

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	0.46	10.91	0.011	0.061	44.04	17.14
BiC	0.46	41.71	0.038	0.218	68.87	46.44
Coil	0.46	104.17	0.066	0.376	67.39	46.34
FOSTER	0.46	37.65	0.040	0.230	72.00	60.38
GEM	0.46	99.93	0.116	0.669	43.87	17.34
iCaRL	0.46	23.78	0.024	0.137	66.66	46.38
PODnet	0.46	28.91	0.028	0.158	66.66	49.90
Replay	0.46	10.72	0.010	0.060	63.81	44.04
RMM-FOSTER	0.46	37.55	0.040	0.227	73.50	63.01
WA	0.46	23.48	0.024	0.136	70.03	56.06
MEMO	7.14	24.73	0.024	0.138	72.78	62.17
DER	9.27	28.50	0.030	0.172	72.93	62.67
CIFNet (Ours)	0.46	13.90	0.010	0.058	71.17	59.98

Table 3: Comparison of energy consumption, environmental impact, and accuracy of various methods on the CIFAR-100 dataset with a 20 classes per task setting. The best results are marked in bold.

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	0.46	13.00	0.013	0.077	77.66	38.22
BiC	0.46	26.57	0.026	0.151	75.56	63.97
Coil	0.46	78.68	0.050	0.289	77.12	59.86
FOSTER	0.46	26.26	0.029	0.164	77.96	68.22
GEM	0.46	66.35	0.078	0.448	76.22	37.48
iCaRL	0.46	19.31	0.020	0.114	77.24	57.88
PODnet	0.46	21.05	0.021	0.118	76.62	62.73
Replay	0.46	12.39	0.013	0.073	76.68	49.91
RMM-FOSTER	0.46	35.65	0.039	0.225	77.96	68.40
WA	0.46	18.92	0.020	0.113	77.24	63.33
MEMO	7.14	13.30	0.013	0.076	78.36	66.88
DER	9.27	19.24	0.020	0.114	77.24	67.42
CIFNet (Ours)	0.46	13.31	0.010	0.055	67.36	61.31

Table 4: Comparison of energy consumption, environmental impact, and accuracy of various methods on the CIFAR-100 dataset with a 50 classes per task setting. The best results are marked in bold.

cost of nearly triple the training time (37.55 vs. 13.90 minutes) and over four times the energy consumption. For any practical application, such a trade-off is indefensible. It highlights that CIFNet’s core design -decoupling the stable feature extractor from the highly adaptive, non-iterative classifier- is fundamentally more efficient.

To effectively compare the trade-off between accuracy and training speed, Figure 4 visualises their relationship across different incremental learning settings. Rather than evaluating each method in isolation, this plot provides a holistic view of how models balance these two factors across all CIL configurations (5, 10, 20, and 50 classes per task) using their final accuracy values. By analysing these trends, we gain deeper insights into the efficiency of each approach, highlighting the strengths and limitations of different methods in a unified framework. The compactness of CIFNet’s curve across the various increment settings shows that, regardless of the number of increments, both training speed and accuracy remain consistently stable. This consistency reinforces the robustness of our approach, as it is less dependent on specific configurations or the number of training phases. In contrast, alternative methods generally exhibit a pronounced drop in accuracy when using smaller incremental steps and a substantial increase in training time as the number of increments grows, making them less suitable for real-world incremental learning scenarios.

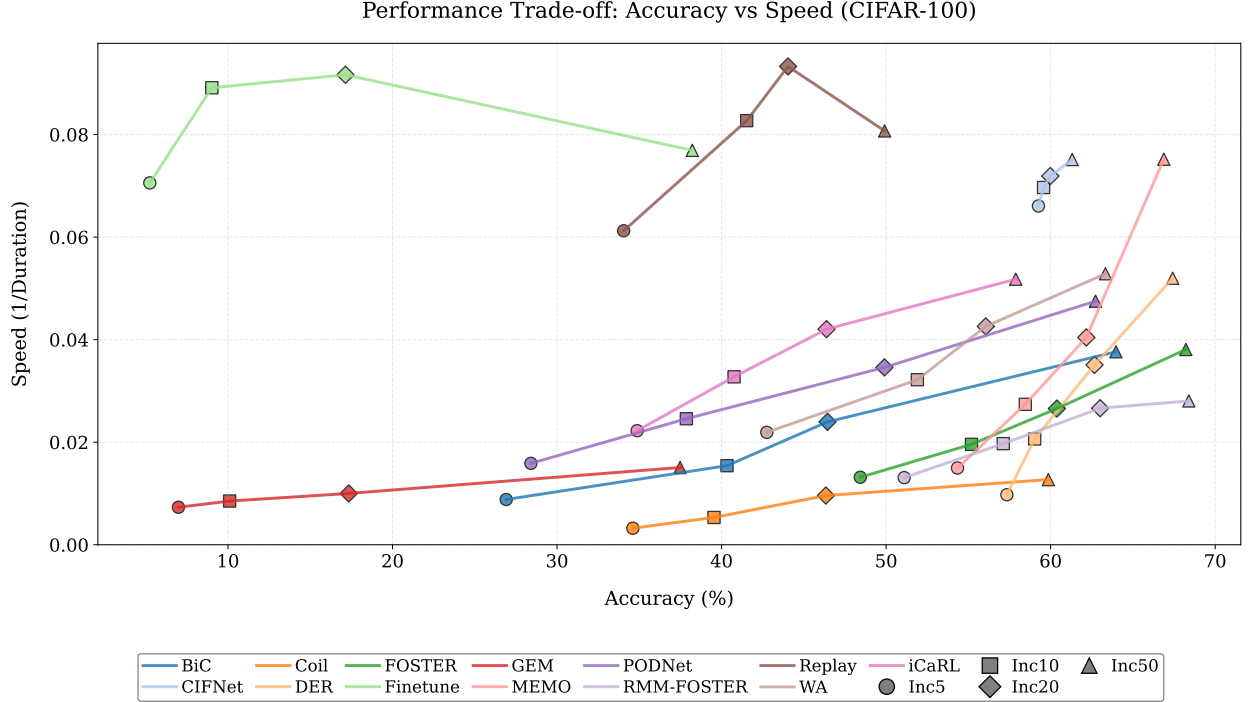


Figure 4: Comparison of final accuracy and (inverted) training duration for different models under various incremental learning configurations on the CIFAR-100 dataset. The optimal region is the upper-right corner, representing both high accuracy and fast training.

A critical test for any CIL method is its robustness as the number of tasks increases (i.e., as the increment size decreases). Figure 5 illustrates this clearly. The performance curves of methods like BiC, iCaRL, and PODnet show a pronounced degradation as the scenario moves from 20 to 10 to 5 classes per task. In contrast, CIFNet’s accuracy curve remains much more stable across all settings.

This consistency underscores the robustness of our approach because we do not fine-tune the backbone; our model is less susceptible to the compounding errors and feature drift that other methods have over longer learning sequences. One could hypothesize that the same level of stability could be achieved through alternative methods of backbone freezing. However, these methods are not designed for fixed features and rely heavily on backbone adaptation to maintain feature quality. Freezing the backbone in such methods generally results in a decline in sharpness accuracy. This is because their classifiers and regularization schemes are not calibrated for non-adaptive embeddings. In contrast, our design (ROLANN classifier, expansion buffer, and oversampling) is explicitly optimized for this fixed-feature regime, which explains its robustness as the number of tasks increases.

5.3 Comparison on ImageNet-100 Dataset

Following the experimental guidelines established for the previous dataset, experiments were conducted on ImageNet-100. However, due to the impractical computational and temporal requirements for large-scale datasets like ImageNet, we omitted the 20 and 50-class increments per task. This highlights a fundamental scalability limitation of existing CIL paradigms, a challenge that CIFNet inherently addresses through its design.

Tables 5 and 6 present the number of parameters (#P, in millions), accuracy, training time, and sustainability metrics, for different task configurations with 5 and 10 classes per task, respectively. In the 5-class setting, our method achieves a final accuracy (A_K) of 78.10%, representing a remarkable 15-point lead over the second-best method, DER (63.66%). Furthermore, as shown in the tables, CIFNet is the only evaluated method to surpass the 70% accuracy mark, setting it apart from the others. Figure 7 clearly illustrates the performance gap, showing that CIFNet’s accuracy far exceeds that of other methods. It maintains its lead consistently across all incremental tasks, demonstrating superior knowledge retention.

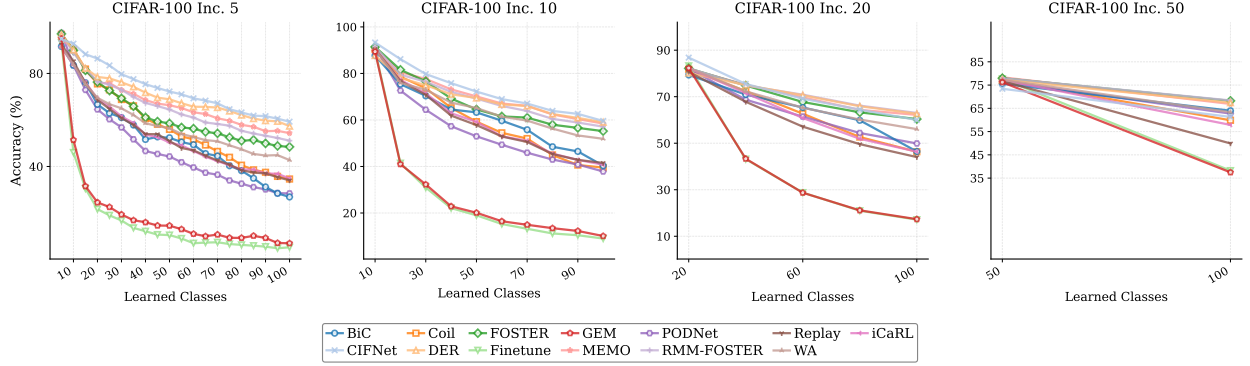


Figure 5: Incremental learning performance (A_k) of various methods on CIFAR-100 under different class-increment settings (5, 10, 20, and 50 classes per task).

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	11.17	188.93	0.198	1.136	17.60	4.64
BiC	11.17	1109.27	1.013	5.821	51.96	23.44
Coil	11.17	1812.03	1.323	7.599	57.29	33.60
FOSTER	11.17	953.89	1.119	6.428	66.01	53.50
GEM*	-	-	-	-	-	-
iCaRL	11.17	536.53	0.599	3.440	54.97	33.06
PODnet	11.17	721.67	0.720	4.138	55.83	37.58
Replay	11.17	226.92	0.227	1.303	56.28	35.38
RMM-FOSTER	11.17	935.67	1.089	6.255	71.73	59.46
WA	11.17	544.00	0.602	3.461	63.21	46.72
MEMO	170.60	758.81	0.847	4.864	68.42	55.52
DER	223.40	1354.40	1.736	9.972	73.88	63.66
CIFNet (Ours)	11.17	71.50	0.047	0.271	85.55	78.10

* Results could not be obtained due to execution issues

Table 5: Comparison of energy consumption, environmental impact, and accuracy of various methods on the ImageNet-100 dataset with a 5 classes per task setting. The best results are marked in bold.

Beyond accuracy, the most striking feature of CIFNet is its remarkable efficiency. In contrast to the extensive time requirements of other methods, which can be up to 25 times longer than CIFNet’s, our approach is notably more efficient, requiring only ~ 70 minutes for training. For instance, Coil and DER demand over 1,800 and 1,300 minutes, respectively. This substantial improvement in efficiency is especially noteworthy in real-world applications, where computational efficiency is a critical factor.

Figure 6 visualizes the trade-off between accuracy and training speed, showcasing CIFNet’s superior balance between these two metrics. Moreover, our approach also stands out for its minimal environmental impact. Compared to methods such as DER (which consumes 9.972 kWh and emits 1.736 kg CO₂-eq), our method requires significantly less energy (only 0.271 kWh) and generates proportionally lower emissions (as low as 0.047 kg CO₂-eq). This represents approximately a 37-fold reduction in energy consumption compared to DER, as detailed in Tables 5 and 6, highlighting the sustainability benefits of our method, making it a viable alternative for large-scale deployments where energy efficiency is a key concern.

To sum up, the robustness of our method is reflected in its consistent performance and training efficiency across a wide range of incremental learning configurations. Regardless of the number of steps or task splits, CIFNet maintains both high accuracy and low computational cost. This stability under varying conditions makes it especially suitable for practical, real-world deployments, where scalability and reliability are essential.

Model	#P	Training Time (min.) ↓	Emissions (kg CO ₂ -eq) ↓	Energy Consumed (kWh) ↓	$\bar{A} \uparrow$	$A_K \uparrow$
Finetune	11.17	184.20	0.201	1.152	28.40	9.34
BiC	11.17	829.46	0.804	4.619	65.93	37.00
Coil	11.17	1249.63	0.977	5.614	64.66	41.60
FOSTER	11.17	767.35	0.918	5.275	70.99	61.06
GEM*	-	-	-	-	-	-
iCaRL	11.17	449.84	0.511	2.936	61.59	40.18
PODnet	11.17	539.05	0.566	3.253	65.80	46.78
Replay	11.17	200.77	0.207	1.191	62.75	42.96
RMM-FOSTER	11.17	757.55	0.903	5.186	76.49	66.10
WA	11.17	453.38	0.512	2.944	70.92	55.44
MEMO	170.60	505.52	0.562	3.231	72.66	60.94
DER	223.40	742.22	0.927	5.328	76.04	64.40
CIFNet (Ours)	11.17	70.02	0.046	0.264	85.20	77.92

* Results could not be obtained due to execution issues

Table 6: Comparison of energy consumption, environmental impact, and accuracy of various methods on the ImageNet-100 dataset with a 10 classes per task setting. The best results are marked in bold.

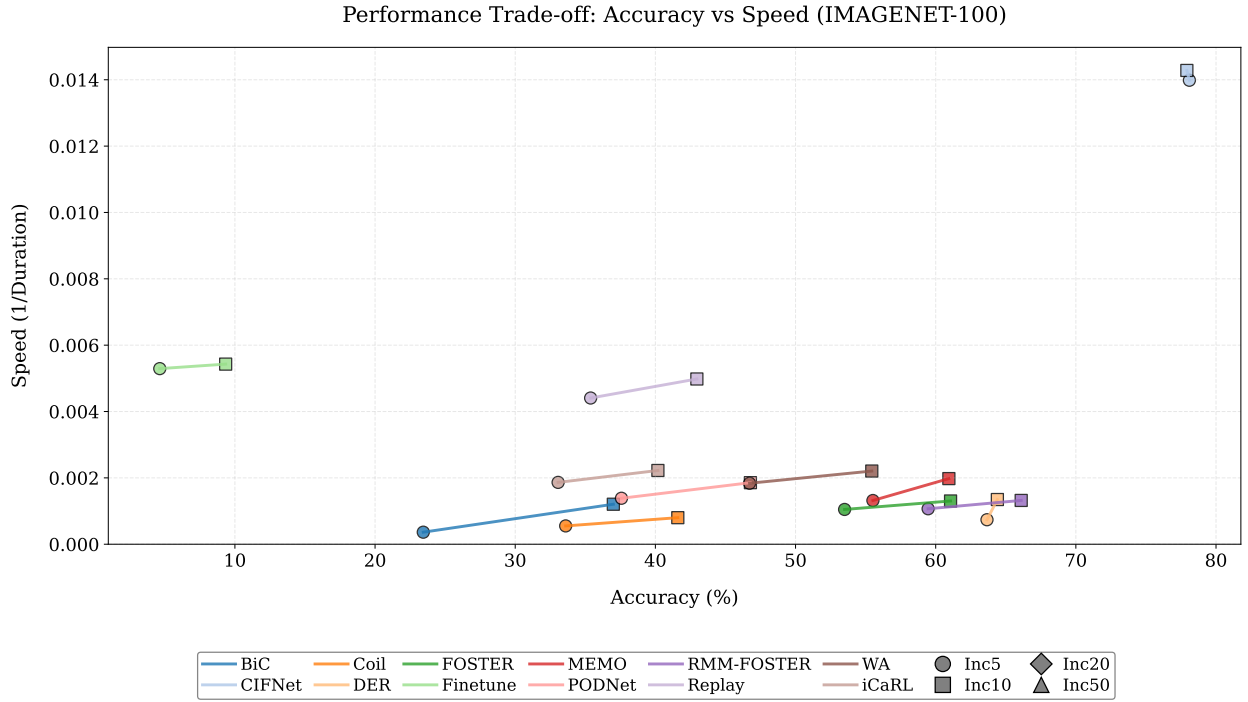


Figure 6: Comparison of final accuracy and (inverted) training duration for different models under various incremental learning configurations on the ImageNet-100 dataset. The optimal region is the upper-right corner, representing both high accuracy and fast training.

5.4 Ablation Study

To validate the design choices of CIFNet, we conducted an ablation study on the CIFAR-100 dataset with class increments of 5 and 10, corresponding to 20 and 10 tasks, respectively. We systematically evaluated the impact of two key components: the expansion buffer and the oversampling strategy for stored samples.

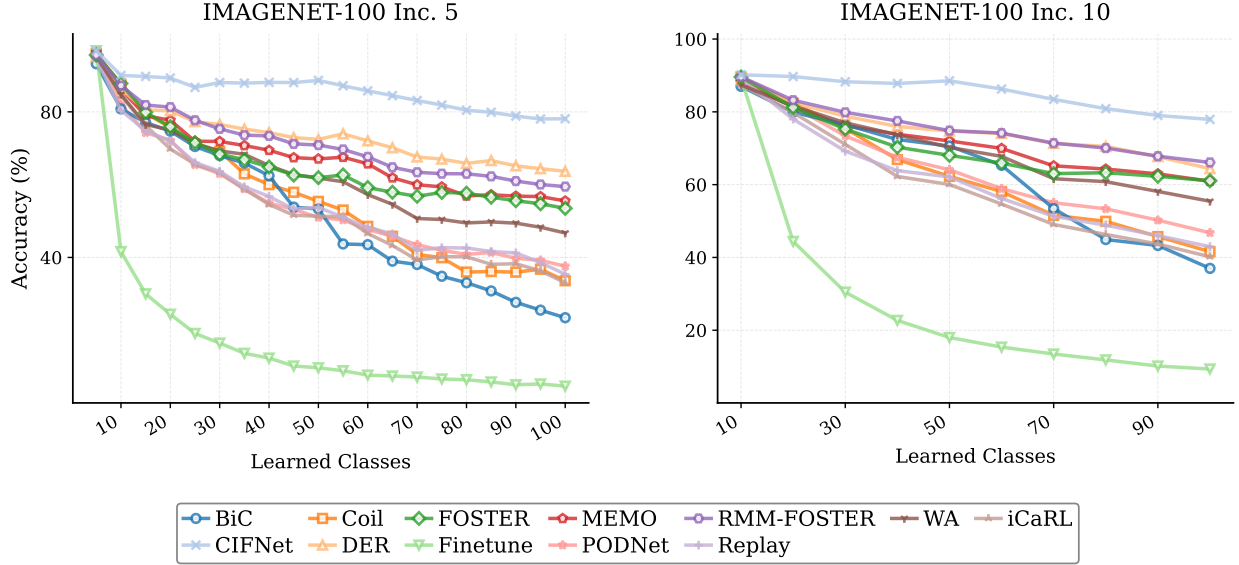


Figure 7: Incremental learning performance (A_k) of various methods on the ImageNet-100 dataset under different class-increment settings.

Table 7 reports the results, while Figure 8 illustrates the incremental learning performance across tasks. The results reveal that removing either component leads to a notable performance drop. Specifically, omitting the expansion buffer yields the largest decrease, reducing \bar{A} by more than 30 points in the 5-class setting and nearly 25 points in the 10-class setting. This highlights the buffer’s crucial role in mitigating catastrophic forgetting and maintaining long-term knowledge retention. Excluding the oversampling strategy also negatively impacts performance, as it fails to adequately balance class distributions during training, leading to biased updates and poorer final accuracy.

In contrast, the full method, which integrates both the expansion buffer and oversampling, consistently achieves the highest overall accuracy and final performance (A_K). These results confirm that both components are indispensable for ensuring stable incremental learning and a uniform representation of past and new classes.

	Increment 5		Increment 10	
	$\bar{A} \uparrow$	$A_K \uparrow$	$\bar{A} \uparrow$	$A_K \uparrow$
Without Expansion Buffer	42.18	19.42	50.36	24.74
Without Oversampling	61.74	36.62	61.63	37.55
Full Method	73.56	59.26	72.93	59.57

Table 7: Impact of CIFNet components on accuracy (\bar{A} and A_K) using the CIFAR-100 dataset with 5- and 10-class increments.

6 Conclusions and Future Work

We presented CIFNet, a novel class incremental learning (CIL) framework that shifts the focus from continual backbone adaptation to efficient classifier learning. By employing a frozen pre-trained feature extractor and the ROLANN layer, CIFNet achieves non-iterative, single-step weight updates per task, eliminating costly fine-tuning while maintaining competitive accuracy. The use of compressed embedding buffer further improves scalability by reducing memory demands.

Experiments on standard benchmarks show that CIFNet matches or surpasses state-of-the-art accuracy while offering orders-of-magnitude reductions in training time, energy consumption, and estimated carbon emissions. These properties make it well suited for deployment in resource-constrained environments such as edge devices and mobile platforms.

A current limitation is the reliance on fixed pre-trained features, which may underperform in domains with large distribution shifts from the pre-training data (e.g., medical or satellite imagery). Future work will explore lightweight

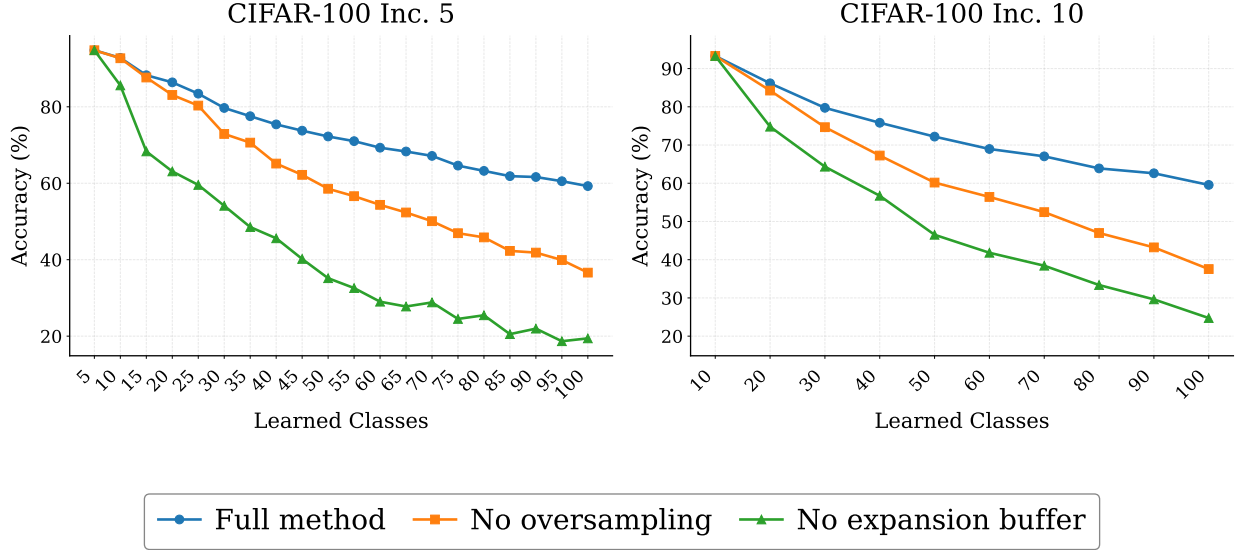


Figure 8: Incremental learning performance (A_k) of CIFNet on CIFAR-100 under different class-increment settings, comparing variants with key components removed.

adaptations on the feature extractor, such as adapter modules, prompt-based tuning, or selective fine-tuning of higher layers, to extend CIFNet’s applicability while preserving its efficiency-oriented design.

Acknowledgements

This work has been funded from Horizon Europe (GA 101070381), MCIN/ AEI/10.13039/501100011033/ERDF, UE (project PID2023-147404OB-I00), and the Ministry for Digital Transformation and Civil Service and Next-GenerationEU/ PRTR (TSI-100925-2023-1). CITIC, as a center accredited for excellence within the Galician University System and a member of the CIGUS Network, receives subsidies from the “Consellería de Cultura, Educación e Universidade” of the Xunta de Galicia, supported in an 80% through ERDF Operational Programme Galicia 2021 - 2027 (Grant ED431G 2023/01).

References

- [1] Jeffrey C Schlimmer and Douglas Fisher. A case study of incremental concept induction. In *Proceedings of the Fifth AAAI National Conference on Artificial Intelligence*, pages 496–501, 1986.
- [2] Sylvestre Alvisé Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning. *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, 2017-January:5533–5542, 11 2017.
- [3] Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. *Psychology of Learning and Motivation - Advances in Research and Theory*, 24:109–165, 1 1989.
- [4] Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. In *International conference on machine learning*, pages 3925–3934. PMLR, 2019.
- [5] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- [6] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.

- [7] Da-Wei Zhou, Zi-Wen Cai, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-incremental learning with pre-trained models: Generalizability and adaptivity are all you need. *International Journal of Computer Vision*, 133(3):1012–1032, 2025.
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [9] Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- [10] German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review, 5 2019.
- [11] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16*, pages 699–715. Springer, 2020.
- [12] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *arXiv preprint arXiv:1708.01547*, 2017.
- [13] Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost Van De Weijer. Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533, 2022.
- [14] Daniel N Barry and Bradley C Love. A neural network account of memory replay and knowledge consolidation. *Cerebral Cortex*, 33(1):83–95, 2023.
- [15] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *Advances in neural information processing systems*, 32, 2019.
- [16] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *Advances in neural information processing systems*, 30, 2017.
- [17] Gido M Van de Ven and Andreas S Tolias. Generative replay with feedback connections as a general strategy for continual learning. *arXiv preprint arXiv:1809.10635*, 2018.
- [18] Gido M Van De Ven, Zhe Li, and Andreas S Tolias. Class-incremental learning with generative classifiers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3611–3620, 2021.
- [19] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [20] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7765–7773, 2018.
- [21] Steven C. Y. Hung, Jia-Hong Lee, Timmy S. T. Wan, Chein-Hung Chen, Yi-Ming Chan, and Chu-Song Chen. Increasingly packing multiple facial-informatics modules in a unified deep-learning model via lifelong learning. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, page 339–343, 2019.
- [22] Dahuin Jung, Dongyoon Han, Jihwan Bang, and Hwanjun Song. Generating instance-level prompts for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11847–11857, 2023.
- [23] Oscar Fontenla-Romero, Bertha Guijarro-Berdiñas, and Beatriz Pérez-Sánchez. Regularized one-layer neural networks for distributed and incremental environments. In *Advances in Computational Intelligence: 16th International Work-Conference on Artificial Neural Networks, IWANN 2021, Virtual Event, June 16–18, 2021, Proceedings, Part II*, page 343–355, Berlin, Heidelberg, 2021. Springer-Verlag.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [25] Da-Wei Zhou, Qi-Wei Wang, Zhi-Hong Qi, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Class-incremental learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(12):9851–9873, 2024.
- [26] Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and forgetting functions. *Psychological review*, 97(2):285, 1990.
- [27] Yaoyao Liu, Bernt Schiele, and Qianru Sun. Rmm: Reinforced memory management for class-incremental learning. *Advances in Neural Information Processing Systems*, 34:3478–3490, 2021.
- [28] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- [29] Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *European conference on computer vision*, pages 398–414. Springer, 2022.

- [30] Da-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards memory-efficient class-incremental learning. *arXiv preprint arXiv:2205.13218*, 2022.
- [31] Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3014–3023, 2021.
- [32] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer vision—ECCV 2020: 16th European conference, Glasgow, UK, August 23–28, 2020, proceedings, part XX 16*, pages 86–102. Springer, 2020.
- [33] Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Co-transport for class-incremental learning. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1645–1654, 2021.
- [34] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 374–382, 2019.
- [35] Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Maintaining discrimination and fairness in class incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13208–13217, 2020.

A ROLANN training algorithm

The training procedure described here extends the ROLANN framework to support incremental and distributed learning. The pseudocode 2 summarizes the training mechanism.

Algorithm 2 ROLANN for incremental/distributed learning.

Input: $X_p \in \mathbb{R}^{m \times n}$ (new training data with m variables and n samples), $Y_p \in \mathbb{R}^{n \times C}$ (desired outputs for C output neurons), G_ϕ (ROLANN layer containing:

- Knowledge Components: $\{M_k, U_k, S_k\}_{c=1}^C$ for each output neuron c ,
- Trainable Parameters: $\{w_k\}_{c=1}^C$ for each output neuron c ,
- Global hyperparameters f, λ shared across all outputs).

Output: Updated ROLANN layer G_ϕ with:

- Updated weights $\{w_k \in \mathbb{R}^{m \times 1}\}_{c=1}^C$ for each output neuron,
- Updated matrices $\{M_k, U_k, S_k\}_{c=1}^C$ for each output neuron.

```

1: function TRAINROLANN( $X_p, Y_p, G_\phi$ )
2:    $f, \lambda \leftarrow G_\phi$ 
3:   for  $c = 1$  to  $C$  do
4:      $M_{k^c}, U_{k^c}, S_{k^c} \leftarrow G_{\phi^c}$ 
5:      $\mathbf{X}_p \leftarrow [\mathbf{1}; X_p]$ 
6:      $\bar{Y}_p \leftarrow f^{-1}(Y_p)$ 
7:      $f' \leftarrow f'(Y_p)$ 
8:      $F_p \leftarrow \text{diag}(f')$ 
9:     if  $M_{k^c}, U_{k^c}, S_{k^c}$  are empty then
10:       $[U_{k^c}, S_{k^c}, \sim] \leftarrow \text{SVD}(X_p F_p)$ 
11:       $M_{k^c} \leftarrow X_p (f' \odot f' \odot Y_p)$ 
12:     else
13:       $M_{k^c} \leftarrow M_{k^c} + X_p (f' \odot f' \odot Y_p)$ 
14:       $[U_{p^c}, S_{p^c}, \sim] \leftarrow \text{SVD}(X_p F_p)$ 
15:       $[U_{k^c}, S_{k^c}, \sim] \leftarrow \text{SVD}\left(\begin{bmatrix} U_{k^c} S_{k^c} & U_{p^c} S_{p^c} \end{bmatrix}\right)$ 
16:     end if
17:      $w_c \leftarrow U_{k^c} \cdot (S_{k^c} \cdot S_{k^c} + \lambda I)^{-1} (U_{k^c}^T M_{k^c})$ 
18:     Update  $G_{\phi^c}$  with  $w_c, M_{k^c}, U_{k^c}, S_{k^c}$ 
19:   end for
20:   return  $G_\phi$ 
21: end function
```

▷ Get global parameters
 ▷ Iterate over all output neurons
 ▷ Neuron c parameters
 ▷ Add bias
 ▷ Inverse activation
 ▷ Derivative of function
 ▷ Diagonal matrix
 ▷ Economy-size
 ▷ Combine with previous knowledge

B Supplementary Material for Buffer Configuration

This appendix provides detailed information on the memory usage of the different replay buffer configurations discussed in the main text. Table 8 illustrates the substantial memory savings achieved by using low-dimensional embeddings instead of raw images.

Buffer Type	Data Configuration	Number of Items	Item Dimensions	Total Memory (MB)
Embeddings	Embeddings (float32)	2,000	512	4.10
Raw Images	CIFAR-100 (int8)	2,000	$32 \times 32 \times 3$	6.14
	CIFAR-100 (float32)	2,000	$32 \times 32 \times 3$	24.56
	ImageNet-100 (int8)	2,000	$224 \times 224 \times 3$	301.06
	ImageNet-100 (float32)	2,000	$224 \times 224 \times 3$	1,206.40

Table 8: Memory Usage Comparison: Replay Buffer with Embeddings vs. Raw Images for Various Configurations

The memory usage calculations are based on standard data type sizes (4 bytes for float32, 1 byte for int8) and provide a conservative estimate for training-time memory.