# Forget Less, Retain More: A Lightweight Regularizer for Rehearsal-Based Continual Learning

#### **Anonymous Author(s)**

Affiliation Address email

### **Abstract**

Deep neural networks suffer from catastrophic forgetting, where performance on previous tasks degrades after training on a new task. We present a novel approach to address this challenge, focusing on the intersection of memory-based methods and regularization approaches. We formulate a regularization strategy, termed Information Maximization (IM) regularizer, for memory-based continual learning methods, which is based exclusively on the expected label distribution, thus making it class-agnostic. As a consequence, IM regularizer can be directly integrated into rehearsal-based continual learning methods, reducing forgetting and favoring faster convergence. Our empirical validation shows that, across datasets and regardless of the number of tasks, our proposed regularization strategy consistently improves baseline performance at the expense of a minimal computational overhead. Finally, we demonstrate the data-agnostic nature of our regularizer by applying it to video data, which presents additional challenges due to its temporal structure and higher memory requirements. Despite the significant domain gap, our experiments show that IM regularizer also improves the performance of video continual learning methods.

## 1 Introduction

2

3

5

6

8

9

10

11

12

13

14

15

16

Continual learning (CL) aims to develop models that can learn from evolving data distributions 18 with minimal forgetting [26]. Due to the high computational and financial costs of training deep 19 neural network models and growing concerns over privacy regulations, the applicability of CL in 20 various real-world scenarios has become increasingly critical. For instance, video-sharing platforms 21 such as YouTube and TikTok receive millions of newly uploaded videos daily, each introducing new trends, visual concepts, and styles. In these dynamic environments, traditional training algorithms for deep learning models struggle to keep pace due to the necessity of frequent retraining, which is 24 resource-intensive and impractical at scale. CL can significantly enhance the effectiveness of models 25 designed for such dynamic data streams by continually adapting previously trained models, rather 26 than retraining them from scratch as new data arrives. 27

In recent years, memory-based methods [26, 5] have emerged as the front-runners in CL, demonstrating better performance at mitigating forgetting compared to their regularization-based counterparts. This superior performance of rehearsal methods is attributed to the use of a memory buffer, a dedicated storage that retains a subset of training data from previously learned tasks. By having access to a subset of past samples, the model can estimate class prototypes effectively and alleviate forgetting despite distribution shifts [26, 25]. This ability to retain past information enables rehearsal-based approaches to maintain stability in long-term learning while adapting to new tasks.

The performance improvement of rehearsal methods over regularization methods comes at the cost of increased memory requirements and greater computational overhead. Experience replay methods

retrain for several epochs on both the current task data and memory buffer samples, thus effectively approximating a joint distribution whenever new data becomes available. This continuous reprocessing of stored data samples not only increases computational demands but also increases training time, making it less scalable for large datasets. This computational penalty is further emphasized outside the image domain; for example, in video data, a single minute-long video recorded at 30 frames per second occupies as much memory as 1800 individual images. However, despite consuming significant storage, such a video might represent only a single instance of a class within the memory buffer, limiting the diversity of stored information and further challenging the efficiency of experience replay in temporal data.

In this paper, we investigate the synergy between memory-based methods and regularization tech-46 niques for CL in a class incremental setup, aiming to leverage the strengths of both approaches to 47 mitigate forgetting while maintaining computational efficiency. Furthermore, we propose a class-48 agnostic regularization strategy for CL, which targets the distribution of the network predictions [18]. 49 Such regularization enables us to learn improved feature representations across several distribution shifts, thereby enhancing the model's generalization to previously seen tasks, while simultaneously 51 minimizing the memory footprint and computational overhead. Despite its simplicity, extensive 52 empirical evaluation shows that the proposed Information Maximization (IM) regularizer emerges 53 as a consistently effective regularization technique, outperforming current regularization strategies 54 tailored for the CL setting in both accuracy and retention of past knowledge. In fact, our proposed 55 approach is not specific to the image CL domain; we further validate its effectiveness by applying the 56 IM regularizer to a video continual learning setup, where it demonstrates similarly improved results 57 in handling the challenges posed by temporal dependencies and increased data complexity.

Contributions: Our work makes two key contributions: (i) We conduct an experimental evaluation of several regularizers, including Elastic Weight Consolidation (EWC), Synaptic Intelligence (SI), Information Maximization (IM), and Entropy Minimization (EM), applied to image continual learning. This evaluation highlights the advantages of the proposed IM regularizer, demonstrating its superiority in terms of both performance and overall reduction in catastrophic forgetting. (ii) We extend our analysis beyond image-based settings by demonstrating the applicability of IM within the context of video continual learning. Given the additional complexity of temporal dependencies and larger data volumes in videos, our results show that IM maintains its effectiveness, achieving substantial gains over traditional memory-based baselines while preserving computational efficiency.

## 2 Related Work

59

60

62

63

64

65

66

**Image Continual Learning.** In the field of image-based continual learning, numerous innovative 69 approaches have been proposed to address catastrophic forgetting. Memory-based methods, such 70 as iCaRL [25], utilize incremental classifiers and representation learning to balance new and old 71 knowledge, while GEM [20] and its more efficient variant A-GEM [8] optimize gradient-based episodic memory to mitigate forgetting. Other approaches, including DER [5], enhance rehearsal 73 74 by incorporating logit distillation, while CoPE [11] leverages class prototypes to structure the latent space, and ER-ACE [7] modifies cross-entropy loss to address task imbalance. Recent work includes 75 Refresh Learning [32], which unifies rehearsal with selective unlearning to refresh model knowledge, 76 and STAR [12], a plug-and-play regularizer that leverages stability-inducing weight perturbations 77 during rehearsal to mitigate forgetting. Regularization-based methods aim to preserve past knowledge 78 by constraining weight updates, typically by identifying the importance of parameters, like Elastic 79 Weight Consolidation (EWC) [15] and Synaptic Intelligence (SI) [36]. Architectural innovations 80 also play a crucial role in continual learning, with L2P [34] demonstrating the effectiveness of 81 learnable prompts in guiding pre-trained models without relying on a rehearsal buffer. More recently, 82 DualPrompt [33] introduced a two-level prompting mechanism for transformer-based architectures. 83 These diverse approaches underscore the rapid advancements in continual learning, paving the way 84 for more scalable and adaptable models in real-world applications. 85

Video Continual Learning. To mitigate catastrophic forgetting in video data, various strategies have been developed, which can be broadly categorized into regularization and memory-based techniques. While regularization methods apply constraints to preserve previous knowledge, memory-based approaches leverage data or representations from past tasks. When analyzing video continual learning, the importance of memory becomes even more pronounced due to the temporal complexity

and higher dimensionality of video data. SMILE [2] underscores this by proposing an efficient 91 replay mechanism that stores a single frame per video, emphasizing video diversity over temporal 92 information. This approach addresses memory constraints effectively, showcasing the critical role of 93 memory in video continual learning. vCLIMB [29] and PIVOT [28] introduce novel benchmarks and 94 methods focusing on class incremental learning and the use of prompting mechanisms, respectively, 95 pushing the boundaries of current methodologies. Utilizing Winning Subnetworks for efficient 96 learning [14], and creating multi-modal datasets for egocentric activity recognition [35] illustrate the expanding scope of continual learning in video domains. Additionally, Continual Predictive Learning [10] and approaches to Video Object Segmentation as a continual learning task [21] represent 99 significant advancements in handling non-stationary environments and long video sequences. Finally, 100 efforts to learn new class representations while preserving old ones through time-channel importance 101 maps [23]. 102

**Test-Time Adaptation.** Test-Time Adaptation (TTA) aims to alleviate performance drop of pre-103 trained models at test time when exposed to domain shifts [27, 1]. Earlier works augmented the 104 training objective with a self-supervised loss function that is later leveraged at test time to combat 105 domain shifts [27, 19]. More recent TTA methods optimize an unsupervised loss function at test-time 106 on the received unlabeled data to improve performance under domain shifts [22]. This includes simple 107 adjustments to the statistics of normalization layers [17], entropy minimization [30], information 108 maximization [18], among others [4, 31]. However, most TTA methods are proposed to combat 109 covariate domain shifts at test time. In this work, we get inspiration from the source hypothesis 110 adaptation method [18] to propose an effective regularizer for continual learning. We also analyze 111 the effectiveness of other adaptation methods such as entropy minimization in mitigating catastrophic 112 forgetting in continual learning. 113

This work aims to enhance continual learning performance by introducing a cost-effective regularizer that improves results even in memory-constrained scenarios. Such scenarios are particularly important when dealing with memory-intensive data, such as videos, or when sample storage is restricted due to privacy concerns. We investigate a class-independent regularizer designed to facilitate the learning of generalizable features.

# 3 Methodology

119

126

127

128

129

130

131

132

133

134

In this section, we formalize the problem of continual learning, with a particular focus on classincremental learning in visual recognition tasks. We define the underlying framework and introduce
the necessary notation to describe the incremental learning process. Additionally, we present the
formulation of the proposed regularizer, Information Maximization (IM), along with the selected
baseline regularizers: Elastic Weight Consolidation (EWC), Synaptic Intelligence (SI), and Entropy
Minimization (EM).

We focus on the offline continual learning problem for visual recognition tasks, where a classifier  $f_{\theta}: \mathcal{X} \to \mathcal{P}(\mathcal{Y})$  (a DNN parameterized by  $\theta$ ) maps an input  $x \in \mathcal{X}$  into the probability simplex  $\mathcal{P}(\mathcal{Y})$ , with  $\mathcal{Y} = \{1, 2, \ldots, K\}$ . In continual learning,  $f_{\theta}$  is presented with a sequence of T tasks  $\{(X_1, Y_1), (X_2, Y_2), \ldots, (X_T, Y_T)\}$  where  $X_i \subset \mathcal{X}$  and  $Y_i \subset \mathcal{Y} \ \forall i$  [29]. Furthermore, we consider the class-incremental problem setup, where the labels presented in each individual task are mutually exclusive  $(Y_i \cap Y_j = \phi \ \forall i \neq j)$ . The main objective of the learner is to maximize its performance (e.g. Accuracy) on all observed tasks. This objective is often hindered by the catastrophic forgetting problem: while learning task i,  $f_{\theta}$  tends to forget previously learned tasks i, significantly dropping its performance for any i and i and i are the following problem: while learning task i, i and i are the following problem: while learning task i, i and i are the following problem: while learning task i, i are the following problem: while learning task i, i and i are the following problem: while learning task i, i and i are the following problem: while learning task i and i are the following problem: i

For our baseline, we consider rehearsal-based continual learning methods where the learner is allowed to store up to N training examples from previous tasks into a replay memory buffer M [9]. Let  $M_t$  denote the replay buffer at task t containing examples from the tasks i < t. Rehearsal-based methods update the parameter set  $\theta$  at task t in the following form:

$$\theta_t^* = \underset{\theta}{\operatorname{arg\,min}} \ \mathbb{E}_{(x,y)\sim(X_t,Y_t)} \mathcal{L}(f_{\theta}(x),y) + \mathbb{E}_{(u,v)\sim M_t} \mathcal{L}(f_{\theta}(u),v). \tag{1}$$

<sup>&</sup>lt;sup>1</sup>e.g. the network's output after Softmax.

That is, for each batch sampled from the newly available data on the  $t^{th}$  task, the learner samples another batch from memory  $\mathcal{M}_t$  and updates the model on the combined loss.

#### 3.1 Regularizing Replay Methods with Information Maximization

Inspired by the work of Liang *et al.* [18] in the domain of test-time adaptation, we hypothesize that in a continual learning setup,  $f_{\theta}$  should output confident predictions that distinctly separate all previously seen classes. This means that for any given input, the model should assign a high probability to a single class. Since the memory buffer  $\mathcal{M}_t$  contains a subset of past examples and new tasks introduce distribution shifts, the model must remain adaptable while preserving knowledge from earlier tasks. To achieve this, we propose a regularizer that encourages the model to make confident predictions across all encountered classes without biasing toward recent task data. By maximizing information in the logits, our approach helps reinforce discriminative representations for all learned classes, improving robustness against distribution shifts. Our proposed regularizer ( $\mathcal{R}_{\text{IM}}$ ) takes the following form:

$$\mathcal{R}_{\text{IM}}(\theta, X_t) = \mathcal{L}_{\text{ent}}(\theta, X_t) + \mathcal{L}_{\text{div}}(\theta, X_t)$$
with 
$$\mathcal{L}_{\text{ent}}(\theta, X_t) = -\mathbb{E}_{x \sim X_t} \sum_{k=1}^K f_{\theta}^k(x) \log f_{\theta}^k(x)$$

$$\mathcal{L}_{\text{div}} = \sum_{k=1}^K \hat{f}_{\theta}^k(x) \log \hat{f}_{\theta}^k(x),$$
(2)

where  $\hat{f}_{\theta}(x) = \mathbb{E}_{x \sim X_t}[f_{\theta}(x)]$  and  $f_{\theta}^k(x)$  is the  $k^{th}$  element in the vector  $f_{\theta}(x)$ . Note that optimizing  $\mathcal{L}_{\text{ent}}$  increases the model's confidence on the prediction, while  $\mathcal{L}_{\text{div}}$  promotes diverse label predictions on  $f_{\theta}$ . Our regularized rehearsal-based method follows the formulation:

$$\min_{\theta} \mathbb{E}_{(x,y)\sim(X_t,Y_t)} \mathcal{L}(f_{\theta}(x),y) + \mathbb{E}_{(u,v)\sim M_t} \mathcal{L}(f_{\theta}(u),v) + \mathcal{R}_{\text{IM}}(\theta,X_t). \tag{3}$$

Our proposed regularizer has the following advantages: (i) It is orthogonal to the most critical design 156 choices of continual learning algorithms, as it can operate regardless of the choice of  $f_{\theta}$ , the replay-157 based method, the size of the memory buffer, and the number of tasks. (ii) Efficient computation 158 of  $\mathcal{R}_{IM}$ : where both  $\mathcal{L}_{\text{ent}}$  and  $\mathcal{L}_{\text{div}}$  depend exclusively on the output predictions of the model and can be computed in  $\mathcal{O}(n)$ . This aspect is essential when dealing with memory-intensive setups. 160 For example, on video data, our regularizer estimates  $\mathcal{L}_{ent}$  and  $\mathcal{L}_{div}$  over clip predictions instead 161 of per-frame estimates. (iii) Our formulation is agnostic to the type of data used in the continual 162 learning problem. Without any modifications, our formulation can be applied to both image-based or 163 video-based continual learning problems. 164

## 3.2 Baseline Regularizers

141

142 143

144

145

146

150

151

152

165 166

167

168

169

170

We compare our proposal against different regularizers to assess its effectiveness in mitigating forgetting and improving continual learning performance. We follow the formulation in Equation (3), and study alternatives to  $\mathcal{R}_{IM}(\theta, X_t)$ . In particular, we analyze different regularizers from the continual learning literature, namely Elastic Weight Consolidation [15] and Synaptic Intelligence [36]. Furthermore, we explore Entropy Minimization [30] from the test-time adaptation literature.

Elastic Weight Consolidation (EWC). Kirkpatrick *et.al* proposed to regularize the parameter update during continual learning to prevent catastrophic forgetting by constraining changes to important weights. The key idea behind EWC is to estimate the importance of each parameter for previously learned tasks and penalize deviations from their learned values. We analyze the effectiveness of combining EWC [15] with rehearsal-based methods by replacing  $\mathcal{R}_{\text{IM}}$  in Equation (3) with  $\mathcal{R}_{\text{EWC}}$ , defined as:

$$\mathcal{R}_{\text{EWC}}(\theta) = \sum_{i} \frac{\lambda}{2} F_i (\theta^i - \theta^i_{t-1})^2,$$

where F is the Fisher information matrix, which quantifies the importance of each parameter based on how sensitive the loss function is to changes in that parameter, and  $\lambda$  is a hyper-parameter balancing the relative importance of the old tasks with respect to the current task. Synaptic Intelligence (SI). It is a biologically inspired regularizer from the continual learning literature. It follows a similar principle to EWC but determines weight importance using a different approach. Instead of using the Fisher Information Matrix, SI tracks the contribution of each parameter during training by accumulating an importance measure based on changes in loss. This adaptive tracking mechanism allows the model to selectively constrain updates to crucial parameters while remaining flexible for learning new tasks. We replace  $\mathcal{R}_{\text{IM}}$  in Equation (3) with  $\mathcal{R}_{\text{SI}}$  which takes the following form:

$$\mathcal{R}_{\mathrm{SI}}(\theta) = \sum_{t}^{T} \frac{\omega_{t}^{k}}{(\Delta \theta_{k}^{t})^{2} + \xi},$$

where  $\Delta \theta_k^t = \theta_k^t - \theta_k^{t-1}$  and the damping parameter  $\xi$  avoids division by zero.

Entropy Minimization (EM). Following the self-supervised spirit of our proposed regularization approach, we include one self-supervised regularizer that encourages the model to produce more confident predictions. In particular, we follow Wang *et.al* [30] and apply entropy minimization to regularize the output distribution, reducing the model's uncertainty when making predictions. Entropy minimization replaces  $\mathcal{R}_{\text{IM}}$  in Equation (3) with  $\mathcal{R}_{\text{EM}}$  where:

$$\mathcal{R}_{EM}(\theta, X_t) = -\mathbb{E}_{x \sim X_t} \sum_{k=1}^K f_{\theta}^k(x) \log f_{\theta}^k(x). \tag{4}$$

Entropy minimization encourages the model to assign higher confidence to its predictions, effectively suppressing uncertain outputs. This can be beneficial in a continual learning setup, where distribution shifts can lead to increased uncertainty.

## 4 Experiments

196

202

203

204

205

209

210

211

213 214

216

In this section, we proceed with the empirical assessment of our proposed approach to validate its effectiveness. For completeness, we first evaluate several rehearsal-based continual learning (CL) methods when paired with the regularizers IM, EWC, SI, and EM. We then extend the analysis by applying IM to additional rehearsal approaches, showing that its benefits generalize consistently across different methods and datasets.

**Datasets.** Following the image CL literature, we focus on two main datasets: Split-CIFAR100 [36] and Split-Tiny ImageNet [16]. Split-CIFAR100 contains a total of 100 classes and 6000 images per class. It is divided into 10 tasks, each containing 10 classes. Split-Tiny ImageNet consists of 200 classes with 500 images per class, and is divided into 10 tasks of 20 classes each.

Evaluation Metrics. To evaluate the performance of CL methods, we consider two metrics: Average Accuracy, which is defined as the average performance across all tasks, and Forgetting Rate that measures the impact of the learned task on the performance of the previous tasks [8].

• Average Accuracy (ACC) quantifies the model's overall performance across all tasks it has encountered. It is defined as:

$$ACC = \frac{1}{T} \sum_{i=1}^{T} a_i, \tag{5}$$

where T is the total number of tasks and  $a_i$  represents the accuracy of the model on the i-th task after it has been trained on all T tasks. ACC provides a comprehensive measure of how well the model learns and retains knowledge across a full sequence of tasks.

• Forgetting Rate (FR) measures the decrease in performance on past tasks after a model has been trained on new ones. It directly measures catastrophic forgetting. FR is defined as:

$$FR = \frac{1}{T - 1} \sum_{i=1}^{T-1} \max_{j < T} (a_{ij} - a_{iT}), \tag{6}$$

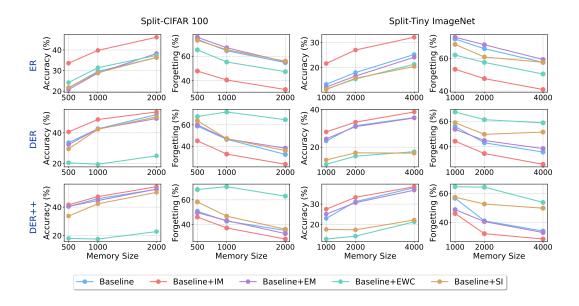


Figure 1: **Results of Integrating Different Regularizers on Split-CIFAR100 and Split-Tiny ImageNet.** This figure plots the average accuracy and forgetting rate of three baseline methods (ER, DER, and DER++) across various memory sizes, in combination with the analyzed regularizers (IM, EM, EWC, and SI), and across two datasets (Split-CIFAR100 and Split-Tiny ImageNet).

where  $a_{ij}$  is the accuracy on task i immediately after training on task j, and  $a_{iT}$  is the accuracy on task i after the final task T has been learned. A lower FR indicates better retention of previously learned knowledge, while a higher FR points to significant forgetting.

**Implementation Details.** We train a ResNet18 [13] model from scratch, following the training scheme and hyper-parameters of each paper. Following our limited memory setting, we define a budget of 5, 10 and 20 samples per class as the maximum allowed in  $\mathcal{M}_t$  at any moment. To balance the loss terms, we multiply the regularization term by  $\lambda$ =0.5 and the cross-entropy loss with 1- $\lambda$ .

**Baselines.** We consider four regularizers: EWC [15], SI [36], EM [30], and IM [18], applied on top of three memory-based continual learning methods: ER [26], DER [5], and DER++ [5]. For the IM regularizer, we further extend the evaluation to include Refresh Learning [32], implemented on top of DER++, and STAR [12], implemented on top of ER.

#### 4.1 Regularized Rehearsal Methods Results

Figure (1) summarizes the performance of rehearsal-based methods ER, DER, and DER++ on Split-CIFAR100 (left columns) and Split-Tiny ImageNet (right columns) for the selected memory sizes. We include the baseline performance (light blue) and outline the impact of incorporating regularization techniques on top of these rehearsal methods.

Our results demonstrate that introducing IM on top of rehearsal-based methods consistently leads to improvements across all memory sizes. For instance, when IM is applied to ER on Split-CIFAR100, we observe an enhancement of 10-13% in performance across all memory sizes. The improvement is slightly lower for DER and DER++, ranging around 2-7% for DER and 1-2% for DER++. In contrast, other regularizers like EWC, SI, and EM generally do not improve the baseline methods, and can even degrade performance in some cases, as is the case for EWC on DER and DER++, where the model's accuracy drops by nearly half.

We can also observe in Figure (1) that forgetting is significantly reduced when the rehearsal methods are paired with IM. For ER (paired with IM) applied on Split-CIFAR100, the reduction in forgetting is around (22-25%) across all memory sizes. For DER and DER++ on Split-CIFAR100, our results show a smaller reduction compared to ER with about (9-13%) and (4-6%), respectively. On the other

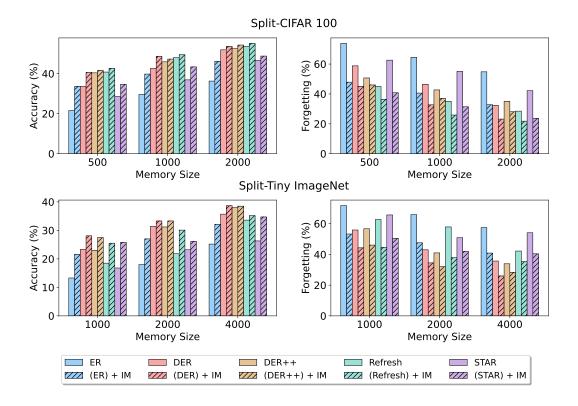


Figure 2: **Results on Split-CIFAR100 and Split-Tiny ImageNet.** This figure presents the average accuracy and forgetting rate of five rehearsal-based methods (ER, DER, DER++, Refresh Learning, and STAR) across different memory buffer sizes, both in their baseline form and when combined with Information Maximization (IM).

hand, EM and SI do not generally reduce forgetting compared to the original baselines. This is since EM aims at increasing the model's confidence in predicting samples from the current task. While this approach might accelerate the learning process over tasks, it does not promote the retention of previously learned information.

For EWC, we observe reduction in forgetting when paired with ER, but not with the remaining baselines. More detailed and comprehensive results can be found in the **Appendix**.

To further validate the effectiveness of IM, we expand our analysis to include two additional rehearsal-based baselines: Refresh Learning and STAR. Figure (2) presents the performance gains when IM is integrated into all five rehearsal methods (ER, DER, DER++, Refresh Learning, STAR) across both Split-CIFAR100 and Split-Tiny ImageNet.

On Split-CIFAR100, IM consistently improves the performance of Refresh Learning and STAR, but with varying impact. When paired with Refresh Learning, the improvements are more modest—about 1–2% in accuracy—yet forgetting is reduced by (7–9%) across all memory sizes. This indicates that while Refresh Learning already stabilizes training to some degree, IM provides an additional layer of retention without significantly altering the learning dynamics. For STAR, the effect is more pronounced: we observe accuracy gains of 2–6%, along with a substantial reduction in forgetting (18–24%). On Split-Tiny ImageNet, IM continues to yield consistent improvements. For Refresh, the accuracy gains are larger than on Split-CIFAR100, ranging from 2–8%, with forgetting reduced by (7–20%). STAR also benefits considerably, with accuracy improvements of 3–9% and forgetting reductions between (9–15%).

**Conclusion.** These results reveal that Information Maximization (IM) serves as an effective regularization technique, consistently improving the performance of image continual learning baselines

Table 1: **Ablation Study on Compute Budget.** This table presents the performance of ER and DER, with and without the Information Maximization (IM) regularizer, on Split-CIFAR100 and Split-Tiny ImageNet datasets. For Split-CIFAR100, the baselines are run with 10 epochs, while for Split-Tiny ImageNet, the experiments are run with 50 epochs.

	Split-Cifar100			Split-Tiny ImageNet			
<b>Buffer Size</b>	500	1000	2000	1000	2000	4000	
ER	20.8	26.8	35.6	13.2	18.7	25.6	
ER + IM	<b>28.8</b>	<b>35.7</b>	<b>40.1</b>	<b>21.5</b>	<b>26.6</b>	<b>31.8</b>	
DER	24.1	21.1	19.9	21.6	26.3	24.2	
DER + IM	<b>33.5</b>	<b>37.3</b>	<b>34.4</b>	<b>27.9</b>	<b>32.1</b>	<b>33.3</b>	

266 across various memory budgets. By encouraging confident yet balanced predictions, IM enhances both accuracy and knowledge retention, thereby mitigating catastrophic forgetting.

## 4.2 Ablation Analysis

To explore the limitations of Information Maximization (IM) as a regularizer for continual learning methods, we conduct two ablation experiments aimed at understanding its performance under various conditions. First, we assess the impact of IM when the computational budget is reduced to determine whether it improves the convergence of the baseline methods. Second, we evaluate if the improvement obtained by using IM diminishes with additional tasks.

**Computational Budget.** In Section (4.1), we followed Mammoth [3, 6] defaults of 50 epochs per task for Split-CIFAR100 and 100 for Split-Tiny ImageNet, ensuring stable training. However, compute efficiency is increasingly critical in continual learning, as resources are costly compared to storage. To assess IM under computational constraints, we ran ablations with 10 epochs per task on Split-CIFAR100 and 50 on Split-Tiny ImageNet.

The results in Table (1) show that, even with a lower computational budget of 10 and 50 epochs per task on Split-CIFAR100 and Split-Tiny ImageNet, respectively, the proposed ER+IM and DER+IM methods outperform their counterparts without IM regularizer. For instance, on the Split-CIFAR100 dataset with a buffer size of 1000, ER+IM achieves an accuracy of 35.7%, significantly higher than ER at 26.8%. Similarly, DER+IM attains 37.3% accuracy, surpassing DER's 21.1% by a large margin. These trends hold for different buffer sizes and datasets, highlighting the effectiveness of the proposed regularization in low-compute regimes.

**Number of Tasks.** In the experiments presented in Section (4.1), we used the conventional 10-tasks split for Split-CIFAR100 and Split-Tiny ImageNet, which is commonly used in continual learning studies. However, as shown in [29, 24], performance may vary when more tasks are introduced. More tasks can make the problem harder because the model has to remember more information and avoid forgetting earlier tasks while learning new information. Consequently, we reran the experiments in Section (4.1), doubling the number of tasks from 10 to 20. This allows us to evaluate whether Information Maximization (IM) regularizer remains effective when the continual learning problem becomes more challenging due to having more tasks to learn. The results presented in Table (2) show that incorporating IM into ER and DER can significantly improve their performance on longer sequences of tasks. For example, ER+IM shows an improvement of (4-7%) and (6-8%) on Split-CIFAR100 and Split-Tiny ImageNet, respectively. On the other hand, DER+IM shows a (4-7%) and (4-5%) improvement on Split-CIFAR100 and Split-Tiny ImageNet, respectively.

## 4.3 Generalization to Video Continual Learning

To further validate the effectiveness of Information Maximization (IM) as a cost-effective regularization technique for continual learning methods, we extend our analysis to experiments in the

Table 2: **Ablation Study on Number of Tasks.** This table presents the performance of ER and DER , with and without Information Maximization (IM) regularizer on a sequence of 20 tasks for both Split-CIFAR100 and Split-Tiny ImageNet datasets.

	Sp	lit-CIF	AR100	Split-Tiny ImageNet			
<b>Buffer Size</b>	500	1000	2000	1000	2000	4000	
ER	16.6	25.8	34.7	9.1	14.5	22.1	
ER + IM	<b>23.4</b>	<b>31.3</b>	<b>38.9</b>	<b>18.3</b>	<b>23.4</b>	<b>29.1</b>	
DER	25.1	35.8	38.9	18.0	22.5	27.8	
DER + IM	<b>32.8</b>	<b>39.8</b>	<b>45.0</b>	<b>23.2</b>	<b>27.0</b>	<b>32.1</b>	

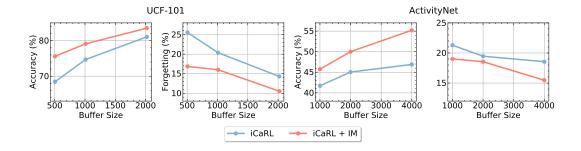


Figure 3: **Application of Information Maximization to Video Continual Learning.** This figure illustrates the average accuracy and forgetting rates of the iCARL video continual learning variant, introduced by vCLIMB [29], with and without our Information Maximization (IM) regularizer.

video domain. Specifically, we experiment with the iCaRL approach as part of the popular vCLIMB 303 framework [25, 29] to assess IM's performance in this video CL context. We evaluate the use of 304 IM regularizer on two widely recognized datasets in the video domain: UCF-101, consisting of 101 305 classes split into 10 tasks, and ActivityNet, comprising 200 classes also divided into 10 tasks. We set 306 the memory size to 5, 10 and 20 samples per class and adopt the training hyperparameters from [29]. 307 Upon applying iCARL to UCF-101 with IM, we achieve an improvement of 2-8% across all memory 308 sizes. Similarly, on ActivityNet, the accuracy gain is between 4-8% with the incorporation of IM. 309 These results highlight the potential of IM as a valuable regularization technique for enhancing 310 performance in the video continual learning scenarios. Note that video data has an additional 311 temporal dimension compared to images, which requires more memory to store. Being able to improve continual learning performance with small memory buffer sizes, as shown in Figure 3, is crucial for facilitating the development of memory-efficient approaches for video continual learning.

#### 5 Conclusion

315

316

317

318

319

320

321

322

323

324

325

326

In conclusion, this paper explores the combined potential of memory-based methods and regularization techniques in the context of Continual Learning (CL), specifically within a class incremental setup. We introduce a novel, class-agnostic regularization strategy for CL, which focuses on the distribution of the network's predictions. This strategy, termed Information Maximization (IM) regularization, facilitates the learning of enhanced feature representations across multiple distribution shifts, while simultaneously minimizing memory requirements and computational overhead. Our extensive empirical evaluation underscores the effectiveness of the proposed IM regularizer. Furthermore, the simplicity and versatility of our approach allow it to be applied across different input domains, as evidenced by its successful application in the video continual learning setup. Unlike traditional image-based settings, video CL presents additional challenges due to its temporal structure and higher memory demands. Despite these complexities, our method demonstrates strong performance, reinforcing its applicability to real-world, resource-constrained scenarios.

#### References

- [1] Motasem Alfarra, Hani Itani, Alejandro Pardo, Shyma Alhuwaider, Merey Ramazanova, Juan C. Pérez,
   Zhipeng Cai, Matthias Müller, and Bernard Ghanem. Revisiting test time adaptation under online evaluation,
   2023.
- [2] Lama Alssum, Juan León Alcázar, Merey Ramazanova, Chen Zhao, and Bernard Ghanem. Just a glimpse:
   Rethinking temporal information for video continual learning. In CVPRW, pages 2473–2482, 2023.
- [3] Matteo Boschini, Lorenzo Bonicelli, Pietro Buzzega, Angelo Porrello, and Simone Calderara. Class incremental continual learning into the extended der-verse. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online test-time
   adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
   pages 8344–8353, 2022.
- [5] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience
   for general continual learning: a strong, simple baseline. *NeurIPS*, 33:15920–15930, 2020.
- Fietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience
   for general continual learning: a strong, simple baseline. In H. Larochelle, M. Ranzato, R. Hadsell,
   M. F. Balcan, and H. Lin, editors, Advances in Neural Information Processing Systems, volume 33, pages
   15920–15930. Curran Associates, Inc., 2020.
- [7] Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky.
   New insights on reducing abrupt representation change in online continual learning. arXiv preprint
   arXiv:2104.05025, 2021.
- 349 [8] Arslan Chaudhry, Marc'Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-GEM. In *ICLR*, 2019.
- [9] 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.
- [10] Geng Chen, Wendong Zhang, Han Lu, Siyu Gao, Yunbo Wang, Mingsheng Long, and Xiaokang Yang.
   Continual predictive learning from videos. In CVPR, pages 10728–10737, 2022.
- [11] Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning online from non-stationary data streams. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8250–8259, 2021.
- Masih Eskandar, Tooba Imtiaz, Davin Hill, Zifeng Wang, and Jennifer Dy. Star: Stability-inducing weight
   perturbation for continual learning. arXiv preprint arXiv:2503.01595, 2025.
- [13] 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.
- 14] Haeyong Kang, Jaehong Yoon, Sung Ju Hwang, and Chang D Yoo. Continual learning: Forget-free winning subnetworks for video representations. *arXiv* preprint *arXiv*:2312.11973, 2023.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu,
   Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic
   forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- 16 Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. CS 231N, 7(7):3, 2015.
- Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch normalization for practical domain adaptation. *arXiv preprint arXiv:1603.04779*, 2016.
- 18] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning*, pages 6028–6039. PMLR, 2020.
- Yuejiang Liu, Parth Kothari, Bastien Van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre
   Alahi. Ttt++: When does self-supervised test-time training fail or thrive? Advances in Neural Information
   Processing Systems, 34:21808–21820, 2021.

- 377 [20] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *NeurIPS*, 378 30, 2017.
- 279 [21] Amir Nazemi, Zeyad Moustafa, and Paul Fieguth. Clvos23: A long video object segmentation dataset for continual learning. In *CVPR*, pages 2495–2504, 2023.
- [22] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui Tan.
   Efficient test-time model adaptation without forgetting. In *International conference on machine learning*,
   pages 16888–16905. PMLR, 2022.
- Jaeyoo Park, Minsoo Kang, and Bohyung Han. Class-incremental learning for action recognition in videos.
   In Proceedings of the IEEE/CVF international conference on computer vision, pages 13698–13707, 2021.
- Ameya Prabhu, Hasan Abed Al Kader Hammoud, Puneet K Dokania, Philip HS Torr, Ser-Nam Lim,
   Bernard Ghanem, and Adel Bibi. Computationally budgeted continual learning: What does matter? In
   Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 3698–3707,
   2023.
- 390 [25] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental 391 classifier and representation learning. In CVPR, pages 2001–2010, 2017.
- [26] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay
   for continual learning. Advances in Neural Information Processing Systems, 32, 2019.
- Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
   with self-supervision for generalization under distribution shifts. In *International conference on machine* learning, pages 9229–9248. PMLR, 2020.
- [28] Andrés Villa, Juan León Alcázar, Motasem Alfarra, Kumail Alhamoud, Julio Hurtado, Fabian Caba
   Heilbron, Alvaro Soto, and Bernard Ghanem. Pivot: Prompting for video continual learning. In CVPR,
   pages 24214–24223, 2023.
- 400 [29] Andrés Villa, Kumail Alhamoud, Victor Escorcia, Fabian Caba, Juan León Alcázar, and Bernard Ghanem.
   401 vclimb: A novel video class incremental learning benchmark. In CVPR, pages 19035–19044, 2022.
- 402 [30] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 404 [31] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In
   405 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7201–7211,
   406 2022.
- 407 [32] Zhenyi Wang, Yan Li, Li Shen, and Heng Huang. A unified and general framework for continual learning.
   408 arXiv preprint arXiv:2403.13249, 2024.
- Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong
   Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual
   learning. In European conference on computer vision, pages 631–648. Springer, 2022.
- 412 [34] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent
   413 Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In CVPR, pages 139–149,
   414 2022.
- Linfeng Xu, Qingbo Wu, Lili Pan, Fanman Meng, Hongliang Li, Chiyuan He, Hanxin Wang, Shaoxu Cheng, and Yu Dai. Towards continual egocentric activity recognition: A multi-modal egocentric activity dataset for continual learning. *arXiv* preprint arXiv:2301.10931, 2023.
- 418 [36] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *ICML*, pages 3987–3995. PMLR, 2017.

# 420 A Appendix

421

#### A.1 Regularized Rehearsal Methods Results

Tables 3 and 4, contain the numerical results for the average accuracy and forgetting rate metrics, respectively, on three baseline methods (ER, DER, and DER++) in combination with the analyzed regularizers (IM, EM, EW, and SI), as well as for integrating IM into Refresh Learning and STAR. These results were summarized as plots in Figures (1) and (2) of the main paper. We observe that our proposed regularize (IM) consistently outperforms other methods across different memory settings.

Table 3: Average Accuracy on Split-CIFAR100 and Split-Tiny ImageNet. This table shows the average accuracy of three baseline methods (ER, DER, and DER++) across various sizes of memory buffer, in combination with the analyzed regularizers (IM, EM, EW, and SI), as well as for integrating IM into Refresh Learning and STAR.

Dataset	Spli	t-CIFAI	R100	Split-Tiny ImageNet		
Buffer	500	1000	2000	1000	2000	4000
ER	21.6	29.6	36.3	13.3	18.0	25.2
ER (IM)	33.6	39.8	46.2	21.6	27.0	32.1
ER (EM)	20.8	28.8	38.3	12.2	16.6	24.1
ER (EWC)	24.2	31.4	37.4	11.4	15.4	21.3
ER (SI)	21.8	29.0	36.4	11.2	15.7	20.4
DER	33.6	42.6	51.9	23.3	31.4	35.7
DER (IM)	40.6	48.7	53.6	28.1	33.3	38.7
DER (EM)	32.3	42.6	49.3	24.4	30.9	35.5
DER (EWC)	20.5	19.5	25.0	10.9	15.2	17.6
DER (SI)	29.5	42.3	50.4	13.2	17.0	16.8
DER++	40.4	45.9	52.5	23.0	31.2	38.0
DER++(IM)	41.6	47.3	54.3	27.5	33.3	38.5
DER++(EM)	40.5	44.6	52.6	25.1	30.7	36.8
DER++ (EWC)	18.0	17.6	22.8	12.9	14.3	21.3
DER++ (SI)	33.7	42.3	50.3	17.6	17.4	22.2
Refresh	40.8	48.0	53.6	18.5	21.8	33.6
Refresh (IM)	42.7	49.5	55.1	25.5	30.1	35.2
STAR	28.6	36.9	46.6	16.8	23.2	26.3
STAR (IM)	34.7	43.4	48.8	25.8	26.1	34.7

# A.2 Regularization Targets

427

For the experimental assessment in Section (4.1), we apply the regularization loss to current task 428 samples only. This raises the question of how the proposed method would behave if the IM loss were 429 applied exclusively to memory samples or to both memory and current task samples. For this reason, 430 we reran the experiments shown in Section (4.1) for both variants, and the results are summarized 431 in Table (5). We find that applying the IM loss to the current task (CT) is superior to applying it to 432 the memory/buffer samples only (BF) or to both buffer and current task samples (ALL). Notably, 433 this trend remains consistent across various buffer sizes, datasets, and continual learning methods, 434 highlighting the robustness of this strategy. 435

For example, with a buffer size of 500 on the Split-CIFAR100, ER+IM (CT) achieves an accuracy of 33.6%, which is significantly higher than ER+IM (ALL) and ER+IM (BF), which achieve 25.9% and 21.7%, respectively. Similarly, DER+IM (CT) consistently outperforms DER+IM (ALL) and DER+IM (BF) across various buffer sizes and datasets, reinforcing the advantage of applying IM

Table 4: **Forgetting Rate on Split-CIFAR100 and Split-Tiny ImageNet.** This table shows the forgetting rate of three baseline methods (ER, DER, and DER++) across various sizes of memory buffer, in combination with the analyzed regularizers (IM, EM, EW, and SI), as well as the results for integrating IM into Refresh Learning and STAR.

Dataset	Spli	t-CIFAI	R100	Split-Tiny ImageNet		
Buffer	500	1000	2000	1000	2000	4000
ER	73.8	64.5	54.8	71.8	65.9	57.5
ER (IM)	47.9	40.6	32.8	53.3	47.6	40.9
ER (EM)	75.5	66.9	55.5	72.9	68.3	59.3
ER (EWC)	65.2	55.3	47.4	62.0	57.5	50.4
ER (SI)	73.0	65.3	56.0	68.5	60.8	57.6
DER	58.8	46.5	32.4	55.9	43.0	35.7
DER (IM)	45.1	32.7	23.2	44.3	34.5	26.0
DER (EM)	60.0	46.9	38.2	53.7	44.8	38.4
DER (EWC)	67.7	72.0	64.8	67.4	61.3	58.8
DER (SI)	63.8	47.1	36.0	58.9	49.7	51.5
DER++	50.7	42.7	35.1	56.7	41.0	33.9
DER++(IM)	46.0	37.1	28.1	46.1	32.1	28.4
DER++(EM)	49.6	43.1	32.5	48.8	40.5	32.9
DER++ (EWC)	68.4	70.7	63.0	64.7	64.4	53.8
DER++(SI)	58.3	46.7	36.1	57.4	52.7	49.8
Refresh	45.0	35.1	28.5	62.7	57.9	42.2
Refresh (IM)	36.4	25.9	21.8	44.6	37.9	35.4
STAR	62.6	55.1	42.2	65.7	50.9	54.1
STAR (IM)	40.8	31.4	23.7	50.4	42.0	40.4

solely to current task samples. For example, it achieves 48.7% accuracy on Split-CIFAR100 with a buffer size of 1000, while DER +IM (ALL) and DER +IM (BF) achieve 46.0% and 41.0%, respectively. These results indicate that applying IM regularizer to the current task samples is more effective than applying it exclusively to the memory/buffer samples only or to both buffer and current task samples.

## 445 A.3 Generalization to Video Continual Learning

Tables 6 and 7 present the numerical results of average accuracy and forgetting rate, respectively, for iCaRL approach with and without IM on two video datasets (UCF101 and ActivityNet). Combining iCaRL with IM shows improvement in average accuracy and reduces forgetting rate across different memory settings.

Table 5: **Ablation Study on Regularization Target Selection.** In the main results, the Information Maximization (IM) regularizer is applied exclusively to the current task samples (CT). This table presents the results of applying the regularizer to the buffer samples only (BF) and to both current task samples and buffer samples simultaneously (ALL). The findings indicate that applying the regularizer to the current task samples consistently leads to superior performance compared to the other variants.

	Split-CIFAR100			Split-Tiny ImageNet		
<b>Buffer Size</b>	500	1000	2000	1000	2000	4000
ER + IM (ALL)	25.9	34.6	42.7	15.0	20.3	27.8
ER + IM (CT)	<b>33.6</b>	<b>39.8</b>	<b>46.2</b>	<b>21.6</b>	<b>27.0</b>	<b>32.1</b>
ER + IM (BF)	21.7	28.9	38.8	12.7	17.7	24.2
DER + IM (ALL)	33.5	46.0	53.5	27.7	32.7	35.1
DER + IM (CT)	<b>40.6</b>	<b>48.7</b>	<b>53.6</b>	<b>28.1</b>	<b>33.3</b>	<b>38.7</b>
DER + IM (BF)	27.3	41.0	50.2	21.2	27.7	34.6

Table 6: **Average Accuracy on UCF101 and ActivityNet.** This table shows the average accuracy of iCaRL with and without IM across various sizes of memory buffer on two datasets (UCF101 and ActivityNet). The results demonstrate that the proposed information maximization approach (+IM) consistently outperforms iCaRL on both datasets regardless of the memory setting.

Dataset		UCF101		ActivityNet		
Buffer	505	1010	2020	1000	2000	4000
iCaRL	68.44	74.70	81.07	41.72	45.05	46.91
iCaRL (IM)	75.60	79.11	83.53	45.76	50.01	55.20

Table 7: **Forgetting Rate on UCF101 and ActivityNet.** This table shows the forgetting rate of iCaRL with and without IM across various sizes of memory buffer on two datasets (UCF101 and ActivityNet). The results demonstrate that the proposed information maximization approach (+IM) consistently achieves lower forgetting rates on both datasets regardless of the memory setting.

Dataset		UCF101		A	ctivityNe	et
Buffer	505	1010	2020	1000	2000	4000
iCaRL	25.55	20.37	14.31	21.29	19.46	18.55
iCaRL + IM	16.87	16.02	10.54	19.01	18.55	15.48