

ACAE-REMIND for Online Continual Learning with Compressed Feature Replay

Kai Wang¹, Luis Herranz¹, Joost van de Weijer¹

¹ Computer Vision Center, Universitat Autònoma de Barcelona, Barcelona, Spain

{kwang, lherranz, joost}@cvc.uab.es

Abstract

Online continual learning aims to learn from a non-IID stream of data from a number of different tasks, where the learner is only allowed to consider data once. Methods are typically allowed to use a limited buffer to store some of the images in the stream. Recently, it was found that feature replay, where an intermediate layer representation of the image is stored (or generated) leads to superior results than image replay, while requiring less memory. Quantized exemplars can further reduce the memory usage. However, a drawback of these methods is that they use a fixed (or very intransigent) backbone network. This significantly limits the learning of representations that can discriminate between all tasks. To address this problem, we propose an auxiliary classifier auto-encoder (ACAE) module for feature replay at intermediate layers with high compression rates. The reduced memory footprint per image allows us to save more exemplars for replay. In our experiments, we conduct task-agnostic evaluation under online continual learning setting and get state-of-the-art performance on ImageNet-Subset, CIFAR100 and CIFAR10 dataset.

1. Introduction

The vast majority of deep learning papers consider that all training data is available jointly, and the learner can process the data several times (epochs) to learn the optimal parameters to solve the task at hand. However, in many real-world scenarios, this would not be possible, and the learner has only access to data of a single task at the time, before proceeding to learn a new task. This scenario refers to *continual learning* (or *incremental learning*, *lifelong learning*). The main challenge in this scenario is to learn from the current data while preventing forgetting the knowledge of previous tasks. With a naive finetuning approach the model will suffer a drastic drop in performance on previous tasks because the model aims to be optimal for the current tasks, and ignores performance on previous tasks. This

phenomenon is known as *catastrophic forgetting* [17, 27]. The field of continual learning studies methods that prevent forgetting [8, 10, 11, 30, 35, 26].

A challenging setting in continual learning, yet common in practical application, is *online continual learning* of non-iid data streams [1, 5, 22]. In the online setting, each image can only be observed once during model optimization (except exemplars in storage). These applications mainly exist in resource constrained devices, such as mobile phones, robots and other smart devices. The majority of methods in continual learning, known as batch incremental learning methods, allow for several cycles (epochs) over the data [6]. These methods cannot operate in the challenging online continual learning setting. Moreover they take longer to train. In this paper, we focus on online continual learning.

Among the approaches to address catastrophic forgetting, some of the best performing ones are rehearsal-based [30, 33, 10]. Several methods save a small set of exemplar images of previous classes [30, 5, 11, 35]. Retrieving them during future training is a straightforward way to prevent forgetting. For example, GEM [22], A-GEM [5] and MIR [1], which address online continual learning, belong to this type. However, this strategy leads to increased memory usage and the problem of training from imbalanced data (between previous tasks and the current task). An alternative is to generate images via generative models (e.g. GANs) [33, 34]. However, image generation is still a difficult problem in computer vision and requires complex generative models, which would also need to be continually learned, making this method not practical for complex datasets.

To circumvent the difficulties of image replay, recent work has focused on feature replay [10, 21]. In [21] a generator is trained to replay compact feature representations of the images (after the last average pooling layer of a ResNet-18). In addition, a distillation loss was applied to prevent forgetting of the feature extractor. Instead of generating features, it was also observed by Hayes *et al.* [10] that saving them as exemplars is very efficient, since it required less

memory per image. To even further reduce the memory requirements, the REMIND method [10] also applies product quantization [13]. This allows them to save up to 1M compressed feature exemplars, instead of 20K exemplar images saved traditionally, in the same memory buffer. A major drawback of these feature replay methods [21, 10] is that they either allow for very little training [21] or no training at all [10] of the backbone feature extractor (located before the replay layer). As a consequence, if this backbone is not yet optimally trained for future tasks, the performance is sub-optimal.

To address the limitation of feature replay, we propose ACAE-REMIND, an auxiliary classifier auto-encoder (ACAE) that allows for compressed feature replay (as in REMIND) at intermediate layers of the network. This contrasts with current feature replay methods that focus on replaying the features in the last layers. The principal advantage of our method is that we can jointly train all layers after the replay layer. This addresses an important problem of feature replay methods, namely the reduced performance because of a large fixed backbone network. Instead of only the 4M parameters that are trained in REMIND when replaying at block 4 of a ResNet-18, we allow to train 9M parameters jointly when replaying on block 3. This leads to feature representations that are more discriminative between the classes of current and previous tasks. We evaluate our method in the challenging, yet more realistic, *online continual learning* setting. From experiments under multiple settings and datasets, we observe state-of-the-art performance in many-task evaluations and competitive in few-task settings.

2. Related work

2.1. Continual learning

Continual learning methods can be categorized into three types which we will shortly comment. For a more elaborate overview see the following surveys [6, 25]).

Regularization-based methods. The first group of techniques is based on regularization. These methods add a regularization term to the loss function which impedes changes to the parameters deemed relevant to previous tasks. The difference depends on how to compute the estimation. From these differences, these methods can be further divided into data-focused [20] and prior-focused [17]. Data-focused methods use knowledge distillation from previously learned models. Prior-focused methods estimate the importance of model parameters as a prior for the new model. However, it has been shown that data-focused methods are vulnerable to domain shifts [2] and prior-focused methods might be not sufficient to restrict the optimization process to keep acceptable performances on previous tasks [9].

Parameter isolation methods. This family focuses on

allocating different model parameters to each task. These models begin with a simplified architecture and updated incrementally with new neurons or network layers in order to allocate additional capacity for new tasks. In Piggyback/PackNet [23, 24], the model learns a separate mask on the weights for each task, whereas in HAT [32], masks are applied to the activations. This method is further developed to the case where no forgetting is allowed in [26]. In general, this branch is restricted to the task-aware (task incremental) setting. Thus, they are more suitable for learning a long sequence of tasks when a task oracle is present and there is no constraint over model capacities.

Replay methods. This type of methods prevent forgetting by including data (real or synthetic) from previous tasks, stored either in an episodic memory or via a generative model. There are two main strategies: exemplar rehearsal [30, 5, 11, 35] and pseudo-rehearsal [33, 34]. The former store a small amount of training samples (also called exemplars) from previous tasks. The latter use generative models learned from previous data distributions to synthesize data.

One of the main drawbacks of exemplar replay is the high memory usage required to store exemplars of previous tasks. REMIND [10] addresses this drawback, instead of saving original data, it saves compressed latent representation of intermediate layer features via product quantization [13]. This is a more efficient usage of memory and computation. However, due to restriction that the backbone is fixed, the majority of feature extraction modules cannot be adopted to later tasks. Therefore, this model has a strong bias towards the first task. Recently, GDumb [28] proposes training a model only from exemplars. The main idea is to balance the sample reservoir in a selection stage, then the model is learnt from scratch on this balanced set. While not designed for any specific continual learning settings, it achieves excellent performance on many. It reveals that sample balancing is crucial for rehearsal-based continual learning methods. While saving images is always expensive compared to saving features, this point is also mentioned in paper on feature adaption [12].

2.2. Auto-encoders and product quantization

Auto-encoders [18] learn representations in an unsupervised way by encouraging the model to reconstruct the input data. An encoder projects the high-dimensional input to a low-dimensional space, and the decoder tries to project back to the original space minimizing the reconstruction error. Product quantization [13] is an effective quantization method that performs a decomposition of a high-dimensional space into the Cartesian product of a series of subspaces, and quantizes them separately.

In our model, we propose an auto-encoder with an auxiliary classifier (ACAE) to force the reconstructions not only

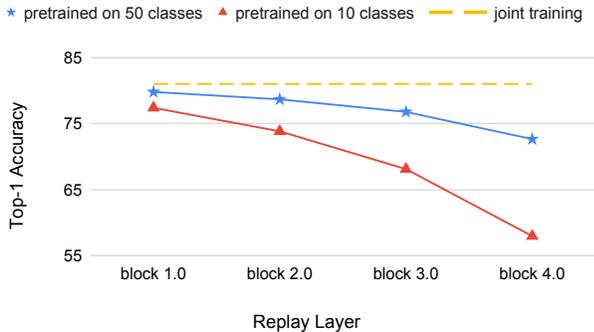


Figure 1: Drop in performance due to frozen backbones (Joint training: 81.0)

to remain close to original inputs but also keeping the classification characteristics. By combining ACAE with Product Quantization (PQ), the feature spaces are decomposed from high dimension to low dimension, from float numbers to integer indexes, which leads to better compression and therefore allows to save more exemplars.

3. Compressed Feature Replay

3.1. Feature replay location

Pseudo-rehearsal methods [33, 34] are limited by the performance of generative models to generate high-quality images. As a result these methods perform poorly on more complex real-world datasets. To address this limitation, Liu *et al.* [21] proposed generative feature replay (GFR) to generate features of an intermediate layer. In their proposal, features before the classifier are generated by a conditional GAN learned in a continual fashion.

REMIND [10] observes that storing features is much more efficient than storing images. They operate on the same block as GFR. With the help of PQ [13] the features are further compressed in REMIND. The features from block 4.0 of ResNet-18 are approximated by a number of codebooks and indexed feature maps. By this means, the floating point values of feature representations are replaced by integer index numbers. This allows them to save $50\times$ more feature exemplars than image exemplars in the same memory space and obtain excellent results for online continual learning.

As noted in the introduction, one of the main drawbacks of REMIND (and feature replay in general) is that these methods freeze the backbone feature extractor (i.e. the layers before the feature replay) after training the first task. They only train the layers that come after the feature replay layer for the remaining tasks. Depending on the continual learning scenario this could lead to a significant drop in per-

Table 1: Comparison of replay methods.

method name	replay layer	compression	online
GFR	last block	GAN	✗
REMIND	last block	PQ	✓
Ours	interm. block	ACAЕ+PQ	✓

formance because of a suboptimal backbone network.

To better understand the impact of freezing the backbone network after the first task we perform an experiment on the ImageNet-Subset dataset [7] with ResNet-18 as the backbone. We consider two scenarios, one with the first task containing 50 classes and the remaining 50 classes divided into five more tasks. In the second scenario, we evenly divide the classes over 10 tasks (each with 10 classes). Clearly, the second scenario is more challenging for REMIND, because now the backbone network can only be trained on the 10 classes of the first task. In Fig. 1 we can see the drop in performance which is caused by freezing the backbone network as a function of the position of the feature replay. The performance of the different backbone networks is computed in the following way: we first train the first task and fix the backbone network, then we jointly train the remaining layers on all training data of all tasks (this can be seen as the upper bound for this continual learning setting, i.e. joint training with the backbone frozen). As can be seen, the drop in performance is significant (by comparing the difference of the blue and red with the yellow line), dropping 8.36% in the first scenario, and 23.04% in the second scenario when replaying the features of block 4.0. As can be seen the drop diminishes considerably by performing the replay at earlier layers. The reason why REMIND chooses to replay at block 4.0 is because the proposed technique does not scale well to lower positions in the network. This is explicitly discussed and they mention that the quantized features would be too large and would significantly increase storage requirements¹.

To overcome the limitations of feature replay, our proposed ACAE-REMIND model aims to apply replay on an intermediate blocks. To reconstruct features in intermediate layers, we introduce a stronger compression module, which achieves dimension reduction, and feature approximation while maintaining the classification characteristics of the replayed features. The method is an extension of REMIND and is based on an Auto-Encoder with Auxiliary Classifier (ACAЕ). We can perform joint training on all layers after the replay layer (and not only the last block as in REMIND). This alleviates the drawback of fixing the backbone neural network. A comparison among the discussed feature replay methods is in Table 1.

¹See Supplementary material S2 [10].

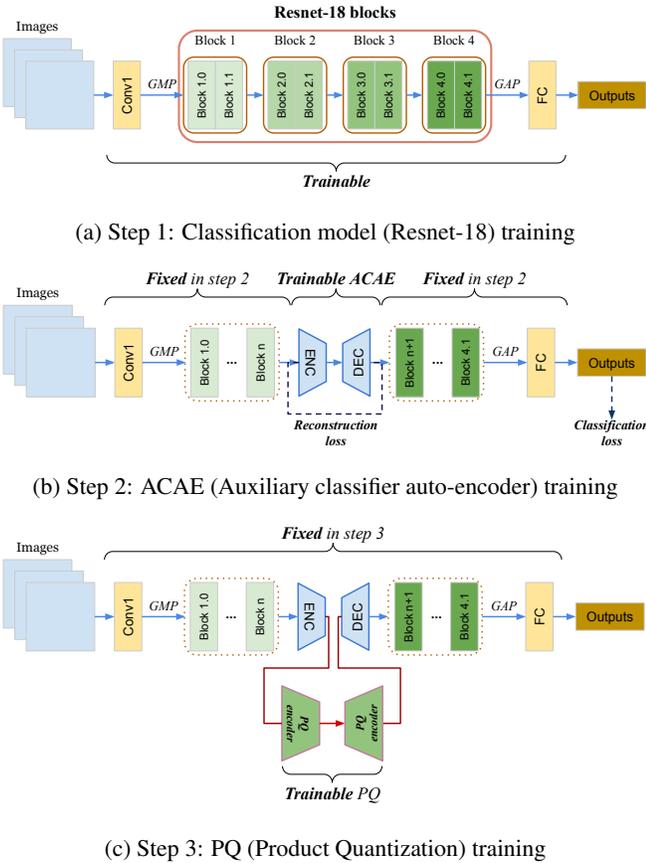


Figure 2: Overview of the initialization stage (trained on first task).

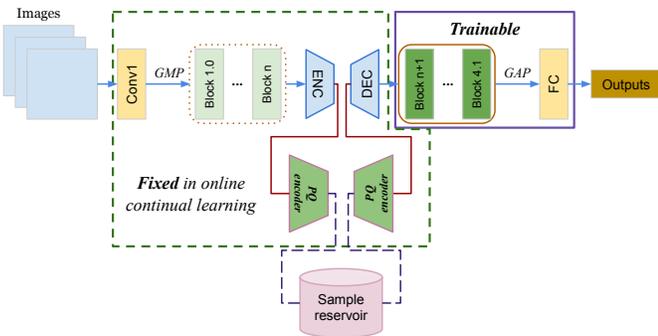


Figure 3: Overview of our online continual learning phase (task $t = 2, \dots, T$).

3.2. Online continual learning setting

Online continual learning is a subarea of incremental learning, where the algorithm is only allowed to make a sin-

gle pass through the data of each task. It is more related to real-life and real-time applications since data comes in sequential streams, and they are not allowed to use the same sample more than one time (unless stick to buffer) in the whole learning process.

Suppose we have a data stream of triplets $\{(x_t^i, y_t^i, t)\}_i$, where x_t^i is the i -th input, y_t^i is the corresponding label and t is the task identifier ($t \in \mathcal{T} = \{1, \dots, T\}$). Each input-label pair is an identical and independent sample drawn from an unknown distribution $P_t(X, Y)$ of task t . We consider the number of tasks T is unknown and the tasks are coming sequentially as $t = 1, \dots, T$. We also assume that data among tasks are disjoint. At inference time, the task-id t is unknown at all time; also referred to as task-agnostic inference. Under this assumption, the resulting model $f(x; \Theta)$, parameterized by Θ , is optimized to minimize a predefined loss $l(x, y; \Theta)$ over new sequential input samples (x_t^i, y_t^i) from current data stream t . And at the same time, the performance on previous tasks should not decrease.

3.3. ACAE-REMIND for compressed feature replay

The ACAE-REMIND model is designed for online task-agnostic continual learning in a memory efficient way. As explained before, we aim to execute feature replay on an intermediate layer. Since all layers after the feature replay can be jointly trained, this strategy can potentially lead to improved performance. Because the distribution in lower layers is more complex, we propose an improved compression mechanism based on an ACAE module. The whole training procedure can be roughly divided into: 1) initialization stage (in Fig. 2) of the classification model, ACAE and PQ modules, 2) online continual learning stage (in Fig. 3).

3.3.1 Initialization

During initialization, the classification model, ACAE and PQ are trained sequentially with data (x_1^i, y_1^i) from the first task $t = 1$. In the first step, the whole classification model is optimized in an offline way. This step aims to learn a robust pretrained model for future tasks (similar as in REMIND). The parameters Θ are updated by minimizing the cross-entropy loss:

$$\text{minimize}_{\Theta} \mathcal{L}_{CE}(y_1^i, \hat{y}_1^i) = -y_1^i \cdot \log \hat{y}_1^i \quad (1)$$

where the prediction is given by $\hat{y}_1^i = f(x_1^i; \Theta)$.

Secondly, the ACAE module is inserted into the layer where we will replay the features. We denote the layers before and after the replay layer as $g(x; \Theta_1)$ and $h(z; \Theta_2)$, where Θ is the union of Θ_1 plus Θ_2 and $z = g(x; \Theta_1)$ in step 1. The encoder and decoder of ACAE are denoted as $u = D_{enc}(z; \Gamma)$ and $\hat{z} = D_{dec}(u; \Pi)$. The ACAE is trained with an auxiliary classification loss and auto-

encoder reconstruction MSE(mean square error) loss on the same data stream (x_1^i, y_1^i) .

Then parameters Γ, Π are computed by minimizing the ACAE loss:

$$\text{minimize}_{\Gamma, \Pi} \mathcal{L}_{ACAE}(y_1^i, x_1^i) = \mathcal{L}_{CE}(y_1^i, \hat{y}_1^i) + \|z_1^i - \hat{z}_1^i\|_2 \quad (2)$$

where

$$\begin{aligned} z_1^i &= g(x_1^i; \Theta_1) \\ \hat{z}_1^i &= D_{dec}(D_{enc}(z_1^i; \Gamma), \Pi) \\ \hat{y}_1^i &= h(\hat{z}_1^i; \Theta_2). \end{aligned} \quad (3)$$

After that, the last step is to train a PQ encoder-decoder pair $P_{enc}(u; \Upsilon)$ and $P_{dec}(v; \Psi)$ to approximate latent representations extracted from ACAE’s encoder. Here Υ, Ψ are learnt from the object function of PQ MSE loss:

$$\begin{aligned} \text{minimize}_{\Upsilon, \Psi} \mathcal{L}_{PQ}(x_1^i) &= \|u_1^i - \hat{u}_1^i\|_2 \\ z_1^i &= g(x_1^i; \Theta_1) \\ u_1^i &= D_{enc}(z_1^i; \Gamma) \\ \hat{u}_1^i &= P_{dec}(P_{enc}(u_1^i; \Upsilon), \Psi) \end{aligned} \quad (4)$$

3.3.2 Online continual learning

In online continual learning, only layers after the ACAE decoder can freely adjust to new tasks (the parameters in Θ_2), other modules (including parameters $\Theta_1, \Gamma, \Pi, \Upsilon, \Psi$) are all fixed during training. The new coming images from the data stream are passed through lower layers, the ACAE encoder and the PQ encoder to get their latent representations v_t^i with corresponding integer indexes. This is computed as:

$$v_t^i = P_{enc}(D_{enc}(g(x_t^i; \Theta_1); \Gamma), \Upsilon) \quad (5)$$

Then its representation is mixed with randomly selected \mathcal{N} previous samples $v_{\bar{t}}^j$ ($\bar{t} < t$) from the reservoir to reconstruct features via the PQ decoder and the ACAE decoder. Those features will be taken to optimize the trainable parameters Θ_2 with the cross-entropy loss $\mathcal{L}_{CE}(y_{\bar{t}}^j, \hat{y}_{\bar{t}}^j)$ and $\hat{y}_{\bar{t}}^j$ is formed as:

$$\hat{y}_{\bar{t}}^j = h(D_{dec}(P_{dec}(v_{\bar{t}}^j; \Psi), \Pi), \Theta_2), \bar{t} \leq t \quad (6)$$

Reservoir sampling After optimization, the new representation will be stored in the reservoir memory. If the reservoir is full, we randomly select a sample to pop up from one of the classes with most samples in the reservoir.

4. Experiments

4.1. Experimental setup

Datasets. Our evaluations are performed on three datasets: ImageNet-Subset [7], CIFAR100 [19] both with 100 classes² and CIFAR10 with 10 classes. We use data augmentation during the initial training of the full model and the ACAE, but removed it to train PQ (we save the representation of the original image - without augmentation -) and during the online continual learning stage. For feature augmentation, we only randomly resize and crop reconstructed features in the online continual learning stage.

Implementation details. We use Resnet-18 as our classification network for ImageNet-Subset. For CIFAR10 and CIFAR100, we use adapted Resnet-18 and Resnet-32 respectively (using only 3 blocks instead of the original 4 blocks). During initialization, the backbone network is learned from scratch with SGD, then the ACAE is trained with Adam [16]. For PQ training, we use the implementation from the Facebook Faiss library [15]. During the online continual learning stage we use SGD.

Evaluation metrics. We consider two widely used metrics: Average of top-1 accuracy over classes (AOC) up to the current task and top-1 accuracy after the last task (LAST).

Experimental settings. We will evaluate our method in five different settings. For the first three settings on ImageNet-Subset and CIFAR100, we use half of the classes as the first task and split the remaining into 5, 25 and 50 tasks with equal split (this setting is widely used [8, 11, 28]). We also refer to these as the 5, 25 and 50 steps setting. The fourth setting is splitting ImageNet-Subset into 10 tasks of the same size (this setting is used in [29, 30, 35]). We compare with several methods: iCaRL [30], BiC [35], UCIR [11], PODNet [8], GDumb [28], RPSnet [29] and REMIND [10]. We note that, except REMIND, the other methods are mainly designed for offline continual learning, which is a simpler setting compared with our online setting.

The fifth setting is on CIFAR10, where we use the commonly used setting from GMED [14], which divides CIFAR10 into 5 tasks equally. And we compare with online continual learning methods: AGEM [5], BGD [36], GEM [22], GSS-Greedy [3], HAL [4], ER [31], MIR [1], and GMED [14].

4.2. Results of online continual learning

Few-task evaluation (5 steps setting). We report the AOC metric on ImageNet-Subset with the 5 steps setting in Table 2. Every time we have a new input image, we randomly sample $\mathcal{N} = 50$ previous latent representations from the

²Samples are presented in a random yet fixed presentation order, as proposed in iCaRL [30], and adopted by others [8, 11, 28, 35].

Table 2: Comparison on Imagenet-Subset, we show the averages over classes (AOC) with 50 classes as the first task and 5/25/50 steps each with 10/2/1 classes. For REMIND and our method, we show the replay layer (block number) in brackets. The highest numbers in each row are highlighted.

On-line	Exemplar info			Methods	AOC over various steps		
	Num.	Shape (CHW)	Mem. (MB)		5	25	50
✗	2K	324 224 (int)	301	iCaRL	65.56	54.56	54.97
				BiC	68.97	59.65	46.49
				UCIR(NME)	69.07	60.81	55.44
				UCIR(CNN)	71.04	62.94	57.25
				PODNet(CNN)	75.54	68.31	62.08
				GDumb	-	-	62.86
✓	130K	87 7 (int)	51	REMIND(3.0)	70.58	67.93	67.35
				REMIND(4.0)	71.02	70.50	70.14
				Ours(1.0)	60.70	56.37	56.10
				Ours(2.0)	70.24	67.65	66.23
				Ours(3.0)	72.58	71.43	70.69
✓	130K	327 7 (int)	204	REMIND(3.0)	72.46	-	-
				REMIND(4.0)	73.98	-	-
				Ours(1.0)	74.08	-	-
				$-\mathcal{L}_{CE}$	70.26	-	-
				Ours(2.0)	73.75	-	-
				Ours(3.0)	73.63	-	-

sample reservoir. It can be seen that our method outperforms REMIND and that, especially for larger memory, we can obtain excellent results by replaying lower layers. For comparison, we have also computed results for block 3.0 for standard REMIND (going to lower blocks did further reduce performance). We observe that our method with 32 codebooks is only 1.46% lower than the state-of-the-art offline PODNet method, and it well outperforms other methods. Even when we have only 8 codebooks, it is still better than all offline algorithms except PODNet. Another interesting phenomenon is that, with 32 codebooks, we get an increase from block 3.0 to block 1.0, but this trend gets reversed with only 8 codebooks. The reason is that in the lower layers, the latent representations contain more information and thus require more codebooks to be represented.

For CIFAR100 with 5 steps, the performance is shown in Table 3. Due to smaller image-size, the compression ratio is not as high as in ImageNet-Subset. In this case, offline continual learning (PODNet) outperforms the online settings by a larger margin (6.1%) under the same memory allocation. It should be noted that PODNet runs for 160 epochs over the data whereas the online methods can only do one epoch. Also, among the online methods, our method performs worse than REMIND(3.0) when considering a memory of 6.4MB. This is because a first task with many classes and more data allows REMIND to also learn a high-quality backbone network. For the larger memory setting (12.8MB) our method performs comparable to REMIND(3.0).

Many-task evaluation (25/50 steps setting). Here the number of rehearsed samples is set to $\mathcal{N} = 200$ because

Table 3: Comparison on CIFAR100 dataset, we show the averages over classes (AOC) with 50 classes as the first task and 5/25/50 steps each with 10/2/1 classes. For REMIND and our method, we show the replay layer (block number) in brackets. The highest numbers in each row are highlighted.

On-line	Exemplar info			Methods	AOC over various steps		
	Num.	Shape (CHW)	Mem. (MB)		5	25	50
✗	2000	332 32 (int)	6.14	iCaRL	58.08	50.60	44.20
				BiC	56.86	48.96	47.09
				UCIR (NME)	63.63	56.82	48.57
				UCIR (CNN)	64.01	57.57	49.30
				PODNet (CNN)	64.83	60.72	57.98
				PODNet (NME)	64.48	62.72	61.40
				GDumb	-	-	58.40
✓	25000	48 8 (int)	6.40	REMIND(2.0)	55.79	58.25	57.93
				REMIND(3.0)	58.71	59.26	58.99
				Ours(1.0)	53.38	53.74	54.25
				Ours(2.0)	57.57	59.28	59.50
✓	50000	48 8 (int)	12.8	REMIND(2.0)	58.51	59.94	59.87
				REMIND(3.0)	61.23	61.02	61.00
				Ours(1.0)	56.40	56.98	57.01
				Ours(2.0)	61.27	62.49	62.30
				$-\mathcal{L}_{CE}$	56.11	60.19	60.07

there are less samples in each step. The results for 25 steps on ImageNet-Subset are shown in Table 2. We obtain state-of-the-art with 3.12%, 0.93% and 4.50% higher than PODNet, REMIND (block 4.0) and REMIND (block 3.0) respectively. The hardest setting is the 50 steps split, where only 1 class is viewed at every time step. We show our performance in Table 2. It is 8.49% higher than GDumb in the offline setting and 1.21% better than REMIND in the online setting.

For the many-task evaluation of CIFAR100 shown in Table 3, we got marginally better than PODNet in 50 steps and worse in 25 steps under higher memory allocation. With lower memory allocation we still got competitive performances. In conclusion, from the many-task evaluation, we observe that bias correction methods suffer from a drop in performance with more time steps, while our model obtains better results and without much drop in performance when the number of tasks increases.

Equal split with 10 tasks (9 steps) on ImageNet-Subset.

In the equal split setting, the 100 classes are divided into 10 tasks with 10 classes each. The top-5 accuracies in each time step are shown in Table 4 (For comparison, top-5 accuracy is adopted here since it is commonly used in this case). For this more challenging setting, REMIND obtains 8.5% lower results than ours, which is due to the less flexible backbone model. Especially here, it makes more sense to perform replay on a lower layer. Also note that in this setting, we get competitive results compared with the methods BiC [35] and RPSnet [29]. However, these methods are offline and perform multiple loops over the data.

Table 4: Comparison on ImageNet-Subset, we show top-5 accuracy with 10 classes as the first task and 9 steps each with 10 classes. For REMIND and our method, we show the replay layer (block number) in brackets. The highest numbers in offline and online settings are highlighted.

Methods		iCaRL	RPSnet	BiC	REMIND (4.0)	Ours (3.0)	Ours (2.0)	Ours (1.0)
Online		✗			✓			
Exem. Info	Num.	2000 (20*100)			130000 (1300*100)			
	Shape (CHW)	3*224*224 (integer)			32*7*7 (integer)			
	Mem. (MB)	301.06			203.84			
Acc.	1	99.3	100	98.4	98.4	98.4	98.4	98.4
	2	97.2	97.4	96.2	91.6	93.3	93.5	94.1
	3	93.5	94.3	94.0	87.1	90.5	91.1	92.7
	4	91.0	92.7	92.9	82.2	87.2	87.7	90.2
	5	87.5	89.4	91.1	79.7	85.3	85.5	89.2
	6	82.1	86.6	89.4	77.7	84.0	85.0	87.8
	7	77.1	83.9	88.1	74.8	81.0	83.7	85.7
	8	72.8	82.4	86.5	72.8	80.9	82.7	85.4
	9	67.1	79.4	85.4	72.2	80.8	83.4	84.7
	10	63.5	74.1	84.4	70.9	79.6	81.8	83.9
	AOC	83.1	88.0	90.6	80.7	86.1	87.1	89.2

Table 5: Comparison on CIFAR10 dataset, we show the LAST accuracy with 2 classes as the first task and 4 steps each with 2 classes. For REMIND and our method, we show the replay layer (block number) in brackets. The highest numbers in each row are highlighted.

Online	Exemplar info			Methods	LAST
	Num.	Shape (C*H*W)	Mem. (MB)		4 steps
✓	500	3*32*32 (integer)	1.536	Finetuning AGEM BGD GEM GSS-Greedy HAL ER MIR GMED(ER) GMED(MIR) GDumb	18.5 18.5 18.2 20.1 28.0 32.1 33.3 34.5 35.0 35.5 45.8
✓	24000	1*8*8 (integer)	1.536	REMIND(3.0) Ours(2.0)	45.2 48.4

Equal split with 5 tasks (4 steps) on CIFAR10. Several existing online methods cannot be straightforwardly applied to large datasets. To be able to compare to them, we also include results on CIFAR10 in Table 5. We divide the 10 classes into 5 tasks with 2 classes each. On this setting, we got 48.4%, which is 2.6% higher than the best reported results of GDumb. We also outperform the REMIND method with more than a 3% margin.

Table 6: Ablation study of classification loss on CIFAR100 and ImageNet-Subset. The features are replayed from block 2.0 for CIFAR100 and block 1.0/2.0/3.0 for ImageNet-Subset.

method name	CIFAR100	ImageNet-Subset		
	block 2.0	block 3.0	block 2.0	block 1.0
Uncompressed features	76.00%	80.96%	80.96%	80.96%
ACAe replay(w/o \mathcal{L}_{CE})	73.31%	78.80%	77.64%	76.52%
ACAe replay(w/ \mathcal{L}_{CE})	74.45%	79.40%	79.76%	80.44%

4.3. Ablation study

One of the key ingredients of the ACAE-REMIND method is the auxiliary classification loss that is used during the training of the auto-encoder (see Eq. 2). This loss ensures that the compression does not remove the features that are crucial for classification. Here we ablate this factor. To show the impact of the classification loss, we evaluate the classification accuracy after compression with and without the loss (directly after Step 2), and compare this to the results that would be obtained with the uncompressed features (see Table 6). The results show that classification loss mitigates the classification drop that occurs due to compression. Finally, we have also ablated the loss in Table 2 and Table 3 on our best performing setting (indicated by rows with $-\mathcal{L}_{CE}$). The results show that the loss does greatly improve results resulting in a performance gain of 2-5%.

To better evaluate the influence of the block number n (the block where we introduce the ACAE-REMIND as seen in Fig. 3) we have included Fig. 4. Here we show a comparison on ImageNet-Subset and CIFAR-100 under the 5, 25 and 50 steps settings. We can conclude that under these settings, the performances are increasing with the feature replay from the first block to the penultimate block, and then decreasing when replaying on the last block. As can be seen the optimal block is relatively stable with respect to the number of steps while keeping the same amount of data for the first task. If we however reduce the number of data for the first task, it becomes more difficult to learn a good backbone network and, as expected, the optimal n decreases: in Fig. 4(c) we can see that $n = 1$ yields optimal results.

5. Conclusions

In this paper, we proposed an extension to the REMIND method, called ACAE-REMIND. We propose a stronger compression module based on an auxiliary classifier auto-encoder that allows to move the feature replay to lower layers. The method is memory efficient and obtains better performance. In evaluation, we perform a comparison over multiple metrics among competitive methods. The strength of our model lies in the fact that with high compression ra-

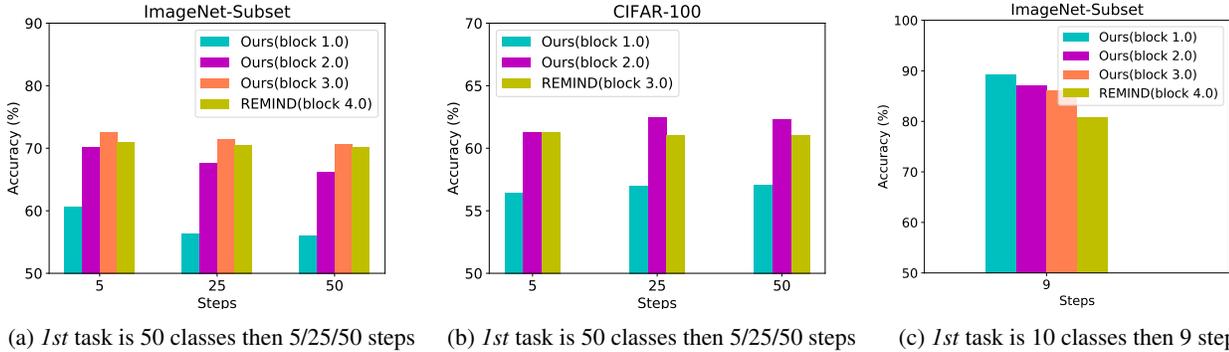


Figure 4: Ablation study on the block number n on ImageNet-Subset and CIFAR100 with various settings. The backbone for ImageNet-Subset is a 4-block Resnet-18 and for CIFAR100 is a 3-block Resnet-32. We show top-1 accuracy in (a) and (b), and top-5 accuracy (c).

tion, we could save more feature exemplars than image exemplars. Especially, when the first task is relatively small (the 10-task scenario in ImageNet-Subset and 5-task in CIFAR10) we outperform REMIND with a large margin. As future work, we are interested in extending this framework to other continual learning problems.

Acknowledgements

We acknowledge the support from Huawei Kirin Solution, the Spanish Government funding for projects PID2019-104174GB-I00 and RTI2018-102285-A-I00, and Kai acknowledges the Chinese Scholarship Council (CSC) No.201706170035. Luis acknowledges the Ramón y Cajal fellowship RYC2019-027020-I.

References

- [1] Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. In *Advances in Neural Information Processing Systems*, pages 11849–11860, 2019. 1, 5
- [2] Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with a network of experts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3366–3375, 2017. 2
- [3] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *Advances in Neural Information Processing Systems*, pages 11816–11825, 2019. 5
- [4] Arslan Chaudhry, Albert Gordo, Puneet K Dokania, Philip Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. *arXiv preprint arXiv:2002.08165*, 2020. 5
- [5] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a gem. *arXiv preprint arXiv:1812.00420*, 2018. 1, 2, 5
- [6] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *arXiv preprint arXiv:1909.08383*, 2019. 1, 2
- [7] 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. 3, 5
- [8] Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Proceedings of the IEEE European Conference on Computer Vision (ECCV)*, volume 2, page 6. Springer, 2020. 1, 5
- [9] Sebastian Farquhar and Yarin Gal. Towards robust evaluations of continual learning. *arXiv preprint arXiv:1805.09733*, 2018. 2
- [10] Tyler L Hayes, Kushal Kafle, Robik Shrestha, Manoj Acharya, and Christopher Kanan. Remind your neural network to prevent catastrophic forgetting. *arXiv preprint arXiv:1910.02509*, 2019. 1, 2, 3, 5
- [11] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via re-balancing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 831–839, 2019. 1, 2, 5
- [12] Ahmet Iscen, Jeffrey Zhang, Svetlana Lazebnik, and Cordelia Schmid. Memory-efficient incremental learning through feature adaptation. *arXiv preprint arXiv:2004.00713*, 2020. 2
- [13] Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2010. 2, 3
- [14] Xisen Jin, Junyi Du, and Xiang Ren. Gradient based memory editing for task-free continual learning. *arXiv preprint arXiv:2006.15294*, 2020. 5

- [15] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*, 2017. 5
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [17] 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. 1, 2
- [18] Mark A Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AICHE journal*, 37(2):233–243, 1991. 2
- [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [20] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. 2
- [21] Xialei Liu, Chenshen Wu, Mikel Menta, Luis Herranz, Bogdan Raducanu, Andrew D Bagdanov, Shangling Jui, and Joost van de Weijer. Generative feature replay for class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 226–227, 2020. 1, 2, 3
- [22] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in neural information processing systems*, 2017. 1, 5
- [23] Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 67–82, 2018. 2
- [24] 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. 2
- [25] Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost van de Weijer. Class-incremental learning: survey and performance evaluation. *arXiv preprint arXiv:2010.15277*, 2020. 2
- [26] Marc Masana, Tinne Tuytelaars, and Joost van de Weijer. Ternary feature masks: continual learning without any forgetting. *arXiv preprint:2001.08714*, 2020. 1, 2
- [27] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. 1
- [28] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *European Conference on Computer Vision*, pages 524–540. Springer, 2020. 2, 5
- [29] Jathushan Rajasegaran, Munawar Hayat, Salman H Khan, Fahad Shahbaz Khan, and Ling Shao. Random path selection for continual learning. In *Advances in Neural Information Processing Systems*, pages 12669–12679, 2019. 5, 6
- [30] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017. 1, 2, 5
- [31] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In *Advances in Neural Information Processing Systems*, pages 350–360, 2019. 5
- [32] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. *arXiv preprint arXiv:1801.01423*, 2018. 2
- [33] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In *Advances in Neural Information Processing Systems*, 2017. 1, 2, 3
- [34] Chenshen Wu, Luis Herranz, Xialei Liu, Joost van de Weijer, Bogdan Raducanu, et al. Memory replay gans: Learning to generate new categories without forgetting. In *Advances in Neural Information Processing Systems*, pages 5962–5972, 2018. 1, 2, 3
- [35] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 374–382, 2019. 1, 2, 5, 6
- [36] Chen Zeno, Itay Golan, Elad Hoffer, and Daniel Soudry. Task agnostic continual learning using online variational bayes. *arXiv preprint arXiv:1803.10123*, 2018. 5