Saliency Guided Experience Packing for Replay in Continual Learning

Gobinda Saha, Kaushik Roy Elmore Family School of Electrical and Computer Engineering Purdue University, West Lafayette, Indiana, USA

gsaha@purdue.edu, kaushik@purdue.edu

Abstract

Artificial learning systems aspire to mimic human intelligence by continually learning from a stream of tasks without forgetting past knowledge. One way to enable such learning is to store past experiences in the form of input examples in episodic memory and replay them when learning new tasks. However, performance of such method suffers as the size of the memory becomes smaller. In this paper, we propose a new approach for experience replay, where we select the past experiences by looking at the saliency maps which provide visual explanations for the model's decision. Guided by these saliency maps, we pack the memory with only the parts or patches of the input images important for the model's prediction. While learning a new task, we replay these memory patches with appropriate zero-padding to remind the model about its past decisions. We evaluate our algorithm on CIFAR-100, miniImageNet and CUB datasets and report better performance than the state-ofthe-art approaches. With qualitative and quantitative analyses we show that our method captures richer summaries of past experiences without any memory increase, and hence performs well with small episodic memory. Codes: https://github.com/sahagobinda/EPR.

1. Introduction

Recent success in deep learning primarily relies on training powerful models with fixed datasets in stationary environments. However, in the non-stationary setting, where data distribution changes over time, artificial neural networks (ANNs) fail to match the efficiency of human learning. In this setup, humans can learn incrementally leveraging and maintaining past knowledge, whereas ANN training algorithms [17] overwrite the representations of the past tasks upon exposure to a new task. This leads to rapid performance degradation on the past tasks - a phenomenon known as 'Catastrophic Forgetting' [29, 32]. Continual Learning (CL) [35, 48] aims to mitigate forgetting while sequentially updating the model on a stream of tasks.

To overcome catastrophic forgetting, an active line of re-

search in continual learning stores a few training samples from the past tasks as experiences in an episodic memory. Variants of experience replay [9, 7, 34, 5] have been proposed, where the model is jointly optimized on the samples from both episodic memory and new task. Such methods provide simple yet effective solutions to the catastrophic forgetting especially in online CL [27] setting where models need to learn from a single pass over the online data stream. However, performance of these methods strongly depends on the size of the episodic memory. The authors in [22] argued that for an optimal performance one needs to store all the past examples in the memory. While experience replay with large memory would yield higher performance, this would essentially mimic the joint optimization process in independent and identically distributed (IID) data setting which puts the effectiveness of CL algorithms into question [31]. Therefore, recent works [7, 9] have explored the idea of designing effective experience replay with tiny episodic memory. However, these methods suffer from high forgetting mainly due to overfitting [49] to the small memory samples, thus show suboptimal performance.

In this paper, we propose a continual learning algorithm that trains a fixed-capacity model on an online stream of tasks using a small episodic memory. Our method, referred to as Experience Packing and Replay (EPR), packs the memory with a more informative summary of the past experiences which improves performance of memory replay by reducing overfitting. To this end, we leverage the tools developed in the field of explainable artificial intelligence (XAI) [45, 57, 1] that shed light on the internal reasoning process of the ANNs. Among various explainability techniques, saliency methods [41, 55] highlight the part of the input data (image) that the model thinks is important for its final decision. Such analyses reveal that ANNs tend to make predictions based on some localized features or objects belonging to a part of the image whereas the rest of the image appears as background information. Thus we hypothesize that storing and replaying only these important parts of the images would be effective in reminding the networks about the past tasks and hence would reduce forgetting.

Therefore, in **EPR**, after learning each task, instead of storing full images, we identify important patches from different images belonging to each class with saliency method [41] and store them in the episodic memory. We introduce Experience Packing Factor (EPF) to set the number of patches kept per class and to determine the size of these patches. Thus, with these patches, we create *composite* images (for each class) that have higher diversity and capture richer summaries of past data distributions without increasing the memory size (Figure 1). While learning a new task, we retrieve these patches from the memory, zero-pad them to match with the original image sizes, and use them for experience replay. We evaluate our algorithm in standard and directly comparable settings [9, 7, 11, 43] on image classification datasets including CIFAR-100, miniImageNet, and CUB. We compare EPR with the state-of-the-art (SOTA) methods for varying memory sizes and report better accuracy with least amount of forgetting. We briefly summarize the contributions of our work as follows:

- We propose a new experience replay method for continual learning EPR, where we select (and store) episodic memory by identifying the parts of the past inputs important for model's prediction from the saliency maps. These parts are then replayed with zero-padding with new data to mitigate catastrophic forgetting.
- We compare EPR with the SOTA methods in both online task-incremental and class-incremental learning [27] and report significantly better performance.
- With comprehensive analyses, we show the effectiveness of saliency methods in selecting the 'informative' memory patches which enables EPR to perform well even with tiny episodic memories.

2. Related Works

Methods for continual learning can be broadly divided into three categories [12]. **Regularization based methods** penalize changes in important parameters for the past tasks to prevent forgetting. Elastic Weight Consolidation (**EWC**) [21] computes such importance from the Fisher information matrix. Other works use selective synaptic plasticity [54, 2], knowledge distillation [24], attention map distillation [13], and variational inference framework [30] for model regularization in continual learning. However, these methods suffer under longer task sequences and perform poorly [7] in the online CL setup considered in this paper.

Parameter isolation methods allocate different subsets of network parameters for each task to overcome forgetting. Some methods [38, 53] under this category expand the network for accommodating new tasks, whereas in other methods [28, 42, 39] a task-specific sub-network is selected by masking out the parameters. Unlike these methods, we train a fixed-sized network in online CL setting.

Memory based methods mitigate forgetting by either storing a subset of old examples in the episodic memory for rehearsal [36, 33], or storing important gradient spaces from the past tasks for constrained optimization [16, 40], or synthesizing old data from generative models for pseudo-rehearsal [44]. Experience Replay (ER) [36, 9] iointly trains the model on the samples from the new tasks and episodic memory. Several recent methods expand on this idea: Meta-Experience Replay (MER) [34] combines episodic memory with meta-learning to maximize knowledge transfer and minimize forgetting; GSS [4] stores examples in the memory for rehearsal based on the gradients; Maximal Interfered Retrieval (MIR) [3] selects a minibatch from the episodic memory for replay that incurs maximum change in loss; Adversarial Shapley value Experience Replay (ASER) [43] uses shapley value scores for episodic memory selection and retrieval; Hindsight Anchor Learning (HAL) [7] improves replay by adding an objective term to minimize forgetting on the meta-learned anchor datapoints; Dark Experience Replay (**DER++**) [5] improves ER by replaying network logits along with the ground truth labels of the memory samples. Gradient Episodic Memory (GEM) [25] and Averaged-GEM (A-GEM) [8] use samples from the memory to compute gradient constraint so that loss on the past task does not increase. Guo et al. [18] improved such methods by proposing a loss balancing update rules in MEGA. Ebrahimi et al. [15] in RRR stores full images with the corresponding saliency maps in the episodic memory and complements conventional experience replay with a saliency-based regularization objective so that model explanations for the past tasks have minimal drift. Our method, EPR, also uses episodic memory for experience replay. However, unlike these methods, we neither store full images nor store any saliency map in the episodic memory. Rather, leveraging the network's reasoning process, we store only a part (patch) of the image in the memory and use them to remind the model about past decisions in replay.

Online continual learning (OCL) [27, 10] builds on top of continual learning with additional set of realistic yet challenging desiderata. Here, model observes data one batch at a time and learns from a single pass over these online data stream. Here, model experiences new tasks (consisting of new classes) or data non-stationary (non-IID input distribution) over time. In this work, we focus on two widely explored sub-classes of OCL namely online task-incremental learning and online class-incremental learning. In the task-incremental setup, an extra supervisory signal (task-ID) is used to select task-specific head of classes over which inference is performed. Whereas in the class-incremental setup inference is performed without any task-ID.

3. Background and Notations

Continual Learning Protocol. For online task-incremental learning we follow the protocol used in [8].

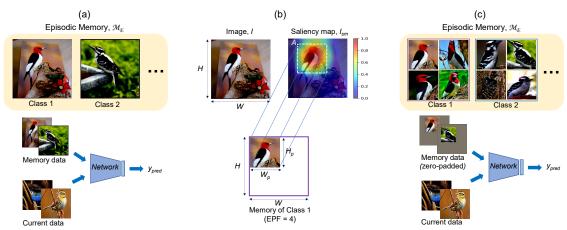


Figure 1. (a) Conventional experience replay method, where full images are stored in the episodic memory and network is jointly trained on these memory data and current task data. (b)-(c) Our experience packing and replay (**EPR**) method. (b) Experience packing, where only the important part of the image for network prediction is selected using saliency method and stored in the episodic memory with the corner coordinate, $A = (x_{cord}, y_{cord})$. (c) Experience packing increases sample diversity per class without memory increase. Stored memory patches are zero-padded and replayed with the current examples.

Here, a model learns from an ordered sequence of dataset, $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_T\}$ consisting of T tasks, where $\mathcal{D}_k = \{(I_i^k, t_i^k, y_i^k)_{i=1}^{n_k}\}$ is the dataset of the k-th task. Each example in these datasets consists of a triplet defined by an input $(I^k \in \mathcal{X})$, an integer task-ID $(t^k \in \mathcal{T})$ and a target vector $(y^k \in \mathbf{y}^k)$, where \mathbf{y}^k is the set of labels specific to task k and $\mathbf{y}^k \subset \mathcal{Y}$. As in [8, 9, 7], we use the first K(<T) Cross-Validation tasks to set the hyperparameters of the continual learning algorithms. The remaining T - K tasks, which is used for training and evaluation, is referred as Training-Evaluation tasks. In this setup, the goal is to learn a neural network, $f_\theta: \mathcal{X} \times \mathcal{T} \to \mathcal{Y}$, parameterized by $\theta \in \mathbb{R}^P$, that maps any input pair (I, t) to its target output y and maintains performance on all the prior tasks. For online class-incremental learning, task-ID (t^k) is omitted.

Saliency Map Generation. Saliency methods provide visual explanations for the model predictions in terms of relevant features in the input. For example, for an input RGB image, $I \in \mathbb{R}^{W \times H \times C}$ belonging to class c, these methods generate a *saliency map*, $I_{sm} \in \mathbb{R}^{W \times H}$ by assigning high intensity values to the relevant image regions that contribute to the model decision. The saliency map is generated by:

$$I_{sm} = XAI (f_{\theta}, I, c)$$
 (1)

where XAI is a saliency method. Simonyan *et al.* [45] first generated saliency maps using a pre-trained neural network. Later works improved the quality of saliency maps [56, 51] and reduced the cost for saliency computation [57, 41]. In this work, we primarily use Gradient-weighted Class Activation Mapping (**Grad-CAM**) [41] as the saliency method. It generates class-specific saliency maps based on gradients back-propagated to the later convolutional layers, given the model prediction. We describe the steps in detail in Appendix A. We also analyze the impact of other saliency

methods such as Grad-CAM++ [6], Smooth-Grad [46] and FullGrad [47] on EPR performance in Section 6.

4. Experience Packing and Replay (EPR)

Experience Replay. Continual learning algorithms, especially OCL methods, have achieved SOTA performance using experience replay [19, 5, 7]. These methods update the model, f_{θ} while storing a few samples from the training data (either randomly [9, 8, 34] or selectively [3, 4, 43]) into a replay buffer called episodic memory, \mathcal{M}_E . When data from a new task becomes available, the model is jointly trained on both the current and the episodic memory examples (Figure 1(a)). Thus, experience replay from \mathcal{M}_E mitigates catastrophic forgetting by reminding the network about the prior tasks. However, performance of these methods shows strong dependence on the number of samples kept in the memory. Though with larger \mathcal{M}_E replay yields better performance, designing effective experience replay with small episodic memory [9, 7] still remains an open research problem. This is because, the model performance becomes highly sensitive to the examples stored in a smallersized \mathcal{M}_E . Moreover, lack of sample diversity (per class) leads to higher overfitting to the memory examples, which causes loss of generalization for the past tasks leading to catastrophic forgetting [49]. To overcome these issues, we propose a method to select and store only patches of images, instead of full images, from the past tasks in \mathcal{M}_E . This enables us to pack diverse experiences from an image class without any memory increase. Next, we introduce the concept of Experience Packing Factor and describe how we select these patches and use them during replay.

Experience Packing Factor (EPF). Let's consider n_{sc} to be the number of (episodic) memory slots assigned for each class. Here, one memory slot can contain one full training image. For a given image, $I \in \mathbb{R}^{W \times H \times C}$ and the

target patch size (from this image) $I_p \in \mathbb{R}^{W_p \times H_p \times C}$ (with $W_p \leq W$ and $H_p \leq H$) Experience Packing Factor (EPF) is defined as the following ratio:

$$EPF = n_{sc} \frac{W \times H}{W_p \times H_p}.$$
 (2)

EPF is integer-valued and it refers to the number of patches one can fit into the given memory slot, n_{sc} , for any particular class. In our design, we consider square images (W=H) and patches $(W_p=H_p)$ and set EPF as a hyperparameter. Thus, for a given EPF we determine the image patch size as:

$$W_p = \sqrt{\frac{n_{sc}}{\text{EPF}}}W. \tag{3}$$

We take the floored integer value of W_p . Equation 3 tells us, for instance, to pack 4 patches (EPF= 4) into 1 memory slot $(n_{sc}=1)$, the patch width (height) should be half of the full image width (height).

Memory Patch Selection and Storage. Explainability techniques [57, 41, 55] reveal that ANN bases its decision on the class-discriminative localized features in the input data (image). Hence, we propose to store only the important part (patch) of that image and use it during replay to remind the network about its past decision. We identify these patches from the saliency maps (Section 3) of the full images. Therefore, while learning each task, we store only a few training images in a small, fixed-sized FIFO 'temporary ring buffer', \mathcal{M}_T [9]. At the end of each task, we extract the desired number of patches from these images using saliency maps and add them to the memory, \mathcal{M}_E . Note that, images from \mathcal{M}_T are only used for patch selection and not used in experience replay in the later tasks. Once the memory patches are stored in \mathcal{M}_E , the temporary buffer, \mathcal{M}_T is freed up to be reused in the next task. If we assume that data from k-th task is available to the model until it sees the next tasks (as in [15]), \mathcal{M}_T is not needed.

Let f_{θ}^k is the trained model after task k. For each examples in \mathcal{M}_T : $(I,t^k,c)\sim \mathcal{M}_T$ we generate the corresponding saliency map, I_{sm} using Equation 1. For given n_{sc} and chosen EPF, we obtain the (square) patch size $(W_p\times W_p)$. Then, we average-pool the saliency map, I_{sm} with kernel size $W_p\times W_p$ and stride (hyperparameter), S_{sm} . We store the top left coordinate (x_{cord},y_{cord}) of the kernel (patch) that corresponds to maximum average-pooled value. In other words, we identify a square region (of size $W_p\times W_p$) in the saliency map that has the maximum average intensity (Figure 1(b)). We obtain the memory patch from the image, I by :

$$I_p = I(x_{cord} : x_{cord} + W_p, y_{cord} : y_{cord} + W_p). \tag{4}$$

In our design, we keep a few more image samples in \mathcal{M}_T per class than the number of image patches we store in \mathcal{M}_E .

Algorithm 1 Experience Packing and Replay (EPR)

Inputs: T: No. of tasks; n_{sc} : No. of memory slots per class; EPF: No. of patches per class; W: image width (height); α : learning rate; \mathcal{D}^{train} : training dataset; f_{θ} : model with param θ ; n: mini-batch size **Output:** Updated model, f_{θ}

```
1: \mathcal{M}_T \leftarrow \{\}
                                                                                     2: \mathcal{M}_E \leftarrow \{\}
                                                                                               ⊳ episodic memory
  3: W_p \leftarrow f_{patch}(n_{sc}, \text{EPF}, W)
                                                                                        ⊳ patch size from Eq. 3
  4: for t \in \{1, 2, ...., T \text{ do} \}
             for \mathcal{B}_t \stackrel{\mathrm{n}}{\sim} \mathcal{D}_t^{train} do
  5:
                                                                                          ⊳ from current dataset
 6:
                    (\mathcal{B}_{\mathcal{M}_E}, \mathbf{x}_{cord}, \mathbf{y}_{cord}) \stackrel{\text{n}}{\sim} \mathcal{M}_E
                                                                                     ⊳ from episodic memory
 7:
                    \mathcal{B}_{\mathcal{M}_E} \leftarrow \text{Zero-pad}(\mathcal{B}_{\mathcal{M}_E}, \mathbf{x}_{cord}, \mathbf{y}_{cord})
 8:
                    \theta \leftarrow \operatorname{SGD}(\theta, \mathcal{B}_t \cup \mathcal{B}_{\mathcal{M}_E}, \alpha) \triangleright update model with replay
 9:
                    \mathcal{M}_T \leftarrow \text{Update-Ring-Buffer}(\mathcal{M}_T, \mathcal{B}_t)
10:
              end for
              \mathcal{M}_E \leftarrow \text{UPDATEMEMORY}(\mathcal{M}_E, \mathcal{M}_T, f_\theta, \text{EPF}, W_p)
11:
       Appendix C
12:
             \mathcal{M}_T \leftarrow \{\}

    ▷ clear 'temporary ring' buffer

13: end for
```

As we will be using these patches with zero-padding (discussed next) for replay, for storage in \mathcal{M}_E , we want to prioritize the patches that after zero-padding gives (or remain closer to) the correct class prediction. Thus we Zero-pad each image patch, I_p and check the model prediction. At first, we populate the memory with the patches for which model gives correct prediction. Then we fill up the remaining slots in \mathcal{M}_E by the patches for which correct class is in model's Top3 predictions. Any remaining memory slot is filled up from the remaining patches irrespective of model predictions. Each selected image patch is then added to $\mathcal{M}_E = \mathcal{M}_E \cup \{(I_p, t^k, c, x_{cord}, y_{cord})\}$, with task-ID, class label and localizable coordinates in the original image.

Replay with Memory Patches. Since the patches stored in \mathcal{M}_E are smaller in size than the original images, we Zero-pad these patches (Figure 1(c)) each time we use them for experience replay. While zero-padding we place these patches in the 'exact' position of their original images using the coordinate values (x_{cord}, y_{cord}). Each sample,

$$I_{rep} = \text{Zero-pad}(I_p, x_{cord}, y_{cord})$$
 (5)

for replay will thus have the same dimensions as the samples of the current task. Throughout the paper, we use zero-padding with the exact placement of the memory patches for replay unless otherwise stated. We discuss other choices for memory patch padding and placement in Section 6. The pseudo-code of our algorithm is given in Algorithm 1.

5. Experimental Setup

Here, we describe the task-incremental learning setup. In Section 6 we discuss the class-incremental learning setup.

Datasets. We evaluate our algorithm on three image classification benchmarks widely used in continual learning. **Split CIFAR** [25] consists of splitting the original CIFAR-100 dataset [23] into 20 disjoint subsets, each of which is considered as a separate task containing 5 classes. **Split miniImageNet** [9, 7, 14] is constructed by splitting

100 classes of miniImageNet [50] into 20 tasks where each task has 5 classes. **Split CUB** [8, 9, 18] is constructed by splitting 200 bird categories from CUB dataset [52] into 20 tasks where each task has 10 classes. The dataset statistics are given in Appendix B. We do not use any data augmentation. All datasets have 20 tasks (T=20), where first 3 tasks (K=3) are used for hyperparameter selection while the remaining 17 tasks are used for training. We report performances on the held-out test sets from these 17 tasks.

Network Architectures and Training. For CIFAR and miniImageNet, we use a reduced ResNet18 [7] with three times fewer feature maps across all layers. For CUB, we use a standard ImageNet pretrained ResNet18 [8, 9]. Similar to [8, 9, 7, 18], we train and evaluate our algorithm in 'multi-head' setting [20] where a task-ID is used to select a task-specific classifier. All the models are trained using SGD with batch size of n=10 for both the current and memory examples. All experiments are averaged over 5 runs using different random seeds, where each seed corresponds to a different model initialization and dataset ordering among tasks. A list of hyperparameters along with the EPFs used in these experiments is given in Appendix D.

Baselines. From memory based methods, we compare with A-GEM [8], MIR [3], MER [34], MEGA-I [18], DER++ [5], ASER [43], HAL [7] and experience replay [9] with ring (ER-RING) and reservoir (ER-Reservoir) buffer. We also compare with EWC [21] which uses regularization and RRR [15] that uses both memory replay and regularization. We include two non-continual learning baselines: Finetune and Multitask. Finetune, where a single model is trained continually without any memory or regularization, gives performance lower bound. Multitask is an oracle baseline where a model is trained jointly on all tasks.

Performance Metrics. We evaluate the classification performance using the **ACC** metric, which is the average test classification accuracy of all tasks. We report backward transfer, **BWT** to measure the influence of new learning on the past knowledge. For instance, negative BWT indicates forgetting. Formally, ACC and BWT are defined as:

$$ACC = \frac{1}{T} \sum_{i=1}^{T} R_{T,i}, \quad BWT = \frac{1}{T-1} \sum_{i=1}^{T-1} -R_i^T \qquad (6)$$

Here, T is the total number of sequential tasks, $R_{T,i}$ is the accuracy of the model on i^{th} task after learning the T^{th} task sequentially [25], and $R_i^T = \max_{l \in \{1, \dots, T-1\}} (R_{l,i} - R_{T,i})$ [7].

6. Results and Analyses

Task-incremental Learning Performance. First, we compare the performance (ACC and BWT) of EPR with the baselines. Table 1 summarizes the results, where for a given n_{sc} , episodic memory (of either ring or reservoir type) can store up to $|\mathcal{M}_E|$ examples. Here, n_{sc} is the number of

memory slots per class and the memory size, $|\mathcal{M}_E|$ is :

$$|\mathcal{M}_E| = n_{sc} \times \text{no. of classes per task} \times \text{no. of tasks.}$$
 (7)

Results in Table 1 show that performance of EWC is almost identical to the 'Finetune' baseline. This indicates that such method is ill-suited for online CL setup. When one memory slot is assigned per class $(n_{sc} = 1)$, our method (EPR) outperforms A-GEM and MEGA-I considerably for all the datasets. Moreover, compared to the other experience replay methods, such as MIR, MER, DER++, and ER, EPR achieves $\sim 2\%$ and $\sim 3\%$ accuracy improvement for CIFAR and miniImageNet respectively with least forgetting. In CUB, EPR obtains $\sim 5\%$ accuracy improvement over these baselines with only $\sim 2\%$ forgetting. For all datasets, EPR considerably outperforms ASER which shows our saliency based memory storage offers better solution for small-memory experience replay compared to the shapely value based memory selection in ASER. Moreover, for CUB dataset, EPR obtains $\sim 9\%$ better accuracy with $\sim 2\%$ less forgetting compared to RRR. This demonstrate the benefit of saliency based input selection (for replay) in EPR over the saliency map regularization in RRR.

Finally, we compare EPR with HAL [7] which holds the SOTA performance in this setup. For the miniImageNet tasks, EPR (with $n_{sc}=1$) achieves slightly better accuracy than HAL, whereas HAL outperforms EPR at the CIFAR tasks. However, in addition to the ring buffer, HAL uses extra memory to store anchor points having the same size of the full images for each class. Thus effectively, HAL uses two memory slots per class ($n_{sc}=2$). In Table 1, we compare EPR with HAL where EPR also uses two memory slots per class. Under this iso-episodic memory condition, EPR has better accuracy and lower forgetting than HAL for both datasets. In this case ($n_{sc}=2$), EPR outperforms all the other methods significantly. For all datasets, amount of forgetting in EPR reduces with increase in memory size.

Class-incremental Learning Performance. In Table 2(a), we compare EPR with SOTA baselines in class-incremental learning setup for varying episodic memory sizes. For CIFAR-100 (20 Tasks) and miniImageNet (10 Tasks) experiments, we used training setup the similar to Continual Prototype Evolution (CoPE) [11] and ASER [43] respectively. In this setup, 'single-head' inference is performed without task-ID and EPR outperforms all the baselines achieving up to 3% ACC gain. All the subsequent analyses are performed in task-incremental setup.

Padding and Placement of Memory Patches. Next, we analyze the impact of different types of padding and placement of the memory patches on the EPR performance. For padding we have two different choices: we can either Zero-pad these patches or we can pad these patches with pixels sampled from normal Gaussian distribution, which we refer to as Random-pad. Similarly, we

Table 1. Performance comparison in task-incremental learning setup. (*) indicates results for CIFAR and miniImageNet are reported from HAL [7] and results for CUB are reported from ER-RING [9]. (†) indicates results are reported from ER-RING. We (re) produced all the other results. Average and standard deviations are computed over 5 runs for different random seeds. No. of memory slots per class, n_{sc} ={1, 2} refers to memory size, $|\mathcal{M}_E|$ ={85, 170} for CIFAR and miniImageNet, and $|\mathcal{M}_E|$ ={170, 340} for CUB.

		Split	CIFAR	Split minil	Split miniImageNet		CUB
n_{sc}	Methods	ACC (%)	BWT	ACC (%)	BWT	ACC (%)	BWT
-	Finetune*	42.9 ± 2.07	-0.25 ± 0.03	34.7 ± 2.69	-0.26 ± 0.03	55.7 ± 2.22	-0.13 ± 0.03
	EWC* [21]	42.4 ± 3.02	- 0.26 ± 0.02	37.7 ± 3.29	- 0.21 ± 0.03	55.0 ± 2.34	- 0.14 ± 0.02
1	RRR [15]	-	-	-	-	62.9 ± 1.33	-0.04 ± 0.01
	A-GEM* [8]	54.9 ± 2.92	- 0.14 ± 0.03	48.2 ± 2.49	- 0.13 ± 0.02	62.1 ± 1.28	- 0.09 ± 0.01
	MIR* [3]	57.1 ± 1.81	- 0.12 ± 0.01	49.3 ± 2.15	- 0.12 ± 0.01	-	-
	MER* [34]	49.7 ± 2.97	- 0.19 ± 0.03	45.5 ± 1.49	- 0.15 ± 0.01	55.4 ± 1.03	- 0.10 ± 0.01
	MEGA-I [18]	55.2 ± 1.21	- 0.14 ± 0.02	48.6 ± 1.11	- 0.10 ± 0.01	65.1 ± 1.30	- 0.05 ± 0.01
	DER++ [5]	54.0 ± 1.18	- 0.15 ± 0.02	48.3 ± 1.44	-0.11 ± 0.01	66.8 ± 1.36	- 0.04 ± 0.01
	ASER [43]	55.4 ± 1.17	- 0.16 ± 0.01	48.2 ± 1.43	- 0.09 ± 0.01	66.2 ± 1.63	- 0.07 ± 0.02
	ER-Reservoir† [9]	53.1 ± 2.66	- 0.19 ± 0.02	44.4 ± 3.22	- 0.17 ± 0.02	61.7 ± 0.62	- 0.09 ± 0.01
	ER-RING* [9]	56.2 ± 1.93	- 0.13 ± 0.01	49.0 ± 2.61	- 0.12 ± 0.02	65.0 ± 0.96	- 0.03 ± 0.01
	EPR (Ours)	$\textbf{58.5} \pm \textbf{1.23}$	- 0.10 ± 0.01	$\textbf{51.9} \pm \textbf{1.57}$	- 0.06 \pm 0.01	$\textbf{72.1} \pm \textbf{0.93}$	- 0.02 ± 0.01
2	RRR	-	-	-	-	67.1 ± 1.27	- 0.03 ± 0.01
	MEGA-I	57.6 ± 0.87	- 0.12 ± 0.01	50.3 ± 1.14	- 0.08 ± 0.01	67.8 ± 1.30	- 0.04 ± 0.01
	DER++	56.3 ± 0.98	- 0.14 ± 0.01	50.1 ± 1.14	- 0.09 ± 0.01	70.7 ± 0.62	- 0.03 ± 0.01
	ASER	57.5 ± 1.21	- 0.13 ± 0.01	50.1 ± 1.07	- 0.08 ± 0.01	69.9 ± 0.85	- 0.05 ± 0.01
	ER-RING	58.6 ± 2.68	- 0.12 ± 0.01	51.2 ± 2.06	- 0.10 ± 0.01	68.3 ± 1.13	- 0.02 ± 0.01
	HAL* [7]	60.4 ± 0.54	- 0.10 ± 0.01	51.6 ± 2.02	- 0.10 ± 0.01	-	-
	EPR (Ours)	$\textbf{60.8} \pm \textbf{0.35}$	- 0.09 ± 0.01	$\textbf{53.2} \pm \textbf{1.45}$	- 0.05 ± 0.01	$\textbf{73.5} \pm \textbf{1.30}$	- 0.01 \pm 0.01
-	MultiTask*	68.3	-	63.5	-	65.6	-

Table 2. (a) ACC(%) comparison in class-incremental learning setup, where results for CIFAR-100 and miniImageNet baselines are taken from [11] and [43] respectively. (b) Impact of padding, placement and selection method of memory patches on EPR performance ($n_{sc} = 1$).

	CIFAR-10	0 (20 Tasks)	miniImageNet (10 Tasks)			
Methods	$ \mathcal{M}_E $ =1k	$ \mathcal{M}_E $ =2k	$ \mathcal{M}_E $ =1k	$ \mathcal{M}_E $ =2k		
GSS [4]	7.6 ± 1.81	9.9 ± 1.17	7.5 ± 0.50	10.7 ± 0.80		
MIR [3]	9.0 ± 1.20	12.0 ± 1.84	8.1 ± 0.30	11.2 ± 0.70		
ER [9]	7.9 ± 1.98	11.9 ± 3.42	8.7 ± 0.40	11.8 ± 0.90		
ASER [43]	-	-	12.2 ± 0.80	14.8 ± 1.10		
CoPE [11]	10.7 ± 1.13	14.8 ± 1.18	-	-		
EPR (Ours)	$\textbf{13.7} \pm \textbf{0.84}$	$\textbf{16.3} \pm \textbf{0.86}$	$\textbf{13.8} \pm \textbf{0.23}$	$\textbf{15.6} \pm \textbf{0.31}$		

	CIFAR		miniIma	miniImageNet		CUB	
Methods	ACC (%)	BWT	ACC (%)	BWT	ACC (%)	BWT	
EPR (Zero-pad, exact)	58.5	- 0.10	51.9	- 0.06	72.1	- 0.02	
EPR (Zero-pad, random)	57.0	- 0.11	51.5	- 0.06	71.9	- 0.02	
EPR (Random-pad, exact)	57.2	- 0.11	49.7	- 0.07	71.5	- 0.02	
Random Snip & Replay	53.6	- 0.14	49.5	- 0.08	67.4	- 0.05	
Input Compression	56.7	- 0.12	49.6	- 0.08	69.6	- 0.03	

can place these patches either in the exact position of their original image using stored coordinate values or we can place them at random positions. Table 2(b) shows that across all datasets, exact placement works slightly better than random placement. This indicate that neural network remembers the past tasks better if it finds the class-discriminative features in their original position during replay. For all the datasets, zero-padding performs better than random-padding (Table 2(b)), which indicates that removing the background information completely serves as a better reminder of past tasks for the network. Thus, in all our experiments we use zero-padding with exact placement. For this, we store a 2D coordinate value per memory patch which has an insignificant overhead compared to $|\mathcal{M}_E|$.

Effectiveness of Saliency Guided Memory Selection. A simple alternative to our saliency guided memory patch selection is to randomly select a patch (of size $W_p \times W_p$)

from the original image and use it for replay with zero-padding. We refer to this method as 'Random Snip & Replay' and compare it with EPR in Table 2(b). For CIFAR and CUB, EPR achieves $\sim 5\%$ and for miniImageNet EPR achieves $\sim 2.5\%$ better accuracy than this baseline. We investigate another simple 'Input Compression' baseline, where we down-sample the images to desired sizes (comparable to EPR) and store in \mathcal{M}_E and for replay up-sample them to the original sizes. Results in Table. 2(b) shows, EPR outperforms this method by up to $\sim 2.5\%$ implying that input compression leads to higher loss in information. These results show that saliency based memory patch selection plays a key role in enabling high performance in EPR.

Experience Replay with Tiny Episodic Memories. Next, we study the impact of buffer size, $|\mathcal{M}_E|$ on the performance of experience replay methods. In Table 1, we reported the results for $n_{sc} = \{1,2\}$ to provide a direct

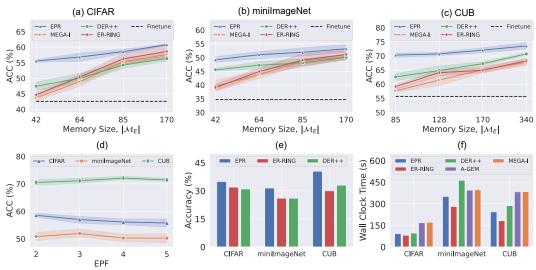


Figure 2. Comparison of ACC for varying memory sizes for (a) Split CIFAR, (a) Split miniImageNet, and (c) Split CUB dataset. (d) ACC for different Experience Packing Factors (EPFs) for different datasets in EPR (for $n_{sc}=1$). (e) Joint training accuracies on episodic memory data compare buffer informativeness (for $n_{sc}=1$). (f) Total wall-clock training time for learning all the tasks sequentially.

comparison to the SOTA works. Here, we analyze whether it is possible to reduce the memory size further and still have an effective experience replay. This means, we consider the fractional values for n_{sc} . In such cases, for instance, $n_{sc} = 0.5$ means only half of the seen classes will have a sample stored in \mathcal{M}_E . Understandably this is a challenging condition for standard experience replay as many classes will not have any representation in the memory, leading to a sharp drop in performance. However, in our method, we can set an appropriate EPF (≥ 1) for any given $n_{sc}(>0)$ and use Equation 3 to get the size of the memory patches. This allows us to pack representative memory patches from each class and preserve performance of experience replay. In Figure 2(a)-(c) we show how the performance (ACC) of different memory replay methods varies with the memory size for different datasets. Here we consider, $n_{sc} = \{0.5, 0.75, 1, 2\}$ which corresponds to $|\mathcal{M}_E| = \{42, 64, 85, 170\}$ for CIFAR and miniImageNet and $|\mathcal{M}_E| = \{85, 128, 170, 340\}$ for CUB. We also provide the results in tabular form in Appendix E (Table E.1).

The 'Finetune' baselines in these figures correspond to $n_{sc}=0$ case, and hence, serve as lower bounds to the performance. These figures show that ACC of memory replay methods such as ER-RING and MEGA-I falls sharply and approaches 'Finetune' baselines as we reduce the memory size. DER++, which uses both stored labels and logits during replay, performs slightly better than these methods. However, it still exhibits high accuracy drop (up to $\sim 9\%$) when memory size is reduced. In contrast, EPR shows high resilience under extreme memory reduction. For example, accuracy drop is only about $\sim 5\%$ for CIFAR, $\sim 4\%$ for miniImageNet, and $\sim 3\%$ for CUB dataset when memory size is reduced by a factor of 4. Thus, among the memory

replay methods for CL, EPR promises to be the best option, especially in the tiny episodic memory regimes.

EPF vs. Performance. In our design, EPF determines how many image patches we can store per class for a given n_{sc} . A higher EPF would select smaller patches (Equation 3), and hence, increase the sample quantity (or diversity) per class. However, a large number of memory patches, unlike full images, does not imply better performance from experience replay. In this regard, feature localization quality in the saliency map gives us a better picture about the quality of these patches for experience replay. Figure 3 shows the saliency maps of different classes for different datasets from our experiments. For larger-sized and better quality images of CUB dataset, we observe that Grad-CAM localizes the object better within small regions of the images (Figure 3(c)). This gives us an impression that important part of the image for network decision can be captured with a smaller patch. Thus a higher EPF can be chosen to select a larger number of high quality patches, which would improve the performance of experience replay. In contrast, for smaller-sized and low quality images of CI-FAR, we observe that network's decisions are distributed over a large portion of the images (Figure 3(a)). Thus, a smaller patch (for high EPF) here may not capture enough information to be fully effective in experience replay. Figure 2(d) shows the impact of EPF on the performance of our method for different datasets. Since, for $n_{sc} = 1$, our methods with EPF = 1 is similar to ER-RING, we consider the cases with EPF ≥ 2 . For CUB, accuracy of EPR improves as we increase EPF from 2 to 4. Beyond that point the accuracy drops which indicates that the memory patches are too small to capture all the relevant information in the given images. For CIFAR, we obtain the best performance for

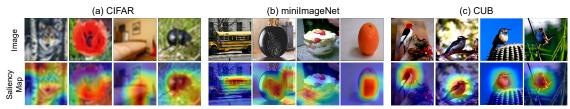


Figure 3. (a) Saliency maps of images from (a) Split CIFAR, (b) Split miniImageNet, (c) Split CUB dataset. CIFAR images have the lowest resolution whereas CUB images have the highest. With increasing image resolution we observe better object localization with Grad-CAM.

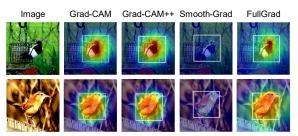


Figure 4. Patch selection from CUB dataset using various saliency methods in EPR.

 ${\sf EPF}=2$ and as we increase EPF we observe drop in accuracy. For miniImageNet, we obtain the best performance for ${\sf EPF}=3$. These results support our observations that link the size and quality of the memory patches to the quality of object localization in saliency maps for a given dataset.

Informativeness of Memory Buffer. Generalization capability of a model trained on the samples from the memory buffer, \mathcal{M}_E can reveal the informativeness of the buffer. Thus, following Buzzega et al. [5], we compare the informativeness of the EPR buffer with the buffers used in DER++ and ER-RING. For each dataset, we train the corresponding model jointly on all the buffer data from all the tasks. This training does not correspond to CL, rather it mimics the multitasks learning. For EPR buffer, we train the model with zero-padded memory patches. Figure 2(e) shows the average (multitask) accuracy on the test set. For all the datasets, models trained on EPR buffer achieve the highest accuracy (better generalization). Thus, our proposed experience packing method captures richer summaries of underlying data distribution (without any memory increase) compared to the other buffers. This reduces overfitting to the memory buffer, which in turn improves accuracy and reduces forgetting (Table 1, Appendix Table E.1).

Impact of Saliency Methods on EPR. To assess the sensitivity of EPR to the choice of saliency algorithms, we also used Grad-CAM++ [6], Smooth-Grad [46] and Full-Grad [47] as saliency method in EPR. Results are compared in Table 3, where we observe only a marginal variations ($\sim 1\%$) in performance. Corresponding saliency maps (Figure 4) show similar patches are being selected in EPR by these methods. These findings indicate our method can work with a wide variety of saliency methods.

Training Time Analysis. Finally, following SOTA works, in Figure 2(f) we provide training time analysis of

Table 3. Impact of various saliency methods on EPR performance.

	CIFA	R	miniIma	geNet	CUB	
Saliency Method	ACC (%)	BWT	ACC (%)	BWT	ACC (%)	BWT
Grad-CAM [41]	58.5	- 0.10	51.9	- 0.06	72.1	- 0.02
Grad-CAM++ [6]	58.0	- 0.10	51.6	- 0.07	72.4	- 0.02
Smooth-Grad [46]	57.5	- 0.12	51.1	- 0.07	72.0	- 0.02
FullGrad [47]	57.8	- 0.12	51.2	- 0.07	72.8	- 0.02

the algorithms where algorithm-specific compute overheads are captured in terms of extra time spent. We measured time on a single NVIDIA GeForce GTX 1060 GPU. Compared to the standard replay (ER-RING), EPR only takes up to $\sim 30\%$ extra time for training. This extra time is spent on saliency based patch selection and zero-padding. EPR trades-off this computational overhead by gaining $\sim 7\%$ ACC improvement over ER-RING. Compared to other recent works such as DER++ and MEGA-I, EPR trains faster and performs better. HAL did not report training time for the datasets under consideration, and hence, we could not provide a comparison. Since, HAL and MER both have meta-optimization steps (higher compute overhead), they are expected to require much larger training time [7].

7. Conclusions

In this paper, we propose a new experience replay method with small episodic memory for continual learning. Using saliency maps, our method identifies the parts of the input images that are important for model's prediction. We store these patches, instead of full images, in the memory and use them with appropriate zero-padding for replay. Our method thus packs the memory with diverse experiences, and hence captures the past data distribution better without memory increase. Comparison with the SOTA methods on diverse image classification tasks shows that our method is simple, fast, and achieves better accuracy with least amount of forgetting. We believe that the work opens up rich avenues for future research. Firstly, better understanding of the model's decision process and better feature localization with saliency methods would improve the quality of the memory patches and hence improve experience replay. Secondly, new replay techniques for the patches can be explored to reduce the memory overfitting further. Finally, future studies can explore possible applications of our concept in other domains such as in reinforcement learning [37].

Acknowledgements. This work was supported in part by the National Science Foundation, Vannevar Bush Faculty Fellowship, Army Research Office, MURI, and by Center for Brain-Inspired Computing (C-BRIC), one of six centers in JUMP, a SRC program sponsored by DARPA.

References

- Amina Adadi and Mohammed Berrada. Peeking inside the black-box: A survey on explainable artificial intelligence (XAI). *IEEE Access*, 6:52138–52160, 2018.
- [2] Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *The European Conference on Computer Vision (ECCV)*, pages 139–154, 2018.
- [3] 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*, volume 32. Curran Associates, Inc., 2019.
- [4] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019.
- [5] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. In *Advances in Neural Information Processing Systems*, volume 33, pages 15920–15930. Curran Associates, Inc., 2020.
- [6] Aditya Chattopadhay, Anirban Sarkar, Prantik Howlader, and Vineeth N Balasubramanian. Grad-CAM++: Generalized gradient-based visual explanations for deep convolutional networks. In 2018 IEEE Winter Conference on Applications of Computer Vision (WACV), pages 839–847, 2018.
- [7] Arslan Chaudhry, Albert Gordo, Puneet K Dokania, Philip HS Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In AAAI, 2021.
- [8] Arslan Chaudhry, Marc' Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with A-GEM. In *International Conference on Learning Representa*tions, 2019.
- [9] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K. Dokania, Philip H. S. Torr, and Marc'Aurelio Ranzato. Continual learning with tiny episodic memories. ArXiv, abs/1902.10486, 2019.
- [10] Aristotelis Chrysakis and Marie-Francine Moens. Online continual learning from imbalanced data. In *Proceedings* of the 37th International Conference on Machine Learning, volume 119 of Proceedings of Machine Learning Research, pages 1952–1961. PMLR, 13–18 Jul 2020.
- [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 (ICCV)*, pages 8250–8259, October 2021.
- [12] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars. A continual

- learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [13] Prithviraj Dhar, Rajat Vikram Singh, Kuan-Chuan Peng, Ziyan Wu, and Rama Chellappa. Learning without memorizing. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 5133–5141, 2019.
- [14] Sayna Ebrahimi, Franziska Meier, Roberto Calandra, Trevor Darrell, and Marcus Rohrbach. Adversarial continual learning. In *The European Conference on Computer Vision* (ECCV), 2020.
- [15] Sayna Ebrahimi, Suzanne Petryk, Akash Gokul, William Gan, Joseph E. Gonzalez, Marcus Rohrbach, and trevor darrell. Remembering for the right reasons: Explanations reduce catastrophic forgetting. In *International Conference on Learning Representations*, 2021.
- [16] Mehrdad Farajtabar, Navid Azizan, Alex Mott, and Ang Li. Orthogonal gradient descent for continual learning. In Silvia Chiappa and Roberto Calandra, editors, Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics, volume 108 of Proceedings of Machine Learning Research, pages 3762–3773. PMLR, 26–28 Aug 2020.
- [17] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- [18] Yunhui Guo, Mingrui Liu, Tianbao Yang, and Tajana Rosing. Improved schemes for episodic memory-based lifelong learning. Advances in Neural Information Processing Systems, 33, 2020.
- [19] Tyler L Hayes, Nathan D Cahill, and Christopher Kanan. Memory efficient experience replay for streaming learning. In *International Conference on Robotics and Automation* (ICRA). IEEE, 2019.
- [20] Yen-Chang Hsu, Yen-Cheng Liu, Anita Ramasamy, and Zsolt Kira. Re-evaluating continual learning scenarios: A categorization and case for strong baselines. In *NeurIPS Continual learning Workshop*, 2018.
- [21] James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114:3521 – 3526, 2017.
- [22] Jeremias Knoblauch, H. Husain, and Tom Diethe. Optimal continual learning has perfect memory and is np-hard. In *International Conference on Machine Learning (ICML)*, 2020.
- [23] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- [24] Zhizhong Li and Derek Hoiem. Learning without forgetting. IEEE Transactions on Pattern Analysis and Machine Intelligence, 40:2935–2947, 2018.
- [25] David Lopez-Paz and Marc' Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [26] Aravindh Mahendran and Andrea Vedaldi. Visualizing deep convolutional neural networks using natural pre-images. *Int. J. Comput. Vision*, 120(3):233–255, Dec. 2016.

- [27] Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022.
- [28] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 7765–7773, 2018.
- [29] Michael Mccloskey and Neil J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. *The Psychology of Learning and Motivation*, 24:104–169, 1989.
- [30] Cuong V. Nguyen, Yingzhen Li, Thang D. Bui, and Richard E. Turner. Variational continual learning. In *International Conference on Learning Representations*, 2018.
- [31] Ameya Prabhu, Philip Torr, and Puneet Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *The European Conference on Computer Vision* (ECCV), August 2020.
- [32] Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and forgetting functions. *Psychological review*, 97 2:285–308, 1990.
- [33] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. iCaRL: Incremental classifier and representation learning. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 5533-5542, 2017.
- [34] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. In *In International Conference on Learning Representations (ICLR)*, 2019.
- [35] Mark B. Ring. Child: A first step towards continual learning. In *Learning to Learn*, 1998.
- [36] Anthony V. Robins. Catastrophic forgetting, rehearsal and pseudorehearsal. *Connect. Sci.*, 7:123–146, 1995.
- [37] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [38] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *ArXiv*, abs/1606.04671, 2016.
- [39] Gobinda Saha, Isha Garg, Aayush Ankit, and Kaushik Roy. SPACE: Structured compression and sharing of representational space for continual learning. *IEEE Access*, pages 1–1, 2021.
- [40] Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. In *International Con*ference on Learning Representations, 2021.
- [41] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE*

- International Conference on Computer Vision (ICCV), Oct 2017.
- [42] Joan Serrà, Dídac Surís, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings* of *Machine Learning Research*, pages 4548–4557. PMLR, 10–15 Jul 2018.
- [43] Dongsub Shim, Zheda Mai, Jihwan Jeong, Scott Sanner, Hyunwoo Kim, and Jongseong Jang. Online class-incremental continual learning with adversarial shapley value. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 9630–9638, 2021.
- [44] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. In Advances in Neural Information Processing Systems, volume 30, 2017.
- [45] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *In Workshop at International Conference on Learning Representations*, 2014.
- [46] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. SmoothGrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825, 2017.
- [47] Suraj Srinivas and François Fleuret. Full-gradient representation for neural network visualization. In Advances in Neural Information Processing Systems (NeurIPS), 2019.
- [48] Sebastian Thrun and Tom M. Mitchell. Lifelong robot learning. *Robotics Auton. Syst.*, 15:25–46, 1995.
- [49] Eli Verwimp, Matthias De Lange, and Tinne Tuytelaars. Rehearsal revealed: The limits and merits of revisiting samples in continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9385–9394, October 2021.
- [50] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In Advances in Neural Information Processing Systems, volume 29, 2016.
- [51] Lijun Wang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Deep networks for saliency detection via local estimation and global search. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 3183– 3192, 2015.
- [52] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010.
- [53] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. In 6th International Conference on Learning Representations, 2018.
- [54] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3987–3995. PMLR, 06–11 Aug 2017.
- [55] Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neu-

- ral attention by excitation backprop. *Int. J. Comput. Vision*, 126(10):1084–1102, Oct. 2018.
- [56] Rui Zhao, Wanli Ouyang, Hongsheng Li, and Xiaogang Wang. Saliency detection by multi-context deep learning. In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 1265–1274, 2015.
- [57] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 2921–2929, 2016

Saliency Guided Experience Packing for Replay in Continual Learning

Appendix

Section A describes the steps of saliency map generation using Grad-CAM. Section B provides the dataset statistics used in different experiments. Pseudo-code of the episodic memory update in EPR is given in Section C. List of hyperparameters used for the baseline algorithms and our method is given in Section D. Additional results are provided in Section E.

A. Saliency Method: Grad-CAM

Gradient-weighted Class Activation Mapping (**Grad-CAM**) [41] is a saliency method that uses gradients to determine the impact of specific feature map activations on a given prediction. Since later layers in the convolutional neural network capture high-level semantics [26], taking gradients of a model output with respect to the feature map activations from one such layers identifies which high-level semantics are important for the model prediction. In our analysis, we select this layer and refer to as *target layer* [15]. List of *target layer* for different experiments is given in Table A.1.

Table A.1. Target layer names in PyTorch package for saliencies generated by different network architectures in Grad-CAM for different datasets.

Dataset	Network	Target Layer		
Split CIFAR	ResNet18 (reduced)	layer4.1.shortcut		
Split miniImageNet	ResNet18 (reduced)	layer4.1.shortcut		
Split CUB	ResNet18	net.layer4.1.conv2		

Let's consider the target layer has M feature maps where each feature map, $A^m \in \mathbb{R}^{u \times v}$ is of width u and height v. Also consider, for a given image $(I \in \mathbb{R}^{W \times H \times C})$ belonging to class c, the pre-softmax score of the image classifier is y_c . To obtain the class-discriminative saliency map, Grad-CAM first takes derivative of y_c with respect to each feature map A^m . These gradients are then global-average-pooled over u and v to obtain importance weight, α_m^c for each feature map:

$$\alpha_m^c = \frac{1}{uv} \sum_{i=1}^u \sum_{j=1}^v \frac{\partial y_c}{\partial A_{ij}^m},\tag{A.1}$$

where A_{ij}^m denotes location (i,j) in the feature map A^m . Next, these weights are used for computing linear combination of the feature map activations, which is then followed by ReLU to obtain the localization map:

$$L_{Grad-CAM}^{c} = \text{ReLU}\left(\sum_{m=1}^{M} \alpha_{m}^{c} A^{m}\right) \tag{A.2}$$

This map is of the same size $(u \times v)$ of A^m . Finally, saliency map, $I_{sm} \in \mathbb{R}^{W \times H}$ is generated by upsampling $L^c_{Grad-CAM}$ to the input image resolution using bilinear interpolation.

$$I_{sm} = \text{Upsample}\left(L_{Grad-CAM}^{c}\right) \tag{A.3}$$

B. Dataset Statistics

Table B.1. Statistics of the CIFAR-100, miniImageNet and CUB datasets used in task-incremental learning experiments.

	Split CIFAR	Split miniImageNet	Split CUB
num. of tasks	20	20	20
input size $(W \times H \times C)$	$32\times32\times3$	$84 \times 84 \times 3$	$224\times224\times3$
num. of classes/task	5	5	10
num. of training samples/tasks	2,500	2,500	300
num. of test samples/tasks	500	500	290

C. Memory Update Algorithm

Algorithm 2 Procedure for saliency guided episodic memory update in EPR

```
1: procedure UPDATEMEMORY(\mathcal{M}_E, \mathcal{M}_T, f_\theta, \text{EPF}, W_p)
             XAI: Procedure for saliency map generation; S_{sm}: stride; t^k: task-ID
             \text{Initialize: } \mathbf{I}_p \leftarrow [\ ]; \mathbf{c} \leftarrow [\ ]; \mathbf{x}_{cord} \leftarrow [\ ]; \mathbf{y}_{cord} \leftarrow [\ ]; \mathbf{P}_{pred} \leftarrow [\ ]
 3:
                                                                                                                                                                                 ▶ Initialize for memory selection
                                                                                                                            \triangleright Sample one example at a time without replacement from \mathcal{M}_T
 4:
             for (I, k, c) \sim \mathcal{M}_T do
 5:
                   I_{sm} \leftarrow \mathtt{XAI}(f_{\theta}, I, c)
                                                                                                                                                                 ⊳ generate saliency map using Equation 1
                   x_{cord}, y_{cord} \leftarrow \text{average-pool}(I_{sm}, W_p, S_{sm})
                                                                                                                                  \triangleright get corner coordinates of the most salient region in input, I
 6:
                   I_p \leftarrow I(x_{cord}: x_{cord} + W_p, y_{cord}: y_{cord} + W_p)
                                                                                                                                                                                         ⊳ get patch from Equation 4
 7:
                   I_p' \leftarrow \text{Zero-pad}(I_p, x_{cord}, y_{cord})
 8:
 9:
                   pred \leftarrow f_{\theta}(I_{p}^{'})
                                                                                                                                                             ⊳ check model prediction after zero-padding
10:
                   \mathbf{I}_p \leftarrow [\mathbf{I}_p, I_p]
                                                                                                                                                                                                                     ⊳ add patch
                   \mathbf{P}_{pred} \leftarrow [\mathbf{P}_{pred}, pred]
11:

    b add prediction

                                                                                                                                                                                                            ⊳ add class label
                   \mathbf{c} \leftarrow [\mathbf{c}, c]
12:
                   \mathbf{x}_{cord} \leftarrow [\mathbf{x}_{cord}, x_{cord}]
                                                                                                                                                                                                                    \triangleright add x_{cord}
13:
14:
                   \mathbf{y}_{cord} \leftarrow [\mathbf{y}_{cord}, y_{cord}]
                                                                                                                                                                                                                    \triangleright add y_{cord}
15:
             (\underline{\textbf{I}}_p, \textbf{c}, \textbf{x}_{cord}, \textbf{y}_{cord}) \leftarrow \texttt{select-patches}(\underline{\textbf{I}}_p, \textbf{c}, \textbf{x}_{cord}, \textbf{y}_{cord}, \underline{\textbf{P}}_{pred}, \texttt{EPF})
                                                                                                                                                                    ⊳ see section 6: memory patch selection
16:
17:
             \mathcal{B}_{\mathcal{M}_E} \leftarrow (\mathbf{I}_p, t^k, \mathbf{c})
18:
             \mathcal{M}_E \leftarrow \mathcal{M}_E \cup \{(\mathcal{B}_{\mathcal{M}_E}, \mathbf{x}_{cord}, \mathbf{y}_{cord})\}

    □ update episodic memory

19:
             return \mathcal{M}_E
20:
21: end procedure
```

D. List of Hyperparameters

List of hyperparameters used for both baseline methods and our approach is provided in Table D.1. EPF values used in different experiments in our method are given in Table D.2.

Table D.1. Hyperparameters grid considered for the baselines and our approach. The best values are given in parentheses. Here, 'lr' represents learning rate. In the table, we represent Split CIFAR as 'cifar', Split miniImageNet as 'minImg' and Split CUB as 'cub'. EPF is experience packing factor and \mathcal{M}_T is the temporary ring buffer in EPR.

Methods	Hyperparameters
Finetune	lr: 0.003, 0.01, 0.03 (cifar, minImg, cub), 0.1, 0.3, 1.0
EWC	$lr: 0.003, 0.01, 0.03$ (cifar, minImg, cub), 0.1, 0.3, 1.0 regularization, $\lambda: 0.1, 1, 10$ (cifar, minImg, cub), 100, 1000
RRR	<i>lr</i> : 0.003, 0.01 (cub), 0.03, 0.1, 0.3, 1.0 regularization : 10, 100 (cub), 1000
A-GEM	lr: 0.003, 0.01, 0.03 (cifar, minImg, cub), 0.1, 0.3, 1.0
MER	$lr: 0.003, 0.01, 0.03$ (cifar, minImg), 0.1 (cub), 0.3, 1.0 with in batch meta-learning rate, $\gamma: 0.01, 0.03, 0.1$ (cifar, minImg, cub), 0.3, 1.0 current batch learning rate multiplier, $s: 1, 2, 5$ (cifar, minImg, cub), 10
MEGA-I	lr : 0.003, 0.01, 0.03 (cifar, minImg, cub), 0.1, 0.3, 1.0 sensitivity parameter, ϵ : $1e^{-5}$, $1e^{-4}$, 0.001, 0.01 (cifar, minImg, cub), 0.1
DER++	$lr: 0.003, 0.01, 0.03 \ (minImg, cub), 0.1 \ (cifar), 0.3, 1.0$ regularization $\alpha: 0.1 \ (minImg), 0.2 \ (cifar), 0.5 \ (cub), 1.0$ regularization, $\beta: 0.5 \ (cifar, minImg, cub), 1.0$
ASER	lr: 0.003, 0.01, 0.03 (cub), 0.1 (cifar, minImg), 0.3, 1.0 $K: 3$ (cifar, minImg, cub); $N_c: 100$ (cifar, miniImg), 150 (cub), 250
ER-Reservoir	lr: 0.003, 0.01, 0.03 (cub), 0.1 (cifar, minImg), 0.3, 1.0
ER-RING	lr: 0.003, 0.01, 0.03 (cifar, minImg, cub), 0.1, 0.3, 1.0
HAL	$lr: 0.003, 0.01, 0.03$ (cifar, minImg), 0.1, 0.3, 1.0 regularization, $\lambda: 0.01, 0.03, 0.1, 0.3$ (minImg), 1 (cifar), 3, 10 mean embedding strength, $\gamma: 0.01, 0.03, 0.1$ (cifar, minImg), 0.3, 1, 3, 10 decay rate, $\beta: 0.5$ (cifar, minImg) gradient steps on anchors, $k: 100$ (cifar, minImg)
Multitask	lr: 0.003, 0.01, 0.03 (cifar, minImg, cub), 0.1, 0.3, 1.0
EPR (ours)	lr (task-incremental) : 0.01, 0.03 (cub), 0.05 (minImg), 0.1 (cifar), 0.3, 1.0 lr (class-incremental) : 0.01, 0.05 (cifar, minImg), 0.1 examples per class temporarily stored in \mathcal{M}_T : $\gamma \times \text{EPF}$; γ : 2(cub), 5 (cifar,minImg) stride, S_{sm} : 1 (cifar, minImg), 2, 3 (cub)

Table D.2. Experience Packing Factor (EPF) for different n_{sc} used in our (a) task-incremental learning and (b) class-incremental learning experiments. Input image width, W for CIFAR, miniImageNet and CUB dataset are 32,84 and 224 respectively. For given n_{sc} , EPF and W, corresponding memory patch sizes (W_p) are also given in the table.

(a)								
Split CIFAR Split miniImageNet Split CUB								
EPF	W_p	EPF	W_p	EPF	W_p			
3	26	5	53	7	119			
2	22	3	48	4	112			
1	27	2	51	3	112			
1	22	2	42	2	112			
	EPF 3	EPF W_p 3 26 2 22 1 27						

			(b)		
CIFAR-1 (20 Task				miniIma (10 Tas	_
$ M_E $	n_{sc}	EPF	W_p	EPF	W_p
2k	20	25	28	25	75
1k	10	13	28	13	73

E. Additional Results

Table E.1. Performance comparison of different experience replay methods for different memory sizes in task-incremental learning setup. Number of memory slots per class, n_{sc} ={0.5, 0.75} refers to memory size, $|\mathcal{M}_E|$ ={42, 64} for CIFAR and miniImageNet, and $|\mathcal{M}_E|$ ={85, 128} for CUB. Average and standard deviations are computed over 5 runs for different random seeds.

		Split CIFAR		CIFAR Split miniImageNet			Split CUB		
n_{sc}	Methods	ACC (%)	BWT	ACC (%)	BWT	ACC (%)	BWT		
-	Finetune	42.9 ± 2.07	-0.25 ± 0.03	34.7 ± 2.69	-0.26 ± 0.03	55.7 ± 2.22	-0.13 ± 0.03		
	MEGA-I	48.9 ± 1.68	- 0.21 ± 0.01	43.8 ± 1.58	-0.14 ± 0.01	61.5 ± 2.08	- 0.08 ± 0.01		
0.75	DER++	50.0 ± 1.81	- 0.19 ± 0.02	47.2 ± 1.54	-0.12 ± 0.01	64.8 ± 1.61	- 0.06 ± 0.01		
	ER-RING	50.4 ± 0.85	- 0.21 ± 0.02	44.9 ± 1.49	- 0.14 ± 0.02	64.0 ± 1.29	- 0.05 ± 0.01		
	EPR (Ours)	$\textbf{56.8} \pm \textbf{1.59}$	- 0.12 ± 0.02	$\textbf{51.1} \pm \textbf{1.47}$	- 0.06 \pm 0.01	$\textbf{70.7} \pm \textbf{0.72}$	- 0.03 \pm 0.01		
	MEGA-I	43.7 ± 1.26	-0.26 ± 0.02	39.6 ± 2.35	-0.18 ± 0.02	57.7 ± 0.62	-0.11 ± 0.01		
0.5	DER++	47.5 ± 1.58	- 0.21 ± 0.01	45.6 ± 0.56	- 0.13 ± 0.01	62.5 ± 1.45	- 0.08 ± 0.01		
	ER-RING	44.6 ± 0.84	- 0.27 ± 0.01	39.1 ± 1.38	- 0.20 ± 0.02	59.2 ± 0.97	- 0.10 ± 0.01		
	EPR (Ours)	$\textbf{55.6} \pm \textbf{0.54}$	- 0.13 ± 0.02	$\textbf{49.2} \pm \textbf{1.20}$	- 0.07 \pm 0.01	$\textbf{70.3} \pm \textbf{0.91}$	- 0.03 \pm 0.01		