# 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, $\alpha_m^c$ 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 $L_{Grad-CAM}^c$ 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 ($W \times H \times C$) | $32 \times 32 \times 3$ | $84 \times 84 \times 3$ | $224 \times 224 \times 3$ |
| 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($\mathcal{M}_E, \mathcal{M}_T, f_\theta, \text{EPF}, W_p$)
2:    XAI: Procedure for saliency map generation; $S_{sm}$: stride; $t^k$: task-ID
3:    Initialize: $\mathbf{I}_p \leftarrow [\,]; \mathbf{c} \leftarrow [\,]; \mathbf{x}_{cord} \leftarrow [\,]; \mathbf{y}_{cord} \leftarrow [\,]; \mathbf{P}_{pred} \leftarrow [\,]$         ▷ Initialize for memory selection
4:    **for** $(I, k, c) \sim \mathcal{M}_T$ **do**         ▷ Sample one example at a time without replacement from $\mathcal{M}_T$
5:       $I_{sm} \leftarrow \text{XAI}(f_\theta, I, c)$         ▷ generate saliency map using Equation 1
6:       $x_{cord}, y_{cord} \leftarrow \text{average-pool}(I_{sm}, W_p, S_{sm})$         ▷ get corner coordinates of the most salient region in input,$I$
7:       $I_p \leftarrow I(x_{cord} : x_{cord} + W_p, y_{cord} : y_{cord} + W_p)$         ▷ get patch from Equation 4
8:       $I_p' \leftarrow \text{Zero-pad}(I_p, x_{cord}, y_{cord})$
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
11:      $\mathbf{P}_{pred} \leftarrow [\mathbf{P}_{pred}, pred]$         ▷ add prediction
12:      $\mathbf{c} \leftarrow [\mathbf{c}, c]$         ▷ add class label
13:      $\mathbf{x}_{cord} \leftarrow [\mathbf{x}_{cord}, x_{cord}]$         ▷ add $x_{cord}$
14:      $\mathbf{y}_{cord} \leftarrow [\mathbf{y}_{cord}, y_{cord}]$         ▷ add $y_{cord}$
15:    **end for**
16:    $(\mathbf{I}_p, \mathbf{c}, \mathbf{x}_{cord}, \mathbf{y}_{cord}) \leftarrow \text{select-patches}(\mathbf{I}_p, \mathbf{c}, \mathbf{x}_{cord}, \mathbf{y}_{cord}, \mathbf{P}_{pred}, \text{EPF})$         ▷ see section 6: memory patch selection
17:    $t^k \leftarrow k$
18:    $\mathcal{B}_{\mathcal{M}_E} \leftarrow (\mathbf{I}_p, t^k, \mathbf{c})$
19:    $\mathcal{M}_E \leftarrow \mathcal{M}_E \cup \{(\mathcal{B}_{\mathcal{M}_E}, \mathbf{x}_{cord}, \mathbf{y}_{cord})\}$         ▷ update episodic memory
20:    **return** $\mathcal{M}_E$
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 <br> regularization, $\lambda$ : 0.1, 1, 10 (cifar, minImg, cub), 100, 1000 |
| RRR | $lr$ : 0.003, 0.01 (cub), 0.03, 0.1, 0.3, 1.0 <br> 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 <br> with in batch meta-learning rate, $\gamma$ : 0.01, 0.03, 0.1 (cifar, minImg, cub), 0.3, 1.0 <br> 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 <br> sensitivity parameter, $\epsilon$ : $1e^{-5}$, $1e^{-4}$, 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 <br> regularization $\alpha$ : 0.1 (minImg), 0.2 (cifar), 0.5 (cub), 1.0 <br> 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 <br> $K$ : 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 <br> regularization, $\lambda$ : 0.01, 0.03, 0.1, 0.3 (minImg), 1 (cifar), 3, 10 <br> mean embedding strength, $\gamma$ : 0.01, 0.03, 0.1 (cifar, minImg), 0.3, 1, 3, 10 <br> decay rate, $\beta$ : 0.5 (cifar, minImg) <br> 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 <br> $lr$ (class-incremental) : 0.01, 0.05 (cifar, minImg), 0.1 <br> examples per class temporarily stored in $\mathcal{M}_T$ : $\gamma \times$ EPF; $\gamma$ : 2(cub), 5 (cifar,minImg) <br> 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)

| $n_{sc}$ | Split CIFAR | | Split miniImageNet | | Split CUB | |
|---|---|---|---|---|---|---|
| | EPF | $W_p$ | EPF | $W_p$ | EPF | $W_p$ |
| 2 | 3 | 26 | 5 | 53 | 7 | 119 |
| 1 | 2 | 22 | 3 | 48 | 4 | 112 |
| 0.75 | 1 | 27 | 2 | 51 | 3 | 112 |
| 0.5 | 1 | 22 | 2 | 42 | 2 | 112 |

(b)

| $|M_E|$ | $n_{sc}$ | CIFAR-100 (20 Tasks) | | miniImageNet (10 Tasks) | |
|---|---|---|---|---|---|
| | | 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 | | Split miniImageNet | | Split CUB | |
|---|---|---|---|---|---|---|---|
| $n_{sc}$ | **Methods** | ACC (%) | BWT | ACC (%) | BWT | ACC (%) | BWT |
| - | Finetune | $42.9 \pm 2.07$ | $- 0.25 \pm 0.03$ | $34.7 \pm 2.69$ | $- 0.26 \pm 0.03$ | $55.7 \pm 2.22$ | $- 0.13 \pm 0.03$ |
| 0.75 | MEGA-I | $48.9 \pm 1.68$ | $- 0.21 \pm 0.01$ | $43.8 \pm 1.58$ | $- 0.14 \pm 0.01$ | $61.5 \pm 2.08$ | $- 0.08 \pm 0.01$ |
| | DER++ | $50.0 \pm 1.81$ | $- 0.19 \pm 0.02$ | $47.2 \pm 1.54$ | $- 0.12 \pm 0.01$ | $64.8 \pm 1.61$ | $- 0.06 \pm 0.01$ |
| | ER-RING | $50.4 \pm 0.85$ | $- 0.21 \pm 0.02$ | $44.9 \pm 1.49$ | $- 0.14 \pm 0.02$ | $64.0 \pm 1.29$ | $- 0.05 \pm 0.01$ |
| | **EPR (Ours)** | $\mathbf{56.8 \pm 1.59}$ | $\mathbf{- 0.12 \pm 0.02}$ | $\mathbf{51.1 \pm 1.47}$ | $\mathbf{- 0.06 \pm 0.01}$ | $\mathbf{70.7 \pm 0.72}$ | $\mathbf{- 0.03 \pm 0.01}$ |
| 0.5 | MEGA-I | $43.7 \pm 1.26$ | $- 0.26 \pm 0.02$ | $39.6 \pm 2.35$ | $- 0.18 \pm 0.02$ | $57.7 \pm 0.62$ | $- 0.11 \pm 0.01$ |
| | DER++ | $47.5 \pm 1.58$ | $- 0.21 \pm 0.01$ | $45.6 \pm 0.56$ | $- 0.13 \pm 0.01$ | $62.5 \pm 1.45$ | $- 0.08 \pm 0.01$ |
| | ER-RING | $44.6 \pm 0.84$ | $- 0.27 \pm 0.01$ | $39.1 \pm 1.38$ | $- 0.20 \pm 0.02$ | $59.2 \pm 0.97$ | $- 0.10 \pm 0.01$ |
| | **EPR (Ours)** | $\mathbf{55.6 \pm 0.54}$ | $\mathbf{- 0.13 \pm 0.02}$ | $\mathbf{49.2 \pm 1.20}$ | $\mathbf{- 0.07 \pm 0.01}$ | $\mathbf{70.3 \pm 0.91}$ | $\mathbf{- 0.03 \pm 0.01}$ |