

Transfer without Forgetting

Supplementary Materials

Matteo Boschini¹, Lorenzo Bonicelli¹, Angelo Porrello¹,
Giovanni Bellitto², Matteo Pennisi², Simone Palazzo²,
Concetto Spampinato², and Simone Calderara¹

¹ AImageLab, University of Modena and Reggio Emilia, Italy
firstname.lastname@unimore.it

² PeRCeiVe Lab, University of Catania, Italy
firstname.lastname@unict.it

A Additional Details on the Model

In this section, we report some additional details on the inner workings of the model which were omitted in the main paper for the sake of brevity.

A.a Further details on \mathbb{M}_{Sp}

The spatial attention map \mathbb{M}_{Sp} is computed on top of the activations of a given layer of the fixed sibling network $\hat{h} \in \mathbb{R}^{b \times c \times h \times w}$, processed through a ResNet-inspired bottleneck structure [1,2]. In detail, we expand and detail Eq. 5 in the main paper:

$$\mathbb{M}_{\text{Sp}} \triangleq C_{1 \times 1}^C \circ \text{ReLU} \circ \text{BN} \circ C_{3 \times 3}^B \circ \text{ReLU} \circ \text{BN} \circ C_{3 \times 3}^B \circ \text{ReLU} \circ \text{BN} \circ C_{1 \times 1}^A, \quad (1)$$

where ReLU denotes a ReLU activation, BN indicates a Batch Normalization layer (conditioned on the task-identifier) and C indicates a Convolutional layer. More specifically, $C_{1 \times 1}^A$ is a 1×1 convolution, projecting from c channels to $c/4$; $C_{3 \times 3}^B$ is a 3×3 dilated convolution with dilation factor 2 and adequate padding to maintain the same spatial resolution as the input, with $c/4$ channels both as input and output; $C_{1 \times 1}^C$ is a 1×1 convolution projecting from $c/4$ channels to 1 channel. This results in \mathbb{M}_{Sp} having shape $b \times 1 \times h \times w$.

A.b Scaling of \mathbb{M}

The second distillation term in Eq. 9 requires storing the binary attention maps \mathbb{M} computed for each sample stored in the memory buffer. While this implies a memory overhead, we point out that this is limited by two factors:

- The binary nature of \mathbb{M} means its elements can be saved using the smallest supported data-type (usually 1 byte due to hardware constraints);
- As \mathbb{M} usually encodes low level features, it contains several redundancies that can be exploited by (a) using lossless compression algorithms, or (b) down-sampling its spatial dimensions before saving.

In TwF we save the feature maps M as bytes and apply down-scaling – with *nearest neighbor* rule – with a factor of 2 if the spatial dimensions are over 16×16 . We use the same strategy to up-scale the maps before computing Eq. 9.

B Hyperparameters

For the experiments of Sec. 4, we employed a choice of hyperparameters validated by grid-search on a random split of 10% of the training set. In the following, we list the values resulting from this process, which can be used to replicate our result. For the sake of fairness, we initialize all models from the same pre-training weights and fix the allowance in terms of iterations and sample efficiency by excluding the number of epochs, lr decay schedule and batch size from the grid-search³.

Split CIFAR-10 - Class-IL	
<i>shared</i>	Eps : 50 bs : 32 Eps _{pretr} : 200 lr _{decay} : no
JOINT	- lr : 0.1
SGD	- lr : 0.1
oEwC	- lr : 0.1 λ : 10 γ : 1
LwF	- lr : 0.1 α : 0.3 τ : 2 wd : 0.0001
ER	500 lr : 0.1
	5120 lr : 0.1
CO ² L	500 lr : 0.5 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
	5120 lr : 0.5 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	500 lr : 0.1 wd : 10^{-5}
	5120 lr : 0.03 wd : 10^{-5}
DER++	500 lr : 0.03 α : 0.2 β : 0.5
	5120 lr : 0.03 α : 0.1 β : 1
ER-ACE	500 lr : 0.03
	5120 lr : 0.03
TwF	500 lr : 0.03 α : 0.3 β : 0.9 λ : 0.1 λ_{FP} : 5×10^{-3} λ_{FP}^{repl} : 0.1
	5120 lr : 0.1 α : 0.3 β : 0.9 λ : 0.1 λ_{FP} : 5×10^{-3} λ_{FP}^{repl} : 0.3

³ It must be noted that, to allow for its regular operation, CO²L demands a larger batch size. All results for this method are influenced by this advantage.

Split CIFAR-10 - Task-IL

<i>shared</i>	Eps : 50 bs : 32 Eps _{pretr} : 200 lr _{decay} : no
JOINT	- lr : 0.1
SGD	- lr : 0.1
oEwC	- lr : 0.03 λ : 0.5 γ : 1
LwF	- lr : 0.01 α : 0.3 τ : 2 wd : 0.0001
ER	500 lr : 0.1 5120 lr : 0.1
CO ² L	500 lr : 0.5 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07 5120 lr : 0.5 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	500 lr : 0.1 wd : 10^{-5} 5120 lr : 0.03 wd : 10^{-5}
DER++	500 lr : 0.03 α : 0.2 β : 0.5 5120 lr : 0.03 α : 0.1 β : 1
ER-ACE	500 lr : 0.03 5120 lr : 0.03
TwF	500 lr : 0.03 α : 0.3 β : 0.9 λ : 0.1 λ_{FP} : 5×10^{-3} $\lambda_{\text{FP}}^{\text{repl}}$: 0.1 5120 lr : 0.1 α : 0.3 β : 0.9 λ : 0.1 λ_{FP} : 5×10^{-3} $\lambda_{\text{FP}}^{\text{repl}}$: 0.3

Split CIFAR-100 - Class-IL

<i>shared</i>	Eps : 50 bs : 64 Eps _{pretr} : 200 lr _{decay} : 0.1 lr _{decay} ^{steps} : [35, 45]
JOINT	- lr : 0.1
SGD	- lr : 0.1
oEwC	- lr : 0.1 λ : 5 γ : 1
LwF	- lr : 0.03 α : 0.3 τ : 2 wd : 0.0005
ER	500 lr : 0.01 2000 lr : 0.01
CO ² L	500 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07 2000 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	500 lr : 1 wd : 10^{-5} 2000 lr : 1 wd : 10^{-5}
DER++	500 lr : 0.1 α : 0.3 β : 0.3 2000 lr : 0.1 α : 0.1 β : 0.5
ER-ACE	500 lr : 0.1 2000 lr : 0.1
TwF	500 lr : 0.03 α : 0.3 β : 1.2 λ : 0.3 λ_{FP} : 0.03 $\lambda_{\text{FP}}^{\text{repl}}$: 1.5 2000 lr : 0.1 α : 0.3 β : 1.2 λ : 0.3 λ_{FP} : 5×10^{-3} $\lambda_{\text{FP}}^{\text{repl}}$: 1.2

Split CIFAR-100 - Task-IL

<i>shared</i>	Eps : 50 bs : 64 Eps _{pretr} : 200 lr _{decay} : 0.1 lr _{decay} ^{steps} : [35, 45]
JOINT	- lr : 0.1
SGD	- lr : 0.01
oEwC	- lr : 0.01 λ : 0.5 γ : 0.7
LwF	- lr : 0.03 α : 0.3 τ : 2 wd : 0.0005
ER	500 lr : 0.01
	2000 lr : 0.01
CO ² L	500 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
	2000 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	500 lr : 1 wd : 10 ⁻⁵
	2000 lr : 1 wd : 10 ⁻⁵
DER++	500 lr : 0.1 α : 0.3 β : 1.2
	2000 lr : 0.1 α : 0.1 β : 0.5
ER-ACE	500 lr : 0.1
	2000 lr : 0.1
TwF	500 lr : 0.03 α : 0.3 β : 1.2 λ : 0.3 λ_{FP} : 0.03 $\lambda_{FP}^{\text{repl}}$: 1.5
	2000 lr : 0.1 α : 0.3 β : 0.8 λ : 0.3 λ_{FP} : 0.03 $\lambda_{FP}^{\text{repl}}$: 0.3

Split CUB-200 - Class-IL

<i>shared</i>	Eps : 50 bs : 64 Eps _{pretr} : 50
JOINT	- lr : 0.1
SGD	- lr : 0.1
oEwC	- lr : 0.01 λ : 1 γ : 1
LwF	- lr : 0.1 α : 1 τ : 2 wd : 0.0005
ER	400 lr : 0.03
	1000 lr : 0.1
CO ² L	400 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
	1000 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	400 lr : 0.1 wd : 10 ⁻⁵
	1000 lr : 0.1 wd : 10 ⁻⁵
DER++	400 lr : 0.1 α : 1 β : 0.5
	1000 lr : 0.1 α : 0.5 β : 0.5
ER-ACE	400 lr : 0.1
	1000 lr : 0.1
TwF	400 lr : 0.03 α : 1 β : 1 λ : 0.3 λ_{FP} : 5 \times 10 ⁻⁴ $\lambda_{FP}^{\text{repl}}$: 0.1
	1000 lr : 0.03 α : 1 β : 1.2 λ : 0.3 λ_{FP} : 5 \times 10 ⁻⁴ $\lambda_{FP}^{\text{repl}}$: 0.1

Split CUB-200 - Task-IL

<i>shared</i>	Eps : 50 bs : 64 Eps _{pretr} : 50
JOINT	- lr : 0.1
SGD	- lr : 0.1
oEwC	- lr : 0.1 λ : 0.5 γ : 0.9
LwF	- lr : 0.1 α : 1 τ : 2 wd : 0.0005
ER	400 lr : 0.1
	1000 lr : 0.1
CO ² L	400 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
	1000 lr : 0.1 bs : 256 κ : 0.2 λ : 1 lr _{lin} : 1 lr _{decay} ^{lin} : 0.2 κ^* : 0.01 τ : 0.07
iCaRL	400 lr : 0.1 wd : 10^{-5}
	1000 lr : 0.1 wd : 10^{-5}
DER++	400 lr : 0.1 α : 0.5 β : 0.5
	1000 lr : 0.1 α : 0.5 β : 0.5
ER-ACE	400 lr : 0.1
	1000 lr : 0.1
TwF	400 lr : 0.03 α : 0.3 β : 1 λ : 0.3 λ_{FP} : 5×10^{-4} λ_{FP}^{repl} : 0.1
	1000 lr : 0.03 α : 1 β : 1 λ : 0.3 λ_{FP} : 5×10^{-4} λ_{FP}^{repl} : 0.1

References

1. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (2016)
2. Park, J., Woo, S., Lee, J.Y., Kweon, I.S.: Bam: Bottleneck attention module. In: British Machine Vision Conference (2018)