

# ICICLE: Interpretable Class Incremental Continual Learning – Supplementary Materials

Dawid Rymarczyk<sup>1,2,3,\*</sup>

Joost van de Weijer<sup>4,5</sup>

Bartosz Zieliński<sup>1,3,6</sup>

Bartłomiej Twardowski<sup>4,5,6</sup>

<sup>1</sup> Faculty of Mathematics and Computer Science, Jagiellonian University

<sup>2</sup> Doctoral School of Exact and Life Sciences, Jagiellonian University

<sup>3</sup> Ardigen SA

<sup>4</sup> Autonomous University of Barcelona

<sup>5</sup> Computer Vision Center

<sup>6</sup> IDEAS NCBR

\*dawid.rymarczyk@doctoral.uj.edu.pl

## Additional experimental setup details

Here we present additional details on the experimental setup. We performed a hyperparameter search for  $\lambda_{dist}$  ( $\lambda_{dist} \in \{10.0, 5.0, 0.0, 0.5, 0.1, 0.05, 0.01, 0.005, 0.001, 0.0005, 0.0001\}$ ). We use Adam optimizer with a learning rate of 0.001 and parameters  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . We set the batch size to 75 and use input images of resolution  $224 \times 224 \times 3$ . The weights of the network are initialized with Xavier’s normal initializer.

We perform a warmup training where the weights of  $f$  are frozen for 10 epochs, and then we train the model until it converges with 12 epochs early stopping. We use the learning schema presented in Table 1. Depending on the number of tasks, we perform warm-up training with  $\{5, 5, 4\}$  epochs and joint training phase for  $\{21, 15, 11\}$ , longer with fewer tasks. Similarly, we perform prototype projection every  $\{10, 7, 5\}$  epoch. So with more tasks, we perform fewer training epochs (Table 1).

## Results on Stanford Cars

Table 2 depicts how the ICICLE method works on the Stanford Cars dataset compared to other baseline methods. Results are consistent with the ones on CUB-200-2011 and show that ICICLE outperforms all standard CL learning methods adapted to ProtoPNet architecture.

## Comparison to GDumb

Additionally, we compared our approach with GDumb [1], a baseline method in CL learning, in scenarios involving 3, 5, and 10 images per class with 4 tasks learning achieving 20.3, 34.2, 57.6, and 13.0, 26.7, 48.8 for task-aware and task-agnostic respectively. ICICLE outperformed GDumb with a small number of examples, and task-aware for GDumb-10 was the only exception

where GDumb achieved a higher accuracy score.

## Distant initialization

In Table 5, we showed that proximity-based initialization is most beneficial. However, here in Figure 1, we show how initialization of prototypical parts at a distance from already existing once generates concepts that are too general or carry information about background.

## Task-recency bias compensation

Here, in Table 3 we show what is the influence of task-recency bias compensation on each task accuracy for task agnostic scenario. After compensation, the accuracy on task 1 increased the most, but at the same time accuracy on all other tasks was sacrificed. However, the average accuracy after the compensation is increased.

## Detailed results after learning each task

In Table 4, and Table 5 we show detailed accuracies of the ICICLE after learning each task for CUB-200-2011 on a single run with the same seed in four and ten tasks learning scenarios.

## Analysis of the hyperparameters for baselines

In Table 6, Table 7, and Table 8 we show the influence of the hyperparameters for each of the baseline methods, EWC, LWF, and LWM, respectively. Based on that, the parameters of these methods were chosen for comparison with ICICLE.

Phase	Model layers	Learning rate	Scheduler	Weight decay	Duration
Warm-up	add-on $1 \times 1$ convolution	$1 \cdot 10^{-3}$	None	None	5, 5, 4 epochs
	prototypical layer	$1 \cdot 10^{-3}$			
Joint	convolutions $f$	$1 \cdot 10^{-4}$	by half every 5 epochs	$10^{-4}$	21, 15, 10 epochs early stopping
	add-on $1 \times 1$ convolution prototypical layer	$1 \cdot 10^{-3}$			

Table 1. Learning schema for the ICICLE method.

METHOD	AVG. INC. TASK-AWARE ACCURACY			AVG. INC. TASK-TASK AGNOSTIC ACCURACY		
	4 TASKS	7 TASKS	14 TASKS	4 TASKS	7 TASKS	14 TASKS
FREEZING	$0.572 \pm 0.031$	$0.518 \pm 0.041$	$0.486 \pm 0.026$	$0.309 \pm 0.012$	$0.155 \pm 0.031$	$0.092 \pm 0.014$
FINETUNING	$0.216 \pm 0.009$	$0.167 \pm 0.011$	$0.149 \pm 0.012$	$0.182 \pm 0.006$	$0.124 \pm 0.013$	$0.057 \pm 0.001$
EWC	$0.456 \pm 0.021$	$0.315 \pm 0.037$	$0.287 \pm 0.041$	$0.258 \pm 0.019$	$0.152 \pm 0.022$	$0.011 \pm 0.009$
LWM	$0.459 \pm 0.072$	$0.416 \pm 0.048$	$0.305 \pm 0.022$	$0.233 \pm 0.026$	$0.171 \pm 0.016$	$0.080 \pm 0.008$
LWF	$0.375 \pm 0.021$	$0.356 \pm 0.024$	$0.250 \pm 0.020$	$0.230 \pm 0.011$	$0.171 \pm 0.005$	$0.092 \pm 0.008$
ICICLE	<b><math>0.654 \pm 0.014</math></b>	<b><math>0.645 \pm 0.003</math></b>	<b><math>0.583 \pm 0.048</math></b>	<b><math>0.335 \pm 0.005</math></b>	<b><math>0.203 \pm 0.010</math></b>	<b><math>0.116 \pm 0.018</math></b>

Table 2. Average incremental accuracy comparison for different numbers of tasks on Stanford Cars, demonstrating the negative impact of the high number of tasks to be learned on models' performance. Despite this trend, ICICLE outperforms the baseline methods across all task numbers.

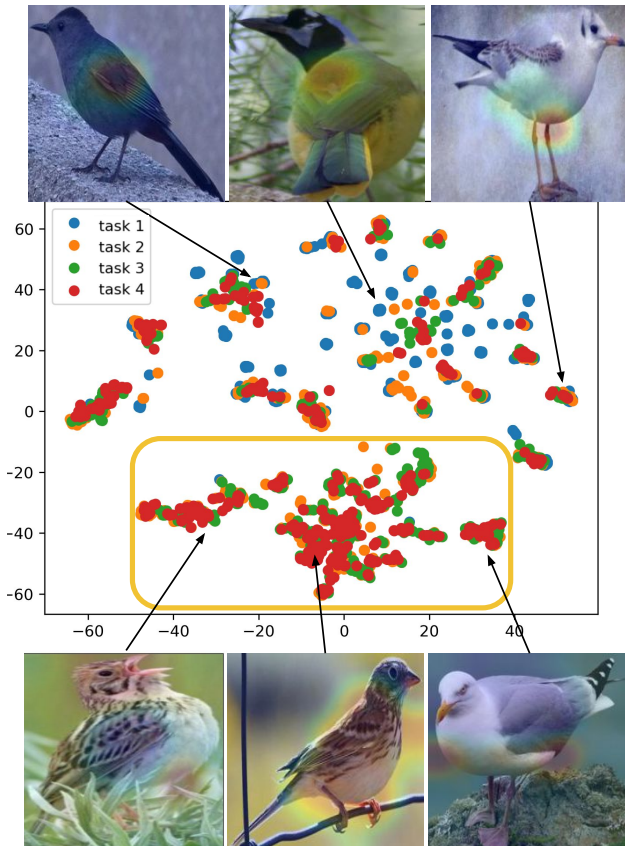


Figure 1. Image depicts the 2D TSNE projection of prototypes. One can observe that there is a cluster of prototypes from all tasks except the first one (yellow box). This is observable with distant initialization and prototype visualizations show that those prototypical parts are representing the background or vague concepts.

	TASK 1	TASK 2	TASK 3	TASK 4	AVG
TASK AWARE	0.514	0.717	0.725	0.698	0.663
BEFORE	0.028	0.301	<b>0.434</b>	<b>0.575</b>	0.335
AFTER	<b>0.233</b>	<b>0.365</b>	0.314	0.486	<b>0.350</b>

Table 3. Task-recency bias compensation influence of a single run. Results show that our compensation method balance more the results per each task in task agnostic scenario.

After	TASK-AWARE ACCURACY					TASK-TASK AGNOSTIC ACCURACY				
	TASK 1	TASK 2	TASK 3	TASK 4	AVG	TASK 1	TASK 2	TASK 3	TASK 4	AVG
TASK 1	0.806	NA	NA	NA	0.806	0.806	NA	NA	NA	0.806
TASK 2	0.740	0.759	NA	NA	0.750	0.089	0.747	NA	NA	0.418
TASK 3	0.622	0.736	0.759	NA	0.706	0.033	0.633	0.549	NA	0.404
TASK 4	0.514	0.717	0.725	0.698	0.663	0.028	0.484	0.378	0.505	0.349

Table 4. Results of the ICICLE method before task-recency compensation for four task learning scenario after each learning episode.

After	TASK-AWARE ACCURACY										
	TASK 1	TASK 2	TASK 3	TASK 4	TASK 5	TASK 6	TASK 7	TASK 8	TASK 9	TASK 10	AVG
TASK 1	0.920	NA	NA	NA	NA	NA	NA	NA	NA	NA	0.920
TASK 2	0.666	0.869	NA	NA	NA	NA	NA	NA	NA	NA	0.767
TASK 3	0.462	0.818	0.858	NA	NA	NA	NA	NA	NA	NA	0.713
TASK 4	0.420	0.751	0.774	0.774	NA	NA	NA	NA	NA	NA	0.680
TASK 5	0.314	0.625	0.680	0.672	0.784	NA	NA	NA	NA	NA	0.615
TASK 6	0.268	0.538	0.627	0.617	0.760	0.747	NA	NA	NA	NA	0.593
TASK 7	0.265	0.476	0.584	0.617	0.713	0.706	0.769	NA	NA	NA	0.590
TASK 8	0.258	0.413	0.551	0.555	0.667	0.634	0.741	0.764	NA	NA	0.573
TASK 9	0.253	0.398	0.494	0.492	0.598	0.587	0.701	0.745	0.852	NA	0.569
TASK 10	0.244	0.371	0.462	0.441	0.573	0.560	0.667	0.729	0.816	0.803	0.567

	TASK-AGNOSTIC ACCURACY										
	TASK 1	TASK 2	TASK 3	TASK 4	TASK 5	TASK 6	TASK 7	TASK 8	TASK 9	TASK 10	AVG
TASK 1	0.920	NA	NA	NA	NA	NA	NA	NA	NA	NA	0.920
TASK 2	0.010	0.869	NA	NA	NA	NA	NA	NA	NA	NA	0.439
TASK 3	0.0	0.060	0.854	NA	NA	NA	NA	NA	NA	NA	0.305
TASK 4	0.0	0.0	0.339	0.751	NA	NA	NA	NA	NA	NA	0.273
TASK 5	0.0	0.0	0.030	0.323	0.746	NA	NA	NA	NA	NA	0.220
TASK 6	0.0	0.0	0.004	0.090	0.451	0.684	NA	NA	NA	NA	0.205
TASK 7	0.0	0.0	0.0	0.020	0.193	0.432	0.712	NA	NA	NA	0.194
TASK 8	0.0	0.0	0.0	0.020	0.073	0.233	0.497	0.643	NA	NA	0.181
TASK 9	0.0	0.0	0.0	0.0	0.035	0.138	0.338	0.435	0.676	NA	0.180
TASK 10	0.0	0.0	0.0	0.0	0.016	0.070	0.214	0.285	0.484	0.643	0.171

Table 5. Results of the ICICLE method before task-recency compensation for ten task learning scenario after each learning episode.

$\alpha$	TASK-AWARE ACCURACY					TASK-TASK AGNOSTIC ACCURACY				
	0.01	0.1	1.0	5.0	10.0	0.01	0.1	1.0	5.0	10.0
	0.185	0.329	0.441	0.197	0.167	0.170	0.185	0.213	0.168	0.144

Table 6. Influence of the  $\alpha$  parameter in EWC on the accuracy of ProtoPNet architecture in four task learning scenario.

$\gamma$	TASK-AWARE ACCURACY					TASK-TASK AGNOSTIC ACCURACY				
	0.001	0.01	0.1	1.0	10.0	0.001	0.01	0.1	1.0	10.0
	0.240	0.240	0.431	0.355	0.231	0.209	0.209	0.212	0.209	0.209

Table 7. Influence of the  $\gamma$  parameter in LWM on the accuracy of ProtoPNet architecture in four task learning scenario.

$\lambda$	TASK-AWARE ACCURACY					TASK-TASK AGNOSTIC ACCURACY				
	0.001	0.01	0.1	1.0	10.0	0.001	0.01	0.1	1.0	10.0
	0.232	0.232	0.238	0.359	0.249	0.207	0.210	0.210	0.231	0.221

Table 8. Influence of the  $\lambda$  parameter in LWF on the accuracy of ProtoPNet architecture in four task learning scenario.

## References

- [1] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II* 16, pages 524–540. Springer, 2020. [1](#)