

# Dual Cognitive Architecture: Incorporating Biases and Multi-Memory Systems for Lifelong Learning

Anonymous authors

Paper under double-blind review

## Abstract

Artificial neural networks (ANNs) exhibit a narrow scope of expertise on stationary independent data. However, the data in the real world is continuous and dynamic, and ANNs must adapt to novel scenarios while also retaining the learned knowledge to become lifelong learners. The ability of humans to excel at these tasks can be attributed to multiple factors ranging from cognitive computational structures, cognitive biases, and the multi-memory systems in the brain. We incorporate key concepts from each of these to design a novel framework, *Dual Cognitive Architecture (DUCA)*, which includes multiple sub-systems, implicit and explicit knowledge representation dichotomy, inductive bias, and a multi-memory system. DUCA shows improvement across different settings and datasets, and it also exhibits reduced task recency bias, without the need for extra information. To further test the versatility of lifelong learning methods on a challenging distribution shift, we introduce a novel domain-incremental dataset *DN4IL*. In addition to improving performance on existing benchmarks, DUCA also demonstrates superior performance on this complex dataset.<sup>1</sup>

## 1 Introduction

Deep learning has seen rapid progress in recent years, and supervised learning agents have achieved superior performance on perception tasks. However, unlike a supervised setting, where data is static and independent and identically distributed, real-world data is changing dynamically. Continual learning (CL) aims to learn multiple tasks when data is streamed sequentially (Parisi et al., 2019). This is crucial in real-world deployment settings, as the model needs to adapt quickly to novel data (plasticity), while also retaining previously learned knowledge (stability). Artificial neural networks (ANN), however, are still not effective lifelong learners, as they often fail to generalize to small changes in distribution and also suffer from forgetting old information when presented with new data (catastrophic forgetting) (McCloskey & Cohen, 1989).

Humans, on the other hand, show a better ability to acquire new skills while also retaining previously learned skills to a greater extent. This intelligence can be attributed to different factors in human cognition. Multiple theories have been proposed to formulate an overall cognitive architecture, which is a broad domain-generic cognitive computation model that captures the essential structure and process of the mind. Some of these theories hypothesize that, instead of a single standalone module, multiple modules in the brain share information to excel at a particular task. CLARION (Connectionist learning with rule induction online) (Sun & Franklin, 2007) is one such theory that postulates an integrative cognitive architecture, consisting of a number of distinct subsystems. It predicates a dual representational structure (Chaiken & Trope, 1999), where the top level encodes conscious explicit knowledge, while the other encodes indirect implicit information. The two systems interact, share knowledge, and cooperate to solve tasks. Delving into these underlying architectures and formulating a new design can help in the quest to build intelligent agents.

Multiple modules can be instituted instead of a single feedforward network. Explicit module that learns from the standard visual input and an implicit module that shares indirect contextual knowledge. The implicit module can be further divided into more submodules, each providing different information. Inductive biases

<sup>1</sup>Code and the *DN4IL* dataset will be made public upon acceptance.

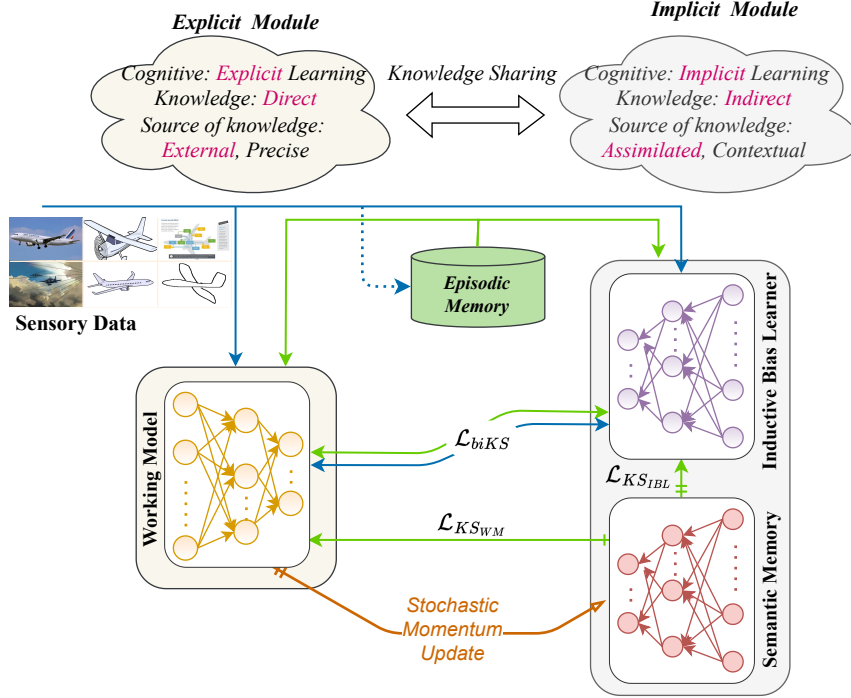


Figure 1: Schematic of *Dual Cognitive Architecture (DUCA)*. The working model in the explicit module learns direct sensory data. Within the implicit module, the inductive bias learner encodes the prior shape knowledge and the semantic memory consolidates information from the explicit module. Only one network (semantic memory) is used during inference as it includes consolidated knowledge across all tasks.

and semantic memories can act as different kinds of implicit knowledge. Inductive biases are pre-stored templates or knowledge that provide some meaningful disposition toward adapting to the continuously evolving world (Chollet, 2019). Theories postulate that after rapidly learning information, a gradual consolidation of knowledge transpires in the brain for slow learning of structured information (Kumaran et al., 2016). Thus, the new design incorporates multiple concepts of cognition architectures, the dichotomy of implicit and explicit representations, inductive biases, and multi-memory systems theory.

To this end, we propose *Dual Cognitive Architecture (DUCA)*, a multi-module architecture for CL. The explicit working module processes the standard input data. Two different submodules are introduced for the implicit module. The inductive bias learner embeds relevant prior information, and as networks are shown to be biased toward textural information (unlike humans that are more biased toward global semantics) (Geirhos et al., 2018), we propose to utilize global shape information as the prior. Both texture and shape are present in the original image, but ANNs tend to rely more on texture and ignore semantic information. Hence, we utilize the implicit shape information and share it with the explicit module to learn more generic and high-level representations. Further, to emulate the consolidation of information in the slow-fast multi-memory system, a gradual accumulation of knowledge from the explicit working module is embedded in the second semantic memory submodule. We show that interacting and leveraging information between these modules can help alleviate catastrophic forgetting, while also increasing the robustness to the distribution shift.

DUCA achieves superior performance across all CL settings on various datasets. DUCA outperforms SOTA CL methods on Seq-CIFAR10 and Seq-CIFAR100 in class incremental settings. Furthermore, in more realistic general class incremental settings where the task boundary is blurry and the classes are not disjoint, DUCA shows significant gains. The addition of inductive bias and semantic memory helps to achieve a better balance between the plasticity-stability trade-off. The prior in the form of shape helps produce generic representations, and this results in DUCA exhibiting a reduced task-recency bias. Furthermore, DUCA also

shows greater robustness against natural corruptions. Finally, to test the capability of CL methods against the distribution shift, we introduce a domain incremental learning dataset, *DN4IL*, which is a carefully designed subset of the DomainNet dataset (Peng et al., 2019). DUCA shows considerable robustness across all domains on these challenging data, thus establishing the efficacy of our cognitive-inspired CL architecture. Our contributions are as follows:

- *Dual Cognitive Architecture (DUCA)*, a novel method that incorporates aspects of cognitive architectures, multi-memory systems, and inductive bias into the CL framework.
- Introducing *DN4IL*, a challenging domain incremental learning dataset.
- Benchmark across different CL settings: class-, task-, generalized class-, and domain-incremental learning.
- Analyses on the plasticity-stability trade-off, task recency bias, and robustness to natural corruptions.

## 2 Methodology

### 2.1 Cognitive Architectures

Cognitive architectures refer to computational models that encapsulate the overall structure of the cognitive process in the brain. The underlying infrastructure of such a model can be leveraged to develop better intelligent systems. Global workspace theory (GWT) (Juliani et al., 2022) postulates that human cognition is composed of a multitude of special-purpose processors and is not a single standalone module. Different sub-modules might encode different contextual information which, when activated, can transfer knowledge to the conscious central workspace to influence and help make better decisions. Furthermore, CLARION (Sun & Franklin, 2007) posits a dual-system cognitive architecture with two levels of knowledge representation. The explicit module encodes direct knowledge that is externally accessible. The implicit module encodes indirect knowledge that is not directly accessible, but can be obtained through some intermediate interpretive or transformational steps. These two modules interact with each other by transferring knowledge between each other.

Inspired by these theories, we formulate a method that incorporates some of the key aspects of cognitive architecture into the CL method. A working module, which encodes the direct sensory data, forms the explicit module. A second module that encodes indirect and interpretive information forms the implicit module. The implicit module further includes multiple sub-modules to encode different types of knowledge.

### 2.2 Inductive Bias

The sub-modules in the implicit module need to encapsulate implicit information that can provide more contextual and high-level supervision. One of such knowledge can be prior knowledge or inductive bias. Inductive biases are pre-stored templates that exist implicitly even in earlier stages of the human brain (Pearl & Mackenzie, 2018). For instance, cognitive inductive bias may be one of the reasons why humans can focus on the global semantics of objects to make predictions. ANNs, on the other hand, are more prone to rely on local cues and textures (Geirhos et al., 2018). Global semantics or shape information already exists in the visual data, but in an indirect way. Hence, we utilize shape as indirect information in the implicit module. The sub-module uses a transformation step to extract the shape and share this inductive bias with the working module. As the standard (RGB) image and its shape counterpart can be viewed as different perspectives/modalities of the same data, ensuring that the representation of one modality is consistent with the other increases robustness to spurious correlations that might exist in only one of them.

### 2.3 Multi Memory System

Moreover, many theories have postulated that an intelligent agent must possess differentially specialized learning memory systems (Kumaran et al., 2016). While one system rapidly learns the individual experience, the other gradually assimilate the knowledge. To emulate this behavior, we establish a second sub-module that slowly consolidates the knowledge from the working module.

## 2.4 Formulation

To this end, we propose a novel method *Cognitive Continual Learner (CCL)*, which incorporates all these concepts into the CL paradigm. DUCA consists of two modules, the explicit module, and the implicit module. The explicit module has a single working model and processes the incoming direct visual data. The implicit module further consists of two submodules, namely the inductive bias learner and the semantic memory. They share relevant contextual information and assimilated knowledge with the explicit module, respectively. Figure 1 shows the overall architecture.

In the implicit module, semantic memory  $N_{SM}$ , consolidates knowledge at stochastic intervals from the working model  $N_{WM}$ , in the explicit module. The other submodule, the inductive bias learner  $N_{IBL}$ , processes the data and extracts shape information (Section F).  $N_{WM}$  processes the RGB data,  $N_{SM}$  consolidates the information from the working module at an update frequency in a stochastic manner, and  $N_{IBL}$  learns from the shape data.  $f$  represents the combination of the encoder and the classifier, and  $\theta_{WM}$ ,  $\theta_{SM}$ , and  $\theta_{IBL}$  are the parameters of the three networks.

A CL classification consists of a sequence of  $T$  tasks and, during each task  $t \in 1, 2 \dots T$ , samples  $x_c$  and their corresponding labels  $y_c$  are drawn from the current task data  $D_t$ . Furthermore, for each subsequent task, a random batch of exemplars is sampled from episodic memory  $B$  as  $x_b$ . An inductive bias (shape) filter is applied to generate shape samples,  $x_{cs} = \mathbb{IB}(x_c)$  and  $x_{bs} = \mathbb{IB}(x_b)$ . Reservoir sampling (Vitter, 1985) is incorporated to replay previous samples. Each of the networks  $N_{WM}$  and  $N_{IBL}$  learns in its own modality with a supervised cross-entropy loss on both the current samples and the buffer samples:

$$\mathcal{L}_{Sup_{WM}} = \mathcal{L}_{CE}(f(x_c; \theta_{WM}), y_c) + \mathcal{L}_{CE}(f(x_b; \theta_{WM}), y_b) \quad (1)$$

$$\mathcal{L}_{Sup_{IBL}} = \mathcal{L}_{CE}(f(x_{cs}; \theta_{IBL}), y_c) + \mathcal{L}_{CE}(f(x_{bs}; \theta_{IBL}), y_b) \quad (2)$$

The Knowledge Sharing (KS) objectives are designed to transfer and share information between all modules. KS occurs for current samples and buffered samples. We employ the mean squared error as the objective function for all KS losses. To provide shape supervision to the working model and vice versa, a bidirectional decision space similarity constraint ( $\mathcal{L}_{biKS}$ ) is enforced to align the features of the two modules.

$$\mathcal{L}_{biKS} = \mathbb{E}_{x \sim D_t \cup B} \|f(x_s; \theta_{IBL}) - f(x; \theta_{WM})\|_2^2 \quad (3)$$

The consolidated structural information in semantic memory is transferred to both the working model and the inductive bias learner by aligning the output space on the buffer samples, which further helps in information retention. The loss functions  $\mathcal{L}_{KS_{WM}}$  and  $\mathcal{L}_{KS_{IBL}}$  are as follows;

$$\begin{aligned} \mathcal{L}_{KS_{WM}} &= \mathbb{E}_{x_b \sim B} \|f(x_b; \theta_{SM}) - f(x_b; \theta_{WM})\|_2^2 \\ \mathcal{L}_{KS_{IBL}} &= \mathbb{E}_{x_b \sim B} \|f(x_b; \theta_{SM}) - f(x_{bs}; \theta_{IBL})\|_2^2 \end{aligned} \quad (4)$$

Thus, the overall loss functions for the working model and the inductive bias learner are as follows;

$$\begin{aligned} \mathcal{L}_{WM} &= \mathcal{L}_{Sup_{WM}} + \lambda(\mathcal{L}_{biKS} + \mathcal{L}_{KS_{WM}}) \\ \mathcal{L}_{IBL} &= \mathcal{L}_{Sup_{IBL}} + \gamma(\mathcal{L}_{biKS} + \mathcal{L}_{KS_{IBL}}) \end{aligned} \quad (5)$$

The semantic memory of the implicit module is updated with a stochastic momentum update (SMU) of the weights of the working model at rate  $r$  with a decay factor of  $d$ ,

$$\theta_{SM} = d \cdot \theta_{SM} + (1 - d) \cdot \theta_{WM} \text{ if } s \sim U(0, 1) < r \quad (6)$$

More details are provided in Algorithm 1. We discuss the computational aspect in Section E. Note that we use semantic memory ( $\theta_{SM}$ ) for inference, as it contains consolidated knowledge across all tasks.

Table 1: Comparison of different methods on standard CL benchmarks (Class-IL, Task-IL and GCIL settings). DUCA shows a consistent improvement over all methods for both buffer sizes.

$ \mathcal{B} $	METHOD	SEQ-CIFAR10		SEQ-CIFAR100		GCIL-CIFAR100	
		CLASS-IL	TASK-IL	CLASS-IL	TASK-IL	UNIFORM	LONGTAIL
-	JOINT	92.20 $\pm$ 0.15	98.31 $\pm$ 0.12	70.62 $\pm$ 0.64	86.19 $\pm$ 0.43	60.45 $\pm$ 1.65	60.10 $\pm$ 0.42
	SGD	19.62 $\pm$ 0.05	61.02 $\pm$ 3.33	17.58 $\pm$ 0.04	40.46 $\pm$ 0.99	10.36 $\pm$ 0.13	9.62 $\pm$ 0.21
200	ER	44.79 $\pm$ 1.86	91.19 $\pm$ 0.94	21.40 $\pm$ 0.22	61.36 $\pm$ 0.39	16.52 $\pm$ 0.10	16.20 $\pm$ 0.30
	DER++	64.88 $\pm$ 1.17	91.92 $\pm$ 0.60	29.60 $\pm$ 1.14	62.49 $\pm$ 0.78	27.73 $\pm$ 0.93	26.48 $\pm$ 2.04
	Co <sup>2</sup> L	65.57 $\pm$ 1.37	93.43 $\pm$ 0.78	31.90 $\pm$ 0.38	55.02 $\pm$ 0.36	-	-
	ER-ACE	62.08 $\pm$ 1.44	92.20 $\pm$ 0.57	32.49 $\pm$ 0.95	59.77 $\pm$ 0.31	27.64 $\pm$ 0.76	25.10 $\pm$ 2.64
	CLS-ER <sup>†</sup>	66.19 $\pm$ 0.75	93.90 $\pm$ 0.60	43.80 $\pm$ 1.89	73.49 $\pm$ 1.04	35.88 $\pm$ 0.41	35.67 $\pm$ 0.72
	DUCA	<b>70.04<math>\pm</math>1.07</b>	<b>94.49<math>\pm</math>0.38</b>	<b>45.38<math>\pm</math>1.28</b>	<b>76.62<math>\pm</math>0.16</b>	<b>38.61<math>\pm</math>0.83</b>	<b>37.11<math>\pm</math>0.16</b>
500	ER	57.74 $\pm$ 0.27	93.61 $\pm$ 0.27	28.02 $\pm$ 0.31	68.23 $\pm$ 0.16	23.62 $\pm$ 0.66	22.36 $\pm$ 1.27
	DER++	72.70 $\pm$ 1.36	93.88 $\pm$ 0.50	41.40 $\pm$ 0.96	70.61 $\pm$ 0.11	35.80 $\pm$ 0.62	34.23 $\pm$ 1.19
	Co <sup>2</sup> L	74.26 $\pm$ 0.77	95.90 $\pm$ 0.26	39.21 $\pm$ 0.39	62.98 $\pm$ 0.58	-	-
	ER-ACE	68.45 $\pm$ 1.78	93.47 $\pm$ 1.00	40.67 $\pm$ 0.06	66.45 $\pm$ 0.71	30.14 $\pm$ 1.11	31.88 $\pm$ 0.73
	CLS-ER	75.22 $\pm$ 0.71	94.94 $\pm$ 0.53	51.40 $\pm$ 1.00	78.12 $\pm$ 0.24	38.94 $\pm$ 0.38	38.79 $\pm$ 0.67
	DUCA	<b>76.20<math>\pm</math>0.70</b>	<b>95.95<math>\pm</math>0.14</b>	<b>54.27<math>\pm</math>1.09</b>	<b>79.80<math>\pm</math>0.32</b>	<b>43.34<math>\pm</math>0.32</b>	<b>41.44<math>\pm</math>0.22</b>

### 3 Experimental Settings

ResNet-18 (He et al., 2016) architecture is used for all experiments. All networks are trained using the SGD optimizer with standard augmentations of random crop and random horizontal flip. The different hyperparameters, tuned per dataset, are provided in D. The different CL settings are explained in detail in Section C. We consider CClass-IL, Domain-IL, and also report the Task-IL settings. Seq-CIFAR10 and Seq-CIFAR100 (Krizhevsky et al., 2009) for the class incremental learning (Class-IL) settings, which are divided into 5 tasks each. In addition to Class-IL, we also consider and evaluate General Class-IL (GCIL) (Mi et al., 2020) on the CIFAR100 dataset. For the domain incremental learning (Domain-IL), we propose a novel dataset, *DN4IL*.

### 4 Results

We provide a comparison of our method with standard baselines and multiple other SOTA CL methods. The lower and upper bounds are reported as SGD (standard training) and JOINT (training all tasks together), respectively. We compare with other rehearsal-based methods in the literature, namely ER, DER++ (Buzzega et al., 2020), Co<sup>2</sup>L (Cha et al., 2021), ER-ACE (Caccia et al., 2021), and CLS-ER (Arani et al., 2021). Table 1 shows the average performance in different settings over three seeds. Co<sup>2</sup>L utilizes task boundary information, and therefore the GCIL setting is not applicable. The results are taken from the original works and, if not available, using the original codes, we conducted a hyperparameter search for the new settings (see Section D for details).

DUCA achieves the best performance across all datasets in all settings. In the challenging Class-IL setting, we observe a gain of  $\sim 50\%$  over DER++, thus showing the efficacy of adding multiple modules for CL. Furthermore, we report improvements of  $\sim 6\%$  on both the Seq-CIFAR10 and Seq-CIFAR100 datasets, over CLS-ER, which utilizes two semantic memories in its design. DUCA has a single semantic memory, and the additional boost is obtained by prior knowledge from the inductive bias learner. Improvement is prominent even when the memory budget is low (200 buffer size). GCIL represents a more realistic setting, as the task boundaries are blurry, and classes can reappear and overlap in any task. GCIL-Longtail version also introduces an imbalance in the sample distribution. DUCA shows a significant improvement on both versions of GCIL-CIFAR100. Additional results are provided in Table 3.

Shape information from the inductive bias learner offers the global high-level context, which helps in producing generic representations that are not biased towards learning only the current task at hand. Furthermore, sharing of the knowledge that has been assimilated through the appearance of overlapping classes through the training scheme, further facilitates learning in this general setting. The overall results indicate that the dual knowledge sharing between the explicit working module and the implicit inductive bias and semantic memory modules enables both better adaptation to new tasks and information retention.

## 5 Domain incremental learning

Intelligent agents deployed in real-world applications need to maintain consistent performance through changes in the data and environment. Domain-IL aims to assess the robustness of the CL methods to the distribution shift. In Domain-IL, the classes in each task remain the same, but the input distribution changes, and this makes for a more plausible use case for evaluation. However, the datasets used in the literature do not fully reflect this setting. For instance, the most common datasets used in the literature are different variations (Rotated and Permuted) of the MNIST dataset (LeCun et al., 1998). MNIST is a simple dataset, usually evaluated on MLP networks, and its variations do not reflect the real-world distribution shift challenges that a CL method faces. As is evident from the different CL methods in the literature, the improvement in performance has been saturated on all variants of MNIST. Farquhar & Gal (2018) propose fundamental desiderata for CL evaluations and datasets based on real-world use cases. One of the criteria is to possess cross-task resemblances, which Permuted-MNIST clearly violates. Thus, a different dataset is needed to test the overall capability of a CL method to handle the distributional shift.

### 5.1 DN4IL Dataset

To this end, we propose *DN4IL* (DomainNet for Domain-IL), which is a well-crafted subset of the standard DomainNet dataset (Peng et al., 2019), used in domain adaptation. DomainNet consists of common objects in six different domains - real, clipart, infograph, painting, quickdraw, and sketch. The original DomainNet consists of 59k samples with 345 classes in each domain. The classes have redundancy, and moreover, evaluating the whole dataset can be computationally expensive in a CL setting. *DN4IL* version considers different criteria such as relevance of classes, uniform sample distribution, computational complexity, and ease of benchmarking for CL.

All classes were grouped into semantically similar supercategories. Of these, a subset of classes was selected that had relevance to domain shift, while also having maximum overlap with other standard datasets such as CIFAR, to facilitate out-of-distribution analyses. 20 supercategories were chosen with 5 classes each (resulting in a total of 100 classes). In addition, to provide a balanced dataset, we performed a class-wise sampling. First, we sample images per class in each supercategory and maintain class balance. Second, we choose samples per domain, so that it results in a dataset that has a near-uniform distribution across all classes and domains. The final dataset *DN4IL* is succinct, more balanced, and more computationally efficient for benchmarking, thus facilitating research in CL. Furthermore, the new dataset is deemed more plausible for real-world settings and also adheres to all evaluation desiderata by (Farquhar & Gal, 2018). The challenging distribution shift between domains provides an apt dataset to test the capability of CL methods in the Domain-IL setting. More details, statistics, and visual examples of this crafted dataset are provided in Section G.

### 5.2 DN4IL Performance

Figure 2 (left) reports the results on *DN4IL* for two different buffer sizes (values are provided in Table 6). DUCA shows a considerable performance gain in the average accuracy across all domains, and this can be primarily attributed to the supervision from the shape data. Standard networks tend to exhibit texture bias and learn background or spurious cues (Geirhos et al., 2018) that result in performance degradation when the distribution changes. Learning global shape information of objects, on the other hand, helps in learning generic features that can translate well to other distributions. Semantic memory further helps to consolidate information across domains. Maintaining consistent performance to such difficult distribution

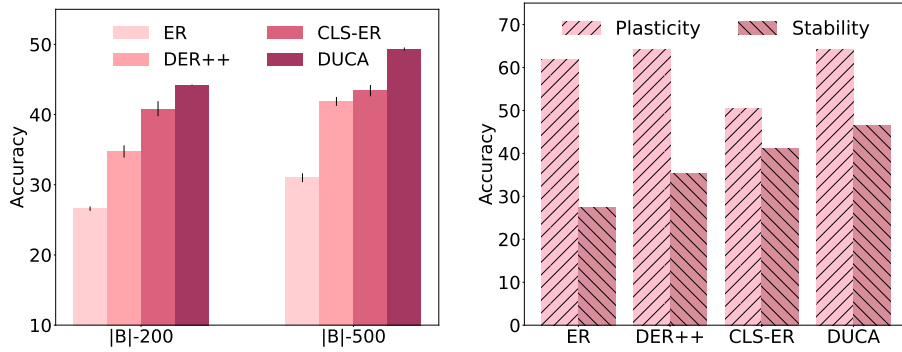


Figure 2: Accuracy (left) and plasticity-stability analysis (right) on *DN4IL* dataset. DUCA substantially outperforms other methods and with better plasticity-stability trade-off.

shifts prove beneficial in real-world applications, and the proficiency of DUCA in this setting can thus open up new avenues for research in cognition-inspired multi-module architectures.

## 6 Analyses

### 6.1 Plasticity-Stability Trade-off

Plasticity refers to the capability of a model to learn new tasks, while stability shows how well it can retain old information. The plasticity-stability dilemma is a long-standing problem in CL, which requires an optimal balance between the two. We measure each of these to assess the competence of the CL methods. Plasticity is computed as the average performance of each task when first learned (e.g., the accuracy of the network trained on task  $T_2$ , evaluated on the test set of  $T_2$ ). Stability is computed as the average performance of all tasks  $1:T-1$ , after learning the final task  $T$ . Figure 2 (right) reports these numbers for the *DN4IL* dataset. As seen, the ER and DER methods exhibit forgetting and show low stability, and focus only on the newer tasks. CLS-ER shows greater stability, but at the cost of reduced plasticity. However, DUCA shows the highest stability while maintaining comparable plasticity. The shape knowledge helps in learning generic solutions that can translate to new tasks, while the semantic consolidation update at stochastic rates acts as a regularization to maintain stable parameter updates. Thus, DUCA strikes a better balance between plasticity and stability.

### 6.2 Task-wise Performance

The average accuracy across all tasks does not provide a complete measure of the ability of a network to retain old information while learning new tasks. To better represent the plasticity-stability measure, we report the task-wise performance at the end of each task. After training each task, we measure the accuracy on the test set of each of the previous tasks. Figure 3 reports this for all tasks of *DN4IL*. The last row represents the performance of each task after the training is completed. ER and DER++ show performance degradation on earlier tasks, as the model continues to train on newer tasks. Both perform well on the last task and display the lowest stability. DUCA reports the highest information retention on older tasks, while also maintaining plasticity. For example, the accuracy on the first task (real) reduces to 27.6 on ER after training the six tasks (domains), while the DUCA maintains the accuracy of 54.9. CLS-ER shows better retention of old information, but at the cost of plasticity. The last task in CLS-ER shows a lower performance compared to DUCA (52.1 vs. 61.0). The performance of the current task in DUCA is relatively lesser and can be attributed to the stochastic update rate.

To shed more light on the performance of each of the modules in DUCA, we also provide the performance of the working model and the inductive bias learner, in Appendix Figure 5. The working model shows better plasticity, while DUCA (semantic memory) displays better stability. Overall, all modules in the proposed approach present unique attributes that improve the learning process and improve performance.

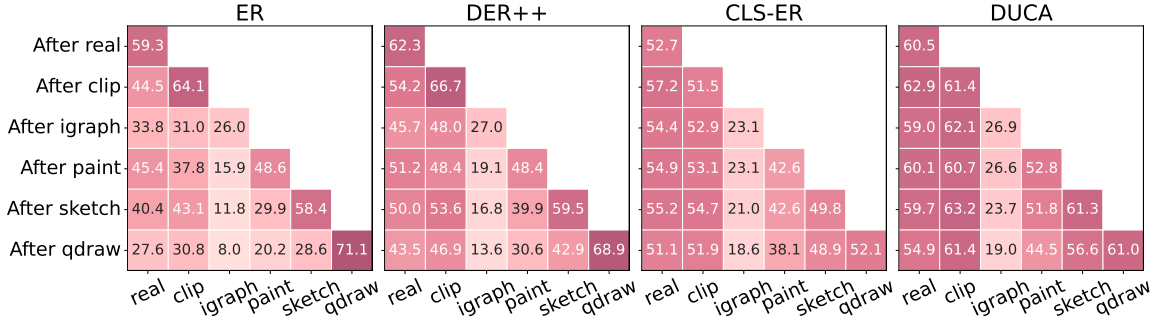


Figure 3: Task-wise performance on  $DN4IL$  ( $|\mathcal{B}|=500$ ), where each task represents a domain. DUCA shows more retention of old information without compromising much on current accuracy.

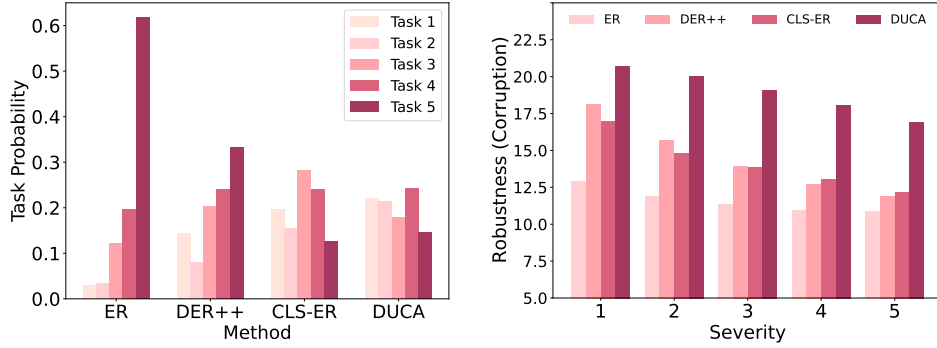


Figure 4: DUCA shows reduced task recency bias (left), as well as higher robustness against natural corruption (right) on Seq-CIFAR10 ( $|\mathcal{B}|=200$ ) dataset.

### 6.3 Recency-Bias Analysis

Recency bias is a behavior in which the model predictions tend to be biased toward the current or the most recent task (Wu et al., 2019). This is undesirable in a CL model, as it results in a biased solution that forgets the old tasks. To this end, after the end of the training, we evaluate the models on the test set (of all tasks) and calculate the probability of predicting each task. The output distribution for each test sample is calculated for all classes and the probabilities are averaged per task.

Figure 4 (left) shows the probabilities for each task on the Seq-CIFAR10 dataset. As shown, the ER and DER++ methods tend to incline most of their predictions toward the classes seen in the last task, thus creating a misguided bias. DUCA shows a lower bias compared to both of these baselines. CLS-ER exhibits reduced bias due to the presence of multiple memories, but the distribution is still relatively skewed (with respect to a probability of 0.2). DUCA shows a more uniform distribution across all tasks. The dual information from the shape data and the consolidated knowledge across tasks helps in breaking away from Occam’s razor pattern of neural networks to default to the easiest solution.

### 6.4 Robustness

Lifelong agents, when deployed in real-world settings, must be resistant to various factors, such as lighting conditions, weather changes, and other effects of digital imaging. Inconsistency in predictions under different conditions might result in undesirable outcomes, especially in safety-critical applications such as autonomous driving. To measure the robustness of the CL method against such natural corruptions, we created a dataset by applying fifteen different corruptions, at varying levels of severity (1- least severe to 5- most severe corruption).



Table 2: Ablation to analyze the effect of each component of DUCA on Seq-CIFAR10 and *DN4IL*.

SM	IBL	KS (WM $\leftrightarrow$ IBL)	Seq-CIFAR10	DN4IL
✓	✓	✓	<b>70.04</b> $\pm$ 1.07	<b>44.23</b> $\pm$ 0.05
✓	✓	✗	69.28 $\pm$ 1.34	40.35 $\pm$ 0.34
✓	✗	-	69.21 $\pm$ 1.46	39.76 $\pm$ 0.56
✗	✓	✓	64.61 $\pm$ 1.22	37.33 $\pm$ 0.01
✗	✗	✗	44.79 $\pm$ 1.86	26.59 $\pm$ 0.31

The performances on the fifteen corruptions are averaged at each severity level and are shown in Figure 4 (right). DUCA outperforms all other techniques at all severity levels. ER, DER++, and CLS-ER show a fast decline in accuracy as severity increases, while DUCA maintains stable performance throughout. Implicit shape information provides a different perspective of the same data to the model, which helps to generate high-level robust representations. DUCA, along with improved continual learning performance, also exhibits improved robustness to corruption, thus proving to be a better candidate for deployment in real-world applications.

## 6.5 Ablation Study

DUCA architecture comprises multiple components, each contributing to the efficacy of the method. The explicit module has the working model, and the implicit module has semantic memory (SM) and inductive bias learner (IBL). Disentangling different components in the DUCA can provide more insight into the contribution of each of them to the overall performance.

Table 2 reports the ablation study with respect to each of these components on both the Seq-CIFAR10 and *DN4IL* datasets. Considering the more complex *DN4IL* dataset, the ER accuracy without any of our components is 26.59. Adding cognitive bias (IBL) improves performance by 40%. Shape information plays a prominent role, as networks need to learn the global semantics of the objects, rather than background or spurious textural information, to translate performance across domains. Adding the dual-memory component (SM) shows an increase of approximately 49% over the vanilla baseline. Furthermore, the KS between explicit and implicit modules on current experiences also plays a key role in performance gain. Combining both of these cognitive components and, in general, following the multi-module design shows a gain of 66%. A similar trend is seen on Seq-CIFAR10.

## 7 Related Works

Rehearsal-based approaches, which revisit examples from the past to alleviate catastrophic forgetting, have been effective in challenging CL scenarios (Farquhar & Gal, 2018). Experience Replay (ER) (Riemer et al., 2018) methods use episodic memory to retain previously seen samples for replay purposes. DER++ (Buzzega et al., 2020) adds a consistency loss on logits, in addition to the ER strategy. CO<sup>2</sup>L (Cha et al., 2021) uses contrastive learning from the self-supervised learning domain to generate transferable representations. ER-ACE (Caccia et al., 2021) targets the representation drift problem in online CL and develops a technique to use separate losses for current and buffer samples. All of these methods limit the architecture to a single stand-alone network, contrary to the biological workings of the brain.

CLS-ER (Arani et al., 2021) proposed a multi-network approach that emulates fast and slow learning systems by using two semantic memories, each aggregating weights at different times. Though CLS-ER utilizes the multi-memory design, sharing of different kinds of knowledge is not leveraged, and hence presents a method with limited scope. DUCA digresses from the standard architectures and proposed a multi-module design that is inspired by the cognitive computational architectures. It incorporates multiple submodules, each sharing different knowledge to develop an effective continual learner that has better generalization and robustness.

## 8 Conclusion

We introduced a novel framework for continual learning which incorporates concepts inspired by cognitive architectures, high-level cognitive biases, and the multi-memory system. *Dual Cognitive Architecture (DUCA)*, includes multiple subsystems with dual knowledge representation. DUCA designed a dichotomy of explicit and implicit modules in which information is selected, maintained, and shared with each other to enable better generalization and robustness. DUCA outperformed on Seq-CIFAR10 and Seq-CIFAR100 on the Class-IL setting. In addition, it also showed a significant gain in the more realistic and challenging GCIL setting. Through different analyses, we showed a better plasticity-stability balance.

Furthermore, shape prior and knowledge consolidation helps to learn more generic solutions, indicated by the reduced problem of task recency bias and greater robustness against natural corruptions. Furthermore, we introduced a challenging domain-IL dataset, *DN4IL*, with six disparate domains. The significant improvement of DUCA on this complex distribution shift demonstrates the benefits of shape context, which helps the network to converge on a generic solution, rather than a simple texture-biased one. The objective of this work was to develop a framework that incorporates elements of cognitive architecture to mitigate the forgetting problem and enhance generalization and robustness. We hope our preliminary work will serve as a foundation for future research in different cognitive biases, and more efficient designs, and pave the way for the advancement of lifelong learning methods for ANNs.

## References

- Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual learning method based on complementary learning system. In *International Conference on Learning Representations*, 2021.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky. New insights on reducing abrupt representation change in online continual learning. In *International Conference on Learning Representations*, 2021.
- Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9516–9525, 2021.
- Shelly Chaiken and Yaacov Trope. *Dual-process theories in social psychology*. Guilford Press, 1999.
- François Chollet. On the measure of intelligence. *arXiv preprint arXiv:1911.01547*, 2019.
- Lijun Ding and Ardeshir Goshtasby. On the canny edge detector. *Pattern Recognition*, 34(3):721–725, 2001.
- Sebastian Farquhar and Yarin Gal. Towards robust evaluations of continual learning. *arXiv preprint arXiv:1805.09733*, 2018.
- Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Arthur Juliani, Kai Arulkumaran, Shuntaro Sasai, and Ryota Kanai. On the link between conscious function and general intelligence in humans and machines. *arXiv preprint arXiv:2204.05133*, 2022.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

- Dharshan Kumaran, Demis Hassabis, and James L McClelland. What learning systems do intelligent agents need? complementary learning systems theory updated. *Trends in cognitive sciences*, 20(7):512–534, 2016.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
- Fei Mi, Lingjing Kong, Tao Lin, Kaicheng Yu, and Boi Faltings. Generalized class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 240–241, 2020.
- German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic books, 2018.
- Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 1406–1415, 2019.
- 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 *International Conference on Learning Representations*, 2018.
- Irwin Sobel and Gary Feldman. A 3x3 isotropic gradient operator for image processing. *a talk at the Stanford Artificial Project in*, pp. 271–272, 1968.
- Ron Sun and Stan Franklin. Computational models of consciousness: A taxonomy and some examples., 2007.
- Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- Jeffrey S Vitter. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)*, 11(1):37–57, 1985.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 374–382, 2019.

## A DUCA Algorithm

---

**Algorithm 1** Dual Cognitive Architecture (DUCA)

---

```

1: Input: Dataset  $\mathcal{D}_t$ , Buffer  $\mathcal{B}$ 
2: Initialize: Three networks: Encoder and classifier  $f$  parameterized by  $\theta_{WM}$ ,  $\theta_{SM}$ , and  $\theta_{IBL}$ 
3: for all tasks  $t \in 1, 2 \dots T$  do
4:   Sample mini-batch:  $(x_c, y_c) \sim \mathcal{D}_t$ 
5:   Extract shape images:  $x_{c_s} = \mathbb{I}\mathbb{B}(x_c)$  where  $\mathbb{I}\mathbb{B}$  is a Sobel filter
6:    $\mathcal{L}_{Sup_{WM}} = \mathcal{L}_{CE}(f(x_c; \theta_{WM}), y_c)$ 
7:    $\mathcal{L}_{Sup_{IBL}} = \mathcal{L}_{CE}(f(x_{c_s}; \theta_{IBL}), y_c)$ 
8:   if  $\mathcal{B} \neq \emptyset$  then
9:     Sample mini-batch:  $(x_b, y_b) \sim \mathcal{B}$ 
10:    Extract shape images:  $x_{b_s} = \mathbb{I}\mathbb{B}(x_b)$ 
11:    Calculate the supervised loss:
12:     $\mathcal{L}_{Sup_{WM}} += \mathcal{L}_{CE}(f(x_b; \theta_{WM}), y_b)$ 
13:     $\mathcal{L}_{Sup_{IBL}} += \mathcal{L}_{CE}(f(x_{b_s}; \theta_{IBL}), y_b)$ 
14:    Knowledge sharing from semantic memory to working model and inductive bias learner:
15:     $\mathcal{L}_{KS_{WM}} = \mathbb{E} \|f(x_b; \theta_{SM}) - f(x_b; \theta_{WM})\|_2^2$ 
16:     $\mathcal{L}_{KS_{IBL}} = \mathbb{E} \|f(x_b; \theta_{SM}) - f(x_{b_s}; \theta_{IBL})\|_2^2$ 
17:    Bidirectional knowledge sharing between working model and inductive bias learner:
18:     $\mathcal{L}_{biKS} = \mathbb{E}_{x \sim \mathcal{D}_t \cup \mathcal{B}} \|f(x; \theta_{WM}) - f(x_s; \theta_{IBL})\|_2^2$ 
19:    Calculate total loss:
20:     $\mathcal{L}_{WM} = \mathcal{L}_{Sup_{WM}} + \lambda(\mathcal{L}_{biKS} + \mathcal{L}_{KS_{WM}})$ 
21:     $\mathcal{L}_{IBL} = \mathcal{L}_{Sup_{IBL}} + \gamma(\mathcal{L}_{biKS} + \mathcal{L}_{KS_{IBL}})$ 
22:    Update both working model and inductive bias learner:  $\theta_{WM}, \theta_{IBL}$ 
23:    Stochastically update semantic memory:
24:    Sample  $s \sim U(0, 1)$ ;
25:    if  $s < r$  then
26:       $\theta_{SM} = d \cdot \theta_{SM} + (1 - d) \cdot \theta_{WM}$ 
27:    Update memory buffer  $\mathcal{B}$ 
28: Return: model  $\theta_{SM}$ 

```

---

## B Additional Results

Figure 5 presents the task-wise performance of all the three networks in the DUCA architecture, on *DN4IL* dataset. Semantic memory helps to retain information by maintaining high accuracy on older tasks and is more stable. The performance of the current task is relatively lower than that of the working model and could be due to the stochastic update rate of this model. The working model has better performance on new tasks and is more plastic. Inductive bias learner is evaluated on the transformed data (shape) and also achieves a balance between plasticity and stability. In general, all modules in our proposed method present unique attributes that improve the learning process by improving performance and reducing catastrophic forgetting. Table 3 reports the results of CIFAR100 on different numbers of tasks. Even when the number of tasks increases, our method consistently improves over the baselines.

Table 3: Comparison on Seq-CIFAR100 dataset for different number of tasks on 500 buffer size.

METHOD	5-TASKS	10-TASKS	20-TASKS
ER	28.02 $\pm$ 0.31	21.49 $\pm$ 0.47	16.52 $\pm$ 0.86
DER++	41.40 $\pm$ 0.96	36.20 $\pm$ 0.52	22.25 $\pm$ 5.87
DUCA	<b>53.23<math>\pm</math>1.62</b>	<b>41.09<math>\pm</math>0.72</b>	<b>33.60<math>\pm</math>0.25</b>

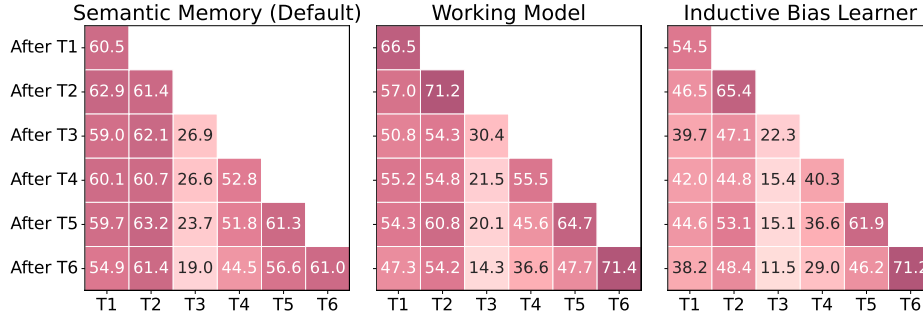


Figure 5: Task probability analysis of all DUCA components on *DN4IL* dataset with 500 buffer size. Semantic memory displays better stability while the working model displays better plasticity.

## C Setting and Datasets

We evaluate all methods in different CL settings. Van de Ven & Tolia (2019) describes three different settings based on increasing difficulty: task incremental learning (Task-IL), domain incremental learning (Domain-IL), and class incremental learning (Class-IL). In Class-IL, each new task consists of novel classes, and the network must learn both new classes while retaining information about the old ones. Task-IL is similar to Class-IL but assumes that task labels are accessible in both training and inference. In Domain-IL, the classes remain the same for each task, but the distribution varies for each task. We report the results for all three settings on the relevant datasets. Class-IL is the most complex setting of the three and is widely studied; however, there are some assumptions that simplify this setting to be realistic. Mi et al. (2020) highlighted some of the limitations of Class-IL, such as the assumption of the same number of classes across different tasks, the absence of reappearance of classes, and the sample distribution per class being well balanced. Hence, Generalized Class-IL (GCIL) was suggested to overcome these limitations and introduce a more realistic setting. GCIL is a more generalized CL setting, where the number of classes in each task is not fixed, and the classes can reappear with varying sample sizes. GCIL samples the number of classes and samples from a probabilistic distribution. The two variations are Uniform (fixed uniform sample distribution over all classes) and Longtail (with class imbalance). We report results on all three settings: Task-IL, Domain-IL, and Class-IL. Furthermore, we also consider the GCIL setting for one of the datasets as an additional evaluation setting. All reported results are averaged over three random seeds.

## D Hyperparameters

For the settings and datasets for which the results are not available in the original papers, using their original codes, we conducted a hyperparameter search for each of the new settings. To this end, we utilize a small validation set from the training set to tune the hyperparameters. For Seq-CIFAR10, the results are taken from the original articles (Buzzega et al., 2020; Cha et al., 2021; Caccia et al., 2021; Arani et al., 2021). For the other datasets (and for each buffer size), we ran a grid search over the hyperparameters reported in the paper. For Seq-CIFAR100 and GCIL-CIFAR100, we based the search range using Seq-CIFAR10 hyperparameters as a reference point.

DN4IL dataset is more complex compared to the CIFAR versions and includes images of larger sizes. Hence, we consider the Seq-TinyImagenet hyperparameters in the respective paper as a reference point for further tuning. The learning rate  $lr$ , the number of epochs, and the batch size are similar across the datasets. The ema update rate  $r$  is lower for more complex datasets, as shown in CLS-ER.  $r$  is chosen in the range of  $[0.01, 0.1]$  with a step size of 0.02 for CLS-ER and DUCA. The different hyperparameters chosen for the baselines, after tuning, are reported in Table 4.

The different hyperparameters chosen for DUCA are shown in Table 5. Hyperparameters:  $lr$ , batch size, number of epochs and decay factor are uniform across all datasets. The stochastic update rate is similar to those of CLS-ER. The hyperparameters are stable across settings and datasets and also complement each

Table 4: Selected hyperparameters for all baselines.

DATASET	$ \mathcal{B} $	METHOD	HYPERPARAMETERS
SEQ-CIFAR100	200	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.1, \beta=0.5$
		CO <sup>2</sup> L	$lr:0.5, \tau:0.5, \kappa:0.2, \kappa^*:0.01, e:100$
		ER-ACE	$lr=0.01$
		CLS-ER	$lr=0.1 \lambda=0.15, r_p=0.1, r_s=0.05, \alpha_p=0.999, \alpha_s=0.999$
	500	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.1, \beta=0.5$
		CO <sup>2</sup> L	$lr:0.5, \tau:0.5, \kappa:0.2, \kappa^*:0.01, e:100$
		ER-ACE	$lr=0.01$
		CLS-ER	$lr=0.1 \lambda=0.15, r_p=0.1, r_s=0.05, \alpha_p=0.999, \alpha_s=0.999$
GCIL-CIFAR100	200	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.5, \beta=0.1$
		ER-ACE	$lr=0.1$
		CLS-ER	$lr=0.1 \lambda=0.1, r_p=0.7, r_s=0.6, \alpha_p=0.999, \alpha_s=0.999$
	500	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.2, \beta=0.1$
		ER-ACE	$lr=0.1$
		CLS-ER	$lr=0.1 \lambda=0.1, r_p=0.7, r_s=0.6, \alpha_p=0.999, \alpha_s=0.999$
DN4IL	200	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.1, \beta=1.0$
		CLS-ER	$lr=0.05 \lambda=0.1, r_p=0.08, r_s=0.04, \alpha_p=0.999, \alpha_s=0.999$
	500	ER	$lr=0.1$
		DER++	$lr=0.03, \alpha=0.5, \beta=0.1$
		CLS-ER	$lr=0.05 \lambda=0.1, r_p=0.08, r_s=0.05, \alpha_p=0.999, \alpha_s=0.999$

Table 5: Selected hyperparameters for DUCA across different settings. The learning rate is set to 0.03, batch size to 32, and epochs to 50 respectively for all the datasets. The decay factor  $d$  is always set to 0.999.

DATASET	$ \mathcal{B} $	$r$	$\lambda$	$\gamma$	DATASET	$ \mathcal{B} $	$r$	$\lambda$	$\gamma$
SEQ-CIFAR10	200	0.2	0.1	0.1	GCIL-CIFAR100	200	0.09	0.1	0.01
	500	0.2	0.1	0.1		500	0.09	0.1	0.01
SEQ-CIFAR100	200	0.1	0.1	0.01	DN4IL	200	0.06	0.1	0.01
	500	0.06	0.1	0.01		500	0.08	0.1	0.01

other. The loss balancing weights ( $\lambda$  and  $\gamma$  respectively) are also similar across all datasets. Therefore, DUCA does not require extensive fine-tuning across different datasets and settings.

## E Complexity

We discuss the computational complexity aspect of our proposed method. DUCA involves three networks during training; however, in inference, only one network is used (SM module). Therefore, for inference purposes, the MAC count, the number of parameters, and the computational capacity remain the same as in the other single-network methods.

The training cost requires three forward passes, as it consists of three different modules. ER, DER++, CO<sup>2</sup>L and ER-ACE have a single network. CLS-ER also has three networks and therefore requires 3 forward passes. DUCA has training complexity similar to that of CLS-ER; however, it outperforms CLS-ER in all provided

**Algorithm 2** Sobel Algorithm - Shape Extraction

- 
- 1: Up-sample the images to twice the original size:  $x_{rgb} = \text{us}(x_{rgb})$
  - 2: Reduce noisy edges:  $x_g = \text{Gaussian\_Blur}(x_{rgb}, \text{kernel\_size} = 3)$
  - 3: Get Sobel kernels:  $S_x = \begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix}$  and  $S_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix}$
  - 4: Apply Sobel kernels:  $x_{dx} = x_g * S_x$  and  $x_{dy} = x_g * S_y$  (\* is 2-D convolution operation)
  - 5: The edge magnitude:  $x_{shape} = \sqrt{x_{dx}^2 + x_{dy}^2}$
  - 6: Down-sample to original image size:  $x_{shape} = \text{ds}(x_{shape})$
- 

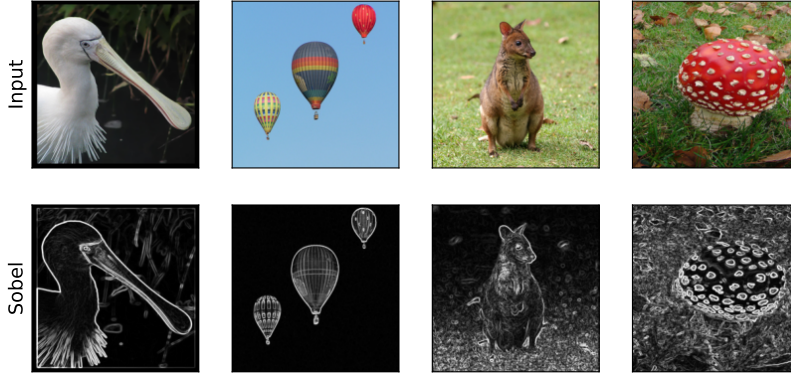


Figure 6: Visual examples of the shape images using Sobel operator

metrics. On the memory front, similar to all methods, we save memory samples based on the memory budget allotted (200 and 500 in the experiments). There are no additional memory requirements, as we do not save any extra information (such as logits in DER++) to be used later in our objectives. Additionally, there are no additional data requirements.

We considered the human brain, or cognition, as the most intelligent agent and wanted to incorporate some of the underlying workings into the neural network architecture. Therefore, our goal was to design a framework inspired by elements of cognitive architecture that reduces forgetting while exhibiting better generalization and robustness. Our work is a preliminary attempt to incorporate elements based on cognitive architecture into the CL algorithm to gauge the gain in reducing forgetting. We hope the promising results result in future work that explores different kinds of cognitive biases and interactions and moves towards more efficient designs.

## F Inductive Bias

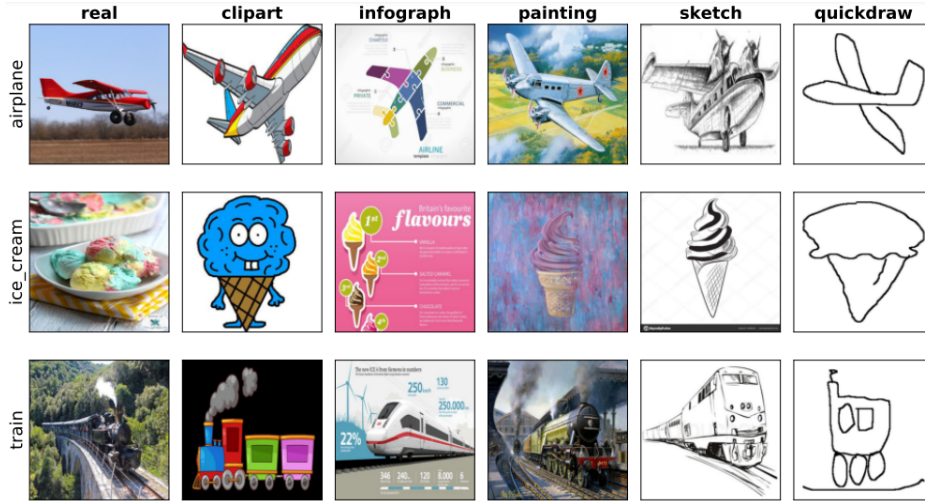
The shape extraction is performed by applying a filter on the input image. Multiple filters (such as Canny (Ding & Goshtasby, 2001), Prewitt) were considered, but the Sobel filter (Sobel & Feldman, 1968) was chosen because it produces a more realistic output by being precise and also smoothing the edges (see Algorithm 2). Figure 6 shows a few examples of applying the Sobel operator on the original RGB images. The Sobel output is fed to the IBL model.

## G DN4IL

We introduce a new dataset for the Domain-IL setting. It is a subset of the standard DomainNet dataset (Peng et al., 2019) used in domain adaptation. It consists of six different domains: real, clipart, infograph, painting, quickdraw, and sketch. The shift in distribution between domains is challenging. A few examples can be seen in Figure 7.

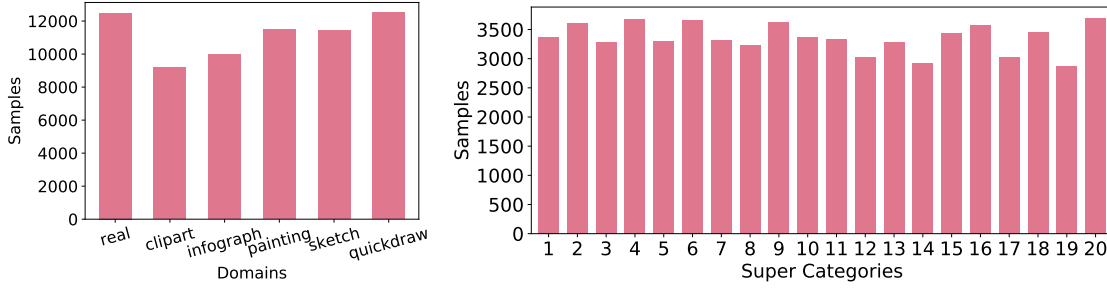
Table 6: Accuracy on the proposed *DN4IL* dataset for the Domain-IL setting. DUCA shows a significant improvement in all disparate and challenging domains.

$ \mathcal{B} $	METHOD	REAL	CLIPART	INFOGRAPH	PAINTING	SKETCH	QUICKDRAW	ACC
-	JOINT SGD	9.98 $\pm$ 0.54	19.97 $\pm$ 0.31	2.32 $\pm$ 0.20	6.58 $\pm$ 0.34	14.91 $\pm$ 0.04	71.23 $\pm$ 0.17	59.93 $\pm$ 1.07 20.83 $\pm$ 0.24
200	ER	20.08 $\pm$ 0.45	26.37 $\pm$ 0.35	5.56 $\pm$ 0.39	13.92 $\pm$ 0.91	23.69 $\pm$ 1.54	69.95 $\pm$ 0.56	26.59 $\pm$ 0.31
	DER++	33.66 $\pm$ 1.65	37.24 $\pm$ 0.64	9.80 $\pm$ 0.63	24.16 $\pm$ 1.17	34.37 $\pm$ 2.00	69.26 $\pm$ 0.79	34.75 $\pm$ 0.87
	CLS-ER	45.53 $\pm$ 0.88	49.17 $\pm$ 1.12	15.79 $\pm$ 0.48	35.80 $\pm$ 0.64	48.03 $\pm$ 0.85	54.40 $\pm$ 1.25	40.83 $\pm$ 1.07
	DUCA	47.52 $\pm$ 0.25	54.69 $\pm$ 0.10	15.70 $\pm$ 0.33	37.54 $\pm$ 0.30	51.98 $\pm$ 0.96	58.80 $\pm$ 0.18	<b>44.23</b> $\pm$ 0.05
500	ER	27.54 $\pm$ 0.05	31.89 $\pm$ 0.93	7.89 $\pm$ 0.45	19.39 $\pm$ 1.02	28.36 $\pm$ 1.35	70.96 $\pm$ 0.10	31.01 $\pm$ 0.62
	DER++	44.49 $\pm$ 1.39	46.17 $\pm$ 0.35	14.01 $\pm$ 0.23	33.44 $\pm$ 0.90	43.59 $\pm$ 1.11	69.53 $\pm$ 0.29	41.87 $\pm$ 0.63
	CLS-ER	49.85 $\pm$ 0.88	51.41 $\pm$ 0.34	18.17 $\pm$ 0.08	37.94 $\pm$ 0.94	49.02 $\pm$ 1.57	55.63 $\pm$ 0.71	43.41 $\pm$ 0.80
	DUCA	54.77 $\pm$ 0.15	60.37 $\pm$ 0.75	19.35 $\pm$ 0.39	44.50 $\pm$ 0.43	56.34 $\pm$ 0.53	60.61 $\pm$ 1.73	<b>49.32</b> $\pm$ 0.23

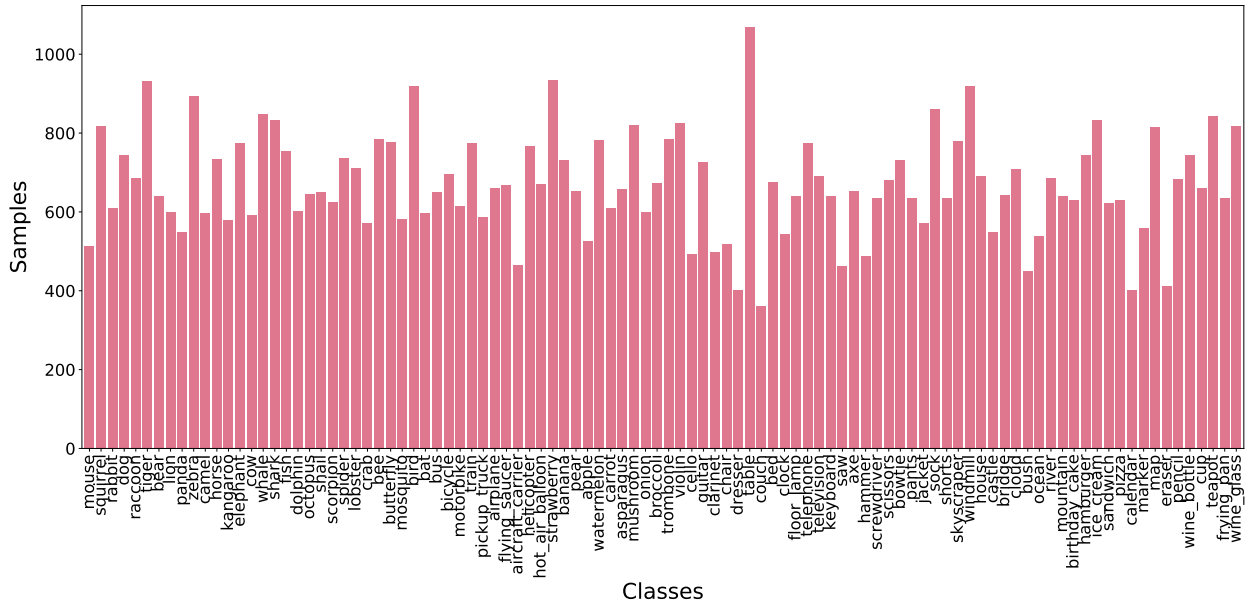
Figure 7: Visual examples of *DN4IL* dataset

Each domain includes 345 classes, and the overall dataset consists of  $\sim 59000$  samples. The classes have redundancy, and also evaluating on the whole dataset can be computationally expensive for CL settings. Therefore, we create a subset by grouping semantically similar classes into 20 supercategories (considering the class overlap between other standard datasets can also facilitate OOD analysis). Each super category has five classes each, which results in a total of 100 classes. The specifications of the classes are given in Table 7. The dataset consists of 67080 training images and 19464 test images. The image size for all experiments is chosen as  $64 \times 64$  (the normalize transform is not applied in augmentations).



Figure 8: Number of samples per domain and per supercategory in *DN4IL* dataset.Table 7: Details on supercategory and classes in *DN4IL* dataset.

SUPERCATEGORY	CLASS
1 SMALL ANIMALS	MOUSE SQUIRREL RABBIT DOG RACCOON
2 MEDIUM ANIMALS	TIGER BEAR LION PANDA ZEBRA
3 LARGE ANIMALS	CAMEL HORSE KANGAROO ELEPHANT COW
4 AQUATIC MAMMALS	WHALE SHARK FISH DOLPHIN OCTOPUS
5 NON-INSECT INVERTEBRATES	SNAIL SCORPION SPIDER LOBSTER CRAB
6 INSECTS	BEE BUTTERFLY MOSQUITO BIRD BAT
7 VEHICLE	BUS BICYCLE MOTORBIKE TRAIN PICKUP_TRUCK
8 SKY-VEHICLE	AIRPLANE FLYING_SAUCER AIRCRAFT_CARRIER HELICOPTER HOT_AIR_BALLOON
9 FRUITS	STRAWBERRY BANANA PEAR APPLE WATERMELON
10 VEGETABLES	CARROT ASPARAGUS MUSHROOM ONION BROCCOLI
11 MUSIC	TROMBONE VIOLIN CELLO GUITAR CLARINET
12 FURNITURE	CHAIR DRESSER TABLE COUCH BED
13 HOUSEHOLD ELECTRICAL DEVICES	CLOCK FLOOR_LAMP TELEPHONE TELEVISION KEYBOARD
14 TOOLS	SAW AXE HAMMER SCREWDRIVER SCISSORS
15 CLOTHES & ACCESSORIES	BOWTIE PANTS JACKET SOCK SHORTS
16 MAN-MADE OUTDOOR	SKYSCRAPER WINDMILL HOUSE CASTLE BRIDGE
17 NATURE	CLOUD BUSH OCEAN RIVER MOUNTAIN
18 FOOD	BIRTHDAY_CAKE HAMBURGER ICE_CREAM SANDWICH PIZZA
19 STATIONARY	CALENDAR MARKER MAP ERASER PENCIL
20 HOUSEHOLD ITEMS	WINE_BOTTLE CUP TEAPOT FRYING_PAN WINE_GLASS

Figure 9: Number of overall samples per class in *DN4IL* dataset.

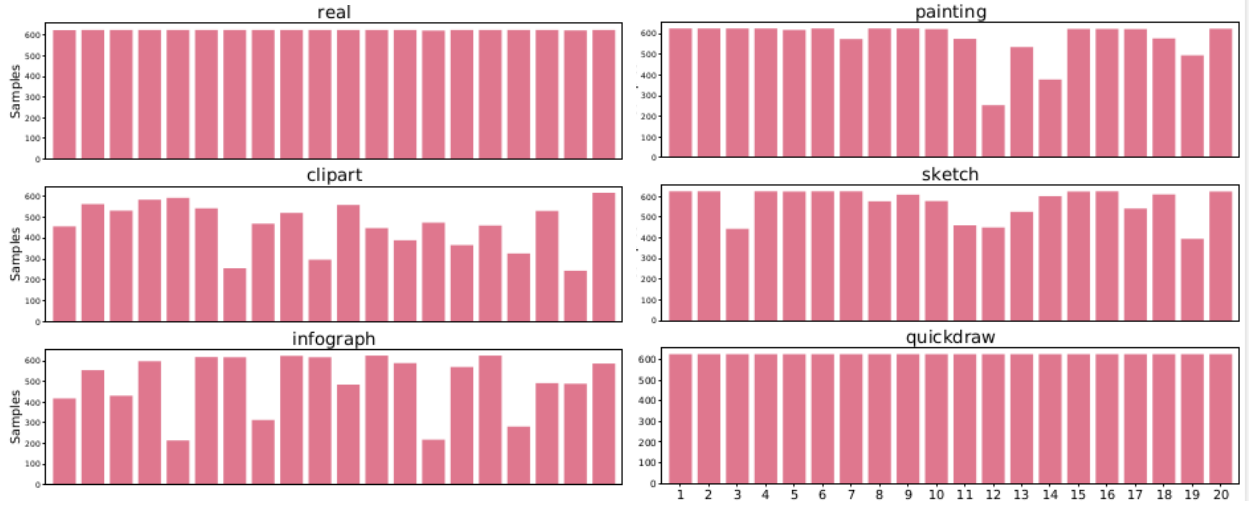


Figure 10: Number of samples per supercategory for each domain in  $DN_4IL$  dataset.