

Continual Knowledge Consolidation LORA for Domain Incremental Learning

Naeem Paeedeh, Mahardhika Pratama, *Senior member, IEEE*, Weiping Ding *Senior member, IEEE*, Jimmy Cao, Wolfgang Mayer, Ryszard Kowalczyk

Abstract—Domain Incremental Learning (DIL) is a continual learning sub-branch that aims to address never-ending arrivals of new domains without catastrophic forgetting problems. Despite the advent of parameter-efficient fine-tuning (PEFT) approaches, existing works create task-specific LoRAs overlooking shared knowledge across tasks. Inaccurate selection of task-specific LoRAs during inference results in significant drops in accuracy, while existing works rely on linear or prototype-based classifiers, which have suboptimal generalization powers. Our paper proposes continual knowledge consolidation low rank adaptation (CONEC-LoRA) addressing the DIL problems. CONEC-LoRA is developed from consolidations between task-shared LORA to extract common knowledge and task-specific LORA to embrace domain-specific knowledge. Unlike existing approaches, CONEC-LoRA integrates the concept of a stochastic classifier whose parameters are sampled from a distribution, thus enhancing the likelihood of correct classifications. Last but not least, an auxiliary network is deployed to optimally predict the task-specific LoRAs for inferences and implements the concept of a different-depth network structure in which every layer is connected with a local classifier to take advantage of intermediate representations. This module integrates the ball-generator loss and transformation module to address the synthetic sample bias problem. Our rigorous experiments demonstrate the advantage of CONEC-LoRA over prior arts in 4 popular benchmark problems with over 5% margins.

Index Terms—Continual Learning, Incremental Learning, Domain Incremental Learning

I. INTRODUCTION

CONTINUAL learning (CL) constitutes a research area of growing interests where the main goal is to develop a learning agent that can accumulate knowledge overtime [1], [2], [3], [4]. This task is challenging for a deep neural network (DNN) because of the catastrophic forgetting (CF) problem [5] where learning a new task over-writes previously valid parameters thus suffering from significant performance drops on previous tasks. In other words, a plastic model is capable of mastering new tasks, but its performance on previous tasks is compromised. In contrast, a stable model retains its old knowledge but fails to learn new tasks. Hence, the key is to balance the issue of plasticity and stability [5].

There exist three classical approaches in combating the CF problem: regularization-based approach [6], [7], [8], [9], [10], [11], [12], [13], memory-based approach [14], [15],

[16], [17], [18], [19], [20], [21], [22], [23] and architecture-based approach [24], [25], [26], [27], [11], [28]. Although the memory-based approaches generally outperform the other two approaches, they impose storage and privacy concerns. Recently, there is a trend toward a rehearsal-free continual learning approach benefiting from strong generalization power of pre-trained model (PTM) [29], [30], [31], [32], [33]. That is, these approaches fix the backbone network to avoid the CF problem while adapting to new tasks via parameter-efficient fine-tuning (PEFT) approaches such as prompts [29], [30], [31], LoRA [34], [35] and adapter [36]. Such an approach is effective because it involves a small number of trainable parameters and is powerful because it enjoys generalized features of the foundation model. Nevertheless, most of these efforts are devoted to addressing the class-incremental learning (CIL) problems [4] where each task features disjoint classes and a model is expected to handle all classes of all tasks in the testing phase.

Domain Incremental Learning (DIL) is a CL sub-problem where the goal is to handle a sequence of varying domains. It differs from the CIL problem because each task is drawn from different domains while sharing the same label space [37]. That is, a model is trained to be robust against various changes, such as data style shifts, data quality degradation, environmental changes, etc., rather than to recognize new classes. In [38], Compositional Prompt (C-Prompt) is proposed as a rehearsal-free solution to the DIL problem, where it puts forward the notion of a domain-specific prompt pool. [39] proposes the concept of dual-level concept prototypes (DualCP), which consists of coarse-grained and fine-grained prototypes for each class. The idea of SOYO is presented in [40], where the key idea lies in the Gaussian Mixture Compressor and Domain Feature Resampler. [41] offers the idea of dual consolidations in the feature level and in the classifier level. Although DIL has rapidly grown, existing works suffer from at least three bottlenecks: 1) they overlook shared knowledge across tasks. Notwithstanding that [38] also makes use of the global prompt to capture shared knowledge. We offer an alternative approach here, where the idea of low-rank adaptation (LoRA) is implemented. In addition, [38] imposes considerable complexities because of the use of domain-specific prompt pools; 2) they suffer from low parameter selection accuracies, leading to incorrect parameters being utilized for inferences. That is, existing approaches often apply the matching degree between the prompt key and query to select task-specific prompts during inferences, resulting in inaccurate parameter selections; 3) they adopt the linear

Naeem Paeedeh, Mahardhika Pratama, Jimmy Cao, Wolfgang Mayer, Ryszard Kowalczyk are with the STEM, University of South Australia, Australia (E-Mail: Naeem.Paeedeh@mymail.unisa.edu.au, {Dhika.Pratama, Wolfgang.Mayer, Jimmy.Cao, Ryszard.Kowalczyk}@unisa.edu.au)

Weiping Ding is with the School of Artificial Intelligence and Computer Science, Nantong University, Nantong, China. (E-Mail: dwp9988@163.com)

or prototype-based classifiers having subpar generalization power.

To correct this shortcoming, we propose continual knowledge consolidation low-rank adaptation (CONEC-LoRA) to deal with the DIL problems. CONEC-LoRA features task-specific LoRAs and task-shared LoRAs in which the first l blocks are assigned to the task-shared LoRA, leaving the remainder $L - l$ blocks for the task-specific LoRA. Our parameter selection strategy for inference is driven by an auxiliary network predicting the domain ID of a testing sample. The auxiliary network is trained during the training process with the projection-based Gaussian mixture model (PGMM) to address the class imbalance problem caused by the absence of previous samples. The idea of the projection or transformation module is introduced to prevent the synthetic sample-bias problem. Another innovation lies in the different-depth network structure, where every layer is assigned a local classifier. The final prediction is drawn from a local classifier maximizing the logit. Last but not least, the concept of a stochastic classifier is integrated to improve the model's generalizations. Unlike the popular prototype-based classifier, the stochastic classifier benefits from the mean and variance vectors sampled from a distribution. This technique guarantees the presence of an infinite number of classifiers, promoting correct classifications.

The key differences between our approach and [34] lie in the task-specific LoRA selection strategy, where the auxiliary network is implemented to select the task-specific LoRA rather than the similarity-based weighting scheme. Such a strategy guarantees the isolation of task-specific LoRAs, allowing for the retention of task-specific details, i.e., robust against the CF problem. By extension, [34] applies the prototype-based classifier, whereas ours is underpinned by the stochastic classifier. We also address the DIL problem here rather than the CIL problem. Note that the adapter fusion strategy in [34] undermines the domain-specific knowledge because all task-specific LoRAs are aggregated. Consequently, this issue leads to the use of orthogonality loss in their approach. In addition, such an approach incurs considerable complexities because multiple forward passes need to be carried out for all task-specific LoRA. On the contrary, our approach is much simpler than that because it predicts the correct task-specific LoRA followed by the network inference. This paper conveys at least four major contributions.

- This paper proposes continual knowledge consolidation low rank adaptation (CONEC-LoRA) for domain incremental learning (DIL) problems, where it features general LoRAs to extract cross-task knowledge and task-specific LoRAs to capture domain-specific knowledge. The building block of CONEC-LoRA is constructed under a stacked pretrained transformer block of L layers, where the first l blocks are designated for the task-shared LoRA, while the remaining $L - l$ blocks are reserved for the task-specific LoRA.
- A parameter selection strategy is proposed based on an auxiliary network predicting the domain label for the task-specific LoRAs for inferences. The auxiliary network receives an input of the frozen backbone and is trained

with the projection-based GMM under a joint loss to address the class imbalance issue as a result of the absence of previous samples. The concept of the projection or transformation module is applied to random inputs of the GMM to prevent the synthetic sample bias problem. In addition, we introduce the ball-generator loss to prevent the synthetic sample bias problem. The auxiliary network is structured under a different-depth configuration, where each layer is connected to a local classifier with an early exit strategy to reduce computational complexities. The local classifier with the most confident prediction is selected for inferences.

- The concept of the stochastic classifier is put forward for domain incremental learning (DIL) problems. This idea extends [42] devised for the CIL problem. This classifier is represented by learnable mean and covariance vectors, with the classifier weight defined as a distribution. This trick enables the creation of an arbitrary number of classifiers, increasing the chance of correct predictions.
- The advantage of CONEC-LoRA has been rigorously evaluated under four benchmark DIL problems. It is compared with prominent DIL algorithms where CONEC-LoRA outperforms them by over 5% margin. In addition, the source code of CONEC-LoRA is made publicly available in <https://github.com/Naeem-Paeedeh/CONEC-LoRA> for reproducibility and convenient further study.

The remainder of this paper is structured as follows: Section 2 discusses the related works; Section 3 outlines preliminaries, including a problem definition and basic concepts; Section 4 describes our algorithm, namely CONEC-LoRA; Section 5 offers our numerical study; and some concluding remarks are drawn in the last section of this paper.

II. RELATED WORKS

Continual learning (CL) aims to address non-stationary learning problems, where a model cannot be fixed once deployed because its predictions quickly become outdated due to changing learning environments. That is, a CL agent is faced with never-ending environments under dynamic conditions [1], [2], [3], [4], resulting in the stability-plasticity dilemma. The CL problem can be divided into three sub-problems [37]: task-incremental learning (TIL), class-incremental learning (CIL), and domain-incremental learning (DIL). The TIL and the CIL are inherently identical, where the difference lies in only the presence of task identifiers during the evaluation phase in the TIL. Therefore, the CIL is deemed more challenging than the TIL.

Class-Incremental Learning (CIL) problem is formulated as a learning problem of sequentially arriving classes. That is, each task presents a set of classes disjoint across tasks. Learning a new task, thus, induces the CF problem because previously valid parameters are catastrophically erased with new ones, i.e., parameter drifts. The regularization-based approach [6], [7], [8], [9], [10], [11], [12], [13] deals with the CF problem via integration of a regularization term to avoid the parameter drift problem. However, the regularization-based approach usually doesn't scale well for large-scale problems

because it is difficult to find an overlapping region of all tasks. The architecture-based approach [24], [25], [26], [27], [11], [28] adds new components while isolating old components when learning new tasks to combat the CF problem. Nevertheless, the architecture-based approach usually calls for the existence of the task IDs, thus being incompatible with the CIL problem, otherwise, it requires old samples to be stored. The memory-based approach [14], [15], [16], [17], [18], [19], [20], [21], [22], [23] utilizes a memory buffer storing old samples. Old samples can then be interleaved with new samples for experience replay steps to prevent the CF problem. Although it is evident that the memory-based approach outperforms the other two approaches, it raises storage and privacy concerns. This issue motivates the development of rehearsal-free approaches benefiting from the strong generalization power of pretrained model (PTM) [33]. This approach is combined with the parameter-efficient fine-tuning (PEFT) approaches such as prompts [29], [30], [31], LoRA [34], [35], and adapter [36]. Such an approach offers strong performances because the CF problem can be effectively alleviated to a very low level by freezing the backbone networks while only adjusting a small number of external parameters. Although there exist numerous contributions in the realm of the CIL problem, the DIL problem remains a relatively uncharted territory.

Domain Incremental Learning (DIL) problem is unlike the CIL problem, where each task carries a unique domain while sharing the same label space. That is, a model is trained to be robust against various changes rather than to recognize new classes. [38] proposes the idea of compositional prompts, where each domain is assigned a prompt pool. In addition, the global prompt is injected and shared across all tasks. [39] proposes the idea of dual prototypes, where coarse-grained and fine-grained prototypes are put forward. In [40], the concepts of Gaussian Mixture Compressor and Domain Feature Resampler are introduced to address inaccurate parameter selections in the DIL context. [41] proposes dual knowledge consolidations in the feature level as well as in the classifier level. The class imbalance problem in the DIL context is analyzed and addressed in [43] by using the multi-expert concept. In [44], the KA-prompt is proposed, where the key idea lies in the initialization of the new prompt with the most compatible previous prompt, thus promoting cross-domain knowledge. Nonetheless, we note at least three shortcomings in existing works: 1) they focus on the domain-specific knowledge while overlooking the domain-shared knowledge, resulting in sub-optimal performances. Note that [38] is based on the idea of the prompt, whereas our approach is constructed under the LoRA concept. In addition, [38] imposes prohibitive complexities due to the use of prompt pools for every domain; 2) existing approaches suffer from inaccurate parameter selection during inferences, significantly affecting the model's performance because wrong components are associated when producing an output. Although this problem is partially addressed in [44], this approach relies on the greedy search technique, incurring prohibitive complexities; 3) existing approaches are built upon a linear or prototype-based classifier having limited generalization power.

III. PRELIMINARIES

A. Problem Definition

Domain-incremental learning (DIL) is formulated as a learning problem of sequentially arriving domains $\{\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_B\}$, where B is the number of domains. Each domain comprises the training set \mathcal{T}_{tr}^b and the testing set \mathcal{T}_{te}^b , where $\mathcal{D}_b = \{\mathcal{T}_{tr}^b, \mathcal{T}_{te}^b\}$. When encountering the b -th domain, a model is only presented with the training set of the b -th domain \mathcal{T}_{tr}^b with the absence of any training samples of the previous domains $\mathcal{T}_{tr}^{1 \sim (b-1)} = \cup_{t=1}^{b-1} \mathcal{T}_{tr}^t$ and any exemplars of the previous domains leading to the catastrophic forgetting (CF) problem. The evaluation is performed against the testing samples of all already seen domains $\mathcal{T}_{te}^{1 \sim b} = \cup_{t=1}^b \mathcal{T}_{te}^t$. The training set of the b -th domain is formed by N_b tuples $\mathcal{T}_{tr}^b = \{(x_i, y_i)\}_{i=1}^{N_b}$, where x_i is an i -th image and y_i is its corresponding label. Suppose that \mathcal{Y}_b presents the label set of the b -th domain, each domain shares the same label set $\forall i \in [1, N_b], y_{i,b} \in \mathcal{Y}_b$, where $\forall b, b' \in [1, B], \mathcal{Y}_b = \mathcal{Y}_{b'}$ but features different characteristics such as styles and/or environmental changes. That is, there exists the presence of concept drifts $P(\mathcal{X}, \mathcal{Y})_b \neq P(\mathcal{X}, \mathcal{Y})_{b'}$. No domain identifiers (IDs) are offered for inferences. In other words, a model is supposed to infer the domain IDs b by itself.

Our model $f_\theta = h_\phi \circ g_\psi$ relies on a pre-trained vision transformer (ViT) model, where $h_\phi : \mathcal{Z} \rightarrow \mathcal{Y}$ constitutes a classifier and $g_\psi : \mathcal{X} \rightarrow \mathcal{Z}$ denotes a feature extractor frozen during the training process. The training process of the b -th domain is to adapt lightweight parameters g_W inserted into the query and key projections of the attention module of all transformer blocks of the ViT backbone. Specifically, this is written as follows:

$$z^{(\ell)} = g_\psi^{(\ell)}(x) + g_W^{(\ell)}(x), \quad (1)$$

where $z^{(\ell)}$ is an embedded feature of layer ℓ for query or values. In addition, the classifier h_ϕ is defined as the stochastic classifier here.

B. Low-Rank Adaptation (LoRA)

LoRA is a parameter-efficient fine-tuning (PEFT) approach [45], [34], which is capable of adapting a foundation model to downstream tasks. The key idea is seen in the use of a pair of rank decomposition matrices, namely a down-projection matrix $B \in \mathbb{R}^{r \times k}$ and an up-projection matrix $A \in \mathbb{R}^{d \times r}$ with rank $r \ll \min(d, k)$. Hence, it achieves learning efficiency by only learning $r \times (d + k)$, which can be expressed as follows:

$$z = g_\psi(x) + \Delta W x, \quad (2)$$

where $\Delta W = AB$. The low-rank matrices can then be attached to the transformer block. For DIL, these matrices can be set to be domain-specific $\{A_b, B_b\}$. When adapting to the b -th domain, only ΔW is learnable, leaving other parameters fixed, thus combating the CF problem. However, such naïve approach of using the same LoRAs in all continual learning tasks loses shared knowledge across domains. Besides, it risks inaccurate parameter selection during inference, which can lead to a loss of performance because no oracle for domain IDs is provided during the testing process.

C. Stochastic Classifier

Given a classifier $h_\phi : \mathcal{Z} \rightarrow \mathcal{Y}$ mapping the latent space to the label space, the prototype-based classifier relies on a set of prototypes defined as an empirical mean of each class $\phi_m = \{\mu_m\}$.

$$\mu_m = \frac{1}{N_m} \sum_{i=1}^{N_m} \mathbb{I}_{y_i=m} z_i, \quad (3)$$

where M denotes the dimension of the label space and $\mathbb{I}_{y_i=m}$ is an indicator function being true if $y_i = m$. Suppose that $\langle \cdot \rangle$ stands for the cosine similarity, the predicted output of the prototype-based classifier is expressed as follows:

$$P(y = m|x) = \frac{\exp(\eta \langle \phi_m, z \rangle)}{\sum_{i=1}^M \exp(\eta \langle \phi_m, z \rangle)}, \quad (4)$$

where η is a temperature controlling the smoothness of a distribution. Such classification implies the use of a single static classifier, leading to suboptimal generalizations. The stochastic classifier [46], [42] offers a different perspective where the classifier weights are sampled from a distribution $\phi_m = \{\mu_m, \sigma_m\}$, paving the way for the applications of multiple classifiers at once. The stochastic classifier is formalized:

$$\langle \phi_m, z \rangle \quad \text{s.t.} \quad \phi_m = \mu_m + \mathcal{N}(0, 1) \odot \sigma_m, \quad (5)$$

where μ_m, σ_m are mean and variance vector automatically learned in an end-to-end manner and \odot is an element-wise multiplication. Such a trick enables the creation of an arbitrary number of classifiers, thereby increasing the number of correct classifications. Fig. 1 visualizes the stochastic classifier.

IV. METHODOLOGY

This paper proposes continual knowledge consolidation low-rank adaptation (CONEC-LORA) built upon a synergy between task-shared LoRA and task-specific LoRA. That is, the first l transformer blocks are attached with task-shared LoRAs while the remainder $L - l$ blocks are assigned with task-specific LoRAs. The classification decision is navigated with the stochastic classifier. For selecting the task-specific LoRAs during inference, the auxiliary network is applied to predict the domain label and is trained with the projection-based Gaussian mixture model (GMM) to address the class imbalance issue resulting from the absence of previous samples and any exemplars. The projection or transformation concept is implemented to cope with the synthetic sample bias problem while introducing a different-depth network structure for inference. An overview of CONEC-LORA's learning policy is outlined in Fig. 2.

A. Network Architecture

Suppose the ViT backbone with L blocks, the shared LoRAs $\{A_s, B_s\}$ are inserted to the first l blocks where $l \leq L$ while the remaining $L - l$ blocks are designated to the task-specific LoRA $\{A_b, B_b\}$. Therefore, the output of the i -th transformer block z_i is written:

$$z_i = g_{\psi_i}(z_{i-1}) + \begin{cases} A_s^i B_s^i z_{i-1} & i \leq l \\ A_b^i B_b^i z_{i-1} & l < i \leq L \end{cases}, \quad (6)$$

where z_{i-1} denotes the output of the previous transformer block $z_0 = x$ while g_{ψ_i} represents the i -th transformer block frozen during the training process. b is the domain label to be predicted during the inference. Note that early layers (1 to l) focus on general patterns to be shared across domains while deep layers ($l+1$ to L) capture domain-specific details.

Motivated by [34], the domain-shared LoRA applies a fixed random orthogonal down-projection matrix $B_s \in \mathbb{R}^{r \times k}$ and a trainable up-projection matrix $A_s \in \mathbb{R}^{d \times r}$ initialized as zeros. It satisfies the following property

$$\Delta W x = A_s B_s x, \quad B_s B_s^\top = I, \quad (7)$$

where $I \in \mathbb{R}^{r \times r}$ is an identity matrix. [47] shows that adjusting A is more effective than adjusting B . We also confirm that a random B can compete to a fully trained one. First, a random matrix $M \in \mathbb{R}^{r \times k}$ is generated from a normal distribution $\mathcal{N}(0, 1)$. A random orthogonal matrix can be obtained via the singular value decomposition (SVD) [48].

$$M = U \Sigma V^\top, B_s = U V_r^\top, \quad (8)$$

where v_r is the first r rows of V . In a nutshell, our effort is devoted to learn A_s during the training process while keeping the random orthogonal matrix B_s frozen.

CONEC-LoRA relies on the stochastic classifier during the training process, whose parameters are drawn from a distribution under the cosine classification strategy as per (5). This strategy allows proper training process where $\phi_m = \{\mu_m, \sigma_m\}$ are trained in an end-to-end fashion. Our investigation reveals that replacing each μ_m with the corresponding prototype after training leads to performance improvements at the inference phase. Note that the stochastic classifier with learnable μ_m must still be applied during the mini-batch training process, as using a prototype classifier during training results in performance degradation. In other words, the use of multiple classifiers as per the stochastic classifier is useful during the training process to boost the learning performance, but it may not surpass the flexibility of the prototype classifier during the inference phase. Therefore, μ_m vectors of the stochastic classifier are replaced with prototypes for the inference phase.

B. Loss Function

CONEC-LORA is learned by minimizing a joint loss function consisting of the conventional cross-entropy loss and the knowledge distillation (KD) loss.

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda_1 \mathcal{L}_{kd}, \quad (9)$$

where λ_1 is a trade-off coefficient controlling the strength of the KD loss function. The KD strategy is omnipresent in the CL domain to combat the CF problem but this strategy is often too strict and leads to degraded plasticity. Hence, an early exit mechanism is implemented here where the KD approach is applied at the transition point between task-shared and task-specific LoRA, i.e., the l -th block. Specifically, using a local classifier h_ϕ , the KD mechanism is done by extracting the [CLS] token representation through the shared LoRAs, i.e., 1

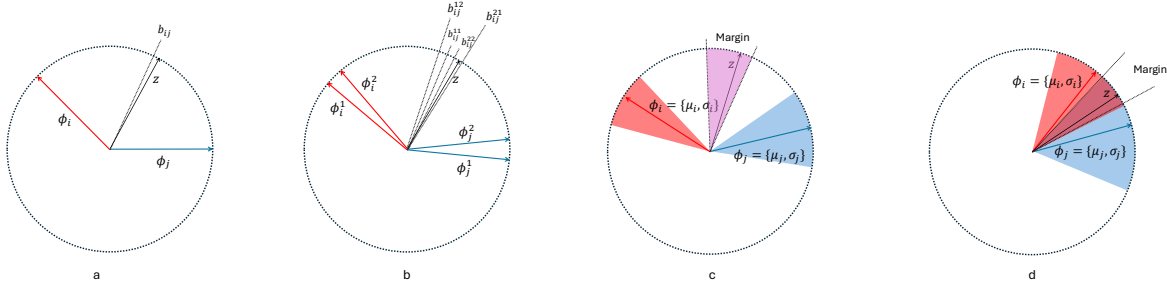


Fig. 1: A stochastic classifier. a) The decision boundary for a single cosine classifier with two weight vectors for two classes. The feature vector belongs to the blue class. b) The illustration of the effect of having two weight vectors per class, resulting in four decision boundaries. c) It shows the effect of having an infinite number of weight vectors on the decision boundary. d) The stochastic nature of the classifier improves the feature separation of the model and classifier by resolving the confusion when the margin is unclear.

to l , for both current domain $z_b^l[\text{CLS}]$ and previous domain $z_{b-1}^l[\text{CLS}]$.

$$\mathcal{L}_{\text{kd}} = \sum_{i \in D_b} s_{b-1,i}^\tau \log s_{b,i}^\tau, \quad (10)$$

where $s_b^\tau = \text{Softmax}(h_\phi^b(z_b^l[\text{CLS}]))/\tau$ and $\tau = 2$ is a temperature. Because A_s is initialized as zero and learned during the training process, its norms mirror the contribution of each element, implying which parts need further tuning. This fact enables redistribution of the gradient from \mathcal{L}_{kd} based on the weight importance in the previous training session. Suppose that $\|a_{s,j}^{b-1}\|_2$ denotes the L_2 norm of the j -th element of the up-projection matrix A_s^{b-1} of the previous domain, the gradient of A_s^b of the current domain is redistributed:

$$\nabla_{A_s^b} \mathcal{L}_{\text{kd}} = \nabla_{A_s^b} \mathcal{L}_{\text{kd}} \odot \sigma(\{\|a_{s,j}^{b-1}\|\}_{j=1}^d), \quad (11)$$

where \odot is Hadamard product, $\sigma(w) = d \times w / \sum_{i=1}^d w_i$ is the dimension-preserving normalization function, and d is the embedding dimension. This redistribution strategy assures that each element is adaptively adjusted, i.e., the essential dimension for the previous task is preserved while allowing other dimensions to move freely in respect to the new task.

C. Auxiliary Domain Classifiers

CONEC-LORA is guided by an auxiliary network that predicts the domain label during inference. The predicted domain label is thus exploited to select the task-specific LoRAs without the CF problem. Since the auxiliary network is trained, it is expected to be more reliable than a non-parametric approach, such as the proposed method in [42] that uses the Mahalanobis distance. The representations of the intermediate layers of Neural Networks (NNs) are more robust to the changes during the continual learning and suffer less from forgetting [49]. Furthermore, the classification can potentially be finished on earlier layers with even higher accuracy. In CONEC-LoRA, we keep the backbone frozen without utilizing the LoRAs and train independent domain classifiers at different levels of abstraction in the selected layers of a ViT. This ensures that the domain classification module would not be affected during the weight updates.

The domain label is predicted by a local linear classifier for every layer. That is, every layer is assigned with a local classifier $h_{\phi_{\text{aux}}^l}(\cdot)$. We also insert a transformation module constructed as a two-layer MLP $\kappa_{\gamma^l}(\cdot)$ to prevent the synthetic sample bias problem, where γ^l denotes the parameters of the transformation module of the l -th layer. Given the embedding of the l -th layer with the frozen backbone network and without any LoRAs z_l , the local output of the l -th layer of the domain classification module is expressed as follows:

$$b = h_{\phi_{\text{aux}}^l}(\kappa_{\gamma^l}(z_l)), l \in [1, L], \quad (12)$$

where $b \in [0, 1]$ denotes the output logit of the l -th local classifier and L is the number of layers. The aggregation is performed by applying the maximum operator to each local classifier in every layer.

$$\hat{b} = \arg \max_{1 \leq l \leq L} h_{\phi_{\text{aux}}^l}(\kappa_{\gamma^l}(z_l)), \quad (13)$$

where $\hat{b} \in B$ and B is the number of domains. In other words, the local classifier having the most confident prediction is selected for inferences. The different-depth network structure is designed to benefit from the intermediate representation, making it more robust to dynamic conditions than a single classifier [49]. During the inference, we also implement the early exit mechanism to relieve the computational burden. That is, we set a confidence threshold ς and examine whether the confidence of local classifiers exceeds the threshold. The earliest layer passing the threshold is selected for inference. This strategy avoids going through all layers for inference.

$$h_{\phi_{\text{aux}}^l}(\kappa_{\gamma^l}(z_l)) \geq \varsigma \rightarrow \hat{b} = \arg \max_b h_{\phi_{\text{aux}}^l}(z_l) \quad (14)$$

The earliest layer satisfying (14) is selected for inference, or the early exit mechanism is implemented. Nonetheless, the training process of the auxiliary network is non-trivial since there exists a class imbalance problem between the current domain and the previous domains. As with [40], we model the previous data distributions using a mixture of C multivariate Gaussian distributions $\mathcal{N}(\nu, K)$.

$$\mathcal{N}(z|\nu, K) = \frac{1}{(2\pi)^{\frac{d}{2}}|K|^{\frac{1}{2}}} \exp(-\frac{1}{2}(z - \nu)^\top K^{-1}(z - \nu)), \quad (15)$$

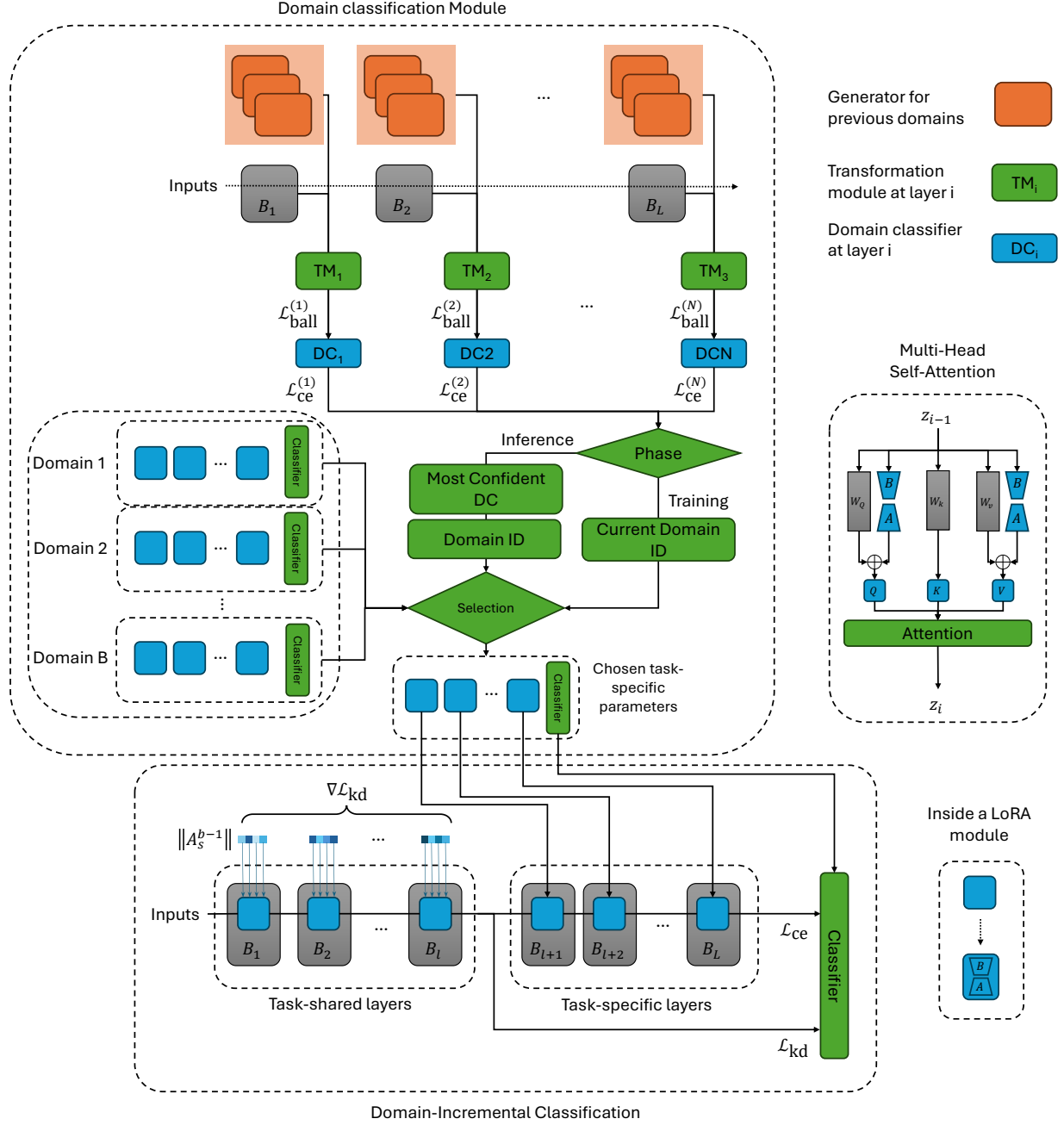


Fig. 2: The architecture of the model and pipeline of CONEC-LoRA. CONEC-LoRA consists of two main stages. First, the input is given to the frozen backbone. The logits of the domain classifiers (DCs) are being calculated. The domain is detected based on the decision of the most confident domain classifier. Second, the domain-specific LoRAs are chosen for use in the domain-specific blocks. The gray rectangles represent the frozen layers of the ViT.

where d is the dimension of embeddings. That is, the mean and covariance matrix of every domain (ν_b, K_b) in the embedding space \mathcal{Z} is calculated for each one of the C components.

The probability density function of the mixture of the Gaussian distributions can be defined as:

$$p(z) = \sum_{c=1}^C \omega_c \mathcal{N}(z | \nu_c, K_c), \quad (16)$$

where $0 \leq \omega_c \leq 1$ and $\sum_{c=1}^C \omega_c = 1$. The GMM parameters are trained with the Expectation-Maximization (EM) algorithm [50]. By storing the $\{\omega_c, \nu_c, K_c\}_{c=1}^C$, we can generate embed-

dings of the previously seen domain from the GMMs:

$$q_l \sim \text{Categorical}(\omega_l), \quad (17)$$

$$\hat{z}_l|q_l \sim \mathcal{N}(\nu_{q_l}^l, K_{q_l}^l). \quad (18)$$

Nevertheless, this naïve sampling strategy causes the synthetic sample bias problem. To remedy this problem, the transformation module $\kappa_{\gamma^l}(\cdot)$ is integrated. Specifically, we feed $\{\tilde{z}_\ell\} = \{\hat{z}_\ell\} \cup \{z_\ell\}$ set to the transformation module. For $\tilde{z}_\ell^i \in \{\tilde{z}_\ell\}$ we calculate $\hat{z}_\ell^i = \kappa_{\gamma^l}(\tilde{z}_\ell^i)$, where $\hat{z}_\ell \in \mathbb{R}^d$. The ball-generator loss [51] is introduced and inspired by the triplet and hinge losses.

$$\mathcal{L}_{\text{ball}}^l = \sum_{\hat{z}_l \in \bar{y}_i, \hat{z}_l \notin \bar{y}_j} \max\{0, d(\hat{z}_l, \mu_i) + r - d(\hat{z}_l, \mu_j)\}, \quad (19)$$

where r is a predefined margin constant, \bar{y} is the domain label, and $d(\cdot)$ is a distance measure. The goal of the ball-generator loss is to pull the synthetic sample towards and close to its center while keeping it away from other centers. Hence, the local auxiliary network $h_{\phi_{\text{aux}}^l}(\cdot)$ of the l -th layer is trained to minimize the following joint loss.

$$\mathcal{L}_{\text{aux}}^l = \mathcal{L}_{\text{ce}}^l + \lambda_2 \mathcal{L}_{\text{ball}}^l, \quad (20)$$

where λ_2 is a tradeoff coefficient controlling the strength of the ball generator loss. Note that $\mathcal{L}_{\text{ce}}^l$ is subject to both the current samples $z \sim \mathcal{D}_b$ and the synthetic samples $\hat{z}_l|q_l \sim \mathcal{N}(\nu_{q_d,k}^l, K_{q_d,k}^l)$, $q_l \sim \text{Categorical}(\omega_l)$ representing the previous domains. The pseudo-codes of different phases of the CONEC-LoRA are shown in Algorithms 1 to 3.

D. Time Complexity Analysis

The proposed method requires two forward passes during training or inference. In one pass, we forward the inputs without LoRA modules to obtain the embeddings for the intermediate layers for the intermediate domain classifiers. In another pass, we should obtain the network's outputs by using both the shared and domain-specific LoRA modules. The crucial parts of the network that their calculations primarily affect in terms of time complexity include the two forward passes of the ViT and LoRA modules, the auxiliary components of the network, and fitting the GMM. In the following, we first calculate the time complexity of those components. Next, we calculate the overall complexity.

The major components of ViT calculations are the MLPs and attention mechanisms within each block. In a ViT with the embedding dimension of d and sequence length of s , the query, key, and value projections require $O(sd^2)$. The attention mechanism requires $O(ds^2 + sd^2)$ operations. If a weight matrix in one layer of the MLP is $\hat{h}_{\ell-1} \times \hat{h}_\ell$ and the width of an input is n_0 , an MLP with L_{MLP} layers requires $O(s \sum_1^{L_{\text{MLP}}} (\hat{h}_{\ell-1} \hat{h}_\ell))$. Since a shallow two-layer MLP is commonly used in the ViTs, the time complexity becomes $O(sd\hat{h}_1)$. Each LoRA adapter has the same time complexity as a shallow two-layer MLP. Two LoRA adapters for the Q and V with rank r add $O(srd)$ operations ($r \ll d$) overhead to each block of the ViT. Therefore, the total time complexity

Algorithm 1 CONEC-LoRA: Training the LoRAs and classifiers

Require: Training sets $\{\mathcal{T}_{\text{tr}}^b\}_{b=1}^B$ for B domains, where $\mathcal{T}_{\text{tr}}^b = \{(x_i, y_i)\}_{i=1}^{N_b}$, backbone g_ψ with L layers, classifiers set $\{h_{\phi_b}\}_{b=1}^B$, task-shared LoRAs $\{g_{W_\ell^s}\}_{\ell=1}^L$, task-specific LoRAs $\{g_{W_\ell^b}\}_{\ell=l+1}^L$, M as the number of epochs, B domains, and \mathcal{C} classes for every domain.

```

1: ► Training:
2: for each  $b \in \{1, 2, \dots, B\}$  do
3:    $h_{\text{temp}} \leftarrow \text{TEMPORARY STOCHASTIC CLASSIFIER}$ 
4:   for each  $i \in \{1, 2, \dots, M\}$  do
5:     for each  $(x_j, y_j) \in \mathcal{T}_{\text{tr}}^b$  do
6:        $z_j \leftarrow g_\psi(x_j; \{g_{W_\ell^s}\}_{\ell=1}^L, \{g_{W_\ell^b}\}_{\ell=l+1}^L)$ 
7:        $z_j \leftarrow h_{\text{temp}}(z_j)$ 
8:        $\mathcal{L}_{\text{ce}} \leftarrow \text{LOSS}_{\text{CE}}(z, y_j)$  {Classification loss}
9:        $z_j^b \leftarrow x_j$ 
10:      if  $b > 1$  then
11:         $z_j^b \leftarrow g_\psi(x_j; \{g_{W_\ell^s}\}_{\ell=1}^L)$ 
12:         $z_j^{b-1} \leftarrow g_\psi(x_j; \{g_{W_\ell^s}\}_{\ell=1}^L)$ 
13:         $\mathcal{L}_{\text{kd}}^j \leftarrow \mathcal{L}_{\text{kd}}^j + s_{b-1,j}^T \log s_{b,j}^T$  {Eq. (10)}
14:        for each  $\ell \leftarrow \{1, 2, \dots, l\}$  do
15:           $\nabla_{A_s^{b,\ell}} \mathcal{L}_{\text{kd}}^j = \nabla_{A_s^{b,\ell}} \mathcal{L}_{\text{kd}}^j \odot \sigma(\|a_{s,k}^{b-1,\ell}\|_{k=1}^d)$ 
16:        end for
17:      end if
18:    end for
19:    Minimize the  $\mathcal{L} = \mathcal{L}_{\text{ce}}(z, y) + \lambda_1 \mathcal{L}_{\text{kd}}$  {Eq. (9)}
20:  end for
21:  {Setting the weight of the classifier to prototypes}
22:  for each  $c \in \mathcal{C}$  do
23:     $S = \{x_j | y_j = c\}, (x_j, y_j) \in \mathcal{T}_{\text{tr}}^b$ 
24:     $\mu_c \leftarrow \frac{1}{|S_c|} \sum_{x \in S_c} g_\psi(x_j; \{g_{W_\ell^s}\}_{\ell=1}^L, \{g_{W_\ell^b}\}_{\ell=l+1}^L)$ 
    {Prototypes}
25:     $\sigma_c \leftarrow h_{\text{temp}}$ 
26:     $h_{\phi_b} \leftarrow (\mu_c, \sigma_c)$ 
27:  end for
28: end for

```

of a ViT with LoRAs is $O(Ls(ds + d^2 + d\hat{h}_1 + rd))$ or an L -layer ViT.

In GMM calculations, the most complex operations are the Mahalanobis distance calculations and matrix inversion. If the number of samples is n and the number of components is c , the Mahalanobis distance operations are of $O(ncd^2)$. The matrix inversion requires $O(cd^3)$. However, since c is usually small (2 in our experiments) and the dimension of embeddings in the network is fixed, this operation is not affected by the number of samples; therefore, it is negligible. Overall, the time complexity of the GMM is $O(I_G Lncd^2)$, where I_G represents the maximum number of iterations for fitting the GMM. The GMM training occurs only during the training.

Finally, there are L MLPs as transformation modules and domain classifiers after each layer of the ViT. Each MLP as a transformation module with L_{TM} layers requires $O(s \sum_1^{L_{\text{TM}}} (\hat{h}_{\ell-1} \hat{h}_\ell))$ if a weight matrix in an MLP layer is $\hat{h}_{\ell-1} \times \hat{h}_\ell$. Moreover, each domain classifier requires $O(sd^2)$ operations. Therefore, the total time complexity of

Algorithm 2 CONEC-LoRA: Training the Domain Classification Module

Require: Training sets $\{\mathcal{T}_{tr}^b\}_{b=1}^B$ for B domains, where $\mathcal{T}_{tr}^b = \{(x_i, y_i)\}_{i=1}^{N_b}$, backbone g_ψ with L layers, domain classifiers set $\{h_{\phi_{aux}^\ell}\}_{\ell=1}^L$, M as the number of epochs, \mathcal{C} classes for every domain, a set of transformation modules $\{\kappa_\gamma^\ell\}$, C components for GMMs, and N as the number of synthetic embeddings to generate.

- 1: **► Training the Domain Classification Module:**
- 2: **for** each $b \in \{1, 2, \dots, B\}$ **do**
- 3: {Obtaining the embeddings for every layer of the current domain}
- 4: **for** each $x_j \in \mathcal{T}_{tr}^b$ **do**
- 5: $z_0^j \leftarrow x_j$
- 6: **for** each $\ell \in \{1, 2, \dots, L\}$ **do**
- 7: $z_\ell^j \leftarrow g_{\psi_{\ell-1}}(z_{\ell-1}^j)$ {Without LoRAs}
- 8: $\tilde{y}_{b,j} \leftarrow b$
- 9: **end for**
- 10: **end for**
- 11: {Calculating the centers for ball-generator loss}
- 12: **for** each $\ell \in \{1, 2, \dots, L\}$ **do**
- 13: $\mu_\ell^b \leftarrow \frac{1}{|\mathcal{Z}_\ell|} \sum_{z \in \mathcal{Z}_\ell} z$ {Center for layer ℓ }
- 14: **end for**
- 15: {Finding the GMM parameters for the following tasks}
- 16: **for** each $\ell \in \{1, 2, \dots, L\}$ **do**
- 17: $\{\omega_c^{b,\ell}, \nu_c^{b,\ell}, K_c^{b,\ell}\}_{c=1}^C \leftarrow \text{GMM}_{EM}(\mathcal{Z}_\ell)$
- 18: **end for**
- 19: {Generating the synthetic embeddings}
- 20: **for** each $d \in \{1, 2, \dots, b-1\}$ **do**
- 21: **for** each $k \in \{1, 2, \dots, N\}$ **do**
- 22: $q_{d,k} \sim \text{Categorical}(\omega^d)$
- 23: $\hat{z}_{d,k} | q_{d,k} \sim \mathcal{N}(\nu_{q_{d,k}}^d, K_{q_{d,k}}^d)$
- 24: $\tilde{y}_{d,k} \leftarrow d$
- 25: **end for**
- 26: **end for**
- 27: **for** each $m \in \{1, 2, \dots, M\}$ **do**
- 28: **for** each $\ell \in \{1, 2, \dots, L\}$ **do**
- 29: {Calculation of the ball-generator loss}
- 30: **if** $b > 1$ **then**
- 31: $\{\tilde{z}_\ell\} \leftarrow \{\hat{z}_\ell\} \cup \{z_\ell\}$
- 32: **for** each $\tilde{z}_\ell^j \in \{\tilde{z}_\ell\}$ **do**
- 33: $\hat{\tilde{z}}_\ell^j \leftarrow \kappa_{\gamma^\ell}(\tilde{z}_\ell^j)$
- 34: **end for**
- 35: $\mathcal{L}_{ball}^\ell = \sum_{\hat{z}_\ell \in \tilde{y}_i, \hat{\tilde{z}}_\ell \notin \tilde{y}_j} \max\{0, d(\hat{\tilde{z}}_\ell, \mu_\ell^i) + r - d(\hat{\tilde{z}}_\ell, \mu_\ell^j)\}$
- 36: $\mathcal{L}_{aux}^\ell \leftarrow \mathcal{L}_{ce}^\ell(h_{\phi_{aux}^\ell}(\hat{\tilde{z}}_\ell)) + \lambda_2 \mathcal{L}_{ball}^\ell$
- 37: **else**
- 38: $\mathcal{L}_{aux}^\ell \leftarrow \mathcal{L}_{ce}^\ell(h_{\phi_{aux}^\ell}(z_\ell))$
- 39: **end if**
- 40: **end for**
- 41: Minimize the \mathcal{L}_{aux}
- 42: **end for**
- 43: **end for**

Algorithm 3 CONEC-LoRA: Inference phase

Require: Test set $\mathcal{T}_{te}^{1 \sim b} = \cup_{t=1}^b \mathcal{T}_{te}^t$, backbone g_ψ , classifiers set $\{h_{\phi_b}\}_{b=1}^B$, domain classifiers set $\{h_{\phi_{aux}^\ell}\}_{\ell=1}^L$, a set of transformation modules $\{\kappa_\gamma^\ell\}$, task-shared LoRAs $\{g_{W_\ell^s}\}_{\ell=1}^L$, task-specific LoRAs $\{g_{W_\ell^b}\}_{\ell=l+1}^L$, \mathcal{C} classes for every domain, and threshold value τ .

- 1: **► Inference phase:**
- 2: **for** each $x_i \in \mathcal{T}_{te}^b$ **do**
- 3: {Obtaining the embeddings for every layer}
- 4: $z_0^i \leftarrow x_i$
- 5: **for** each $\ell \in \{1, 2, \dots, L\}$ **do**
- 6: {Embedding vectors from the frozen backbone for every layer without LoRAs}
- 7: $z_\ell^i \leftarrow g_{\psi_{\ell-1}}(z_{\ell-1}^i)$
- 8: $\tilde{z}_\ell^i \leftarrow h_{\phi_{aux}^\ell}(\kappa_{\gamma^\ell}(z_\ell^i))$
- 9: {Confidence of every domain classifier}
- 10: $\rho_\ell^i \leftarrow \max \tilde{z}_\ell^i$
- 11: **end for**
- 12: $S_i \leftarrow \{\ell | \rho_\ell^i \geq \tau\}, 1 \leq \ell \leq L$
- 13: **if** $S \neq \emptyset$ **then**
- 14: {The index of the first classifier that meets the criterion}
- 15: $\gamma_i \leftarrow \min(S_i)$
- 16: **else**
- 17: {The index of the most confident classifier}
- 18: $\gamma_i \leftarrow \arg \max_\ell (\rho_\ell^i)$
- 19: **end if**
- 20: $\hat{b}_i \leftarrow \arg \max(\rho_{\gamma_i}^i)$ {Predicted domain label}
- 21: $z_i \leftarrow g_\psi(x_i; \{g_{W_\ell^s}\}_{\ell=1}^L, \{g_{W_\ell^b}\}_{\ell=l+1}^L)$
- 22: $z_i \leftarrow h_{\phi_{\hat{b}_i}}(z_i)$
- 23: $\hat{y}_i \leftarrow \arg \max(z_i)$
- 24: **end for**

$$O\left(Ls(ds + d^2 + \sum_1^{L_{TM}}(\hat{h}_{\ell-1}\hat{h}_\ell))\right).$$

Considering all the complexities of the components, with a batch size of B , I_{TC} iterations for training for the classification, and I_{TD} iterations for the domain classification phase of the training, the time complexity of the training is

$$O\left(BL((I_{TC} + I_{TD})(ds^2 + sd^2 + sd\tilde{h}_1) + I_{TC}srd + I_{TD}s \sum_1^{L_{TM}}(\hat{h}_{\ell-1}\hat{h}_\ell)) + I_G Lncd^2\right), \quad (21)$$

and the time complexity of the inference for a sample is

$$O\left(Ls(ds + d^2 + d\tilde{h}_1 + \sum_1^{L_{TM}}(\hat{h}_{\ell-1}\hat{h}_\ell))\right). \quad (22)$$

V. EXPERIMENTS

A. Datasets

We evaluate the effectiveness of CONEC-LoRA on four established DIL datasets: DomainNet [52], CORE50 [53], CDDb-Hard [54], and Office-Home [55]. These datasets are

the domain classification phase (parameter selection) is of

executed with a random seed across five domain orders. The final numerical results are taken from the average across five different domain orders.

B. Baseline Methods

We compare our method comprehensively with DCE [43], DualCP [39], S-iPrompt [56], CODA-Prompt [57], RanPAC [58], EASE [59], DualPrompt [30], L2P [29], and SimpleCIL [60]. Moreover, we report the results for the replay-based method, including MEMO [61], iCaRL [62], and Replay [63].

C. Implementation Details

The CDDb, Office-Home, and CORe50 experiments are conducted on an NVIDIA RTX 4090 GPU, while the DomainNet experiments are performed on an RTX A5000 GPU. The images are resized to 224x224 pixels for all experiments. The transformation modules are two-layer MLPs with a hidden layer dimension of 1024. Regarding to the hyper-parameters, the rank for LoRA matrices is set to 8, λ_1 and λ_2 are set to 5 and 2, respectively. The margin is set to 1. ς is set to 0.9. The learning rate for the LoRAs and the temporary classifier is set to 0.02. The learning rates for the domain classifiers and transformation modules are set to 2×10^{-3} and 1×10^{-4} , respectively. We use the SGD optimizer for training the model with a batch size of 64. The first 6 blocks of the network are dedicated to the task-shared LoRAs, and the next 6 layers are utilized for the task-specific LoRAs.

D. Numerical Results

Table I reports the accuracy of CONEC-LoRA and other consolidated algorithms. It is clearly seen that our algorithm outperforms other algorithms across all datasets with significant margins. In the CDDb case, our algorithm outperforms DualCP by 6%, while CONEC-LoRA surpasses DCE by 2% in the office-home problem. Ours is also superior in the CORe50 problem, with an over 3% margin, despite the fact that the test samples belong to 3 unseen domains rather than the 8 domains the network is trained on, which demonstrates the generalization capability of CONEC-LoRA. Last but not least, our method outperforms DCE with over 2% margins in the challenging DomainNet problem. This finding clearly demonstrates the advantage of our approach based on the combination of task-shared and task-specific LoRAs, the concept of the stochastic classifier, and the introduction of an auxiliary network for domain ID predictions. On the other hand, we find that the gap between that with and without an oracle is negligible, i.e., $< 1\%$. This confirms the efficacy of our auxiliary network, which produces reliable domain ID predictions for the selection and isolation of task-specific LoRA. Such a parameter isolation approach assures that task-specific knowledge can be maintained and effectively combined with shareable information to deliver accurate classifications.

E. Domain Classification Accuracies

Since our approach relies on the auxiliary network for the parameter selection strategy, Table II compares the domain classification accuracy of CONEC-LoRA, SOYO [40], and S-iPrompt [56]. Note that both SOYO and S-iPrompt utilize an external network for parameter isolation strategies. The advantage of our auxiliary network is clearly reported in Table II, where it outperforms other methods with significant margins. It beats SOYO by a 5% margin in the CDDb problem while surpassing SOYO with 3% gap in the domainNet problem. We don't show the domain classification accuracies of other methods in the Office-Home problem because they are absent in their original papers. On the other hand, our auxiliary network performs well in the Office-Home problem, where it attains 83.72% accuracy. Note that the novelty of our auxiliary network is found in at least 3 facets: 1) it applies the concept of ball-generator loss to cope with the synthetic sample-bias problem; 2) it is structured under the different-depth network structure to exploit the strength of intermediate representations; 3) it makes use of the transformation module under a fixed backbone network.

F. Ablation Studies

Table III reports the numerical results of our ablation studies, where the contribution of different components is analyzed. It is observed that the use of the cosine classifier in lieu of the stochastic classifier reduces the domain incremental learning accuracy. A significant performance deterioration is observed when changing the classifier to the linear classifier, i.e., CONEC-LoRA suffers from over 10% loss of performance. This finding confirms our claim that the stochastic classifier plays a key role in our algorithm. On the other hand, the importance of the shared LoRAs is perceived by configuring our method with only the task-specific LoRA. CONEC-LoRA suffers from performance degradation using only the task-specific LoRA without the shared LoRA, i.e., the final accuracy deteriorates, showing the CF problem due to the absence of task-invariant features.

We also evaluate the performance of our auxiliary classifier without the ball generator loss. The removal of ball generator loss results in a decline in domain classification accuracy by about 4%, ultimately affecting the domain incremental learning accuracy. On the other hand, using a single classifier instead of multiple local classifiers, i.e., the different-depth network structure, reduces domain classification accuracy. The different-depth network structure fully exploits the intermediate representations in identifying the correct domain labels.

G. Robustness Analysis

We evaluate the robustness of the CONEC-LoRA to the variations of λ_1 and λ_2 , which are crucial coefficients to control the contributions of loss functions. Table IV shows the robustness of the proposed method to the variations of each coefficient while keeping the other fixed. It is seen that CONEC-LoRA is not sensitive to different trade-off constants λ_1, λ_2 where variations of these parameters do not significantly affect the performances of CONEC-LoRA. This finding

TABLE I: Comparison of the average and last accuracies of different methods. The ViT-B/16 IN1K is used as a backbone. Other backbones are mentioned after the method name. Methods with † indicate implemented with exemplars (10 per class). Note that since the train and test domains of the CORE50 dataset do not overlap, the ground truth domain labels are not available, and the accuracy with the oracle cannot be calculated.

Method	Venue	Office-Home		DomainNet		CORE50		CDDb-Hard	
		Avg	Last	Avg	Last	Avg	Last	Avg	Last
Fine-tune [41]	CVPR 2025	78.32 \pm 3.28	76.16 \pm 1.39	28.17 \pm 6.47	38.82 \pm 7.65	75.44 \pm 1.68	76.19 \pm 2.36	52.08 \pm 1.35	50.11 \pm 1.62
Replay† [63]	Psychology Review	84.23 \pm 2.31	83.75 \pm 0.68	64.78 \pm 2.98	61.16 \pm 1.19	85.56 \pm 0.38	92.21 \pm 0.63	66.91 \pm 18.0	63.21 \pm 11.6
iCaRL† [62]	CVPR 2017	83.2 \pm 3.2	83.5 \pm 0.7	58.4 \pm 5.3	55.9 \pm 1.7	71.0 \pm 3.7	76.0 \pm 2.7	57.5 \pm 9.9	54.4 \pm 9.9
MEMO† [61]	ICLR 2022	78.7 \pm 3.7	79.1 \pm 1.2	57.8 \pm 6.7	56.4 \pm 1.4	66.0 \pm 2.7	68.2 \pm 1.7	66.0 \pm 2.7	68.2 \pm 1.7
SimpleCIL [60]	IJCV 2025	75.2 \pm 4.9	76.2 \pm 0.0	41.1 \pm 6.8	40.6 \pm 0.0	62.5 \pm 1.6	67.2 \pm 0.0	65.5 \pm 3.2	64.1 \pm 1.7
L2P [29]	CVPR 2022	78.7 \pm 4.2	80.5 \pm 0.6	48.5 \pm 6.8	45.2 \pm 2.3	72.3 \pm 1.3	81.7 \pm 0.5	67.3 \pm 5.3	65.0 \pm 6.5
DualPrompt [30]	ECCV 2022	77.2 \pm 3.1	79.1 \pm 0.6	52.3 \pm 9.7	50.9 \pm 4.3	71.0 \pm 4.5	77.7 \pm 1.0	66.9 \pm 5.8	65.6 \pm 2.6
CODA-Prompt [57]	CVPR 2023	82.4 \pm 3.8	83.3 \pm 0.3	47.6 \pm 6.1	45.1 \pm 1.2	72.8 \pm 1.2	81.4 \pm 1.1	67.9 \pm 6.5	66.6 \pm 2.8
EASE [59]	CVPR 2024	81.16 \pm 3.52	76.33 \pm 2.16	50.50 \pm 2.27	43.72 \pm 1.70	86.30 \pm 0.04	87.02 \pm 1.21	67.78 \pm 2.44	64.96 \pm 8.36
RanPAC [58]	NeurIPS 2023	83.4 \pm 3.8	83.3 \pm 0.3	57.8 \pm 5.5	56.1 \pm 0.6	76.7 \pm 1.3	78.4 \pm 1.7	61.7 \pm 4.3	62.5 \pm 2.8
S-iPrompt [56]	NeurIPS 2022	81.40 \pm 3.3	80.80 \pm 0.20	59.00 \pm 6.8	57.90 \pm 0.2	62.7 \pm 2.2	65.8 \pm 0.9	64.2 \pm 5.8	63.4 \pm 2.8
DCE [43]	ICML 2025	84.6 \pm 3.0	84.4 \pm 0.2	64.3 \pm 6.0	63.5 \pm 0.5	80.1 \pm 0.70	84.8 \pm 0.30	74.6 \pm 6.5	71.8 \pm 4.2
DualCP [39]	AAAI 2025	N/A	N/A	N/A	60.13	N/A	88.10	N/A	82.16
CONEC-LoRA		86.29 \pm 1.53	86.43 \pm 0.38	66.79 \pm 2.94	66.42 \pm 0.38	88.88 \pm 0.54	90.24 \pm 1.64	88.43 \pm 2.21	88.21 \pm 0.88
CONEC-LoRA with the oracle		86.98 \pm 1.55	87.63 \pm 0.20	68.16 \pm 2.96	68.48 \pm 0.22	N/A	N/A	88.98 \pm 2.24	89.68 \pm 0.68

TABLE II: Domain classification accuracy

Method	Venue	CDDb-Hard	DomainNet	Office-Home
S-iPrompts [56], [40]	NeurIPS 2022	81.79	80.35	N/A
PINA-D+SOYO [40]	CVPR 2025	89.84	85.37	N/A
CONEC-LoRA		94.70 \pm 2.09	88.52 \pm 2.58	83.72 \pm 5.74

TABLE III: Ablation study on the CDDb-Hard dataset. The columns show the average and last domain-incremental learning (DIL) accuracies for modifications to the domain-incremental classification module, as well as the domain classification accuracy (DC) for modifications to the domain classification module, across five different orders of domains.

Description	DIL Average	DIL Last	Description	DC Average
CONEC-LoRA	88.43 \pm 2.21	88.21 \pm 0.88	CONEC-LoRA	94.73 \pm 2.08
Cosine classifier instead of the stochastic classifier	87.42 \pm 3.18	87.80 \pm 2.40	Without the ball-generator loss	90.53 \pm 3.66
Linear classifier instead of the stochastic classifier	81.26 \pm 3.47	77.96 \pm 6.74	Without the intermediate domain classifiers	94.31 \pm 2.46
Only task-specific LoRAs	88.79 \pm 2.54	87.62 \pm 3.80		

TABLE IV: Robustness of CONEC-LoRA to the variations of λ_1 coefficient on the domain-incremental learning (DIL) accuracy and λ_2 coefficient on the domain classification (DC) accuracy on the CDDb-Hard dataset.

λ_1	DIL Average	DIL Last	λ_2	DC Average
3	88.36 \pm 2.24	88.23 \pm 1.12	1	94.23 \pm 2.32
4	88.43 \pm 2.21	88.23 \pm 0.94	1.5	94.66 \pm 2.09
5	88.43 \pm 2.21	88.21 \pm 0.88	2	94.73 \pm 2.08
6	88.49 \pm 2.24	88.20 \pm 1.03	2.5	94.83 \pm 2.05
7	88.57 \pm 2.22	88.34 \pm 1.08	3	94.96 \pm 2.00

confirms the robustness of CONEC-LoRA against different trade-off constants steering the contributions of loss functions. For the sake of convenience, the parameters are simply set to $\lambda_1 = 5$, $\lambda_2 = 2$ for our experiments.

H. UMAP Analysis

Fig. 3 exhibits the UMAP [64] graphs of the network on 5 domains of the CDDb dataset before training and after training with CONEC-LoRA. It is clearly seen that before the training process, the embeddings are scattered in the feature space, and the pre-trained model seems to be confused in distinguishing samples of different domains and classes. After the training process, CONEC-LoRA exhibits strong discriminative features, where samples of different domains are projected into

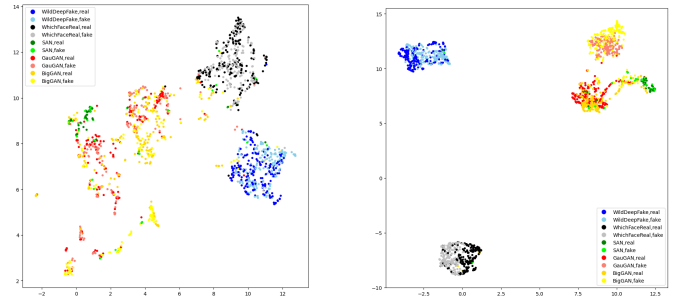


Fig. 3: UMAP plots on the embedding space of the model before training on the left and after training on the right on five domains of the CDDb-Hard dataset.

different locations of the latent spaces. Besides, our algorithm is capable of separating different classes within the domains. This finding also confirms that our approach is relatively robust to the CF problem, where our model recognizes all sequentially presented domains well, i.e., the UMAP plot refers to the model after training on the last domain.

VI. CONCLUSION

This paper proposes continual knowledge consolidation low rank adaptation (CONEC-LoRA) as a solution to domain incremental learning (DIL) problems. CONEC-LoRA resolves

the three shortcomings of existing methods: parameter isolation, inaccurate selection of task-specific parameters, and the use of traditional classifiers. CONEC-LoRA features a synergy between shared LoRA and task-specific LoRA while putting forward the auxiliary domain classifier, which includes the ball-generator loss, the different-depth network structure, and the transformation module. It integrates the concept of stochastic classifiers having class mean and variance vectors drawn from a distribution. The efficacy of CONEC-LoRA has been rigorously evaluated with the four benchmark problems of DIL, where it beats recently published works across all problems with over 5% margins. Its domain classification accuracies are also superior to existing works, while the ablation study substantiates the advantage of each learning component. The robustness analysis confirms that our algorithms are not sensitive to variations of loss coefficients. Our future work is devoted to studying the data scarcity problem in DIL.

REFERENCES

- [1] Z. Chen and B. Liu, *Lifelong machine learning*. Springer, 2018, vol. 1.
- [2] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural networks : the official journal of the International Neural Network Society*, vol. 113, pp. 54–71, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:73497737>
- [3] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. van de Weijer, "Class-incremental learning: Survey and performance evaluation on image classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, pp. 5513–5533, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:234353728>
- [4] D.-W. Zhou, Q. Wang, Z. Qi, H.-J. Ye, D. chuan Zhan, and Z. Liu, "Class-incremental learning: A survey," *IEEE transactions on pattern analysis and machine intelligence*, vol. PP, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:256627357>
- [5] D. Kim and B. Han, "On the stability-plasticity dilemma of class-incremental learning," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20 196–20 204, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:257921372>
- [6] J. Kirkpatrick, R. Pascanu, N. C. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, pp. 3521 – 3526, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:4704285>
- [7] F. Mao, W. Weng, M. Pratama, and E. K. Y. Yapp, "Continual learning via inter-task synaptic mapping," *ArXiv*, vol. abs/2106.13954, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:233710658>
- [8] I. Paik, S. Oh, T. Kwak, and I. Kim, "Overcoming catastrophic forgetting by neuron-level plasticity control," in *AAAI Conference on Artificial Intelligence*, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:199001153>
- [9] F. Zenke, B. Poole, and S. Ganguli, "Continual learning through synaptic intelligence," *Proceedings of machine learning research*, vol. 70, pp. 3987–3995, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:10409742>
- [10] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuytelaars, "Memory aware synapses: Learning what (not) to forget," pp. 139–154, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:4254748>
- [11] X. Li, Y. Zhou, T. Wu, R. Socher, and C. Xiong, "Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting," *ArXiv*, vol. abs/1904.00310, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:90259576>
- [12] J. Schwarz, W. M. Czarnecki, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R. Hadsell, "Progress & compress: A scalable framework for continual learning," *ArXiv*, vol. abs/1805.06370, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:21718339>
- [13] S. Cha, H. Hsu, F. du Pin Calmon, and T. Moon, "Cpr: Classifier-projection regularization for continual learning," *ArXiv*, vol. abs/2006.07326, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:219636462>
- [14] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny, "Efficient lifelong learning with a-gem," *ArXiv*, vol. abs/1812.00420, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:54443381>
- [15] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5533–5542, 2016. [Online]. Available: <https://api.semanticscholar.org/CorpusID:206596260>
- [16] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. S. Torr, and M. Ranzato, "On tiny episodic memories in continual learning," *arXiv: Learning*, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:173188188>
- [17] A. Chaudhry, A. Gordo, P. K. Dokania, P. H. S. Torr, and D. Lopez-Paz, "Using hindsight to anchor past knowledge in continual learning," in *AAAI Conference on Artificial Intelligence*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:210957697>
- [18] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara, "Dark experience for general continual learning: a strong, simple baseline," *ArXiv*, vol. abs/2004.07211, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:215768806>
- [19] T. Dam, M. Pratama, M. M. Ferdaus, S. G. Anavatti, and H. Abbas, "Scalable adversarial online continual learning," in *ECML/PKDD*, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:252089827>
- [20] M. V. de Carvalho, M. Pratama, J. Zhang, and Y. San, "Class-incremental learning via knowledge amalgamation," in *ECML/PKDD*, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:252090326>
- [21] M. A. Ma'sum, M. Pratama, E. D. Lughofer, W. Ding, and W. Jatmiko, "Assessor-guided learning for continual environments," *Inf. Sci.*, vol. 640, p. 119088, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:257636622>
- [22] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," *Advances in neural information processing systems*, vol. 30, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:37308416>
- [23] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *Neural Information Processing Systems*, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:1888776>
- [24] J. Xu, J. Ma, X. Gao, and Z. Zhu, "Adaptive progressive continual learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, pp. 6715–6728, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:235767757>
- [25] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong learning with dynamically expandable networks," *ArXiv*, vol. abs/1708.01547, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:3693512>
- [26] A. Ashfahani and M. Pratama, "Unsupervised continual learning in streaming environments," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, pp. 9992–10 003, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:237572299>
- [27] M. Pratama, A. Ashfahani, and E. D. Lughofer, "Unsupervised continual learning via self-adaptive deep clustering approach," in *CSSL*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:235658071>
- [28] A. Rakaraddi, S.-K. Lam, M. Pratama, and M. V. de Carvalho, "Reinforced continual learning for graphs," *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:252089650>
- [29] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. G. Dy, and T. Pfister, "Learning to prompt for continual learning," *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 139–149, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:245218925>
- [30] Z. Wang, Z. Zhang, S. Ebrahimi, R. Sun, H. Zhang, C.-Y. Lee, X. Ren, G. Su, V. Perot, J. G. Dy, and T. Pfister, "Dualprompt: Complementary prompting for rehearsal-free continual learning," *ArXiv*, p. 631–648, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:248085201>
- [31] M. A. Ma'sum, M. Pratama, S. Ramasamy, L. Liu, Habibullah, and R. Kowalczyk, "Vision and language synergy for rehearsal free continual learning," in *International Conference on Learning Representations*,

2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:278602871>
- [32] S. Momeni, S. Mazumder, and B. Liu, "Continual learning using a kernel-based method over foundation models," in *AAAI Conference on Artificial Intelligence*, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:274965301>
- [33] D.-W. Zhou, H.-L. Sun, J. Ning, H.-J. Ye, and D. chuan Zhan, "Continual learning with pre-trained models: A survey," in *International Joint Conference on Artificial Intelligence*, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:267312447>
- [34] J. He, Z. Duan, and F. M. Zhu, "CI-lora: Continual low-rank adaptation for rehearsal-free class-incremental learning," *ArXiv*, vol. abs/2505.24816, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:279070575>
- [35] X. Liu and X. Chang, "Lora subtraction for drift-resistant space in exemplar-free continual learning," *ArXiv*, vol. abs/2503.18985, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:277313510>
- [36] T. Fukuda, H. Kera, and K. Kawamoto, "Adapter merging with centroid prototype mapping for scalable class-incremental learning," *ArXiv*, vol. abs/2412.18219, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:274992202>
- [37] G. M. van de Ven and A. S. Tolias, "Three scenarios for continual learning," *ArXiv*, vol. abs/1904.07734, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:119309522>
- [38] Z. Liu, Y. Peng, and J. Zhou, "Compositional prompting for anti-forgetting in domain incremental learning," *Int. J. Comput. Vis.*, vol. 132, pp. 5783–5800, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:270778693>
- [39] Q. Wang, Y. He, S. Dong, X. Song, J. Han, H. Luo, and Y. Gong, "Dualcp: Rehearsal-free domain-incremental learning via dual-level concept prototype," in *AAAI Conference on Artificial Intelligence*, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:277272564>
- [40] Q. Wang, X. Song, Y. He, J. Han, C. Ding, X. Gao, and Y. Gong, "Boosting domain incremental learning: Selecting the optimal parameters is all you need," *ArXiv*, vol. abs/2505.23744, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:278996449>
- [41] D.-W. Zhou, Z.-W. Cai, H.-J. Ye, L. Zhang, and D.-C. Zhan, "Dual consolidation for pre-trained model-based domain-incremental learning," *ArXiv*, vol. abs/2410.00911, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:273022628>
- [42] N. Paedeh, M. Pratama, S. Wibirama, W. Mayer, Z. Cao, and R. Kowalczyk, "Few-shot class incremental learning via robust transformer approach," *Inf. Sci.*, vol. 675, p. 120751, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:269741067>
- [43] L. Li, D.-W. Zhou, H.-J. Ye, and D.-C. Zhan, "Addressing imbalanced domain-incremental learning through dual-balance collaborative experts," in *Forty-second International Conference on Machine Learning*, 2025. [Online]. Available: <https://openreview.net/forum?id=dwjwvTwV3V>
- [44] K. Xu, X. Zou, G. Hua, and J. Zhou, "Componential prompt-knowledge alignment for domain incremental learning," in *Forty-second International Conference on Machine Learning*, 2025. [Online]. Available: <https://openreview.net/forum?id=QWCdBzLOsk>
- [45] J. E. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *ArXiv*, vol. abs/2106.09685, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:235458009>
- [46] J. Kalla and S. Biswas, "S3c: Self-supervised stochastic classifiers for few-shot class-incremental learning," in *European Conference on Computer Vision*, 2022.
- [47] J. Zhu, K. H. Greenewald, K. Nadjahi, H. S. de Oc'ariz Borde, R. B. Gabrielsson, L. Choshen, M. Ghassemi, M. Yurochkin, and J. Solomon, "Asymmetry in low-rank adapters of foundation models," *ArXiv*, vol. abs/2402.16842, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268033026>
- [48] G. Golub and W. Kahan, "Calculating the singular values and pseudo-inverse of a matrix," *Journal of the Society for Industrial and Applied Mathematics, Series B: Numerical Analysis*, vol. 2, no. 2, pp. 205–224, 1965.
- [49] F. Szatkowski, Y. Zheng, F. Yang, B. Twardowski, T. Trzeciński, and J. van de Weijer, "Improving continual learning performance and efficiency with auxiliary classifiers," 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268363428>
- [50] C. M. Bishop and N. M. Nasrabadi, *Pattern recognition and machine learning*. Springer, 2006, vol. 4, no. 4.
- [51] M. A. Ma'sum, M. Pratama, L. Liu, E. D. Lughofer, Habibullah, and R. Kowalczyk, "Few-shot continual learning via flat-to-wide approaches," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 36, pp. 8966–8978, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259251661>
- [52] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, "Moment matching for multi-source domain adaptation," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1406–1415.
- [53] V. Lomonaco and D. Maltoni, "Core50: a new dataset and benchmark for continuous object recognition," in *Conference on robot learning*. PMLR, 2017, pp. 17–26.
- [54] C. Li, Z. Huang, D. P. Paudel, Y. Wang, M. Shahbazi, X. Hong, and L. Van Gool, "A continual deepfake detection benchmark: Dataset, methods, and essentials," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2023, pp. 1339–1349.
- [55] H. Venkateswara, J. Eusébio, S. Chakraborty, and S. Panchanathan, "Deep hashing network for unsupervised domain adaptation," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5385–5394, 2017. [Online]. Available: <https://api.semanticscholar.org/CorpusID:2928248>
- [56] Y. Wang, Z. Huang, and X. Hong, "S-prompts learning with pre-trained transformers: An occam's razor for domain incremental learning," *ArXiv*, vol. abs/2207.12819, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:251066766>
- [57] J. Smith, L. Karlinsky, V. Gutta, P. Cascante-Bonilla, D. Kim, A. Arbelle, R. Panda, R. S. Feris, and Z. Kira, "Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning," *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11909–11919, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:253801593>
- [58] M. D. McDonnell, D. Gong, A. Parvaneh, E. Abbasnejad, and A. van den Hengel, "Ranpac: Random projections and pre-trained models for continual learning," *ArXiv*, vol. abs/2307.02251, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259342681>
- [59] D.-W. Zhou, H.-L. Sun, H.-J. Ye, and D. chuan Zhan, "Expandable subspace ensemble for pre-trained model-based class-incremental learning," *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 23554–23564, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268531460>
- [60] D.-W. Zhou, Z.-W. Cai, H.-J. Ye, D.-C. Zhan, and Z. Liu, "Revisiting class-incremental learning with pre-trained models: Generalizability and adaptivity are all you need," *International Journal of Computer Vision*, vol. 133, no. 3, pp. 1012–1032, 2025.
- [61] D.-W. Zhou, Q.-W. Wang, H.-J. Ye, and D.-C. Zhan, "A model or 603 exemplars: Towards memory-efficient class-incremental learning," *arXiv preprint arXiv:2205.13218*, 2022.
- [62] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.
- [63] R. Ratcliff, "Connectionist models of recognition memory: constraints imposed by learning and forgetting functions," *Psychological review*, vol. 97, no. 2, p. 285, 1990.
- [64] L. McInnes, J. Healy, and J. Melville, "Umap: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.