

# CONTINUAL LEARNING VIA LEARNING CONTINUAL MEMORY IN VISION TRANSFORMER

**Anonymous authors**  
 Paper under double-blind review

## ABSTRACT

This paper explores continual learning (CL) using Vision Transformer (ViT) in streaming tasks under the challenging exemplar-free class-incremental (ExfCCL) setting. We formulate ExfCCL as a learning problem consisting of two key sub-systems: (i) task ID inference for test data, which selects appropriate task-specific head classifiers to accounting for varying class distributions across tasks and streams, and (ii) a dynamic learning-to-grow feature backbone that balances stability and plasticity, mitigating catastrophic forgetting through task synergies. Following the common protocol that the first task can train a ViT sufficiently well as the base model, we address these sub-systems from a continual memory learning perspective. To support task ID inference, we utilize an external memory mechanism that maintains task centroids computed by the base ViT throughout CL. For the feature backbone, we identify optimal placements for internal (parameter) memory to enable a dynamic, task-synergy guided growing feature backbone. We propose a Hierarchical Exploration-Exploitation (HEE) sampling-based neural architecture search (NAS) method that effectively learns task synergies by continually and structurally updating internal memory with four basic operations: *reuse*, *adapt*, *new*, and *skip*. Our approach, dubbed **Continual Hierarchical-Exploration-Exploitation Memory (CHEEM)**, is evaluated on the challenging Visual Domain Decathlon (VDD) and ImageNet-R benchmarks, demonstrating its effectiveness.

## 1 INTRODUCTION

Developing continual learning machines is a key objective in Artificial Intelligence (AI), aiming to replicate human-like adaptability and the ability to learn-to-learn, enabling proficiency in streaming tasks. Despite their advances, state-of-the-art Deep Neural Networks (DNNs) still lack true biological intelligence in the realm of continual learning and are particularly hindered by the critical issue of *catastrophic forgetting* when exposed to streaming tasks in dynamic environments (McCloskey & Cohen, 1989; Thrun & Mitchell, 1995).

To address catastrophic forgetting, two primary categories of continual learning methods have emerged: exemplar-based methods (Aljundi et al., 2019b; Hayes et al., 2019; Wu et al., 2019) and exemplar-free methods (Kirkpatrick et al., 2017; Li et al., 2019; Wang et al., 2022d;c;a). While both have shown promising progress, exemplar-free methods are especially appealing due to their ability to learn new tasks without retaining any data from previous tasks. Furthermore, continual learning has evolved from the traditional task-incremental protocol, where task IDs of test data are available during inference, to the more challenging class-incremental protocol, in

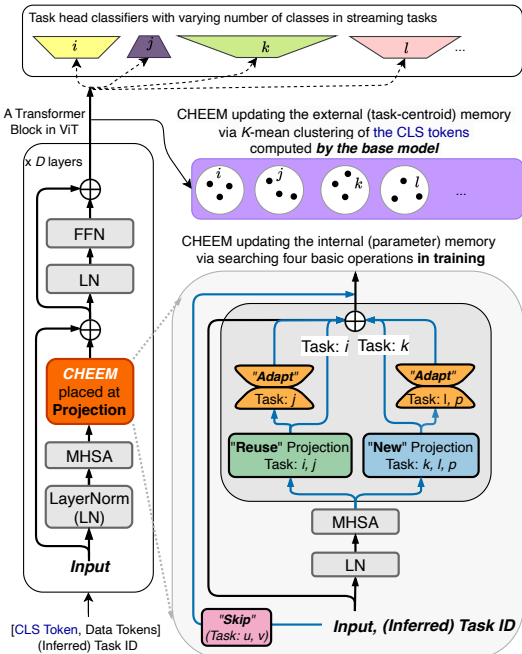


Figure 1: Illustration of the proposed CHEEM.

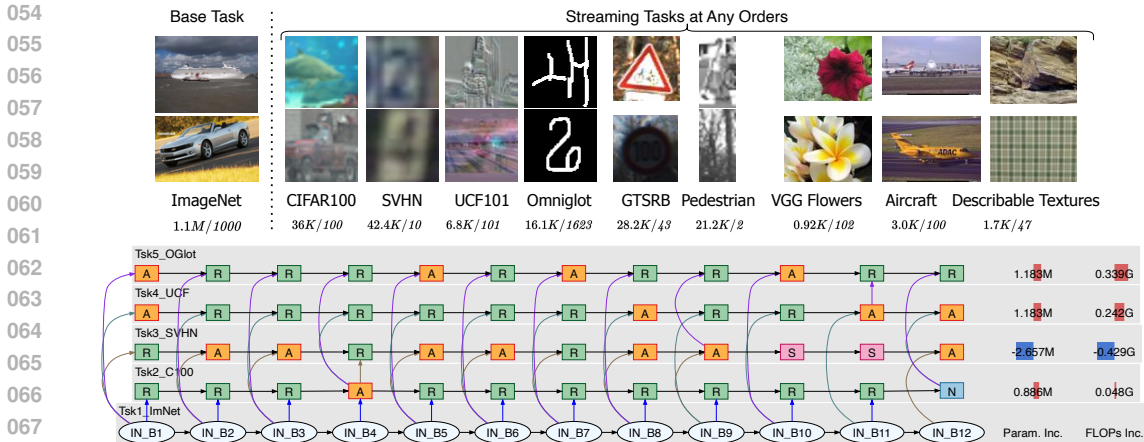


Figure 2: **Top:** The challenging VDD benchmark (Rebuffi et al., 2017a) consisting of 10 tasks of different nature with #training images/#classes significantly varying across different tasks. **Bottom:** Starting from the base task (Tsk1\_ImNet) trained ViT (Dosovitskiy et al., 2021) (e.g., the 12-layer ViT-B model consisting of IN\_B1 to IN\_B12), together with the best performance on VDD in experiments, our CHEEM can learn sensible memory structures: e.g., Tsk3\_SVHN learns to skip both B10 and B11 with both mode complexity and FLOPs reduced. Tsk4\_UCF and Tsk5\_OGlot learn to adapt the first layers to account for the domain shifts. We show the first 5 tasks for clarity (see more in the supplementary). The last two columns show the number of new task-specific parameters and added FLOPs, in comparison with the base model, subject to **R** reuse, **A** dapt, **N** ew and **S** kip.

which task IDs are unknown during inference. This paper focuses on exemplar-free class-incremental continual learning (ExfCCL).

Recently, the emergence of powerful pretrained Transformer models (Vaswani et al., 2017; Dosovitskiy et al., 2021; Radford et al., 2021) has driven significant interest in ExfCCL using pretrained and frozen Vision Transformers (ViTs) (Dosovitskiy et al., 2021), primarily explored through the lens of prompt-tuning or prefix-tuning (Wang et al., 2022d;c;a; Smith et al., 2023a). However, this approach presents two main drawbacks: (i) A single frozen pretrained Transformer backbone cannot accommodate all streaming tasks of diverse nature, such as those found in the VDD benchmark (Rebuffi et al., 2017a), shown in Fig. 2. While this method maximizes the stability of the feature backbone, it relies on input or prefix prompts to address plasticity by explicitly leveraging the self-attention mechanisms. (ii) Due to the quadratic complexity of ViTs concerning the number of input tokens, incorporating task-specific input prompts or layer-wise prefix prompts into a frozen pretrained Transformer backbone significantly increases computational demands. Continual learning where computational requirements always grow for new tasks by design regardless of task complexity fails to reflect true intelligence.

The challenge of structurally and dynamically updating a pretrained Transformer backbone for ExfCCL remains an open problem, as it requires achieving a balance between the backbone’s stability and plasticity while enabling dynamic computation and mitigating catastrophic forgetting.

Moreover, ExfCCL involves tackling the challenge of designing task head classifiers, which can either explicitly infer task IDs to select task-specific head classifiers (Wang et al., 2022a) or share a growing, common head classifier (Wang et al., 2022d;c; Smith et al., 2023a). The former approach is straightforward and can help in learning task-specific components and enhancing the plasticity of the backbone, but may suffer from low average precision of the task ID inference for certain tasks. In contrast, the latter approach resorts to more sophisticated prefix prompt tuning to induce plasticity in the feature backbone and can create a discrepancy between training, where head segments from previous tasks are masked out in the softmax function of loss computation, and testing, where the entire head is used to infer class labels. This discrepancy can lead to instability when dealing with streaming tasks that involve a varying number of classes among tasks, such as those encountered in the VDD benchmark.

**Our aim** in this paper is to study ExfCCL using ViTs. As illustrated in Fig. 1, we formulate it as a problem of learning continual memory in ViT, which has two components: (i) The internal

memory enabling a dynamic learning-to-grow feature backbone that balances stability and plasticity, mitigating catastrophic forgetting through **task synergies, in which a new task learns automatically to reuse/adapt modules from previous similar tasks, or to introduce new modules when needed, or to skip some modules when it appears to be an easier task** (see the bottom of 2). The internal parameter memory learning presents alternative perspectives to the input and prefix prompting based methods (Wang et al., 2022d;c; Smith et al., 2023a; Wang et al., 2022a). (ii) The external memory enabling task ID inference for test data, for which we adopt a method proposed in (Wang et al., 2022a) for its simplicity.

We follow the common protocol in the prior art that the first task (e.g., ImageNet-1k (Russakovsky et al., 2015)) can train a ViT sufficiently well as the base model. To enable a dynamic, learning-to-grow feature backbone starting from the base model, we identify and provide simple yet effective solutions for two key challenges: **(i) Which modules in a ViT should be reused, adapted, newly created, or skipped for a new task in ExfCCL?** It is computationally impractical, and counterproductive to the stability-plasticity trade-off, to make all ViT components dynamic (as opposed to a fully frozen model). We designate the output projection layer following multi-head self-attention (MHSA) as *the task-synergy internal (parameter) memory* that will be structurally and dynamically updated. **(ii) How can we represent and continually learn this task-synergy internal memory to enable dynamic memory structures across streaming tasks?** We propose organizing the memory using a mixture of experts (MoEs) similar in spirit to (Riquelme et al., 2021), as depicted in Fig. 1 and illustrated through the “task-factorized” visualization in Fig. 2. To learn to select the optimal synergy operation from the four choices (*reuse*, *adapt*, *new*, and *skip*), we introduce an effective hierarchical exploration-exploitation (HEE) sampling-based neural architecture search (NAS) method. This approach is inspired by the task-incremental learn-to-grow method (Li et al., 2019). Our proposed method is termed **CHEEM** (*Continual Hierarchical-Exploration-Exploitation Memory*).

**Our Contributions.** This paper makes three main contributions to the field of exemplar-free class-incremental continual learning (ExfCCL) using ViT. (i) It presents a hierarchical task-synergy exploration-exploitation sampling based NAS method for learning task-aware dynamic models continually with respect to four operations: *Skip*, *Reuse*, *Adapt*, and *New* to mitigate catastrophic forgetting. (ii) It identifies a “sweet spot” in ViT as the task-synergy internal (parameter) memory, i.e., the output projection layers after MHSA in ViT. It also presents a new usage for the class-token CLS in ViT as the internal memory updating guidance, in addition to leveraging it in maintaining the external (task-centroid) memory for task ID inference on the fly. (iii) It is the first work, to the best of our knowledge, to evaluate continual learning with ViTs on the large-scale, diverse and imbalanced VDD benchmark (Rebuffi et al., 2017a), with better performance than the prior art.

## 2 APPROACH

### 2.1 PROBLEM FORMULATION OF CHEEM IN EXFCCL

We start with a vanilla  $D$ -layer ViT model (e.g., the 12-layer ViT-Base) (Dosovitskiy et al., 2021). The left of Fig. 1 shows a ViT block. Denote by  $x_{L,d}$  an input sequence consisting of  $L$  tokens encoded in a  $d$ -dimensional space. In ViTs, the first token is the so-called class-token, CLS. The remaining  $L - 1$  tokens are formed by patchifying an input image and then embedding patches, together with additive positional encoding. A ViT block is defined by,

$$z_{L,d} = x_{L,d} + \text{Proj}(\text{MHSA}(\text{LN}_1(x_{L,d}))), \quad (1)$$

$$y_{L,d} = z_{L,d} + \text{FFN}(\text{LN}_2(z_{L,d}))), \quad (2)$$

where  $\text{LN}(\cdot)$  represents the layer normalization (Ba et al., 2016), and  $\text{Proj}(\cdot)$  is a linear transformation fusing the multi-head outputs from MHSA module. The FFN is often implemented by a multi-layer perceptron (MLP) with a feature expansion layer  $\text{MLP}^u$  and a feature reduction layer  $\text{MLP}^d$  with a nonlinear activation function (such as the GELU (Hendrycks & Gimpel, 2016)) in the between.

Denote by  $\mathcal{T} = \{T_1, \dots, T_t, \dots, T_N\}$  a stream of tasks in continual learning, where each task  $T_t$  consists of a training set  $D_t^{\text{train}}$  and a testing set  $D_t^{\text{test}}$ . We make no restrictive assumptions regarding the nature, order, or number of classes in streaming tasks, either per task or in total. For example, there are 9 diverse, streaming tasks in the VDD benchmark (2).

Denote by  $(f_1, C_1)$  the ViT (Dosovitskiy et al., 2021) base model (e.g., the 12-layer ViT-Base) trained on the first task  $T_1$  (e.g., ImageNet-1k (Russakovsky et al., 2015)), where  $f_1$  is the backbone and  $C_1$  the head classifier. Denote by  $(\mathcal{F}_t, C_t)$  the sequentially and continually learned model after task  $T_t$

(for  $t \geq 1$ ).  $\mathcal{F}_t$  is the super ViT backbone structurally and dynamically updated from the base model  $f_1 = \mathcal{F}_1$ , and  $\mathcal{C}_t = \{C_1, \dots, C_t\}$  consisting of task-specific head classifiers each of which is trained from scratch. Let  $\Theta_t = \mathcal{F}_t \setminus \mathcal{F}_{t-1}$  be task-specific backbone parameters, which we term **the internal parameter memory** in ExfCCL to exploit task synergies. We note that ExfCCL requires no use of exemplar data of previous tasks in any forms in learning  $(\Theta_t, C_t)$  for task  $T_t$ .

In inference, since the task IDs of test data are unknown. For a test sample  $x$ , we will need to infer its task ID on the fly. Following the prior art (Wang et al., 2022d;c; Smith et al., 2023a; Wang et al., 2022a), we use the base model  $f_1(\cdot)$  as the query function  $q(\cdot)$ , and use  $q(x)$  to retrieve the task ID from **the external memory** for which we adopt the method proposed in (Wang et al., 2022a) for its simplicity. With task ID available (provided in training or inferred in testing), we can allocate the task-specific backbone  $f_t \subset \mathcal{F}_N$  for task  $T_t$ . The task-specific model is then specified by  $(f_t, C_t)$ .

## 2.2 CONSTRUCTING THE EXTERNAL MEMORY

We choose to explicitly infer task IDs for test data (see Appendix A for the analyses). We adopt the method proposed in S-Prompts (Wang et al., 2022a). For a task  $t$ , we leverage  $K$ -mean clustering of the CLS tokens of its training images computed by the query function  $q(\cdot)$  (i.e., the base model  $f_0$ ). Denote by  $Z_t = \{z_t^1, \dots, z_t^K\}$  the clustered task centroids for task  $T_t$ , where for simplicity we use the same  $K = 5$  across tasks. We have the external memory,  $\Psi = \cup_{t=1}^N Z_t$  after training. In inference, for a testing sample  $x$ , we first compute its CLS token using the same query function  $q(x)$ , denoted by  $z$ . The task ID is then inferred via  $K$ -NN retrieval,  $KNN(z, \Psi)$ , either by retrieving the class ID of the Top-1 NN centroid in the external memory  $\Psi$ , or by majority voting of the class ID from the  $K$ -NN centroids.

We note that the such constructed external memory does not break the exemplar-free protocol since we do not use exemplars in any forms from previous tasks. We also note that we focus on the learning of the internal parameter memory. The decoupled design between the external memory and the internal memory will enable us to integrate and/or develop more advanced task ID inference approaches in future work, while potentially promoting our proposed CHEEM in the field of ExfCCL.

## 2.3 IDENTIFYING THE TASK-SYNERGY INTERNAL MEMORY

The proposed identification process is straightforward. Without introducing any modules handling forgetting, we compare both the task-to-task forward transferrability and the sequential forgetting for different components in a ViT block. **Our intuition is that a desirable component for placing the task-synergy parameter memory must enable strong transferrability with manageable forgetting, while being lightweight to account for the trade-off between stability and plasticity.**

To that end, we use the VDD benchmark (Rebuffi et al., 2017a) (see Fig. 2). We first train a ViT-Base (Dosovitskiy et al., 2021) on the first task, ImageNet (Russakovsky et al., 2015), as the base model  $f_1(\cdot)$ . To measure the task-to-task transferability, we *individually fine-tune*  $f_1$  in a task-to-task transfer learning manner for the remaining 9 streaming tasks. Let  $f_{t|1}$  be the backbone fine-tuned for task  $T_t$  (for  $t \geq 1$ ), and  $C_t$  the head classifier trained from scratch. The average Top-1 accuracy is:

$$AA = \frac{1}{N} \sum_{t=1}^N \text{Acc}(T_t; f_{t|1}, C_t) \quad (3)$$

where  $\text{Acc}()$  uses the Top-1 classification accuracy.

To measure the sequential forgetting, we *continually fine-tune* the backbone started from  $f_1$  on the 9 tasks in a randomly sampled and fixed streaming order (as shown in Fig. 2). Let  $f_{1:t}$  be the backbone trained sequentially and continually after task  $T_t$  and  $C_t$  is its head classifier. The average forgetting (Chaudhry et al., 2018) on the first  $N - 1$  streaming tasks is,

$$AF = \frac{1}{N - 1} \sum_{t=1}^{N-1} \left( \max_{j \in [t, N]} a_{j,t} - a_{N,t} \right), \quad (4)$$

where  $a_{j,t} = \text{Acc}(T_t; f_{1:j}, C_t)$ .

Table 1: Results of identifying the optimal placement of our CHEEM in ViT by testing 11 components (Eqns. 1 and 2).

Index	Finetuned Component	AA (Eqn. 3)	AF (Eqn. 4)
1	LN <sub>1</sub> + LN <sub>2</sub>	81.76	21.24
2	FFN	84.20	44.76
3	MLP <sup>d</sup>	83.66	37.99
4	LN <sub>2</sub>	80.04	16.35
5	MHSA + LN <sub>1</sub>	85.26	54.38
6	LN <sub>1</sub>	81.18	19.04
7	Query	81.57	19.69
8	Key	81.56	19.19
9	Query+Key	81.49	31.10
10	Value	84.99	37.58
11	Projection ( <b>CHEEM</b> )	85.11	30.50
Classifier w/ Frozen Backbone		70.78	-

As shown in Table 1, we compare 11 components or composite components in ViT. Consider the strong forward transfer ability, manageable forgetting, maintaining simplicity and for less invasive implementation in practice, **we select the Projection layer after the MHSA as the task-synergy internal (parameter) memory** (Fig. 1) to realize our proposed CHEEM for ExfCCL. We note that the last row of Table 1 shows the result of a conventional transfer learning setting in which only the head classifier is trained with the backbone frozen, which clearly shows that one frozen backbone can not fit all streaming tasks, as aforementioned and entailing dynamic updating of the backbone in continual learning.

## 2.4 LEARNING THE TASK-SYNERGY INTERNAL MEMORY

The proposed internal memory of our CHEEM is represented by a MoEs. Starting with the base ViT model  $f_1$ , the internal memory at the  $l$ -th layer in ViT consists of a single expert defined by a tuple,

$$E_l^{(1,\cdot)} = (\theta_l^{(1,\cdot)}, \mu_l^1), \tag{5}$$

where the subscript represents the layer index and the list-based superscript shows which task(s) use this expert.  $\theta_l^{(1,\cdot)}$  are the parameters of the projection layer and  $\mu_l^1 \in R^d$  is the associated mean class-token (CLS) pooled from the training dataset after the model is trained, which is task specific (as indicated by the superscript). For example, if an expert is reused by another task (say, 3) in continual learning, we will have  $E_l^{(1,3,\cdot)} = (\theta_l^{(1,3,\cdot)}, \mu_l^1, \mu_l^3)$ .

As shown in Fig. 3, for a new task  $t$ , **learning to update CHEEM consists of:** i) the Supernet construction, ii) the Supernet training, and iii) the target network selection and finetuning.

### 2.4.1 SUPERNET CONSTRUCTION VIA Reuse, Adapt, New AND Skip

For clarity, we consider a single layer  $l$  for a new task, and the current memory consists of two experts,  $\{E_l^{(1,\cdot)}, E_l^{(2,\cdot)}\}$  (Fig. 3, left). We utilize four operations in the Supernet construction:

- **R**euse the projection layer from a previous (similar) task for a new task.
- **A**dapt by adding a new lightweight layer on top of the projection layer of a previous task, implemented by a MLP with one squeezing hidden layer.
- **N**ew by adding a new projection layer, which enables the model to handle dissimilar situations.
- **S**kip the entire MHSA block, which encourages the adaptivity accounting for the diverse nature of tasks.

The Supernet is constructed by reusing and adapting each existing expert at layer  $l$ , and adding a new and a skip expert (the bottom of Fig. 3). The newly added adapt MLPs and new projection layers will be trained from scratch using the data of a new task only. The right-top of Fig. 3 shows the Adapt operation on top of  $E_l^{(2,\cdot)}$  is learned and added,  $E_l^{(3,\cdot)} = (\theta_l^{(3,\cdot)}, \mu_l^3)$  where  $\theta_l^{(3,\cdot)}$  represents parameters of the adapt MLPs learned for the task 3, and  $\mu_l^3$  is the mean CLS token pooled for the task 3. The expert  $E_l^{(2,\cdot)}$  is updated to  $E_l^{(2,3,\cdot)} = (\theta_l^{(2,3,\cdot)}, \mu_l^2)$  indicating its weights are shared with task 3. We provide implementation details for Adapt in the Appendix D.5.

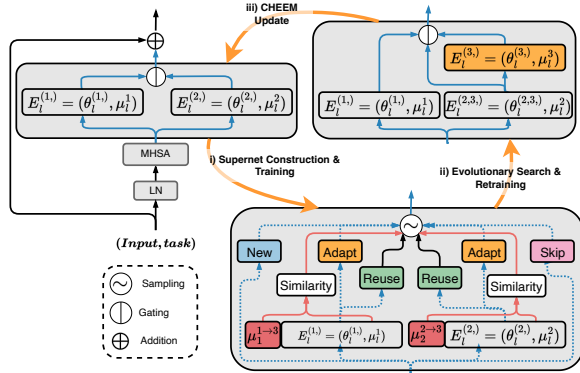


Figure 3: Illustration of CHEEM learning via NAS.

### 2.4.2 SUPERNET TRAINING VIA THE PROPOSED HEE SAMPLING-BASED NAS

To train the Supernet constructed for a new task  $t$ , we build on the SPOS method (Guo et al., 2020) due to its efficiency. The basic idea of SPOS is to train a single-path sub-network from the Supernet by sampling an expert at every layer in each mini-batch of training. One key aspect is the sampling strategy. The vanilla SPOS method uses uniform sampling (i.e., the *pure exploration* (PE) strategy, Fig. 4 left). **We propose an exploitation strategy** (Fig. 4 right), which utilizes a hierarchical sampling



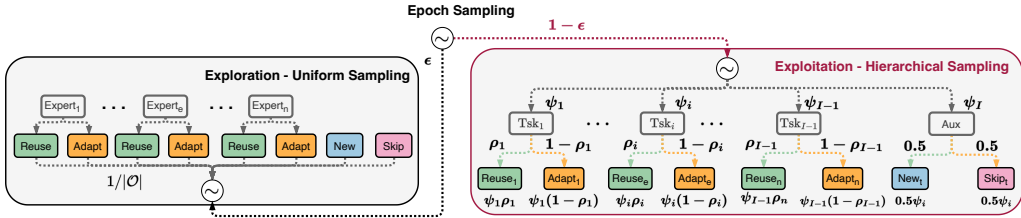


Figure 4: Illustration of the proposed hierarchical task-synergy exploration-exploitation (HEE) sampling based NAS. It integrates the vanilla exploration strategy (left) and the proposed exploitation strategy (right) with an epoch-wise scheduling.

method that forms the categorical distribution over the operations in the search space *explicitly based on task synergies computed based on the pooled task-specific CLS tokens*.

As illustrated in the right of Fig. 4, at each layer  $l$  in the Supernet, for a new task  $t$ , our proposed 3-level HEE sampling is realized by:

- **The epoch-level sampling:** At the beginning of an epoch in the Supernet training, we choose the pure exploration strategy with a probability of  $\epsilon_1$  (e.g., 0.3), and the exploitation strategy with a probability of  $1 - \epsilon_1$ .
- **The task-level sampling:** There are  $t - 1$  previous tasks, some of which may use `Skip` at the  $l$ -th layer and will not be used in sampling. We introduce an auxiliary task to handle sampling `New` and `Skip` for the current task  $t$ . *The task-level sampling is based on a categorical distribution*  $(\psi_1, \dots, \psi_i, \dots, \psi_{I-1}, \psi_I)$ , which is computed by the Softmax function over the similarity scores defined below, where  $I \leq t$ .
- **The operation-level sampling:** With a previous task  $i$  sampled with the probability  $\psi_i$ , we sample the `Reuse` and `Adapt` operation using a Bernoulli distribution  $(\rho_i, 1 - \rho_i)$ , where  $\rho_i$  is the success rate computed by the Sigmoid function of the task similarity score (defined below).

**The Task Similarity Score.** Consider one expert  $E_l^{(i,j)}$  at the  $l$ -th layer which is shared by two previous tasks  $i$  and  $j$  with their mean CLS tokens  $\mu_l^i$  and  $\mu_l^j$  respectively, we compute the mean CLS tokens using the models for task  $i$  and  $j$  on the training dataset  $D_t^{train}$  of the current task  $t$ ,  $\mu_l^{i \rightarrow t}$  and  $\mu_l^{j \rightarrow t}$ . The task similarity between task  $i$  and  $t$  is computed by,

$$S_l^{i,t} = \text{NormCosine}(\mu_l^i, \mu_l^{i \rightarrow t}), \quad (6)$$

where  $\text{NormCosine}(\cdot, \cdot)$  is the Normalized Cosine Similarity, which is calculated by scaling the Cosine Similarity score between  $-1$  and  $1$  using the minimum and the maximum Cosine Similarity scores from all the experts in all the MHSA blocks of the ViT. This normalization is necessary to increase the difference in magnitudes of the similarities between tasks, which results in better Expert sampling distributions during the sampling process in our experiments.

**The Auxiliary Task.** For a new task  $t$ , to handle the `New` and `Skip` experts at each layer  $l$  for which we do not have direct similarity scores. Instead, we introduce an auxiliary task, `Aux` (see the right of Fig. 4) which gives equally-likely chance to select `New` or `Skip` expert. For the `Aux` task itself, the similarity score between it and task  $t$  is defined by,

$$S_l^{aux,t} = -\max_{i=1}^{t-1} S_l^{i,t}, \quad (7)$$

which intuitively means we probabilistically resort to the `New` operation or the `Skip` operation when other experts turn out not “helpful” for the task  $t$ .

#### 2.4.3 TARGET NETWORK SELECTION AND FINETUNING

After the Supernet is trained, we adopt the same evolutionary search used in the SPOS method (Guo et al., 2020) based on the proposed HEE sampling strategy with a different epoch sampling probability  $\epsilon_2$  ( $\epsilon_2 > \epsilon_1$ , e.g.,  $\epsilon_2 = 0.5$ , to encourage more exploration during the evolutionary search). The evolutionary search is performed on the validation set to select the path which gives the best validation accuracy. After the target network for a new task is selected, we retrain the newly added layers by the `New` and `Adapt` operations from scratch (random initialization), rather than keeping the weights from the Supernet training. This is based on the observations in network pruning that it is the neural architecture topology that matters and that the warm-up weights may not need to be preserved to ensure good performance on the target dataset (Liu et al., 2019b).

Table 2: Results on the VDD benchmark (Rebuffi et al., 2017a) using ViT-B/8 (Dosovitskiy et al., 2021) over 3 different seeds. The last two columns show the performance under the task-incremental continual learning (TCL) setting.

Method	C100 36k/100	SVHN 42.4k/10	UCF 6.8k/101	OGIt 16.1k/1623	GTSR 28.2k/43	DPed 21.2k/2	Flwr 0.92k/102	Airc. 3.0k/100	DTD 1.7k/47	AA (Eqn. 8)	AF (Eqn. 4)	TCL AA	TCL AF
L2P++	79.16	0.62	5.40	14.05	62.82	4.09	2.97	2.73	2.98	19.42 ± 0.09	5.49 ± 0.37	68.68 ± 0.67	10.25 ± 0.77
DualPrompt	82.92	15.27	12.19	33.71	80.12	7.65	3.17	3.47	3.23	26.86 ± 0.40	3.54 ± 0.27	76.57 ± 0.23	2.57 ± 0.26
CODA-Prompt	88.06	18.09	12.16	53.95	95.65	19.42	6.83	6.34	5.62	34.01 ± 0.99	8.0 ± 0.72	75.83 ± 0.60	6.25 ± 0.60
S-Prompts	87.38	88.53	64.05	72.17	98.53	<b>99.65</b>	<b>96.63</b>	45.49	<b>58.07</b>	78.95 ± 0.07	0.32 ± 0.32	79.71 ± 0.08	0.0 ± 0.0
EWC	79.15	85.62	47.03	64.57	90.65	54.56	55.26	34.87	54.49	62.91 ± 0.81	19.82 ± 0.69	75.00 ± 0.67	6.33 ± 0.30
L2 Regularization	82.59	65.57	51.57	26.83	95.62	98.50	83.04	33.80	51.67	65.47 ± 0.33	7.23 ± 0.51	69.63 ± 0.17	3.00 ± 0.47
Experience Replay	56.68	6.23	33.45	74.33	7.00	0.02	64.25	27.03	37.55	34.06 ± 0.69	44.89 ± 0.72	54.55 ± 2.67	22.67 ± 3.01
<b>Our CHEEM</b>	<b>88.56</b>	<b>95.63</b>	<b>75.05</b>	<b>83.81</b>	<b>99.15</b>	99.64	90.92	<b>55.53</b>	56.42	<b>82.74 ± 0.54</b>	<b>0.33 ± 0.00</b>	<b>84.65 ± 0.33</b>	<b>0.0 ± 0.0</b>

### 3 EXPERIMENTS

**Data:** We use the VDD benchmark (Rebuffi et al., 2017a) and the Split ImageNet-R (attention) benchmark introduced in Wang et al. (2022c), derived from (Hendrycks et al., 2021). The VDD benchmark (Fig. 2) consists of 10 tasks with significantly varying distributions of classes per task. The vanilla ImageNet-R dataset consists of 200 classes and 30k images in total. The Split ImageNet-R benchmark is constructed by randomly splitting the ImageNet-R dataset into 10 tasks each of which has 20 classes. To test the effects of imbalance class distributions in streaming tasks, we construct another Split ImageNet-R benchmark consisting of 6 tasks with imbalance number of classes (5, 10, 15, 20, 50, 100) randomly sampled from the 200 total classes without replacement for each task. More details are provided in Appendix. E.

**Baselines of ExfCCL.** We compare with S-Prompts (Wang et al., 2022a), Learning to Prompt (L2P) Wang et al. (2022d), DualPrompt Wang et al. (2022c) and CODA-Prompt Smith et al. (2023a). On the VDD benchmark, we also test two regularization based methods, the L2 Parameter Regularization (Smith et al., 2023b) and the EWC (Kirkpatrick et al., 2017), and a strong baseline of Experience Replay (with iCARL (Rebuffi et al., 2017b) for buffer updates) using the total buffer size 20k and 10 exemplars per class. For EWC, L2 and iCarl, we freeze the backbone model and only keep the output projection layers (where CHEEM is placed) in the MHSA blocks trainable for better understanding the advantages of structurally updating them in our CHEEM. More details are in Appendix F and G.

**Metric:** In addition to average forgetting (Eqn. 4), we compare average accuracy in ExfCCL setting,

$$AA = \frac{1}{N} \sum_{t=1}^N \text{Acc}(T_t; f_{1:N}, C_t). \quad (8)$$

#### 3.1 RESULTS ON THE VDD BENCHMARK

Table 2 shows the results. We use the ImageNet-1k (the 1st task in the VDD) trained ViT-B/8 as the base model. For fair comparisons with prompting based baselines, we evaluate average accuracy and forgetting on the remaining 9 tasks in the VDD. **Our proposed CHEEM obtains significantly better performance than all the baselines.** We analyze the results as follows.

**i) The Advantage of Task-Aware Dynamic Backbone Over Frozen Backbone: CHEEM vs S-Prompts.** CHEEM adopts the same task ID inference method (i.e., the external memory) proposed by S-Prompts (Wang et al., 2022a). The improvement by CHEEM with above 3% absolute average accuracy increase clearly shows the advantage of our proposed internal parameter memory for maintaining task-aware dynamic backbones, against the method of prepending retrieved task-specific prompts in S-Prompts. Table 2 also shows a potential for harnessing prompt-based methods in CHEEM. S-Prompts performs better on tasks that are similar to the first ImageNet task and have less training data such as Flwr (918 training images) and DTD (1692 training images).

**ii) Explicit Task ID Inference: With vs Without.** The three baselines, L2P (Wang et al., 2022d), DualPrompt (Wang et al., 2022c) and CODA-Prompt (Smith et al., 2023a), use a shared head classifier in inference, which, as aforementioned in the introduction, suffers from the discrepancy between training and testing. Due to the significant imbalance class distributions in the VDD, their average accuracy undergo catastrophic drops (19.42% by L2P, 26.86% by DualPrompt and 34.01% by CODA-Prompt, vs 78.95% by S-Prompts and 82.74% by our CHEEM).

To understand the effects of head classifiers better, we also compare the performance when task IDs are provided in inference (i.e., Task-incremental CL, TCL), so the three methods can use task-specific

head segments of the shared classifier, and both S-Prompts and our CHEEM do not need the external task-centroid memory. The last two columns of Table 2 show the results. *We have three observations:*

- The performance of L2P, DualPrompt and CODA-Prompt are drastically boosted (e.g., from 34.01% to 75.83% by CODA-Prompt), showing the importance of explicitly inference task IDs for imbalance class distributions of streaming tasks in continual learning. Our CHEEM under CCL is still significantly better than the three methods under TCL (by more than 6%), **which reinforces the importance of the task-aware dynamic backbone learned via our CHEEM.**
- The performance difference between CCL and TCL by our CHEEM is small, around 2%, showing that the task ID inference method (Wang et al., 2022a) is effective in handling the imbalance class distributions in the VDD benchmark.
- In terms of average forgetting under CCL and TCL, we should compare the ratio between average forgetting and average accuracy, rather than the average forgetting rates in isolation, for which the first three methods undergo severer forgetting under CCL as intuitively expected. Our CHEEM does not have forgetting under TCL.

**iii) Stability vs Plasticity of the Backbone.** To further show the advantage of our CHEEM learning to update the feature backbone from the base model, we compare with EWC (Kirkpatrick et al., 2017), L2 Regularization (Smith et al., 2023b) and Experience Replay (Rebuffi et al., 2017b). The three methods regularize the weights, including the shared head classifier, in learning new tasks using different formulations. Comparing with L2P, DualPrompt and CODA-Prompt, under CCL, the regularization based methods work better, showing the negative effects of the discrepancy of handling the shared head in training and testing by the three prompting based methods, as well as a potential of integrating them. Comparing with S-Prompts, the plasticity of updating the entire backbone with regularization is less effective than the plasticity introduced by prompts while keeping the backbone frozen in S-Prompts. Comparing with our CHEEM, the three regularization based methods suffer from significant forgetting, highlighting the balanced stability and plasticity via our CHEEM.

Table 3: Results on two constructed Split ImageNet-R benchmarks, averaged across 3 different task orders. The base model is ImageNet-1k trained ViT-B/16 sourced from (Wightman, 2019).

Method	Imbalanced (6 tasks with (5, 10, 15, 20, 50, 100) classes per task)				Balanced (10 tasks with 20 classes per task)			
	AA (Eqn. 8)	AF (Eqn. 4)	TCL AA	TCL AF	AA (Eqn. 8)	AF (Eqn. 4)	TCL AA	TCL AF
L2P++	64.44 ± 1.38	8.62 ± 4.02	88.21 ± 0.28	0.50 ± 0.36	73.01 ± 0.57	5.80 ± 0.29	89.08 ± 0.38	0.48 ± 0.02
DualPrompt	67.20 ± 1.38	6.95 ± 3.96	89.46 ± 0.57	0.45 ± 0.24	73.01 ± 0.55	3.35 ± 0.36	90.07 ± 0.19	0.29 ± 0.19
CODAPrompt	<b>67.66</b> ± 3.09	7.06 ± 2.97	90.69 ± 0.44	0.46 ± 0.23	<b>76.48</b> ± 0.25	5.04 ± 0.12	91.86 ± 0.34	0.45 ± 0.10
Our CHEEM	66.97 ± 0.43	8.44 ± 3.55	<b>91.48</b> ± 0.63	0.0 ± 0.0	64.72 ± 0.28	8.73 ± 0.41	<b>92.47</b> ± 0.64	0.0 ± 0.0

### 3.2 RESULTS ON THE SPLIT IMAGENET-R BENCHMARK

The ImageNet-R dataset is created to challenge models trained using ImageNet. So, by design, the external memory of our CHEEM that is based on ImageNet-1k trained base model will not work well. Table 3 shows the results which we analyze as follows:

**i) Randomly Assigned and Imbalanced Classes in Streaming Task Challenge All of ExFCCL Methods We Test.** Under CCL, the three prompting based methods suffer from catastrophic drop of performance from balanced to imbalanced settings (e.g., from  $76.48 \pm 0.25\%$  to  $67.66 \pm 3.09\%$  by CODA-Prompt, see Fig. 5 and analyses in Appendix A ). Our CHEEM is on-par with DualPrompt and CODA-Prompt. Our CHEEM is also stable from balanced to imbalanced scenarios. Under TCL, all methods obtain significant performance boost, and our CHEEM is better than all others. Similarly, we envision that leveraging the dynamically learned backbone by our CHEEM in constructing the external memory will be a promising direction to be studied in future work.

**ii) Randomly Assigned Yet Balanced Classes in Streaming Tasks Challenge the Task-Centroid based Task ID Inference.** On the one hand, for streaming tasks with balanced but randomly sampled classes per task (without replacement), our CHEEM under CCL has the worst performance ( $64.72\%$  vs  $76.48\%$  by CODA-Prompt), which was caused by the low average precision of task ID inference (see Fig. 6 and analyses in Appendix B). Our CHEEM under TCL obtains the best performance ( $92.47\%$ ), which shows that the learned task-aware backbone is expressive. Overall, task centroids computed using the CLS tokens in the base model are not able to distinguish tasks from each other with high accuracy. One potential solution is to leverage the dynamically learned backbone by our CHEEM in constructing the external memory, considering its superior performance under TCL, at the expense of more costly task ID inference. On the other hand, *although it is interesting to test continual learning approaches using the Split ImageNet-R benchmark, the random composition of*



432 *classes in a task is not natural in comparison with scenarios in natural human learning. Unlike*  
 433 *the Split ImageNet-R, the VDD benchmark may suit continual learning better from perspective the*  
 434 *perspective of real-world scenarios, for which our CHEEM works the best.*  
 435

#### 436 4 VISUALIZATION AND ABLATION STUDIES

437 Due to space limit, we present visualizations of the learned CHEEM (examples like Fig. 2) in  
 438 Appendix C. We also conduct ablation studies including: **(i)** The proposed HEE sampling significantly  
 439 outperforms pure exploration (PE) sampling in Appendix D.1. **(ii)** The proposed HEE-empowered  
 440 SPOS NAS significantly outperforms the DARTS (Liu et al., 2019a) based NAS used in the learn-to-  
 441 grow method (Li et al., 2019) in Appendix D.2. **(iii)** The proposed internal memory learning method  
 442 outperforms three state-of-the-art methods: SupSup (Wortsman et al., 2020), EFT (Verma et al., 2021)  
 443 and LL (Ge et al., 2023) under the task-to-task transfer learning setting in Appendix D.3. And, more  
 444 are in Appendix D.3, D.4 and D.6.

#### 445 5 RELATED WORK

447 *Experience Replay Based approaches* aim to retain some exemplars, in the form of either raw data or  
 448 latent features, from the previous tasks and replay them to the model along with the data from a new  
 449 task (Aljundi et al., 2019b; Rebuffi et al., 2017b; Aljundi et al., 2019a; Balaji et al., 2020; Bang et al.,  
 450 2021; Chaudhry et al., 2021; Pham et al., 2021), or *generative replay methods* (Shin et al., 2017; Cong  
 451 et al., 2020) which replay exemplars sampled from the learned generators along with the data from  
 452 the current task. For exemplar-free continual learning, *Regularization-based Approaches* explicitly  
 453 prevent model parameters from deviating too far from their stable values learned on previous tasks  
 454 when learning a new task (Aljundi et al., 2018; 2019c; Douillard et al., 2020; Nguyen et al., 2018;  
 455 Kirkpatrick et al., 2017; Li & Hoiem, 2018; Zenke et al., 2017; Schwarz et al., 2018).

456 *Dynamic Models* aim to use different parameters for each task to eliminate the use of stored exemplars.  
 457 Dynamically Expandable Network (Yoon et al., 2018) adds neurons to a network based on learned  
 458 sparsity constraints and heuristic loss thresholds. PathNet (Fernando et al., 2017) and RPSNet  
 459 (Rajasegaran et al., 2019) learn task-specific submodules or paths. Progressive Neural Networks (Rusu  
 460 et al., 2016) learn a new network per task and adds lateral connections to the previous tasks’ networks.  
 461 The L2G (Li et al., 2019) uses Differentiable Architecture Search (DARTS) (Liu et al., 2019a) to  
 462 determine if a layer can be reused, adapted, or renewed for a task, which is tested for ConvNets  
 463 and the learning-to-grow operations are applied uniformly at each layer in a ConvNet. Our method  
 464 is motivated by the L2G method, but with significant differences. NAS has also been explored by  
 465 (Gao et al., 2023a) for continual learning on small scale tasks. Approaches that learn task-specific  
 466 components on top of a backbone include (Ge et al., 2023; Verma et al., 2021; Wortsman et al., 2020;  
 467 Abati et al., 2020).

468 Recently, there has been increasing interest in continual learning using Vision Transformers (Wang  
 469 et al., 2022d;c; Xue et al., 2022; Ermis et al., 2022; Douillard et al., 2022; Pelosin et al., 2022; Yu  
 470 et al., 2021; Li et al., 2022; Iscen et al., 2022; Wang et al., 2022a;b; Mohamed et al., 2023; Gao et al.,  
 471 2023b; Zhou et al., 2023). *Prompt Based approaches* learn external parameters appended to the data  
 472 tokens that encode task-specific information useful for classification (Wang et al., 2022d;a; Douillard  
 473 et al., 2022; Smith et al., 2023a; Jung et al., 2023; Tang et al., 2023; Zhou et al., 2024). Our proposed  
 474 method is complementary to prompting-based based approaches.

#### 475 6 CONCLUSION

476 This paper presents a method of transforming Vision Transformers (ViTs) for exemplar-free class-  
 477 incremental continual learning (ExfCCL), dubbed **CHEEM** (Continual Hierarchical-Exploration-  
 478 Exploitation Memory). Our CHEEM consists of the external (task-centroid) memory and the internal  
 479 (parameter) memory. The former is for task ID inference for test data based on clustered task  
 480 centroids in training. The latter is realized by a proposed Hierarchical-Exploration-Exploitation  
 481 (HEE) sampling based neural architecture search algorithm. The external and internal memory are  
 482 maintained in a decoupled way. Our CHEEM is tested on two challenging benchmarks, the VDD  
 483 benchmark and the continual ImageNet-R benchmark. It obtains state-of-the-art performance on  
 484 the VDD outperforming the prior art by a large margin, with sensible CHEEM structures continually  
 485 learned. It obtains on-par performance on the imbalanced ImageNet-R benchmark, while performs  
 worse on the balanced one, for which we analyze the reason thoroughly due to the external memory.

## REFERENCES

- 486  
487  
488 Davide Abati, Jakub Tomczak, Tijmen Blankevoort, Simone Calderara, Rita Cucchiara, and  
489 Babak Ehteshami Bejnordi. Conditional channel gated networks for task-aware continual  
490 learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR*  
491 *2020, Seattle, WA, USA, June 13-19, 2020*, pp. 3930–3939. Computer Vision Foundation / IEEE,  
492 2020. doi: 10.1109/CVPR42600.2020.00399. URL [https://openaccess.thecvf.com/  
493 content\\_CVPR\\_2020/html/Abati\\_Conditional\\_Channel\\_Gated\\_Networks\\_  
494 for\\_Task-Aware\\_Continual\\_Learning\\_CVPR\\_2020\\_paper.html](https://openaccess.thecvf.com/content_CVPR_2020/html/Abati_Conditional_Channel_Gated_Networks_for_Task-Aware_Continual_Learning_CVPR_2020_paper.html).
- 495 Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars.  
496 Memory aware synapses: Learning what (not) to forget. In Vittorio Ferrari, Martial Hebert, Cristian  
497 Sminchisescu, and Yair Weiss (eds.), *Computer Vision - ECCV 2018 - 15th European Conference,*  
498 *Munich, Germany, September 8-14, 2018, Proceedings, Part III*, volume 11207 of *Lecture Notes*  
499 *in Computer Science*, pp. 144–161. Springer, 2018. doi: 10.1007/978-3-030-01219-9\_9. URL  
500 [https://doi.org/10.1007/978-3-030-01219-9\\_9](https://doi.org/10.1007/978-3-030-01219-9_9).
- 501 Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin, and  
502 Lucas Page-Caccia. Online continual learning with maximal interfered retrieval. In Hanna M. Wal-  
503 lach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garn-  
504 nett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural*  
505 *Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC,*  
506 *Canada*, pp. 11849–11860, 2019a. URL [https://proceedings.neurips.cc/paper/  
507 2019/hash/15825aee15eb335cc13f9b559f166ee8-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/15825aee15eb335cc13f9b559f166ee8-Abstract.html).
- 508 Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample  
509 selection for online continual learning. In Hanna M. Wallach, Hugo Larochelle, Alina  
510 Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances*  
511 *in Neural Information Processing Systems 32: Annual Conference on Neural Information*  
512 *Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*,  
513 pp. 11816–11825, 2019b. URL [https://proceedings.neurips.cc/paper/2019/  
514 hash/e562cd9c0768d5464b64cf61da7fc6bb-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/e562cd9c0768d5464b64cf61da7fc6bb-Abstract.html).
- 515 Rahaf Aljundi, Marcus Rohrbach, and Tinne Tuytelaars. Selfless sequential learning. In *7th*  
516 *International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA,*  
517 *May 6-9, 2019*. OpenReview.net, 2019c. URL [https://openreview.net/forum?id=  
518 Bkxbrn0cYX](https://openreview.net/forum?id=Bkxbrn0cYX).
- 519 Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *CoRR*,  
520 abs/1607.06450, 2016. URL <http://arxiv.org/abs/1607.06450>.
- 521  
522 Yogesh Balaji, Mehrdad Farajtabar, Dong Yin, Alex Mott, and Ang Li. The effectiveness of  
523 memory replay in large scale continual learning. *CoRR*, abs/2010.02418, 2020. URL [https:  
524 //arxiv.org/abs/2010.02418](https://arxiv.org/abs/2010.02418).
- 525 Jihwan Bang, Heesu Kim, Youngjoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow  
526 memory: Continual learning with a memory of diverse samples. In *IEEE Conference on*  
527 *Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pp. 8218–  
528 8227. Computer Vision Foundation / IEEE, 2021. doi: 10.1109/CVPR46437.2021.00812.  
529 URL [https://openaccess.thecvf.com/content/CVPR2021/html/Bang\\_  
530 Rainbow\\_Memory\\_Continual\\_Learning\\_With\\_a\\_Memory\\_of\\_Diverse\\_  
531 Samples\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Bang_Rainbow_Memory_Continual_Learning_With_a_Memory_of_Diverse_Samples_CVPR_2021_paper.html).
- 532  
533 Yoshua Bengio, Nicholas Léonard, and Aaron C. Courville. Estimating or propagating gradients  
534 through stochastic neurons for conditional computation. *CoRR*, abs/1308.3432, 2013. URL  
535 <http://arxiv.org/abs/1308.3432>.
- 536 Hakan Bilen, Basura Fernando, Efstratios Gavves, Andrea Vedaldi, and Stephen Gould. Dynamic  
537 image networks for action recognition. In *2016 IEEE Conference on Computer Vision and Pattern*  
538 *Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pp. 3034–3042. IEEE Com-  
539 puter Society, 2016. doi: 10.1109/CVPR.2016.331. URL [https://doi.org/10.1109/  
CVPR.2016.331](https://doi.org/10.1109/CVPR.2016.331).

- 540 Arslan Chaudhry, Puneet Kumar Dokania, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss (eds.), Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XI, volume 11215 of Lecture Notes in Computer Science, pp. 556–572. Springer, 2018. doi: 10.1007/978-3-030-01252-6\_33. URL [https://doi.org/10.1007/978-3-030-01252-6\\_33](https://doi.org/10.1007/978-3-030-01252-6_33).
- 547 Arslan Chaudhry, Albert Gordo, Puneet K. Dokania, Philip H. S. Torr, and David Lopez-Paz. Using hindsight to anchor past knowledge in continual learning. In Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021, pp. 6993–7001. AAAI Press, 2021. doi: 10.1609/AAAI.V35I8.16861. URL <https://doi.org/10.1609/aaai.v35i8.16861>.
- 554 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2014, Columbus, OH, USA, June 23-28, 2014, pp. 3606–3613. IEEE Computer Society, 2014. doi: 10.1109/CVPR.2014.461. URL <https://doi.org/10.1109/CVPR.2014.461>.
- 559 Yulai Cong, Miaoyun Zhao, Jianqiao Li, Sijia Wang, and Lawrence Carin. GAN memory with no forgetting. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/bf201d5407a6509fa536afc4b380577e-Abstract.html>.
- 566 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- 572 Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (eds.), Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XX, volume 12365 of Lecture Notes in Computer Science, pp. 86–102. Springer, 2020. doi: 10.1007/978-3-030-58565-5\_6. URL [https://doi.org/10.1007/978-3-030-58565-5\\_6](https://doi.org/10.1007/978-3-030-58565-5_6).
- 578 Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 9275–9285. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00907. URL <https://doi.org/10.1109/CVPR52688.2022.00907>.
- 583 Beyza Ermis, Giovanni Zappella, Martin Wistuba, Aditya Rawal, and Cédric Archambeau. Continual learning with transformers for image classification. In IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2022, New Orleans, LA, USA, June 19-20, 2022, pp. 3773–3780. IEEE, 2022. doi: 10.1109/CVPRW56347.2022.00422. URL <https://doi.org/10.1109/CVPRW56347.2022.00422>.
- 588 Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. CoRR, abs/1701.08734, 2017. URL <http://arxiv.org/abs/1701.08734>.
- 592 Qiang Gao, Zhipeng Luo, Diego Klabjan, and Fengli Zhang. Efficient architecture search for continual learning. IEEE Trans. Neural Networks Learn. Syst., 34(11):8555–8565, 2023a. doi: 10.1109/TNNLS.2022.3151511. URL <https://doi.org/10.1109/TNNLS.2022.3151511>.

- 594 Qiankun Gao, Chen Zhao, Yifan Sun, Teng Xi, Gang Zhang, Bernard Ghanem, and Jian Zhang.  
595 A unified continual learning framework with general parameter-efficient tuning. In IEEE/CVF  
596 International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023, pp.  
597 11449–11459. IEEE, 2023b. doi: 10.1109/ICCV51070.2023.01055. URL [https://doi.org/](https://doi.org/10.1109/ICCV51070.2023.01055)  
598 [10.1109/ICCV51070.2023.01055](https://doi.org/10.1109/ICCV51070.2023.01055).
- 599 Yunhao Ge, Yuecheng Li, Di Wu, Ao Xu, Adam M. Jones, Amanda Sofie Rios, Iordanis Fostitropoulos,  
600 shixian wen, Po-Hsuan Huang, Zachary William Murdock, Gozde Sahin, Shuo Ni, Kiran Lekkala,  
601 Sumedh Anand Sontakke, and Laurent Itti. Lightweight learner for shared knowledge lifelong  
602 learning. Transactions on Machine Learning Research, 2023. ISSN 2835-8856. URL [https:](https://openreview.net/forum?id=Jj12c8kWUc)  
603 [//openreview.net/forum?id=Jj12c8kWUc](https://openreview.net/forum?id=Jj12c8kWUc).
- 604 Zichao Guo, Xiangyu Zhang, Haoyuan Mu, Wen Heng, Zechun Liu, Yichen Wei, and Jian Sun. Single  
605 path one-shot neural architecture search with uniform sampling. In Andrea Vedaldi, Horst Bischof,  
606 Thomas Brox, and Jan-Michael Frahm (eds.), Computer Vision - ECCV 2020 - 16th European  
607 Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XVI, volume 12361 of Lecture  
608 Notes in Computer Science, pp. 544–560. Springer, 2020. doi: 10.1007/978-3-030-58517-4\_32.  
609 URL [https://doi.org/10.1007/978-3-030-58517-4\\_32](https://doi.org/10.1007/978-3-030-58517-4_32).
- 610 Tyler L. Hayes, Nathan D. Cahill, and Christopher Kanan. Memory efficient experience replay  
611 for streaming learning. In International Conference on Robotics and Automation, ICRA 2019,  
612 Montreal, QC, Canada, May 20-24, 2019, pp. 9769–9776. IEEE, 2019. doi: 10.1109/ICRA.2019.  
613 8793982. URL <https://doi.org/10.1109/ICRA.2019.8793982>.
- 614 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
615 recognition. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR  
616 2016, Las Vegas, NV, USA, June 27-30, 2016, pp. 770–778. IEEE Computer Society, 2016. doi:  
617 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.
- 618 Dan Hendrycks and Kevin Gimpel. Bridging nonlinearities and stochastic regularizers with gaussian  
619 error linear units. CoRR, abs/1606.08415, 2016. URL [http://arxiv.org/abs/1606.](http://arxiv.org/abs/1606.08415)  
620 [08415](http://arxiv.org/abs/1606.08415).
- 621 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul  
622 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer.  
623 The many faces of robustness: A critical analysis of out-of-distribution generalization. In 2021  
624 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada,  
625 October 10-17, 2021, pp. 8320–8329. IEEE, 2021. doi: 10.1109/ICCV48922.2021.00823. URL  
626 <https://doi.org/10.1109/ICCV48922.2021.00823>.
- 627 Ahmet Iscen, Thomas Bird, Mathilde Caron, Alireza Fathi, and Cordelia Schmid. A memory  
628 transformer network for incremental learning. In 33rd British Machine Vision Conference 2022,  
629 BMVC 2022, London, UK, November 21-24, 2022, pp. 388. BMVA Press, 2022. URL [https:](https://bmvc2022.mpi-inf.mpg.de/388/)  
630 [//bmvc2022.mpi-inf.mpg.de/388/](https://bmvc2022.mpi-inf.mpg.de/388/).
- 631 Dahuin Jung, Dongyoon Han, Jihwan Bang, and Hwanjun Song. Generating instance-level prompts  
632 for rehearsal-free continual learning. In IEEE/CVF International Conference on Computer  
633 Vision, ICCV 2023, Paris, France, October 1-6, 2023, pp. 11813–11823. IEEE, 2023. doi:  
634 10.1109/ICCV51070.2023.01088. URL [https://doi.org/10.1109/ICCV51070.2023.](https://doi.org/10.1109/ICCV51070.2023.01088)  
635 [01088](https://doi.org/10.1109/ICCV51070.2023.01088).
- 636 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua  
637 Bengio and Yann LeCun (eds.), 3rd International Conference on Learning Representations, ICLR  
638 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015. URL [http:](http://arxiv.org/abs/1412.6980)  
639 [//arxiv.org/abs/1412.6980](http://arxiv.org/abs/1412.6980).
- 640 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A.  
641 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis,  
642 Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural  
643 networks. Proceedings of the National Academy of Sciences, 114(13):3521–3526, 2017. doi:  
644 10.1073/pnas.1611835114. URL [https://www.pnas.org/doi/abs/10.1073/pnas.](https://www.pnas.org/doi/abs/10.1073/pnas.1611835114)  
645 [1611835114](https://www.pnas.org/doi/abs/10.1073/pnas.1611835114).

- 648 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.  
649
- 650 Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. Human-level concept learning  
651 through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015. doi: 10.1126/  
652 science.aab3050. URL [https://www.science.org/doi/abs/10.1126/science.  
653 aab3050](https://www.science.org/doi/abs/10.1126/science.aab3050).
- 654 Duo Li, Guimei Cao, Yunlu Xu, Zhanzhan Cheng, and Yi Niu. Technical report for ICCV 2021  
655 challenge sslad-track3b: Transformers are better continual learners. *CoRR*, abs/2201.04924, 2022.  
656 URL <https://arxiv.org/abs/2201.04924>.  
657
- 658 Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual  
659 structure learning framework for overcoming catastrophic forgetting. In Kamalika Chaudhuri  
660 and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine  
661 Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings  
662 of Machine Learning Research*, pp. 3925–3934. PMLR, 2019. URL [http://proceedings.  
663 mlr.press/v97/li19m.html](http://proceedings.mlr.press/v97/li19m.html).
- 664 Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE Trans. Pattern Anal. Mach. Intell.*,  
665 40(12):2935–2947, 2018. doi: 10.1109/TPAMI.2017.2773081. URL [https://doi.org/10.  
666 1109/TPAMI.2017.2773081](https://doi.org/10.1109/TPAMI.2017.2773081).
- 667 Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. In  
668 *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA,  
669 May 6-9, 2019*. OpenReview.net, 2019a. URL [https://openreview.net/forum?id=  
670 S1eYHoC5FX](https://openreview.net/forum?id=S1eYHoC5FX).
- 671  
672 Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value of  
673 network pruning. In *7th International Conference on Learning Representations, ICLR 2019, New  
674 Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019b. URL [https://openreview.  
675 net/forum?id=rJlnB3C5Ym](https://openreview.net/forum?id=rJlnB3C5Ym).  
676
- 677 Subhansu Maji, Esa Rahtu, Juho Kannala, Matthew B. Blaschko, and Andrea Vedaldi. Fine-grained  
678 visual classification of aircraft. *CoRR*, abs/1306.5151, 2013. URL [http://arxiv.org/abs/  
679 1306.5151](http://arxiv.org/abs/1306.5151).
- 680 Michael McCloskey and Neal J. Cohen. Catastrophic interference in connectionist networks: The  
681 sequential learning problem. volume 24 of *Psychology of Learning and Motivation*, pp. 109–165.  
682 Academic Press, 1989. doi: [https://doi.org/10.1016/S0079-7421\(08\)60536-8](https://doi.org/10.1016/S0079-7421(08)60536-8). URL [https://  
683 www.sciencedirect.com/science/article/pii/S0079742108605368](https://www.sciencedirect.com/science/article/pii/S0079742108605368).
- 684  
685 Abdelrahman Mohamed, Rushali Grandhe, K. J. Joseph, Salman H. Khan, and Fahad Shahbaz Khan.  
686  $D^3$ former: Debaised dual distilled transformer for incremental learning. In *IEEE/CVF Conference  
687 on Computer Vision and Pattern Recognition, CVPR 2023 - Workshops, Vancouver, BC, Canada,  
688 June 17-24, 2023*, pp. 2421–2430. IEEE, 2023. doi: 10.1109/CVPRW59228.2023.00240. URL  
689 <https://doi.org/10.1109/CVPRW59228.2023.00240>.  
690
- 691 Stefan Munder and Darius M. Gavrilă. An experimental study on pedestrian classification. *IEEE  
692 Trans. Pattern Anal. Mach. Intell.*, 28(11):1863–1868, 2006. doi: 10.1109/TPAMI.2006.217. URL  
693 <https://doi.org/10.1109/TPAMI.2006.217>.
- 694 Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading  
695 digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning  
696 and Unsupervised Feature Learning 2011*, 2011. URL [http://ufldl.stanford.edu/  
697 housenumbers/nips2011\\_housenumbers.pdf](http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf).  
698
- 699 Cuong V. Nguyen, Yingzhen Li, Thang D. Bui, and Richard E. Turner. Variational continual  
700 learning. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver,  
701 BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net, 2018. URL  
<https://openreview.net/forum?id=BkQqq0gRb>.



- 702 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of  
703 classes. In Sixth Indian Conference on Computer Vision, Graphics & Image Processing, ICVGIP  
704 2008, Bhubaneswar, India, 16-19 December 2008, pp. 722–729. IEEE Computer Society, 2008.  
705 doi: 10.1109/ICVGIP.2008.47. URL <https://doi.org/10.1109/ICVGIP.2008.47>.
- 706  
707 Francesco Pelosin, Saurav Jha, Andrea Torsello, Bogdan Raducanu, and Joost van de Weijer. Towards  
708 exemplar-free continual learning in vision transformers: an account of attention, functional and  
709 weight regularization. In IEEE/CVF Conference on Computer Vision and Pattern Recognition  
710 Workshops, CVPR Workshops 2022, New Orleans, LA, USA, June 19-20, 2022, pp. 3819–3828.  
711 IEEE, 2022. doi: 10.1109/CVPRW56347.2022.00427. URL [https://doi.org/10.1109/](https://doi.org/10.1109/CVPRW56347.2022.00427)  
712 [CVPRW56347.2022.00427](https://doi.org/10.1109/CVPRW56347.2022.00427).
- 713 Quang Pham, Chenghao Liu, and Steven C. H. Hoi. Dualnet: Continual learning, fast and slow. In  
714 Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman  
715 Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on  
716 Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp.  
717 16131–16144, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/86a1fa88adb5c33bd7a68ac2f9f3f96b-Abstract.html)  
718 [86a1fa88adb5c33bd7a68ac2f9f3f96b-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/86a1fa88adb5c33bd7a68ac2f9f3f96b-Abstract.html).
- 719 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agar-  
720 wal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya  
721 Sutskever. Learning transferable visual models from natural language supervision. In Marina Meila  
722 and Tong Zhang (eds.), Proceedings of the 38th International Conference on Machine Learning,  
723 ICML 2021, 18-24 July 2021, Virtual Event, volume 139 of Proceedings of Machine Learning  
724 Research, pp. 8748–8763. PMLR, 2021. URL [http://proceedings.mlr.press/v139/](http://proceedings.mlr.press/v139/radford21a.html)  
725 [radford21a.html](http://proceedings.mlr.press/v139/radford21a.html).
- 726 Jathushan Rajasegaran, Munawar Hayat, Salman H. Khan, Fahad Shahbaz Khan, and Ling Shao.  
727 Random path selection for continual learning. In Hanna M. Wallach, Hugo Larochelle, Alina  
728 Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances  
729 in Neural Information Processing Systems 32: Annual Conference on Neural Information  
730 Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp.  
731 12648–12658, 2019. URL [https://proceedings.neurips.cc/paper/2019/hash/](https://proceedings.neurips.cc/paper/2019/hash/83da7c539e1ab4e759623c38d8737e9e-Abstract.html)  
732 [83da7c539e1ab4e759623c38d8737e9e-Abstract.html](https://proceedings.neurips.cc/paper/2019/hash/83da7c539e1ab4e759623c38d8737e9e-Abstract.html).
- 733 Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual do-  
734 mains with residual adapters. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio,  
735 Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances  
736 in Neural Information Processing Systems 30: Annual Conference on Neural  
737 Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp.  
738 506–516, 2017a. URL [https://proceedings.neurips.cc/paper/2017/hash/](https://proceedings.neurips.cc/paper/2017/hash/e7b24b112a44fdd9ee93bdf998c6ca0e-Abstract.html)  
739 [e7b24b112a44fdd9ee93bdf998c6ca0e-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/e7b24b112a44fdd9ee93bdf998c6ca0e-Abstract.html).
- 740  
741 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl:  
742 Incremental classifier and representation learning. In 2017 IEEE Conference on Computer Vision  
743 and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pp. 5533–5542.  
744 IEEE Computer Society, 2017b. doi: 10.1109/CVPR.2017.587. URL [https://doi.org/10.](https://doi.org/10.1109/CVPR.2017.587)  
745 [1109/CVPR.2017.587](https://doi.org/10.1109/CVPR.2017.587).
- 746 Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Su-  
747 sano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. In  
748 Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman  
749 Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on  
750 Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp.  
751 8583–8595, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/48237d9f2dea8c74c2a72126cf63d933-Abstract.html)  
752 [48237d9f2dea8c74c2a72126cf63d933-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/48237d9f2dea8c74c2a72126cf63d933-Abstract.html).
- 753 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,  
754 Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet  
755 large scale visual recognition challenge. Int. J. Comput. Vis., 115(3):211–252, 2015. doi: 10.  
[1007/S11263-015-0816-Y](https://doi.org/10.1007/s11263-015-0816-y). URL <https://doi.org/10.1007/s11263-015-0816-y>.

- 756 Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick,  
757 Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. CoRR,  
758 abs/1606.04671, 2016. URL <http://arxiv.org/abs/1606.04671>.  
759
- 760 Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye  
761 Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual  
762 learning. In Jennifer G. Dy and Andreas Krause (eds.), Proceedings of the 35th International  
763 Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July  
764 10-15, 2018, volume 80 of Proceedings of Machine Learning Research, pp. 4535–4544. PMLR,  
765 2018. URL <http://proceedings.mlr.press/v80/schwarz18a.html>.
- 766 Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning  
767 with deep generative replay. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio,  
768 Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.),  
769 Advances in Neural Information Processing Systems 30: Annual Conference on Neural  
770 Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp.  
771 2990–2999, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/  
772 0efbe98067c6c73dba1250d2beaa81f9-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/0efbe98067c6c73dba1250d2beaa81f9-Abstract.html).
- 773 James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf  
774 Arbelle, Rameswar Panda, Rogério Feris, and Zsolt Kira. Coda-prompt: Continual decomposed  
775 attention-based prompting for rehearsal-free continual learning. In IEEE/CVF Conference on  
776 Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24,  
777 2023, pp. 11909–11919. IEEE, 2023a. doi: 10.1109/CVPR52729.2023.01146. URL <https://doi.org/10.1109/CVPR52729.2023.01146>.
- 779 James Seale Smith, Junjiao Tian, Shaunak Halbe, Yen-Chang Hsu, and Zsolt Kira. A closer look  
780 at rehearsal-free continual learning. In IEEE/CVF Conference on Computer Vision and Pattern  
781 Recognition, CVPR 2023 - Workshops, Vancouver, BC, Canada, June 17-24, 2023, pp. 2410–  
782 2420. IEEE, 2023b. doi: 10.1109/CVPRW59228.2023.00239. URL [https://doi.org/10.  
783 1109/CVPRW59228.2023.00239](https://doi.org/10.1109/CVPRW59228.2023.00239).
- 784  
785 Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human  
786 actions classes from videos in the wild. CoRR, abs/1212.0402, 2012. URL [http://arxiv.  
787 org/abs/1212.0402](http://arxiv.org/abs/1212.0402).
- 788 Johannes Stalldkamp, Marc Schlipfing, Jan Salmen, and Christian Igel. Man vs. computer: Bench-  
789 marking machine learning algorithms for traffic sign recognition. Neural Networks, 32:323–  
790 332, 2012. doi: 10.1016/J.NEUNET.2012.02.016. URL [https://doi.org/10.1016/j.  
791 neunet.2012.02.016](https://doi.org/10.1016/j.neunet.2012.02.016).
- 792  
793 Yu-Ming Tang, Yi-Xing Peng, and Wei-Shi Zheng. When prompt-based incremental learning does  
794 not meet strong pretraining. In IEEE/CVF International Conference on Computer Vision, ICCV  
795 2023, Paris, France, October 1-6, 2023, pp. 1706–1716. IEEE, 2023. doi: 10.1109/ICCV51070.  
796 2023.00164. URL <https://doi.org/10.1109/ICCV51070.2023.00164>.
- 797  
798 Sebastian Thrun and Tom M. Mitchell. Lifelong robot learning. Robotics Auton. Syst., 15(1-2):  
799 25–46, 1995. doi: 10.1016/0921-8890(95)00004-Y. URL [https://doi.org/10.1016/  
800 0921-8890\(95\)00004-Y](https://doi.org/10.1016/0921-8890(95)00004-Y).
- 801  
802 Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of  
803 Machine Learning Research, 9(86):2579–2605, 2008. URL [http://jmlr.org/papers/  
v9/vandermaaten08a.html](http://jmlr.org/papers/v9/vandermaaten08a.html).
- 804  
805 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez,  
806 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von  
807 Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman  
808 Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on  
809 Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp.  
5998–6008, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/  
3f5ee243547dee91fbd053c1c4a845aa-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html).

- 810 Vinay Kumar Verma, Kevin J. Liang, Nikhil Mehta, Piyush Rai, and Lawrence Carin. Efficient  
811 feature transformations for discriminative and generative continual learning. In IEEE Conference  
812 on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, pp.  
813 13865–13875. Computer Vision Foundation / IEEE, 2021. doi: 10.1109/CVPR46437.2021.01365.  
814 URL [https://openaccess.thecvf.com/content/CVPR2021/html/Verma\\_](https://openaccess.thecvf.com/content/CVPR2021/html/Verma_Efficient_Feature_Transformations_for_Discriminative_and_Generative_Continual_Learning_CVPR_2021_paper.html)  
815 [Efficient\\_Feature\\_Transformations\\_for\\_Discriminative\\_and\\_](https://openaccess.thecvf.com/content/CVPR2021/html/Verma_Efficient_Feature_Transformations_for_Discriminative_and_Generative_Continual_Learning_CVPR_2021_paper.html)  
816 [Generative\\_Continual\\_Learning\\_CVPR\\_2021\\_paper.html](https://openaccess.thecvf.com/content/CVPR2021/html/Verma_Efficient_Feature_Transformations_for_Discriminative_and_Generative_Continual_Learning_CVPR_2021_paper.html).
- 817 Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained trans-  
818 formers: An occam’s razor for domain incremental learning. In Sanmi Koyejo, S. Mo-  
819 hamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural  
820 Information Processing Systems 35: Annual Conference on Neural Information Processing  
821 Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9,  
822 2022, 2022a. URL [http://papers.nips.cc/paper\\_files/paper/2022/hash/](http://papers.nips.cc/paper_files/paper/2022/hash/25886d7a7cf4e33fd44072a0cd81bf30-Abstract-Conference.html)  
823 [25886d7a7cf4e33fd44072a0cd81bf30-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/25886d7a7cf4e33fd44072a0cd81bf30-Abstract-Conference.html).
- 824 Zhen Wang, Liu Liu, Yiqun Duan, Yajing Kong, and Dacheng Tao. Continual learning with life-  
825 long vision transformer. In IEEE/CVF Conference on Computer Vision and Pattern Recognition,  
826 CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 171–181. IEEE, 2022b. doi: 10.  
827 1109/CVPR52688.2022.00027. URL [https://doi.org/10.1109/CVPR52688.2022.](https://doi.org/10.1109/CVPR52688.2022.00027)  
828 [00027](https://doi.org/10.1109/CVPR52688.2022.00027).
- 829 Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren,  
830 Guolong Su, Vincent Perot, Jennifer G. Dy, and Tomas Pfister. Dualprompt: Complementary  
831 prompting for rehearsal-free continual learning. In Shai Avidan, Gabriel J. Brostow, Moustapha  
832 Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), Computer Vision - ECCV 2022 - 17th  
833 European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXVI, volume  
834 13686 of Lecture Notes in Computer Science, pp. 631–648. Springer, 2022c. doi: 10.1007/  
835 978-3-031-19809-0\_36. URL [https://doi.org/10.1007/978-3-031-19809-0\\_](https://doi.org/10.1007/978-3-031-19809-0_36)  
836 [36](https://doi.org/10.1007/978-3-031-19809-0_36).
- 837 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent  
838 Perot, Jennifer G. Dy, and Tomas Pfister. Learning to prompt for continual learning. In IEEE/CVF  
839 Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA,  
840 June 18-24, 2022, pp. 139–149. IEEE, 2022d. doi: 10.1109/CVPR52688.2022.00024. URL  
841 <https://doi.org/10.1109/CVPR52688.2022.00024>.
- 842 Ross Wightman. Pytorch image models. [https://github.com/rwightman/](https://github.com/rwightman/pytorch-image-models)  
843 [pytorch-image-models](https://github.com/rwightman/pytorch-image-models), 2019.
- 844 Mitchell Wortsman, Vivek Ramanujan, Rosanne Liu, Aniruddha Kembhavi, Mohammad  
845 Rastegari, Jason Yosinski, and Ali Farhadi. Supermasks in superposition. In Hugo  
846 Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien  
847 Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference  
848 on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020,  
849 virtual, 2020. URL [https://proceedings.neurips.cc/paper/2020/hash/](https://proceedings.neurips.cc/paper/2020/hash/ad1f8bb9b51f023cdc80cf94bb615aa9-Abstract.html)  
850 [ad1f8bb9b51f023cdc80cf94bb615aa9-Abstract.html](https://proceedings.neurips.cc/paper/2020/hash/ad1f8bb9b51f023cdc80cf94bb615aa9-Abstract.html).
- 851 Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and  
852 Yun Fu. Large scale incremental learning. In IEEE Conference on Computer Vision  
853 and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pp.  
854 374–382. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00046.  
855 URL [http://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Wu\\_Large\\_](http://openaccess.thecvf.com/content_CVPR_2019/html/Wu_Large_Scale_Incremental_Learning_CVPR_2019_paper.html)  
856 [Scale\\_Incremental\\_Learning\\_CVPR\\_2019\\_paper.html](http://openaccess.thecvf.com/content_CVPR_2019/html/Wu_Large_Scale_Incremental_Learning_CVPR_2019_paper.html).
- 857 Mengqi Xue, Haofei Zhang, Jie Song, and Mingli Song. Meta-attention for vit-backed continual  
858 learning. In IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022,  
859 New Orleans, LA, USA, June 18-24, 2022, pp. 150–159. IEEE, 2022. doi: 10.1109/CVPR52688.  
860 [2022.00025](https://doi.org/10.1109/CVPR52688.2022.00025). URL <https://doi.org/10.1109/CVPR52688.2022.00025>.
- 861 Peng Ye, Baopu Li, Yikang Li, Tao Chen, Jiayuan Fan, and Wanli Ouyang.  $\beta$ -darts: Beta-decay  
862 regularization for differentiable architecture search. In IEEE/CVF Conference on Computer  
863 and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp. 150–159. IEEE, 2022. doi: 10.1109/CVPR52688.2022.00025.

864 Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, pp.  
865 10864–10873. IEEE, 2022. doi: 10.1109/CVPR52688.2022.01060. URL [https://doi.org/](https://doi.org/10.1109/CVPR52688.2022.01060)  
866 10.1109/CVPR52688.2022.01060.  
867

868 Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically  
869 expandable networks. In 6th International Conference on Learning Representations, ICLR 2018,  
870 Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net,  
871 2018. URL <https://openreview.net/forum?id=Sk7KsfW0->.

872 Pei Yu, Yinpeng Chen, Ying Jin, and Zicheng Liu. Improving vision transformers for incremental  
873 learning. CoRR, abs/2112.06103, 2021. URL <https://arxiv.org/abs/2112.06103>.  
874

875 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence.  
876 In Doina Precup and Yee Whye Teh (eds.), Proceedings of the 34th International Conference  
877 on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017, volume 70 of  
878 Proceedings of Machine Learning Research, pp. 3987–3995. PMLR, 2017. URL [http://](http://proceedings.mlr.press/v70/zenke17a.html)  
879 [proceedings.mlr.press/v70/zenke17a.html](http://proceedings.mlr.press/v70/zenke17a.html).

880 Da-Wei Zhou, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-incremental learning  
881 with pre-trained models: Generalizability and adaptivity are all you need. CoRR, abs/2303.07338,  
882 2023. doi: 10.48550/ARXIV.2303.07338. URL [https://doi.org/10.48550/arXiv.](https://doi.org/10.48550/arXiv.2303.07338)  
883 [2303.07338](https://doi.org/10.48550/arXiv.2303.07338).

884 Da-Wei Zhou, Hai-Long Sun, Han-Jia Ye, and De-Chuan Zhan. Expandable subspace ensemble  
885 for pre-trained model-based class-incremental learning. In IEEE/CVF Conference on Computer  
886 Vision and Pattern Recognition, CVPR 2024, Seattle, WA, USA, June 16-22, 2024, pp. 23554–  
887 23564. IEEE, 2024. doi: 10.1109/CVPR52733.2024.02223. URL [https://doi.org/10.](https://doi.org/10.1109/CVPR52733.2024.02223)  
888 [1109/CVPR52733.2024.02223](https://doi.org/10.1109/CVPR52733.2024.02223).  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

## 918 APPENDIX

919

## 920 OUTLINE

921

922

923 In this Appendix, we elaborate on the following aspects that are not presented in the submission due  
924 to the space limit:

925

- 926 • **Section A – Analyses of head classifier design choices** in continual learning.
- 927 • **Section B – Analysis of the discriminative capacity of centroids for task ID inference**
- 928 • **Section C – More examples of learned CHEEM** by our proposed hierarchical exploration-  
929 exploitation (HEE) sampling scheme and by the vanilla pure exploration (PE) sampling scheme.
- 930 • **Section D - Ablation Studies:**
  - 931 – **Section D.1 – Comparing Pure Exploration vs. our proposed Hierarchical Exploration-**  
932 **Exploitation sampling scheme for SPOS NAS**
  - 933 – **Section D.2 – Comparing proposed HEE empowered SPOS NAS with DARTS:** We  
934 compare with DARTS, which has been used by Learn-to-Grow and show that our proposed  
935 method outperforms L2G with DARTS.
  - 936 – **Section D.3 – Can CHEEM be used for task-to-task transfer?**
  - 937 – **Section D.4 – Can we use other components in ViTs as CHEEM?:** We compare with  
938 the Query/Key/Value layer and the FFN layer, and verify the effectiveness of the proposed  
939 identification in the main paper.
  - 940 – **Section D.5 – How is the `Adapt` operation implemented?:** We elaborate two implementa-  
941 tion methods of the `Adapt` operation and compare their performance.
  - 942 – **Section D.6 - Effects of task orders.**
- 943 • **Section E – Details of the two benchmarks** tested in the experiments: the Visual Domain  
944 Decathlon (VDD) (Rebuffi et al., 2017a) benchmark and the ImageNet-R (Hendrycks et al., 2021)  
945 benchmark.
- 946 • **Section F – The base model and its training details:** the Vision Transformer (ViT) model  
947 specification (ViT-B/8) and how it is initially training on the ImageNet in our experiments.
- 948 • **Section G – Experimental settings and training hyperparameters in our implementation** on  
949 the VDD benchmark and the 5-dataset benchmark.
- 950 • **Section H - Details of Modifying SupSup, EFT and LL on the VDD Benchmark to work with**  
951 **Vision Transformers:**

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971



## A TASK-SPECIFIC HEAD SELECTION VIA TASK ID INFERENCE VS SHARING A GROWING, COMMON TASK HEAD

As aforementioned in the introduction of main text, explicitly inferring task IDs from the external memory for test data is to select the task head classifiers as done in (Wang et al., 2022a), while also enabling selecting task-specific feature backbones (exploited in our CHEEM). An alternative approach is to share a growing, common task head as done in (Wang et al., 2022d;c; Smith et al., 2023a). The former can suffer from low average precision, while for the latter we observe serious instability caused by imbalanced classes in stream tasks due to the discrepancy in training and testing (see Table 2 and Table 3).

To show the instability, let  $\mathcal{C} = (z_1, z_2, \dots, z_K)$  be the logits output of the common head. During training, the task ID for each  $z_k \in \mathcal{C}$  is known, and the logits of previous tasks are masked out (using  $-\infty$ ) before computing the softmax (i.e., task-specific heads are used in training). During testing, the label for a testing sample is predicated by  $\arg \max_{k \in [1, K]} z_k$ , across the entire head. This discrepancy between training and testing imposes a very strong assumption of the stability of logits in training and testing, which entails that the task-specific head segment for a testing example will underlyingly resemble the 1-hot vector, such that it can hold across the entire head. Although (Wang et al., 2022d;c; Smith et al., 2023a) show strong performance using the same number of classes per task (20) in ImageNet-R (Hendrycks et al., 2021), we observe that slightly varying the number of classes per task, while keeping everything else the same, can lead to around 10% drop of performance in ImageNet-R. For more challenging task-class distributions as in the VDD (Fig. 2), (Wang et al., 2022d;c; Smith et al., 2023a) have a catastrophic drop of performance.

To analyze this further, we probe the model after the training on the final task is complete. For each task, we first make predictions using the task specific part of the head (i.e., the task-incremental setting) and calculate the entropy of predicted probabilities of all the samples with the task-specific head. Then, we make predictions with the full head (i.e., the class-incremental setting), including all the classes in the previous tasks, and track samples for which the logit with the maximum values changes to one of the classes from the previous tasks, which we refer to as **switched samples**.

As shown in Fig. 5, we observe that the entropy for the switched samples is consistently higher than those of the unswitched samples. Furthermore, the variance of the the entropy consistently goes down as the tasks progress, which indicates that for later tasks, only the predictions that are highly confident remain unswitched. However, the entropy for the switched samples shows a large variance. This indicates that sharp predictions are not a sufficient condition to avoid switching.

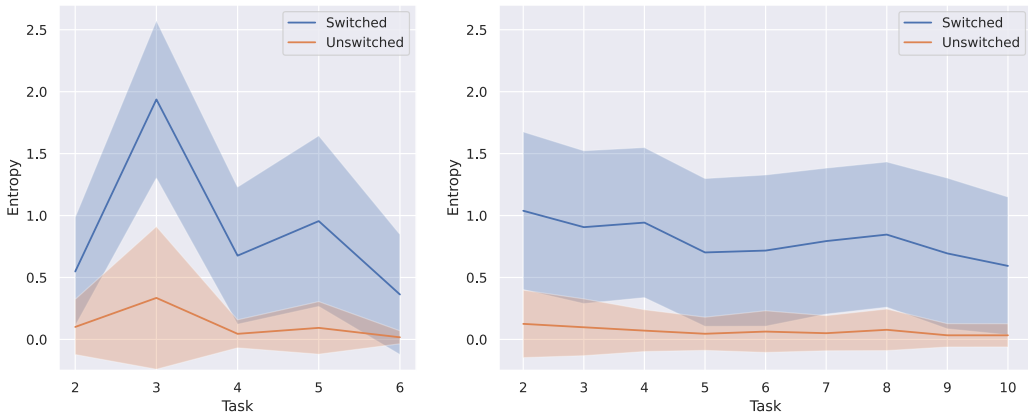


Figure 5: Mean entropy for the switched and unswitched samples from each task on the Split ImageNet-R benchmarks using CODA-Prompt. *Left*: Imbalanced ImageNet-R with varying number of classes (5,10,15,20, 50, 100) per task and 6 tasks in total. *Right*: Balanced ImageNet-R with same number (20) of classes per task and 10 tasks in total.

## B THE PROS AND CONS OF BASE MODEL COMPUTED TASK CENTROIDS AS THE EXTERNAL MEMORY FOR TASK ID INFERENCE

To understand effects of the external task-centroid memory we adopted from S-Prompts (Wang et al., 2022a), we analyze the separability of the K-Mean clusters of different tasks on the VDD (Rebuffi et al., 2017a) (see results in Table 2) and Balanced ImageNet-R benchmarks (see results in Table 3) by visualizing them using t-SNE (van der Maaten & Hinton, 2008) in Fig. 6.

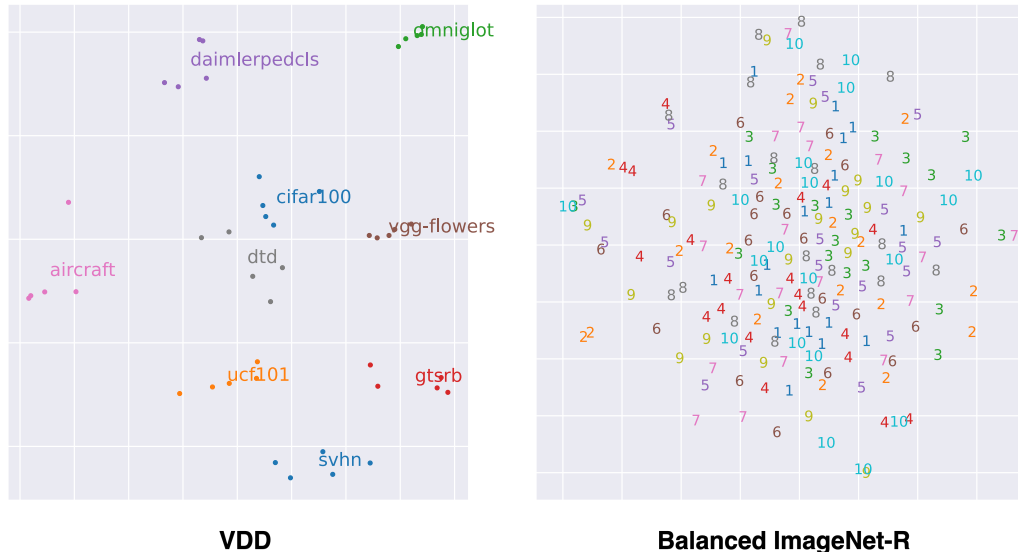


Figure 6: The t-SNE visualization of task centroids using the frozen base model CLS tokens of training examples of each task. We run K-Mean on the VDD benchmark with  $K = 5$  per task. We use the mean CLS token per class (20 centroids per task) on the balanced Imagenet-R benchmark, which we found works better.

For our proposed CHEEM, we can observe that the external task-centroid memory is effective on the VDD, providing sufficiently high recall rates for the internal memory of CHEEM and resulting in overall high average accuracy and negligible average forgetting. On the contrary, the external task-centroid memory of the balanced Imagenet-R benchmark explains why our CHEEM performs much worse under CCL. We note that the external memory is constructed based on the frozen base model. As discussed in the main text, a potential improvement is to leverage dynamically updated backbones via the internal memory of our CHEEM in computing the task centroids at the expense of increased task ID inference cost. That is also to shift from our current decoupled external and internal memory design to integrative / coupled external and internal memory.

Together with the experimental results that the prompting based methods (L2P, DualPrompt and CODA-Prompt) can handle balanced ImageNet-R better in terms of preventing switched samples as analyzed in Appendix A, and that regularized based method (EWC and L2 Regularization in Table 2) can handle VDD better than prompting based methods, we envision there are opportunities of exploring the integration between them and the internal memory of our CHEEM for more robust ExfCCL across scenarios including both VDD and balanced/imbalanced ImageNet-R.

## C EXAMPLES OF CHEEM LEARNED SEQUENTIALLY AND CONTINUALLY

Fig. 7 shows the CHEEM learned sequentially and continually via the proposed HEE-based NAS on the VDD benchmark (Rebuffi et al., 2017a) with three different random seeds.

As comparisons, Fig. 8 shows the structure updates learned using the vanilla PE-based NAS. It can be seen that pure exploration does not reuse components from similar tasks. The pure exploration based method adds many unnecessary `Adapt` and `New` operations even though the tasks are similar (e.g., ImNet  $\rightarrow$  C100), verifying the effectiveness of the proposed sampling method. While the pure

1080 exploration scheme adds many `Skip` operations, thereby reducing the overall FLOPs, the average  
1081 accuracy is low by a large margin, about 6%. This shows that the pure exploration scheme cannot  
1082 learn to choose operations in a task synergy aware way.  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

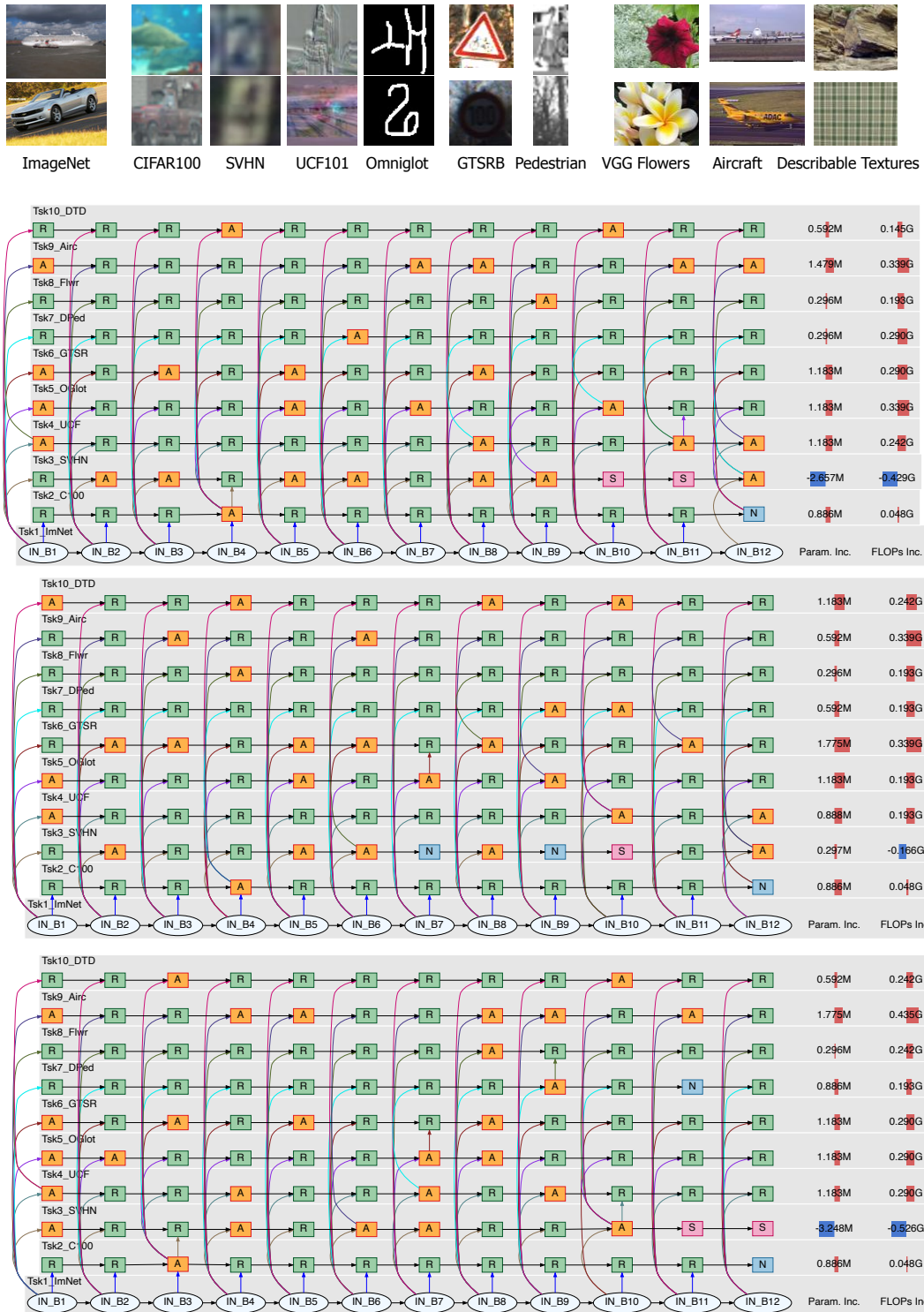
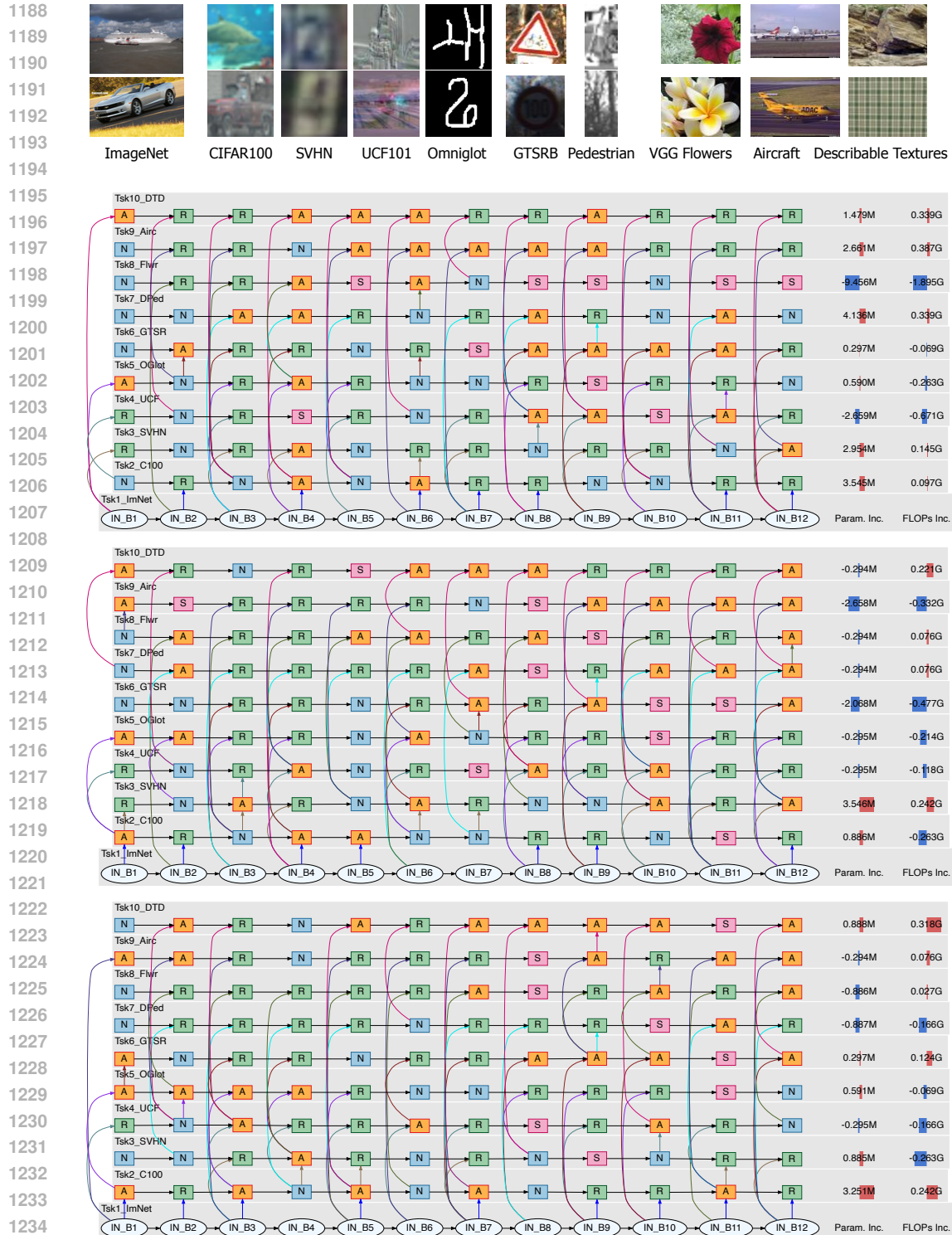


Figure 7: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi et al., 2017a) with the task sequence shown in the top **using our proposed HEE-based NAS** and three different random seeds. The overall performance is reported in Table 2 in the main paper. **S**, **R**, **A** and **N** represent Skip, Reuse, Adapt and New respectively. The last two columns show the number of new task-specific parameters and added FLOPs respectively, in comparison with the first task, ImNet model. Overall, the learned task synergies make intuitive sense and remain relatively stable across different random seeds.



1235 Figure 8: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi  
1236 et al., 2017a) with the task sequence shown in the top using the vanilla PE-based NAS and three  
1237 different random seeds. S, R, A and N represent Skip, Reuse, Adapt and New respectively.  
1238 The overall performance is reported in Table 4 in this appendix and our proposed HEE-based NAS  
1239 significantly improves the performance by an absolute 6% average accuracy. In terms of the learned  
1240 CHEEM, the PE-based NAS leads to much more New operations, which shows it is less effective in  
1241 terms of leveraging task synergies.



## D ABLATION STUDIES

### D.1 HIERARCHICAL EXPLORATION AND EXPLOITATION (HEE) VS. PURE EXPLORATION (PE)

We verify the effectiveness of our HEE sampling (Table 4), which is significantly better by a large margin, 6% absolute average accuracy increase.

Table 4: The effectiveness of our proposed hierarchical exploration and exploitation sampling empowered SPOS NAS (Fig. 4). The structure updates learned by the PE strategy are visualized in Fig. 8 in the Appendix.

Method	C100	SVHN	UCF	OGIt	GTSR	DPed	Flwr	Airc.	DTD	Avg. Accuracy	Avg. Forgetting
PE	80.34	94.75	72.63	81.43	99.13	99.59	72.45	38.27	37.22	75.09 $\pm$ 0.81	0.32 $\pm$ 0.00
HEE	88.56	95.63	75.05	83.81	99.15	99.64	90.92	55.53	56.42	82.74 $\pm$ 0.54	0.33 $\pm$ 0.0

### D.2 COMPARISON WITH DARTS

Table 5 shows the effectiveness of the proposed HEE empowered SPOS NAS against DARTS, as used in the original learn-to-grow method. Our proposed HEE empowered SPOS NAS outperforms DARTS, and the more advanced  $\beta$ -DARTS (Ye et al., 2022) by a large margin. We compare with Learn-to-Grow under a task-incremental continual learning (TCL) setting as used in the original formulation of L2G by (Li et al., 2019).

Table 5: The effectiveness of our proposed HEE sampling empowered SPOS NAS (Fig. 4). The two L2G variants use the same ViT-B base model as our HEE.

Method	C100	SVHN	UCF	OGIt	GTSR	DPed	Flwr	Airc.	DTD	Avg. Accuracy
L2G (DARTS)	88.47	85.20	79.22	80.19	99.28	98.06	76.14	39.29	46.01	77.45 $\pm$ 2.41
L2G ( $\beta$ -DARTS)	88.95	94.73	75.31	79.76	99.84	99.76	78.86	34.50	47.09	78.14 $\pm$ 0.54
HEE	90.93	95.96	80.74	83.25	99.94	99.96	94.12	58.90	60.05	84.65 $\pm$ 0.33

### D.3 RESULTS BY TASK-TO-TASK TRANSFER BASED CONTINUAL LEARNING

To evaluate the transfer learning capabilities of CHEEM, we compare with three state-of-the-art methods: Supermasks in Superposition (SupSup) (Wortsman et al., 2020), Efficient Feature Transformation (EFT) (Verma et al., 2021) and Lightweight Learner (LL) (Ge et al., 2023) that modify the backbone model separately for each task, which we refer to as a task-to-task transfer setting. We modify all the methods to work with ViTs. We provide the details of the modifications in Section H. All three methods are fine-tuned on a new task for 50 epochs. For CHEEM, we use 50 epochs for the supernet training and 30 epochs for fine-tuning.

Table 6 demonstrates that CHEEM achieves the highest average accuracy on the VDD benchmark. Although its performance on individual tasks may be lower, CHEEM exhibits robustness to domain variations, as reflected in its overall average accuracy. In contrast, other methods show varying performance depending on the domain. SupSup performs well on tasks that are significantly different from the base ImageNet task, such as UCF 101, Omniglot, and SVHN, but performs poorly on tasks closely related to ImageNet. Conversely, EFT and LL excel on tasks similar to ImageNet but struggle with out-of-distribution tasks. Despite its lower individual task performance, CHEEM’s consistent performance across diverse downstream tasks highlights its adaptability and makes it more general.

Table 6: Results on the VDD benchmark (Rebuffi et al., 2017a) using ViT-B/8 (Dosovitskiy et al., 2021) under the task-to-task transfer based continual learning protocol. The learned CHEEM are visualized in Fig. 10 in the Appendix.

Method	ImNet	C100	SVHN	UCF	OGIt	GTSR	DPed	Flwr	Airc.	DTD	Avg. Accuracy
SupSup	82.65	89.96	<b>96.05</b>	<b>81.68</b>	<b>84.60</b>	<b>99.97</b>	99.97	78.76	44.18	51.60	81.14 $\pm$ 0.04
EFT	82.65	91.86	93.51	73.89	75.62	99.58	<b>99.98</b>	96.34	48.17	<b>64.40</b>	82.60 $\pm$ 0.07
LL	82.65	<b>91.92</b>	93.90	75.63	77.07	99.71	99.96	<b>96.47</b>	49.33	64.34	83.10 $\pm$ 0.02
Our CHEEM	82.65	90.93	95.96	80.74	83.25	99.94	99.96	94.12	<b>58.90</b>	60.05	<b>84.65 <math>\pm</math> 0.33</b>

#### D.4 CHEEM PLACED AT OTHER ViT COMPONENTS

Table 7 shows the performance comparisons with other four different components in the ViT (the Query/Key/Value linear projection layer and the FFN block) used in realizing the proposed CHEEM. The Query/Key/Value component as the CHEEM does not perform as well as the Projection component. The FFN block as the CHEEM performs only slightly better than the Projection layer, but at the expense of a much larger parameter cost. This reinforces our identification above.

Table 7: Results of ablation study on other components of the ViT used for realizing the CHEEM. The results have been averaged over 3 different seeds.

Component	ImNet	C100	SVHN	UCF	OGIt	GTSR	DPed	Flwr	Airc.	DTD	Avg. Accuracy	Avg. Param. Inc./task (M)
Projection	82.65	90.54	96.12	75.53	83.81	99.93	99.88	91.21	55.59	59.18	83.44 ± 0.50	1.06 ± 0.04
Query	82.65	89.66	93.74	71.53	82.02	99.87	99.89	90.03	49.57	59.40	81.84 ± 0.32	2.38 ± 0.12
Key	82.65	89.29	94.77	72.25	81.86	99.86	99.90	88.86	51.72	60.46	82.16 ± 0.17	2.41 ± 0.03
Value	82.65	84.94	95.90	75.85	84.68	99.89	99.89	86.54	48.83	55.37	81.46 ± 0.25	1.70 ± 0.11
FFN	82.65	91.05	96.08	76.96	85.22	99.94	99.94	93.79	56.74	59.61	84.20 ± 0.28	2.31 ± 0.28

#### D.5 IMPLEMENTATION DETAILS OF THE Adapt OPERATION

**How to Adapt in a sustainable way?** The proposed Adapt operation will effectively increase the depth of the network in a plain way. In the worst case, if too many tasks use Adapt on top of each other, we will end up stacking too many MLP layers together. This may lead to unstable training due to gradient vanishing. Shortcut connections (He et al., 2016) have been shown to alleviate the gradient vanishing and exploding problems, making it possible to train deeper networks. We introduce the shortcut connection in adding a MLP Adapt operation. We test two different implementations: with shortcut in all the three components (supernet training, target network selection and target network finetuning) versus with shortcut only in target network finetuning (i.e., without shortcut in the NAS including both supernet training and target network selection).

Table 8: Results of the ablation study on the implementation of the Adapt operation: with (w/) vs without (w/o) shortcut connection for the MLP Adapt layer in NAS. The first two rows are for the sequential and continual paradigm and the last two rows for the task-to-task (T2T) transfer based paradigm.

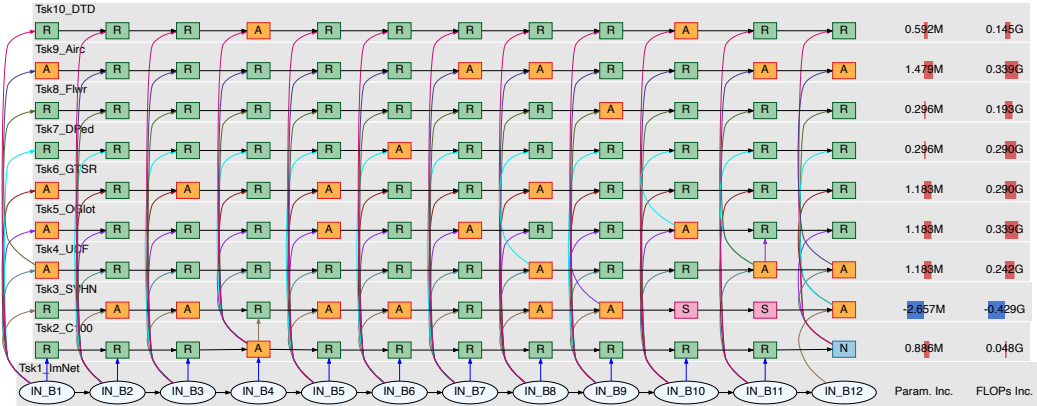
Shortcut in NAS	ImNet	C100	SVHN	UCF	OGIt	GTSR	DPed	Flwr	Airc.	DTD	Avg. Accuracy	Avg. Param. ↑/task (M)	Avg. FLOPs ↑/task (G)
w/o	82.65	90.54	96.12	75.53	83.81	99.93	99.88	91.21	55.59	59.18	83.44 ± 0.50	1.06 ± 0.04	0.17 ± 0.01
w/	82.65	91.18	96.18	82.34	86.03	99.91	99.95	91.60	58.90	58.56	84.73 ± 0.19	2.01 ± 0.18	0.49 ± 0.13
w/o (T2T)	82.65	90.93	95.96	80.74	83.25	99.94	99.96	94.12	58.90	60.05	84.65 ± 0.33	2.61 ± 0.15	-0.19 ± 0.09
w/ (T2T)	82.65	91.24	99.25	84.14	85.99	99.97	99.95	94.64	60.34	61.63	85.68 ± 0.16	3.23 ± 0.12	0.01 ± 0.02

Table 8 shows the performance comparisons on the VDD Benchmark under the continual learning paradigm. In terms of sequentially introduced complexities, a more compact model is learned without the shortcut in the Adapt during NAS (supernet training and target network selection) as evidenced by the number of additional parameters and the increase in FLOPs. Using the shortcut in both supernet training and target network selection results in twice the parameter increase, and almost 4× increase in FLOPs. Fig. 9 and Fig. 10 show comparisons of the learned CHEEM by the two implementation methods under the two paradigms respectively.

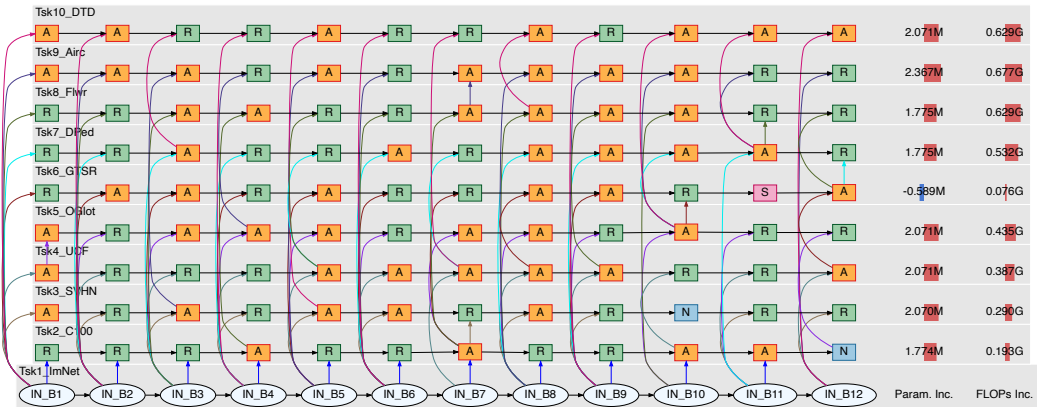
**Remarks.** We have two remarks as follows.

- We use the more parameter-efficient implementation (i.e., w/o shortcut for the Adapt in NAS) in the main paper for both the continual learning and the task-to-task transfer learning paradigms, even though the counterparts have better performance.
- We note that although the T2T paradigm results in larger parameter increase per task, its computational costs are relatively lower due to either more Skip operations learned and/or the fact that there is no Adapt-on-Adapt operations since it is task-to-task transfer based learning.

1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403



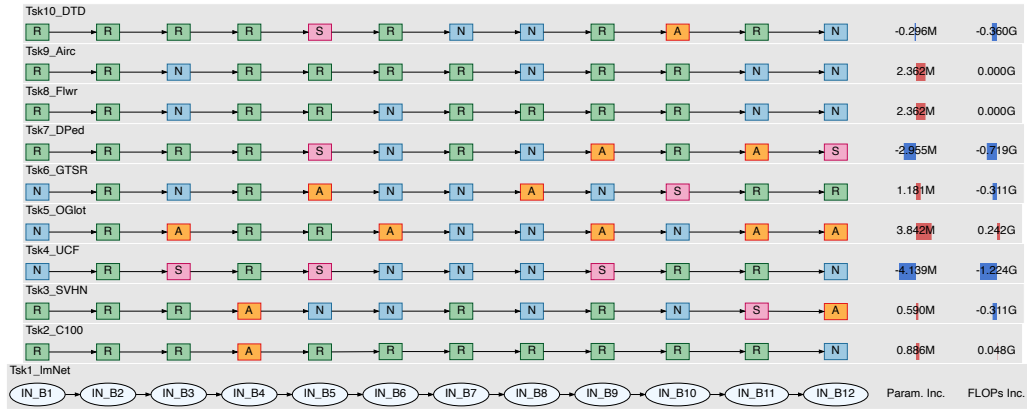
(a) An example of CHEEM learned **without** the shortcut for the MLP Adapt layer in NAS (the same one as the 2nd row in Fig. 7).



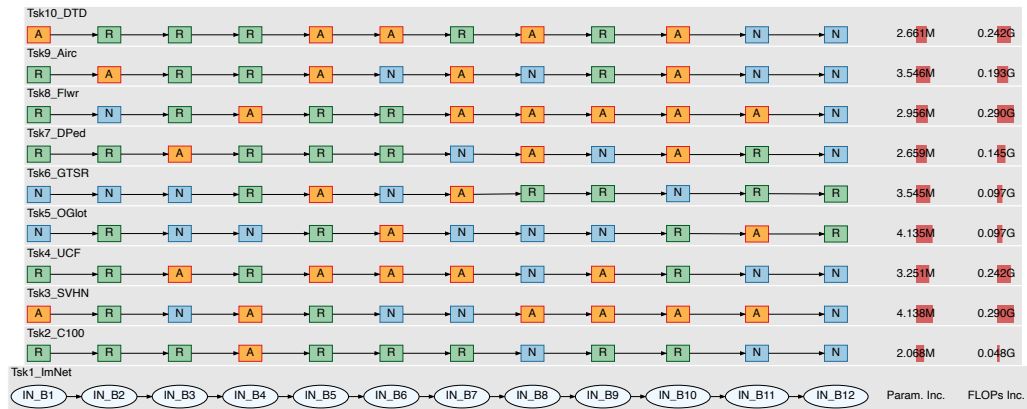
(b) An example of CHEEM learned **with** the shortcut for the MLP Adapt layer in NAS. More Adapt on top of Adapt operations are learned.

Figure 9: Comparisons between CHEEM learned by two different implementations of the MLP Adapt operation under the sequential and continual learning paradigm. S, R, A and N represent Skip, Reuse, Adapt and New respectively.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457



(a) An example of CHEEM learned **without** the shortcut for the MLP Adapt layer in NAS.



(b) An example of CHEEM learned **with** the shortcut for the MLP Adapt layer in NAS. More Adapt operations are learned.

Figure 10: Comparisons between CHEEM learned by two different implementations of the Adapt operation under the task-to-task (T2T) transfer based lifelong learning setting. Here all the 9 tasks are transferred from the base Tsk1\_ImNet model, so we omit the arrows linking the blocks for clarity. S, R, A and N represent Skip, Reuse, Adapt and New respectively.

## D.6 EFFECTS OF TASK ORDERS

We investigate the effects of task orders of the 9 tasks in the VDD benchmark. We test four more task sequences in addition to the one presented in the main paper. Overall, The CHEEM learned by our proposed HEE-based NAS achieve similar performance across different task orders, and consistently significantly outperform those learned by the vanilla PE-based NAS. We note that since the task ID inference is based on a frozen backbone, it is independent of the task order. Hence, for simplicity, we evaluate the effect of task order in a task-incremental setting.

Table 9 reports the performance. Fig. 11, Fig. 12, Fig. 13 and Fig. 14 show the learned CHEEM using our proposed HEE-based NAS.

Table 9: Results of ablation study on CHEEM learning with four different task orders using both our proposed HEE-based NAS and the vanilla PE-based NAS. The results have been averaged over 3 different seeds.

NAS	ImNet	OGIt	UCF	Airc.	Flwr	SVHN	DTD	GTSR	DPed	C100	Avg. Accuracy	Avg. Param. Inc./task (M)
HEE	82.65	84.32	75.27	54.32	90.29	95.83	57.89	99.92	99.72	89.96	83.02 ± 0.31	1.25 ± 0.15
PE	82.65	77.41	70.12	39.40	64.35	94.12	37.02	99.83	99.41	70.78	73.51 ± 0.80	2.86 ± 0.14
NAS	ImNet	DPed	SVHN	DTD	Airc.	OGIt	C100	GTSR	Flwr	UCF	Avg. Accuracy	Avg. Param. Inc./task (M)
HEE	82.66	99.94	95.83	58.56	42.43	83.55	89.98	99.95	91.99	75.67	82.06 ± 1.28	1.42 ± 0.05
PE	82.66	99.65	95.04	45.66	35.87	77.62	71.51	99.85	66.11	63.99	73.79 ± 0.50	2.85 ± 0.12
NAS	ImNet	UCF	C100	OGIt	GTSR	DTD	Flwr	SVHN	DPed	Airc.	Avg. Accuracy	Avg. Param. Inc./task (M)
HEE	82.66	79.73	90.75	84.93	99.90	58.14	91.27	96.05	99.89	54.06	83.74 ± 0.51	1.37 ± 0.05
PE	82.66	74.49	74.17	78.76	99.91	41.01	70.49	94.15	99.25	37.77	75.27 ± 2.41	2.75 ± 0.16
NAS	ImNet	Flwr	UCF	OGIt	GTSR	DPed	C100	Airc.	DTD	SVHN	Avg. Accuracy	Avg. Param. Inc./task (M)
HEE	82.66	87.52	77.17	84.20	99.92	99.80	90.30	54.49	56.83	96.03	82.89 ± 0.58	1.41 ± 0.09
PE	82.66	72.75	76.31	78.47	99.89	99.42	70.09	34.95	39.89	93.72	74.81 ± 1.40	2.70 ± 0.11

1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

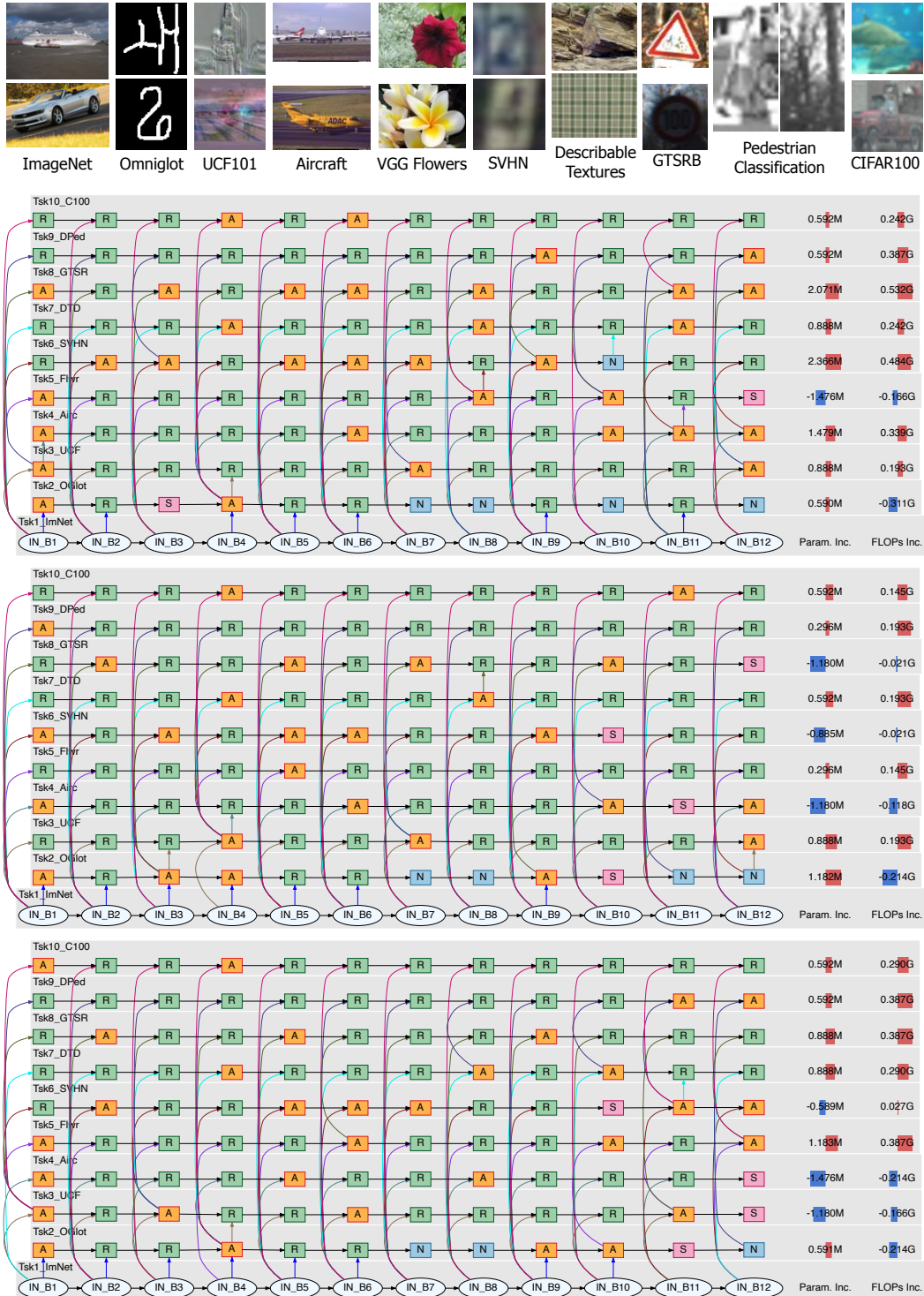


Figure 11: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi et al., 2017a) with the task sequence shown in the top **using our proposed HEE-based NAS** and three different random seeds. The overall performance is reported in Table 9. S, R, A and N represent Skip, Reuse, Adapt and New respectively. The last two columns show the number of new task-specific parameters and added FLOPs respectively, in comparison with the first task, ImNet model. Overall, the learned task synergies make intuitive sense and remain relatively stable across different random seeds.



1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

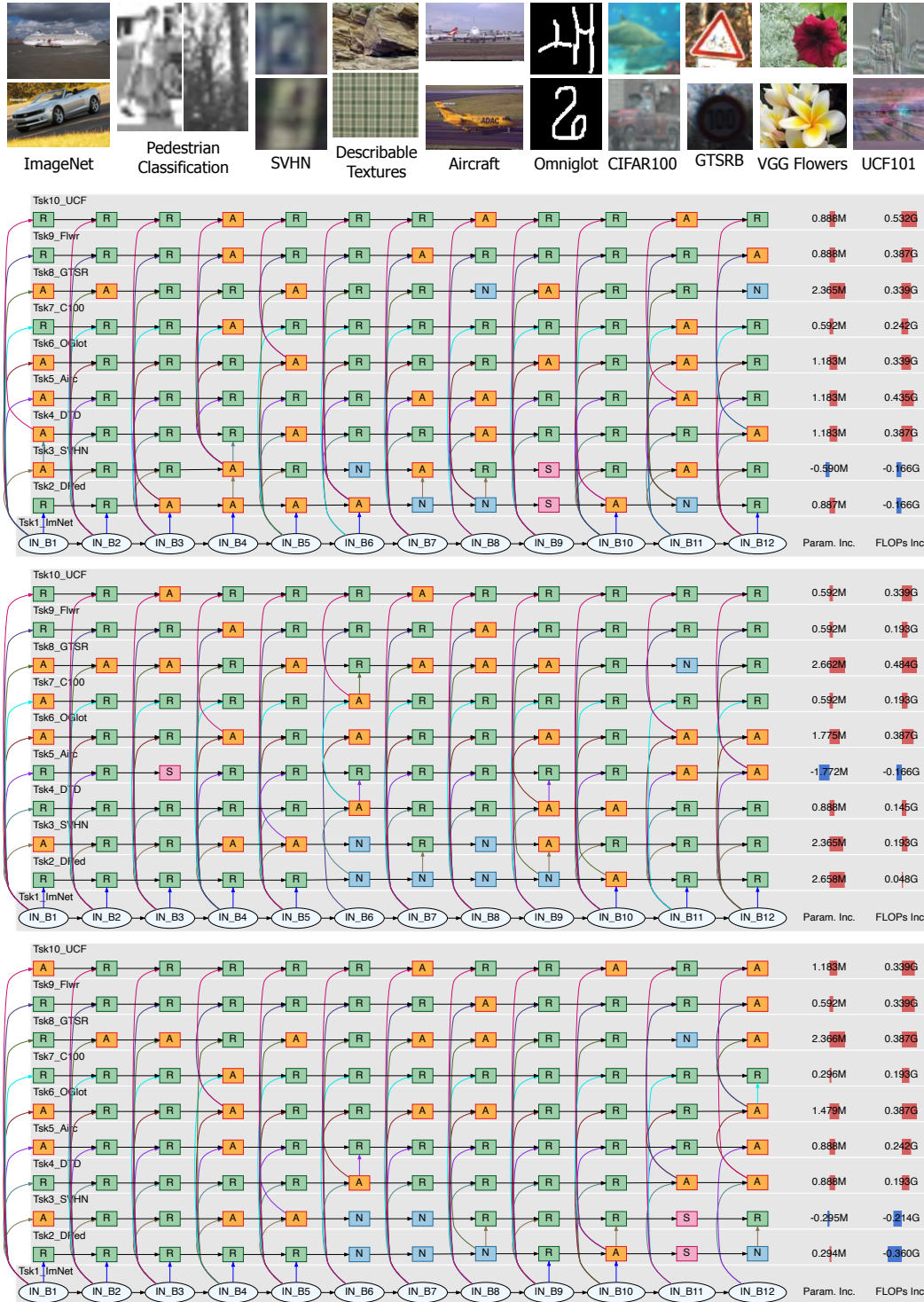


Figure 12: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi et al., 2017a) with the task sequence shown in the top using our proposed HEE-based NAS and three different random seeds. The overall performance is reported in Table 9. S, R, A and N represent Skip, Reuse, Adapt and New respectively. The last two columns show the number of new task-specific parameters and added FLOPs respectively, in comparison with the first task, ImNet model. Overall, the learned task synergies make intuitive sense and remain relatively stable across different random seeds.



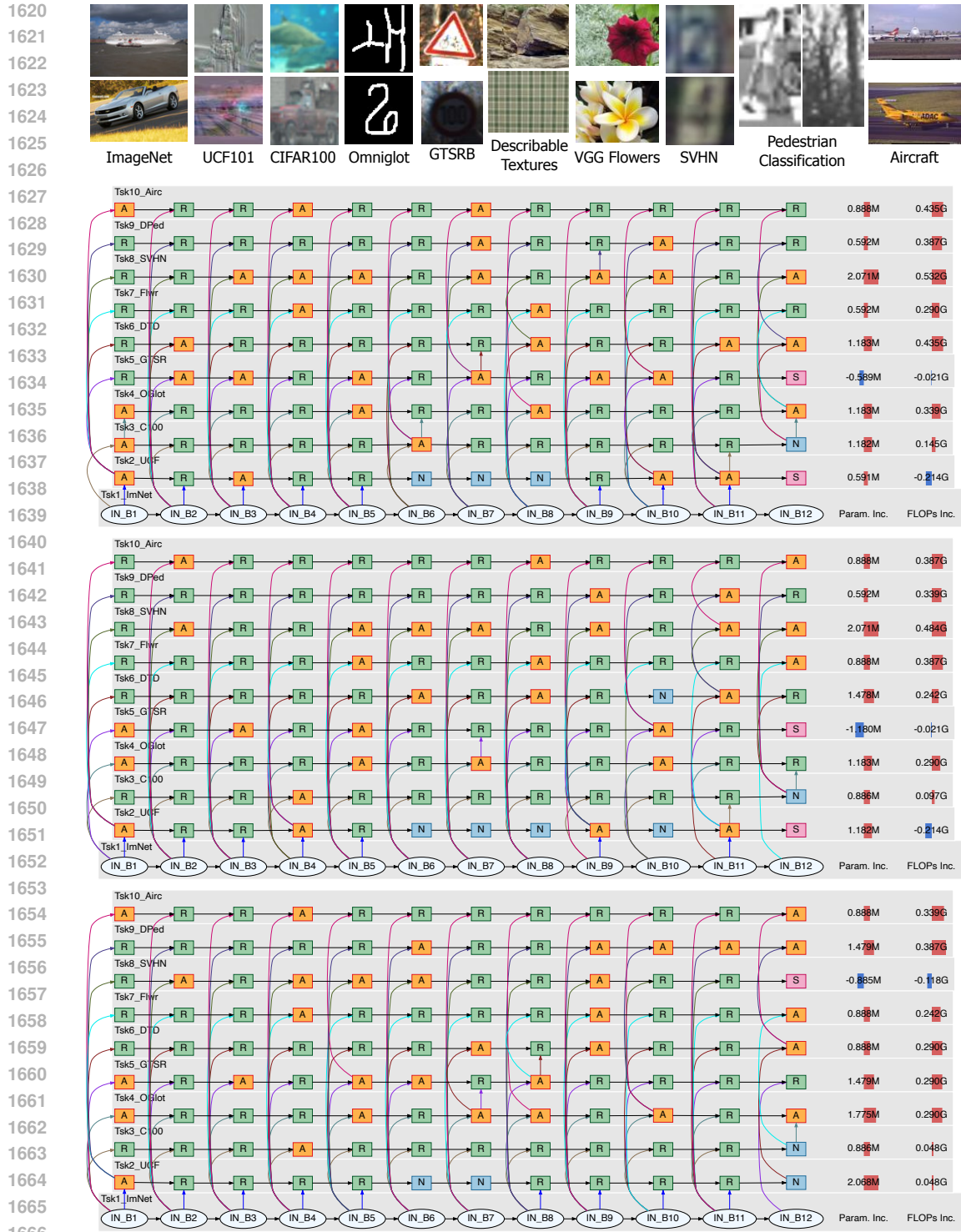


Figure 13: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi et al., 2017a) with the task sequence shown in the top **using our proposed HEE-based NAS** and three different random seeds. The overall performance is reported in Table 9. S, R, A and N represent Skip, Reuse, Adapt and New respectively. The last two columns show the number of new task-specific parameters and added FLOPs respectively, in comparison with the first task, ImNet model. Overall, the learned task synergies make intuitive sense and remain relatively stable across different random seeds.

1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727

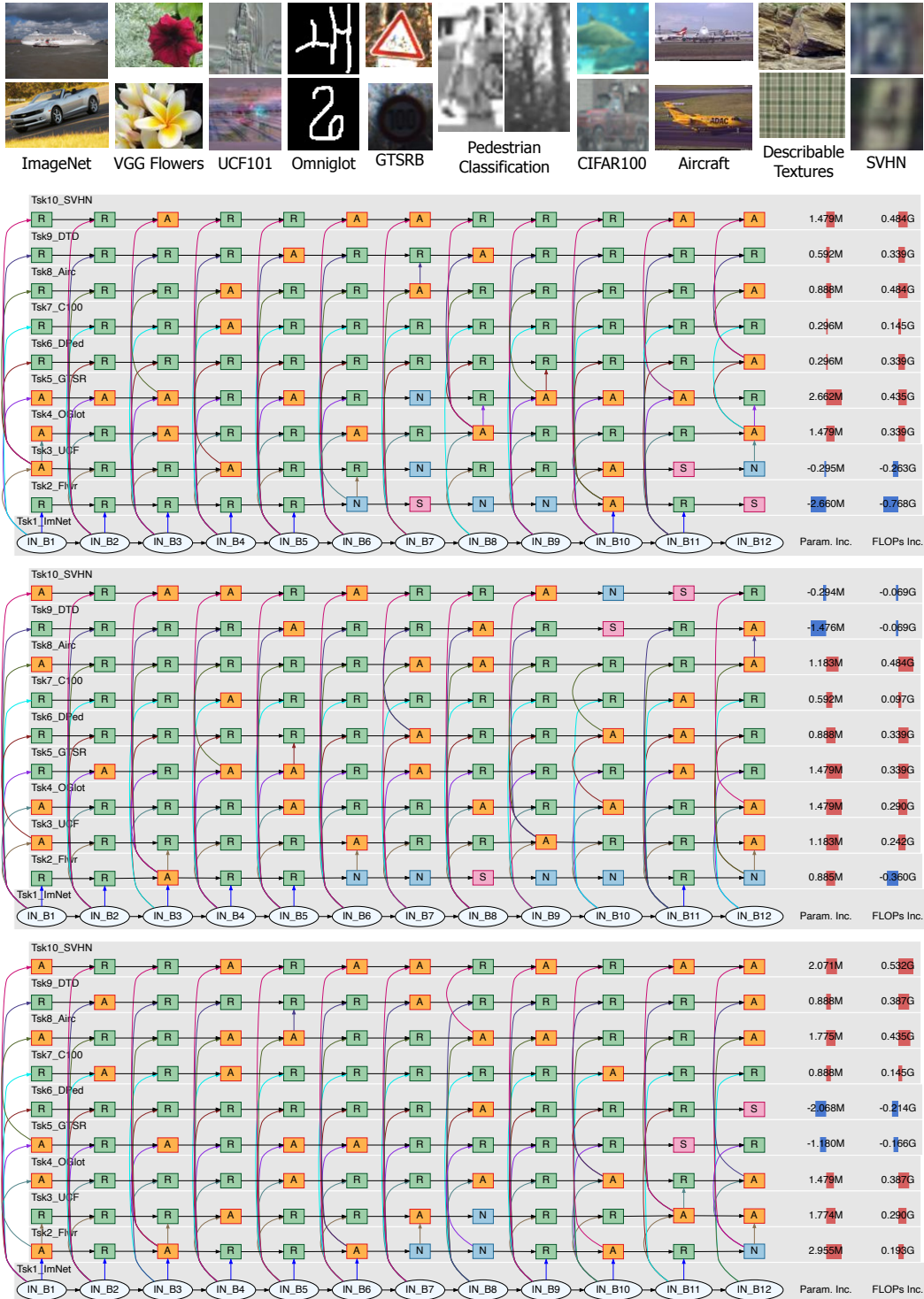


Figure 14: Examples of the task-synergy memory (CHEEM) learned on the VDD benchmark (Rebuffi et al., 2017a) with the task sequence shown in the top using our proposed HEE-based NAS and three different random seeds. The overall performance is reported in Table 9. S, R, A and N represent Skip, Reuse, Adapt and New respectively. The last two columns show the number of new task-specific parameters and added FLOPs respectively, in comparison with the first task, ImNet model. Overall, the learned task synergies make intuitive sense and remain relatively stable across different random seeds.

## E DATASET DETAILS

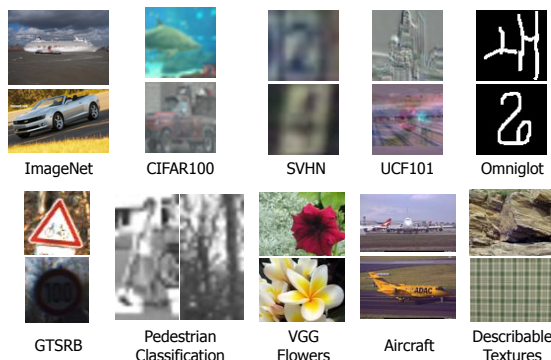


Figure 15: Example images from the VDD benchmark (Rebuffi et al., 2017a). Each task has a significantly different domain than others, making VDD a challenging benchmark for lifelong learning.

Table 10: The number of samples in training, validation and testing sets per task used in our experiments on the VDD benchmark (Rebuffi et al., 2017a).

Task	Train	Validation	Test	#Categories
ImageNet12	1108951	123216	49000	1000
CIFAR100	36000	4000	10000	10
SVHN	42496	4721	26040	10
UCF	6827	758	1952	101
Omniglot	16068	1785	6492	1623
GTSR	28231	3136	7842	43
DPed	21168	2352	5880	2
VGG-Flowers	918	102	1020	102
Aircraft	3001	333	3333	100
DTD	1692	188	1880	47

### E.1 THE VDD BENCHMARK

It consists of 10 tasks: ImageNet-1k (Russakovsky et al., 2015), CIFAR100 (Krizhevsky et al., 2009), SVHN (Netzer et al., 2011), UCF101 Dynamic Images (UCF) (Soomro et al., 2012; Bilen et al., 2016), Omniglot (Lake et al., 2015), German Traffic Signs (GTSR) (Stallkamp et al., 2012), Daimler Pedestrian Classification (DPed) (Munder & Gavrilu, 2006), VGG Flowers (Nilsback & Zisserman, 2008), FGVC-Aircraft (Maji et al., 2013), and Describable Textures (DTD) (Cimpoi et al., 2014). All the images in the VDD benchmark have been scaled such that the shorter side is 72 pixels. Table 10 shows the number of samples in each task. Fig. 15 shows examples of images from each task of the VDD benchmark. In our experiments, we use 10% of the official training data from each of the tasks for validation (e.g., used in the target network selection in Section 3.2.3 in main text), and report the accuracy on the official validation set due to the unavailability of the ground-truth labels for the official test data. In Table 10, the train, validation and test splits are thus referred to 90% of the official training data, 10% of the official training data, and the entire official validation data respectively. When finetuning the learned architecture (i.e., the searched target network) for each task, we use the entire official training data to train and report results on the official validation set.

### E.2 IMAGENET-R BENCHMARK

The ImageNet-R(etention) (Hendrycks et al., 2021) dataset contains art, cartoons, deviantart, graffiti, embroidery, graphics, origami, paintings, patterns, plastic objects, plush objects, sculptures, sketches, tattoos, toys, and video game renditions of ImageNet classes. It has renditions of 200 ImageNet classes resulting in 30,000 images. The Split ImageNet-R benchmark proposed by (Wang et al., 2022c) uses the ImageNet-R dataset to build a continual learning benchmark specifically for studying methods that use model pretrained on ImageNet (Russakovsky et al., 2015). ImageNet-R poses challenges for such methods because of the diversity within the same class. The Split-ImageNet benchmark proves to be challenging for experience-replay based approaches because of this intra-class variance, as well as replay-free methods that use a frozen backbone from ImageNet as the accuracy of the standard models on ImageNet-R is low. We use the same training and validation splits as those used by (Smith et al., 2023a). For the balanced evaluation, we divide the data set into 10 tasks with 20 classes each, and report results across 3 runs with random class orders. For the imbalanced evaluation, we construct 6 tasks with 5, 10, 15, 20, 50 and 100 classes. We report results across 3 runs with tasks, with varying task orders.

## F THE BASE VISION TRANSFORMER: ViT-B/8

We use the base Vision Transformer (ViT) model, with a patch size of  $8 \times 8$  (ViT-B/8) model from (Dosovitskiy et al., 2021). The base ViT model contains 12 Transformer blocks. A Transformer block is defined by stacking a Multi-Head Self-Attention (MHSA) block and a Multi-Layer Perceptron (MLP) block with residual connections for each block. ViT-B/8 uses 12 attention heads in each of the MHSA blocks, and a feature dimension of 768. The MLP block expands the dimension size to 3072 in the first layer and projects it back to 768 in the second layer. For all the experiments, we use an image size of  $72 \times 72$  following the VDD setting. We base the implementation of the ViT on the `timm` package (Wightman, 2019).

**Training the Base Model** To train the ViT-B/8 model, we use the ImageNet data provided by the VDD benchmark (the `train` split in Table 10). To save the training time, we initialize the weights from the ViT-B/8 trained on the full resolution ImageNet dataset ( $224 \times 224$ ) and available in the `timm` package, and finetune it for 30 epochs on the downsized version of ImageNet ( $72 \times 72$ ) in the VDD benchmark. We use a batch size of 2048 split across 4 Nvidia Quadro RTX 8000 GPUs. We follow the standard training/finetuning recipes for ViT models. The file `cheem/artifacts/imagenet_pretraining/args.yaml` in our code folder provides all the training hyperparameters used for training the the ViT-B/8 model on ImageNet. During testing, we take a single center crop of  $72 \times 72$  from an image scaled with the shortest side to scaled to 72 pixels.

## G SETTINGS AND HYPERPARAMETERS IN LEARNING CHEEM

Starting with the ImageNet trained ViT-B/8, the proposed CHEEM learning consists of three components: *supernet training*, *evolutionary search for target network selection*, and *target network finetuning*.

Table 11: Data augmentations for the 9 tasks in the VDD benchmark.

Task	Scale and Crop	Hor. Flip	Ver. Flip
CIFAR100	Yes	p=0.5	No
Aircraft	Yes	p=0.5	No
DPed	Yes	p=0.5	No
DTD	Yes	p=0.5	p=0.5
GTSR	Yes	p=0.5	No
OGIt	Yes	No	No
SVHN	Yes	No	No
UCF101	Yes	p=0.5	No
Flwr.	Yes	p=0.5	No

Table 12: Data augmentations used for each task in the 5-Datasets benchmark.

Task	Scale and Crop	Hor. Flip
MNIST	Yes	No
not-MNIST	Yes	No
SVHN	Yes	No
CIFAR100	Yes	p=0.5
Fashion MNIST	Yes	No

**Data Augmentations** A full list of data augmentations used for the VDD benchmark is provided in Table 11, and the data augmentations used for the tasks in the 5-datasets benchmark is provided in Table 12. The augmentations are chosen so as not to affect the nature of the data. Scale and Crop transformation scales the image randomly between 90% to 100% of the original resolution and takes a random crop with an aspect ratio sampled from a uniform distribution over the original aspect ratio  $\pm 0.05$ . In evaluating the supernet and the finetuned model on the validation set and test set respectively, images are simply resized to  $72 \times 72$  with bicubic interpolation.

**Supernet Training** *The VDD Benchmark:* For each task, we train the supernet for 100 epochs, unless otherwise stated. We use a label smoothing of 0.1. We use a learning rate of 0.001 and the Adam optimizer (Kingma & Ba, 2015) with a Cosine Decay Rule. We use a batch size of 512, and ensure the minimum number of batches in an epoch is 15 (via repeatedly sampling when the number of total samples of a task is not sufficient). As stated in the paper, for the Exploration-Exploitation sampling scheme, we use an exploration probability  $\epsilon = 0.3$ .

*The 5-datasets Benchmark:* We use the same hyperparameters as those used in the VDD Benchmark, but train the supernet for 50 epochs to account for its relatively lower complexity.

1836 *L2G with DARTS and  $\beta$ -DARTS*: We train the supernet of the Learn-to-Grow (L2G) (Li et al., 2019)  
1837 for 50 epochs on the VDD benchmark and 25 epochs on the 5-datasets benchmark, since DARTS  
1838 simultaneously trains all sub-networks (i.e. the entire supernet) at each epoch. We use a weight of 1  
1839 for the beta loss in all the experiments with  $\beta$ -DARTS.

1840  
1841 **Evolutionary Search** The evolutionary search is run for 20 epochs. We use a population size of  
1842 50. We use 25 candidates both in the mutation stage and the crossover stage. The top 50 candidates  
1843 are retained. The crossover is performed among the top 10 candidates, and the top 10 candidates  
1844 are mutated with a probability of 0.1. For the Exploration-Exploitation sampling scheme, we use an  
1845 exploration probability  $\epsilon = 0.5$  when generating the initial population.

1846  
1847 **Finetuning** The target network for a task selected by the evolutionary search is finetuned for 30  
1848 epochs with a learning rate of 0.001, Adam optimizer, and a Cosine Learning Rate scheduler. Drop  
1849 Path of 0.25 and label smoothing of 0.1 are used for regularization. We use a batch size of 512, and a  
1850 minimum of 30 batches are drawn.

1851 We use a single Nvidia A100 GPU for all the experiments.

1852

## 1853 H MODIFYING SUPSUP, EFT AND LL TO WORK WITH ViTs

1854

1855 In the main paper, we compare with Supermasks in Superposition (SupSup) (Wortsman et al., 2020),  
1856 Efficient Feature Transformation (EFT) (Verma et al., 2021), and Lightweight Learner (LL) (Ge et al.,  
1857 2023) in Table 5 under the task-to-task transfer learning paradigm. The three methods are originally  
1858 developed for Convolutional Neural Networks. We modify them to be compatible with ViTs for a fair  
1859 comparison with our CHEEM.

1860 We use the same ViT-B/8 base model (Sec. F) for SupSup, EFT and LL. For the SupSup  
1861 method (Wortsman et al., 2020), we learn masks for the weights of the final linear projection  
1862 layer of the Multi-Head Self-Attention block using the straight through estimator (Bengio et al.,  
1863 2013). We apply the EFT (Verma et al., 2021) on all the linear layers in the ViT-B/8 (i.e., all the  
1864 Query/Key/Value projection layers, the final projection layer, and the FFN layers) by scaling their acti-  
1865 vation maps via the Hadamard product with learnable scaling vectors, following the original proposed  
1866 formulation for fully-connected layers in the EFT (Verma et al., 2021). For the LL method (Ge et al.,  
1867 2023) which learns a task-specific bias vector that is added to all the feature maps of convolutional  
1868 layers, we learn a similar bias vector and add it to the output of all the linear layers of the ViT.

1869

1870

1871

1872

1873

1874

1875

1876

1877

1878

1879

1880

1881

1882

1883

1884

1885

1886

1887

1888

1889