

# HYPERADAPTER: GENERATING ADAPTERS FOR PRE-TRAINED MODEL-BASED CONTINUAL LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Humans excel at leveraging past experiences to learn new skills, while artificial neural networks suffer from the phenomenon of catastrophic forgetting during sequential learning. Efforts have been made to alleviate forgetting by introducing a rehearsal buffer into the model, but this way is impractical in real-world scenarios with data privacy. Recently, pre-trained model-based continual learning methods have provided new insights into addressing this issue by effectively utilizing the powerful representational capabilities of pre-trained models to avoid catastrophic forgetting without a rehearsal buffer. In this work, we propose a novel pre-trained model-based continual learning framework, HyperAdapter, which utilizes a hyper-network to generate adapters based on the current input, adapting the pre-trained model to the corresponding task. This paradigm requires fewer additional parameters as the number of tasks increases, which is a critical advantage for scaling to long sequences continual learning. Unlike methods that partition task-related knowledge into relatively independent subspaces, it promotes positive knowledge transfer across tasks. Comprehensive experiments across various datasets demonstrate that HyperAdapter consistently outperforms all existing methods and even exceeds the upper bounds of multi-task learning, establishing a new state-of-the-art for pre-trained model-based continual learning. Our code will be released.

## 1 INTRODUCTION

Humans exhibit remarkable abilities for constantly acquiring new knowledge in the dynamically changing real world, which helps them grasp new skills more easily with a richer knowledge base. However, when trained on successive task stages, neural networks tend to overfit on the current task and perform poorly on previous ones, a problem known as catastrophic forgetting (McCloskey & Cohen, 1989). Continual learning aims to learn new tasks while retaining past knowledge (De Lange et al., 2021; Mai et al., 2022). Inspired by the replay process in the hippocampus, some approaches rely on a rehearsal buffer to store samples from previous tasks (Buzzega et al., 2020; Cha et al., 2021; Chaudhry et al., 2019). These samples are then combined with data from the current task to mitigate forgetting. While these methods have shown promising results, they face challenges in real-world scenarios with privacy concerns (Shokri & Shmatikov, 2015) or memory constraints (Smith et al., 2021). The need for rehearsal-free methods to address continual learning in practical settings remains.

Without a rehearsal buffer, methods need to focus on the parameters and architecture of the model. Drawing inspiration from synaptic consolidation in the neocortex, EWC (Kirkpatrick et al., 2017) prevents forgetting of critical past knowledge by applying regularization to the model weights. While regularization-based methods (Kirkpatrick et al., 2017; Li & Hoiem, 2017) offer new insights into rehearsal-free continual learning, they alleviate forgetting at the expense of model plasticity, leading to suboptimal performance when learning new tasks. Recent advances in pre-training (Chen et al., 2020; He et al., 2020; Bao et al., 2021; Xie et al., 2022; He et al., 2022a) inspire the community to integrate pre-trained models into continual learning, aiming to prevent forgetting while efficiently learning new tasks. In particular, prompt-based methods (Wang et al., 2022b;a; Smith et al., 2023) can even outperform rehearsal-based methods. These methods maintain a prompt pool for the pre-trained backbone, selecting prompts based on input samples to instruct the learning of corresponding tasks. However, a small pool may lead to forgetting, while a large one burdens memory and hinders positive knowledge transfer between tasks, presenting the stability-plasticity dilemma (Jung et al., 2023). EASE (Zhou et al., 2024) introduces the adapter into continual learning to enhance the adaptability

054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

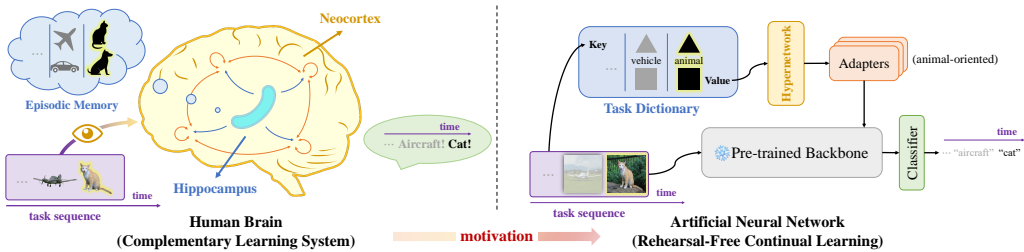


Figure 1: Motivation of HyperAdapter. Inspired by the CLS of human brain, we utilize a hypernetwork to generate adapters to adapt the pre-trained backbone to different tasks. In our framework, the task dictionary functions akin to episodic memory in the hippocampus, while the hypernetwork represents the neocortex, storing past knowledge. The rapidly updating task-specific embeddings and the slowly updating general hypernetwork work together to achieve rehearsal-free continual learning.

of pre-trained models across different tasks. However, it is impractical for learning large numbers of tasks, as the number of adapters involved during inference is proportional to the number of tasks.

In this paper, we propose a novel pre-trained model-based framework named HyperAdapter, for rehearsal-free continual learning. Our method consists of a pre-trained backbone, a hypernetwork, and a set of task embeddings. For any given input, representative features from the pre-trained model are leveraged as queries to identify the most similar task embedding, thus eliminating the necessity of knowing the task identities during inference. Subsequently, the hypernetwork generates a series of adapter parameters based on the obtained task embedding to adapt the pre-trained model to the corresponding tasks. Our design elegantly addresses the problems of existing methods while inheriting all their advantages. The frozen pre-trained model effectively prevents catastrophic forgetting, adapters enhance the model’s expressive ability and adaptive capacity compared to prompts, and the design of hypernetwork avoids the excessive number of adapters in EASE.

According to the Complementary Learning Systems (CLS) theory (Kumaran et al., 2016; McClelland et al., 1995), humans achieve continual learning through the synergy of two systems: the rapidly updating hippocampus focuses on learning task-specific representations, while the slowly updating neocortex specializes in learning more general representations based on past experiences. As shown in Figure 1, task embeddings can be likened to episodic memories in the hippocampus, selectively invoked based on different inputs, while the hypernetwork acts as the neocortex of brain, storing past knowledge in the form of neural connections and updating slowly through the indirect optimization of its generated adapters. Finally, the pre-trained model represents the prior knowledge acquired before learning, aiding the model in better acquiring new tasks. Without a fixed-size prompt pool, the hypernetwork effectively expands model capacity, facilitating positive knowledge transfer. Furthermore, for each new task, our method requires only the addition of a learnable task embedding vector, making it highly suitable for continual learning scenarios involving long sequences and large task numbers. Our empirical results demonstrate that HyperAdapter surpasses the performance of all existing works, offering a significant step forward in rehearsal-free continual learning. In our proposed CL-100 benchmark, HyperAdapter outperforms the previously best DAP and EASE by 2.24% and 3.07% in average final accuracy, respectively. On the larger ImageNet-R and DomainNet, HyperAdapter surpasses the previous SOTA CODA-P and EASE, by 1.06% and 4.80% respectively.

Our main contributions can be summarized as follows:

1. We propose HyperAdapter, a novel rehearsal-free continual learning framework. The method leverages a hypernetwork to generate adapters, adapting the pre-trained model to each task effectively, thereby mitigating forgetting and facilitating positive knowledge transfer.
2. Extensive experiments on various datasets demonstrate that HyperAdapter consistently outperforms all existing methods and even surpasses the multi-task learning upper bound in some cases, establishing a new SOTA for rehearsal-free continual learning. Moreover, the hypernetwork design makes it suitable for continual learning with longer task sequences.
3. To the best of our knowledge, this is the first work leveraging hypernetworks to unlock the potential of pre-trained models in the field of continual learning. We expect that our approach provides a novel perspective on continual learning of pre-trained models.

## 2 RELATED WORK

### 2.1 PRE-TRAINED MODEL-BASED CONTINUAL LEARNING

In recent years, continual learning has emerged as a focal point in the field of machine learning, with the primary challenge being how to incorporate new information effectively without forgetting prior knowledge. Traditional continual learning methods (Rolnick et al., 2019; Rebuffi et al., 2017; Lopez-Paz & Ranzato, 2017; Chaudhry et al., 2018b; Aljundi et al., 2019; Chaudhry et al., 2019) typically rely on the rehearsal of old data, which can raise privacy and storage issues. Although carefully designed regularizations (Li & Hoiem, 2017; Lopez-Paz & Ranzato, 2017) have addressed forgetting to some extent, regularization-based methods still underperform rehearsal-based ones.

With advancements in model pre-training, an increasing number of studies begin to explore integrating knowledge from pre-trained models into continual learning, achieving comparable results even without a rehearsal buffer. Some approaches (Wang et al., 2022b;a; Smith et al., 2023; Jung et al., 2023) facilitate continual learning by providing appropriate visual prompts (Jia et al., 2022) to pre-trained models. However, these prompt-based methods typically focus on instance-level enhancements, which offer limited overall benefits to the model. Recently, adapter-based techniques (Zhou et al., 2024; Gao et al., 2024) have shown promising results, outperforming prompt-based ones. The modular design of adapters (Pfeiffer et al., 2021; Chen et al., 2022) allows to retain and leverage the knowledge from pre-trained models more effectively, although they also suffer from limited adaptability and scalability. Leveraging a hypernetwork to generate adapters, HyperAdapter proposes a novel rehearsal free mechanism with the capable of generalization.

### 2.2 COMPLEMENTARY LEARNING SYSTEMS

The Complementary Learning Systems (CLS) theory reveals that humans achieve continual learning through the synergy of two systems that update at different frequencies. Inspired by CLS, several methods (Parisi et al., 2018; Pham et al., 2021; Arani et al., 2022) incorporate multiple networks along with rehearsal buffers, regularization constraints, or other components that expand with the task number. FearNet (Kemker & Kanan, 2017), for instance, employs a three-network structure: one hippocampal network for recalling recent instances, one PFC network for long-term memories, and one additional network for deciding between the two for specific cases. Gomez-Villa et al. (2024) proposes to train an expert network that, unburdened by the task of retaining prior knowledge, focuses on excelling in new tasks. Closer to our approach, Gurbuz et al. (2024) introduces a contextual encoding for rehearsal-free scenarios. However, all the methods above do not take pre-trained models into consideration, thus lack of powerful representational capabilities and outstanding performance.

### 2.3 HYPERNETWORKS

Hypernetworks (Ha et al., 2016) were initially proposed for network compression, aimed at generating smaller sets of weights to reduce the computational and storage demands. In efficient fine-tuning (Mahabadi et al., 2021; Zhmoginov et al., 2022; He et al., 2022b; Zhang et al., 2022; Üstün et al., 2022), federated learning (Zhang et al., 2023) and few-shot learning tasks (Sendera et al., 2023), hypernetworks have been employed to dynamically generate parameters, achieving significant performance. These applications demonstrate that hypernetworks can not only adjust models for new tasks but also retain memory of old knowledge. Following this strategy, some hypernet-based continual learning methods (Chandra et al., 2023; Książek & Spurek, 2023; Hemati et al., 2023) generate new parameters for the entire network to adapt to new tasks, but these methods may not fully exploit the prior knowledge of pre-trained models. Different from previous works, HyperAdapter leverages representative features from pre-trained models, thus eliminating necessity of knowing the task identities during inference or the dependence on any rehearsal buffers.

## 3 PRELIMINARIES

### 3.1 CONTINUAL LEARNING

Continual learning requires the model to learn a sequence of tasks in order of arrival. For a sequence of  $T$  tasks  $\mathcal{D} = \{\mathcal{D}_1, \dots, \mathcal{D}_T\}$ , where each task  $\mathcal{D}_t = \{(x, y)\}$  contains tuples of the input sample

$x \in \mathcal{X}$  and its corresponding label  $y \in \mathcal{Y}$ , a single model  $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$  parameterized by  $\theta$  needs to learn on  $\mathcal{D}$  sequentially and predicts the label  $y = f(x; \theta)$  given input sample  $x$  from arbitrary task. Data from previous tasks will not be seen anymore when training future tasks, requiring the model to avoid forgetting old knowledge while learning new tasks. Depending on the differences in task transition, continual learning can be further categorized into multiple settings. In this work, we focus on the more challenging class-incremental learning setting, where task identities are only known during training, which is more common in real-world scenarios.

### 3.2 ADAPTER TUNING

Adapter Tuning, which first emerged in NLP (Houlsby et al., 2019), achieves efficient adaptation to downstream tasks by freezing pre-trained weights and inserting lightweight bottleneck modules into the model. AdaptFormer (Chen et al., 2022) later introduces adapters to visual recognition tasks. The vanilla Vision Transformer consists of a series of blocks, each containing a multi-head self-attention layer (MHSA) and a feed-forward network (FFN). The processing of FFN can be formalized as:

$$\mathbf{x} = \text{MLP}(\text{LN}(\mathbf{x}')) + \mathbf{x}' \quad (1)$$

where  $\mathbf{x}' \in \mathbb{R}^d$  is the output of MHSA in the same block. AdaptFormer replaces the MLP in FFN with AdaptMLP, which includes an adapter with a down-projection matrix  $\mathbf{D} \in \mathbb{R}^{d \times r}$  and an up-projection matrix  $\mathbf{U} \in \mathbb{R}^{r \times d}$ , where  $r$  is the bottleneck dimension satisfying  $r \ll d$ . The process of extracting adapted features can be represented as:

$$\tilde{\mathbf{x}} = \text{GELU}(\mathbf{x}' \cdot \mathbf{D}) \cdot \mathbf{U} \quad (2)$$

All features are then passed through residual connection to obtain the final output of AdaptMLP:

$$\mathbf{x} = \text{MLP}(\text{LN}(\mathbf{x}')) + s \cdot \tilde{\mathbf{x}} + \mathbf{x}' \quad (3)$$

where  $s$  is a scale factor to ensure convergence.

## 4 METHOD

We propose HyperAdapter as a hypernetwork-based method for rehearsal-free continual learning. The overall framework is illustrated in Figure 2. We first introduce task-conditional embeddings and the query-key matching mechanism in Section 4.1. Then, we explain how the hypernetwork utilizes these embeddings to generate task-oriented adapters in Section 4.2. And we further extend the hypernetwork to a block-wise implementation in Section 4.3, significantly reducing forgetting and improving the performance. Finally, we present the overall optimization objective in Section 4.4.

### 4.1 TASK-CONDITIONAL EMBEDDINGS

The hypernetwork requires task-conditional embeddings as inputs to generate task-specific adapter parameters. These embeddings serve as episodic memories in the hippocampus, selectively activated based on different inputs, and fed into the hypernetwork storing past knowledge to accomplish any task. To achieve input-dependent selection of task embeddings, we maintain a dictionary  $\mathcal{C} = \{(\mathbf{k}_1, \mathbf{z}_1), (\mathbf{k}_2, \mathbf{z}_2), \dots, (\mathbf{k}_T, \mathbf{z}_T)\}$  of length  $T$ , where each entry contains a key  $\mathbf{k} \in \mathbb{R}^d$  and a value  $\mathbf{z} \in \mathbb{R}^e$ , both being learnable parameter vectors. For task embedding selection, we employ a query-key matching mechanism following the previous work. Given a query  $\mathbf{q} \in \mathbb{R}^d$ , we search for the closest key in the dictionary:

$$\tilde{\mathbf{k}} = \arg \min_{\mathbf{k} \in \mathbf{K}} \gamma(\mathbf{q}, \mathbf{k}) \quad (4)$$

where  $\mathbf{K} = \{\mathbf{k}_i\}_{i=1}^T$  is the set of all keys in the dictionary  $\mathcal{C}$  and  $\gamma$  is the distance metric (we use the negative cosine similarity in our experiments). The corresponding  $\tilde{\mathbf{z}}$  for the retrieved  $\tilde{\mathbf{k}}$  is then the task embedding corresponding to that query. This design decouples the learning of keys and values and allows adjusting the total parameters of the hypernetwork through the dimension of task embeddings, providing greater flexibility to the entire framework.

For any given input, we aim for the generated query to be relatively fixed, as this facilitates the learning of task keys. Hence, we simply utilize the pre-trained model as a frozen feature extractor

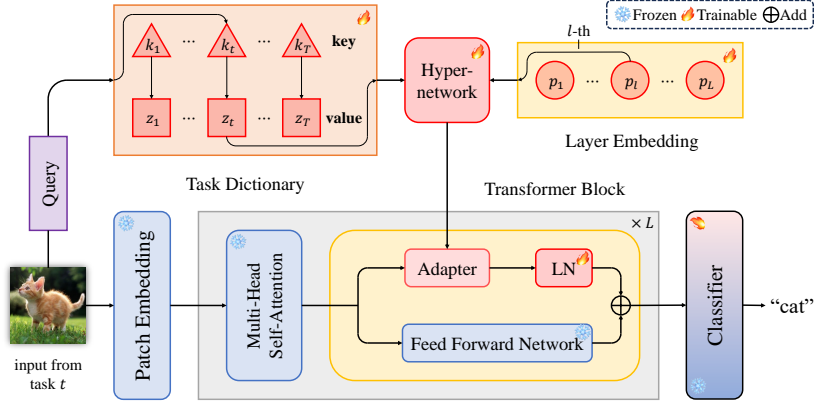


Figure 2: Framework of HyperAdapter. First, inputs from any task are matched to the most similar task embedding through a query. It is then fed into the hypernetwork along with the layer position embedding to generate a set of adapters. The input image is processed through the model with these adapters, and the output is finally obtained via the classifier. We do not update the classifier weights related to previous tasks. And layer embeddings are not included in the block-wise HyperAdapter.

$q(\mathbf{x}) = f(\mathbf{x})[\text{CLS}]$ , where  $[\text{CLS}]$  represents the class token. In this way, the same input consistently yields the same query, preventing forgetting during the query generation. Moreover, a powerful pre-trained feature extractor also ensures that similar inputs produce similar queries, which aids the model in obtaining the correct task embeddings. And the matching loss can be formulated as:

$$\mathcal{L}_{\text{match}}(\mathbf{x}, \mathbf{k}_t) = \gamma(q(\mathbf{x}), \mathbf{k}_t), \quad \mathbf{x} \in \mathcal{D}_t \quad (5)$$

During training, since the task identities are known, we use ground truth task embeddings  $z_t$  as input to the hypernetwork. In inference, the trained task dictionary  $\mathcal{C}$  can be used to obtain input-dependent task embeddings  $\tilde{z}$ , eliminating the need for task identities and making the entire method applicable to the more challenging class-incremental learning.

## 4.2 TASK-ORIENTED HYPER-ADAPTERS

With the task embeddings, we can adapt the pre-trained model to each task for continual learning using task-oriented hyper-adapters. According to the theory of CLS, humans learn continually through the synergy of two systems. Here we use a hypernetwork to play the role of the neocortex, taking the task embedding  $z_t$  as input and generating a set of adapter parameters  $\{(D_t^l, U_t^l)\}_{l=1}^L$  for the task  $\mathcal{D}_t$ . This design enables the hypernetwork to continuously retain past knowledge and promote positive knowledge transfer through information sharing across different tasks.

To increase the model capacity, we introduce a set of positional embeddings  $\mathcal{P} = \{\mathbf{p}^l \in \mathbb{R}^e\}_{l=1}^L$  for each layer. In practice, we simply add these positional embeddings to the task embedding:

$$\mathbf{I}_t^l = z_t + \mathbf{p}^l \quad (6)$$

With this combined embedding  $\mathbf{I}_t^l \in \mathbb{R}^e$ , the hypernetwork can generate different parameters for each layer, helping the model better adapt to specific tasks. Furthermore, to prevent forgetting, considering that different  $\mathbf{p}^l$  are used only to help the hypernetwork distinguish between different layers, we freeze all the positional embeddings  $\mathcal{P}$  after learning the first task  $\mathcal{D}_1$ .

The neocortex stores knowledge in the form of neural connections. Inspired by this, we use a large linear layer as the hypernetwork  $h$ . Specifically, for the down- and up-projection layers of the adapter, we use two linear layers  $\mathbf{W}_D \in \mathbb{R}^{(d \times r) \times e}$  and  $\mathbf{W}_U \in \mathbb{R}^{(r \times d) \times e}$  to generate the adapter parameters:

$$(\mathbf{D}_t^l, \mathbf{U}_t^l) = h(\mathbf{I}_t^l) = (\mathbf{W}_D, \mathbf{W}_U)\mathbf{I}_t^l \quad (7)$$

Additionally, we include an extra layer normalization  $\text{LN}^l$  after the adapter to enhance the model's ability to learn new tasks, which is crucial for the performance of continual learning, as detailed in the ablation study. Overall, our task-oriented hyper-adapters can be represented as:

$$\tilde{\mathbf{x}} = \text{LN}^l(\text{GELU}(\mathbf{x}' \cdot \mathbf{D}_t^l) \cdot \mathbf{U}_t^l) \quad (8)$$

### 270 4.3 BLOCK-WISE HYPER-ADAPTERS

271  
272 While incorporating a single hypernetwork for the entire model is consistent with the original  
273 intention of the hypernetwork design, which is to compress the parameters of the network, we find  
274 that introducing a separate hypernetwork for each layer significantly enhances the continual learning  
275 performance. Since continual learning itself does not impose strict requirements on the number  
276 of parameters, this approach effectively decouples task-specific knowledge from position-related  
277 knowledge, further mitigating forgetting. Particularly, we introduce a hypernetwork  $h^l$  for each  
278 layer of the model, comprising two linear layers  $\mathbf{W}_D^l \in \mathbb{R}^{(d \times r) \times e}$  and  $\mathbf{W}_U^l \in \mathbb{R}^{(r \times d) \times e}$ , to generate  
279 adapter parameters for a specific layer:

$$280 (\mathbf{D}_t^l, \mathbf{U}_t^l) = h^l(\mathbf{z}_t) = (\mathbf{W}_D^l, \mathbf{W}_U^l) \mathbf{z}_t \quad (9)$$

281  
282 The block-wise hyper-adapters improve performance at the cost of increased parameters. We refer to  
283 this method as  $\text{HA}_{\text{block}}$ , whereas the method described in Section 4.2 is referred to as  $\text{HA}_{\text{model}}$ .

### 284 4.4 OPTIMIZATION OBJECTIVE

285  
286 During training, for an input  $\mathbf{x}$  from task  $\mathcal{D}_t$ , we use the ground truth task key  $\mathbf{k}_t$  to compute the  
287 matching loss in the class-incremental setting. Then we combine the input with the corresponding task  
288 embedding  $\mathbf{z}_t$  and positional embeddings  $\mathcal{P}$  (excluded in  $\text{HA}_{\text{block}}$ ) and pass it through the pre-trained  
289 model  $f_\theta$  and the hypernetwork  $h$ . The output is then fed into the classifier  $g_\phi$  for prediction, are  
290 compared with the labels  $y$  to calculate the task loss. In summary, the overall optimization objective  
291 for HyperAdapter can be expressed as:

$$292 \min_{\mathcal{C}, \mathcal{P}, \mathbf{W}_D, \mathbf{W}_U, \phi_t} \mathcal{L}(g_{\phi_t} \circ f(\mathbf{x}; \theta, h(\mathbf{I}_t; \mathbf{W}_D, \mathbf{W}_U)), y) + \lambda \mathcal{L}_{\text{match}}(\mathbf{x}, \mathbf{k}_t), \quad (\mathbf{x}, y) \in \mathcal{D}_t \quad (10)$$

293  
294 where  $\mathcal{L}$  represents the task loss,  $\mathcal{L}_{\text{match}}$  is the matching loss defined in Section 4.1, and  $\lambda$  is a scalar  
295 loss weight balancing term (we set it to 0.1 by default). The parameters to be updated include the  
296 task dictionary  $\mathcal{C}$ , the positional embeddings  $\mathcal{P}$  (frozen after learning the first task), the hypernetwork  
297  $\mathbf{W}_D, \mathbf{W}_U$ , and the classifier  $\phi$ . Notably, consistent with previous works, we do not update the  
298 weights in the classifier related to past tasks, with the remaining parts for the  $t$ -th task denoted as  $\phi_t$ .

## 300 5 EXPERIMENTS

### 301 5.1 EXPERIMENTAL SETUP

302  
303 **Datasets.** We conduct extensive experiments on seven datasets with varying data scales and task  
304 sequence lengths. The first five datasets are uniformly restricted to 100 classes each and combined to  
305 form the CL-100 benchmark. More details of the datasets are provided in the appendix.

- 306  
307
- 308 • **CL-100 Benchmark:** We select 5 commonly used datasets and combined them into a new  
309 benchmark CL-100, for a more comprehensive evaluation. These datasets vary in scale,  
310 listed from small to large: Oxford Flowers 102 (Nilsback & Zisserman, 2008), Caltech-  
311 101 (Fei-Fei et al., 2006), Stanford Dogs (Dataset, 2011), CIFAR-100 (Krizhevsky et al.,  
312 2009), and Food-101 (Bossard et al., 2014). We restrict these datasets to 100 classes each,  
313 with each class containing between 20 to 750 training images. For continual learning, the  
314 CL-100 benchmark is splitted into 10 tasks by default, each containing 10 classes.
- 315 • **Large Benchmarks:** To validate the effectiveness of methods in more challenging scenarios,  
316 we include two larger benchmarks: Split ImageNet-R (Hendrycks et al., 2021) and Split  
317 DomainNet (Peng et al., 2019). The Split ImageNet-R benchmark is build upon the test  
318 set of the ImageNet-R, containing a total of 200 classes and 30,000 images. It is splitted  
319 into 10 tasks by default, each with 20 classes. The Split ImageNet-R dataset presents a  
320 substantial divergence from the dataset used for backbone pre-training (i.e. ImageNet-21k),  
321 and the significant intra-class diversity imposes higher demands on rehearsal-free continual  
322 learning methods. The Split DomainNet benchmark comprises over a hundred thousand  
323 images across 345 classes from 6 different domains. It is splitted into 15 tasks by default,  
each containing 23 classes. The Split DomainNet dataset includes images in various styles,  
which is beneficial for evaluating the generalization ability of model across distinct domains.

**Comparing Baselines.** We compare our method against several rehearsal-free baselines and SOTA methods, including regularization- and architecture-based (e.g. prompt and adapter) approaches. All methods use the same ViT-B/16 backbone pre-trained on ImageNet-21K (Deng et al., 2009). For completeness, we also include two naive baselines. More details can be found in the appendix.

- **Naive Baselines:** We present Full-seq and Linear-seq as two naive baselines. Full-seq denotes the fully sequential training, while Linear-seq refers to a version based on the pre-trained backbone, where only the classifier is updated during sequential training.
- **Regularization-Based Methods:** We choose the classical EWC (Kirkpatrick et al., 2017) and LwF (Li & Hoiem, 2017) as our regularization-based baselines, both of which introduce certain regularization to model parameters during sequential training to prevent forgetting.
- **Prompt-Based Methods:** There are many prompt-based methods utilizing pre-trained models for continual learning. We select L2P (Wang et al., 2022b), DualPrompt (Wang et al., 2022a), CODA-Prompt (Smith et al., 2023), and DAP (Jung et al., 2023) as the current SOTA methods for comparison.
- **Adapter-Based Methods:** We choose EASE (Zhou et al., 2024), the only adapter-based continual learning method prior to our work for a fair comparison.
- **Upper Bound:** Following previous works, we use the multi-task learning results as the upper bound to demonstrate the exceptional performance of our method.

**Evaluation Metrics.** Following DAP, we employ three widely used metrics for method evaluation: final accuracy (FnI. Acc.  $\uparrow$ ) of the accuracy after the last task as final performance, forgetting rate (Forgetting  $\downarrow$ ) of the ability to alleviate forgetting (negative transfer inhibition), and learning accuracy (Lrn. Acc.  $\uparrow$ ) of the ability to acquire new information (positive transfer promotion). We repeat each experiment three times and report the average values with standard errors. Please refer to the appendix for more details on the evaluation metrics.

**Implementation Details.** We train HyperAdapter using Adam with  $\beta_1, \beta_2$  of 0.9, learning rate of 0.01 and batch size of 128. All input images are resized to  $224 \times 224$ . For smaller datasets, we train every task for 30 epochs to ensure convergence, while for others, we train each task for 10 epochs. In the task dictionary, we maintain only 1 key and 1 embedding for each task. Following previous works, we use a single classifier and do not update the weights corresponding to classes of past tasks. For model-wise HyperAdapter, we set the dimension of both task embeddings and adapter bottleneck to 32, while for block-wise HyperAdapter, we set them to 16. The hyperparameters  $s$  in Equation 3 and  $\lambda$  in Equation 10 are both set to 0.1. All experiments are conducted on NVIDIA A100 GPUs.

## 5.2 MAIN RESULTS

**CL-100 Benchmark.** Table 1 presents a comprehensive comparison between HyperAdapter and other methods on CL-100 benchmark. Our method consistently outperforms all baselines, establishing a new SOTA for rehearsal-free continual learning. In Table 1(a), which shows the final accuracy, HyperAdapter leads all other methods. Specifically, it surpasses the regularization-based methods EWC and LwF by 37.20% and 32.39% in average accuracy. Moreover, compared to more advanced prompt-based and adapter-based methods, HyperAdapter exceeds the previous best DAP and EASE by 2.24% and 3.07%. The block-wise HyperAdapter further achieves a 1.77% improvement over model-wise HyperAdapter through more parameters. Notably, compared to the multi-task learning upper bound, HyperAdapter even achieved an improvement of 0.78-2.94% on larger datasets (with more than 10,000 training images).

Table 1(b) shows the forgetting performance. With the help of a pre-trained backbone and task dictionary, the model-wise HyperAdapter significantly reduced the forgetting rate by more than 20% compared to regularization-based methods. Benefiting from the decoupling of task and position information, block-wise HyperAdapter further reduced the forgetting rate by 1.11%, reaching an astonishing 1.30%, significantly lower than the previous best DAP (2.26%) and EASE (3.18%). Table 1(c) reflects the ability to learn new tasks based on old knowledge. The design of the hyper-network enables HyperAdapter to better facilitate positive knowledge transfer between tasks. Both versions of HyperAdapter surpass DAP and EASE, as well as Full-seq, in learning accuracy.

Table 1: Results on CL-100 benchmark. Datasets are sorted with their scales in ascending order.

| Dataset                   | Split Flowers-100                                    | Split Caltech-100                  | Split Dogs-100                     | Split CIFAR-100                    | Split Food-100                     |                                       |
|---------------------------|--|------------------------------------|------------------------------------|------------------------------------|------------------------------------|---------------------------------------|
| <b>Method</b>             | <b>(a) Final Accuracy (<math>\uparrow</math>)</b>    |                                    |                                    |                                    |                                    | <b>Mean (<math>\uparrow</math>)</b>   |
| Full-seq                  | 35.64 $\pm$ 1.92                                     | 27.04 $\pm$ 1.25                   | 22.67 $\pm$ 0.96                   | 30.39 $\pm$ 1.92                   | 26.90 $\pm$ 0.51                   | 28.53                                 |
| Linear-seq                | 78.95 $\pm$ 0.58                                     | 76.17 $\pm$ 0.08                   | 66.22 $\pm$ 0.16                   | 68.43 $\pm$ 0.09                   | 60.58 $\pm$ 0.32                   | 70.07                                 |
| EWC                       | 69.79 $\pm$ 1.76                                     | 57.96 $\pm$ 1.83                   | 45.72 $\pm$ 1.26                   | 59.60 $\pm$ 1.27                   | 55.27 $\pm$ 1.06                   | 57.67                                 |
| LwF                       | 71.78 $\pm$ 1.98                                     | 63.26 $\pm$ 1.37                   | 48.97 $\pm$ 1.23                   | 68.22 $\pm$ 1.63                   | 60.15 $\pm$ 0.66                   | 62.48                                 |
| L2P                       | 94.53 $\pm$ 1.23                                     | 89.34 $\pm$ 1.78                   | 76.59 $\pm$ 1.32                   | 83.05 $\pm$ 1.02                   | 70.48 $\pm$ 1.38                   | 82.80                                 |
| DualPrompt                | 95.25 $\pm$ 0.82                                     | 91.52 $\pm$ 0.87                   | 78.36 $\pm$ 0.63                   | 84.77 $\pm$ 0.69                   | 75.31 $\pm$ 0.57                   | 85.04                                 |
| CODA-P                    | 97.02 $\pm$ 0.39                                     | 91.44 $\pm$ 0.34                   | 84.41 $\pm$ 1.44                   | 86.25 $\pm$ 0.74                   | 77.58 $\pm$ 1.18                   | 87.34                                 |
| DAP                       | 96.49 $\pm$ 0.05                                     | 97.23 $\pm$ 0.33                   | 87.02 $\pm$ 1.12                   | 94.05 $\pm$ 1.19                   | 88.37 $\pm$ 0.63                   | 92.63                                 |
| EASE                      | 97.53 $\pm$ 0.16                                     | 96.54 $\pm$ 0.47                   | 88.31 $\pm$ 0.82                   | 87.81 $\pm$ 0.23                   | 88.82 $\pm$ 1.47                   | 91.80                                 |
| <b>HA<sub>model</sub></b> | 97.57 $\pm$ 0.68                                     | 96.44 $\pm$ 0.15                   | 87.14 $\pm$ 0.29                   | 95.85 $\pm$ 0.37                   | 88.64 $\pm$ 0.54                   | 93.13                                 |
| <b>HA<sub>block</sub></b> | <b>98.04 <math>\pm</math> 0.42</b>                   | <b>97.66 <math>\pm</math> 0.24</b> | <b>89.32 <math>\pm</math> 0.68</b> | <b>97.14 <math>\pm</math> 0.35</b> | <b>92.20 <math>\pm</math> 0.45</b> | <b>94.87</b>                          |
| Upper Bound               | 98.83 $\pm$ 0.75                                     | 98.32 $\pm$ 0.16                   | 88.54 $\pm$ 0.35                   | 94.20 $\pm$ 0.80                   | 90.40 $\pm$ 1.24                   | 94.06                                 |
| <b>Method</b>             | <b>(b) Forgetting Rate (<math>\downarrow</math>)</b> |                                    |                                    |                                    |                                    | <b>Mean (<math>\downarrow</math>)</b> |
| Full-seq                  | 67.53 $\pm$ 0.53                                     | 68.97 $\pm$ 1.14                   | 73.12 $\pm$ 1.27                   | 69.21 $\pm$ 0.18                   | 75.70 $\pm$ 0.55                   | 70.91                                 |
| Linear-seq                | 21.93 $\pm$ 0.14                                     | 23.66 $\pm$ 0.26                   | 25.78 $\pm$ 0.35                   | 17.34 $\pm$ 0.57                   | 20.69 $\pm$ 0.61                   | 21.88                                 |
| EWC                       | 23.73 $\pm$ 3.00                                     | 26.34 $\pm$ 3.39                   | 29.26 $\pm$ 1.64                   | 24.65 $\pm$ 0.07                   | 23.67 $\pm$ 1.12                   | 25.53                                 |
| LwF                       | 25.14 $\pm$ 2.42                                     | 24.63 $\pm$ 0.75                   | 33.42 $\pm$ 1.86                   | 15.44 $\pm$ 1.48                   | 17.15 $\pm$ 0.62                   | 23.16                                 |
| L2P                       | 2.12 $\pm$ 0.14                                      | 4.86 $\pm$ 0.49                    | 6.93 $\pm$ 0.44                    | 7.21 $\pm$ 0.16                    | 9.09 $\pm$ 0.44                    | 6.04                                  |
| DualPrompt                | 1.40 $\pm$ 0.32                                      | 2.85 $\pm$ 0.31                    | 4.84 $\pm$ 0.11                    | 5.60 $\pm$ 0.40                    | 8.76 $\pm$ 0.21                    | 4.69                                  |
| CODA-P                    | 1.16 $\pm$ 0.18                                      | 2.73 $\pm$ 0.48                    | 4.03 $\pm$ 2.11                    | 4.67 $\pm$ 0.26                    | 7.58 $\pm$ 1.23                    | 4.03                                  |
| DAP                       | 0.41 $\pm$ 0.07                                      | 0.92 $\pm$ 0.09                    | 2.52 $\pm$ 0.80                    | 2.28 $\pm$ 0.96                    | 5.19 $\pm$ 0.52                    | 2.26                                  |
| EASE                      | 0.52 $\pm$ 0.04                                      | 1.41 $\pm$ 0.58                    | 3.48 $\pm$ 1.52                    | 5.40 $\pm$ 0.96                    | 5.07 $\pm$ 0.37                    | 3.18                                  |
| <b>HA<sub>model</sub></b> | 1.25 $\pm$ 0.79                                      | 1.45 $\pm$ 0.14                    | 3.56 $\pm$ 0.25                    | 2.08 $\pm$ 0.37                    | 3.73 $\pm$ 1.17                    | 2.41                                  |
| <b>HA<sub>block</sub></b> | <b>0.40 <math>\pm</math> 0.12</b>                    | <b>0.54 <math>\pm</math> 0.07</b>  | <b>2.08 <math>\pm</math> 0.46</b>  | <b>0.80 <math>\pm</math> 0.44</b>  | <b>2.66 <math>\pm</math> 0.49</b>  | <b>1.30</b>                           |
| <b>Method</b>             | <b>(c) Learning Accuracy (<math>\uparrow</math>)</b> |                                    |                                    |                                    |                                    | <b>Mean (<math>\uparrow</math>)</b>   |
| Full-seq                  | 97.93 $\pm$ 1.84                                     | 94.87 $\pm$ 2.87                   | 89.48 $\pm$ 0.18                   | 97.75 $\pm$ 1.39                   | 94.04 $\pm$ 0.03                   | 94.81                                 |
| Linear-seq                | 96.01 $\pm$ 0.20                                     | 92.58 $\pm$ 0.05                   | 83.20 $\pm$ 1.24                   | 89.50 $\pm$ 1.38                   | 80.20 $\pm$ 0.24                   | 88.12                                 |
| EWC                       | 93.47 $\pm$ 1.98                                     | 86.29 $\pm$ 0.89                   | 74.23 $\pm$ 1.56                   | 81.78 $\pm$ 1.29                   | 75.49 $\pm$ 0.14                   | 82.25                                 |
| LwF                       | 97.14 $\pm$ 0.17                                     | 88.01 $\pm$ 0.08                   | 89.51 $\pm$ 0.57                   | 82.05 $\pm$ 0.07                   | 75.59 $\pm$ 0.21                   | 86.46                                 |
| L2P                       | 97.37 $\pm$ 0.19                                     | 93.54 $\pm$ 0.54                   | 89.78 $\pm$ 0.81                   | 89.13 $\pm$ 0.07                   | 78.65 $\pm$ 0.45                   | 89.69                                 |
| DualPrompt                | 97.72 $\pm$ 0.16                                     | 96.52 $\pm$ 0.41                   | 90.11 $\pm$ 0.23                   | 90.61 $\pm$ 0.13                   | 82.82 $\pm$ 0.03                   | 91.56                                 |
| CODA-P                    | 98.06 $\pm$ 0.24                                     | 96.89 $\pm$ 1.42                   | 90.61 $\pm$ 0.51                   | 91.79 $\pm$ 0.68                   | 87.80 $\pm$ 1.56                   | 93.03                                 |
| DAP                       | 96.74 $\pm$ 0.11                                     | 97.87 $\pm$ 0.28                   | 89.17 $\pm$ 0.48                   | 96.37 $\pm$ 0.74                   | 93.03 $\pm$ 0.58                   | 94.64                                 |
| EASE                      | 98.07 $\pm$ 0.12                                     | 96.67 $\pm$ 1.07                   | 90.36 $\pm$ 0.24                   | 92.60 $\pm$ 1.54                   | 93.67 $\pm$ 1.24                   | 94.27                                 |
| <b>HA<sub>model</sub></b> | 98.24 $\pm$ 0.56                                     | 97.70 $\pm$ 0.21                   | 90.26 $\pm$ 0.07                   | 97.71 $\pm$ 0.06                   | 94.83 $\pm$ 0.02                   | 95.75                                 |
| <b>HA<sub>block</sub></b> | <b>98.32 <math>\pm</math> 0.49</b>                   | <b>98.04 <math>\pm</math> 0.20</b> | <b>91.10 <math>\pm</math> 0.42</b> | <b>97.82 <math>\pm</math> 0.07</b> | <b>94.59 <math>\pm</math> 0.12</b> | <b>95.97</b>                          |

Table 2: Results on Split ImageNet-R and Split DomainNet.

| Dataset<br>Method         | Split ImageNet-R                   |                                   |                                    | Split DomainNet                    |                                   |                                    |
|---------------------------|------------------------------------|-----------------------------------|------------------------------------|------------------------------------|-----------------------------------|------------------------------------|
|                           | FnL. Acc. ( $\uparrow$ )           | Forgetting ( $\downarrow$ )       | Lrn. Acc. ( $\uparrow$ )           | FnL. Acc. ( $\uparrow$ )           | Forgetting ( $\downarrow$ )       | Lrn. Acc. ( $\uparrow$ )           |
| Full-seq                  | 21.09 $\pm$ 3.45                   | 54.89 $\pm$ 2.31                  | 75.96 $\pm$ 1.32                   | 27.89 $\pm$ 3.21                   | 72.89 $\pm$ 2.72                  | 89.58 $\pm$ 1.98                   |
| Linear-seq                | 55.21 $\pm$ 1.59                   | 19.89 $\pm$ 0.45                  | 74.32 $\pm$ 1.45                   | 72.15 $\pm$ 1.98                   | 12.15 $\pm$ 0.89                  | 84.15 $\pm$ 2.72                   |
| L2P                       | 60.98 $\pm$ 0.70                   | 9.93 $\pm$ 0.43                   | 69.23 $\pm$ 0.78                   | 80.67 $\pm$ 0.85                   | 5.33 $\pm$ 0.87                   | 85.14 $\pm$ 0.99                   |
| DualPrompt                | 68.97 $\pm$ 2.87                   | 4.66 $\pm$ 2.15                   | 72.85 $\pm$ 2.27                   | 81.89 $\pm$ 0.63                   | 5.21 $\pm$ 1.17                   | 87.27 $\pm$ 1.80                   |
| CODA-P                    | 75.45 $\pm$ 0.56                   | <b>1.64 <math>\pm</math> 0.10</b> | 77.90 $\pm$ 1.97                   | 76.52 $\pm$ 0.78                   | 6.83 $\pm$ 0.10                   | 82.53 $\pm$ 0.24                   |
| DAP                       | 70.12 $\pm$ 2.24                   | 2.90 $\pm$ 2.70                   | 73.24 $\pm$ 2.81                   | 83.51 $\pm$ 1.07                   | 5.30 $\pm$ 0.52                   | 88.77 $\pm$ 0.79                   |
| EASE                      | 73.23 $\pm$ 1.65                   | 6.79 $\pm$ 1.19                   | 76.97 $\pm$ 1.25                   | 86.76 $\pm$ 1.29                   | 4.67 $\pm$ 0.92                   | 91.74 $\pm$ 1.28                   |
| <b>HA<sub>model</sub></b> | 75.65 $\pm$ 2.67                   | 5.78 $\pm$ 2.94                   | <b>79.22 <math>\pm</math> 0.96</b> | 89.20 $\pm$ 0.54                   | 4.18 $\pm$ 0.49                   | 93.10 $\pm$ 0.24                   |
| <b>HA<sub>block</sub></b> | <b>76.51 <math>\pm</math> 1.32</b> | 3.83 $\pm$ 1.46                   | 77.97 $\pm$ 0.81                   | <b>91.56 <math>\pm</math> 0.11</b> | <b>2.18 <math>\pm</math> 0.10</b> | <b>93.58 <math>\pm</math> 0.12</b> |
| Upper Bound               | 77.13 $\pm$ 1.54                   | -                                 | -                                  | 90.65 $\pm$ 0.98                   | -                                 | -                                  |

Interestingly, the sequential learning method that only updates the classifier (known as linear probing) also reaches an impressive final accuracy of 70.07%, surpassing the classic regularization-based methods EWC and LwF, which fully demonstrates the tremendous potential of pre-trained models.

**ImageNet-R and DomainNet.** To further demonstrate the performance of HyperAdapter, we conduct experiments on two larger benchmarks: Split ImageNet-R and Split DomainNet, with results shown in Table 2. On both datasets, HyperAdapter surpasses the previous best methods, CODA-P and EASE, by 1.06% and 4.80%, respectively. This fully demonstrates the capability of HyperAdapter to handle more categories and longer task sequences. Moreover, HyperAdapter closely approaches the upper bound of multi-task learning on ImageNet-R and even surpasses the upper bound on the larger DomainNet dataset. Experimental results on longer sequences are provided in the appendix.



Table 3: LoRA results on CL-100 benchmark. HLA stands for the LoRA version of HyperAdapter.

| Dataset              | Split Flowers-100                    | Split Caltech-100 | Split Dogs-100   | Split CIFAR-100  | Split Food-100   |                       |
|----------------------|--------------------------------------|-------------------|------------------|------------------|------------------|-----------------------|
| Method               | (a) Final Accuracy ( $\uparrow$ )    |                   |                  |                  |                  | Mean ( $\uparrow$ )   |
| HLA <sub>model</sub> | 95.86 $\pm$ 0.00                     | 93.13 $\pm$ 0.49  | 86.07 $\pm$ 0.08 | 92.20 $\pm$ 0.69 | 87.64 $\pm$ 0.18 | 90.98                 |
| HLA <sub>block</sub> | 96.21 $\pm$ 1.21                     | 97.93 $\pm$ 0.14  | 90.04 $\pm$ 0.13 | 97.80 $\pm$ 0.01 | 93.72 $\pm$ 0.14 | 95.14                 |
| Method               | (b) Forgetting Rate ( $\downarrow$ ) |                   |                  |                  |                  | Mean ( $\downarrow$ ) |
| HLA <sub>model</sub> | 2.24 $\pm$ 0.03                      | 0.45 $\pm$ 1.08   | 5.64 $\pm$ 0.28  | 2.06 $\pm$ 0.93  | 3.73 $\pm$ 0.21  | 2.82                  |
| HLA <sub>block</sub> | 1.18 $\pm$ 0.33                      | 0.54 $\pm$ 0.03   | 2.06 $\pm$ 0.07  | 0.40 $\pm$ 0.00  | 2.14 $\pm$ 0.14  | 1.30                  |
| Method               | (c) Learning Accuracy ( $\uparrow$ ) |                   |                  |                  |                  | Mean ( $\uparrow$ )   |
| HLA <sub>model</sub> | 99.22 $\pm$ 0.01                     | 97.85 $\pm$ 0.08  | 91.03 $\pm$ 0.10 | 97.60 $\pm$ 0.00 | 94.11 $\pm$ 0.03 | 95.96                 |
| HLA <sub>block</sub> | 99.32 $\pm$ 0.38                     | 98.12 $\pm$ 0.06  | 91.72 $\pm$ 0.00 | 98.18 $\pm$ 0.00 | 95.65 $\pm$ 0.02 | 96.60                 |

Table 4: Ablation studies on core designs.

| Dataset<br>Method      | Split CIFAR-100          |                             |                          | Split ImageNet-R         |                             |                          |
|------------------------|--------------------------|-----------------------------|--------------------------|--------------------------|-----------------------------|--------------------------|
|                        | FnL. Acc. ( $\uparrow$ ) | Forgetting ( $\downarrow$ ) | Lrn. Acc. ( $\uparrow$ ) | FnL. Acc. ( $\uparrow$ ) | Forgetting ( $\downarrow$ ) | Lrn. Acc. ( $\uparrow$ ) |
| HA <sub>model</sub>    | 95.85 $\pm$ 0.37         | 2.08 $\pm$ 0.37             | 97.71 $\pm$ 0.06         | 75.65 $\pm$ 2.67         | 5.78 $\pm$ 2.94             | 79.22 $\pm$ 0.96         |
| w/o task dictionary    | 66.32 $\pm$ 1.61         | 33.46 $\pm$ 1.84            | 96.43 $\pm$ 0.06         | 45.48 $\pm$ 5.35         | 35.06 $\pm$ 5.35            | 77.19 $\pm$ 5.58         |
| w/o position embedding | 94.43 $\pm$ 0.11         | 2.30 $\pm$ 0.17             | 96.57 $\pm$ 0.11         | 69.18 $\pm$ 1.20         | 9.09 $\pm$ 0.74             | 76.39 $\pm$ 2.29         |
| w/o adapter LN         | 29.58 $\pm$ 2.73         | 22.84 $\pm$ 1.49            | 50.13 $\pm$ 3.89         | 13.01 $\pm$ 0.46         | 6.80 $\pm$ 0.88             | 18.85 $\pm$ 0.69         |

### 5.3 HYPERADAPTER WITH LORA

To verify that the hypernetwork-based design of our HyperAdapter is not confined to any specific adapter structure, we choose Low-Rank Adaptation or LoRA (Hu et al., 2021), an adapter structure initially proposed for efficiently fine-tuning large language models. We refer to the LoRA version of our approach as HyperLoRA or HLA, which is also divided into model-wise and block-wise versions. The results on CL-100 benchmark are shown in Table 3. In this experiment, we followed the original settings in Hu et al. (2021), applying LoRA to the query and value matrices in ViT’s attention mechanism, with the rank set to half of the bottleneck dimension in HA’s adapter to ensure that the learnable parameters of these two versions are consistent. All other hyperparameters remained the same as in the HA experiments. The results show that our method also achieves excellent performance when applied to LoRA, even outperforming the HA version without hyperparameter tuning, fully demonstrating the versatility of our design.

### 5.4 ABLATION STUDIES

**Core Design.** To validate the effectiveness of the core designs in HyperAdapter, we conduct ablation studies on the Split CIFAR-100 and ImageNet-R datasets, all based on the model-wise version. As shown in Table 4, row 1 represents the complete version of HyperAdapter. In row 2, we remove the task dictionary and use an additional MLP to project queries to the specified dimension, providing instance embeddings to the hypernetwork as input. This MLP is frozen after learning the first task. The results indicate that removing task embeddings has a minor impact on learning accuracy, but overly diverse inputs easily lead to forgetting in the hypernetwork, resulting in  $\sim 30\%$  drop in final accuracy. Row 3 removes the position embedding from each layer, causing the hypernetwork to always generate the same adapter. While repeated adapters have a minimal effect on model performance, due to the minimal parameter count of position embeddings, we ultimately chose the performance improvement brought by diverse adapters. In row 4, we remove the layer normalization (LN) after each adapter. The results demonstrate that this design has a significant impact on the final performance, with removing it resulting in over 60% performance loss. Since learning accuracy also decreases, it indicates that the LN is more for stabilizing the training of the adapters.

**Parameter Scale.** In Figure 3, we show the impact of parameter scale on performance. Figure 3(a) and (b) represent two different versions of HyperAdapter. It can be observed that model performance is positively correlated with parameter scale, which means more parameters lead to higher accuracy. However, as the parameter scale continues to increase, this marginal benefit gradually diminishes. To strike a balance between parameters and performance, we set both the task embedding and adapter bottleneck dimension of model-wise HyperAdapter to 32, while both dimensions of block-wise HyperAdapter are set to 16. Figure 3(c) demonstrates the performance change when fixing one

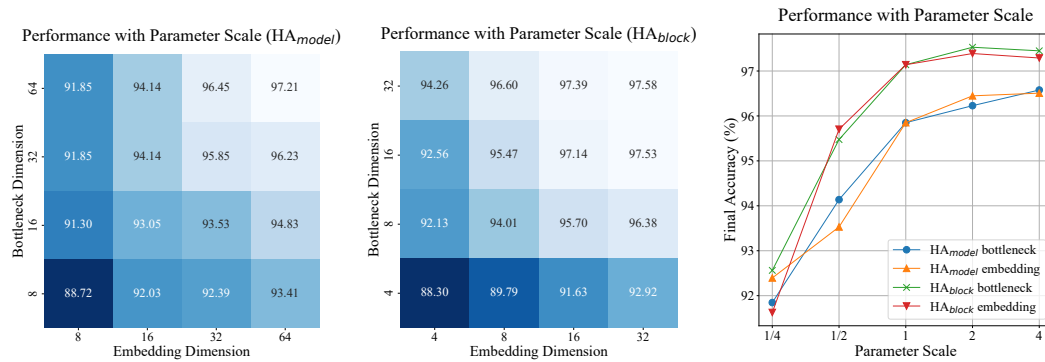


Figure 3: Further analysis on parameter scales on Split CIFAR-100.

dimension and only altering the other one. It can be seen that when the parameter scale reaches a certain threshold, the final performance no longer increases and may even decrease.

## 5.5 VISUALIZATIONS

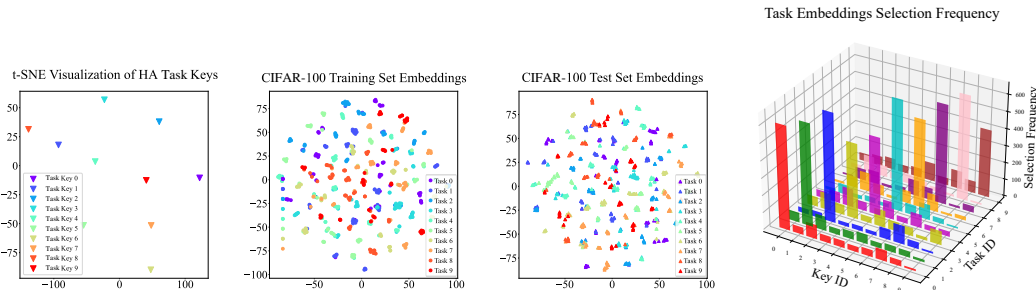


Figure 4: Visualizations on task embeddings and the selection on Split CIFAR-100.

**t-SNE.** In Figure 4(a), we present t-SNE results of input queries and task keys. To enable the hypernetwork to generate corresponding adapters for different tasks, we adopt a query-key matching mechanism to obtain task embeddings, which can be viewed as clustering based on the pre-trained model. Due to the gap between pre-training data and downstream tasks, the pre-trained model may not always distinguish inputs from different classes well. Moreover, since a task comprises data from different classes, the large inter-class gap can make the clustering process more challenging. Nevertheless, this matching mechanism still plays a crucial role in our method, as shown in Table 4.

**Task Embedding Selection.** In Figure 4(b), we further show the instance-level selection frequency of different task embeddings. Most tasks correctly select their corresponding embeddings, and the noise from a few incorrect selections do not significantly impact model performance, benefiting from the similarity-based matching mechanism. Even if the matching is incorrect, the matched task is still highly similar, and the embedding can provide enough information for the current task to achieve correct classification. Improving this selection mechanism is left as a direction for future work.

## 6 CONCLUSION

In this paper, we propose HyperAdapter as a novel rehearsal-free continual learning method, utilizing a hypernetwork to generate adapters to adapt the pre-trained model to different tasks. This work represents the first attempt to employ a hypernetwork for continual learning based on pre-trained models, providing a feasible approach for the continual learning of large-scale pre-trained models. We establish a new SOTA on our newly curated CL-100 benchmark and two commonly used large benchmarks, even surpassing the multi-task learning upper bound in some cases. This highlights the immense potential of pre-trained models in continual learning and marks a significant step forward in the application of neural networks in real-world scenarios.

## REFERENCES

- 540  
541  
542 Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for  
543 online continual learning. *Advances in neural information processing systems*, 32, 2019.
- 544  
545 Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual  
546 learning method based on complementary learning system. *arXiv preprint arXiv:2201.12604*,  
547 2022.
- 548  
549 Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: Bert pre-training of image transformers.  
550 In *International Conference on Learning Representations*, 2021.
- 551  
552 Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101 – mining discriminative compo-  
553 nents with random forests. In *European Conference on Computer Vision*, 2014.
- 554  
555 Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark  
556 experience for general continual learning: a strong, simple baseline. *Advances in neural information  
557 processing systems*, 33:15920–15930, 2020.
- 558  
559 Hyuntak Cha, Jaeho Lee, and Jinwoo Shin. Co2l: Contrastive continual learning. In *Proceedings of  
560 the IEEE/CVF International conference on computer vision*, pp. 9516–9525, 2021.
- 561  
562 Dupati Srikar Chandra, Sakshi Varshney, PK Srijith, and Sunil Gupta. Continual learning with  
563 dependency preserving hypernetworks. In *Proceedings of the IEEE/CVF Winter Conference on  
564 Applications of Computer Vision*, pp. 2339–2348, 2023.
- 565  
566 Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian  
567 walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the  
568 European conference on computer vision (ECCV)*, pp. 532–547, 2018a.
- 569  
570 Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient  
571 lifelong learning with a-gem. In *International Conference on Learning Representations*, 2018b.
- 572  
573 Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K  
574 Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual  
575 learning. *arXiv preprint arXiv:1902.10486*, 2019.
- 576  
577 Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu Wang, Yibing Song, Jue Wang, and Ping Luo.  
578 Adaptformer: Adapting vision transformers for scalable visual recognition. *Advances in Neural  
579 Information Processing Systems*, 35:16664–16678, 2022.
- 580  
581 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
582 contrastive learning of visual representations. In *International conference on machine learning*, pp.  
583 1597–1607. PMLR, 2020.
- 584  
585 E Dataset. Novel datasets for fine-grained image categorization. In *First Workshop on Fine Grained  
586 Visual Categorization, CVPR. Citeseer. Citeseer. Citeseer*, 2011.
- 587  
588 Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory  
589 Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification  
590 tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021.
- 591  
592 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale  
593 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,  
pp. 248–255. Ieee, 2009.
- 594  
595 Li Fei-Fei, Robert Fergus, and Pietro Perona. One-shot learning of object categories. *IEEE transac-  
596 tions on pattern analysis and machine intelligence*, 28(4):594–611, 2006.
- 597  
598 Xinyuan Gao, Songlin Dong, Yuhang He, Qiang Wang, and Yihong Gong. Beyond prompt learning:  
599 Continual adapter for efficient rehearsal-free continual learning. *arXiv preprint arXiv:2407.10281*,  
2024.

- 594 Alex Gomez-Villa, Bartłomiej Twardowski, Kai Wang, and Joost Van de Weijer. Plasticity-optimized  
595 complementary networks for unsupervised continual learning. In *Proceedings of the IEEE/CVF*  
596 *Winter Conference on Applications of Computer Vision*, pp. 1690–1700, 2024.
- 597 Mustafa Burak Gurbuz, Jean Michael Moorman, and Constantine Dovrolis. Nice: Neurogenesis  
598 inspired contextual encoding for replay-free class incremental learning. In *Proceedings of the*  
599 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23659–23669, 2024.
- 600 David Ha, Andrew Dai, and Quoc V Le. Hypernetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- 602 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
603 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*  
604 *computer vision and pattern recognition*, pp. 9729–9738, 2020.
- 605 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked  
606 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer*  
607 *vision and pattern recognition*, pp. 16000–16009, 2022a.
- 609 Yun He, Steven Zheng, Yi Tay, Jai Gupta, Yu Du, Vamsi Aribandi, Zhe Zhao, YaGuang Li, Zhao  
610 Chen, Donald Metzler, et al. Hyperprompt: Prompt-based task-conditioning of transformers. In  
611 *International conference on machine learning*, pp. 8678–8690. PMLR, 2022b.
- 612 Hamed Hemati, Vincenzo Lomonaco, Davide Bacciu, and Damian Borth. Partial hypernetworks for  
613 continual learning. In *Conference on Lifelong Learning Agents*, pp. 318–336. PMLR, 2023.
- 614 Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul  
615 Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer.  
616 The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*,  
617 2021.
- 618 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,  
619 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for  
620 nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- 621 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,  
622 et al. Lora: Low-rank adaptation of large language models. In *International Conference on*  
623 *Learning Representations*, 2021.
- 624 Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and  
625 Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision*, 2022.
- 626 Dahuin Jung, Dongyoon Han, Jihwan Bang, and Hwanjun Song. Generating instance-level prompts  
627 for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF International Conference*  
628 *on Computer Vision*, pp. 11847–11857, 2023.
- 629 Ronald Kemker and Christopher Kanan. Fearnnet: Brain-inspired model for incremental learning.  
630 *arXiv preprint arXiv:1711.10563*, 2017.
- 631 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A  
632 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming  
633 catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114  
634 (13):3521–3526, 2017.
- 635 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 636 Kamil Książek and Przemysław Spurek. Hypermask: Adaptive hypernetwork-based masks for  
637 continual learning. *arXiv preprint arXiv:2310.00113*, 2023.
- 638 Dharshan Kumaran, Demis Hassabis, and James L McClelland. What learning systems do intelligent  
639 agents need? complementary learning systems theory updated. *Trends in cognitive sciences*, 20(7):  
640 512–534, 2016.
- 641 Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis*  
642 *and machine intelligence*, 40(12):2935–2947, 2017.

- 648 David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning.  
649 *Advances in neural information processing systems*, 30, 2017.
- 650
- 651 Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-  
652 efficient multi-task fine-tuning for transformers via shared hypernetworks. *arXiv preprint*  
653 *arXiv:2106.04489*, 2021.
- 654 Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online  
655 continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51,  
656 2022.
- 657
- 658 James L McClelland, Bruce L McNaughton, and Randall C O’Reilly. Why there are complementary  
659 learning systems in the hippocampus and neocortex: insights from the successes and failures of  
660 connectionist models of learning and memory. *Psychological review*, 102(3):419, 1995.
- 661
- 662 Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The  
663 sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165.  
664 Elsevier, 1989.
- 665
- 666 Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number  
667 of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pp.  
668 722–729. IEEE, 2008.
- 669
- 670 German I Parisi, Jun Tani, Cornelius Weber, and Stefan Wermter. Lifelong learning of spatiotemporal  
671 representations with dual-memory recurrent self-organization. *Frontiers in neurorobotics*, 12:78,  
672 2018.
- 673
- 674 Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching  
675 for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on*  
676 *Computer Vision*, pp. 1406–1415, 2019.
- 677
- 678 Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfu-  
679 sion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference*  
680 *of the European Chapter of the Association for Computational Linguistics: Main Volume*, pp.  
681 487–503, 2021.
- 682
- 683 Quang Pham, Chenghao Liu, and Steven Hoi. Dualnet: Continual learning, fast and slow. *Advances*  
684 *in Neural Information Processing Systems*, 34:16131–16144, 2021.
- 685
- 686 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:  
687 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on*  
688 *Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- 689
- 690 Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald  
691 Tesauero. Learning to learn without forgetting by maximizing transfer and minimizing interference.  
692 In *International Conference on Learning Representations*. International Conference on Learning  
693 Representations, ICLR, 2019.
- 694
- 695 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience  
696 replay for continual learning. *Advances in neural information processing systems*, 32, 2019.
- 697
- 698 Marcin Sendera, Marcin Przewieźlikowski, Konrad Karanowski, Maciej Zieba, Jacek Tabor, and  
699 Przemysław Spurek. Hypershot: Few-shot learning by kernel hypernetworks. In *Proceedings of*  
700 *the IEEE/CVF winter conference on applications of computer vision*, pp. 2469–2478, 2023.
- 701
- 702 Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd*  
703 *ACM SIGSAC conference on computer and communications security*, pp. 1310–1321, 2015.
- 704
- 705 James Smith, Jonathan Balloch, Yen-Chang Hsu, and Zsolt Kira. Memory-efficient semi-supervised  
706 continual learning: The world is its own replay buffer. In *2021 International Joint Conference on*  
707 *Neural Networks (IJCNN)*, pp. 1–8. IEEE, 2021.

- 702 James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf  
703 Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed  
704 attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF*  
705 *Conference on Computer Vision and Pattern Recognition*, pp. 11909–11919, 2023.
- 706 Ahmet Üstün, Arianna Bisazza, Gosse Bouma, Gertjan van Noord, and Sebastian Ruder. Hyper-x: A  
707 unified hypernetwork for multi-task multilingual transfer. In *Proceedings of the 2022 Conference*  
708 *on Empirical Methods in Natural Language Processing*, pp. 7934–7949, 2022.
- 709 Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren,  
710 Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for  
711 rehearsal-free continual learning. In *European Conference on Computer Vision*, pp. 631–648.  
712 Springer, 2022a.
- 713 Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent  
714 Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings*  
715 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 139–149, 2022b.
- 716 Zhen-da Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu.  
717 Simmim: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF*  
718 *conference on computer vision and pattern recognition*, pp. 9653–9663, 2022.
- 719 Rongyu Zhang, Yun Chen, Chenrui Wu, Fangxin Wang, and Jiangchuan Liu. Optimizing efficient  
720 personalized federated learning with hypernetworks at edge. *IEEE Network*, 37(4):120–126, 2023.
- 721 Zhengkun Zhang, Wenya Guo, Xiaojun Meng, Yasheng Wang, Yadao Wang, Xin Jiang, Qun Liu, and  
722 Zhenglu Yang. Hyperpelt: Unified parameter-efficient language model tuning for both language  
723 and vision-and-language tasks. *arXiv preprint arXiv:2203.03878*, 2022.
- 724 Andrey Zhmoginov, Mark Sandler, and Maksym Vladymyrov. Hypertransformer: Model generation  
725 for supervised and semi-supervised few-shot learning. In *International Conference on Machine*  
726 *Learning*, pp. 27075–27098. PMLR, 2022.
- 727 Da-Wei Zhou, Hai-Long Sun, Han-Jia Ye, and De-Chuan Zhan. Expandable subspace ensemble for  
728 pre-trained model-based class-incremental learning. *arXiv preprint arXiv:2403.12030*, 2024.
- 729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755