

# SELF-ORGANIZING PATHWAY EXPANSION FOR NON-EXEMPLAR INCREMENTAL LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Non-exemplar class-incremental learning aims to recognize both the old and new classes without access to old class samples. The conflict between old and new class optimization is exacerbated since the shared neural pathways can only be differentiated by the incremental samples. To address this problem, we propose a novel self-organizing pathway expansion scheme. Our scheme consists of a class-specific pathway organization strategy that decouples the optimization pathway of different classes to enhance the independence of the feature representation, and a pathway-guided feature optimization mechanism to mitigate the update interference between the old and new classes. Extensive experiments on four datasets demonstrate superior incremental performance, outperforming the state-of-the-art methods by a margin of 1%, 3%, 2% and 2%, respectively.

## 1 INTRODUCTION

Since deep neural networks have achieved good performance in fully supervised scenarios, how to extend this learning capability to open environment has attracted great attention. Particularly, it is essential to ensure that the network can continuously learn new knowledge while maintaining the abilities to identify old tasks (*i.e.*, incremental learning (Rebuffi et al., 2017; Douillard et al., 2020)). Fine-tuning the network directly with new data can lead to a serious bias of the representation and classifier, which is often referred to as catastrophic forgetting. Due to privacy and hardware limits, old samples are usually unavailable for joint training, making it more difficult to maintain the old class performance in the subsequent optimization process. In this paper, we focus on this ability to continuously learn new tasks without any old samples or exemplars, which is called non-exemplar class-incremental learning (NECIL) (Zhu et al., 2021b;a; 2022; Yu et al., 2020b; Yin et al., 2020).

Most methods maintain the feature representation of old classes by means of various distillation loss functions (Douillard et al., 2020; Hu et al., 2021). Although catastrophic forgetting is somewhat mitigated, incremental performance still suffers from the confusion between the old and new class in the feature space. Furthermore, in the absence of old class samples, the degree of forgetting is only related to the initial model and incremental samples (Zhu et al., 2021b). Existing NECIL works (Zhu et al., 2021a; Yin et al., 2020) mainly focus on enhancing the overall performance by improving the discrimination and generalization of the initial model, which brings a significant improvement on the incremental performance.

Instead, we focus on the impact of incremental samples on the optimization process. Intuitively, since different incremental classes cause disparate feature confusion, the interference on the old class performance is also different even if initialized from the same model (Zhu et al., 2022; 2021b). To further explore the association, we estimate the inter-class confusion by measuring the status of feature activation (Zhou et al., 2016) in existing incremental model. As shown in Fig. 1 (b), we filter out the positions of strongly activated modules as the class-specific pathways, and find that the pathway of incremental class is commonly confused with the previous ones in the baseline. Furthermore, it can be seen in Fig. 1 (a) that the degree of pathway overlap (*i.e.*, similarity) between the old class and incremental class is positively correlated with the forgetting degree, which motivates us to address the interference problem from the perspective of pathway optimization.

Based on the above observation, we propose a self-organizing pathway expansion scheme to learn a pathway-aware representation, mitigating the feature interference during the subsequent incremental process. The scheme is mainly manifested in two aspect. Firstly, during the initial phase, we adopt

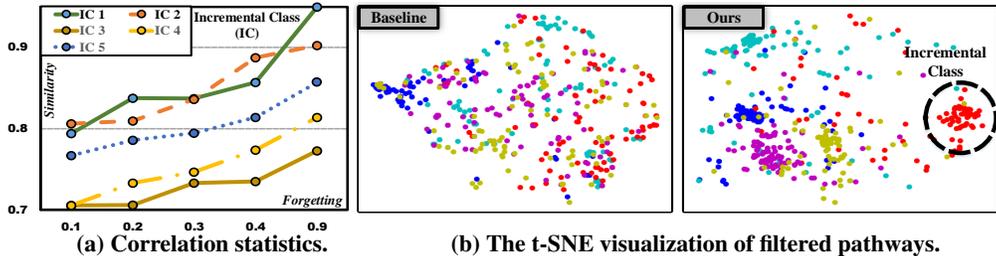


Figure 1: Motivation of our method. (a) The accuracy degradation of old classes (*i.e.* forgetting rates in the horizontal coordinate) is positively correlated to the corresponding pathway similarity with incremental classes. The concept of pathway is formed from the aggregation of important modules, which are filtered out by the contribution to the final recognition performance. (b) Compared to the standard classification method (*i.e.*, baseline in Sec. 3.2), the discriminative pathway in our method brings out lower inter-class overlap, which benefits the mitigation of feature confusion.

the class-specific pathway organization strategy to enhance the independence of feature representation by forcing the optimization pathways specific to different classes. A global pathway planner is utilized to explicitly select the most relevant modules, facilitating the pathway identification. It is noted that we do not modify the network structure, but only divide the output channels of each convolution module to match the output of the pathway planner. Secondly, during the incremental phases, we introduce a pathway-guided feature update mechanism to promote the effectiveness of new classes involved in incremental optimization by adjusting the classification weight with the pathway similarity. Since the pathway value is either 0 or 1, we calculate the intersection of union (*i.e.*, IoU) value to better measure the class relevance, reducing the interference of vector normalization. Furthermore, an incremental pathway update mechanism is proposed to ensure the long-term effect by alternating the optimization of the pathway planner and feature representation. To summarize, our main contributions are as follows:

- 1) A self-organizing pathway expansion scheme is proposed for non-exemplar incremental learning, in which a progressive decoupling optimization is accomplished by a class-specific pathway organization strategy, resulting in a pathway-aware representation.
- 2) A pathway-guided feature update mechanism is proposed, which utilizes the similarity of pathways to guide the optimization of incremental samples.
- 3) Extensive experiments are performed on benchmark including CIFAR-100, TinyImageNet, ImageNet-Subset and ImageNet-Full datasets, and the results demonstrate the superiority of our method over the state-of-the-art.

## 2 RELATED WORK

### 2.1 INCREMENTAL LEARNING

As deep learning research advances, there is a growing demand for continual learning (Kirkpatrick et al., 2017; Zenke et al., 2017; Aljundi et al., 2018), which requires the network to learn new tasks without forgetting the old knowledge to achieve the stability-plasticity trade-off. Class-incremental learning (CIL (Rebuffi et al., 2017; Wu et al., 2019; Hou et al., 2019; Douillard et al., 2020; Yan et al., 2021)), a difficult type in continual learning, has attracted much attention due to the agnosticism to task identity (van de Ven & Tolias, 2019).

Recently, some works (Yu et al., 2020b; Zhu et al., 2021b;a; Yin et al., 2020) focus on a challenging but practical non-exemplar class-incremental learning (NECIL) problem, where no past data can be stored due to equipment limits or privacy security. Yu et al. (2020b) estimates the semantic drift of the initial model inherited from the base phase, and compensates the prototypes in each test phase. Yin et al. (2020) inverts the old samples from the initial model for the joint distillation process. Zhu et al. (2021b;a) consider to enhance the generalization of the representation to learn more transferable features for future tasks. We follow their NECIL settings. However, different from their work focusing on the utilization and enhancement of the initial model, we mainly consider the rectification of the incremental samples on joint classification and distillation process.

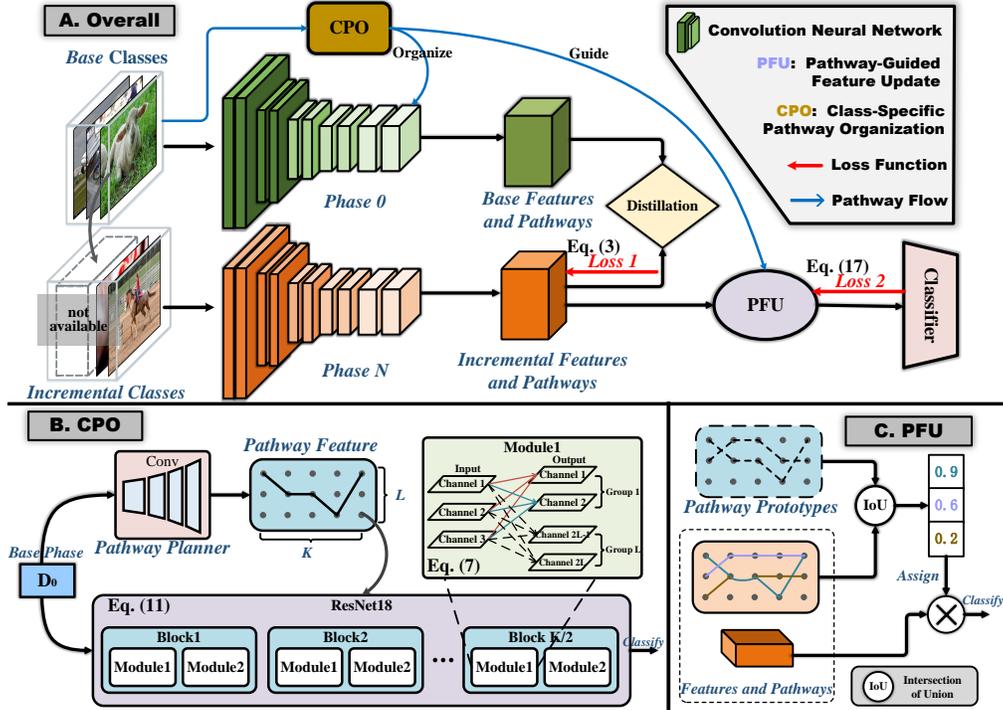


Figure 2: Our proposed self-organized pathway expansion scheme for NECIL. (A) Overall pipeline. (B) During the base phase, a CPO strategy is proposed to mitigate the incremental interference, in which the pathway feature extracted by a pathway planner is utilized to organize the class-specific learning in feature extractor (e.g., ResNet18). (C) During the incremental phase, the similarity scores (i.e. IoU) between the pathway feature and saved pathway prototypes are assigned to the optimization process as loss weights, facilitating the pathway-guided feature update (PFU).

## 2.2 NEURAL PATHWAYS

Recently, large language models have been scaled up with pipelining rather than pure data-parallelism (Zhang et al., 2022), demonstrating the potential of pathways. To enhance the adaptation of the network to new tasks, several continual learning methods (Chen et al., 2020; Rajasegaran et al., 2019) have been proposed to decouple the learning process from the perspective of pathway. However, the targeted models are continuously expanded with the update of pathway, which is difficult to adapt to the standard classification network (e.g., ResNet (He et al., 2016)). The expansion direction of pathway tends to be selected randomly, making it hard to search for an explanation. In this paper, we target on the pathway learning on the standard network without changing the structure, and guiding the incremental optimization based on the pathway relationship.

## 3 METHODOLOGY

### 3.1 PROBLEM DESCRIPTION

The NECIL problem is defined as follows. Here we denote  $D_t$  as the training set at the current phase  $t$ , which consists of the sample set  $X_t$  and label set  $Y_t$ . Our task is to train the model from a continuous data stream, i.e., training sets  $D_0, D_1, \dots, D_T$ , where labels of a set  $X_i$  ( $0 \leq i \leq T$ ) are from the set  $Y_i$ , and  $T$  represents the number of incremental phases. It should be mentioned that all the incremental classes are disjoint, that is,  $Y_i \cap Y_j = \emptyset (i \neq j)$ . At the current phase  $t$ , there are no old training sets (i.e.,  $D_{0:t-1}$ ) in memory, but incremental samples (i.e.,  $D_t$ ) for the current phase. To measure the performance of models at current phase  $t$ , we calculate the classification accuracy on the test set  $Z_t$ , in which the classes are from all the seen label sets  $Y_0 \cup Y_1 \dots \cup Y_t$ .

### 3.2 BASELINE FOR NECIL

Following the paradigm of existing NECIL works (Zhu et al., 2021b;a; 2022; Yin et al., 2020), we adapt the distillation-based CIL methods (Rebuffi et al., 2017) to the NECIL setting as the baseline. Specifically, at the incremental phase (*i.e.*,  $t > 0$ ), a standard classification model that consists of the feature extractor  $f_{\theta_t}$  and classifier  $g_{\phi_t}$  should be optimized under the full supervision (*i.e.*,  $D_{0:t}$ ),

$$\min_{\theta_t, \phi_t} \mathcal{L}_t = \mathcal{L}_{cls}(\theta_t, \phi_t; D_{0:t}) = \mathcal{L}_{cls}(\theta_t, \phi_t; D_{0:t-1}) + \mathcal{L}_{cls}(\theta_t, \phi_t; D_t), \quad (1)$$

$$\mathcal{L}_{cls}(\theta_t, \phi_t; D_t) = \sum_{x \in X_t} \sum_{y \in Y_t} y \cdot \log(g_{\phi_t}(f_{\theta_t}(x))), \quad (2)$$

where  $\mathcal{L}_t$  represents the overall loss function for feature optimization. However in the NECIL setting, since the previous training sets are unavailable, the corresponding loss  $\mathcal{L}_{cls}(\theta_t, \phi_t; D_{0:t-1})$  for both the feature extractor and classifier is missing, leading to a serious bias to current classes. To solve the problem, existing methods (Hou et al., 2019; Douillard et al., 2020) replace the old classification supervision with the feature distillation and classifier correction. Specifically, the parameters  $\theta_{t-1}$  of the old feature extractor from previous phase  $t-1$  is frozen and saved during each incremental phase  $t$ . To maintain the old informative feature, the knowledge distillation  $\mathcal{L}_{kd}$  is used to ensure the similarity between the current representation  $f_{\theta_t}(x)$  and the previous one  $f_{\theta_{t-1}}(x)$ :

$$\min_{\theta_t} \mathcal{L}_{kd}(\theta_t; \theta_{t-1}, D_t) = \sum_{x \in X_t} \|f_{\theta_t}(x) - f_{\theta_{t-1}}(x)\|_2, \quad (3)$$

where  $\|\cdot\|_2$  denotes Euclidean Norm. As there are no exemplars for balanced classifier optimization in NECIL, we turn to consider the class-representative prototypes  $P_{0:t-1}$  (Zhu et al., 2021b) in the deep feature space. Specifically, we compute and memorize one prototype  $p^c \in P_{0:t-1}$  for each class  $c$  as:

$$p^c = \mathbb{E}_{(x,y) \sim D_{0:t-1}} [f_{\theta_t}(x) \mid y = c]. \quad (4)$$

In each training iteration, we choose to oversample (Chawla et al., 2002) memorized prototypes  $P_{0:t-1}$  as training prototypes  $\tilde{P}_{0:t-1}$  by the ratio of batch size. Training prototypes are directly involved in the standard classification optimization, achieving the augmentation of the classifier, which is consistent with the baseline in PASS (Zhu et al., 2021b) and IL2A (Zhu et al., 2021a):

$$\min_{\phi_t} \mathcal{L}_{aug}(\phi_t; \tilde{P}_{0:t-1}) = \sum_{p^c \in \tilde{P}_{0:t-1}} \sum_{y \in Y_{0:t-1}} y \cdot \log(g_{\phi_t}(p^c)). \quad (5)$$

In conclusion, the overall feature optimization problem for the baseline method can be written as follows,

$$\min_{\theta_t, \phi_t} \mathcal{L}_t = \mathcal{L}_{cls}(\theta_t, \phi_t; D_t) + \mathcal{L}_{kd}(\theta_t; \theta_{t-1}, D_t) + \mathcal{L}_{aug}(\phi_t; \tilde{P}_{0:t-1}). \quad (6)$$

### 3.3 SELF-ORGANIZING PATHWAY EXPANSION

Our proposed self-organizing pathway expansion scheme consists of a class-specific pathway organization strategy that reduces the pathway overlap during the base phase to mitigate the overall feature confusion, and a pathway-guided feature optimization mechanism to refine the incremental optimization guided by the inter-class pathway correlation. The main procedures are summarized in Algorithms 1 and 2 respectively, and the specific implementation is described below.

**Class-Specific Pathway Organization.** To mitigate the interference during the feature optimization process, we perform a structural decomposition on the feature extractor and organize the class-specific pathway adaptively. As shown in Fig. 2, each standard convolution module consists of a  $3 \times 3$  convolution layer and a BatchNorm layer. We firstly reorganize  $K$  convolution modules, each of which is equally divided into  $L$  groups along the output channels. We define  $\theta_t^k \in \mathbb{R}^{C_{in} \times C_{out}}$  as the parameters of  $k_{th}$  convolution module  $\text{Conv}_{\theta_t^k}$  of the feature extractor  $f_{\theta_t}$ , in which  $\theta_t^{k,l} \in \mathbb{R}^{C_{in} \times C_{out}/L}$  denotes the parameters of  $l_{th}$  group  $\text{Conv}_{\theta_t^{k,l}}$ .  $C_{in}$  and  $C_{out}$  represent the number of input and output channels. Let  $z_t^{k-1}$  be the input feature of  $\text{Conv}_{\theta_t^k}$ , the convolution operation is organized as follows,

$$z_t^k = \text{Conv}_{\theta_t^k}(z_t^{k-1}) = \text{Concat}[\text{Conv}_{\theta_t^{k,1}}(z_t^{k-1}) \dots \text{Conv}_{\theta_t^{k,L}}(z_t^{k-1})], (1 < k \leq K), \quad (7)$$

**Algorithm 1** Class-Specific Pathway Organization

- 1: **Input:** Feature extractor  $f_{\theta_t}$ , pathway planner  $f_{\alpha_t}$ , base set  $D_0$  and maximum sparse rate  $\zeta_{max}$ ,
- 2: **Initialize:** Reorganize the structure  $f_{\theta_t}$  by Eq. (7);
- 3: **for all**  $(x, y) \in D_0$  **do**
- 4:   Extract the pathway score  $\mathbf{S} = f_{\alpha_t}(x)$ ;
- 5:   Compute the specific sparse rate  $\zeta (\leq \zeta_{max})$  in each epoch by Eq. (10);
- 6:   Confirm the position  $(l, k)$  of filtered pathway with the soft threshold  $\varepsilon$  by Eq. (8);  
 $\hat{\mathbf{S}} \leftarrow \{s^{l,k} \mid s^{l,k} > \varepsilon, s^{l,k} \in \mathbf{S}\}$
- 7:   Guide the forward feature optimization with filtered pathway by Eq. (11);
- 8:   Update  $\theta_t$  and  $\alpha_t$  by taking a SGD step on the image and pathway loss (Eqs. (13) and (14));
- 9: **end for**
- 10: **Output:** Calculated feature prototypes  $\mathbf{P}_0$  and pathway prototypes  $\mathbf{A}_0$  by Eqs. (4) and (15).

**Algorithm 2** Pathway-Guided Feature Update

- 1: **Input:** Old  $f_{\theta_{t-1}}$  and new feature extractor  $f_{\theta_t}$ , old  $f_{\alpha_{t-1}}$  and new pathway planner  $f_{\alpha_t}$ , incremental set  $D_t (t > 0)$ , feature prototypes  $\mathbf{P}_{0:t-1}$  and pathway prototypes  $\mathbf{A}_{0:t-1}$ .
- 2: **Initialize:** Freeze the parameters of  $f_{\alpha_t}$ ;
- 3: **for all**  $(x, y) \in D_t$  **do**
- 4:   Filter the pathway with  $\zeta_{max}$  by Eq. (9);
- 5:   Compute feature classification and distillation loss weighted with pathway similarity by Eq. (17);
- 6:   Compute the augmentation loss by Eq. (5);
- 7:   Update  $\theta_t$  and  $\alpha_t$  based on above losses;
- 8: **end for**
- 9: Freeze the parameters of  $f_{\theta_t}$ , and unfreeze  $f_{\alpha_t}$ ;
- 10: **for all**  $(x, y) \in D_t$  **do**
- 11:   Update incremental pathway planner with pathway update loss  $\mathcal{L}_t^{path}$  by Eq. (18)
- 12: **end for**

where Concat denotes the concatenation along the output channels. The output feature  $\mathbf{z}_t^k$  is the same as the that of standard convolution module before reorganization.

Then, we introduce a pathway planner  $f_{\alpha_t}$ , which consists of several standard convolution blocks. It receives the image  $x$  as input, and output a probability score  $\mathbf{S} \in \mathbb{R}^{K \times L} = f_{\alpha_t}$ , representing the pathway importance of  $K$  modules and  $L$  groups in the feature extractor. According to the obtained score, a gradually decreasing sparse rate is adopted to filter the most adequate components of the global pathway to guide the feature optimization. Specifically, given a target sparse rate  $\zeta$ , we solve the minimum pathway threshold  $\varepsilon$  from the equation,

$$1 - \zeta = \frac{|\{s^{k,l} \mid s^{k,l} > \varepsilon, s^{k,l} \in \mathbf{S}\}|}{|\{s^{k,l}, s^{k,l} \in \mathbf{S}\}|}, \quad (8)$$

where  $|\cdot|$  means the element number. The pathway score can be filtered by the calculated threshold:

$$\hat{\mathbf{S}} = \text{Filter}(\mathbf{S}, \zeta) = \mathbf{S} * \text{Bool}(\mathbf{S} - \varepsilon > 0), \quad (9)$$

where  $*$  represents the element-wise multiplication, and Bool denotes the element-wise boolean operation. As the threshold  $\varepsilon$  is not a given hard value (Csordás et al., 2020) but a filtered soft one in Eq. (8), no special gradient correction is required. To stabilize the optimization process with the threshold, we use a three-step strategy to jointly optimize features and pathways in which different values of sparse rate are adopted at different epoch  $e$ :

$$\zeta = \begin{cases} 0, & e < e_1 \\ \frac{e-e_1}{e_2-e_1} \zeta_{max}, & e_1 \leq e < e_2 \\ \zeta_{max}, & e \geq e_2, \end{cases} \quad (10)$$

where  $e_1$  and  $e_2$  are two hyper-parameters.  $\zeta_{max}$  is another hyper-parameter that defines the maximum value of sparse rate. According to the filtered scores  $\hat{\mathbf{S}}$ , we reorganize the pathway of the network, and Eq. (7) can be rewritten as follows,

$$\mathbf{z}_t^k = \text{Conv}_{\theta_t^k}(\mathbf{z}_t^{k-1}, \hat{\mathbf{S}}) = \text{Concat}[\hat{s}^{k,1} * \text{Conv}_{\theta_t^{k,0}}(\mathbf{z}_t^{k-1}) \dots \hat{s}^{k,L} * \text{Conv}_{\theta_t^{k,L}}(\mathbf{z}_t^{k-1})], \quad (11)$$

$$\mathbf{z}_t^K = f_{\theta_t}(\mathbf{z}_t^0, \hat{\mathbf{S}}) = f_{\theta_t}(x; \hat{\mathbf{S}}) = f_{\theta_t}(x; f_{\alpha_t}(x)) = f_{\theta_t, \alpha_t}(x), \quad x \in X_t, \quad (12)$$

where  $\hat{s}^{k,l}$  denotes the element in  $\hat{\mathbf{S}}$  at the  $(k, l)$  position. Eq. (2) can be rewritten as follows,

$$\mathcal{L}_{cls}(\theta_t, \phi_t, \alpha_t; D_t) = \sum_{x \in X_t} \sum_{y \in Y_t} y \cdot \log(g_{\phi_t}(f_{\theta_t, \alpha_t}(x))). \quad (13)$$

Finally, we binarize the filtered pathway and improve inter-class discriminability with a learnable pathway classifier  $g_{\beta_t}$ :

$$\min_{\alpha_t, \beta_t} \mathcal{L}_{cls}^{path}(\alpha_t, \beta_t; D_t) = \sum_{x \in X_t} \sum_{y \in Y_t} y \cdot \log(g_{\beta_t}(\delta(f_{\alpha_t}(x)))), \quad (14)$$

where  $\delta$ ,  $\alpha_t$  and  $\mathcal{L}_{cls}^{path}$  denotes the gate function (Serra et al., 2018), the learnable parameters in the pathway planner  $f_{\alpha_t}$  and the overall pathway classification loss, respectively.

**Pathway-Guided Feature Update.** To promote the efficiency of incremental learning, we adopt a pathway-guided feature update mechanism in the incremental phase. Specifically, we involve new samples into the classification process according to the pathway overlap with old ones. We preserve the class-specific pathway prototype  $\mathbf{a}^c \in \mathbf{A}_{0:t-1}$  for class  $c$  at the phase end,

$$\mathbf{a}^c = \text{Filter}(\mathbb{E}_{(x,y) \sim D_{0:t-1}} [f_{\alpha_t}(x) \mid y = c], \zeta), \quad (15)$$

where Filter is the same as that in Eq. (9). The binarized pathway score  $\delta(f_{\alpha_t}(x))$  is compared to the saved pathway prototype with intersection over union (IoU (Long et al., 2015)), thus measuring the relevance  $\lambda$  of the corresponding samples to the previous parameter space:

$$\lambda(x) = \frac{1}{C} \sum_{c=1}^C (\text{IoU}(\delta(f_{\alpha_t}(x)), \mathbf{a}^c)), \quad (16)$$

where  $C$  represents the number of pathway prototypes. To ensure the stability of the incremental representation optimization, we freeze the parameters of the pathway planner (*i.e.*,  $\alpha_t$ ). More relevant samples are assigned smaller weights to optimize the novel classes, thus the classification loss in Eq. (6) (*i.e.*,  $\mathcal{L}_t$ ) can be rewritten as follows,

$$\min_{\theta_t, \phi_t} \mathcal{L}_{cls}(\theta_t, \phi_t; \alpha_t, D_t, \mathbf{A}_{0:t-1}) = \sum_{x \in X_t} \lambda(x) \sum_{y \in Y_t} y \cdot \log(g_{\phi_t}(f_{\theta_t, \alpha_t}(x))). \quad (17)$$

**Incremental Pathway Update.** To enhance the effectiveness of the pathway planner, we then freeze the parameters of feature extractor  $\theta_t$ , and adopt the incremental pathway update mechanism, which is similar to the optimization process of incremental feature in Eq. (6),

$$\min_{\alpha_t, \beta_t} L_t^{path} = \mathcal{L}_{cls}^{path}(\alpha_t, \beta_t; D_t) + \mathcal{L}_{kd}^{path}(\alpha_t; \alpha_{t-1}, D_t) + \mathcal{L}_{aug}^{path}(\beta_t, \mathbf{A}_{0:t-1}). \quad (18)$$

The old pathway planner with frozen parameters  $\alpha_{t-1}$  is utilized to distill with the current planner, and the pathway prototypes are oversampled to correct the pathway classifier bias to the old class. Overall, the loss functions  $L_t$  and  $L_t^{path}$  are utilized sequentially in the incremental phase  $t$ .

## 4 EXPERIMENTS

### 4.1 DATASETS AND SETTINGS

**Datasets.** Following the setting in Zhu et al. (2021b), we conduct comprehensive experiments on four datasets CIFAR-100 (Krizhevsky, 2009), TinyImageNet (Le & Yang, 2015), ImageNet-Subset and ImageNet-Full. CIFAR-100 contains 60,000 images of  $32 \times 32$  size from 100 classes, and each class includes 500 training images and 100 test images. TinyImageNet contains 200 classes, and each class contains 500 training images, 50 validation images and 50 test images. It provides more incremental phases and classes for the sensitivity analysis on different methods. ImageNet-Subset is a 100-class subset of ImageNet-Full (Deng et al., 2009), which provides a large-scale evaluation scenery. Except for 40 base classes in 20 incremental phases setting of CIFAR-100, we train the model on half of classes for the base phase, and equal classes in the rest incremental phases. We conduct different incremental settings (5, 10 and 20 phases) for both CIFAR-100 and TinyImageNet, and 10 incremental phases setting for the rest datasets, which is consistent with Zhu et al. (2021b).

**Settings and Metric.** For a fair comparison with (Zhu et al., 2021b), we adopt the same backbone network (*i.e.*, ResNet-18), and maintain the same accuracy at the first phase for all datasets. We report average incremental accuracy and average forgetting (Zhu et al., 2021b). Average incremental accuracy  $A_T$  is computed as the average accuracy of all incremental phases  $a_t$  (including the first phase), which compares the overall performance of different methods fairly,

$$A_T = \frac{1}{T} \sum_{t=0}^T a_t. \quad (19)$$

Average forgetting is computed as the average forgetting throughout the incremental process, which directly measures the ability of different methods to resist catastrophic forgetting. The forgetting at phase  $t$  ( $t > 0$ ) is calculated as  $F_t = \frac{1}{t} \sum_{j=1}^{t-1} f_j^t$ , where  $f_j^t$  denotes the performance drop of classes:

$$f_j^t = \max_{i \in \{j, \dots, t-1\}} a_{i,j} - a_{t,j}, \quad (20)$$

where  $a_{i,j}$  represents the accuracy of classes first encountered in phase  $j$  after the model has been incrementally trained up to phase  $i$  ( $i > j$ ). Other implementation details on the settings are available in the Appendix (*i.e.*, Appendix A.1.2).

+CPO	+PFU	+IPU	5	10	20	Method	5	10	20
			48.51	46.66	40.29	Rps	63.74	62.71	59.06
✓			51.55	49.87	48.60	Hat	59.44	57.69	55.68
✓	✓		52.61	51.97	51.17	Iap	56.00	55.11	52.79
✓	✓	✓	53.69	52.88	51.94	Piggy	55.79	54.36	38.78
						Ours	66.64	65.84	61.83

Table 1: Ablation study of our method on TinyImageNet. CPO, PFU and IPU represent the proposed components in Sec. 3.3. 5, 10 and 20 represents the number of incremental phases (*i.e.*, P).

Table 2: The impact of the pathway structure on CIFAR-100. Rps, Hat, Iap and Piggy are detailed in Sec. 4.3.

## 4.2 ABLATION STUDY

To prove the effectiveness of our proposed method, we conduct several ablation experiments on TinyImageNet. The performance of our scheme is mainly attributed to three prominent components: the class-specific pathway organization strategy (CPO), the pathway-guided feature update (PFU) and the incremental pathway update (IPU) mechanism. Since the three components are sequential, we add them gradually for comparison. As can be seen in Tab. 1, CPO bring a 3.04%, 3.21% and 8.31% improvement in overall performance. It demonstrates that the initial pathway decoupling plays an important role in mitigating the interference during the incremental process, especially in the case of longer phases. IPU and PFU also achieves average improvement of 1 and 2 points, facilitating the rectification of features and pathways during the subsequent incremental processes.

## 4.3 ANALYSIS

**The impact of the pathway optimization strategy.** To explore the impact of pathway optimization strategy on the incremental representation learning, we compare our methods to some classical pathway-related ones. Since most of methods are not designed for class-incremental learning, we adapt their core strategies in our settings. As shown in Tab. 2, our method is obviously superior to other ones in three settings. Piggy (Mallya et al., 2018) simply optimizes the mask of parameters on the basis of the initial model, which is not sufficient to handle the complex incremental process. The hard threshold adopted in Hat (Serra et al., 2018) and Iap (Chen et al., 2020) brings great optimization difficulty. Although the RPS (Rajasegaran et al., 2019) achieves good results, its complex network structure and random path search strategy are not efficient.

**The impact of the numbers of divided groups (*i.e.*,  $L$ ).** To explore the sensitivity of divided groups on the incremental performance, we design the following experiments. We divide the output channels into different channels equally. If the channels are not divisible, we round down it. It can be seen in Fig. 3 (a) that the performance fluctuates little except for exceptionally few divisions, demonstrating the stability of our pathway learning. When the number of division is equal to 2, the overall decoupling space for pathways is too small to promote sparse learning.

**The impact of the maximum sparse rate (*i.e.*,  $\zeta_{max}$ ).** To explore the effect of sparse rate on the incremental performance, we conduct multiple experiments with different sparse rates on CIFAR-100. As shown in Fig. 3 (b), the performance with high sparse rate is obviously worse than that with other values. In this case, due to the increase of difficulty of pathway independence, the initial classification accuracy is greatly disturbed. When the sparse rate is too low (*e.g.*, 0.2), the initial accuracy is obviously higher, bring the overall improvement of the incremental performance. When the sparsity value is between 0.3 and 0.45, the initial accuracy is consistent and the incremental performance gets better with heavier sparsity, demonstrating the effectiveness of the pathway decoupling.

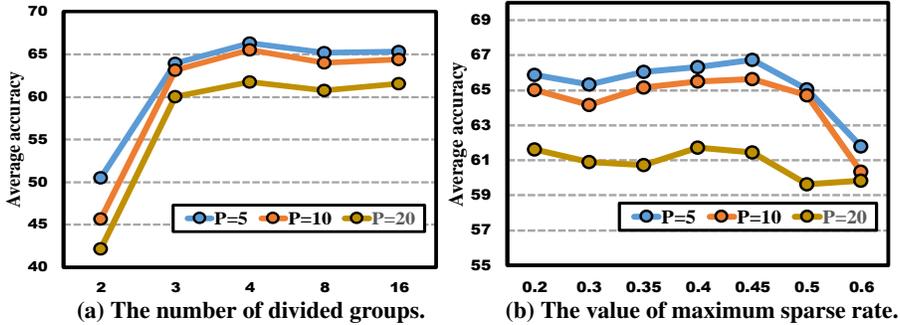


Figure 3: The impact of the values of divided groups (*i.e.*,  $L$  in Sec. 3.3) and the maximum sparse rate (*i.e.*,  $\zeta_{max}$  in Sec. 3.3) on the incremental performance.

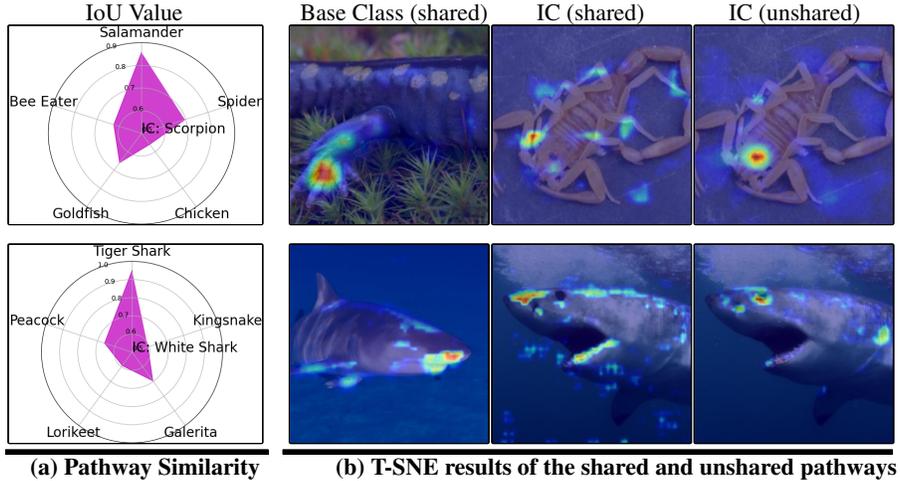


Figure 4: Effect of our scheme on the pathway learning. (a) CPO realizes the organization of distinguishable pathways, thus mitigating the overlap between the incremental classes (*i.e.*, IC) and old ones. (b) PFU promotes the pathway expansion of similar classes. The first two columns represent the activated features of shared pathways, and the last represents the unshared ones.

#### 4.4 VISUALIZATION

To better demonstrate the role of CPO and PFU during optimization, we show the corresponding visualization results. In Fig. 4 (a), the center of the circle represents the incremental class, and the surrounding represents the five different base classes. The middle values represent the intersection of union (IoU) of pathways between the new and old classes. It can be seen the pathways are class-specific, and the similarity is also positively related to the class relationship. For example, the pathway of white sharp is closer to the one of tiger sharp. As shown in Fig. 4 (b), for the incremental class, the features of shared and unshared pathways are visualized by t-SNE (Maaten & Hinton, 2008). For example, the white sharp and tiger sharp are discriminatory to other classes due to the features of teeth. To further distinguish between these two ones, the white shark expands new pathways to learn the texture features on their bodies. Owing to our PFU, the incremental pathways are promoted to differentiate from the old ones, thus improving the separation of novel clusters.

#### 4.5 COMPARISON WITH SOTA

To better assess the overall performance, we compare it to the SOTA of NECIL (LwF\_MC (Rebuffi et al., 2017), MUC, SDC, PASS, IL2A, ABD and SSRE) and some classical exemplar-based CIL methods (iCART (Rebuffi et al., 2017), EEIL, UCIR and PODNet (Douillard et al., 2020)).

As shown in Tab. 3, compared to the SOTA of non-exemplar methods (*i.e.*, E=0), our method achieves average improvement of about 1 point and 2 points on the average accuracy and aver-

Methods		Average Accuracy ( $\uparrow$ )			Average Forgetting ( $\downarrow$ )		
		$P=5$	$P=10$	$P=20$	$P=5$	$P=10$	$P=20$
(1) $E=20$	iCaRL-CNN*	51.07	48.66	44.43	42.13	45.69	43.54
	iCaRL-NCM*	58.56	54.19	50.51	24.90	28.32	35.53
	EEIL* (Castro et al., 2018)	60.37	56.05	52.34	23.36	26.65	32.40
	UCIR* (Hou et al., 2019)	63.78	62.39	59.07	21.00	25.12	28.65
	PODNet $\ddagger$	64.88	63.05	61.62	19.12	22.55	25.64
(2) $E=0$	LwF_MC	45.93	27.43	20.07	44.23	50.47	55.46
	MUC (Yu et al., 2020a)	49.42	30.19	21.27	40.28	47.56	52.65
	SDC $\ddagger$ (Yu et al., 2020b)	56.77	57.00	58.90	6.96	7.50	10.77
	PASS (Zhu et al., 2021b)	63.47	61.84	58.09	25.20	30.25	30.61
	IL2A $\ddagger$ (Zhu et al., 2021a)	65.72	62.69	59.90	27.25	37.35	39.27
	ABD $\ddagger$ (Yin et al., 2020)	63.85	62.46	57.40	23.12	27.34	33.42
	SSRE (Zhu et al., 2022)	65.88	65.04	61.70	18.37	19.48	19.00
	Ours	<b>66.64+0.76</b>	<b>65.84+0.80</b>	<b>61.83+0.13</b>	<b>6.50+0.46</b>	<b>3.30+4.20</b>	<b>9.14+1.63</b>

Table 3: Comparisons with other methods on CIFAR-100 dataset. P represents the number of phases and E represents the number of exemplars. Models with an asterisk \* represent the reproduced results in (Zhu et al., 2021b). Models with a marker  $\ddagger$  represent the reproduced results by this paper. The red footnotes in the last row represent the relative improvement compared with the results of SOTA.

Methods		TinyImageNet			ImageNet-Subset
		$P=5$	$P=10$	$P=20$	$P=10$
(1) $E=20$	iCaRL-CNN*	34.64	31.15	27.90	50.53
	iCaRL-NCM* (Rebuffi et al., 2017)	45.86	43.29	38.04	60.79
	EEIL* (Castro et al., 2018)	47.12	45.01	40.50	63.34
	UCIR* (Hou et al., 2019)	49.15	48.52	42.83	66.16
(2) $E=0$	LwF_MC (Rebuffi et al., 2017)	29.12	23.10	17.43	31.18
	MUC (Yu et al., 2020a)	32.58	26.61	21.95	35.07
	MAS (Aljundi et al., 2018)	18.97	11.82	7.17	19.11
	EWC (Kirkpatrick et al., 2017)	19.64	16.18	17.09	27.32
	PASS (Zhu et al., 2021b)	49.55	47.29	42.07	61.80
	SSRE (Zhu et al., 2022)	50.39	48.93	48.17	67.69
Ours	<b>53.69+3.30</b>	<b>52.88+3.95</b>	<b>51.94+3.77</b>	<b>69.22+1.53</b>	

Table 4: Comparisons of the average incremental accuracy (%) with other methods on TinyImageNet and ImageNet-Subset. P represents the number of phases and E represents the number of exemplars. Models with an asterisk \* represent the reproduced results in Zhu et al. (2021b).

age forgetting of CIFAR-100 dataset, respectively. The performance of our method is comparable to the classical exemplar-based methods (*i.e.*,  $E=20$ ), which shows that our method further mitigate the gap between the two settings. To provide further insight into the behaviors of different methods on larger benchmarks, we compare their average accuracy on TinyImageNet, ImageNet-Subset and ImageNet-Full. As shown in Tabs. 4 and 5, our method achieves average improvement of 3 points. Due to the larger size images in these datasets, the pathway independence during the feature optimization is clearer, bringing greater performance improvement.

Methods	iCaRL-NCM $\ddagger$	UCIR $\ddagger$	PODNet $\ddagger$	PASS $\ddagger$	SSRE $\ddagger$	Ours
$E=0, P=10$	32.43	53.27	50.67	55.90	58.12	<b>60.20+2.08</b>

Table 5: Comparisons of the average incremental accuracy (%) with other methods on ImageNet-Full. Models with an asterisk  $\ddagger$  represent the reproduced results by this paper.

## 5 CONCLUSION

In this paper, a novel self-organized pathway expansion scheme is presented for the NECIL task. A class-specific pathway organization strategy is first proposed to mitigate the feature interference during the optimization of pathway-aware representation. Based on the learnable pathway planner, a pathway-guided feature update mechanism is introduced to adjust the involvement in joint training of classification and distillation. Experimental results show that our method is superior in both performance and adaptability to the state-of-the-art methods, especially on larger datasets.

## REFERENCES

- Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 139–154, 2018.
- Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 233–248, 2018.
- Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.
- Hung-Jen Chen, An-Chieh Cheng, Da-Cheng Juan, Wei Wei, and Min Sun. Mitigating forgetting in online continual learning via instance-aware parameterization. *Advances in Neural Information Processing Systems*, 33:17466–17477, 2020.
- Róbert Csordás, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Are neural nets modular? inspecting functional modularity through differentiable weight masks. In *International Conference on Learning Representations*, 2020.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255. Ieee, 2009.
- Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XX 16*, pp. 86–102. Springer, 2020.
- Xitong Gao, Yiren Zhao, Łukasz Dudziak, Robert Mullins, and Cheng-zhong Xu. Dynamic channel pruning: Feature boosting and suppression. In *International Conference on Learning Representations*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 831–839, 2019.
- Xinting Hu, Kaihua Tang, Chunyan Miao, Xiansheng Hua, and Hanwang Zhang. Distilling causal effect of data in class-incremental learning. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3956–3965, 2021.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- Ya Le and Xuan Yang. Tiny imagenet visual recognition challenge. *CS 231N*, 7(7):3, 2015.
- Lucas Liebenwein, Cenk Baykal, Harry Lang, Dan Feldman, and Daniela Rus. Provable filter pruning for efficient neural networks. In *International Conference on Learning Representations*, 2019.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3431–3440, 2015.

- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.
- Arun Mallya, Dillon Davis, and Svetlana Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning to mask weights. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 67–82, 2018.
- Jathushan Rajasegaran, Munawar Hayat, Salman Khan, Fahad Shahbaz Khan, and Ling Shao. Random path selection for incremental learning. *Advances in Neural Information Processing Systems*, 2019.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pp. 618–626, 2017.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pp. 4548–4557. PMLR, 2018.
- Yang Sui, Miao Yin, Yi Xie, Huy Phan, Saman Aliari Zonouz, and Bo Yuan. Chip: Channel independence-based pruning for compact neural networks. *Advances in Neural Information Processing Systems*, 34, 2021.
- Gido M. van de Ven and A. Tolia. Three scenarios for continual learning. *ArXiv*, abs/1904.07734, 2019.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 374–382, 2019.
- Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3014–3023, 2021.
- Hongxu Yin, Pavlo Molchanov, Jose M Alvarez, Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K Jha, and Jan Kautz. Dreaming to distill: Data-free knowledge transfer via deepinversion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8715–8724, 2020.
- L. Yu, S. Parisot, G. Slabaugh, J. Xu, and T. Tuytelaars. More classifiers, less forgetting: A generic multi-classifier paradigm for incremental learning. *European Conference on Computer Vision*, 2020a.
- Lu Yu, Bartłomiej Twardowski, Xialei Liu, Luis Herranz, Kai Wang, Yongmei Cheng, Shangling Jui, and Joost van de Weijer. Semantic drift compensation for class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6982–6991, 2020b.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pp. 3987–3995. PMLR, 2017.
- Fan Zhang, Duyu Tang, Yong Dai, Cong Zhou, Shuangzhi Wu, and Shuming Shi. Skillnet-nlu: A sparsely activated model for general-purpose natural language understanding. *arXiv e-prints*, pp. arXiv–2203, 2022.
- Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2921–2929, 2016.

Fei Zhu, Zhen Cheng, Xu-yao Zhang, and Cheng-lin Liu. Class-incremental learning via dual augmentation. *Advances in Neural Information Processing Systems*, 34, 2021a.

Fei Zhu, Xu-Yao Zhang, Chuang Wang, Fei Yin, and Cheng-Lin Liu. Prototype augmentation and self-supervision for incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5871–5880, 2021b.

Kai Zhu, Wei Zhai, Yang Cao, Jiebo Luo, and Zheng-Jun Zha. Self-sustaining representation expansion for non-exemplar class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9296–9305, 2022.

## A APPENDIX

### A.1 DETAILED EXPLANATION

#### A.1.1 STANDARD DEVIATION OF THE INCREMENTAL PERFORMANCE (ERROR BARS)

All results of the average incremental accuracy and average forgetting are evaluated on three different runs. To show the stability of our method, we report its standard deviation on three runs. As shown in Tabs. 6 and 7, random factors have little impact on our scheme.

Methods		Average Accuracy ( $\uparrow$ )			Average Forgetting ( $\downarrow$ )		
		$P=5$	$P=10$	$P=20$	$P=5$	$P=10$	$P=20$
(1) $E=20$	iCaRL-CNN*	51.07	48.66	44.43	42.13	45.69	43.54
	iCaRL-NCM*	58.56	54.19	50.51	24.90	28.32	35.53
	EEIL*	60.37	56.05	52.34	23.36	26.65	32.40
	UCIR* (Hou et al., 2019)	63.78	62.39	59.07	21.00	25.12	28.65
	PODNet $\ddagger$	64.88	63.05	61.62	19.12	22.55	25.64
(2) $E=0$	LwF_MC	45.93	27.43	20.07	44.23	50.47	55.46
	MUC (Yu et al., 2020a)	49.42	30.19	21.27	40.28	47.56	52.65
	SDC $\ddagger$ (Yu et al., 2020b)	56.77	57.00	58.90	6.96	7.50	10.77
	PASS (Zhu et al., 2021b)	63.47	61.84	58.09	25.20	30.25	30.61
	IL2A $\ddagger$ (Zhu et al., 2021a)	65.72	62.69	59.90	27.25	37.35	39.27
	ABD $\ddagger$ (Yin et al., 2020)	63.85	62.46	57.40	23.12	27.34	33.42
	Ours	<b>66.64<math>\pm</math>0.01</b>	<b>65.84<math>\pm</math>0.07</b>	<b>61.83<math>\pm</math>0.12</b>	<b>6.50<math>\pm</math>0.13</b>	<b>3.30<math>\pm</math>0.39</b>	<b>9.14<math>\pm</math>1.42</b>

Table 6: Comparisons with other methods on CIFAR-100 dataset. P represents the number of phases and E represents the number of exemplars. Models with an asterisk \* represent the reproduced results in (Zhu et al., 2021b). Models with a marker  $\ddagger$  represent the reproduced results by this paper. The blue footnotes in the last row represent the values of error bars.

Methods		TinyImageNet			ImageNet-Subset
		$P=5$	$P=10$	$P=20$	$P=10$
(1) $E=20$	iCaRL-CNN*	34.64	31.15	27.90	50.53
	iCaRL-NCM*	45.86	43.29	38.04	60.79
	EEIL* (Castro et al., 2018)	47.12	45.01	40.50	63.34
	UCIR* (Hou et al., 2019)	49.15	48.52	42.83	66.16
(2) $E=0$	LwF_MC (Rebuffi et al., 2017)	29.12	23.10	17.43	31.18
	MUC (Yu et al., 2020a)	32.58	26.61	21.95	35.07
	MAS (Aljundi et al., 2018)	18.97	11.82	7.17	19.11
	EWC (Kirkpatrick et al., 2017)	19.64	16.18	17.09	27.32
	PASS (Zhu et al., 2021b)	49.55	47.29	42.07	61.80
Ours	<b>53.69<math>\pm</math>0.14</b>	<b>52.88<math>\pm</math>0.05</b>	<b>51.94<math>\pm</math>0.28</b>	<b>69.22<math>\pm</math>0.05</b>	

Table 7: Comparisons of the average incremental accuracy (%) with other methods on TinyImageNet and ImageNet-Subset. P represents the number of phases and E represents the number of exemplars. Models with an asterisk \* represent the reproduced results in (Zhu et al., 2021b). The blue footnotes in the last row represent the values of error bars.

### A.1.2 DETAILED SETTING

We use an Adam optimizer, in which the initial learning rate is set to 0.001 and the attenuation rate is set to 0.0002. The batch size is set to 128. The model stops training after 160 epochs and 60 epochs during the initial phase and incremental phases, respectively. We adopt ResNet18 and 3 standard convolution blocks as the backbone of feature extractor  $f_{\theta_t}$  and pathway planner  $f_{\alpha_t}$ , respectively.

In the main text, the maximum value of sparse rate in Eq. (8) is set to 0.4. The values of  $e_1$  and  $e_2$  in Eq. (10) are set to 0.08 and 0.75, respectively. The values of  $L$  and  $K$  in Section 3.3 are set to 4 and 16, respectively. One NVIDIA GTX2080Ti gpu is utilized for CIFAR-100 and TinyImageNet datasets. Two NVIDIA GTX3090 and eight NVIDIA Tesla A100 gpu are utilized for ImageNet-Sub and ImageNet-Full datasets, respectively. All datasets adopted in this paper are open to the public.

### A.1.3 OPTIMIZATION EXPLANATION IN OUR SCHEME

The optimization of feature representation  $\theta_t$  is mainly guided by the classification loss function  $\mathcal{L}_{cls}$  and feature distillation loss function  $\mathcal{L}_{kd}$ . Assuming that the optimal solution at the incremental phase  $t - 1$  is taken when  $\theta_{t-1} = \theta_{t-1}^*$ . As  $\theta_t$  is initialized by the value of  $\theta_{t-1}^*$ , it can be assumed that  $\theta_t$  is close to  $\theta_{t-1}^*$ . Then the Taylor expansion on  $\theta_t$  can be written as follows,

$$f(\theta_t) = f(\theta_{t-1}^*) + \left(\frac{\partial f(\theta)}{\partial \theta} \Big|_{\theta=\theta_{t-1}^*}\right)(\theta_t - \theta_{t-1}^*) + \frac{1}{2}(\theta_t - \theta_{t-1}^*)^T \left(\frac{\partial^2 f(\theta)}{\partial^2 \theta} \Big|_{\theta=\theta_{t-1}^*}\right)(\theta_t - \theta_{t-1}^*) + o(\theta_{t-1}^*). \quad (21)$$

The first order component is constrained to zero by the gradient descent, and the ones higher than second order can be ignored. The subscript  $t$  can be omitted for brevity, and Eq. (21) can be approximated as follows:

$$f(\theta) = f(\theta^*) + \frac{1}{2}(\theta - \theta^*)^2 f''(\theta^*) = f(\theta^*) + \frac{1}{2}\Omega(\theta - \theta^*)^2 = f(\theta^*) + \frac{1}{2}(\Omega_{cls} + \Omega_{kd})(\theta - \theta^*)^2, \quad (22)$$

where  $\Omega_{cls}$  and  $\Omega_{kd}$  represents the importance of parameter space on the classification and distillation tasks, which is commonly estimated in different incremental methods (Kirkpatrick et al., 2017; Aljundi et al., 2018). To mitigate the interference between the two objectives, we can improve their respective weight sparsity (*i.e.*, the sparsity of  $\Omega_{cls}$  and  $\Omega_{kd}$ ), and reduce the shared space of important parameters.

### A.1.4 DETAILED VALUES OF THE CURVES

To facilitate the fair comparison of subsequent work, we report the detailed values of incremental accuracy for each phase in Tabs. 8 to 10. The average accuracy is consistent with the one in Table 3 and 4 of the main text.

Dataset	Phase									
	0	1	2	3	4	5	6	7	8	9
A	82.40	78.23	74.87	72.34	68.62	67.96	65.52	64.84	62.57	60.83
B	62.70	59.92	58.55	57.01	55.25	54.42	53.18	52.74	52.27	51.70

Dataset	Phase									
	10	11	12	13	14	15	16	17	18	19
A	59.76	58.85	57.39	55.72	54.66	53.54	53.21	52.55	52.24	51.54
B	51.25	50.76	50.19	49.25	48.71	47.95	47.66	47.09	46.66	45.57

Dataset	Phase
	20
A	50.82
B	45.03

Table 8: Detailed values of classification accuracy under the setting of 20 incremental phases. A and B represent the CIFAR-100 and TinyImageNet datasets, respectively.

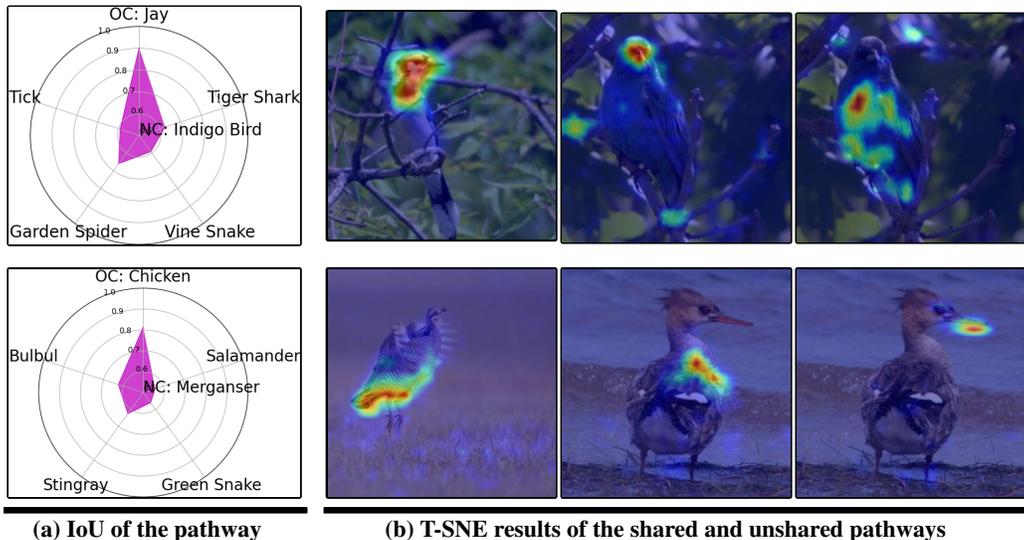


Figure 5: Effect of our scheme on the pathway learning. (a) CPO realizes the organization of distinguishable pathways, thus mitigating the overlap between the incremental classes and the old ones. NC represents the novel classes. (b) PFU promotes the pathway expansion of similar classes. The first two columns represent the shared pathways, and the last represents the unshared ones.

Dataset	Phase										
	0	1	2	3	4	5	6	7	8	9	10
CIFAR-100	80.90	76.15	73.13	69.58	66.73	64.59	63.01	60.81	58.68	56.16	54.50
TinyImageNet	62.70	59.05	56.00	54.25	53.03	52.37	51.21	50.13	48.73	47.81	46.41
ImageNet-Subset	83.40	77.29	73.93	71.75	69.64	68.56	67.61	65.07	63.01	61.22	59.91
ImageNet-Full	76.46	67.59	64.92	62.89	60.61	58.72	57.12	55.75	54.17	52.30	51.65

Table 9: Detailed values of classification accuracy under the setting of 10 incremental phases.

### A.1.5 MORE RESULTS ON VISUALIZATION

To better demonstrate the role of CPO and PFU during optimization, we show more corresponding visualization results. In Fig. 5 (a), the center of the circle represents the novel class, and the surrounding represents the five different base classes. The middle values represent the intersection of union (IoU) of pathways between the new and old classes. It can be seen the pathways are class-specific, and the similarity is also positively related to the class relationship. As shown in Fig. 5 (b), the features of shared and unshared pathways are visualized by Grad-CAM (Selvaraju et al., 2017). To further distinguish between the old and novel class, the novel one expands new pathways to learn representative features.

### A.1.6 CONFUSION MATRIX

To evaluate performance of both old and new classes during training, we compare their accuracy on two setting (*i.e.* 5 and 10 incremental phases). As shown in Fig. 6, our method achieves similar performance between the old and new classes without favoring one side due to overfitting, which is a prerequisite for a good incremental learning system.

### A.1.7 RELATED WORK ON FILTER PRUNING METHODS

Network pruning (Liebenwein et al., 2019; Sui et al., 2021; Gao et al., 2018) is an important technology to reduce memory size and bandwidth. Recently, various network pruning techniques have been proposed, which can be classified from the structural aspect, *i.e.*, the structured and unstructured pruning. Specifically, structured methods remove parameters in groups by pruning neurons, filters, or channels. Classical filter pruning methods (Sui et al., 2021; Gao et al., 2018) make up the

Dataset	Phase					
	0	1	2	3	4	5
CIFAR-100	80.90	72.63	65.87	62.94	60.76	56.76
TinyImageNet	62.70	57.35	54.30	52.04	48.96	46.79

Table 10: Detailed values of classification accuracy under the setting of 5 incremental phases.

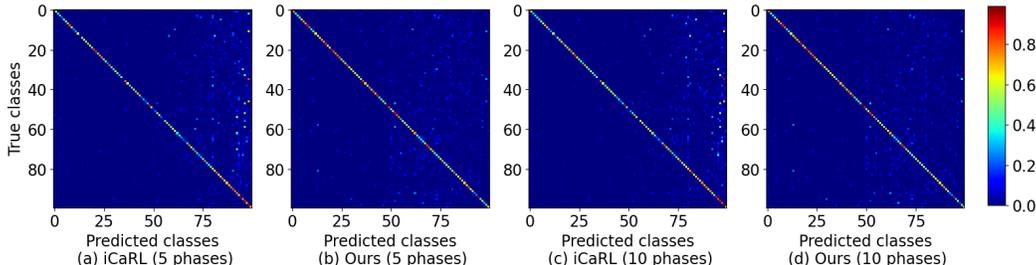


Figure 6: Confusion matrices of different methods on CIFAR-100. 5 phases and 10 phases settings are considered to evaluate the stability of our method on the old and novel classes.

prominent family of structured methods for CNNs. Different from the pruning methods designed for the network efficiency, our scheme aims at the mitigation of update interference. The concept of group in this paper is slightly different as only the output channels of features are divided for the organization of pathway.

#### A.1.8 LIMITATION AND SOCIETAL IMPACT

The division way of the standard classification model in our method is too simple, which constrains the adjustment of some factors. As shown in Fig. 3, the maximum sparsity can only be kept below 0.5, which deserves the further improvement. Our non-exemplar method avoids the issue of privacy but an old model needs to be maintained during the training, which poses a risk of information leak. This calls for future research that addresses this aspect.

## A.2 ADDITIONAL RESULTS

### A.2.1 PATHWAY VISUALIZATION

To better demonstrate the role of self-organizing pathway expansion scheme during optimization, we show more visualization results on the pathways of different classes. For the simplicity of viewing, we plot the most important group (*i.e.*, vertical coordinate) in each module (*i.e.*, horizontal coordinate). As shown in Fig. 7 (a), at the initial phase, two different classes (*i.e.*, the classes in the first and second columns) tend to utilize different pathways to extract the corresponding features. At the incremental phase, the novel class (*i.e.*, the class in the third column) tends to utilize the novel pathway to optimize the incremental features, and the whole pathway is similar to the semantically close class (*i.e.*, the class in the second column). The observation is consistent with the one from Fig. 4 in the main text.

### A.2.2 THE IMPACT OF THE SET EPOCHS

To explore the effect of the set epochs on the incremental performance, we conduct multiple experiments with different start epochs (*i.e.*  $e_1$  in Eq. (10) of the main text) and end epochs (*i.e.*  $e_2$  in Eq. (10) of the main text) on CIFAR-100. As shown in Fig. 8 (b), the performance with larger start epoch is obviously worse than those with other values. In this case, due to the increase of initial classification accuracy with full pathways, it is more difficult to separate the class-specific pathway, which influences the overall performance. At the same time, the values of the end epoch have almost no effect on the incremental performance, which is more robust to the optimization process.

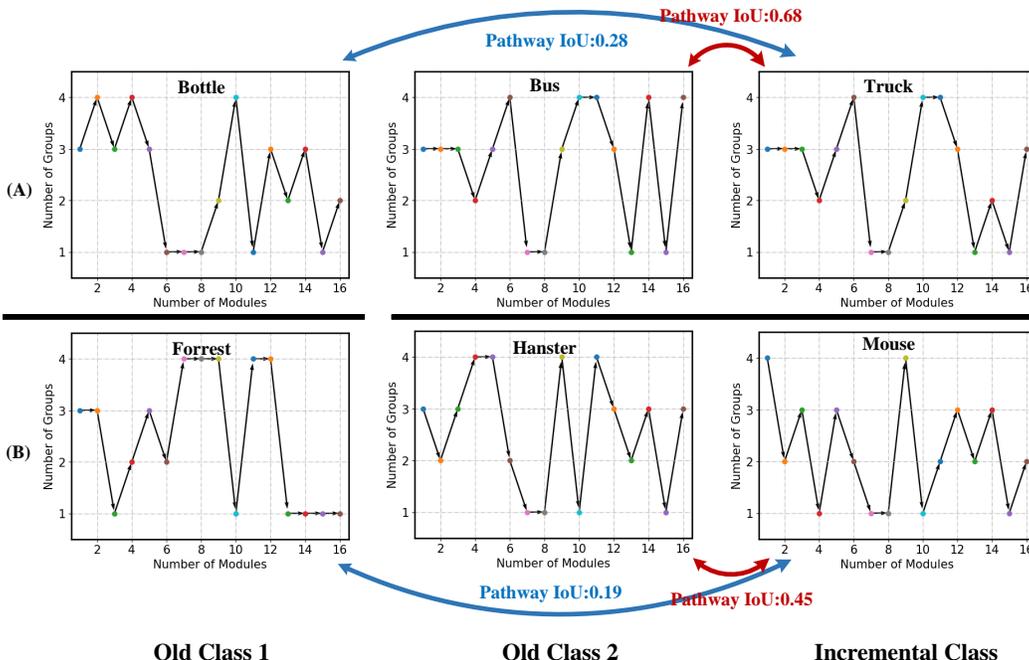


Figure 7: More visualization results on the pathway. The first and second columns represent the pathways of two old classes, which differ significantly in semantics. The third column represents the pathway of the incremental class, which is semantically closer to the one in the second column. For the simplicity of viewing, we plot the most important group (*i.e.*,  $L$ ) in each module (*i.e.*,  $K$ ). Pathway IoU represents the overlap rate of corresponding class-specific pathways.

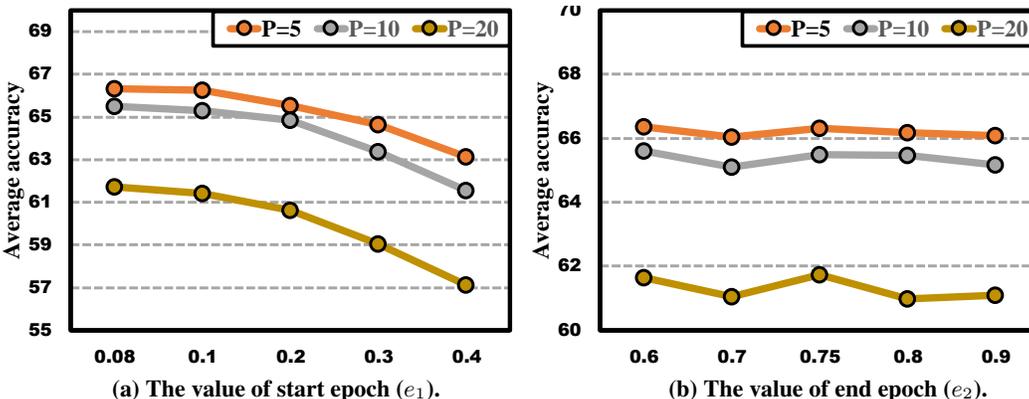


Figure 8: The impact of the values of the set epochs (*i.e.*,  $e_1$  and  $e_2$  in Eq. (10) of the main text) on the incremental performance in the three-step strategy.

### A.2.3 FURTHER ANALYSIS OF MORE CIL METHODS

In the comparative experiments of the main text, we compare with some classical CIL methods at two different settings, demonstrating that our method reduces the gap between the two settings. At the same time, most of the classical methods are not applicable to the NECIL settings, let alone the latest CIL methods. For example, we adapt the latest CIL method dynamic expandable network (Yan et al., 2021) to the NECIL setting (*i.e.* NDER), and its performance is poor as shown in Tab. 11. Due to the lack of old samples, it is difficult to perform effective optimization with such large expanding parameters.

Method	CIFAR-100		
	5 phases	10 phases	20 phases
NDER	29.08	21.13	13.10
Ours	66.64	65.84	61.83

Table 11: Further analysis on the CIL method.

#### A.2.4 GENERALIZATION TO THE CIL SETTING

To further prove the effectiveness and generalization of our method, we introduce it into the CIL setting. As (Douillard et al., 2020) is one of the SOTA methods in CIL setting, we modify its implementation with our self-organized pathway expansion scheme directly. As shown in Tab. 12, our method achieves average improvement of 2 points. Even if the effect of incremental samples on the overall performance is weakened by exemplars in CIL setting, our scheme still brings a boost to the existing method (Douillard et al., 2020). It can be seen that our method has great potential for the CIL setting, which will serve as our future work.

#### A.2.5 COMPARISON WITH SOTA ON IMAGENET-FULL DATASET

To better assess the overall performance of our scheme on larger dataset, we compare it to the SOTA of NECIL (PASS) and some classical methods of exemplar-based CIL (iCaRL, UCIR and PODNet) on ImageNet-Full.

As shown in Tab. 13, compared to the SOTA of non-exemplar methods (*i.e.*,  $E=0$ ), our method achieves average improvement of 2 points on the average accuracy. The performance of our method is comparable to the classical exemplar-based methods (*i.e.*,  $E=20$ ), which shows that our method further mitigate the gap between the two settings on larger dataset.

Method	CIFAR-100 (B50)	
	5 phases	10 phases
Podnet	64.88	63.05
Ours	66.64	65.84

Table 12: Comparisons of the average incremental accuracy (%) under the CIL setting.

Methods		ImageNet-Full
		$P=10$
$E=20$	iCaRL (Rebuffi et al., 2017)	46.72
	UCIR (Hou et al., 2019)	63.27
	PODNet (Douillard et al., 2020)	64.17
$E=0$	iCaRL <sup>‡</sup> (Rebuffi et al., 2017)	32.43
	UCIR <sup>‡</sup> (Hou et al., 2019)	53.27
	PODNet <sup>‡</sup> (Douillard et al., 2020)	50.67
	PASS <sup>‡</sup> (Douillard et al., 2020)	55.90
	SSRE <sup>‡</sup> (Zhu et al., 2022)	58.12
	Ours	<b>60.20</b>

Table 13: Comparisons of the average incremental accuracy (%) with other methods on ImageNet-Full. P represents the number of phases and E represents the number of exemplars. Models with an asterisk <sup>‡</sup> represent the reproduced results by this paper.