# Consistent Prompting for Rehearsal-Free Continual Learning

Zhanxin Gao[1], Jun Cen[2], Xiaobin Chang[1,3*]

[1]School of Artificial Intelligence, Sun Yat-sen University, China
[2]Cheng Kar-Shun Robotics Institute, The Hong Kong University of Science and Technology, China
[3]Key Laboratory of Machine Intelligence and Advanced Computing, Ministry of Education, China

gaozhx27@mail2.sysu.edu.cn, jcenaa@connect.ust.hk, changxb3@mail.sysu.edu.cn

## Abstract

*Continual learning empowers models to adapt autonomously to the ever-changing environment or data streams without forgetting old knowledge. Prompt-based approaches are built on frozen pre-trained models to learn the task-specific prompts and classifiers efficiently. Existing prompt-based methods are inconsistent between training and testing, limiting their effectiveness. Two types of inconsistency are revealed. Test predictions are made from all classifiers while training only focuses on the current task classifier without holistic alignment, leading to Classifier inconsistency. Prompt inconsistency indicates that the prompt selected during testing may not correspond to the one associated with this task during training. In this paper, we propose a novel prompt-based method, Consistent Prompting (CPrompt), for more aligned training and testing. Specifically, all existing classifiers are exposed to prompt training, resulting in classifier consistency learning. In addition, prompt consistency learning is proposed to enhance prediction robustness and boost prompt selection accuracy. Our Consistent Prompting surpasses its prompt-based counterparts and achieves state-of-the-art performance on multiple continual learning benchmarks. Detailed analysis shows that improvements come from more consistent training and testing. Our code is available at https://github.com/Zhanxin-Gao/CPrompt.*

## 1. Introduction

Continual learning [2, 6, 26, 37, 39] aims to equip deep models with the capacity to continuously acquire new knowledge, e.g., learn to recognize new object categories, while handling catastrophic forgetting [18, 27, 28, 32] of the old knowledge. Rehearsal-based methods [3, 34, 49] alleviate the forgetting problem with a small number of exemplars of previous tasks stored in the memory buffer and
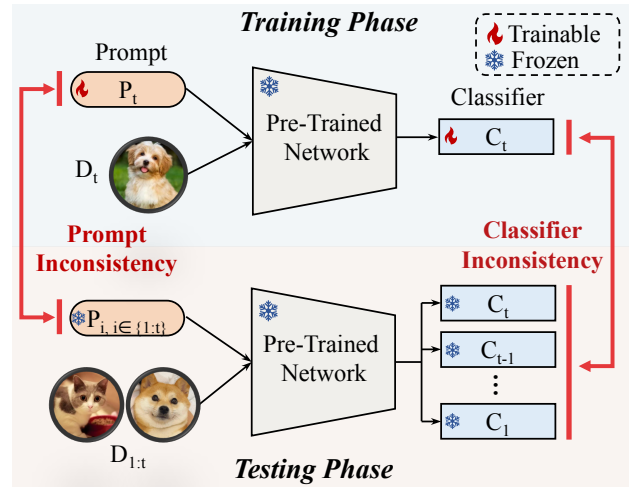
*indicates corresponding author.



Figure 1. Existing prompt-based methods are inconsistent between training and testing, and prompt inconsistency and classifier inconsistency are illustrated. $P_t$ and $C_t$ represent the prompt and the classifier of the current task t, respectively.

replayed with the new task data during training. However, exemplars of previous tasks may not be available due to constraints such as data privacy or memory limitations. Therefore, rehearsal-free continual learning methods [9, 23, 25, 35, 43, 51] have attracted much attention.

Prompt-based models [30, 36, 38, 45, 46] have shown exceptional results in rehearsal-free continual learning. Based on the frozen backbone network pre-trained on a large-scale dataset, such models can efficiently adapt to the new task by only training a tiny set of parameters, i.e., the prompts and fully connected (FC) classifiers. However, the training and testing of existing approaches lack consistency, as the prompt and the classifier are optimized solely within the current task. As shown in Figure 1, two types of such inconsistencies are discussed. Firstly, classifier inconsistency indicates that test predictions are made from all classifiers rather than only from the training one. Secondly, prompt
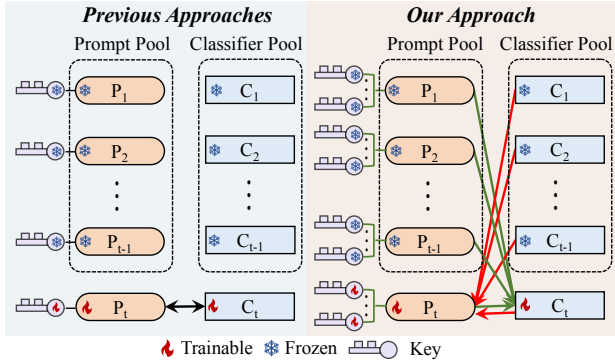
Figure 2. Previous approaches typically train the current task prompt and classifier in isolation. Our Consistent Prompting leverages all existing prompts and classifiers to instruct the training of the current task prompt and classifier. Meanwhile, we suggest using multiple keys to map each task prompt, instead of relying on a single key, to adapt to the diverse nature of each task.

inconsistency refers to the mismatch of prompts between training and testing. Existing prompt-based methods do not identify or handle such inconsistent issues, leading to suboptimal performance.

To align the training of prompts and classifiers with testing, we propose a novel training method, called Consistent Prompting (CPrompt), which surpasses existing prompt-based ones. CPrompt consists of two modules: Classifier Consistency Learning (CCL) and Prompt Consistency Learning (PCL). Specifically, CCL is proposed to address the classifier inconsistency issue by exposing the current task prompt training to all classifiers seen so far. On the other hand, PCL is proposed to handle the prompt inconsistency problem. The classifier of the current task is trained with prompts randomly selected from the pool. It enables the classifier to better adapt to different prompts for more robust prediction. Moreover, a multi-key mechanism is proposed to boost the prompt selection accuracy in PCL. The comparison between existing prompt-based approaches with our Consistent Prompting is illustrated in Figure 2. We evaluate the proposed method on four challenging continual learning benchmark datasets. Extensive analysis is also provided to demonstrate the superior performance of CPrompt mainly comes from its consistent training and testing. The contributions of this work are summarized in three-fold:

1. The inconsistency issues between the training and testing of prompt-based rehearsal-free continual learning methods are identified and discussed for the first time in this work. A novel Consistent Prompting (CPrompt) is then proposed for better consistency.

2. To maintain the classifier consistency at testing, the prompt should be exposed to all seen classifiers during training, as proposed in Classifier Consistency Learning

(CCL) of CPrompt.

3. We propose Prompt Consistency Learning (PCL) in CPrompt with two complementary purposes. For more robust testing predictions, the current classifier should be trained under different prompts. For more precise prompt selection, a multi-key mechanism is exploited.

## 2. Related Work

### 2.1. Continual Learning

Continual learning methods aim to reduce catastrophic forgetting [18, 27] of the old knowledge when adapting to new one. Three kinds of pipelines are proposed based on distinctive perspectives. Regularization-based approaches [1, 15, 21, 50] aim to prevent significant changes to important attributes to protect previously learned knowledge from excessive interference. As a prevalent regularization method, knowledge distillation [4, 8, 47] is widely adopted to transfer the retained knowledge in the previous model (as a teacher) to the current student model. Parameter isolation methods [13, 20, 48] freeze specific parameters while allocating the rest for subsequent tasks or expanding the network for new knowledge learning. Such approaches are inherently intuitive and can yield promising results when an ample number of parameters are extended. However, these techniques often induce model complexity, posing maintenance challenges. Rehearsal [3, 5, 10, 12, 22] is a popular strategy in continual learning, allowing the model to partially access previous exemplars. These approaches can effectively mitigate the forgetting of prior knowledge. However, they require additional memory and computational overhead and raise data privacy concerns as well. The aforementioned paradigms are highly complementary and can be combined to enhance continual learning performance [5, 40, 41, 48].

### 2.2. Prompt-based Methods in Continual Learning

Prompt [19, 24] is a fundamental technique used in Natural Language Processing (NLP). It serves as a transfer approach or provides specific instructions for downstream tasks. Recent prompt-based continual learning methods [14, 29, 33] handle the rehearsal-free setting by encoding task-specific knowledge within prompts. It enables the network to efficiently retrieve previous information by querying the appropriate prompts and eliminates the need for rehearsal buffers. L2P [46] proposes to learn a prompt pool and use a query-key mechanism to select a prompt. However, since the entire prompt pool is always trainable, forgetting of prior knowledge within the prompt pool is inevitable. Instead of using the same prompt pool across tasks, two complementary pools, G-Prompt and E-Prompt, are proposed by DualPrompt [45] to encode task-invariant and tasks-specific knowledge, respectively. CODAPrompt [36] introduces
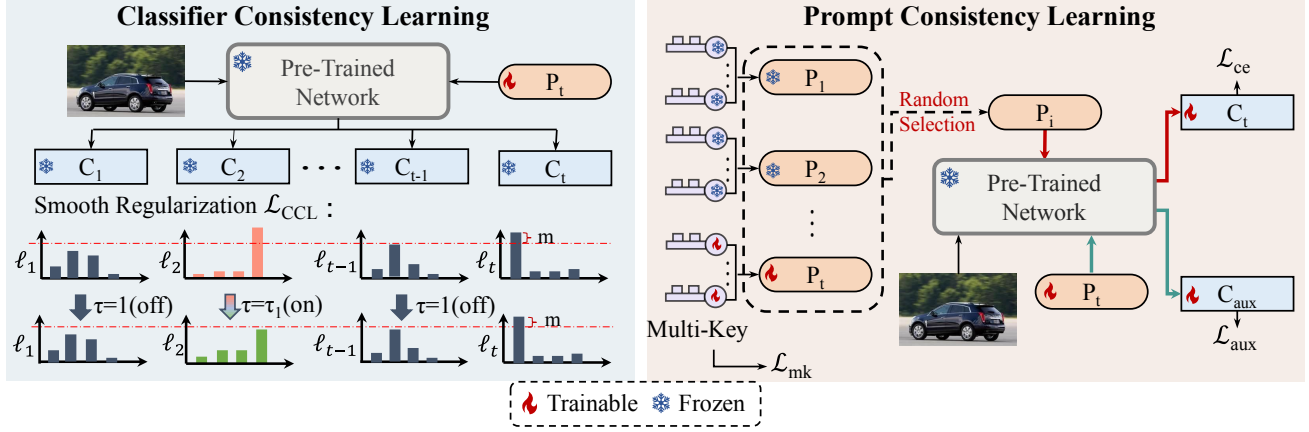
Figure 3. The illustration of the proposed consistent prompting (CPrompt). CPrompt aims to align the training of prompts and classifiers with testing for more consistency. It consists of two main modules: classifier consistency learning (CCL, detailed in Section 3.2) and prompt consistency learning (PCL, detailed in Section 3.3).

a decomposed attention-based prompting method and expands the prompt component according to different tasks. To handle the expansion of classifiers in continual learning, ESN [44] proposes a prompt-based method that employs energy self-normalization. This is achieved by adding a self-attention block, which produces consistent and high confidence scores for in-task data.

Existing prompt-based methods face difficulties in clarifying the relationship between the newly optimized prompts and the ones learned previously. This leads to interference among different prompts or components during the training process. Additionally, these prompts are not fully aligned with either input images or overall classifiers, causing the training-testing inconsistency. To address these challenges, we propose consistent prompting as a solution to enhance the effectiveness of continual learning.

## 3. Consistent Prompting

### 3.1. Prerequisites

**Continual Learning Setting** Continual learning is designed to acquire knowledge from a data stream comprising T non-overlapping sequential datasets, denoted as $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_T\}$. Each dataset $\mathcal{D}_t$ corresponds to a specific task $t$ and can be represented as the union of individual class datasets, i.e., $\mathcal{D}_t = \bigcup_j \mathcal{D}_{t,j}$, where $j$ denotes the $j$th class within task $t$. The objective is to train a mapping function $f$ capable of predicting $f(x) \in Y_{t,j}$ for every input $x \in \mathcal{D}_{t,j}$, where $Y_{t,j}$ represents the $j$th class within task $t$. In this research, we focus on class incremental learning (CIL) where task identities are not provided during testing. Moreover, adopting the rehearsal-free CIL setting prohibits the use of any exemplar from prior tasks.

**Prompt-based Method** A pre-trained vision transformer

(ViT) [7] comprises an embedding layer $f_e$ and multiple self-attention layers $f_i$, $i = 1, 2, ..., N$. Each image $x \in \mathcal{D}_{t,k}$ is initially processed by the embedding layer $f_e$, yielding a sequential feature $\xi_e = f_e(x) \in \mathbb{R}^{L \times D}$, where $L$ refers to the number of patches, i.e., token length, and $D$ denotes the embedding dimension. During training, we keep the parameters of the pre-trained model $f_i$ and $f_e$ frozen and exclusively update additional learnable parameters, known as the prompt $P \in \mathbb{R}^{L_P \times D}$, where $L_P$ indicates the length of the prompt. Specifically, a prompt $P$ can simultaneously fit with $s$ self-attention layers by splitting $P$ into $s$ segments and each segment $p_i \in \mathbb{R}^{(L_P/s) \times D}$, $i = 1, 2, ..., s$. To insert a segment $p_j$ into its corresponding self-attention layer $f_i$, we extend the output $\xi_{i-1}$ of the preceding layer $f_{i-1}$ with $p_j$, resulting the input $[p_j; \xi_{i-1}]$ for $f_i$. During testing, the pre-trained model collaborates with the learned prompt. However, task-specific prompts are learned across tasks in continual learning. Therefore, an additional prompt selection mechanism is required.

### 3.2. Classifier Consistency Learning

Prompt-based approaches typically utilize all task classifiers for inference during test time. However, different classifiers are learned within the corresponding task. The alignment among their behaviors is not guaranteed. The proposed Classifier Consistency Learning (CCL) aims to handle this issue by exposing the current task prompt training to all classifiers encountered, which is straightforward and effective. To achieve accurate predictions, we need to ensure that each input has the highest logit response for its corresponding class among all encountered classifiers. However, we observe that such classifiers without regularization are naturally biased towards the previous parts, resulting in higher logit values for old task classifiers than for the cur-

rent one. Therefore, a new regularization loss is proposed to encourage the maximum logit value of the old task classifiers to be lower than that of the current class by a pre-defined margin. This regularization is applied to the current task prompt to mitigate bias and enable model training. Further details are provided below.

Assuming the pre-trained network as $f_\theta$, we can extract the image feature with respect to a specific class using the token [class] as in [7]:

$$h = f_\theta(x, P_t)[0], \qquad (1)$$

where $x \in D_t$ is an input image of the current training task $t$, 0 means indexing the first [class] token. The extracted feature $h$ is then fed into all task classifiers, and the respective logit values are obtained via,

$$\ell_i = C_i(h), i \in \{1, ..., t\}. \qquad (2)$$

The proposed smooth regularization is built on the entropy of previous tasks. For the $i$th task, $i \in \{1, ..., t-1\}$, an adaptive entropy is calculated based on logit $\ell_i$,

$$\mathcal{L}_e(i) = - < \sigma(\ell_i/\tau), \, \log(\sigma(\ell_i)) >, \qquad (3)$$

where $\sigma$ is the softmax function and $<,>$ is the inner product operator. We also find that blocking the gradients from $\sigma(\ell_i/\tau)$ stabilizes the optimization. The temperature $\tau$ is chosen accordingly,

$$\tau = \begin{cases} \tau_1, & \max(\ell_i) + m \geq \max(\ell_t), \\ 1, & \text{otherwise.} \end{cases} \qquad (4)$$

where smooth regularization is needed when the maximum logit of the previous task classifier $\max(\ell_i)$ exceeds the maximum logit of the current task classifier $\max(\ell_t)$ by a margin of $m \geq 0$. Therefore, $\tau$ equals to a predefined $\tau_1 > 1$ and enables the regularization. Otherwise, smooth regularization is not necessary. Simply setting $\tau$ to 1 yields zero gradients can turns it off. The supplementary material provides relevant theoretical analysis and proof. An illustration of the smooth regularization process is depicted in Figure 3 for better understanding.

The Classifier Consistency Learning loss $\mathcal{L}_{CCL}$ is given by,

$$\mathcal{L}_{CCL} = \frac{\alpha}{t-1} \sum_{i=1}^{t-1} \mathcal{L}_e(i), \qquad (5)$$

where $\alpha$ indicates the strength of the regularization and $t-1$ normalizes the growth of tasks in continual learning.

### 3.3. Prompt Consistency Learning

Another significant problem of prompt-based continual learning is the model's uncertainty in selecting the correct prompt for inference. Therefore, the proposed Prompt Consistency Learning (PCL) aims to establish a more robust prompt-classifier relationship, ensuring correct output even with incorrect prompts. During training, a task-specific prompt $P_i$ is randomly selected from the current prompt pool,

$$P_i, i \sim \text{Uni}(1, t), \qquad (6)$$

where $\text{Uni}(1, t)$ is a uniform distribution on the integers $1, 2, ..., t$. All prompts except $P_t$ are frozen to preserve the encoded knowledge and defy catastrophic forgetting. The output logit of the current task classifier $C_t$ can then be obtained via,

$$\ell_t = C_t(f_\theta(x, P_i)[0]). \qquad (7)$$

The corresponding loss is calculated,

$$\mathcal{L}_{ce} = \text{CrossEntropy}(\ell_t, y), \qquad (8)$$

where $y$ is the ground truth class label of the input image $x$. In this way, the classifier is trained to make the correct predictions even based on the wrong prompts, which is more consistent with testing. This training process is illustrated in Figure 3.

It is noteworthy that the current task-specific prompt $P_t$ is selected with a probability of only $1/t$, which results in inadequate training for $P_t$. Therefore, we employ an auxiliary classifier, $C_{aux}$, to assist in the training of $P_t$ as follows:

$$\mathcal{L}_{aux} = \text{CrossEntropy}(C_{aux}(h), y). \qquad (9)$$

Here, $h$ represents the extracted feature calculated by Eq. (1). It is noted that we employ the auxiliary classifier on $C_{aux}$ instead of $C_t$. Employing the auxiliary classifier on $C_t$ would lead $C_t$ to lean towards adapting more to the current prompt $P_t$ rather than adapting to any task prompt with equal probability as in Eq. (8). This could result in $C_t$ lacking robustness when selecting other prompts during testing and harming the training-testing consistency. More discussions and results about the necessity of $C_{aux}$ is presented in Section 4.4.

**Multi-Key for Prompt Selection** A new multi-key mechanism is proposed for more accurate prompt selection and thus enhances the continual learning performance. The query feature is extracted from the pre-trained network, which (the feature vector corresponding to [class] token [7]) can be expressed as,

$$q = f_\theta(x)[0]. \qquad (10)$$

Query features extracted from the pre-trained network of the same class tend to be similar, while they may exhibit diversity across different classes within the same task. Therefore, relying on a single key to represent each task, as the previous prompt-based methods did, is not sufficient. The proposed multi-key mechanism employs multiple keys in a

prompt to map each task, resulting in more precise representations of the various categories within each task, as illustrated in Figure 3. Specifically, each class within a task is assigned a unique key, resulting in the same number of keys as classes. The cosine similarity is then employed to measure the discrepancy between a query and its corresponding key,

$$d_{i,j} = \cos(q, k_{i,j}), \tag{11}$$

$k_{i,j}$ indicates the key of the $j$th class within task $i$. The query retrieves the closest key and the corresponding prompt is selected,

$$i = \underset{i \in \{1:t\}, j \in \{1:|Y_i|\}}{\arg\max} d_{i,j}, \tag{12}$$

where $|Y_i|$ is the number of classes in task $i$. During training, the softmax cross-entropy is employed to maximize the similarity between the query and its corresponding key while minimizing other similarities,

$$\mathcal{L}_{mk} = -\log\left(\frac{e^{d_{t,y}(x)}}{\sum_{i \in \{1:t\}, j \in \{1:|Y_i|\}} e^{d_{i,j}(x)}}\right), \tag{13}$$

where $y$ represents the class label of the input image $x$ at current task $t$. The total loss function of Prompt Consistency Learning is,

$$\mathcal{L}_{PCL} = \mathcal{L}_{ce} + \mathcal{L}_{aux} + \mathcal{L}_{mk}. \tag{14}$$

Finally, the overall learning objective of the proposed CPrompt is,

$$\mathcal{L} = \mathcal{L}_{CCL} + \mathcal{L}_{PCL}. \tag{15}$$

## 4. Experiments

### 4.1. Experimental Details

**Datasets and Protocols** Extensive experiments on four benchmark datasets are conducted for thorough comparisons among different continual learning methods. Following the class incremental setting, all test samples are without task identity.
- **StanfordCars** [16] is a fine-grained car dataset comprising 196 classes and 16,185 images. There are 8,144 training images and the rest for testing. All classes are randomly divided into the 10-task (20 classes per task) and 20-task (10 classes per task) continual learning setting, denoted Split StanfordCars. It poses a challenge due to the distinct image styles and the difficulty in distinguishing between fine-grained classes.
- **ImageNet-R** [11] has 30,000 images of 200 ImageNet classes. The images of each class exhibit various styles, including art, cartoons, Deviant-Art, graffiti, and hard examples sourced from the original ImageNet dataset. This

benchmark is challenging because the various styles significantly differ from the pre-training data. The continual learning benchmark, Split ImageNet-R, consists of the 10-task (20 classes per task) and 20-task (10 classes per task) settings.
- **DomainNet** [31] dataset is a cross-domain dataset, including 345 common objects from 6 diverse domains, including Clipart, Real, Sketch, Infograph, Painting, and Quickdraw. Due to the significant disparity in image counts of different classes, the 200 categories with the most images are selected. Its continual learning setting consists of 10 tasks (10 classes per task). Different classes are distributed into different tasks at random. Moreover, each task comprises images from multiple domains rather than a single one as in the existing split [44].
- **CIFAR-100** [17] consists of 60,000 32 × 32 colour images of 100 classes. There are 500 training images and 100 test images per class. It has a widely adopted continual learning benchmark with 10 tasks (10 classes per task).

**Evaluation Metrics** The continual learning performance of classification models is mainly evaluated by two metrics [42]: the average accuracy of all classes after learning the last task (denoted as Last-acc) and the averaged incremental accuracy over all learned tasks (denoted as Avg-acc). The average forgetting [42], denoted as FF, provides additional context about the performance drops over tasks. We give more emphasis to Last-acc and Avg-acc as they are more comprehensive metrics [36].

**Implementation Details** Our model training employs the SGD optimizer with a momentum of 0.9 and an initial learning rate of 0.01. The initial learning rate gradually diminishes to zero following a cosine annealing scheduler. The mini-batch size is 16. In practice, the balancing hyper-parameters of different losses are all fixed at 1 without tuning. The hyper-parameters $(\tau_1, m)$ in Eq. (4) are $(1.02, 0.05)$ for DomainNet, $(1.2, 0)$ for CIFAR-100 and $(1.15, 0.1)$ otherwise. These values are chosen via cross-validation. Following the previous setting [45, 46], the ImageNet pre-trained ViT-B/16 [7] is used as the backbone. Moreover, the ViT-B/16 architecture encompasses 12 self-attention layers. We divided each task-specific prompt into two segments and inserted them into the first and the middle 6th self-attention layers, respectively. All results are averaged over 3 runs with the corresponding standard deviations reported to mitigate the influence of random factors.

**Competitors** We focus on comparing our model with the state-of-the-art prompt-based methods, i.e., L2P [46], DualPrompt [45], ESN [44] and CODAprompt [36]. For fair comparisons, their best results are also reproduced under our experimental settings. The upper-bound (UB) performance is achieved by fine-tuning the prompt and classifier with all task data collectively.

Table 1. The continual learning results on Split StanfordCars under 10-task and 20-task settings. Best results are marked in **bold**.

| Method | 10-task | | | 20-task | | |
|--------|---------|---|---|---------|---|---|
| | Last-acc↑ | Avg-acc ↑ | FF ↓ | Last-acc ↑ | Avg-acc ↑ | FF ↓ |
| UB | 83.96 | - | - | 83.96 | - | - |
| L2P | 60.39±1.99 | 71.92±1.12 | **13.00±0.12** | 45.14±4.33 | 60.03±3.56 | **15.22±1.28** |
| DualPrompt | 57.27±0.34 | 70.36±2.33 | 16.31±0.97 | 43.99±1.55 | 60.22±0.86 | 18.25±1.45 |
| ESN | 56.91±0.56 | 72.82±0.79 | 13.50±1.64 | 46.53±2.02 | 62.54±1.81 | 15.96±0.76 |
| CODAprompt | 62.24±0.14 | 73.28±0.93 | 15.08±0.89 | 48.94±1.77 | 63.78±1.46 | 17.38±2.16 |
| Ours | **66.77±0.37** | **76.81±0.27** | 13.95±0.46 | **55.16±0.19** | **68.74±0.83** | 20.02±1.78 |

Table 2. The continual learning results on Split ImageNet-R under 10-task and 20-task settings. Best results are marked in **bold**.

| Method | 10-task | | | 20-task | | |
|--------|---------|---|---|---------|---|---|
| | Last-acc ↑ | Avg-acc ↑ | FF ↓ | Last-acc ↑ | Avg-acc ↑ | FF ↓ |
| UB | 80.27 | - | - | 80.27 | - | - |
| L2P | 74.60±0.90 | 80.83±1.39 | **4.52±0.50** | 72.09±1.12 | 78.39±0.94 | **4.86±1.37** |
| DualPrompt | 74.87±0.85 | 81.39±1.25 | 7.18±0.15 | 71.69±1.06 | 79.12±1.27 | 7.68±0.96 |
| ESN | 75.11±0.36 | 81.63±1.10 | 5.68±0.77 | 70.57±0.62 | 77.95±0.76 | 6.84±0.36 |
| CODAprompt | 75.51±0.81 | 81.32±1.01 | 5.91±1.36 | 72.25±0.78 | 78.07±0.40 | 6.65±0.31 |
| Ours | **77.14±0.11** | **82.92±0.70** | 5.97±0.68 | **74.79±0.28** | **81.46±0.93** | 7.34±0.65 |

Table 3. The continual learning results on Split CIFAR-100 under the 10-task setting. Best results are marked in **bold**.

| Method | Last-acc ↑ | Avg-acc ↑ | FF ↓ |
|--------|-----------|-----------|------|
| UB | 91.99 | - | - |
| L2P | 86.38±0.31 | 91.45±0.19 | 5.88±0.76 |
| DualPrompt | 86.61±0.22 | 90.82±1.47 | 5.86±0.62 |
| ESN | 86.42±0.80 | 91.65±0.67 | 6.08±0.48 |
| CODAprompt | 85.73±0.14 | 91.03±0.57 | 7.13±0.44 |
| Ours | **87.82±0.21** | **92.53±0.23** | **5.06±0.50** |

Table 4. The continual learning results on Split DomainNet under the 10-task setting. Best results are marked in **bold**.

| Method | Last-acc ↑ | Avg-acc ↑ | FF ↓ |
|--------|-----------|-----------|------|
| UB | 89.15 | - | - |
| L2P | 81.17±0.83 | 87.43±0.95 | 8.98±1.25 |
| DualPrompt | 81.70±0.78 | 87.80±0.99 | 8.04±0.31 |
| ESN | 79.22±2.04 | 86.69±1.18 | 10.62±2.12 |
| CODAprompt | 80.04±0.79 | 86.27±0.82 | 10.16±0.35 |
| Ours | **82.97±0.34** | **88.54±0.41** | **7.45±0.93** |

## 4.2. Main Results

Comparisons between our Consistent Prompting (CPrompt) with its SOTA competitors on the Split StanfordCars benchmarks are shown in Table 1. The Split StanfordCars benchmark is a challenging task for continual learning methods, as their results are far behind the upper bound of joint learning and non-neglectable forgetting occurs. However, the proposed CPrompt achieves the best performance among its counterparts. Specifically, CPrompt outperforms the SOTA CODAprompt significantly, achieving 4.53% and 3.53% higher Last-acc and Avg-acc respectively under the 10-task setting. Such improvements become more apparent with additional tasks. Under the 20-task setting, CPrompt achieves 6.22% higher Last-acc and 4.96% higher Avg-acc than those of CODAprompt.

Our CPrompt and SOTA methods are also compared in Split ImageNet-R, as shown in Table 2. Under the 10-task setting, our method still outperforms CODAprompt, showing about 1.60% improvements in both Last-acc and

Avg-acc. Larger improvements, 2.54% higher Last-acc and 3.33% higher Avg-acc, achieved by CPrompt are also observed under the more challenging 20-task setting, similar to the trends above.

Table 3 and Table 4 present the results of Split CIFAR-100 and Split DomainNet both under 10-task continual learning. The proposed CPrompt consistently outperforms other prompt-based methods on both benchmarks under all criteria. Furthermore, CPrompt's improvements are non-trivial, mostly resulting in a greater than 1% increase in Last-acc and Avg-acc compared to all competitors.

Detailed comparisons between different methods along the continual learning procedure are illustrated in Figure 4. It demonstrates the superiority of our method. The curves of CPrompt are consistently above those of its counterparts across different tasks. Typically, these performance gaps tend to widen as tasks become more complex and involve greater numbers of classes. This finding suggests that CPrompt can be more resistant to catastrophic forgetting
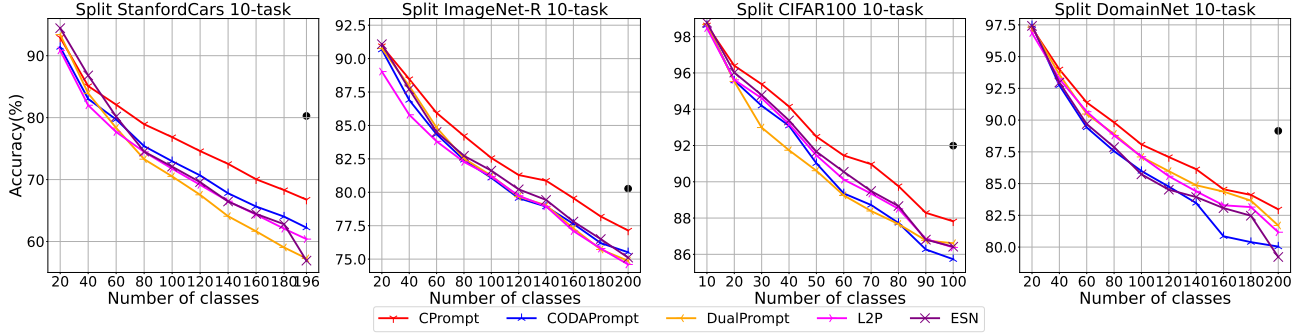
Figure 4. Illustrations of continual learning performance at each task. Each dot indicates the accuracy of the seen classes. The result of the upper-bound (UB) is represented by a dot on the overall classes.

Table 5. Ablation study of the proposed CPrompt with three components: Classifier Consistency Learning (CCL), Prompt Consistency Learning without Multi-Key mechanism (PCL w/o MK) , and Multi-Key (MK) mechanism.

| CCL | PCL w/o MK | MK | Split StanfordCars | | | | Split ImageNet-R | | | |
| | | | 10-task | | 20-task | | 10-task | | 20-task | |
| | | | Last-acc ↑ | Avg-acc ↑ | Last-acc ↑ | Avg-acc ↑ | Last-acc ↑ | Avg-acc ↑ | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 61.96 | 73.10 | 49.23 | 64.10 | 74.75 | 81.19 | 72.44 | 79.80 |
| ✔ | | | 62.52 | 73.42 | 51.89 | 66.11 | 76.37 | 82.19 | 73.63 | 80.46 |
| | ✔ | | 65.18 | 75.82 | 52.55 | 67.96 | 75.93 | 81.93 | 73.83 | 81.02 |
| | | ✔ | 63.66 | 74.26 | 51.65 | 65.91 | 75.36 | 82.01 | 73.14 | 80.22 |
| ✔ | ✔ | | 66.37 | 76.58 | 53.44 | 67.16 | 76.53 | 82.61 | 74.51 | 81.43 |
| ✔ | ✔ | ✔ | 66.77 | 76.81 | 55.16 | 68.74 | 77.14 | 82.92 | 74.79 | 81.46 |

Table 6. Detail analysis of CCL on 10-task continual learning of Split StanfordCars. Reported results are obtained by matching each input with the corresponding prompt. w/o CCL means CPrompt without CCL.

| | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|
| w/o CCL | 67.78 | 77.27 |
| CPrompt | 71.14+3.36 | 80.44+3.17 |

compared to its SOTA competitors.

## 4.3. Ablation Study

The proposed CPrompt consists of two main components: Classifier Consistency Learning (CCL) and Prompt Consistency Learning (PCL) including the Multi-Key (MK) mechanism. The effectiveness of each component is evaluated through ablative experiments on two large-scale continual learning benchmarks, Split StanfordCars and Split ImageNet-R. The results are presented in Table 5. The performance of our full CPrompt (with all three components) is significantly better than that of the vanilla backbone (results in the first row). For example, CPrompt achieves $4.81\%$ higher Last-acc and $3.71\%$ higher Avg-acc than the backbone under the Split StanfordCars 10-task setting. Furthermore, including each component alone can typically result in substantial improvements, exceeding $1\%$, in both criteria.

Table 7. Results of DualPrompt with CCL on 10-task continual learning of Split StanfordCars.

| | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|
| DualPrompt | 57.27 | 70.36 |
| +CCL | 61.79+4.52 | 72.07+1.71 |

Therefore, the effectiveness of each proposed component is demonstrated. Last but not least, combining CCL and PCL w/o MK improves the performance of each individual component. Such a combination clearly boosts the performance of all settings to the second best. It demonstrates that these two components serve distinctive purposes and complement each other.

## 4.4. Detailed Analysis

**Detailed Analysis of CCL** During the training phase, CCL is proposed to regulate the behavior of all classifiers, aiming to achieve more consistency across training and testing. In order to demonstrate the advantages exclusively derived from the consistent classifier, the corresponding prompt of each input is provided during testing. Table 6 presents the corresponding comparison between the CPrompts trained with and without CCL. Our full approach (with CCL) achieves significantly better results, $3.36\%$ higher Last-acc and $3.17\%$ higher Avg-acc, than the

Table 8. Detail analysis of PCL on 10-task continual learning of Split StanfordCars. With the prompt selection fixed at the initial task one, the reported results are averaged from the second to the final task. w/o PCL means CPrompt without PCL.

|  | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|
| w/o PCL | 48.83 | 63.17 |
| CPrompt | 59.75 +10.92 | 70.71 +7.54 |

Table 9. Results of DualPrompt with PCL on 10-task continual learning of Split StanfordCars.

|  | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|
| DualPrompt | 57.27 | 70.36 |
| +PCL | 60.09 +2.82 | 72.12 +1.76 |

CPrompt without CCL. Furthermore, CCL can be applied to other prompt-based continual learning methods and brings them significant improvements, e.g., $4.52\%$ Last-acc, to DualPrompt, as demonstrated in Table 7. More results of CODAPrompt and L2P with our approach are in the supplementary material.

**Detailed Analysis of PCL** To demonstrate the benefits of robust prediction brought by PCL, a new experiment procedure is developed. During testing, we consistently choose the prompt of the initial task, which does not correspond to any of the subsequent tasks. Using this prompt, we calculate the model's performance from the second task to the last one and then take the average. This setting helps to showcase the model's robustness in dealing with mismatched prompts, where existing methods may have difficulties. Table 8 presents the corresponding comparison between the CPrompts trained with and without PCL. Our full CPrompt trained with PCL outperforms its counterpart with large margins, $10.92\%$ and $7.54\%$ on Last-acc and Avg-acc, respectively. This clearly demonstrates the necessity of PCL in the proposed CPrompt. Moreover, combining PCL with the DualPrompt method yields significant improvements, as demonstrated in Table 9.

**Detailed Analysis of MK** The proposed multi-key (MK) mechanism assigns multiple keys to each task prompt to handle the inherent diversity of different classes. This approach leads to higher accuracy, i.e., for more than $20\%$, in selecting task-specific prompts, thus improving continual learning performance compared to the vanilla one-key (OK) mapping per task prompt, as illustrated in Table 10. The importance of accurate prompt selection can also be demonstrated by the upper-bound performance under the setting of task incremental learning (TIL). Appropriate input prompts can always be selected ($100\%$ Task-acc) with the task identities provided in TIL. As a result, its performance is clearly better than the CIL counterparts.

**Detailed Analysis of auxiliary classifier $C_{aux}$** The im-

Table 10. Detail analysis of the Multi-Key (MK) mechanism on 10-task continual learning of Split StanfordCars. The experiment is conducted on the vanilla backbone network with the one-key (OK) mechanism by default. Task-acc refers to the accuracy of selecting an appropriate prompt.

|  | Task-acc ↑ | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|---|
| TIL | 100 | 67.50 | 77.22 |
| OK | 19.44 | 61.96 | 73.10 |
| MK | 39.83 +20.39 | 63.66 +1.70 | 74.26 +1.16 |

Table 11. Detail analysis of auxiliary classifier $C_{aux}$ on 10-task continual learning of Split StanfordCars.

|  | Last-acc ↑ | Avg-acc ↑ |
|---|---|---|
| w/o $\mathcal{L}_{aux}$ | 60.34 | 73.17 |
| w/o $C_{aux}$ | 64.17 | 74.32 |
| CPrompt | 66.77 | 76.81 |

pacts of auxiliary classifier $C_{aux}$ and its corresponding loss $\mathcal{L}_{aux}$ in Eq. (9) are evaluated in Table 11. w/o $\mathcal{L}_{aux}$ refers to training CPrompt without auxiliary classifier. w/o $C_{aux}$ indicates learning $C_t$ instead of $C_{aux}$ with the auxiliary loss $\mathcal{L}_{aux}$. Both variants are inferior to CPrompt, demonstrating the effectiveness of the auxiliary classifier $C_{aux}$ and its objective $\mathcal{L}_{aux}$, as shown in Table 11.

# 5. Conclusion

In this paper, the training-testing inconsistency of existing prompt-based continual learning methods is first revealed and thoroughly discussed. We propose a novel approach, consistent prompting (CPrompt), to handle this important issue. Our CPrompt consists of two complementary components: classifier consistency learning (CCL) and prompt consistency learning (PCL). CCL addresses the classifier inconsistency by training with all classifiers, while PCL enhances the prediction robustness and prompt selection accuracy to alleviate the prompt inconsistency. The effectiveness of the proposed method is demonstrated by its state-of-the-art performance. Extensive analysis is also conducted to show the importance of training-testing consistency in prompt-based methods. **Limitations.** The multi-key mechanism is used to improve prompt selection accuracy during testing. It improves prompt selection accuracy by over $20\%$, as shown in Table 10. However, it still has a relatively low accuracy, i.e., less than $40\%$, achieved. Further improving the prompt selection accuracy can significantly boost the continual learning performance, as suggested by TIL's superior performance in Table 10. We will further enhance the proposed consistent prompting in this direction.

# References

[1] 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)*, pages 139–154, 2018. 2

[2] Eden Belouadah, Adrian Popescu, and Ioannis Kanellos. A comprehensive study of class incremental learning algorithms for visual tasks. *Neural Networks*, 135:38–54, 2021. 1

[3] Lorenzo Bonicelli, Matteo Boschini, Angelo Porrello, Concetto Spampinato, and Simone Calderara. On the effectiveness of lipschitz-driven rehearsal in continual learning. *Advances in Neural Information Processing Systems*, 35: 31886–31901, 2022. 1, 2

[4] Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020. 2

[5] Xiuwei Chen and Xiaobin Chang. Dynamic residual classifier for class incremental learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18743–18752, 2023. 2

[6] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE transactions on pattern analysis and machine intelligence*, 44(7):3366–3385, 2021. 1

[7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 3, 4, 5

[8] Tao Feng, Mang Wang, and Hangjie Yuan. Overcoming catastrophic forgetting in incremental object detection via elastic response distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9427–9436, 2022. 2

[9] Qiankun Gao, Chen Zhao, Bernard Ghanem, and Jian Zhang. R-dfcil: Relation-guided representation learning for data-free class incremental learning. In *European Conference on Computer Vision*, pages 423–439. Springer, 2022. 1

[10] Tyler L Hayes, Giri P Krishnan, Maxim Bazhenov, Hava T Siegelmann, Terrence J Sejnowski, and Christopher Kanan. Replay in deep learning: Current approaches and missing biological elements. *Neural computation*, 33(11):2908–2950, 2021. 2

[11] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8340–8349, 2021. 5

[12] Kishaan Jeeveswaran, Prashant Bhat, Bahram Zonooz, and Elahe Arani. Birt: Bio-inspired replay in vision transformers for continual learning. *arXiv preprint arXiv:2305.04769*, 2023. 2

[13] Zixuan Ke, Bing Liu, and Xingchang Huang. Continual learning of a mixed sequence of similar and dissimilar tasks. *Advances in Neural Information Processing Systems*, 33: 18493–18504, 2020. 2

[14] Muhammad Gul Zain Ali Khan, Muhammad Ferjad Naeem, Luc Van Gool, Didier Stricker, Federico Tombari, and Muhammad Zeshan Afzal. Introducing language guidance in prompt-based continual learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11463–11473, 2023. 2

[15] 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. 2

[16] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *Proceedings of the IEEE international conference on computer vision workshops*, pages 554–561, 2013. 5

[17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5

[18] Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming catastrophic forgetting by incremental moment matching. *Advances in neural information processing systems*, 30, 2017. 1, 2

[19] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021. 2

[20] Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. In *International Conference on Machine Learning*, pages 3925–3934. PMLR, 2019. 2

[21] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017. 2

[22] Huiwei Lin, Baoquan Zhang, Shanshan Feng, Xutao Li, and Yunming Ye. Pcr: Proxy-based contrastive replay for online class-incremental continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24246–24255, 2023. 2

[23] Huan Liu, Li Gu, Zhixiang Chi, Yang Wang, Yuanhao Yu, Jun Chen, and Jin Tang. Few-shot class-incremental learning via entropy-regularized data-free replay. In *European Conference on Computer Vision*, pages 146–162. Springer, 2022. 1

[24] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9): 1–35, 2023. 2

[25] Chunwei Ma, Zhanghexuan Ji, Ziyun Huang, Yan Shen, Mingchen Gao, and Jinhui Xu. Progressive voronoi diagram

subdivision enables accurate data-free class-incremental learning. In *The Eleventh International Conference on Learning Representations*, 2022. 1

[26] Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D Bagdanov, and Joost Van De Weijer. Class-incremental learning: survey and performance evaluation on image classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5513–5533, 2022. 1

[27] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, pages 109–165. Elsevier, 1989. 1, 2

[28] Sanket Vaibhav Mehta, Darshan Patil, Sarath Chandar, and Emma Strubell. An empirical investigation of the role of pre-training in lifelong learning. *arXiv preprint arXiv:2112.09153*, 2021. 1

[29] Jun-Yeong Moon, Keon-Hee Park, Jung Uk Kim, and Gyeong-Moon Park. Online class incremental learning on stochastic blurry task boundary via mask and visual prompt tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11731–11741, 2023. 2

[30] Yixuan Pei, Zhiwu Qing, Shiwei Zhang, Xiang Wang, Yingya Zhang, Deli Zhao, and Xueming Qian. Space-time prompting for video class-incremental learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11932–11942, 2023. 1

[31] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1406–1415, 2019. 5

[32] Vinay Venkatesh Ramasesh, Aitor Lewkowycz, and Ethan Dyer. Effect of scale on catastrophic forgetting in neural networks. In *International Conference on Learning Representations*, 2021. 1

[33] Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi. Progressive prompts: Continual learning for language models. *arXiv preprint arXiv:2301.12314*, 2023. 2

[34] 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*, pages 2001–2010, 2017. 1

[35] James Smith, Yen-Chang Hsu, Jonathan Balloch, Yilin Shen, Hongxia Jin, and Zsolt Kira. Always be dreaming: A new approach for data-free class-incremental learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9374–9384, 2021. 1

[36] James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11909–11919, 2023. 1, 2, 5

[37] Shagun Sodhani, Mojtaba Faramarzi, Sanket Vaibhav Mehta, Pranshu Malviya, Mohamed Abdelsalam, Janarthanan Janarthanan, and Sarath Chandar. An introduction to lifelong supervised learning. *arXiv preprint arXiv:2207.04354*, 2022. 1

[38] Yu-Ming Tang, Yi-Xing Peng, and Wei-Shi Zheng. When prompt-based incremental learning does not meet strong pretraining. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1706–1716, 2023. 1

[39] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019. 1

[40] Fu-Yun Wang, Da-Wei Zhou, Liu Liu, Han-Jia Ye, Yatao Bian, De-Chuan Zhan, and Peilin Zhao. Beef: Bi-compatible class-incremental learning via energy-based expansion and fusion. In *The Eleventh International Conference on Learning Representations*, 2022. 2

[41] Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *European conference on computer vision*, pages 398–414, 2022. 2

[42] Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. A comprehensive survey of continual learning: Theory, method and application. *arXiv preprint arXiv:2302.00487*, 2023. 5

[43] Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical decomposition of prompt-based continual learning: Rethinking obscured suboptimality. *Advances in Neural Information Processing Systems*, 36, 2024. 1

[44] Yabin Wang, Zhiheng Ma, Zhiwu Huang, Yaowei Wang, Zhou Su, and Xiaopeng Hong. Isolation and impartial aggregation: A paradigm of incremental learning without interference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 10209–10217, 2023. 3, 5

[45] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European Conference on Computer Vision*, pages 631–648. Springer, 2022. 1, 2, 5

[46] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022. 1, 2, 5

[47] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 374–382, 2019. 2

[48] 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*, pages 3014–3023, 2021. 2

[49] Jaehong Yoon, Divyam Madaan, Eunho Yang, and Sung Ju

Hwang. Online coreset selection for rehearsal-based continual learning. *arXiv preprint arXiv:2106.01085*, 2021. 1

[50] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International conference on machine learning*, pages 3987–3995. PMLR, 2017. 2

[51] Gengwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learner with classifier alignment for continual learning on a pre-trained model. *arXiv preprint arXiv:2303.05118*, 2023. 1