

# Mind the Gap: Preserving and Compensating for the Modality Gap in CLIP-Based Continual Learning

Linlan Huang<sup>1</sup>, Xusheng Cao<sup>1</sup>, Haori Lu<sup>1</sup>, Yifan Meng<sup>1</sup>, Fei Yang<sup>2,1</sup>, Xialei Liu<sup>2,1\*</sup>

<sup>1</sup>VCIP, CS, Nankai University      <sup>2</sup>NKIARI, Shenzhen Futian

{huanglinlan, caoxusheng, luhaori}@mail.nankai.edu.cn, {feiyang, xialei}@nankai.edu.cn

## Abstract

Continual learning aims to enable models to learn sequentially from continuously incoming data while retaining performance on previously learned tasks. With the Contrastive Language-Image Pre-trained model (CLIP) exhibiting strong capabilities across various downstream tasks, there has been growing interest in leveraging CLIP for continual learning in such scenarios. Most existing works overlook the inherent modality gap in CLIP, a key factor in its generalization and adaptability. In this paper, we analyze the variations in the modality gap during the fine-tuning of vision-language pre-trained models. Our observations reveal that the modality gap effectively reflects the extent to which pre-trained knowledge is preserved. Based on these insights, we propose a simple yet effective method, **MG-CLIP**, that improves CLIP’s performance in class-incremental learning. Our approach leverages modality gap preservation to mitigate forgetting and modality gap compensation to enhance the capacity for new data, introducing a novel modality-gap-based perspective for continual learning. Extensive experiments on multiple benchmarks demonstrate that our method outperforms existing approaches without requiring additional replay data. Our code is available at <https://github.com/linlany/MindtheGap>.

## 1. Introduction

The goal of continual learning is to enable models to continuously acquire new knowledge and adapt to the ever-changing real world [22, 31]. Traditional approaches to continual learning focus on training a model from scratch, aiming to reduce catastrophic forgetting of old task knowledge while maintaining plasticity to adapt to new data [18]. This is often referred to as the stability-plasticity trade-off. However, with the recent development of large-scale pre-trained models, these models offer stronger stability

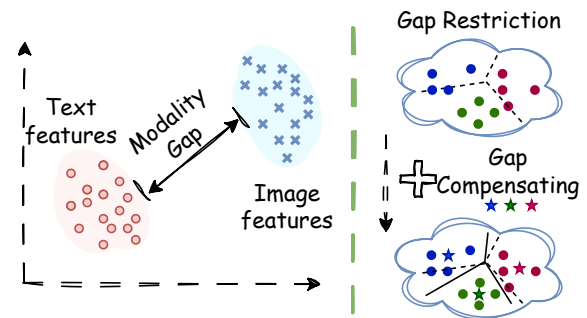


Figure 1. Left: In continual learning, we preserve the modality gap to retain CLIP’s original knowledge and mitigate forgetting. Right: Due to the modality gap, the text classifier’s performance is restricted. We compensate for it through an intra-modal classifier by enhancing its adaptability and refining decision boundaries.

and generalization capabilities for continual learning [15, 23, 36]. Among pre-trained models, visual language pre-trained models such as CLIP demonstrate excellent generalization to downstream tasks, showcasing impressive zero-shot capabilities [16, 17, 26]. Consequently, CLIP-based continual learning (CL) has emerged as a promising new direction, attracting increasing attention from researchers [11, 12, 30, 40].

Existing CLIP-based continual learning methods can be broadly categorized into two approaches: one fine-tunes the backbone to modify feature representations [21, 40], while the other freezes the backbone and introduces learnable modules for continual learning [11, 12, 49]. They often treat CLIP as a feature extractor, approximating it to a text-enhanced visual model [21]. Based on this view, they focus on better feature fusion for improved performance [12, 49] or leveraging textual information to guide visual feature adaptation [11]. However, these methods often overlook CLIP’s unique cross-modal properties, the modality gap, distinguishing it from unimodal systems. The modality gap is observed in previous work [19]. As shown in Fig. 1 (left), the modality gap is that the intra-modal feature distance is small while the inter-modal distance is large. The feature

\*Corresponding author.

distributions of the two modalities are represented in two distinct cones, with an inherent spacing between them. To fully exploit CLIP’s cross-modal properties, it is crucial to explore strategies that preserve the inherent modality gap while addressing its limitations in continual learning.

In this paper, we investigate CLIP-based class-incremental learning from the perspective of the modality gap. In existing studies on multimodal model training and adaptation of downstream tasks, the modality gap is often seen as a source of suboptimal performance. Researchers focus on reducing the modality gap in pre-trained multimodal models [6, 10, 24]. However, in the context of continual learning with CLIP, our objective is to retain the strong generalization capability of the model while learning new data. Thus, how to handle the modality gap in continual learning remains an open question. Our work is based on the assumption that this modality gap reflects the intrinsic knowledge of the pre-trained model.

Downstream datasets are typically much smaller than pretraining datasets, making them insufficient to fully train the model. As a result, modifying the fundamental properties of a pre-trained model, such as the modality gap, based on the downstream data may disrupt the pre-trained knowledge. It may potentially compromise the performance of continual learning. During training on the current task, the cross-entropy optimization objective tends to expand this gap. This phenomenon intensifies as incremental learning data continue to arrive. Furthermore, directly aligning the two modalities for downstream tasks will also significantly alter the pre-trained knowledge. Therefore, we keep the modality gap relatively stable to maintain the stability of the model. We analyze the changes in modality gap during training and propose a modality gap-aware adjustment strategy. By tracking variations in the modality gap, we regulate the training to maintain a stable modality gap, thereby preserving pre-trained knowledge and mitigating forgetting.

Additionally, we analyze the impact of preserving the modality gap. As shown in Fig. 1 (right), when text features are used as classifiers, the modality gap may restrict the model’s ability to learn new data in continual learning. This limits the adaptability of the model, reducing the plasticity. To compensate modality gap, we propose to build a classifier in the visual space, where the modality gap does not pose a restriction. By integrating its output with that of the text classifier, we compensate for the modality gap and improve the learning capacity of CLIP.

The main contributions of this paper are as follows:

- We propose a method to preserve the modality gap that maintains pre-trained knowledge and mitigates forgetting in continual learning.
- We compensate for the limitations of the modality gap by introducing a complementary mechanism, enhancing model plasticity.

- Our approach achieves state-of-the-art results across multiple datasets without replay or complex structures.

## 2. Related Work

### 2.1. Class-Incremental Learning

Class-incremental learning methods are generally grouped into three types [3]. Regularization methods limit changes to the model. Some methods use distillation to reduce the deviation of model features or penalize changes in model parameters [2, 13, 18, 44, 47]. Dynamic network approaches allow the model structure to evolve as new tasks are introduced [5, 32, 37, 38]. Replay-based methods retain original samples or relevant information from them. When learning a new task, these methods replay original samples [20, 27, 29] or recover old samples from the retained information [11, 33, 41, 43]. The latter often involves techniques such as memory compression or feature replay. Dynamic networks or feature replay methods tend to increase both parameter count and memory storage requirements. With the widespread use of pre-trained models, incremental learning methods for visual pre-trained models have also emerged. Parameter-efficient fine-tuning approaches gradually expand a small number of parameters as tasks increase [28, 34, 35]. APER [48] maintains stability by training the model only on the first task. SLCA [43] uses a small learning rate fine-tuning backbone for continual learning. They show that preserving the stability of pre-trained visual models enhances generalization and benefits class-incremental learning on downstream tasks.

### 2.2. Class-Incremental Learning with CLIP

CLIP performs remarkably well in class-incremental learning tasks with downstream data [30]. Continual learning based on CLIP has gained increasing attention. Some approaches add learnable modules to the original CLIP features to modify them for better adaptation to new tasks. PROOF [49] and CLAP [12] add learnable modules to the CLIP output features to facilitate cross-modal interaction. RAPF [11] introduces a linear layer after the CLIP visual encoder for downstream adaptation and uses textual modality information to guide feature replay. Other methods fine-tune CLIP to alter its output features. ZSCL [46] distills additional datasets during training. MOE4CL [40] fine-tunes CLIP with a mixture of experts, introducing a selection mechanism to selectively use the original CLIP model. Magmax [21] sequentially fine-tunes the entire CLIP model, using a task vector algorithm to merge models at inference time. LGVLM [45] trains separate LoRA modules for each task. These methods primarily focus on leveraging the prior knowledge of natural language to assist continual learning, often neglecting the preservation of zero-shot capabilities while emphasizing performance on

sequential tasks. In contrast, we focus on retaining the model’s capabilities while enhancing its continual learning ability, starting from the inherent property of the CLIP.

### 2.3. Modality Gap

Liang *et al.* [19] demonstrated a phenomenon in pre-trained vision language models called the modality gap. In contrastive vision-language models, the modality gap is reflected in the fact that text and image features are located in two separate narrow cones. The features of different modalities exhibit distinct separation, whereas the features within the same modality tend to cluster more closely. Existing studies have explored the impact of the modality gap on tasks such as domain adaptation [10], few-shot classification [39], model pre-training [7] and retrieval task [24]. Reducing the modality gap can improve the model performance in these tasks. However, in continual learning tasks, the stability of pre-trained knowledge is more crucial for downstream tasks over time. Therefore, we propose using the modality gap as an indicator of the changes in the pre-trained model’s feature space, with the goal of preserving the modality gap in continual tasks.

### 3. Preliminaries

**Class incremental learning definition.** We consider a class-incremental learning setup based on a pre-trained CLIP model  $M$ . The objective is to sequentially train the model on a series of classification tasks. Each task  $t$  consists of a set of categories  $C_t$ , and the categories between tasks do not overlap,  $\forall i \neq j, C_i \cap C_j = \emptyset$ . During the training of task  $t$ , the model has no access to information related to the previous  $t - 1$  tasks. After training in the  $t$ -th task, the model  $M_t$  is required to correctly classify all previously learned categories,  $C_1 \cup C_2 \cup \dots \cup C_t$ , without access to task identifiers.

**Modality gap measure.** In the classification task, given  $N$  images and  $K$  class name of text, we measure the inter-modality similarity as the average cosine similarity between all image and text features:  $\frac{1}{N} \sum_{i=1}^N \frac{1}{K} \sum_{j=1}^K \cos(\mathbf{x}_i, \mathbf{t}_j)$ , where  $\mathbf{x}_i$  is an image feature and  $\mathbf{t}_j$  is a text feature. This captures the overall similarity between the image and text modalities, reflecting the modality gap in the feature space.

To analyze the impact of image-to-text similarity in more fine-grained, we define the average similarity of positive image-text pairs as:

$$pos = \frac{1}{N} \sum_i^N \cos(\mathbf{x}_i, \mathbf{t}_{y_i}), \quad (1)$$

where  $\mathbf{t}_{y_i}$  denotes the corresponding class text feature for the image  $\mathbf{x}_i$ . Similarly, we measure the average similarity

of negative image-text pairs as:

$$neg = \frac{1}{N} \sum_{i=1}^N \frac{1}{K-1} \sum_{j=1}^{K-1} \cos(\mathbf{x}_i, \mathbf{t}_j^{neg}), \quad (2)$$

where  $\mathbf{t}_j^{neg}$  represents the text features of non-matching classes for the image  $\mathbf{x}_i$ .

### 4. Method

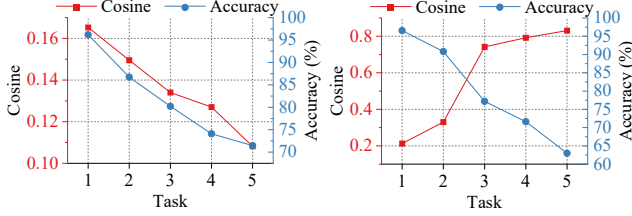
Our method approaches CLIP-based continual learning from the perspective of the modality gap, recognizing its crucial role in preserving CLIP’s generalization capabilities during continual learning. We identify that uncontrolled training can disrupt this inherent modality gap, leading to suboptimal performance in downstream incremental learning. To address this, we propose a two-stage strategy: (1) Preserving the modality gap for model stability. We regulate the training process to ensure a stable modality gap across tasks; (2) Compensating the modality gap for model plasticity. The constraint for the modality gap may limit the model’s adaptability, we introduce a complementary classifier that enhances task-specific performance without altering the preserved modality gap. During inference, we integrate both outputs to balance stability and task plasticity. Next, we discuss the motivation (Sec. 4.1) and details of our method (Sec. 4.2 and Sec. 4.3). The pseudo-code can be found in Sec.5 of the supplementary material.

#### 4.1. Effect of Modality Gap in Continual Learning

In this section, we begin by analyzing the evolution of the modality gap during continual learning.

**Cross-entropy loss expands the modality gap.** The optimization objective of cross-entropy is not consistent with the modality gap of the original CLIP. When optimized via cross-entropy loss, the training objective aligns image-text pairs to a cosine similarity of 1 for matches and -1 for non-matches. However, this optimization contradicts the moderate similarity distribution of the original CLIP model, which ranges from approximately 0 to 0.3 reflecting the inherent modality gap. It causes the modality gap of the trained model to expand. Furthermore, downstream classification tasks introduce a structural imbalance absent in CLIP pre-training. Unlike symmetric text-image pairs in pre-training, classification tasks pair multiple images of the same class with a single text embedding. In a  $C$ -class dataset, each text forms positive pairs with only  $\frac{1}{C}$  of the images while being repelled from the remaining  $\frac{C-1}{C}$  negative pairs. The imbalance leads to the optimization process dominated by repelling non-matching samples, further expanding the modality gap.

The modality gap implicitly reflects pre-trained knowledge, and expanding it may result in forgetting previous



(a) Naive fine-tuning reduces cosine (b) The alignment loss increases co-similarity, widens modality gap, and sine similarity and reduces modal-ity gap, but it disrupts pre-trained knowledge, leading to forgetting.

Figure 2. Average cosine similarity between text and image features, along with accuracy on ImageNet-R during CL.

knowledge. As shown in Fig. 2a, experimental results confirm this effect: the cosine similarity between image and text representations steadily declines as tasks progress, demonstrating the growing modality gap. This coincides with a continuous drop in accuracy, highlighting its adverse impact on CLIP’s pre-trained knowledge and overall performance in class-incremental learning.

**Direct alignment loss reduces modality gap but disrupts pre-trained knowledge.** As shown in Fig. 2b, introducing a loss of alignment [7] in downstream tasks can reduce the modality gap by minimizing the Euclidean distance between matched image and text features. The key issue lies in the mismatch of the subspace: downstream datasets span a limited region of the original CLIP feature space, which makes the direct alignment risk of overspecialization. Directly aligning the two modalities can significantly alter pre-trained representations, leading to forgetting.

**Sustaining CLIP’s lifelong learning capacity.** Previous analyses have revealed a phenomenon: naive fine-tuning tends to expand the modality gap, while direct alignment reduces it. As class-incremental learning progresses, this problem becomes increasingly severe, gradually eroding CLIP’s pretraining capabilities. Since CLIP’s pretraining knowledge plays a crucial role in providing stability for downstream tasks, we need to maintain a relatively stable modality gap. Ensuring this stability is essential for allowing CLIP to incorporate new knowledge without compromising its fundamental vision-language correspondence.

#### 4.2. Adaptive Modality Gap Preservation

In this section, we characterize the asymmetric evolution of the modality gap in continual learning and propose an adaptive preservation strategy to enhance model stability.

**Asymmetry phenomenon of modality gap change during training.** As shown in Fig. 3, we observe an asymmetric

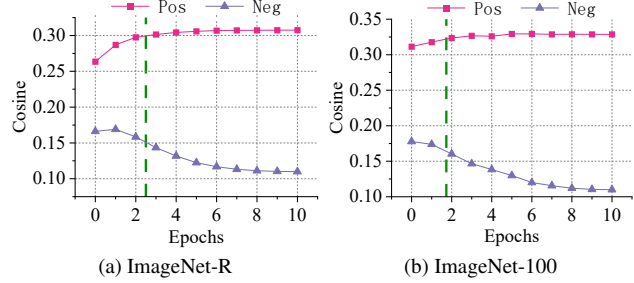


Figure 3. The variation in the mean of positive and negative cosine within a task during training. The cosine change pattern before and after the green dotted line is different.

change in cosine similarity between text and image features during single task training. The ‘pos’ and ‘neg’ use Eq. 1 and Eq. 2 to calculate. Before training, *i.e.*, epoch 0, the ‘pos’ is already larger than the ‘neg’, indicating the zero-shot capability of the model. Even the similarity between an image and its corresponding class text remains relatively low, which is a reflection of the modality gap. During the early phase of training, the output changes are primarily due to an increase in the similarity of positive image-text pairs, indicating the model has learned new knowledge and becomes more confident about the true class. This corresponds to the part before the green dashed line in Fig. 3. As long as the positive output remains larger than the negative outputs, the model can still make correct predictions even if the loss exists. However, in the later stages of training, the positive output stabilizes and the negative outputs begin to decrease. This indicates that the distance between the image and most other text features increases. Therefore, an optimal training stage can retain its learned knowledge while keeping the inter-modal distance relatively stable.

**Adaptive training for preserving modality gap.** Based on these observations, we propose determining the number of training epochs by monitoring the mean negative outputs. As shown in Fig. 4, using the first task data of the dataset, we estimate the required number of epochs for subsequent tasks. Specifically, we first compute the mean of the negative outputs using Eq. 2, denoted as  $neg^0$ , using the original CLIP model for the first task. Then, we train the model using LoRA, and after the  $e$  epoch, we compute the  $neg^e$  of the negative outputs for all data. We calculate the relative difference between the current and the original as follows:

$$\Delta = \frac{|neg^e - neg^0|}{neg^0}. \quad (3)$$

When the difference  $\Delta$  exceeds a predefined threshold  $\alpha$ , we record the last epoch  $e$  where  $\Delta$  was still below  $\alpha$ . Finally, for all tasks in the downstream dataset, we train the model for the  $\max(e, 1)$  epochs.

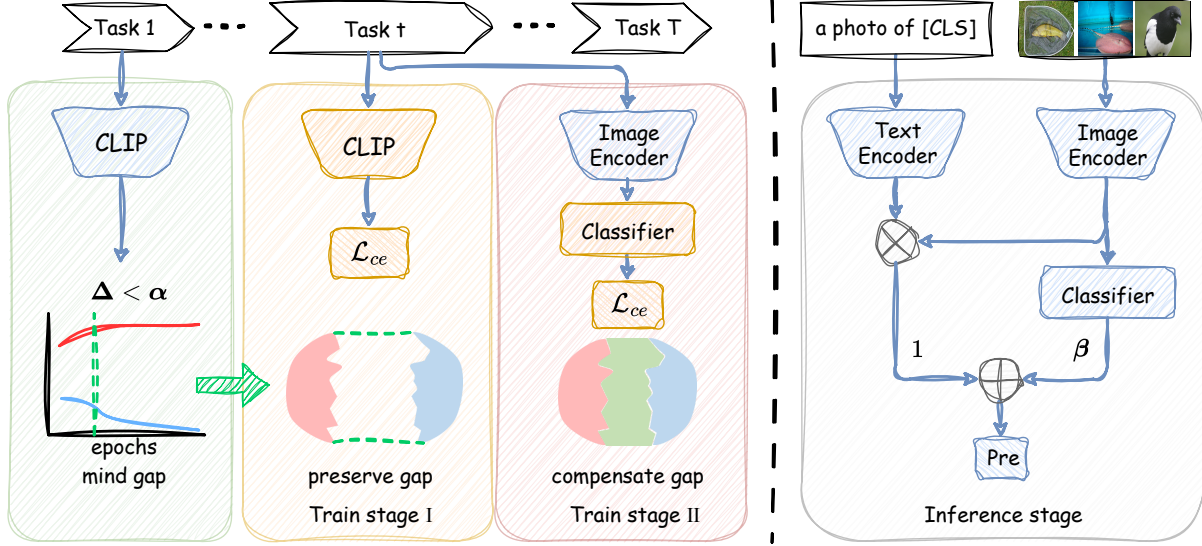


Figure 4. Our method tackles CLIP-based continual learning through dual mechanisms: (1) **Modality Gap Preservation** stops training when modality gap deviation exceeds a stability threshold, preventing cross-modal knowledge distortion; (2) **Gap Compensation** trains a visual-space classifier on frozen features to enhance task-specific plasticity while respecting preserved gap. In inference stage, we combine two different classifier subspaces for prediction.

### 4.3. Intra-modal Compensation for Modality Gap

In this section, we explain how the modality gap limits the model’s capability, as it must be maintained as discussed in the previous section. Consequently, we construct an intra-modal classifier to compensate for the limitations.

**Modal gap restriction on text classifier capability.** To achieve minimal cross-entropy loss, an optimal classifier must exist within the image feature space. Specifically, for any classifier that minimizes classification error, there exists an equivalent form  $\mathbf{W}_{\text{opt}}$  entirely contained in the span of image features. This follows from the decomposition of any classifier into components parallel and orthogonal to the image space, where the orthogonal component does not contribute to classification. A detailed proof of this claim is provided in the supplementary material 2.1.

Given this, we analyze the modality gap’s impact on text classifiers. The text feature matrix  $\mathbf{T}$  decomposes as:  $\mathbf{T} = \mathbf{T}_{\parallel} + \mathbf{T}_{\perp}$  where  $\mathbf{T}_{\parallel}$  lies in the subspace spanned by visual features  $\mathbf{X}$ , and  $\mathbf{T}_{\perp}$  is orthogonal to it. Since text features generally do not fully span the image feature space, the best achievable text classifier is restricted to a low-rank subspace, leading to misalignment error. The lower bound of the distance from  $\mathbf{T}_{\parallel}$  to the optimal image-space classifier is determined by the singular values  $s^2$  outside the text feature subspace:

$$\|\mathbf{T}_{\parallel} - \mathbf{W}_{\text{opt}}\|_F^2 \geq \sum_{i=r+1}^{r'} s_i^2, \quad (4)$$

where  $r'$  is the rank of  $\mathbf{W}_{\text{opt}}$  and  $r$  is the rank of  $\mathbf{T}_{\parallel}$ . The

formal derivation of this bound is provided in the supplementary material 2.2.

This result reveals an inherent limitation: unless the text feature space has sufficient capacity to represent the optimal classifier, perfect alignment is unattainable. Due to the modality gap, text classifiers often operate in a lower-rank subspace, restricting their classification effectiveness.

**Modal gap compensation via intra-modal classifier.** To compensate the modality gap, we introduce an auxiliary classifier in the visual space. The fine-tuned CLIP model  $f_{\text{clip}}(\cdot)$  and the classifier weights for old classes are frozen. For classes introduced in the current task, we initialize their classifier weights within the cosine classifier  $\mathbf{W}_v$  using their class prototypes and train them using image features without text. Since the gradients of the classifier remain in the input space, which is the visual space, this ensures that the classifier operates within the visual subspace.

As shown in Fig. 4, in the model inference stage, we combine the predictions from both the text and visual classifiers. The final predict score is calculated as:

$$\text{pre}(\mathbf{x}) = f_{\text{clip}}(\mathbf{x}, \mathbf{t}) + \beta \cdot \text{softmax}(\mathbf{W}_v^T \mathbf{x}), \quad (5)$$

where  $\beta$  is a constant hyperparameter.

## 5. Experiments

### 5.1. Experimental Setup

**Datasets.** We evaluate our method through continual learning tasks on five benchmark datasets: CIFAR-100 [14],

| Method              | Exemplar | CIFAR-100 |       | ImageNet-R |       | ImageNet100 |       | ImageNet-1K |       | VTAB  |       |
|---------------------|----------|-----------|-------|------------|-------|-------------|-------|-------------|-------|-------|-------|
|                     |          | Avg       | Last  | Avg        | Last  | Avg         | Last  | Avg         | Last  | Avg   | Last  |
| PROOF [49]          | DR       | 84.88     | 76.29 | 82.83      | 77.05 | 84.71       | 72.48 | 76.23       | 65.26 | 89.09 | 83.97 |
| CLAP [12]           |          | 86.13     | 78.21 | 85.77      | 79.98 | 87.76       | 79.16 | 81.72       | 73.19 | 91.37 | 89.67 |
| SLCA [43]           | FR       | 80.53     | 67.58 | 75.92      | 70.37 | 78.63       | 59.92 | 79.10       | 68.27 | 84.25 | 82.54 |
| RAPF [11]           |          | 86.19     | 79.04 | 85.58      | 80.28 | 87.51       | 80.23 | 81.73       | 72.58 | 90.88 | 82.31 |
| L2P++ [35]          | NR       | 81.90     | 73.08 | 81.67      | 75.98 | 80.51       | 67.22 | 79.30       | 69.60 | 63.23 | 38.37 |
| DualPrompt [34]     |          | 81.45     | 72.51 | 82.01      | 75.77 | 80.65       | 67.38 | 79.39       | 69.79 | 61.89 | 37.58 |
| CODA [28]           |          | 76.98     | 62.25 | 78.00      | 67.52 | 64.13       | 34.76 | 76.99       | 66.96 | 62.51 | 38.25 |
| Continual-CLIP [30] |          | 75.15     | 66.68 | 79.12      | 72.00 | 84.98       | 75.40 | 72.96       | 64.44 | 53.64 | 31.50 |
| Aper-Adapter [48]   |          | 75.76     | 65.50 | 78.65      | 71.35 | 85.84       | 76.40 | 76.60       | 68.74 | 80.75 | 71.21 |
| MOE4CL [40]         |          | 85.36     | 78.37 | 85.28      | 80.77 | 86.39       | 76.66 | 81.29       | 72.73 | 68.49 | 61.70 |
| CLAP* [12]          |          | 74.19     | 63.45 | 81.22      | 75.80 | 81.07       | 72.00 | 75.85       | 67.36 | 82.11 | 80.11 |
| MagMax [21]         |          | 85.63     | 79.00 | 87.13      | 80.85 | 86.33       | 75.92 | 80.74       | 71.31 | 64.63 | 53.90 |
| MG-CLIP (Ours)      |          | NR        | 87.00 | 80.57      | 87.58 | 82.67       | 87.31 | 78.38       | 81.88 | 73.68 | 94.67 |

Table 1. Comparison of performance with different methods. DR denotes using real data for replay, FR denotes generating old class features for replay, and NR denotes non-replay. The results are mainly obtained from references [11, 12] and reproduced using their publicly available code. The performance of ours is the average over three different class orders. Except for VTAB, which is divided into 5 tasks, the other datasets are divided into 10 tasks. CLAP\* represents the replay-free version provided by the paper of CLAP.

ImageNet-R [8], ImageNet-100 [4], ImageNet-1K [4], and VTAB [42]. All of these datasets except the VTAB are evenly divided into 10 sequential tasks. Following the previous work [48], we extracted a subset from VTAB, including 5 tasks, each with ten categories. For more experimental details, please refer to Sec.4 of the supplementary material.

**Competing methods.** Our experiments compare two types of methods: (1) Vision-only approaches that include L2P++ [35], DualPrompt [34], CODA [28], SLCA [43], and Aper-Adapter [48]; (2) CLIP-based methods that include PROOF [49], CLAP [12], RAPF [11], MOE4CL [40], MagMax [21] and the zero-shot CLIP baseline Continual-CLIP [30]. All methods adopt the ViT-B/16 weights of OpenAI [26] by default. While original MagMax implementations employ enhanced augmentations and optimized text templates, we follow previous works [11, 40, 40] and standardize evaluations using basic image augmentations and a fixed prompt template for fair comparison.

**Evaluation metrics.** The average accuracy on the test data after training the  $t$ -th task, considering the first through  $t$  tasks, is represented as  $A_t$ . ‘Avg’ is the average of the accuracies of all tasks, *i.e.*,  $\frac{1}{t} \sum_{i=1}^t A_i$ . ‘Last’ indicates the average accuracy after the final task  $T$ , *i.e.*,  $A_T$ .

**Implementation detail.** We develop our method using PyTorch on an A40 GPU. Unless otherwise specified, we use OpenAI’s pre-trained CLIP model, specifically the ViT-B/16 version. Results for another version of CLIP can be found in the supplementary materials. During the training of the backbone, we apply LoRA [9] for model adaptation, with the default rank set to 8. Since the focus of our work is

not on the fine-tuning process itself, we simplify the implementation by applying LoRA only to the key and value parts of the attention module. Different fine-tuning model implementations are left for future work. We use the Adam optimizer with a cosine learning rate scheduler with an initial learning rate of 0.001. The number of epochs is determined by the first part of the method, where the threshold  $\alpha$  is set to 10%. In the training stage for the image-space classifier, we train for 3 epochs, with an initial learning rate of 0.0005. The classifier uses a cosine classifier. During inference, the hyperparameter  $\beta$  to integrate the two classification results is set by default to 4. Hyperparameter effects are described in the supplementary material.

## 5.2. Comparison Results

Table 1 presents a comparison of our approach with other methods. In most datasets, our method outperforms all others, including those relying on replay. Specifically, on CIFAR-100, our method surpasses all competitors in Last accuracy by at least 1.53%. On ImageNet-R, we achieve an improvement of at least 1.82% in Last accuracy. On ImageNet-100, while our approach performs slightly worse than RAPF and CLAP, both of which use replay, we do not rely on replay and still outperform all non-replay methods by at least 1.72% in Last accuracy. Additionally, on ImageNet-100, CLAP\* performs significantly worse than CLAP, indicating that much of CLAP’s performance gain stems from data replay. On the larger-scale ImageNet-1K dataset, our method achieves comparable or even slightly better results than the best-performing alternatives.

The VTAB dataset poses a significant challenge to CLIP, as evidenced by the zeroshot accuracy, *i.e.*, Continual.CLIP, which is only 31.5%. On this challenging benchmark, most methods suffer substantial performance drops. Our method outperforms the best replay-based CLAP approach by 1.86% in Last accuracy, significantly outperforming its non-replay version. It shows a clear advantages of our method in this challenging cross-domain scenario.

### 5.3. Comparison of the Zero-shot Capability

Continual learning aims to allow models to incrementally acquire new knowledge while retaining existing capabilities. CLIP models, compared to traditional pretrained vision models, demonstrate superior zero-shot generalization. Therefore, continual learning with CLIP should not only minimize forgetting on new downstream tasks but also preserve its original zero-shot abilities. This ensures CLIP based continual learning is more than just a method for task-specific initialization that leads to overfitting, a limitation of previous works in this area. These works often follow traditional evaluation protocols, neglecting the preservation of the model’s intrinsic abilities. We propose evaluating zero-shot generalization of CLIP model on independent benchmarks after fine-tuning to address this issue.

As shown in Tab. 2, we assess zero-shot performance on three standard datasets following continual learning on CIFAR-100. Our approach yields slight improvements over the original CLIP on Food101 [1] and Oxford Pets [25], which are clearly distinct from CIFAR-100, while all baselines show performance degradation. On the more similar ImageNet-1K, three methods outperform the original CLIP, with our method achieving the most substantial improvement (a 2.85% gain). These results suggest that our approach successfully preserves pre-trained knowledge while effectively incorporating new information.

Notably, replay-based methods (PROOF, RAPF, and CLAP) experience more significant performance declines in zero-shot tasks compared to non-replay methods. This indicates that replay mechanisms may lead to overfitting to downstream tasks, whereas non-replay methods must preserve original representations to prevent forgetting. This further highlights the potential negative impact of replay strategies on pre-trained models.

### 5.4. Training Cost Analysis.

Using CIFAR-100 as an example, we analyze the parameter overhead of different methods by comparing the additional learnable parameters introduced by our approach and other baselines, as illustrated in Fig. 5. Our method introduces only 0.54M additional trainable parameters. Although the RAPF method requires fewer learnable parameters, it still necessitates storing a covariance matrix for each class, leading to additional storage consumption proportional to the

|        | Food101 | Pets  | ImageNet-1K | Avg   |
|--------|---------|-------|-------------|-------|
| CLIP   | 85.14   | 87.6  | 64.44       | 79.06 |
| PROOF  | 9.75    | 22.81 | 11.18       | 14.58 |
| RAPF   | 17.18   | 28.56 | 15.2        | 20.31 |
| CLAP   | 80.82   | 74.68 | 56.04       | 70.51 |
| MOE4CL | 82.85   | 84.06 | 66.02       | 77.64 |
| MagMax | 81.67   | 85.96 | 66.44       | 78.02 |
| Ours   | 85.70   | 88.17 | 67.29       | 80.39 |

Table 2. Zero-shot performance of the model on different downstream datasets after completing all class-incremental learning tasks on CIFAR-100. CLIP refers to the original CLIP pre-trained model without any downstream task-specific training.

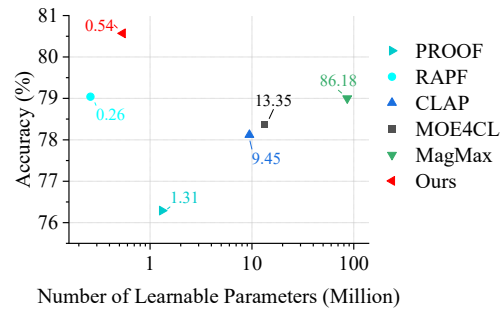


Figure 5. Comparison of accuracy and learnable parameters of different methods.

square of the feature dimension, *i.e.*,  $nd^2$ . For instance, with 100 classes, this results in an extra storage overhead of over 26M parameters. In contrast, our method avoids replay entirely, maintaining minimal storage consumption.

In the training process, the only extra cost in our method occurs after each epoch of the first task, where we perform an additional forward pass on the data of task to evaluate change of the modality gap. However, this step only needs to be done for the first task, does not involve backpropagation, and is minimal compared to the training cost.

### 5.5. Ablation Study

**Module ablation.** Table 3 presents the ablation study of our proposed modules. The baseline setup involves simple fine-tuning of the CLIP model for each task. Following the settings of previous work, we fine-tuned each task for 10 epochs [21] and used the widely adopted cross-entropy loss with the old-class output mask [40]. MGC represents Modality Gap Compensation, while MGP stands for Modality Gap Preservation. As shown in Tab. 3, both components of our method, when used individually, lead to performance improvements, with the modality gap preservation(MGP) providing a more significant boost. The best performance is achieved when both components are used together.

| MGP | MGC | Avg          | Last         |
|-----|-----|--------------|--------------|
| ✗   | ✗   | 84.30        | 72.74        |
| ✗   | ✓   | 85.10        | 74.58        |
| ✓   | ✗   | 86.73        | 76.86        |
| ✓   | ✓   | <b>87.31</b> | <b>78.38</b> |

Table 3. Ablation study of the modules on ImageNet-100. MGP refers to our adaptive modality gap preservation, and MGC denotes the intra-modal Compensation for Modality Gap.

When using MGC alone in the baseline setup, the improvement in last’ accuracy is 1.84%, which is larger than the improvement (1.52%) when adding MGC to MGP. A similar trend is observed for ‘Avg’ accuracy. This suggests that simple fine-tuning in the baseline setup has expanded the modality gap, and therefore the role of MGC in compensating for the gap becomes more pronounced.

**Image space and classifier space.** We investigate the relationships among the linear subspaces spanned by the image features, the text-based classifier, and the learned modality-gap-compensating classifier. To quantify their differences, we apply matrix decomposition to extract the orthonormal bases of each space. For detailed calculations, refer to the supplementary materials Sec.3.  $\mathbf{B}_i$ : Orthonormal basis of image feature space;  $\mathbf{B}_t$ : Orthonormal basis of text classifier;  $\mathbf{B}_{vc}$ : Orthonormal basis of the modality-gap-compensating classifier;  $\mathbf{B}_{t+vc}$ : Orthonormal basis of the combined space spanned by both classifiers. To measure how well the text classifier space  $\mathbf{B}_t$  covers the image feature space  $\mathbf{B}_i$ , we compute the mean norm of the orthogonal component of  $\mathbf{B}_i$  with respect to  $\mathbf{B}_t$ :

$$d(\mathbf{B}_i, \mathbf{B}_t) = \frac{1}{|\mathbf{B}_i|} \sum_{\mathbf{x} \in \mathbf{B}_i} \|\mathbf{x} - \mathbf{B}_t \mathbf{B}_t^\top \mathbf{x}\|. \quad (6)$$

This metric reflects how much of  $\mathbf{B}_i$  lies outside  $\mathbf{B}_t$ : It equals 1 if the spaces are orthogonal and 0 if one is a subspace of the other. Similarly, we compute  $d(\mathbf{B}_i, \mathbf{B}_{vc})$  and  $d(\mathbf{B}_i, \mathbf{B}_{t+vc})$ .

As shown in Tab. 4, due to the modality gap, the text classifier space significantly deviates from the image feature space. In contrast, the classifier that compensates for modality gap better aligns with the image space. Combining both classifiers further improves coverage and highlighting their complementary roles.

**Analyzing the impact of our method on the modality gap.** We analyze the impact of our method on the modality gap. Tab. 5 presents the mean cosine similarity of the positive and negative for the final task under different experimental settings, along with the last accuracy (Last Acc). The ‘Base’ refers to the results of naive fine-tuning. It can be observed that the mean of the negative similarity signif-

|                                      | CIFAR-100     | ImageNet-R    | ImageNet100   |
|--------------------------------------|---------------|---------------|---------------|
| $d(\mathbf{B}_i, \mathbf{B}_t)$      | 0.7732        | 0.7160        | 0.7664        |
| $d(\mathbf{B}_i, \mathbf{B}_{vc})$   | 0.6076        | 0.5843        | 0.5414        |
| $d(\mathbf{B}_i, \mathbf{B}_{t+vc})$ | <b>0.4917</b> | <b>0.3284</b> | <b>0.4431</b> |

Table 4. Difference measurement of image subspace and subspace of different classifiers on different datasets.

|          | CLIP   | Base   | Distill | Ours         |
|----------|--------|--------|---------|--------------|
| pos      | 0.3157 | 0.3037 | 0.3168  | 0.3251       |
| neg      | 0.1719 | 0.0579 | 0.1718  | 0.1520       |
| Last Acc | 75.40  | 72.74  | 75.85   | <b>78.38</b> |

Table 5. Comparison of cosine similarity and last accuracy (Last Acc) across different experimental settings on ImageNet-100.

icantly decreases compared to the original CLIP. This suggests that the modality gap is expanded, leading to a lower last accuracy. The ‘Distill’ represents the traditional distillation method, which explicitly limits the output magnitude of the model. It restricts the changes in the modality gap and has some positive effect. However, a complete restriction of the modality gap evidently leads to a decline in the model’s learning capability. This leads to suboptimal performance. We can see that its performance is close to the original CLIP. In contrast, our method maintains a relatively stable modality gap, allowing for a moderate increase in the positive similarity and a decrease in the negative similarity. This balance ensures that the model can appropriately learn new knowledge while maintaining pre-trained knowledge, ultimately achieving optimal performance.

## 6. Conclusion

This paper investigates the impact of modality gap on the performance of vision-language pre-trained models in class-incremental learning. We find that maintaining a relatively stable modality gap helps preserve pre-trained knowledge and prevents its degradation. Under the condition of a stable modality gap, training a visual space classifier that is not restricted by the modality gap can compensate for some of its negative effects, further enhancing the model’s capabilities. The experimental results demonstrate the effectiveness of our approach.

**Limitations and Future Work.** Our focus has been on the modality gap, and the model has been fine-tuned simply using LoRA without considering other fine-tuning methods. our current method does not incorporate specially designed loss functions or parameter constraint to mitigate forgetting. Future work will explore the integration of this approach with other continual learning methods and investigate suitable distillation strategies.

## Acknowledgment

This work was funded by NSFC (NO. 62206135, 62225604), Young Elite Scientists Sponsorship Program by CAST (2023QNRC001), “Science and Technology Yongjiang 2035” key technology breakthrough plan project (2024Z120), Shenzhen Science and Technology Program (JCYJ20240813114237048), and the Fundamental Research Funds for the Central Universities (Nankai University, 070-63233085). Computation was supported by the Supercomputing Center of Nankai University.

## References

- [1] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014, proceedings, part VI 13*, pages 446–461. Springer, 2014. 7
- [2] Yoojin Choi, Mostafa El-Khamy, and Jungwon Lee. Dual-teacher class-incremental learning with data-free generative replay. In *CVPR*, pages 3543–3552, 2021. 2
- [3] 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 TPAMI*, 44(7):3366–3385, 2021. 2
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. Ieee, 2009. 6
- [5] Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *CVPR*, pages 9285–9295, 2022. 2
- [6] Sedigheh Eslami and Gerard de Melo. Mitigate the gap: Investigating approaches for improving cross-modal alignment in clip. *arXiv preprint arXiv:2406.17639*, 2024. 2
- [7] Abrar Fahim, Alex Murphy, and Alona Fyshe. It’s not a modality gap: Characterizing and addressing the contrastive gap. *arXiv preprint arXiv:2405.18570*, 2024. 3, 4
- [8] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*, 2021. 6
- [9] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6
- [10] Xuefeng Hu, Ke Zhang, Lu Xia, Albert Chen, Jijia Luo, Yuyin Sun, Ken Wang, Nan Qiao, Xiao Zeng, Min Sun, et al. Reclip: Refine contrastive language image pre-training with source free domain adaptation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2994–3003, 2024. 2, 3
- [11] Linlan Huang, Xusheng Cao, Haori Lu, and Xialei Liu. Class-incremental learning with clip: Adaptive representation adjustment and parameter fusion. In *European Conference on Computer Vision*, pages 214–231. Springer, 2024. 1, 2, 6
- [12] Saurav Jha, Dong Gong, and Lina Yao. Clap4clip: Continual learning with probabilistic finetuning for vision-language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 1, 2, 6
- [13] 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
- [14] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 5
- [15] Kuan-Ying Lee, Yuanyi Zhong, and Yu-Xiong Wang. Do pre-trained models benefit equally in continual learning? In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 6485–6493, 2023. 1
- [16] Yunheng Li, Yuxuan Li, Quansheng Zeng, Wenhai Wang, Qibin Hou, and Ming-Ming Cheng. Unbiased region-language alignment for open-vocabulary dense prediction. *arXiv preprint arXiv:2412.06244*, 2024. 1
- [17] Yunheng Li, Zhong-Yu Li, Quan-Sheng Zeng, Qibin Hou, and Ming-Ming Cheng. Cascade-CLIP: Cascaded vision-language embeddings alignment for zero-shot semantic segmentation. In *Proceedings of the 41st International Conference on Machine Learning*, pages 28243–28258. PMLR, 2024. 1
- [18] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE TPAMI*, 40(12):2935–2947, 2017. 1, 2
- [19] Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y Zou. Mind the gap: Understanding the modality gap in multi-modal contrastive representation learning. *Advances in Neural Information Processing Systems*, 35:17612–17625, 2022. 1, 3
- [20] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017. 2
- [21] Daniel Marczak, Bartłomiej Twardowski, Tomasz Trzcíński, and Sebastian Cygert. Magma: Leveraging model merging for seamless continual learning. In *European Conference on Computer Vision*, pages 379–395. Springer, 2024. 1, 2, 6, 7
- [22] 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 TPAMI*, 45(5):5513–5533, 2022. 1
- [23] 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
- [24] Marco Mistretta, Alberto Baldreti, Lorenzo Agnolucci, Marco Bertini, and Andrew D Bagdanov. Cross the gap: Exposing the intra-modal misalignment in clip via modality inversion. *arXiv preprint arXiv:2502.04263*, 2025. 2, 3
- [25] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *2012 IEEE conference on*

- computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012. 7
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1, 6
- [27] 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. 2
- [28] 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 *CVPR*, pages 11909–11919, 2023. 2, 6
- [29] Zhicheng Sun, Yadong Mu, and Gang Hua. Regularizing second-order influences for continual learning. In *CVPR*, pages 20166–20175, 2023. 2
- [30] Vishal Thengane, Salman Khan, Munawar Hayat, and Fahad Khan. Clip model is an efficient continual learner. *arXiv preprint arXiv:2210.03114*, 2022. 1, 2, 6
- [31] Gido M Van de Ven, Tinne Tuytelaars, and Andreas S Tolias. Three types of incremental learning. *Nature Machine Intelligence*, 4(12):1185–1197, 2022. 1
- [32] Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *ECCV*, pages 398–414. Springer, 2022. 2
- [33] Liyuan Wang, Xingxing Zhang, Kuo Yang, Longhui Yu, Chongxuan Li, HONG Lanqing, Shifeng Zhang, Zhenguo Li, Yi Zhong, and Jun Zhu. Memory replay with data compression for continual learning. In *International Conference on Learning Representations*, 2021. 2
- [34] 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 *ECCV*, pages 631–648. Springer, 2022. 2, 6
- [35] 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 *CVPR*, pages 139–149, 2022. 2, 6
- [36] Tz-Ying Wu, Gurumurthy Swaminathan, Zhizhong Li, Avinash Ravichandran, Nuno Vasconcelos, Rahul Bhotika, and Stefano Soatto. Class-incremental learning with strong pre-trained models. In *CVPR*, pages 9601–9610, 2022. 1
- [37] Tz-Ying Wu, Gurumurthy Swaminathan, Zhizhong Li, Avinash Ravichandran, Nuno Vasconcelos, Rahul Bhotika, and Stefano Soatto. Class-incremental learning with strong pre-trained models. In *CVPR*, pages 9601–9610, 2022. 2
- [38] Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *CVPR*, pages 3014–3023, 2021. 2
- [39] Chao Yi, Lu Ren, De-Chuan Zhan, and Han-Jia Ye. Leveraging cross-modal neighbor representation for improved clip classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27402–27411, 2024. 3
- [40] Jiazuo Yu, Yunzhi Zhuge, Lu Zhang, Ping Hu, Dong Wang, Huchuan Lu, and You He. Boosting continual learning of vision-language models via mixture-of-experts adapters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23219–23230, 2024. 1, 2, 6, 7
- [41] Jiang-Tian Zhai, Xialei Liu, Andrew D Bagdanov, Ke Li, and Ming-Ming Cheng. Masked autoencoders are efficient class incremental learners. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19104–19113, 2023. 2
- [42] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. A large-scale study of representation learning with the visual task adaptation benchmark. *arXiv preprint arXiv:1910.04867*, 2019. 6
- [43] 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. 2, 6
- [44] Junting Zhang, Jie Zhang, Shalini Ghosh, Dawei Li, Serafettin Tasci, Larry Heck, Heming Zhang, and C-C Jay Kuo. Class-incremental learning via deep model consolidation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1131–1140, 2020. 2
- [45] Wentao Zhang, Yujun Huang, Weizhuo Zhang, Tong Zhang, Qicheng Lao, Yue Yu, Wei-Shi Zheng, and Ruixuan Wang. Continual learning of image classes with language guidance from a vision-language model. *IEEE Transactions on Circuits and Systems for Video Technology*, 2024. 2
- [46] Zangwei Zheng, Mingyuan Ma, Kai Wang, Ziheng Qin, Xianguyu Yue, and Yang You. Preventing zero-shot transfer degradation in continual learning of vision-language models. In *ICCV*, pages 19125–19136, 2023. 2
- [47] Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Co-transport for class-incremental learning. In *ACM MM*, pages 1645–1654, 2021. 2
- [48] Da-Wei Zhou, Zi-Wen Cai, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-incremental learning with pre-trained models: Generalizability and adaptivity are all you need. *International Journal of Computer Vision*, pages 1–21, 2024. 2, 6
- [49] Da-Wei Zhou, Yuanhan Zhang, Yan Wang, Jingyi Ning, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Learning without forgetting for vision-language models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025. 1, 2, 6