CTP: Towards Vision-Language Continual Pretraining via Compatible Momentum Contrast and Topology Preservation

Hongguang Zhu^{1,2*} Yunchao Wei^{1,2,3} Yao Zhao^{1,2,3†} Xiaodan Liang^{3,4,5} Chuniie Zhang^{1,2}

¹Institute of Information Science, Beijing Jiaotong University ²Beijing Key Laboratory of Advanced Information Science and Network ³Peng Cheng Laboratory ⁴Sun Yat-sen University ⁵MBZUAI

Abstract

Vision-Language Pretraining (VLP) has shown impressive results on diverse downstream tasks by offline training on large-scale datasets. Regarding the growing nature of real-world data, such an offline training paradigm on ever-expanding data is unsustainable, because models lack the continual learning ability to accumulate knowledge constantly. However, most continual learning studies are limited to uni-modal classification and existing multimodal datasets cannot simulate continual non-stationary data stream scenarios. To support the study of Vision-Language Continual Pretraining (VLCP), we first contribute a comprehensive and unified benchmark dataset P9D which contains over one million product image-text pairs from 9 industries. The data from each industry as an independent task supports continual learning and conforms to the real-world long-tail nature to simulate pretraining on web data. We comprehensively study the characteristics and challenges of VLCP, and propose a new algorithm: Compatible momentum contrast with Topology Preservation, dubbed CTP. The compatible momentum model absorbs the knowledge of the current and previous-task models to flexibly update the modal feature. Moreover, Topology Preservation transfers the knowledge of embedding across tasks while preserving the flexibility of feature adjustment. The experimental results demonstrate our method not only achieves superior performance compared with other baselines but also does not bring an expensive training burden. Dataset and codes are available at https://github. com/KevinLight831/CTP.

1. Introduction

Benefiting from the remarkable generalization ability derived from large-scale pretraining, Vision-Language Pretraining (VLP) [50, 33] has emerged as the prevalent approach for addressing downstream vision-language tasks. The recent advancements in artificial intelligence such as CLIP [50] and ChatGPT [1] have further fueled this trend





(b) Vision-Language Continual Pretraining (VLCP).

Figure 1: The traditional Class-Incremental Learning (CIL) is inflexible in the continual learning of visual concepts, which needs ever-expanding classifier parameters and endless human annotation. Moreover, it is difficult for singleclass labeling to cover all visual concepts in an image. e.g., CIL only focuses on the foreground class dog and ignores the background class flower, while Vision-Language Continual Pretraining (VLCP) can flexibly represent the image content by text. Compared with CIL, which fixes the label space and only updates the image encoder, VLCP updates the image and text encoders simultaneously in the fixed dimension. Meanwhile, previous-task data also cannot be used as contrast samples in the continual pretraining.

of using larger models and more data. In the long term, this computational headlong rush does not seem reasonable to move toward sustainable solutions and actually also excludes academic laboratories with limited resources. Current VLP paradigms all train on prepared data in advance. Nevertheless, the world is ever-changing. Offlinetrained models can not evolve in a dynamic environment to continually acquire, integrate and accumulate new knowledge. Moreover, repeated offline pretraining on the everexpanding dataset will impose growing and endless training costs. Only finetuning on new data will also suffer severe

arXiv:2308.07146v1 [cs.CV] 14 Aug 2023

^{*} This work was done when Hongguang Zhu worked as an research intern in Peng Cheng Laboratory. Email: kevinlight831@gmail.com

[†] Corresponding author: yzhao@bjtu.edu.cn

degradation due to *catastrophic forgetting* [45]. Hence, in practical application scenarios, it is significant for VLP to continually integrate knowledge from the incoming data.

Prior studies on continual learning [28, 80, 65, 78, 76] focused on supervised class-incremental learning (CIL), as Figure 1(a), which aims to maintain discriminative features for known classes and expand new classifiers to learn new classes. However, this paradigm is inflexible due to the constant demand for laborious annotation and increasing classifier parameters. In contrast, VLP allows for learning open-world visual patterns without explicit "class" concepts, which can capture more comprehensive visual concepts rather than just category-based features. Moreover, massive web weak-aligned image-text pairs can be used as training data without extra human annotation, and no extra parameters are needed for category expansion as the output dimension is fixed for image and text encoders.

Nevertheless, Vision-Language Continual Pretraining (VLCP) remains understudied due to the lack of datasets that satisfy both massive image-text pairs and continual tasks with discrepant knowledge. Therefore, we contribute the first VLCP dataset P9D, which contains more than 1 million product image-text pairs and over 3800 categories from 9 industries. Different task data are split according to industry (e.g. food and clothing) to support the continual pretraining. P9D not only is larger than previous CIL datasets both in terms of both class number and data size. but also conforms to real-world long-tailed distributions. As shown in the Figure 1(b), VLCP, as a new paradigm, also suffers new challenges compared with traditional CIL. 1) Fixed-dimensional embedding: CIL methods typically address the stability-plasticity dilemma by preserving old logits [37] or freezing old backbone [44] and finetuning new classifiers. Without explicit class supervision and increasing embedding dimension, the VLCP can only adjust the fixed-dimensional shared embedding to incorporate both old and new knowledge. 2) Missing contrast samples: CIL still can use the gradient from negative logits of old classes [37, 51, 24] even if the old data is unseen. But the lack of contrast samples from old tasks leads to suboptimal shared embedding in VLCP. 3) Multi-modal/task optimization: Unlike CIL has fixed label space to optimize image encoder, VLCP involves the complicated joint optimization of image, text, and multi-modal encoders.

Therefore, we propose a simple yet effective method, Compatible momentum contrast with Topology Preservation (CTP), which maintains a compatible momentum model that absorbs both new and old knowledge to separately adjust uni-modal and multi-modal encoders. Moreover, different from CIL methods that distill visual features across tasks, topology preservation keeps consistent sample relationship across tasks. It not only transfers the topology knowledge of the old embedding while preserving the flexibility of feature adjustment. Meanwhile, to systematically investigate the vision-language continual pretraining, we extend a series of traditional CIL methods to VLCP and evaluate them in a unified setting. Interestingly, we find that the multi-modal fusion feature by masked modeling pretraining has a strong anti-forgetting ability, and the performance of continual finetuning approximates that of joint training in multi-modal retrieval. Oppositely, due to the lack of contrastive samples from different tasks, the cross-modal alignment ability suffers serious forgetting in continual pretraining and has a big gap with joint training. Meanwhile, The experimental results show our method is not only compatible with both memory-buffer and memory-free situations, but also achieve leading performance without incurring expensive training time costs.

Our contributions are as follows:

- We build the first Vision-Language Continual Pretraining (VLCP) benchmark dataset P9D to support the study of VLCP, which contains massive image-text pairs and the continual non-stationary task data.
- We systematically study the characteristics and challenges of VLCP, and establish baseline library for VLCP by extending popular continual learning methods and evaluating them in a unified setting.
- We propose a simple yet effective method CTP for VLCP, which achieve both superior performance and efficient training.

2. Related Work

2.1. Vision-language pretraining

Vision-Language Pretraining (VLP) [20] leverages large-scale web image-text pairs as pretraining data and adopts self-supervised learning (contrastive learning [50] or masked modeling [5]) to train the transferable image-text embeddings. The VLP models can coarsely be divided into two paradigms according to architectures: 1) Dual-Encoder and 2) Fusion-Encoder. Dual-Encoder models encode images and texts respectively by separate encoders and employ cosine similarity to build the image-text alignment. The Dual-Encoder models [41, 50, 27] achieve promising results on image-text retrieval with linear time complexity. However, the loose modal interaction by cosine similarity also limits the multi-modal fusion ability [4, 10, 70]. Thus, the other paradigm Fusion-Encoder employs crossmodal attention to jointly encode images and text. The prior works [9, 35, 77] use pretrained detectors to extract regional features and Transformers [61] for multi-modal fusion. However, extracting region features is computationally expensive and the joint transformer requires quadratic time complexity for retrieval tasks. Thus, align before fuse (ALBEF) architecture [34, 33, 73] incorporates the imagetext contrastive loss before multi-modal fusion and replaces

the detector with VIT [16] for the end-to-end pretraining. It not only eliminates the burden of object detection preprocessing but also keeps multi-modal interaction by the top fusion layers and linear retrieval time complexity by the bottom dual encoders.

Nevertheless, current VLP models are all trained in a joint manner using prepared data and keep fixed pretrained parameters. In the long term, they cannot constantly accumulate knowledge and evolve themselves to accommodate the dynamic world. Thus, we concentrate on the vision-language continual pretraining based on the ALBEF architecture and evaluate the cross-modal alignment and multi-modal fusion capabilities in the continual environment.

2.2. Continual Learning

Continual learning aims to overcome catastrophic forgetting and integrate novel knowledge in a sequential fashion where old data are unavailable. Conventional continual learning methods mainly focus on image classification tasks. They can be roughly categorized into three groups: 1) **Regularization-based methods** [28, 2, 74, 64, 7, 37] limit the plasticity of the model to address catastrophic forgetting by regularizing important parameters or knowledge distillation. Although these methods alleviate forgetting to some extent without storing old samples, they cannot get satisfactory performance in some complex datasets [69] and challenging settings [42]. 2) Architecture-based methods [54, 56, 44, 36] keep old parameters fixed while growing and allocating weights for learning new data. These methods can expand sub-networks or focus more on the specific part of network modules. However, these models require task identity to condition the network at test time, which is impractical for more realistic and task-agnostic settings e.g. retrieval tasks. Additionally, as the number of tasks increases, the parameters of the added sub-networks become very huge, which is also not suitable for application deployment. 3) **Replay-based methods** [8, 51, 24] apply extra memory to store a few samples from previous tasks or learn to generate pseudo data and train with the current data together. Based on this simple yet effective idea, recent methods further improve and achieve state-of-the-art performance by involving different sampling strategies. However, the memory size and the training complexity will be enlarged significantly and unaffordable as the growth of tasks, especially for costly large-scale pretraining.

However, It is an inflexible way to get the foundation model with continual learning ability through image classification. Firstly, most real-world data can not be accurately represented semantics by simple category labeling, and class annotation is labor-intensive. Secondly. the everexpanding classifier will also bring endless growth of parameters. In contrast, VLP does not require explicit "class" labeling and can cover wider visual concepts by text, while its fixed output dimensions can continually support downstream tasks without increasing parameters.

3. Dataset

Traditional CIL datasets [32, 31] usually use simple images with single-semantic and have the same sample number for different tasks. This ideal setting is difficult to simulate the real-world data which is noisy and long-tail distributed. Meanwhile, despite massive web data collected by some VLP datasets [57, 6], The simple random partitioning will make each chunk still conform to the overall distribution and unable to simulate the continual environment. Considering the rich product samples in e-commerce websites, which not only conform to the real web-data nature but also have category, industry, and title information as weakly semantic correspondence, we use e-commerce data to construct the first vision-language continual pretraining dataset P9D and establish the unified evaluation benchmark.

3.1. Dataset Split

As the Figure 2(c), P9D includes more than 1 million image-text pairs of real products. According to the industry name of products, P9D is divided into 9 tasks to sequential training (default order) which are Household, Furnishings, Food, Beauty, Clothing, Auto, Parenting, Outdoor, and Electronics. We select 1,014,599 image-text pairs for training, and 2,846 pairs as the test set of cross-modal retrieval. 4,615 and 46,855 pairs as the query set and gallery set of multi-modal retrieval. The quantity distributions of the four sets are consistent and more details can refer to the appendix. Considering the descriptions of similar products may be very similar, to avoid the situation that one image corresponds to multiple captions affecting the overall evaluation of cross-modal retrieval, we save one sample for the categories owning more than 100 samples in the training set. Meanwhile, For the query set, the number of samples per class is about 0.5% of that in the training set.

3.2. Data Analysis.

Real-world Web Data: We collect massive commodity data from e-commerce websites and split task identity according to industry to simulate real-world rich-concept but non-stationary data streams. Different from the existing CIL datasets with clear images and labels, our images present the characteristics of multi-domain mixing (as Figure 2(b)). e.g., multiple backgrounds, amorphous watermarks, occlusion, and multi-view. Meanwhile, although only text can be used in training, P9D also includes product class labels to evaluate the fusion feature clustering by multi-modal retrieval.

Rich Categories: Some recent papers can obtain preferable classification results on traditional CIL datasets by finetun-



Figure 2: The dataset statistic of P9D. (a) The abbreviated hierarchical categories structure of P9D. (b) Some multi-domain mixing examples about the product 'Keyboard and mouse set'. (c) The quantity distribution of P9D train set. (d) The sample number distribution of category in P9D. The red line represents the sample number of each category in decreasing order.

ing [67] and even freezing [58] the weight of pretrained models. We suppose that the traditional benchmark datasets contain limited classes (e.g. MNIST [32] and CIFAR-100 [31]) and thus cannot adequately evaluate the continual learning methods in the face of powerful generalized pre-trained models. Unlike these datasets with the narrow space of category labels, our dataset contains over 3800 categories and divides each task in a way that matches the real-world industry domains. Therefore, our setting is more challenging and realistic. The abbreviated hierarchical categories structure of P9D is shown in the Figure 2(a), and more comparisons with other datasets can refer to the appendix.

Real-World Distribution: The conventional CIL datasets consider a balanced distribution for each task but ignore the nature of long-tailed distributions in the real world. In contrast, our dataset contains a different number of categories for each task and the long-tailed distribution aligns well with real-world scenarios. The Figure 2(d) shows the sample number distribution of the categories.

4. Methodology

4.1. Preliminary

Problem Setting: We propose a vision-language continual pretraining (VLCP) setting, where models are supposed to be sequentially trained on T tasks data $\mathcal{D} = \{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_T\}$. In the t-th task, the whole sub-dataset $\mathcal{D}_t = \{(I_i^t, T_i^t)\}_{i=1}^{N_t} \text{ contains the } N_t \text{ image-text pairs where } I_i^t \text{ and } T_i^t \text{ respectively denote the } i\text{-th image and the corresponding text description of the } t\text{-th task. Because the old task data is unseen in the continual setting, the VLP model trained on the current dataset <math>\mathcal{D}_t$ needs to resist forgetting and performs well for all learned datasets $\{\mathcal{D}_1, \mathcal{D}_2, ..., \mathcal{D}_t\}$.

Model Architecture: As shown in Figure 3, We use a 12layer ViT-B/16 [16] as the image encoder f_v , and initialize it with weights pretrained on ImageNet-1k from [60]. The image I is encoded into patch feature sequence $f_v(I) =$ $\{v^{\text{cls}}, v^1, ..., v^N\}$. Meanwhile, We use the first and the last 6 layers of BERT_{base} [14] model to initialize the text encoder f_t and multi-modal encoder f_m . The text T is encoded into word feature sequence $f_t(T) = \{w^{\text{cls}}, w^1, ..., w^N\}$. Then, the text features $f_t(T)$ will fuse with the image features $f_v(I)$ through cross attention at each layer of multi-modal encoder f_m . Two linear transformations g_v and g_t are map the v^{cls} and w^{cls} to the low-dimensional (256-d) representations for Image-Text Alignment (ITA). It can be denoted $v = g_v(v^{\text{cls}})$ and $w = g_w(w^{\text{cls}})$.

Given a pair of image-text data (I_i, T_j) in the current D_t , the image-text similarity function $s(I_i, T_j)$ is defined as the cosine similarity $s(I_i, T_j) = \frac{v_i^T w_j}{\|v_i\| \|w_j\|}$. Given a batch of Bpairs, the model uses the symmetric cross-entropy loss over the $B \times B$ similarity matrix to optimize the parameters. The image-to-text loss \mathcal{L}_{i2t} and the text-to-image loss \mathcal{L}_{t2i} are



Figure 3: Illustration of the proposed vision-language continual pretraining method CTP. (a) shows the overall continual pretraining model architecture, and (b) shows the interactive adjustment of the compatible momentum model M_c and current training model M_t . The M_c absorb the parameter of current model M_t and previous-task model M_{t-1} , and in turn constrains the update of M_t . (c) shows that the topology relationship of previous task is maintained as much as possible while allowing the overall embedding to be updated.

formulated as:

$$\mathcal{L}_{i2t} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(s(I_i, T_i)/\tau)}{\sum_{j=1}^{B} \exp(s(I_i, T_j)/\tau)},$$

$$\mathcal{L}_{t2i} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(s(T_i, I_i)/\tau)}{\sum_{j=1}^{B} \exp(s(T_i, I_j)/\tau)},$$
(1)

where τ is the temperature parameter. The Image-Text Alignment (ITA) loss is defined as : $\mathcal{L}_{ita} = \frac{\mathcal{L}_{i2t} + \mathcal{L}_{t2i}}{2}$.

In Masked Language Modeling (MLM), give an imagetext pair, we randomly mask out the words with probability of 15% [14, 34], and replace masked ones $w_{\rm m}$ with the special token [MASK]. The goal is to predict these masked tokens based on their surrounding words and the image features, by minimizing the cross-entropy loss:

$$\mathcal{L}_{mlm} = -\mathbb{E}_{(I,\hat{T})\sim D_t} \mathbf{H}(y^{\mathsf{m}}, p_{\theta}^{\mathsf{m}}(I, \hat{T})), \qquad (2)$$

where the $y^{\rm m}$ is the one-hot vocabulary distribution where the masked token has a probability of 1, and $p_{\theta}^{\rm m}(I, \hat{T})$ is the predicted probability of model θ for masked token $w_{\rm m}$. The total VLP loss is defined as $\mathcal{L}_{VLP} = \mathcal{L}_{ita} + \mathcal{L}_{mlm}$.

4.2. Compatible Momentum Contrast

Due to the fixed-dimensional embedding, VLP cannot isolate old and new knowledge like CIL by extending projection parameters. Therefore, both the vision and language embeddings need to be constantly adjusted to simultaneously accommodate the old and new image-text pairs. In order to review old knowledge and adapt to the new task, we use the momentum model M_c initialized by the previous-task trained model M_{t-1} as additional supervision of the current training model M_t . Some single-modal [22] or vison-language [26] pretraining works also adopt momentum model as a temporal ensembling method [19] to smoothly guide the training. However, the traditional momentum model is updated by only parameters of training model. With the accumulation of training steps, it will be gradually affected by the new model and also suffers catastrophic forgetting. Moreover, recklessly maintaining old knowledge will also make the model lose the plasticity to acquire new knowledge. Therefore, we propose the compatible momentum update which synchronously absorbs the old and new knowledge:

$$\theta_c \leftarrow m \cdot \theta_c + \frac{1-m}{2} \cdot \theta_{t-1} + \frac{1-m}{2} \cdot \theta_t,$$
 (3)

where $m \in [0, 1)$ is the momentum coefficient and θ is the model parameters. The adjustment of models is interactive. Compatible momentum model M_c updates parameters θ_c through the previous-task model M_{t-1} and training model M_t . In turn, the training model M_t are optimized by the back-propagation of momentum contrast and affects the parameters θ_t to be passed in the next step. To further steadily update uni-modal encoders f_v and f_t , we maintain two dynamic queues to preserve the K negative image/text features. The image features v_i^c and text feature w_i^c from compatible momentum encoders are constantly pushed into visual queue $Q^I = \{v_1^c, v_2^c, \cdots, v_{N_Q}^c\}$ and text queue $Q^T = \{w_1^c, w_2^c, \cdots, w_{N_Q}^c\}$ which $N_Q = K + B$. The compatible momentum contrastive losses about the vision and language encoders can be formulated as follows:

$$\mathcal{L}_{i2t}^{c} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(s(I_i, Q_i^T)/\tau)}{\sum_{j=1}^{N_Q} \exp(s(I_i, Q_j^T))/\tau)},$$

$$\mathcal{L}_{t2i}^{c} = -\frac{1}{B} \sum_{i=1}^{B} \log \frac{\exp(s(T_i, Q_i^I)/\tau)}{\sum_{j=1}^{N_Q} \exp(s(T_i, Q_j^I)/\tau)},$$
 (4)

Similarly, $\mathcal{L}_{ita}^{c} = \frac{\mathcal{L}_{i2t}^{c} + \mathcal{L}_{t2i}^{c}}{2}$. \mathcal{L}_{ita}^{c} constraints that the contrastive relation of image-text pairs is still workable between the encoders of the compatible momentum model and current model. It allows slow adjustment of uni-modal encoders. Besides, The compatible momentum model also provides the soft predicted probability for masked language modeling loss.

$$\mathcal{L}_{mlm}^{c} = -\mathbb{E}_{(I,\hat{T})\sim D_{t}} \mathbf{H}(p_{\theta_{c}}^{\mathsf{m}}(I,\hat{T}), p_{\theta_{t}}^{\mathsf{m}}(I,\hat{T})), \quad (5)$$

Thus, Compatible Momentum Contrastive loss can be defined as $\mathcal{L}_{CMC} = \mathcal{L}_{ita}^c + \mathcal{L}_{mlm}^c$.

4.3. Topology Preservation

Although the compatible momentum contrast can flexibly adjust the output feature of uni-modal and fusion encoder. It does not directly transfer relationship knowledge of samples across tasks and the model may forget the overall topology of prior embedding to obtain sub-optimal performance. Unlike CIL can receive the gradient of old classes according to labels, VLCP has no gradient from old contrast samples and the image-text encoder is bidirectional synchronization adjustment according to the image-text similarity. To integrate the old and new knowledge while maintaining the topology relationship of prior tasks, we constrained the mini-batch sample relationships of the current and previous-task models to be consistent. Specifically, it is mainly divided into cross-modal topology preservation loss \mathcal{L}_c and same-modal topology preservation loss \mathcal{L}_s . To image $I = \{I_1, I_2, ..., I_B\}$ and text $T = \{T_1, T_2, ..., T_B\}$, we define the image-to-text similarity distribution as follows:

$$\mathcal{P}^{i2t} = \frac{\exp(s(I, T_i)/\tau)}{\sum_{i=1}^{B} \exp(s(I, T_i))/\tau)},$$
(6)

and the definition of text-to-image similarity distribution \mathcal{P}^{t2i} is similar. Thus, the \mathcal{L}_c can be formulated as the crossentropy H between current model M_t and previous-task model M_{t-1} :

$$\mathcal{L}_{c} = \frac{1}{2} \mathbb{E}_{(I,T)\sim D_{t}} [H(\mathcal{P}_{\theta_{t-1}}^{i2t}, \mathcal{P}_{\theta_{t}}^{i2t}) + H(\mathcal{P}_{\theta_{t-1}}^{t2i}, \mathcal{P}_{\theta_{t}}^{t2i})], \quad (7)$$

In the same-modal topology preservation, the similarity of the same sample is 1 under both old and new models. However, such a large similarity would suppress the relational distillation of other unmatched pairs [79]. Thus, we conduct the simple strategy that changes the similarity at the diagonal of s(I, I) from 1 to a minimum (*e.g.*-1000) to exclude the "apical dominance" of matched samples. We formulate the changed matrix as $\hat{s}(I, I)$ and the \mathcal{L}_s can be formulated as:

$$\hat{\mathcal{P}}^{i2i} = \frac{\exp(\hat{s}(I, I_i)/\tau)}{\sum_{i=1}^{B} \exp(\hat{s}(I, I_i))/\tau)},$$
(8)

$$\mathcal{L}_{s} = \frac{1}{2} \mathbb{E}_{(I,T)\sim D_{t}} [H(\hat{\mathcal{P}}_{\theta_{t-1}}^{i2i}, \hat{\mathcal{P}}_{\theta_{t}}^{i2i}) + H(\hat{\mathcal{P}}_{\theta_{t-1}}^{t2t}, \hat{\mathcal{P}}_{\theta_{t}}^{t2t})], \quad (9)$$

Thus, the overall continual pretraining loss is as follows:

$$\mathcal{L} = \mathcal{L}_{VLP} + \mathcal{L}_{CMC} + \mathcal{L}_c + \mathcal{L}_s, \tag{10}$$

5. Experiment

5.1. Experimental Setup

Implementation Details Regarding offline pretrained VLP models have established widespread generalization, direct finetuning may involve knowledge leakage and weaken the accurate evaluation of continual learning ability. Therefore, we do not load the weight of the pretrained VLP model for initialization. For each task, All models are trained for 5 epochs on 4 NVIDIA A100 GPUs with batch size 128 per GPU. We use the AdamW [40] optimizer with a weight decay of 0.05, and the learning rate is set as 1e-4 and decays to 1e-6 following a cosine schedule. We take random image crops of resolution 224×224 during pretraining and also apply RandAugment. The maximum sequence length of tokens is limited to 30. The momentum parameter is set as 0.9, and the queue size K is set as 1024. The cross-modal alignment temperature $\tau = 0.07$.

Evaluation Setting. With the knowledge accumulation in continual learning, the evaluation galleries of cross-modal and multi-modal retrieval also need to expand by merging the data of new tasks. Similar to the standard vision-language modeling setting[50, 34, 33], For cross-modal retrieval, we measure the performance with Recall at K (R@K, K=1,5,10) [18], which is defined as the proportion of ground truth being retrieved at top-K of the ranking list. Also, we use Rm as the overall metric, which is defined as the mean of R@K of both text and image retrieval.

The multi-modal encoder can fuse the multi-modal information and produce a comprehensive representation to predict masked words. Thus, we use the multi-modal output [CLS] token to retrieve samples of the same class to evaluate the zero-sample clustering ability of multi-modal fusion features for similar products. It will not introduce extra interference from finetuning on downstream tasks. (*e.g.* classification needs to train new classification head from

Methods	Training			Cross-modal Retrieval					Multi-modal Retrieval		
	hours	TR@1	TR@5	TR@10	IR@1	IR@5	IR@10	Rm	mAP@1	mAP@5	mAP@10
JointT	_	61.31	87.17	91.67	61.60	86.79	91.95	80.08	63.79	70.10	67.40
Memory-Free											
SeqF	3.4	34.79	62.83	72.03	35.73	63.63	72.14	56.86	62.15	68.03	65.30
SI [74]	4.9	34.96	63.60	73.01	35.84	63.49	73.44	57.39	61.32	67.48	64.78
MAS [2]	4.1	36.65	64.69	73.93	37.25	64.55	74.39	58.57	62.62	68.63	65.79
EWC [28]	4.6	37.28	65.11	74.38	38.05	65.57	75.05	59.24	62.99	68.75	65.62
AFEC [64]	8.7	37.67	66.34	74.67	38.79	66.48	75.16	59.85	62.58	68.44	65.66
LWF [37]	4.0	37.63	66.55	75.09	38.26	66.62	75.26	59.90	62.34	68.64	65.94
RWalk [7]	6.2	37.77	68.10	76.70	38.83	67.39	76.74	60.92	62.41	68.43	65.59
Our:CTP	4.0	43.43	72.10	80.08	43.39	71.15	79.06	64.87	62.64	68.21	65.10
Memory-Buff	er										
MoF [51]	4.8	42.87	71.93	80.60	43.71	72.03	80.67	65.30	62.84	68.79	65.98
LUCIR [24]	5.5	43.82	73.68	82.04	44.38	73.54	80.99	66.41	61.06	67.57	64.74
ER [<mark>8</mark>]	4.4	45.19	73.40	81.97	44.97	72.70	81.10	66.56	62.21	68.39	65.80
Kmeans [8]	4.7	46.17	74.65	82.33	45.92	73.68	81.80	67.26	62.77	68.82	66.04
ICARL [51]	5.3	45.85	74.63	82.57	46.24	73.23	81.83	67.39	63.51	69.20	66.39
Our:CTP+ER	4.7	50.53	77.62	84.57	49.79	76.77	84.47	70.63	62.62	68.68	65.81

Table 1: The final cross-modal and multi-modal retrieval performance comparison with different Memory-Free and Memory-Buffer continual learning baselines.



Figure 4: The performance curve of different methods on the continual learning. For Task *i*, the model θ_i only test on the merged test set of learned tasks $\{\mathcal{D}_0, \mathcal{D}_1, ..., \mathcal{D}_i\}$.

scratch). Because each query here corresponds to multiple targets with the same class, multi-modal retrieval considers both the number and ranking of targets in the top-K retrieved candidates. For multi-modal retrieval, we adopt mean Average Precision (mAP@N) [68, 75, 15] as the evaluation metrics. mAP@N is computed as follows:

$$mAP@N = \frac{1}{Q} \sum_{i=1}^{Q} \frac{1}{\mathbf{m}_{q}} \sum_{k=1}^{N} P_{q}(k) \delta_{q}(k)$$

$$P_{q}(k) = \frac{R_{q}(k) + 1}{k}$$
(11)

where Q is the total number of multi-modal queries and m_q is the total target number of the q-th query in the retrieved

top-N relevant candidates from gallery. $P_q(k)$ is the precision at rank k for the q-th query, and the $\delta_q(k)$ is a binary indicator function that returns 1 when the k-th prediction is correct for the q-th query and 0 otherwise. $R_q(k)$ is the current sum of returned correct predictions at the rank k.

Comparison Methods. To verify the effectiveness of our method, we compare it with several popular continual learning methods. Since all competitors are originally proposed for classification settings. We replace the instance-label cross-entropy loss with image-text contrastive loss to train the VLP model and add the Masked Language Modeling loss to train the multi-modal fusion capability as follows [34]. Besides, We use joint training (JointT) of all seen

samples as the upper-bound performance and sequential finetuning (SeqF) as the lower-bound. All baselines can be separated into Memory-Free and Memory-Buffer methods based on the availability of old samples. For the former. we evaluate representative regularization-based methods such as EWC [28], SI [74], MAS [2], AFEC [64], RWalk [7], and LwF [37]. As for the latter, we evaluate different replay sampling strategies: ER [8], Mean-of-Feature (MoF)[51], and Kmeans [8], and some replay-based methods: ICARL[51] and LUCIR [24]. We must emphasize that although Memory-Buffer methods usually have higher performance than Memory-Free methods in CIL, they bring expensive storage costs in large-scale pretraining, and performance is directly related to the replay sample size. For the comprehensive and unified study of VLCP, to all Memory-Buffer methods, we maintain a fixed-size buffer of 10000 image-text pairs (about 1% dataset) and continually update the stored samples. The detailed introduction for each method can refer to the appendix.

5.2. Experimental Results

Cross-modal vs. Multi-modal. From the Figure 4, we found an interesting phenomenon that SeqF (lower-bound) and JointT (upper bound) have a large gap (23.22% in Rm) in cross-modal retrieval, but a small gap (1.65% in mAP@1) in multi-modal retrieval. This phenomenon shows that multi-modal fusion is stronger anti-forgetting ability than cross-modal alignment in VLCP. We suppose that on the one hand, The redundancy and complementarity of multi-modal information help multi-modal fusion resist the forgetting of class attributes. On the other hand, the pretext task MLM beforehand creates a word prediction classifier that corresponds to each word of the dictionary and keeps consistency across tasks. This predefined and well-initialized label space [71] maybe train models more stably than ever-expanding label space.

Parameter vs. Topology Preservation. From the comparison of Memory-Free methods in Table 1, we observe that the regularization methods [74, 2, 64, 7] represented by EWC [28] perform poorly. Probably because such methods conservatively trust the old model parameters and cannot flexibly update representation to better accommodate new knowledge. Meanwhile, They introduce the extra postprocessing of calculating the fisher matrix which reduces training efficiency. In contrast, benefiting from flexible updating of compatible momentum contrast and soft constraints of topology preservation, Our CTP achieves 3.95% improvement over the most competitive method on Rm while not bring more time costs and extra memory burden.

Memory-Free *vs.* **Memory-Buffer.** It can be noted that the Memory-Buffer methods exhibit superior final performance and smaller performance fluctuations compared to

\mathcal{L}_{CMC}	\mathcal{L}_c	\mathcal{L}_s	ER	TR@1	IR@1	Rm	mAP@1
X	X	X	X	34.79	35.73	56.86	62.15
1	X	X	X	37.81	37.63	59.86	62.10
X	1	X	X	41.45	40.39	62.50	61.32
1	1	X	X	42.69	42.41	64.38	61.89
1	X	1	X	41.92	42.02	63.17	61.11
1	1	1	X	43.43	43.39	64.87	62.64
1	1	1	✓	50.53	49.79	70.63	62.62

Table 2: The ablation study on each component of CTP. the \mathcal{L}_{CMC} represent the compatible momentum contrast.

Method	TR@1	IR@1	Rm	mAP
only θ^{t-1}	41.95	42.23	63.57	62.06
only θ^t	40.41	40.34	62.32	62.32
w/o Q	42.69	40.83	63.94	61.63
w/o $\hat{\mathcal{P}}^{i2i}$ & $\hat{\mathcal{P}}^{t2t}$	40.76	40.86	63.12	62.04
Our:CTP	43.43	43.39	64.87	62.64

Table 3: The more detailed ablation results of CTP.

the Memory-Free methods. This is attributed to the use of old data from the memory buffer as contrast samples, which plays the role of joint optimization in continual pretraining. Meanwhile, we found different replay sampling strategies have similar performance, but Kmeans[8] has a slight advantage, It may be because Kmeans can unsupervised cluster features to sample representative points of current embedding without category prior. However, the Memory-Buffer methods also bring extra storage and time cost causing by old data preservation and retraining process. The results show CTP can be further improved by 5.76% with ER sampling strategy and outperform the second-place method by 3.24% on Rm.

5.3. Ablations

We perform ablation of each module in the CTP method and find that each module effectively improves the crossmodal retrieval performance from the Table 2.

CMC *vs.* **TP.** It seems that cross-modal topology preservation plays a more direct anti-forgetting role than compatible momentum contrast in cross-modal retrieval, which improves 5.64% on Rm but also decreases 0.82% on mAP@1. When they are combined, the cross-modal and multi-modal retrieval performance were all improved by 1.88% and 0.57%. In addition, The same-modal topology preservation further improves by 0.49% and 0.75%.

Compatible Update *vs.* **Single-way Update.** As shown in Table 3, we compare compatible momentum update with the single-way momentum update from the current model θ^t and the previous-step model θ^{t-1} , and find the single-way update can not get well results due to totally relay on

the current or old model. In addition, the momentum update from θ^{t-1} is slightly higher than that from θ^t in Rm. This may be because the model accumulates large knowledge in continual learning, and the non-globally optimal update will lead to more forgetting in the later tasks. Therefore, it should pay more attention to the maintenance of old knowledge in the later periods of continual learning.

Momentum Queue.: We compare with the situation without the momentum queue 'w/o Q' in Table 3 and find that the queue can slightly improve the performance. We suppose the queue stores negative samples of the current task, which plays a role of smoothing training and anti-forgetting to a certain extent.

Suppress Same-modal Maximum Similarity.: We find it results in a negative effect and leads to performance degradation if not suppressing same-modal maximum similarity. It indicates that the maximum similarity will obscure the sample relationship and affect the performance.

6. Conclusion

In this paper, we build the first Vision-Language Continual Pretraining benchmark dataset P9D which contains over 1 million image-text pairs and task data with discrepant knowledge to simulate the continual pretraining environment. Further, we comprehensively study new characteristics and challenges of VLCP, and propose a new approach CTP which combines the compatible momentum contrast and topology preservation to flexibly update model to accommodate the ever-changing embedding. It can achieve the superior performance while ensuring efficient training.

References

- [1] Chatgpt. https://openai.com/blog/chatgpt/. 1
- [2] 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. 3, 7, 8, 13, 14, 16
- [3] Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with a network of experts. In *Proceedings of the IEEE Conference on Computer Vision* and Pattern Recognition, pages 3366–3375, 2017. 14
- [4] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425– 2433, 2015. 2
- [5] Hangbo Bao, Wenhui Wang, Li Dong, and Furu Wei. Vlbeit: Generative vision-language pretraining. arXiv preprint arXiv:2206.01127, 2022. 2
- [6] Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pretraining to recognize long-tail visual concepts. In *Proceed*-

ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 3558–3568, 2021. 3, 14, 15

- [7] Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *Proceedings of the European conference on computer vision* (ECCV), pages 532–547, 2018. 3, 7, 8, 13, 16
- [8] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and M Ranzato. Continual learning with tiny episodic memories. 2019. 3, 7, 8, 13, 16
- [9] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *Proceedings* of the European Conference on Computer Vision., 2020. 2, 15
- [10] Marcella Cornia, Lorenzo Baraldi, and Rita Cucchiara. Show, control and tell: A framework for generating controllable and grounded captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8307–8316, 2019. 2
- [11] Teófilo Emídio De Campos, Bodla Rakesh Babu, Manik Varma, et al. Character recognition in natural images. VIS-APP (2), 7(2), 2009. 14
- [12] 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. 13, 14
- [13] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. *Computer Vision and Pattern Recognition*, 2009. 14
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019. 4, 5
- [15] Xiao Dong, Xunlin Zhan, Yangxin Wu, Yunchao Wei, Michael C Kampffmeyer, Xiaoyong Wei, Minlong Lu, Yaowei Wang, and Xiaodan Liang. M5product: Selfharmonized contrastive learning for e-commercial multimodal pretraining. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21252–21262, 2022. 7
- [16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 3, 4
- [17] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111:98–136, 2015. 14

- [18] Fartash Faghri, David J Fleet, Jamie Ryan Kiros, and Sanja Fidler. Vse++: Improving visual-semantic embeddings with hard negatives. In *Proceedings of the British Machine Vision Conference.*, 2017. 6
- [19] Geoff French, Michal Mackiewicz, and Mark Fisher. Selfensembling for visual domain adaptation. In *International Conference on Learning Representations*, 2018. 5
- [20] Zhe Gan, Linjie Li, Chunyuan Li, Lijuan Wang, Zicheng Liu, Jianfeng Gao, et al. Vision-language pre-training: Basics, recent advances, and future trends. *Foundations and Trends® in Computer Graphics and Vision*, 14(3–4):163–352, 2022.
- [21] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6904–6913, 2017. 14, 15
- [22] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 5, 15
- [23] Micah Hodosh, Peter Young, and Julia Hockenmaier. Framing image description as a ranking task: data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 2013. 15
- [24] Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 831–839, 2019. 2, 3, 7, 8, 13, 16
- [25] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6700–6709, 2019. 14, 15
- [26] Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, Guoxing Yang, Jingyuan Wen, Heng Zhang, Baogui Xu, Weihao Zheng, et al. Wenlan: Bridging vision and language by large-scale multi-modal pre-training. arXiv preprint arXiv:2103.06561, 2021. 5
- [27] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021. 2
- [28] 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, 3, 7, 8, 12, 13, 16
- [29] 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. 14

- [30] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017. 14, 15
- [31] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. **3**, **4**, **14**
- [32] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 3, 4, 14
- [33] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888– 12900. PMLR, 2022. 1, 2, 6, 15
- [34] Junnan Li, Ramprasaath Selvaraju, Akhilesh Gotmare, Shafiq Joty, Caiming Xiong, and Steven Chu Hong Hoi. Align before fuse: Vision and language representation learning with momentum distillation. 2021. 2, 5, 6, 7, 15
- [35] Xiujun Li, Xi Yin, C. Li, X. Hu, Pengchuan Zhang, Lei Zhang, Longguang Wang, H. Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. *Proceedings of the European Conference on Computer Vision.*, 2020. 2
- [36] 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. 3
- [37] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelli*gence, 40(12):2935–2947, 2017. 2, 3, 7, 8, 13, 16
- [38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer, 2014. 14, 15
- [39] Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. In *Conference on Robot Learning*, pages 17–26. PMLR, 2017. 14, 15
- [40] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. arXiv preprint arXiv:1711.05101, 2017. 6
- [41] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. Advances in neural information processing systems, 32, 2019. 2
- [42] Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022. 3
- [43] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. arXiv preprint arXiv:1306.5151, 2013. 14

- [44] Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition, pages 7765–7773, 2018. 2, 3
- [45] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages 109–165. Elsevier, 1989. 2
- [46] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 14
- [47] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In 2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing, pages 722–729. IEEE, 2008. 14
- [48] Vicente Ordonez, Girish Kulkarni, and Tamara Berg. Im2text: Describing images using 1 million captioned photographs. Advances in neural information processing systems, 24, 2011. 14, 15
- [49] Ariadna Quattoni and Antonio Torralba. Recognizing indoor scenes. In 2009 IEEE conference on computer vision and pattern recognition, pages 413–420. IEEE, 2009. 14
- [50] 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, 2, 6, 15
- [51] 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, 3, 7, 8, 13, 16
- [52] Negar Rostamzadeh, Seyedarian Hosseini, Thomas Boquet, Wojciech Stokowiec, Ying Zhang, Christian Jauvin, and Chris Pal. Fashion-gen: The generative fashion dataset and challenge. arXiv preprint arXiv:1806.08317, 2018. 14, 15
- [53] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 14, 15
- [54] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. arXiv preprint arXiv:1606.04671, 2016. 3
- [55] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. arXiv preprint arXiv:2111.02114, 2021. 14, 15
- [56] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pages 4548–4557. PMLR, 2018. 3
- [57] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, im-

age alt-text dataset for automatic image captioning. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2556–2565, 2018. 3, 14, 15

- [58] Vishal Thengane, Salman Khan, Munawar Hayat, and Fahad Khan. Clip model is an efficient continual learner. arXiv preprint arXiv:2210.03114, 2022. 4
- [59] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016. 14, 15
- [60] Hugo Touvron, Matthieu Cord, Douze Matthijs, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. *International Conference on Machine Learning*, 2020. 4
- [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 2
- [62] Jeffrey Scott Vitter. Random sampling with a reservoir. ACM Transactions on Mathematical Software, 1985. 13
- [63] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. 14
- [64] Liyuan Wang, Mingtian Zhang, Zhongfan Jia, Qian Li, Chenglong Bao, Kaisheng Ma, Jun Zhu, and Yi Zhong. Afec: Active forgetting of negative transfer in continual learning. Advances in Neural Information Processing Systems, 34:22379–22391, 2021. 3, 7, 8, 13, 16
- [65] 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. 2
- [66] Xiao Wang, Guangyao Chen, Guangwu Qian, Pengcheng Gao, Xiao-Yong Wei, Yaowei Wang, Yonghong Tian, and Wen Gao. Large-scale multi-modal pre-trained models: A comprehensive survey. arXiv preprint arXiv:2302.10035, 2023. 13
- [67] 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. 4
- [68] Tobias Weyand, Andre Araujo, Bingyi Cao, and Jack Sim. Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2575–2584, 2020. 7
- [69] 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. 3
- [70] Peng Xu, Xiatian Zhu, and David A Clifton. Multimodal learning with transformers: A survey. arXiv preprint arXiv:2206.06488, 2022. 2

- [71] S YANG, P SUN, Y JIANG, X XIA, and P Luo. Objects in semantic topology. In *International Conference on Learning Representation (ICLR)*, 2022. 8
- [72] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. 14, 15
- [73] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. 2
- [74] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International conference on machine learning*, pages 3987–3995. PMLR, 2017. 3, 7, 8, 13, 16
- [75] Xunlin Zhan, Yangxin Wu, Xiao Dong, Yunchao Wei, Minlong Lu, Yichi Zhang, Hang Xu, and Xiaodan Liang. Product1m: Towards weakly supervised instance-level product retrieval via cross-modal pretraining. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11782–11791, 2021. 7
- [76] 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
- [77] Pengchuan Zhang, Xiujun Li, Xiaowei Hu, Jianwei Yang, Lei Zhang, Lijuan Wang, Yejin Choi, and Jianfeng Gao. Vinvl: Revisiting visual representations in vision-language models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5579– 5588, 2021. 2
- [78] Zekang Zhang, Guangyu Gao, Zhiyuan Fang, Jianbo Jiao, and Yunchao Wei. Mining unseen classes via regional objectness: A simple baseline for incremental segmentation. Advances in Neural Information Processing Systems, 35:24340–24353, 2022. 2
- [79] Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. Decoupled knowledge distillation. In *Proceedings of* the IEEE/CVF Conference on computer vision and pattern recognition, pages 11953–11962, 2022. 6
- [80] Da-Wei Zhou, Qi-Wei Wang, Zhi-Hong Qi, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Deep class-incremental learning: A survey. arXiv preprint arXiv:2302.03648, 2023. 2

Appendix

This supplementary document mainly provides more information about our P9D dataset and implementation details of the baseline methods. Besides, we provide the pseudocode of CTP and more experimental studies.

A. P9D Dataset.

A.1. Dataset Split.

Figure 5 shows the quantity distribution of each subset of our P9D. The different subsets have a consistent quantity distribution across tasks, and this consistent distribution ensures the comprehensive and unified evaluation for pretraining. Different from the training set, the test set (cross-modal retrieval evaluation) and query set (multi-modal retrieval evaluation) need to be further filtered by humans. The filter criterion is that the text describes the image content as accurately as possible while ensuring that the test/query set is proportional to the training set for the same category.

A.2. Image-Text Examples.

Figure 6 shows some image-text examples. We show some images of same class and keep one described caption for simplicity. It shows that real-world web data is noisy and multi-domain mixing. There are prevalent and complicated situations in the web image domain, such as complex backgrounds, amorphous watermarks, irrelevant objects, and occlusion.

B. Details of Baseline Methods.

Because these baseline methods are originally proposed for continual learning on the image classification task. Thus, we re-implement them to adapt the setting of visionlanguage pretraining. In addition to the replacement of the main optimization loss, we present the implementation details of each comparison method as follows:

B.1. Memory-Free methods

EWC [28] is the classical regularization methods. It maintains the old model parameters θ_{t-1} and an important matrix Ω with the same scale as the model. EWC builds an additional regularization loss to remember the old parameters according to the important matrix. Because the model θ_{t-1} at the last task stores the old knowledge, consolidating important parameters can fix the knowledge from being forgotten. The training loss can be formulated as:

$$\mathcal{L}_{EWC} = \mathcal{L}_{VLP} + \frac{1}{2}\lambda\Sigma_k\Omega_k(\theta_{t,k} - \theta_{t-1,k})^2, \quad (12)$$

where the $\theta_{t-1,k}$ denotes the k-th parameter after training last task data \mathcal{D}_{t-1} . Ω_k means the important weight of the



Figure 5: The quantity distribution of different task data is consistent for the four subsets of our P9D dataset.

k-th parameter and is calculated by the Fisher Information Matrix (FIM) in the EWC method.

SI [74] considers that the EWC is conducted at the end of each task and will ignore the optimization dynamics over the entire training trajectory. Thus, SI online estimates the importance weight Ω_k by its contribution (backward gradient) to the total loss variation. However, this online strategy need to backpass the gradient twice for each iteration. In the re-implement, we store the gradient of each parameter by retaining the forward graph.

MAS [2] calculate Ω_k by a unsupervised way. Specifically, It accumulates important measures based on the sensitivity of predictive results (output features) to parameter changes. In our re-implement, we sum the norm of the visual, textual, and multi-modal features as the predictive result to calculate the importance.

RWalk [7] combines the regularization terms of SI [74] and EWC [28] to integrate their advantages. In each iteration, Rwalk simultaneously consolidates the parameter by considering the online importance weight from SI method and the offline important weight from EWC method.

AFEC [64] proposes to actively forget the old knowledge that interferes with the learning of new tasks for continual learning. Specifically, It introduces the extra forward-step trained model θ_t^{\star} as the expansion and collaboratively guides the update of the current model with the EWC method. Similar to EWC, the training loss can be formulated as:

$$\mathcal{L}_{\mathcal{AFEC}} = \mathcal{L}_{VLP} + \frac{1}{2}\lambda\Sigma_k\Omega_k(\theta_{t,k} - \theta_{t-1,k})^2 + \frac{1}{2}\lambda_e\Sigma_k\Omega_k^*(\theta_{t,k} - \theta_{t,k}^*)^2,$$
(13)

where θ_t^{\star} is the parameter of forward-step trained model and λ_e is the FIM of θ_t^{\star} .

LWF [37] aligns the representations of previous-step and current models for all new arriving data. We maintain one reference model whose parameters are copied from the previous-step trained model and align the image and text representation of the reference and current model by the cross-entropy loss for each iteration.

B.2. Memory-Buffer methods

For the memory updating processing, the replay buffer will delete some old samples and add some new samples according to the size of the new task data.

ER [8] is a popular sample selection strategy. It uses the reservoir sampling [62] randomly stores a fixed number of training samples for each input batch and each sample has the same probability of being replaced.

Kmeans [8] use the Kmeans clustering to process all samples of the current task and set the number of clusters to the number of corresponding replaced samples. Then the cluster-center samples are chosen to update the buffer.

MoF [51] is first proposed by ICARL [51] and selects samples that are closest to the feature mean of each class. Because vision-language pretraining has no class concept, we choose the samples that are closest to the multi-modal feature mean of the current task.

ICARL [51] perform knowledge distillation on both buffer samples and new samples. The sample selection strategy is Mean-of-Feature (MoF). In our implementation, we combine the LWF term to optimize the current model.

LUCIR [24] proposes to utilize a cosine classifier to avoid the influence of the biased classifier and encourage similar feature orientation of the new and previous-step models. In our implementation, we replace the regular projected linear layer with the cosine normalizing linear layer. In addition, we constrain that the similarity of same-modal embedding from new and old models is big as possible. However, the inter-class distance constraint of the original paper [24] cannot be re-implemented because there is no class label in vision-language pretraining.

C. Dataset Comparison.

In Table 4, we present the comparison of our P9D with popular datasets from the continual learning domain [12] and multi-modal domain [66]. We observe that traditional continual learning datasets have a small number of data samples with limited classes (mostly at the thousand level), and only for single-target class labeling without detailed text description. In addition, although existing multi-modal datasets contain a large number of web image-text pairs,

Dataset	Train Samples	Categories	Modal	Objects	Continual Task Split	URL		
Popular Class-Incremental Learning Dataset								
Oxford Flowers [47]	2,040	102	image	single	yes	Link		
VOC Actions [17]	3,102	11	image	single	yes	Link		
MIT Scenes [49]	5,360	67	image	single	yes	Link		
CUB200-2011 [63]	5,994	200	image	single	yes	Link		
FGVC-Aircraft [43]	6,666	100	image	single	yes	Link		
Letters [11]	6,850	52	image	single	yes	Link		
Stanford Cars [29]	8,144	196	image	single	yes	Link		
SVHN [46]	73,257	10	image	single	yes	Link		
CIFAR10 [31]	50,000	10	image	single	yes	Link		
CIFAR100 [31]	50,000	100	image	single	yes	Link		
MNIST [32]	60,000	10	image	single	yes	Link		
Tiny-ImageNet [12]	80,000	200	image	single	yes	Link		
ImageNet-100 [53]	130,000	100	image	single	yes	Link		
CORe50 [39]	120,000	50	image	single	yes	Link		
Popular Multi-Modal Dataset								
Flickr30K [72]	29,000	_	image-text	multi	no	Link		
COCO [38]	113,287	80	image-text	multi	no	Link		
Visual Genome [30]	108K	_	image-text	multi	no	Link		
FashionGen [52]	325,536	_	image-text	multi	no	Link		
SBU [48]	875K	-	image-text	multi	no	Link		
GQA [25]	1M	_	image-text	multi	no	Link		
VQA v2.0 [21]	1.1M	-	image-text	multi	no	Link		
CC3M [57]	3.1M	-	image-text	multi	no	Link		
CC12M [6]	12M	-	image-text	multi	no	Link		
YFCC-100M [59]	100M	-	image-text	multi	no	Link		
LAION-400M [55]	400M	-	image-text	multi	no	Link		
Our: P9D	1,014,599	3,814	image-text	multi	yes	-		

Table 4: The overview of datasets about continual learning and vision-language pretraining domains. 'Categories' means the number of classes in the corresponding dataset and '-' means not mentioned. 'Objects' means the number of labeled/described objects in images. 'Continual Task Split' means the dataset contains different data chunks with discrepant semantic concepts and supports to simulate the continual environment. 'URL' means the hyperlink of corresponding dataset websites.

their data are too noisy and mixed to conform to the data split for continual tasks. In contrast, our P9D contains abundant image-text pairs to support vision-language pretraining. Besides, Each task contains rich semantic concepts, and different generalized semantic domains. It can support the simulation of continual learning environments.

C.1. Class-Incremental Learning Datasets

Oxford Flowers [47], MIT Scenes [49], CUB200-2011 [63], Stanford Cars [29], FGVC-Aircraft [43], VOC Actions [17], Letters [11], SVHN [46]. Aljundi *et al.* [3, 2] propose to use a sequence of 8 highly diverse recognition tasks as continual tasks. This sequence is composed of 8 different topics, going from flowers, scenes, birds, and cars, to aircrafts, actions, letters, and digits.

CIFAR10/100 [31] consists of 60,000 32×32 color images in 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images. The CIFAR100 dataset has 100 classes containing 600 images each. There are 500 training images and 100 testing images per class and the 100 classes can be grouped into 20 superclasses.

MNIST [32] is a large handwritten digits dataset. It has 60,000 samples as the training set and 10,000 samples as the test set.

Tiny-ImageNet [12] first used in the study of continual learning by Matthi *et al.* [12]. This is a subset of 200 classes from ImageNet [13] and the image size is rescaled to 64×64 . Each class contains 500 samples subdivided into train-

ing (80%) and validation (10%), and 50 samples for evaluation.

ImageNet-100 (SubImageNet) [53] is a 100-class random sample subset of ImageNet. It contains 130,000 images for training and 5,000 images for testing.

CORe50 [39] is a collection of 50 objects collected in 11 distinct domains, where 8 of them (120,000 samples) are used for training, and the rest are used as a single test set (45,000).

C.2. Multi-modal Datasets

Flickr30K [72] is obtained by extending the corpus of Hodosh *et al.* [23] and the image topic contain everyday scenes and activities. There are 31,783 images associated with five manually annotated captions each, and 29,000 images are used for training.

COCO [38] is built based on MSCOCO dataset [38]. It consists of 123,287 images and each image is annotated with 5 captions. There are 113,287 training images, 5000 test images, and 5000 validation images. COCO and Flickr30K datasets are often used as the retrieval evaluation dataset for large-scale vision-language pretraining.

Visual Genome [30] is proposed to help to develop of visual understanding tasks (*i.e.* image caption and visual question answering, *etc.*) by mining the relationships between objects. The dataset contains more than 108K images and each image has about 35 objects, 26 attributes, and 21 pairwise relationships.

FashionGen [52] contains 325,536 1360×1360 fashion images and each image has a paragraph-length caption as the description. Six different angles are photographed for all fashion items.

SBU [48] is collected and filtered from Flickr.com. It is usually used as the subset of vision-language pretraining [34, 33, 9].

GQA [25] is a balanced dataset with 1.7M samples which is mainly proposed for visual reasoning and compositional question answering.

VQA v2.0 [21] is proposed to reduce the language biases that existed in previous VQA datasets. It consists of around 1.1M image-question pairs and 13M corresponding answers based on 200K MSCOCO images.

CC3M [57] is a dataset annotated with conceptual captions and the image-text samples are mainly collected from the web. It contains about 3.3M image-description pairs.

CC12M [6] is a product of the urgent need for large-scale data with rapidly developing vision-language pre-training. The authors of CC3M relax the image-text filters and obtain

Method	TR@1	IR@1	Rm	mAP
only θ^{t-1}	41.95	42.23	63.57	62.06
only θ^t	40.41	40.34	62.32	62.32
m=0.7	43.64	43.04	64.83	62.95
m=0.8	43.68	43.11	64.74	62.36
m=0.9	43.43	43.39	64.87	62.64
m=0.99	43.50	43.01	64.75	62.23
m=0.995	44.27	42.23	65.04	61.63

Table 5: The results of momentum selection experiment.

the larger dataset CC12M.

YFCC-100M [59] totally contains 100 million media objects (99.2 million photos, 0.8 million videos) collected from Flickr.com.

LAION-400M [55] is filtered using pre-trained CLIP [50] and contains 400 million image-text pairs.

D. Algorithm

The Alg. 1 shows the training pipeline of our CTP in the task data $D_t \in \{D_1, D_2, ..., D_T\}$.

E. More Experiments.

E.1. Momentum Setting.

For the first task (t = 0), There is no previous-step model and the current model has not adapted to the product domain. Thus, the momentum m of the first task is set to 0.995, and we keep it the same for all ablation studies. Because we find that the training loss oscillates and fails to converge if m is 0.9 in the first task.

For the following tasks, we set m as 0.9 and we also do the parameter-sensitive study about m. The results of Table 5 show the training on the $\{1, 2, ..., T\}$ task is not sensitive to the setting of compatible momentum m. We suspect this is due to the fact that the model accepts parameters from both the previous-step and current models and is less prone to biased updates. In addition to the main vision-language pretraining loss, the compatible momentum contrastive loss is an auxiliary loss for continual learning. Thus, the model is more robust to momentum parameter selection and does not easily collapse [22].

E.2. Reverse Task Order.

In the main text, all experiments are conducted in the default task order. To study the impact of task order on the performance ranking, we supplement a check experiment with the reversed task order ¹. The Table 6 shows the results of all baselines and our method in the reversed task order.

¹ Electronics, Outdoor, Parenting, Auto, Clothing, Beauty, Food, Furnishings, and Household

Algorithm 1 Pseudocode of CTP in a PyTorch-like style.

```
# F, M, R: training, momentum, and reference (previous-task) model
# m, t, q_v, q_t: momentum, temperature, visual and textual queues
M.params, R.params = F.params, F.params # initialize momentum and reference model
for (image, text) in loader: # load a minibatch with N image-text pairs
feat_v, feat_t = F.get_featuere(image, text)
    ita_loss = CE(S_i2t/t,eye_like(S_i2t))+CE(S_t2i/t,eye_like(S_t2i)) # image-text contrastive loss
    mlm_loss = CE(multimodal_out, labels=text.labels, mask=text.mask) # mask language modeling loss
loss = ita_loss/2 + mlm_loss # conventional loss of the current task
    # compatible momentum update
M.params = m*M.params+(1-m)/2*F.params+(1-m)/2*R.params
    feat_vm, feat_tm = M.get_featuere(image, text)
    multimodal_out_m = M.multimodal_fusion(image, text, text.mask)
    enqueue(q_v, feat_vm.detach(), q_t, feat_tm.detach()) # enqueue current features
S_i2t_m, S_t2i_m = feat_v @ q_t.T, feat_t @ q_v.T
    # compatible momentum contrast
ita_loss_m = CE(S_i2t_m/t, eye_like(S_i2t_m)) + CE(S_t2i_m/t, eye_like(S_t2i_m))
mlm_loss_m = CE(multimodal_out,labels=multimodal_out_m.logits,mask=text.mask)
    loss += ita_loss_m/2 + mlm_loss_m
    dequeue(q_v, q_t) # dequeue earliest features
       topology preservation
    # topology preservation
feat_vr, feat_tr = R.get_featuere(image, text)
S_i2t_r, S_t2t_r = feat_vr @ feat_tr.T, feat_tr @ feat_vr.T
S_i2i_r, S_t2t_r = feat_vr @ feat_vr.T, feat_tr @ feat_tr.T
loss_sm = CE(S_i2i/t, S_i2i_r/t, mask=eye_like(S_i2i)) + CE(S_t2t/t, S_t2t_r/t, mask=eye_like(S_t2t))
    loss_cm = CE(S_i2t/t, S_i2t_r/t) + CE(S_t2i/t, S_t2i_r/t)
    loss += loss_sm + loss_cm
    # parameter updat
```

```
loss.backward()
update(F.params)
```

CE: cross entropy loss; eye_like: create an identity matrix with the same size as input.

Methods	Cross-modal Retrieval							Multi-modal Retrieval			
	TR@1	TR@5	TR@10	IR@1	IR@5	IR@10	Rm	mAP@1	mAP@5	mAP@10	
JointT	60.72	86.05	91.74	61.98	86.82	91.85	79.86	64.07	70.09	67.33	
Memory-Free											
SeqF	37.81	64.69	74.00	38.05	64.23	74.46	58.87	61.86	68.08	65.18	
SI [74]	38.51	64.41	74.84	38.90	65.14	74.14	59.32	61.56	67.63	64.69	
MAS [2]	39.81	66.97	75.86	40.86	66.87	75.93	61.05	61.65	67.78	65.23	
EWC [28]	39.11	67.96	77.09	41.46	68.73	77.30	61.94	62.17	68.11	65.22	
AFEC [64]	40.13	68.17	77.62	41.57	68.69	76.99	62.19	61.60	67.66	64.96	
LWF [37]	41.18	67.29	76.39	39.81	67.81	76.32	61.31	61.76	68.12	65.20	
RWalk [7]	39.04	67.85	78.00	40.20	68.94	77.65	61.95	62.40	68.43	65.60	
Our:CTP	45.96	73.47	80.85	44.98	72.34	80.25	66.31	61.08	67.20	64.19	
Memory-Buffe	er										
MoF [51]	43.92	72.28	80.99	45.01	73.19	81.03	66.07	61.37	67.65	64.79	
LUCIR [24]	45.36	72.91	80.92	45.61	73.68	80.74	66.54	61.89	67.92	65.31	
ER [8]	44.59	72.87	81.20	45.92	73.05	80.89	66.42	62.32	68.35	65.42	
Kmeans [8]	45.19	74.42	81.83	46.03	73.05	80.67	66.53	62.73	68.35	65.39	
ICARL [51]	47.33	75.65	83.63	47.61	75.90	83.24	68.89	62.54	68.54	65.87	
Our:CTP+ER	51.05	79.20	86.23	51.30	78.60	85.45	71.97	61.58	67.65	64.78	

Table 6: The final cross-modal and multi-modal retrieval performance comparison when conducting the reversed task order.

The result shows that although there are some changes in the ranking of some methods with similar performance, the overall performance ranking is still consistent with the performance ranking of default task order. Additionally, our method CTP exhibits superior performance in both continual learning scenarios (Memory-Free and Memory-Buffer), even when the order of tasks is changed. It indicates that the task order can affect the performance value of final result but not the performance ranking of our method. Our method consistently outperforms in different task order settings.

F. License

Our P9D dataset is released under CC BY-NC-SA 4.0 license and can freely be used for non-commercial purposes. The collection of data has obtained permission from the relevant websites. Once a conflict of interest, our group reserves all the rights for the final explanation.



KAMJOVE T-57 thermal thermostatic electric kettle 304 stainless steel tea art special boiling water teapot

Finance office solar real person voice 12 big keystroke calculator Summer sun-shading fashion women's neck protection scarf/ veil/mask

Nuts and preserved fruit assorted snacks original 2 cans SMACO CROSS men's business travel leisure computer bag backpack

Velbon EX-MACRO tripod set, MINI tripod, SLR camera tripod



Figure 6: Some examples of our dataset. The first and second rows are the corresponding image-text pair. For simplicity, the rest rows show the images from the same class.