

# CluMo: Cluster-based Modality Fusion Prompt for Continual Learning in Visual Question Answering

Anonymous authors

Paper under double-blind review

## Abstract

Large vision-language models (VLMs) have shown significant performance boost in various application domains. However, adopting them to deal with several sequentially encountered tasks has been challenging because finetuning a VLM on a task normally leads to reducing its generalization power and the capacity of learning new tasks as well as causing catastrophic forgetting on previously learned tasks. Enabling using VLMs in multimodal continual learning (CL) settings can help to address such scenarios. To improve generalization capacity and prevent catastrophic forgetting, we propose a novel prompt-based CL method for VLMs, namely **Cluster-based Modality Fusion Prompt (CluMo)**. We design a novel **Key-Key-Prompt** pair, where each prompt is associated with a visual prompt key and a textual prompt key. We adopt a two-stage training strategy. During the first stage, the single-modal keys are trained via  $K$ -means clustering algorithm to help select the best semantically matched prompt. During the second stage, the prompt keys are frozen, the selected prompt is attached to the input for training the VLM in the CL scenario. Experiments on two benchmarks demonstrate that our method achieves SOTA performance.

## 1 Introduction

Visual Question Answering (VQA) is a complicated task, where the goal is to answer questions described in natural language (text) about a given input image. Addressing VQA requires understanding and fusion of information from both the visual and textual domains to generate accurate responses. Recently, significant advancements in addressing VQA tasks have emerged due to the development of pre-trained large vision-language models (VLMs) Radford et al. (2021); Kim et al. (2021). Despite these advances, one of the persistent challenges in VQA tasks is the ability to adapt a VLM in CL setting to avoid finetuning a copy of an underlying VLM per task. In a CL setting, we learn new tasks and aim for continuously improving the model performance without forgetting previously learned knowledge, also known as catastrophic forgetting French (1999). To address catastrophic forgetting, a group of CL algorithms are deployed. Regularization-based methods Kirkpatrick et al. (2017); Li & Hoiem (2017) constrain the drastic parameter shift when learning new tasks. Expansion-based methods Douillard et al. (2022); Cai et al. (2023) expand the model with small portion of additional weights and use the expanded weights to learn the new incoming tasks. Rehearsal-based methods Rebuffi et al. (2017); Rolnick et al. (2019) store a representative subset of the training dataset for each task into a small memory buffer and replay them back during the learning of the current task to maintain the encoded knowledge of the previously learned tasks. More recently, prompt-based methods Wang et al. (2022b;a) aim to use prompts that contains task-specific or semantic-specific information to prevent catastrophic forgetting. A prompt is attached to the embedded features of the input to adapt the model to focus on the specific characteristics of the input task that has been learned before.

Most existing CL methods consider unimodal, i.e., vision-only or language-only, settings and hence are inapplicable to address VQA tasks. To tackle this shortcoming, we propose a novel two-stage prompt learning-based CL method, namely cluster-based modality fusion prompt (CluMo). Figure 1 visualizes the high-level idea of our approach. Our method adopts a pre-trained VLM as its backbone and benefits from a clustering-based modal-specific key strategy to boost generalization capacity and minimize catastrophic

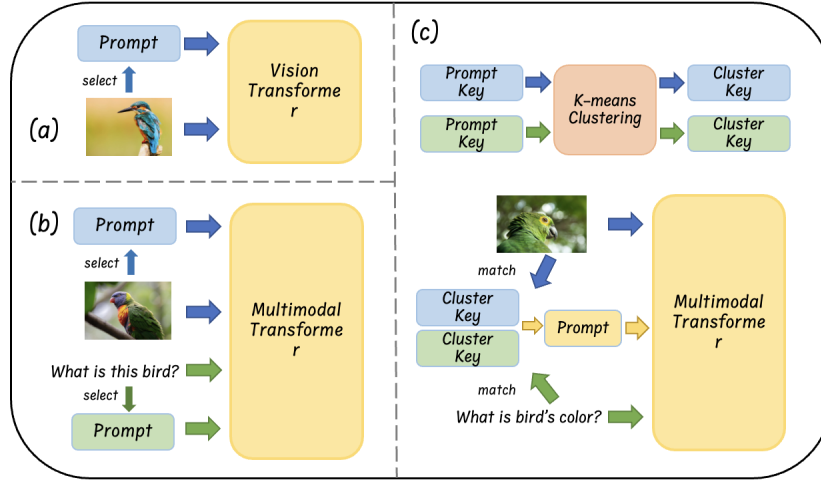


Figure 1: Comparison between existing prompt-based CL methods and our proposed method: **(a)** Uni-modal based methods use image feature to select prompts from a prompt pool. **(b)** Multi-modal based methods use the image features to select image prompts and use the text features to select the text prompts. **(c)** We first train the prompt key using a clustering algorithm to form a cluster key and use the combination of the cluster key from both modalities to select the fusion prompt.

forgetting. More specifically, we use a clustering-based algorithm to train visual-prompt keys and textual-prompt keys during the first stage. During the second stage, we assign each input image-question pair with well-trained prompt keys to its corresponding visual keys and textual keys. We then use the combination of two modal-specific keys to find the best-matched prompt to adapt the model for the input task. We also benefit from knowledge distillation during training to further improve the performance. Our proposed method outperforms existing alternative methods. Our specific contribution includes:

- We propose a novel clustering-based prompt learning method for training VLMs in CL settings to address VQA tasks with vision-and-language inputs.
- We use a two-stage training strategy to train the prompt keys before training the whole model to guarantee optimal prompt selection that is necessary for generalization on the input.
- We offer extensive experiments to demonstrate that the proposed approach achieves SOTA performance against CL existing methods and offer insight about the reason of this improved performance.

## 2 Related Works

**Visual Question Answering** Visual Question Answering (VQA) has been a pivotal task at the intersection of computer vision and natural language processing which led to advances on more complex tasks. Initially, VQA was formulated as a classification task in which answers are selected from a predefined set of answers Agrawal et al. (2016) and was solved by using CNNs for image feature extraction and RNNs for text processing. These models were too simple to be used in most practical cases. With the development of transformer and BERT-like models Lu et al. (2019); Li et al. (2019), performance in VQA tasks has significantly been improved due to the better capacity of capturing the intricate relationship between two modalities and generating the response in the form of meaningful and descriptive texts. Despite these advances, VQA tasks are mostly studied in static settings Goyal et al. (2017); Johnson et al. (2016); Marino et al. (2019), where it is assumed that there is a single input task. As a result, most existing methods are inapplicable in dynamic environments and settings such as continual learning (CL) when we encounter several tasks over time.

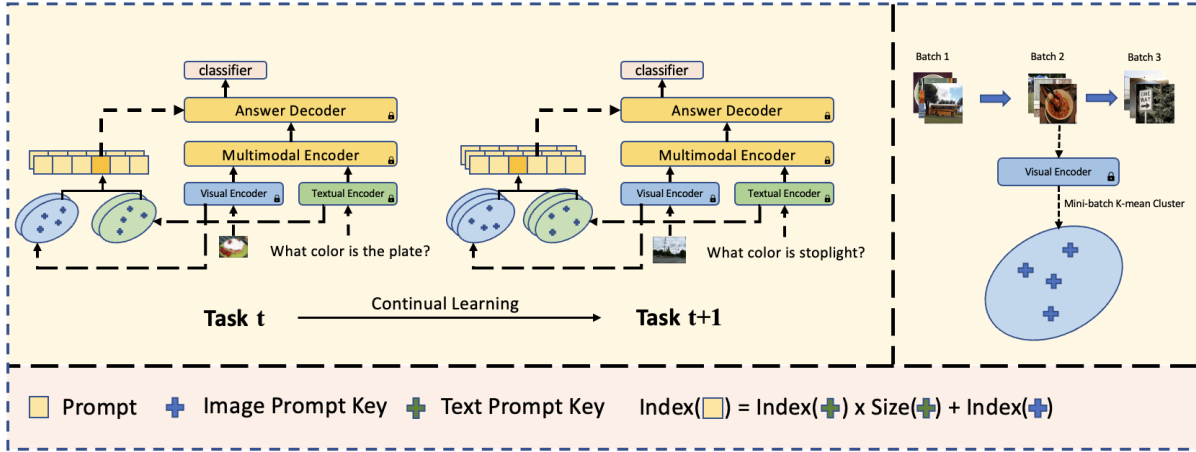


Figure 2: Block diagram of the proposed approach: **Left:** the backbone contains a pre-trained frozen visual encoder, a textual encoder, and a multimodal encoder. The answer decoder shares the same architecture as multimodal encoder. During the training phase, the fixed number of visual prompt keys, textual prompt keys, and a prompt pool will be added for each new task. **Right:** the procedure of visual prompt key training consists of training the modal-specific prompt key by a sequence of randomly selected batches of training data from current task until convergence is reached. Same procedure for textual prompt key.

**Prompt-Based Learning** Prompt learning is a powerful technique for leveraging pre-trained language models to frame downstream tasks in NLP. It is more memory-efficient than using Adapters Pfeiffer et al. (2021) or LoRA Hu et al. (2021) and has been used successfully to guide responses of VLMs for a particular task. The reason is that additional prompt parameters are in the form of small embeddings which are directly added to the input to make the model task. As a result, the prompts are often much smaller in size compared to the layers of the model, leading to a minimal increase in the total number of parameters compared to using adapters or LoRA. Moreover, prompts can be stored in a memory-efficient way and retrieved dynamically based on the observed task in the input. This means that even when the model handles many tasks, the memory overhead remains low. Brown et al. 2020 introduced the concept of prompt for the natural language instruction task to guide the model towards desired outputs. Prompt learning is based on providing a fixed function to condition a model so that it gets extra information token which specializes it to perform the down-stream task. Prompts are mostly considered as trainable parameters, either task-specific or domain-specific, to guide the model by obtaining task-specific knowledge Lester et al. (2021); Li & Liang (2021). Prompt learning has also been found helpful in handling a single VQA task Hu et al. (2023).

**Prompt Learning for Continual Learning** Prompt learning has been used in CL to prevent catastrophic forgetting when a large pre-trained models is trained on a stream of sequentially encountered tasks. It allows a single model to quickly adapt to new tasks in the stream without needing extensive retraining. By using task-specific prompts, the model can retain and recall knowledge from earlier tasks, mitigating the issue of catastrophic forgetting. It also allows scalability to learning a large number of tasks since each task primarily requires learning or generating new prompts rather than retraining the entire model. L2P Wang et al. (2022b) pioneered to connect prompt-based learning and CL. Instead of having a single shared prompt to learn all tasks, L2P introduced the concept of “prompt pool” to maintain prompts for different tasks independently from each other. DualPrompt Wang et al. (2022a) extended the idea of prompt pool in l2p by introducing E-prompt and G-prompt. While E-prompt is task-specific, G-prompt encodes the knowledge used for all tasks to further allow knowledge sharing and transferring while mitigating negative transfer. S-Prompt Wang et al. (2023) applied clustering to build the prompt pool with domain-specific prompts. These prompt learning methods for CL only consider single-modality, i.e., vision-only or text-only, and hence are sub-optimal for tasks with multi-modal inputs such VQA when the modalities are related. Our method benefits from the specific properties of multi-modal data to address VQA in CL settings using prompt learning and leads to performance improvements against these methods.

### 3 Problem Description

Consider a set of VQA tasks,  $\{\mathcal{T}_i\}_{i=1}^T$ , which are encountered sequentially and each of them is from different domain. For each of the tasks, a labeled training dataset  $\mathcal{D}^i = \{(\langle \mathbf{I}_i^j, \mathbf{L}_i^j \rangle^i, y_i^j)_{j=1}^{N_i}\}$  is accessible,  $N_i$  denotes the size of dataset,  $\mathbf{I}_i^j \in \mathbb{R}^{H \times W \times C}$  denotes the input image,  $\mathbf{L}_i^j \in \mathbb{R}^{L \times |V|}$  denotes the input text, and  $y_i^j$  denotes the text-typed discrete label. The order in which the VQA tasks are observed is not known in advance and the training data points are assumed to be drawn iid from a task-specific joint distribution  $p_i^t(\cdot, \cdot, \cdot)$ . Upon learning each task, the model moves forward to learn the next task. Since all the previously learned tasks can be encountered at any time during testing in the future, the model should learn new tasks such that its knowledge of previously learned tasks is maintained, i.e., by preventing catastrophic forgetting.

We formulate our problem in a **domain-incremental** learning setting Van de Ven & Tolias (2019) which assumes all the tasks are from different domains and the boundaries between them are known during learning time. We consider that each task can be learnt individually by adapting a pre-trained large multimodal transformer  $f_{\theta_M}^i(\cdot, \cdot)$  via minimizing a suitable discrimination loss  $\mathcal{L}$ , e.g., cross entropy. In our approach, all the model parameters, except the final classifier layer  $\theta_{cls}$ , are frozen during training to preserve the generalizability of the model. We benefit from prompt learning to enable using a single model to learn all tasks. To prevent catastrophic forgetting, a trainable task-specific prompt pool, which contains several task-specific prompts, is attached to the model  $f_{\theta_M}^i(\cdot, \cdot)$  such that the best-semantically-matched prompt is selected based on image and text inputs for task specialization. The prompt is then pre-pended to the input vectors so that the output is generated based on specialization. Our method is rehearsal-free and does not need any memory buffer similar to prior approaches Lopez-Paz & Ranzato (2022); Rebuffi et al. (2017).

### 4 Proposed Architecture

Our architecture, named cluster-based modality fusion prompt (**CluMo**), contains two unimodal task-specific cluster-based keys for vision and text embeddings and one prompt pool. The combination of the selections from both keys is then used to select the best matched prompt from the prompt pool. A high-level diagram of our approach is presented in Figure 2. In this section, we first introduce the preliminaries such as backbone model and prompt pool-based method in 4.1, and modality fusion prompt in Sec. 4.2, then the cluster-based prompt key is described in Sec. 4.3, and the training and the inference strategy is discussed in Sec. 4.4.

#### 4.1 Preliminary

**Backbone** The base multimodal transformer contains three encoders: the visual encoder  $VE$ , the textual encoder  $TE$ , and the multimodal fusion encoder  $FE$ . Given a visual input  $\mathbf{V}$ , i.e., a single image, and a textual input  $\mathbf{T}$ , i.e., a question, the data processing pipeline for the model is:

$$\hat{y}(\mathbf{V}, \mathbf{T}) = \mathcal{F}(FE([VE(\mathbf{V}); TE(\mathbf{T})])), \quad (1)$$

where  $\mathcal{F}(\cdot)$  is the classifier to predict the answer.

**Prompt Pool** As an adoption of prompt learning in continual learning, a prompt pool is a set of trainable key-value ( $K$ - $P$ ) pair, in which  $K \in \mathbb{R}^{1 \times D}$  denotes the ‘‘prompt key’’, and  $P \in \mathbb{R}^{L_p \times D}$  is the prompt.  $L_p$  and  $D$  denote the length and dimension of the prompt. Given an input image  $\mathbf{V}$ , we compute  $v_I = VE(\mathbf{V}) \in \mathbb{R}^{L_v \times D}$ , where  $L_v$  is the dimension of the features, after passing the image through the visual encoder.  $v_{I_0} = v_I[0]$  is matched with all the keys  $K$  within the prompt pool via similarity score, such as cosine similarity, to find the most similar  $K_i$ . The corresponding  $P_i$  is selected and prepend to  $\mathbf{V}$  as  $\mathbf{V}' = [P_i; \mathbf{V}]$ . Parameters of  $K$  and  $P$  are updated through back-propagation during the training. However, in our setting, we adopt a two stage training strategy that prompt keys are trained before model and prompts.

#### 4.2 Modality Fusion Prompt

Previous prompt-based CL methods such as L2P Wang et al. (2022b) associate each prompt in the prompt pool with a single prompt key to form Key-Value pair. In practice, the prompt keys in prompt-based CL

can be considered as cluster centers. These cluster encode a notion of similarity between the prompts. The input feature vectors that form a cluster in the feature space can be assigned to these cluster centers, which are prompt keys. The intuition behind this idea is that **feature vectors with small geometric distance in the feature space are semantically similar** Wang et al. (2023).

However, such a key-value pair design considers only single modality without tasks with multimodal inputs. The reason is that different input modalities contain different or complementary semantic information. Hence, having prompt keys that associate with each modality can help guiding prompt selection, which is more comprehensive and representative in term of semantic properties of each modality. Thus, we propose a task-specific prompt pool architecture, namely **Modality Fusion Prompt**, which is composed of the **visual prompt keys**  $K_v$ , the **textual prompt keys**  $K_t$ , and the **prompt pool**  $P$  as following:

$$\begin{aligned} K_t &= [K_{t_1}, K_{t_2}, \dots, K_{t_{S_t}}], \\ K_v &= [K_{v_1}, K_{v_2}, \dots, K_{v_{S_v}}], \\ P &= [P_1, P_2, \dots, P_{S_p}], \\ K_{t_m} &\in R^D, K_{v_n} \in R^D, P_l \in R^{L_p \times D}, \end{aligned} \quad (2)$$

where  $S_t$ ,  $S_v$ , and  $S_p$  are the sizes of textual prompt key, the visual prompt key, and the prompt pool, respectively.  $L_p$  is the length of each prompt and  $D$  is the hidden dimension of the transformer backbone. The prompt pool size  $S_p$  is then determined as  $S_p = S_v \times S_t$ . Each prompt is associated with the unique combination of one visual prompt key and one textual prompt key. Given a specific visual prompt key  $K_{v_m}$  and a specific textual prompt key  $K_{t_n}$ , the Key-Key-Value pair is defined as the following:

$$(K_{v_m}, K_{t_n}) \rightarrow P_{m*S_v+n}. \quad (3)$$

As modality fusion prompt is **task-specific**, new visual prompt keys, textual prompt keys and a prompt pool will be initialized for each of the new coming task. The previous ones are frozen during training.

### 4.3 Cluster-based Prompt Key

---

#### Algorithm 1 PromptKeyTraining

---

**Require:** Dataset  $D$ , Image Prompt Key Pool  $K_v$ , Text Prompt Key Pool  $K_t$ , Image Prompt Size  $S_I$ , Text Prompt Size  $S_T$

```

while Not Converge do
  Random Select batch of image  $I$ , text  $T$  from  $D$ 
   $\hat{v}_I = \text{mean}(VE(I), \text{dim} = 1)$ 
   $\hat{v}_T = \text{mean}(VT(T), \text{dim} = 1)$ 
   $Cluster_I = \text{dictionary}()$ 
   $Cluster_T = \text{dictionary}()$ 
  for  $i, t$  in  $v_{I_{mean}}, v_{T_{mean}}$  do
     $Key_{img} = \text{image key with top similarity}(i, K_v)$ 
     $Key_{txt} = \text{text key with top similarity}(t, K_t)$ 
     $Cluster_I[Key_{img}].\text{append}(i)$ 
     $Cluster_T[Key_{txt}].\text{append}(t)$ 
  end for
  for  $i$  in  $S_I$  do
     $K_v[i] = \text{mean}(Cluster_I[i])$ 
  end for
  for  $i$  in  $S_T$  do
     $K_t[i] = \text{mean}(Cluster_T[i])$ 
  end for
end while

```

---

$K$ -means clustering has been widely adopted in machine learning algorithms for semantic separation and understanding, where data from different domains can be explicitly separated via clustering in an unsuper-

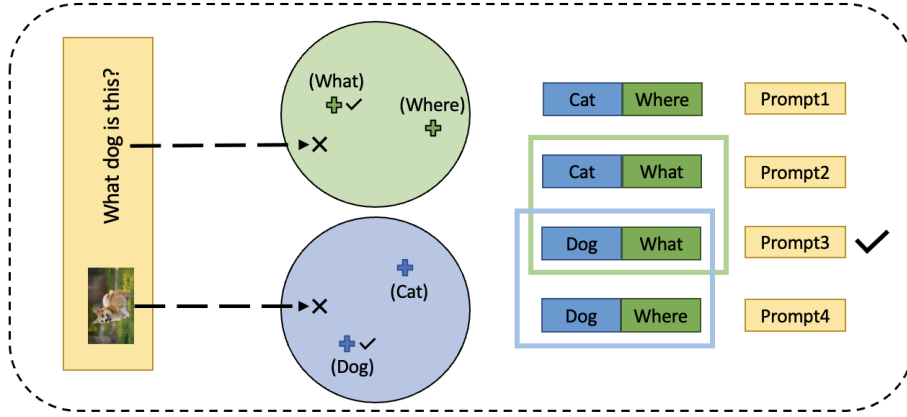


Figure 3: Naive Example of Prompt Selection. Consider a naive animal VQA dataset which only contains dog and cat images and questions only about "what" and "where". During the first stage training, the visual prompt keys and textual prompt keys are learnt to represent "dog"/"cat" and "where"/"what" respectively (in realistic settings the keys will learn the implicit cluster instead of explicit category). Given a test image-question pair, the image and question are projected to their modality-specific feature space through the encoders. The nearest prompt keys, which are keys represent "dog" and "what" are selected. The combination of the two prompt keys lead to "prompt3" which is finally selected.

vised way Wang et al. (2023) Cohn & Holm (2021). However, even though the data from single task belong to the same domain, they can still be further divided into sub-domains based on the semantic property. To make each prompt key be the semantically cluster center of the sub-domains for both vision and text inputs, we adopt mini-batch  $K$ -means clustering algorithm on prompt keys of  $K_v$  and  $K_t$  to make each prompt key diverse and representative. Let  $\mathcal{B} = (I, T)$  be the random batch from the training dataset. We extract the image feature vector  $v_I$  and the text feature vector  $v_T$  as follows:

$$v_I = VE(I), v_T = TE(T), \quad (4)$$

where  $v_I \in R^{B \times L_I \times D}$  and  $v_T \in R^{B \times L_T \times D}$ ,  $B$  is the batch size,  $L_I$  and  $L_T$  are the length of vectors for image and text features, represent the embedded image and text input respectively. For visual prompt key clustering, each image feature vector,  $v_{I_n}$ , is set by taking mean along second the dimension such that  $\hat{v}_{I_n} \in R^{B \times D}$ , and  $\hat{v}_{I_n}$  is used to compare with every prompt key in  $K_v$ :

$$similarity(n, m) = \|\hat{v}_{I_n} - K_{v_m}\|_2, \quad (5)$$

and the prompt key with highest similarity is assigned to match  $v_{I_n}$ . After calculation of the whole batch  $\mathcal{B}$ , the prompt keys are then updated by calculating the mean of all  $\hat{v}_{I_n}$  assigned to that specific visual prompt key. We repeat the above step until the convergence. The procedure of updating the text prompt key  $K_t$  is similar to updating the image prompt keys. Algorithm 1 summarizes our approach for prompt key training.

#### 4.4 Training and Inference

During training, we adopt a **two-stage training** strategy to ensure that the prompt keys are correctly settled before learning the current task. In the first stage of learning each task  $T_i$ , we random select batches from the current task's dataloader to train minibatch  $k$ -means Cluster on the visual and the textual prompt key  $K_v$  and  $K_t$  until reaching the convergence of the clustering algorithm. During the second stage, the trained  $K_v$  and  $K_t$  are frozen. Within the iteration of training dataloader, each training instance is assigned to its nearest prompt key using  $k$ -nearest neighbor (KNN) algorithm to find the best match prompt  $P_k$  from the prompt pool  $P$ .  $P_k$  is then attached to the model pipeline:

$$\hat{y}(\mathbf{V}, \mathbf{T}) = \mathcal{F}(FE([P_k; VE(\mathbf{V}); TE(\mathbf{T})])) \quad (6)$$

**Algorithm 2** Continual Learning Procedure**Require:** Datasets  $\mathbf{D}$ , Pretrained Model  $\mathbf{M}$ 


---

```

Freeze  $\mathbf{M}$  except  $\mathbf{M}$ .classifier
 $\mathbf{M}.\text{CluMoList} = \text{dict}()$ 
for dataset  $D_i$  in  $\mathbf{D}$  do
  Initialize new CluMo  $\mathbf{C}_i$ 
   $\mathbf{M}.\text{CluMoList}[D_i] = \mathbf{C}_i$ 
  PromptKeyTraining( $D_i, \mathbf{C}_i$ )
  Freeze visual key and textual key of  $\mathbf{C}_i$ 
  for Batch  $\mathbf{B}$  in  $D_i$  do
    for img  $i$ , txt  $t$  in  $\mathbf{B}$  do
      Select  $Key_{img}$  based on  $i$  from  $\mathbf{C}_i$ 
      Select  $Key_{txt}$  based on  $t$  from  $\mathbf{C}_i$ 
      Find prompt  $P$  by  $Key_{img}$  and  $Key_{txt}$  from  $\mathbf{C}_i$ 
    end for
    loss = Train( $\mathbf{M}, \mathbf{B}, P$ )
    loss.backward()
  end for
end for

```

---

During the second stage, we also use knowledge distillation to further boost the performance. Before the training of task  $T$ , we keep a frozen copy of model after finishing  $T - 1$ , denoted as  $\mathcal{M}_{T-1}$ . To prevent significant parameter shift, we pass the same input image and question to both  $\mathcal{M}_T$  and  $\mathcal{M}_{T-1}$  and add the difference between the two model’s output to the loss:

$$\mathcal{L}_{KD}(\mathbf{V}, \mathbf{T}) = \text{MSE}(\hat{y}_{\mathcal{M}_T}(\mathbf{V}, \mathbf{T}), \hat{y}_{\mathcal{M}_{T-1}}(\mathbf{V}, \mathbf{T})). \quad (7)$$

The final objective loss function would be:

$$\mathcal{L} = \mathcal{L}_{ce}(\hat{y}(\mathbf{V}, \mathbf{T}), y) + \mathcal{L}_{KD} \quad (8)$$

Where  $\mathcal{L}_{ce}$  is the same cross entropy loss. The overall training procedure for all tasks come in sequence is presented in Algorithm 2.

During inference, the model is frozen and we follow a procedure similar to the second stage of training. For every training image-text pair, the image input is aligned with the best-matched image prompt key while the text input is aligned with the best-matched text prompt key. The combination of prompt keys is deployed to find the corresponding prompt, which is pre-pend to the output of multimodal encoder. To help better understand the procedure of prompt selection, an visualization example is provided in Figure 3.

## 5 Experiments

Our implementation code is available as a supplement. Please refer to the code for reporducing the results.

### 5.1 Experiment Setup

**Backbone** We used the public pre-trained large multimodal transformer, ALBEF Li et al. (2021), as our backbone for VQA task. It consists of an image encoder, a text encoder, a multimodal encoder, which uses cross-attention between the two modalities. Specifically for VQA tasks, an pre-trained answer decoder is append after the multimodal encoder, which has same architecture as multimodal encoder.

**Baselines for comparison** We use seven methods for comparison. We include algorithms from major CL approaches. We include two regularization-based methods: **EWC** Kirkpatrick et al. (2017) and **LwF** Li & Hoiem (2017), two rehearsal-based methods: **ER** Rolnick et al. (2019) and **GEM** Su et al. (2021). We also include three SOTA prompt-based continual learning methods, **L2P** Wang et al. (2022b), **DualPrompt** Wang et al. (2022a), and **S-Prompt** Wang et al. (2023). We also include finetuning to demonstrate the

Method	CLOVE-scene											
	abcdef		dbafec		bdcafe		acbefd		caefdb		bafedc	
	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$
Finetune	34.03	34.28	34.89	34.99	38.83	21.65	34.45	35.14	34.42	34.47	33.95	35.67
EWC	37.49	28.04	37.00	29.10	37.95	27.46	37.99	29.68	37.13	28.63	37.83	27.97
LwF	38.18	26.82	35.03	32.84	37.31	29.11	37.85	29.87	37.94	28.15	38.21	27.48
ER	41.05	19.92	42.09	17.12	42.37	18.09	41.91	20.28	41.11	19.65	42.08	20.52
GEM	41.52	18.33	43.14	14.73	42.89	17.43	42.54	20.13	41.90	20.88	43.11	19.86
L2P	43.01	18.22	45.84	15.03	44.64	17.41	45.63	14.96	44.78	17.99	46.58	14.85
DualPrompt	45.51	15.86	46.58	13.49	45.83	16.48	46.27	15.45	46.21	15.89	47.01	13.16
S-Prompt	45.73	14.11	45.93	14.17	46.99	14.38	46.68	14.77	47.53	12.19	44.09	22.32
<b>CluMo</b>	<b>48.73</b>	<b>10.76</b>	<b>48.8</b>	<b>10.25</b>	<b>48.83</b>	<b>9.73</b>	<b>48.94</b>	<b>10.29</b>	<b>48.26</b>	<b>11.04</b>	<b>48.98</b>	<b>10.23</b>

Method	CLOVE-function											
	soarkl		rsaolk		osrlak		oarlks		skaolr		ksoarl	
	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$	A $\uparrow$	F $\downarrow$
Finetune	31.55	53.76	37.34	39.64	23.34	57.32	24.09	62.79	16.34	74.82	17.50	84.46
EWC	35.70	47.92	37.82	41.55	38.92	40.48	40.74	33.89	37.53	37.22	40.85	32.32
LwF	37.18	46.86	36.81	44.12	39.21	39.81	36.81	41.29	30.49	53.11	29.17	55.84
ER	42.22	32.97	39.78	38.62	41.22	35.79	37.14	33.38	33.41	48.99	38.23	38.01
GEM	44.58	30.87	41.43	29.46	40.87	32.98	39.81	28.77	36.88	39.14	40.26	31.87
L2P	44.80	16.38	43.39	21.26	43.27	21.97	42.54	19.18	40.4	31.92	43.37	24.19
DualPrompt	45.01	15.90	44.26	17.43	44.66	18.50	43.69	15.31	39.32	34.78	45.65	20.54
S-Prompt	45.45	13.47	45.01	14.76	45.27	14.29	42.98	20.20	42.85	25.82	44.09	22.32
<b>CluMo</b>	<b>46.18</b>	<b>10.62</b>	<b>45.36</b>	<b>11.69</b>	<b>46.34</b>	<b>10.22</b>	<b>45.95</b>	<b>9.15</b>	<b>45.66</b>	<b>19.89</b>	<b>46.89</b>	<b>17.41</b>

Table 1: Comparative experimental results: the accuracy and forgetting rate for different task order. For each task sequence, **A**  $\uparrow$  indicates the accuracy of the method, while **F**  $\downarrow$  is the forgetting rate of each.

positive effect of CL. Following the original setting of each method, we leave the whole backbone model unfrozen for non-prompt-based methods and freeze the whole backbone model for prompt-based methods except for the classifier. To make the fair comparison, we fit all the continual learning methods into our backbone, ALBEF, instead of using the original model proposed in each method.

**CL Tasks** We evaluate our method on tasks built using the **CLOVE** Lei et al. (2022) dataset which is a VQA-based continual learning dataset. The benchmark contains two sepecate benchmarks for different scenario, including scene-incremental setting benchmark, **CLOVE-scene**, and function-incremental setting benchmark, **CLOVE-function**. Each of the task sets contains six tasks which are domain-specific and diverse from each other. For more details about **CLOVE** and the tasks we use, please refer to the Appendix.

**Metrics for comparison** We use the average accuracy and the average forgetting rate on all tasks to evaluate the performance of our method and its ability to tackle catastrophic forgetting. Different from dataset such as **VQAv2** Goyal et al. (2017), where each question is paired with different ground truth answers, questions in **GLOVE** dataset only contains exactly one correct answer for each question. Thus, the accuracy is simply calculated by  $y == \hat{y}$  for every training and testing data instance. On the other hand, the forgetting rate is calculated as:

$$F = \frac{A_i - A_{ij}}{A_i} \quad (9)$$

where  $A_i$  is the accuracy of task  $i$ , and  $A_{ij}$  is the accuracy of task  $i$  after the model is trained on task  $j$ . For details about the optimization and implementation processes, please refer to the Appendix.

## 5.2 Comparative Results

We conduct the comparison experiments on both the **CLOVE-scene** and **CLOVE-function** task sets with a randomly selected task order. In table 1, the task order *abcdef* represents the CL tasks:



Table 2: Ablative Experiments

Methods	Accuracy	Forgetting
Full Method	48.73	10.76
Ablative KD	47.36	11.25
Ablative Clustering	46.08	12.86
Ablative Textual Key	46.16	12.49
Ablative Visual Key	46.53	12.22

*ShopAndDining*, *WorkPlace*, *HomeOrHotel*, *Transportation*, *SportAndLeisure Outdoors* in sequence. The *oarlks* in **CLOVE-function** represents tasks: *ObjectRecognition*, *AttributeRecognition*, *RelationReasoning*, *LogicReasoning*, *KnowledgeReasoning* and *SceneTextRecognition*.

We observe in Table 1 that our method outperforms all the baselines across all task order sets in terms of both average accuracy and average forgetting rates. We also observe that the performance of different method within the same group tend to be similar. The regularization-based methods, **EWC** and **LwF**, obtain the sub-optimal accuracy and forgetting rate besides. The reason is that the domain for each task in the dataset is significantly different from the rest of tasks and hence regularization methods fail to capture the common space of the parameter distribution. This challenge makes it difficult to maintain the accuracy of the current task and previous tasks at the same time using regularization. The replay methods, **ER** and **GEM**, achieve better performance than regularization-based methods. This can be explained by the fact that replaying the data from previous task is an efficient way to remind the model and adjust its parameter distribution not too diverse from previous ones. However, because we need to rely on a memory buffer to store samples for replay, these methods are memory-consuming and thus not space-efficient. Moreover, replay-based methods are still limited by the upper-bound of joint training, as they generally can only reduce catastrophic forgetting without boosting the accuracy of individual tasks. On the other hand, the prompt-based methods, namely **L2P**, **DualPrompt**, and **SPrompt**, achieve superior performances compared to more traditional CL methods. Rather than tuning the whole model with regularization, prompt-based methods store the prior knowledge in trainable prompts, which are smaller and more efficient than memory buffer, and keep the main body of backbone model frozen. With the combination of generalization capacity of pre-trained model and specific previous knowledge stored in prompt, prompt-based method can outperform the replay and regularization methods. Among all the methods, our method achieve the best performance.

Compared with the baseline prompt-based methods which only consider visual modality for prompt selecting and updating, **CluMo** takes care of both the visual and textual modalities, as well as the fusion of the two for selecting the prompt which deploys the given information more comprehensively to process the prompt. Our design thus fits better in multimodal learning scenario than other existing continual learning methods.

### 5.3 Ablation Experiments

To offer a better insight about our method, we perform an ablation study for each component of **CluMo** to study the positive contribution of each component. We study the effect of the following:

- **Visual Prompt Key**, key to separate the inner-task image features by their semantic property.
- **Textual Prompt Key**, key to separate the inner-task text features by their semantic property.
- **Minibatch  $k$ -means Clustering** which train the prompt keys as centers of clustering algorithm to better fit the semantic meaning.
- **Knowledge Distillation**, to prevent the drastic parameter shift of unfrozen classifier.

We conduct ablation experiment on **CLOVE-scene** dataset with the task order *abcdef*. We set the size for both the visual prompt key and the textual prompt key to be three. For ablative text experiments, we change

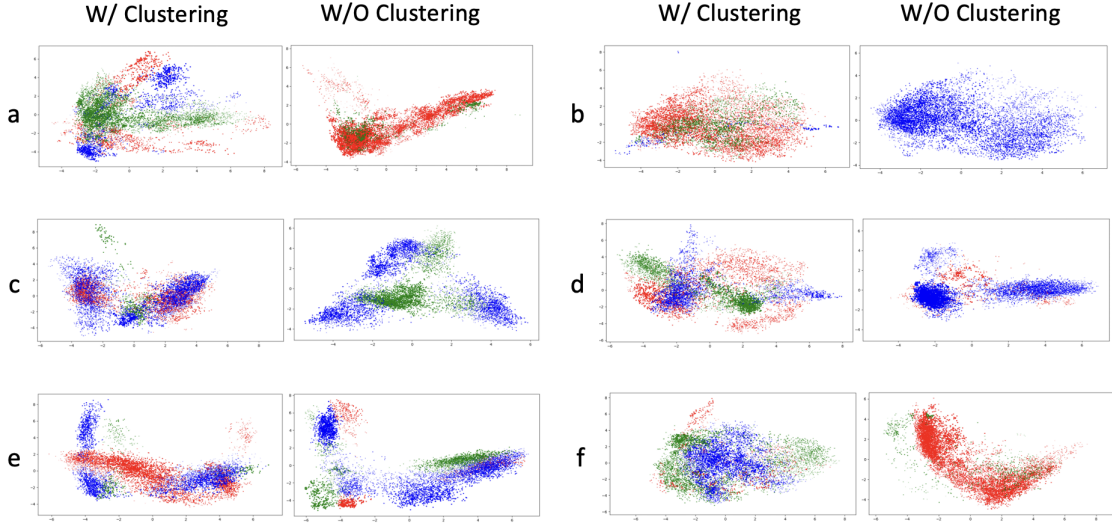


Figure 4: Cluster distribution on all training image data of **CLOVE-Scene**’s six sub-tasks before and after applying mini-batch  $k$ -means clustering algorithm with visual key size of 3 using PCA. **Color of more diversity indicates more even distribution of key selection.**

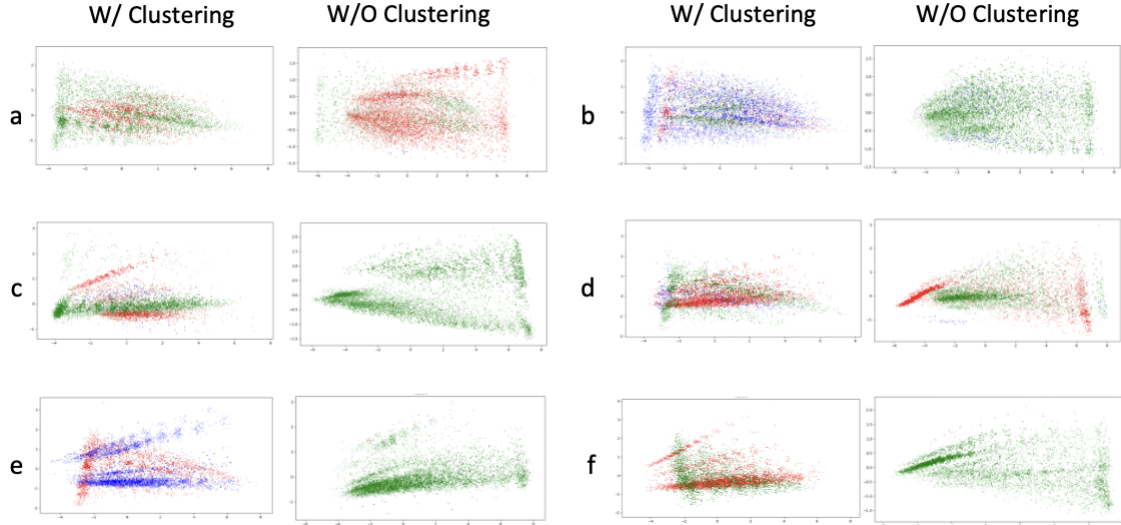


Figure 5: Cluster distribution on all training text data of **CLOVE-Scene**’s six sub-tasks before and after applying mini-batch  $k$ -means clustering algorithm with textual key size of 3 using PCA.

the size of textual prompt key to 9 to achieve the same prompt size. We also removed the visual prompt key which is the same for ablative image experiments. Results for this experiment is presented in Table 2. We observe that despite having the same number of prompts, the performance values of Ablative Textual Key and Ablative Visual Key are lower than our full pipeline. This result verifies our hypothesis that both modalities should be used to guide the prompt selection and the missing of any will cause information lost and lead to sub-optimal performance. In other words, current approaches for unimodal settings do not use all the information we have in multimodal scenarios. We also observe that without the clustering algorithm, the performance of ablative clustering is the lowest among all the settings which indicate the significance of doing cluster training for learning the prompt keys.

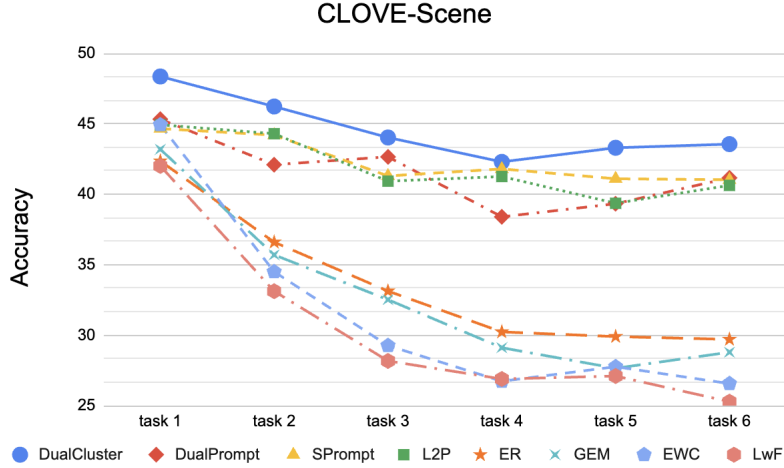


Figure 6: Accuracy on the first task after running task sequence.

#### 5.4 Analytic Experiments

**Effect of clustering** To show the effect of clustering algorithm, we empirically show the correlation between the clustering error and the downstream accuracy. As we apply Euclidean distance as metric to learn the clusters, we record the average distance between each point to its assigned cluster center for every task, and take the average for all the tasks:

$$\mathcal{E} = Avg\left(\sum_{i=1}^N Avg\left(\sum_{j=1}^M ||x_j - c_k||_2\right)\right) \quad (10)$$

where  $i$  represent the number of tasks,  $j$  represent the training data from task  $i$  and  $k$  is the  $k^{th}$  cluster center. We consider both the visual prompt key training and the textual prompt key training in this experiment. Table 3 presents the results. Same as our expectation, we observe a negative correlation between the clustering error and the performance accuracy, in another word, lower  $\mathcal{E}$  for image and text prompt keys leads to a higher accuracy. Without the clustering component, we observe  $\mathcal{E}$  to be as high as 42.38 and 42.8 for image prompt key and text prompt key, respectively. After applying clustering algorithm,  $\mathcal{E}$  drops below 20, which can be considered significant, for both modalities and the accuracy improves 2.97%.

**Cluster Visualization** To show the effect of clustering on prompt key more intuitively, we visualize the visual prompt key selection distribution and textual prompt key selection distribution on the visual and textual portion of the training data for **CLOVE-Scene** in Figure 4 and Figure 5. Since we use three visual prompt keys for each task, the vector features of visual data are split into three groups, which are the green, blue and red points in Figure 4. We observe that without using clustering, visual data are more likely to overlap on the same cluster center which means they would lead to select the same visual prompt key. After performing clustering, we observe that the distribution becomes more evenly, and every cluster of data is diverse and separated from the others which means that the visual data can be separated explicitly. The diversity of prompt key selection indicates each input image can find the "correct" prompt which is semantically closer to it, which contains more specific knowledge about the sub-domain the given image belongs to. The visualization of textual prompt keys indicates similar observation.

**Tracking the Accuracy for the First Task** To take a closer look into the effect on preventing catastrophic forgetting and increasing the accuracy in CL, we track the accuracy of the first task while learning the task sequence. The result is shown in Figure 6. We see that the accuracy drops until task 4, and then slightly increases until task 6. Although it is not our main focus, this behavior shows a trend of forward transfer

Table 3: Acc. with different clustering error

$\mathcal{E}$ . Image	$\mathcal{E}$ . Text	Accuracy
15.40	10.72	48.73
15.74	12.22	48.03
17.21	12.53	47.94
42.38	42.8	47.32

Table 4: Acc. with different prompt key size

$S_{img} \times S_{txt}$	Accuracy
$2 \times 2$	48.51
$3 \times 3$	48.73
$4 \times 4$	48.32
$5 \times 5$	48.32
$10 \times 10$	48.51

Table 5: Acc. with length/number trade-off

$S_v \times S_t \times L_p$	Accuracy
$2 \times 2 \times 22$	47.84
$3 \times 3 \times 10$	48.73
$4 \times 4 \times 6$	46.59

between the tasks. Among all the baseline methods, we notice that prompt-based methods, **SPrompt**, **DualPrompt** and **L2P**, significantly outperform other methods which verifies the SOTA status of prompt learning in CL and its success in preventing catastrophic forgetting. Our method **CluMo**, on the other hand, still outperform all prompt-based baseline methods. We observe that using the cluster-based prompts, the accuracy on the first task is superior compared to the other methods at the very beginning. Similar to other prompt-based method, our method’s accuracy slightly drops until task 4 and improves subsequently. As the accuracy of our proposed method is higher than others at all time steps, our method has the leading performance in terms of both accuracy and backward transfer.

**Effect of Prompt Key Size** We also conduct an experiment to study the effect of prompt pool size to show the stability of our method with respect to this hyper-parameter. In Table 4, we choose different visual prompt key and textual prompt key sizes,  $2 \times 2$ ,  $3 \times 3$ ,  $4 \times 4$ ,  $5 \times 5$ ,  $10 \times 10$ , corresponding to 4, 9, 16, 25, and 100 prompt pool sizes. Although having variant sizes, we only observe minor changes in accuracy in Table 4 when prompt pool size changes, i.e., between 48.32 and 48.73. This observation means that our method is not sensitive to the change of the prompt pool size and hence we don’t need to focus on tuning it.

**Prompt Length and Prompt Size** In the experiment setting, we utilize 3 visual keys and 3 textual keys, and each prompt has length of 10. In sum, we have  $3 \times 3 \times 10 = 90$  total prompt length for each task. We conduct experiments to show the trade-off of prompt length and prompt numbers given the fixed total prompt length. To make a fair comparison, we keep the visual key and textual key the same size, and choose the closest integer prompt length to let the multiplication of the three be around 90. Within table 6, among the three different combination, we find that the current setting, which is the most balanced one, obtains the best result. Regarding the rest of the two, the higher prompt length has the better performance.

## 6 Conclusion

We introduced a novel prompt-based continual learning method for learning multimodal tasks. While most of existing methods apply simple prompts on a single modality, our method proposes modal-specific visual prompt keys and textual prompt keys and train them to capture the semantic properties of the training dataset using K-means clustering algorithm. We use the combination of both the visual prompt key and the textual prompt key to select prompts, which enable the prompt to better boost the performance. Our experiments show that our method achieves the state-of-the-art performance in continual VQA tasks in different domains compared to other regularization-based, rehearsal-based and prompt-based CL methods.

## References

- Aishwarya Agrawal, Jiasen Lu, Stanislaw Antol, Margaret Mitchell, C. Lawrence Zitnick, Dhruv Batra, and Devi Parikh. Vqa: Visual question answering, 2016.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- Yuliang Cai, Jesse Thomason, and Mohammad Rostami. Task-attentive transformer architecture for continual learning of vision-and-language tasks using knowledge distillation, 2023.
- Ryan Cohn and Elizabeth Holm. Unsupervised machine learning via transfer learning and k-means clustering to classify materials image data. *Integrating Materials and Manufacturing Innovation*, 10(2):231–244, 2021.
- Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion, 2022.
- Robert M French. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4): 128–135, 1999.
- 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, 2017.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- Yushi Hu, Hang Hua, Zhengyuan Yang, Weijia Shi, Noah A Smith, and Jiebo Luo. Promptcap: Prompt-guided image captioning for vqa with gpt-3. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2963–2975, 2023.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning, 2016.
- Wonjae Kim, Bokyung Son, and Ildoo Kim. Vilt: Vision-and-language transformer without convolution or region supervision, 2021.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, March 2017. ISSN 1091-6490. doi: 10.1073/pnas.1611835114. URL <http://dx.doi.org/10.1073/pnas.1611835114>.
- Stan Weixian Lei, Difei Gao, Jay Zhangjie Wu, Yuxuan Wang, Wei Liu, Mengmi Zhang, and Mike Zheng Shou. Symbolic replay: Scene graph as prompt for continual learning on vqa task, 2022.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning, 2021.
- Junnan Li, Ramprasaath R. Selvaraju, Akhilesh Deepak Gotmare, Shafiq Joty, Caiming Xiong, and Steven Hoi. Align before fuse: Vision and language representation learning with momentum distillation, 2021.
- Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language, 2019.
- Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation, 2021.
- Zhizhong Li and Derek Hoiem. Learning without forgetting, 2017.

- David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning, 2022.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks, 2019.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. Ok-vqa: A visual question answering benchmark requiring external knowledge, 2019.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. icarl: Incremental classifier and representation learning, 2017.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P. Lillicrap, and Greg Wayne. Experience replay for continual learning, 2019.
- Lin Su, Nan Duan, Edward Cui, Lei Ji, Chenfei Wu, Huaishao Luo, Yongfei Liu, Ming Zhong, Taroon Bharti, and Arun Sacheti. Gem: A general evaluation benchmark for multimodal tasks, 2021.
- Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained transformers: An occam’s razor for domain incremental learning, 2023.
- Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Dualprompt: Complementary prompting for rehearsal-free continual learning, 2022a.
- 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, 2022b.

Table 6: size of each task in CLOVE

task	training size	testing size	image source
ShopAndDining	20K	3K	MS-COCO
WorkPlace	20K	3K	MS-COCO
HomeOrHotel	20K	3K	MS-COCO
Transportation	20K	3K	MS-COCO
SportAndLeisure	20K	3K	MS-COCO
Outdoors	20K	3K	MS-COCO
ObjectRecognition	20K	3K	MS-COCO
AttributeRecognition	20K	3K	MS-COCO
RelationReasoning	20K	3K	MS-COCO
LogicReasoning	20K	3K	MS-COCO
KnowledgeReasoning	20K	3K	MS-COCO
SceneTextRecognition	16.8K	2.4K	VG

## A Appendix

### A.1 Hardware Setup and Hyper-parameter

All experiments were conducted using a single Nvidia A40 GPU. We employed the AdamW optimizer across all experiments, utilizing a cosine learning rate scheduler, and set the initial learning rate to  $lr = 3 \times 10^{-4}$ . The models were trained for 5 epochs with a batch size of 16.

For **EWC**, we set the fisher sample percentage to be 0.1 and ewc loss weight equals to 0.1 as well.

For **Experiment Replay**, we store 1% of data from each tasks into the memory buffer. During the training of the current task, we randomly select a batch of data from the memory buffer to train the model for every 100 batches of current data training.

For **GEM**, we also store 1% of data to memory buffer, which is randomly picked from each tasks.

For **CluMo** framework, we configured the visual prompt key size ( $S_v$ ) and the text prompt key size ( $S_t$ ) both to 3, with the prompt length ( $L_p$ ) set to 10.

For **L2P**, we set the prompt pool size equals to 20 and prompt length to be 5.

For **DuamPrompt**, the G-Prompt was inserted into layers 0 and 1 of the visual encoder, while the E-Prompt was integrated into layers 2, 3, and 4.

For **S-Prompt**, we set the prompt length to be 10 and prompt pool size to be 30.

Furthermore, for prompt-based techniques such as L2P, DualPrompt, and SPrompt, we opted to freeze the entire backbone model, allowing only the final classifier layer to remain trainable. Conversely, for all other baseline methods, no parameters were frozen, ensuring the entire network was fine-tuned during training.

### A.2 CLOVE dataset detail description

In the **CLOVE-Scene** and **CLOVE-Function** datasets, all tasks have a uniform distribution of training and testing data, with the exception of the *SceneTextRecognition* task, which comprises 16.8K training samples and 2.4K testing samples. The remaining tasks within these datasets contain 20K training samples and 3K testing samples each. It is reflected in the Table 5.

To provide a deeper understanding of the **CLOVE** dataset, we offer additional details here. We have included visualizations in Figure 7 and Figure 8 to showcase two sample images from each task within the datasets, emphasizing the distinctiveness of each domain. From these samples, it is evident that the images in the **CLOVE-Scene** dataset vary significantly across tasks, even though the questions associated with them are similar in structure, differing primarily based on the content depicted in the images.

On the other hand, the **CLOVE-Function** dataset presents a different scenario. While the images across various tasks may appear to belong to similar or overlapping domains, making it challenging to distinguish one task from another based solely on visual content, the diversity becomes apparent when considering the questions. Each task within the **CLOVE-Function** dataset involves distinct types of reasoning, as reflected in the varied nature of the questions posed, which are tailored to serve different reasoning purposes.



Q: Is there a sandwich on the tray?  
A: No.

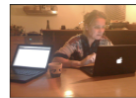


Q: Is the bottle on the couch?  
A: Yes.

#### Shopping and Dining



Q: Which room is it?  
A: Office.



Q: Is a laptop to the left of her?  
A: Yes.

#### Workplace



Q: Does the blanket look red?  
A: Yes.



Q: What toy on the top of sink?  
A: Rubber Duck.

#### Hotel



Q: What is the plane behind the man?  
A: Runway.



Q: What are the letters on?  
A: Stairs.

#### Transportation



Q: Is the man on the right of frisbee wearing a hat?  
A: Yes.



Q: What is the man playing?  
A: Frisbee.

#### SportAndLeisure



Q: Does the zebra nose have the white color?  
A: No.



Q: What is flying in the sky?  
A: Kite.

#### Outdoors

Figure 7: CLOVE-scene dataset samples





Q: What color is the helmet in the middle of the image?  
A: Blue.



Q: Is it indoor or outdoor  
A: Yes.

#### Attribute Recognition



Q: Which red object behind the black piano can keep light out of the room?  
A: Curtain.



Q: What food is next to the object that I can use to for making toast?  
A: Bread.

#### Knowledge Reasoning



Q: What do the marker and the post have in common?  
A: Color.



Q: Is the color of pillow different than that of counter?  
A: Yes.

#### Logic Reasoning



Q: What place is it?  
A: Harbor.



Q: What is this, a couch or a table?  
A: Table.

#### Object Recognition



Q: What color is the jersey the boy is wearing?  
A: Black.



Q: Who is wearing goggles?  
A: Woman.

#### Relation Reasoning



Q: What is the brand of this camera?  
A: Dakota.



Q: What is the brand of the phone?  
A: Nokia.

#### Scene Text Recognition

Figure 8: CLOVE-function dataset samples