# PLUG-AND-PLAY COMPOSITIONALITY FOR BOOSTING CONTINUAL LEARNING WITH FOUNDATION MODELS

**Anonymous authors** 

Paper under double-blind review

#### **ABSTRACT**

Vision learners often struggle with catastrophic forgetting due to their reliance on class recognition by comparison, rather than understanding classes as compositions of representative concepts. This limitation is prevalent even in state-of-the-art continual learners with foundation models and worsens when current tasks contain few classes. Inspired by the recent success of concept-level understanding in mitigating forgetting, we design a universal framework CompSLOT to guide concept learning across diverse continual learners. Leveraging the progress of object-centric learning in parsing semantically meaningful slots from images, we tackle the challenge of learning slot extraction from ImageNet-pretrained vision transformers by analyzing meaningful concept properties. We further introduce a primitive selection and aggregation mechanism to harness concept-level image understanding. Additionally, we propose a method-agnostic self-supervision approach to distill sample-wise concept-based similarity information into the classifier, reducing reliance on incorrect or partial concepts for classification. Experiments show CompSLOT significantly enhances various continual learners and provides a universal concept-level module for the community.

#### 1 Introduction

Artificial intelligence systems mimic the learning behavior of human intelligence by collecting information and managing knowledge pools from continually assigned tasks in the open world. This need to handle non-independent and identically distributed training data has driven research in continual learning (CL) (Zhou et al., 2024c;a; Biesialska et al., 2020), which aims to balance the objectives of overcoming *catastrophic forgetting* (McCloskey & Cohen, 1989) of learned tasks and achieving *efficient adaptation* to future tasks, also known as the *stability-plasticity dilemma* (Grossberg, 2012). Leveraging a powerful pre-trained backbone to ensure a basic understanding of the world, CL methods of foundation models (FMs), including prompt-based methods (Gao et al., 2023; Smith et al., 2023; Wang et al., 2022c;b; 2024; Gao et al., 2024), representation-based methods (Zhou et al., 2025; 2024b; McDonnell et al., 2023; Zhang et al., 2023), and model-mixture-based methods (Gao et al., 2023; Wang et al., 2024; Marouf et al., 2024), have emerged as a popular direction in this field. However, FMs need to be updated when encountering out-of-distribution data in the upcoming tasks (Yang et al., 2025).

The human brain exhibits *compositionality* (Hupkes et al., 2020; Liao et al., 2024) when comprehending the world, decomposing seen concrete objects into abstract concepts. For example, a *Chihuahua* consists of general dog concepts such as *body shapes* and chihuahua-specific concepts like *small size* and *head shapes*. This interpretability is intuitive to humans, enabling them to generalize novel dog species by decomposing them into combinations of existing concepts while learning disentangled new concepts to refine the knowledge base, thus, facilitating efficient reuse (Liao et al., 2024). A common strategy for existing CL methods for FMs to alleviate forgetting is to inherit parameters learned from old tasks when initializing new tasks' models, as done in Wang et al. (2024); Gao et al. (2024). These state-of-the-art (SOTA) approaches generally do not fully exploit cross-task potential correlations (i.e., common concepts shared across tasks). In contrast, learning low-dimensional concept combinations to understand classes does not require establishing class representations from the high-dimensional feature level, as in traditional methods, thereby mitigating catastrophic forgetting and enabling rapid adaptation to novel classes (Liao et al., 2024; Yu et al., 2025; Yang et al., 2024; Kundargi et al., 2025; Lai et al., 2024). Thus, a set of CL methods leverages interpretable tools, e.g., ChatGPT (Brown

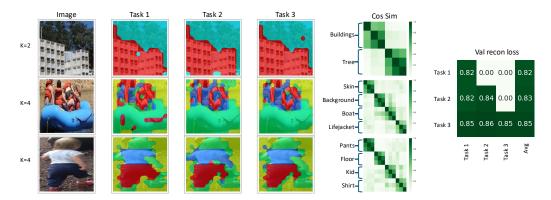


Figure 1: Examples of learned slots by continual COBJ reconstruction tasks and validation reconstruction losses. **Left**: slots are extracted for three example images from the first task on learners after training on the 1-st, 2-nd, and 3-rd tasks. Slots are masked with different colors. **Middle**: corresponding slot cosine similarity matrices grouped by correlated regions. Each group contains slots from three tasks and is identified by the Hungarian matching algorithm. **Right**: validation reconstruction loss matrix. Each row indicates a learner trained after a specific task and evaluated on all seen tasks, respectively. **Takeaway**: learned Slot Attention module enjoys almost no forgetting across compositional-relevant tasks.

et al., 2020), and concept bottleneck models (Yu et al., 2025; Yang et al., 2024; Lai et al., 2024) to bring attention to concepts within images, which achieves great success on boosting CL performance. Another challenge is that canonical benchmarks, such as Split-CIFAR100 (Krizhevsky et al., 2009) and Split-ImageNet-R (Hendrycks et al., 2021), are not specifically designed to evaluate the compositionality of continual models. The CFST evaluation framework (Liao et al., 2024) (including CGQA and COBJ) is the only work, to our knowledge, that systematically studies the compositionality of a continual learner. CFST introduces two component-relevant phases in which the data share a common concept set with different combinations. In the first phase, the dataset is split into several continual tasks, aiming to train a continual learner. Subsequently, the second phase is used to evaluate the learner's compositional generalization performance on unseen concept combinations.

Motivated by the above analysis, we pose the research question: Can the compositionality in concept learning truly enhance the CL performances of SOTA continual learners with FMs? We propose a **Compositional Slot** plug-in (**CompSLOT**) for continual learning to answer the above question. The first step involves extracting concepts from raw images, namely, conducting concept learning. Several studies have shown significant progress in concept learning by utilizing explicit concept-level supervision obtained from segmentation masks (Kirillov et al., 2023; Ravi et al., 2025) or natural language annotations (Ramesh et al., 2021; Yu et al., 2025). Nevertheless, it is crucial to compare with SOTA CL methods of FMs, where only labels from the current CL tasks are available as supervision. Consequently, Slot Attention (Locatello et al., 2020), as an SOTA unsupervised object-centric learning approach (Greff et al., 2020), has effectively emerged as a viable self-supervised solution. A Slot Attention module learns to group and encode spatial features into a set of low-dimensional distinct slots, with each slot representing a disentangled region and binding to an object (i.e., concept) in the image. To avoid additional learning of the encoder, the input to Slot Attention can be specified as semantic patch features provided by a pre-trained vision transformer (ViT), which also serves as the learner's backbone. We present a preliminary experiment demonstrating that the learned Slot Attention module exhibits almost no forgetting across compositional tasks, as shown in Figure 1. Specifically, we train a Slot Attention module on COBJ 3-tasks as a continual reconstruction task (i.e., trained with reconstruction loss). We then extract slots for images from the first task using the modules after training on the first, second, and third tasks. We observe that each corresponding slot consistently represents a human-interpretable concept and remains stable after training on new tasks, maintaining a high cosine similarity.

With the above method to extract the hidden concepts in images, the next step is to introduce concept learning into CL methods with FMs. The challenge is that there is no unified forwarding framework to organize all CL methods with FMs so that we can easily perform concept learning and assist vanilla learning processes of feature extractors. Hence, we propose regularizing the outputs of learners with **sample-wise similarity based on concepts**. This makes our approach a method-agnostic plugin for

any CL method with FM. We first use a learnable aggregation mechanism based on attention to extract class-relevant concepts (i.e., **primitives** (Zou et al., 2024)) as the weighted sum of slots based on their similarity to a learnable task key. The distance of primitives between two images carries information about the similarity in concept level. For example, a *Chihuahua* is close to other dog species (e.g., *German Shepherd*) rather than cat species (e.g., *Siamese*) because they share considerably more concepts (e.g., *dog body*). Subsequently, we propose a method-agnostic primitive-logit alignment plugin to distill the learned sample-wise concept-level similarity into the outputs of models based on a contrastive loss. Our experiments demonstrate that the above procedures successfully select meaningful concepts in images as primitives and ultimately achieve a superior continual learning performance attributed to a better compositional generalization performance.

The contributions of this work are summarized as follows:

- CompSLOT successfully achieved desired concept extraction and task-dependent aggregation mechanisms using an ImageNet-pretrained ViT with an additional primitive loss. This enabled the utilization of hidden concept information to benefit CL tasks without introducing additional powerful pretrained segmentation models, such as Oquab et al. (2024); Kirillov et al. (2023).
- CompSLOT designed a method-agnostic primitive-logit alignment plugin that distills concept similarity into the outputs of models. This allowed learners to intentionally discover shared and distinct concepts among classes, guiding the decision-making process of classifiers.
- The experimental results showed that our mechanisms successfully leverage concept-wise compositionality to significantly enhance a large range of continual learners.

#### 2 RELATED WORKS

**Continual Learning of Foundation Models** Benefiting from the rich knowledge in large-scale pre-trained ViT, CL methods with FMs (Zhou et al., 2024a) greatly mitigate forgetting previously learned classification tasks and achieve fast adaptation to new ones. The community has mainly developed three families of approaches, according to the way of utilizing the pre-trained knowledge: 1) Prompt-based methods (Gao et al., 2023; Smith et al., 2023; Wang et al., 2022c;b; Gao et al., 2024; Liang & Li, 2024; Le et al., 2024) efficiently tune prompts for tasks rather than fine-tune the backbone; 2) Representation-based methods(Zhou et al., 2025; 2024b; McDonnell et al., 2023; Zhang et al., 2023) involve leveraging the advantages of representations from the pre-trained backbone with a class prototype-based classifier; 3) Model-mixture-based methods (Gao et al., 2023; Wang et al., 2024; Marouf et al., 2024) utilize hybrid techniques such as model fusion (Wang et al., 2024; Marouf et al., 2024) and model ensemble (Gao et al., 2023) to query a set of models, thus, making the prediction more robust; Moreover, rehearsing old samples is an effective way to alleviate forgetting old tasks. Several methods (Wang et al., 2022a; Yan et al., 2021; Zhou et al., 2023) contribute to efficient sample storage mechanisms and auxiliary supervision to address class imbalance, achieving a better stability-plasticity trade-off. However, the above methods ignore hidden conceptual relationships among classes, limiting their significance on handling compositionally relevant tasks.

Compositionality Compositionality has been extensively studied in natural language processing (Biesialska et al., 2020; Kaushik & Martin, 2020; Lake & Baroni, 2018; Keysers et al., 2020). To achieve a compositional learner, methods include the introduction of sparse coding (Murphy et al., 2012), regularization (Sun et al., 2016; Luo et al., 2015), and applying independent component analysis (Musil & Mareček, 2022; Yamagiwa et al., 2023). In Hupkes et al. (2020), the authors summarize five types of tests for language compositionality, which are further extended to vision in Liao et al. (2024). Meanwhile, researchers in vision utilize compositional information between objects and attributes to boost zero-shot inference through regularization (Nagarajan & Grauman, 2018), separate learning (Ruis et al., 2021), causal reasoning (Atzmon et al., 2020), self-attention (Khan et al., 2023), and uniting energy-based modules (Wu et al., 2022). Common strategies to learn hidden concepts among continual tasks are external interpretability tools (Yang et al., 2024), learnable mapping (Lai et al., 2024), CLIP (Kundargi et al., 2025), ChatGPT (Yu et al., 2025), and assigning different module paths for tasks (Rajasegaran et al., 2019; Ostapenko et al., 2021). Our work, instead, does not require prior concept-level supervision for training or an extra concept bottleneck model (Yu et al., 2025), making it more adaptable and easier to integrate with different methods.

Object-centric Learning We adopt object-centric learning to autonomously extract concept information directly from images. The introduction of Slot Attention (Locatello et al., 2020) marked the emergence of a new paradigm for disentangling objects (i.e., concepts) within a scene. Subsequent research has focused on improving its robustness in complex environments—primarily through encoder enhancements like covariance regularization (Stange et al., 2023) and bi-level optimization (Jia et al., 2023; Chang et al., 2022). Other efforts have explored advanced decoders to refine decomposition. For example, SLATE (Singh et al., 2022) uses an autoregressive transformer decoder, while Wu et al. (2023); Jiang et al. (2023) propose diffusion-based approaches. Kakogeorgiou et al. (2024) leverages distillation to refine object segmentation via decoder-guided encoder training, and Kori et al. (2023) introduces conditional Slot Attention with a foundational slot dictionary to address specialization limitations. Our method, instead, employs a lightweight MLP decoder to minimize computational cost while preserving effectiveness. Experiments show that this simple design can still significantly benefit continual learning.

# 3 PRELIMINARIES

Class-incremental vision continual classification tasks We consider T sequential vision classification tasks with a dataset  $\mathcal{D} = [\mathcal{D}^1, \dots, \mathcal{D}^T]$ , where each  $\mathcal{D}^t$  consists of image samples  $\boldsymbol{x} \in \mathcal{X}^t$  with corresponding labels  $y \in \mathcal{Y}^t$ . Here,  $\mathcal{Y}^t$  is a subset of the global label set  $\mathcal{Y}$ , and  $\forall \mathcal{Y}^t \cap \mathcal{Y}^k = \emptyset$  for  $t \neq k$ , with task identity unknown during inference, i.e., class-incremental learning (CIL) setting. A general model-based continual learner includes a Vision Transformer (ViT)-based backbone  $f(\cdot|\theta_f)$  and classification heads  $h_t(\cdot|\theta_{h_t})$ , where t is the task identity. Each head is trained separately for the corresponding task, but the outputs from all heads are concatenated for final inference:  $\boldsymbol{H_{te}} = f(\boldsymbol{x_{te}}|\theta_f)[0]$ , where [0] indicates the [CLS] token (i.e., the first dimension of the output of f), and  $\operatorname{pred}(\boldsymbol{x_{te}}) = \arg\min\left[h_1\left(\boldsymbol{H_{te}}|\theta_{h_1}\right);\dots;h_T\left(\boldsymbol{H_{te}}|\theta_{h_T}\right)\right]$ , where  $[\cdot;\cdot]$  denotes concatenation.

Slot attention (Locatello et al., 2020) As the state-of-the-art object-centric plug-in, slot attention aims to decompose a single image into a set of K disentangled slots  $\mathbf{S} \in \mathbb{R}^{K \times D_s}$ , each encoding one compositional component of the image.  $D_s$  is the dimension of slot representation. The output  $f(\mathbf{x}|\theta_f)$  from a pre-trained ViT backbone consists of two parts: the uninstructed image feature  $\mathbf{H} = f(\mathbf{x}|\theta_f)[0] \in \mathbb{R}^D$  with the token [CLS] and the semantic patch features  $\mathbf{E} = f(\mathbf{x}|\theta_f)[1:] \in \mathbb{R}^{N \times D}$ , where N is the patch number. These N patches are further encoded into the slot space and refined into K slots through an iterative attention procedure. The K slots are first initialized with a learnable Gaussian distribution. In each refinement iteration, slots collect soft assignment information from each patch with an attention mask  $\mathbf{A} \in \mathbb{R}_+^{K \times N}$ . The weighted mean K is then computed along the patch dimension, and a Gated Recurrent Unit (GRU) (Cho et al., 2014) aggregates the patch information into the assigned slots, as follows:  $\mathbf{A} = \sigma\left(\frac{q(\mathbf{S})k(\mathbf{E})^\top}{\sqrt{D_s}}\right)$ ,  $A_{i,n} \leftarrow \frac{A_{i,n}}{\sum_{j=1}^N A_{i,j}}$ ,  $\mathbf{S} \leftarrow \mathrm{GRU}(\mathbf{S}, \mathbf{A}v(\mathbf{E}))$ , where  $q(\cdot), k(\cdot), v(\cdot)$  are learnable query, key, value projections, respectively, and  $\sigma(\cdot)$  is the softmax function.

# 4 METHODS

We present our CompSLOT framework in Figure 2. For each continual task  $\mathcal{D}^t$ , we first perform **concept learning** (detailed in section 4.1) through a **slot decomposition** and a **primitive selection** mechanism, and then distill the pair-wise similarity statistic of the extracted primitives to model outputs (detailed in section 4.2) in a method-agnostic manner.

Unless otherwise stated, the proposed slot attention and primitive selection modules are **globally shared** across tasks, and no parameters except the ViT backbone are frozen. They are initialized at the beginning of the first CL task. In future CL tasks, their architectures (e.g., the number of slots) remain fixed, while parameters will be fine-tuned throughout all CL tasks. This design prevents parameter explosion and supports long-sequence tasks, as demonstrated in Figure 3b.

#### 4.1 Concept Learning

Firstly, we define **concepts** as the ground truth slot decomposition of an image. Since slot attention exhibits permutation equivalence w.r.t. the order of the slots (and masks) (Locatello et al., 2020),

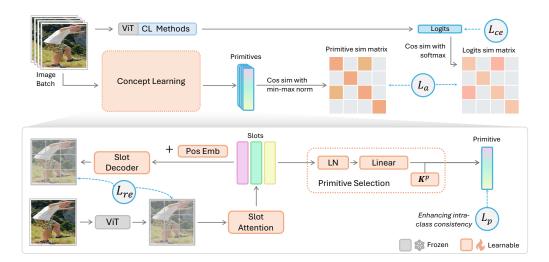


Figure 2: Proposed CompSLOT framework. Given an image batch, we extract primitives for each image with a concept learning procedure. Then, we can distill the sample-wise conceptual similarity from the primitive representations of the image batch into logits to enhance compositionality of CL methods.

we regard  $\{S, A\}$  as the corresponding set representations of  $\{S, A\}$ , where  $S = \{s_i\}_{i=1}^K$  and  $A = \{a_i\}_{i=1}^K$ , with  $s_i \in \mathbb{R}^{D_s}$  and  $a_i \in \mathbb{R}^N_+$  being the *i*-th row of S and A, respectively.

**Definition 1** (Concept & Disentanglement). Let x be an image, then  $\{S, A\}$  is a *disentangled* decomposition of x (a.k.a., *concepts* and corresponding *attention regions*), if 1)  $a_i, a_j \in A, a_i, a_j \in \mathbb{R}^N$ ,  $a_i \perp a_j$ , and 2) S satisfies  $\arg\min_{S} \sum_{s_i, s_j \in S} |\sin(s_i, s_j)|$ , where  $s_i, s_j \in \mathbb{R}^{D_s}$ , and  $|\cdot|$  is the absolute value,  $\perp$  is orthogonal symbol, and  $\sin(\cdot, \cdot)$  is a similarity score function, e.g., cosine similarity.

Remark 1. The examples of concepts are Chihuahua's head w.r.t. Chihuahua objects in section 1 and buildings w.r.t. images in Figure 1. In Figure 1,  $\mathcal A$  corresponds to patch regions and  $\mathcal S$  corresponds to slot representations used to calculate cosine similarity.

To train the slot attention and primitive selection modules, we resort to continually reconstructing  $\mathcal{D}$  and an additional contrastive primitive loss.

**Continual image reconstruction** For clarity, we denote the forward path of slot attention as  $\{S,A\} = s(E|\theta_s), S \in \mathbb{R}^{K \times D_s}, A \in \mathbb{R}^{K \times N}$ . We augment the position embedding into slots S when reconstructing the image, as the ViT does.  $S'_n = S \oplus pos_n$ , where  $pos_n \in \mathbb{R}^{D_s}$  is the learnable position embedding at patch n and  $\oplus$  is the element-wise addition with broadcasting. Next,  $S' \in \mathbb{R}^{N \times K \times D_s}$ , a collection of N position-augmented slots, are mapped back individually to the D-dim patch space with an MLP slot decoder  $d(\cdot|\theta_d)$ . Subsequently, we apply weighted-sum with the attention mask A and finally get the reconstructed patch features E. The reconstruction loss  $L_{re}$  is the MSE loss between the ground truth patch features E and the reconstructed E, as follows:

$$\tilde{\boldsymbol{E}} = \boldsymbol{A}^{\top} d(\boldsymbol{S'}|\theta_d) \in \mathbb{R}^{N \times D}, \quad L_{re} = ||\boldsymbol{E} - \tilde{\boldsymbol{E}}||_2.$$
 (1)

When describing the object *Chihuahua*, some concepts (like *Chihuahua's head*) are class-relevant, while others (like *sky*) are class-irrelevant. We name such class-relevant concepts as **primitives**, containing information to identify the desired classes.

**Definition 2** (Primitives). Let  $\mathcal{X}^y$ ,  $\mathcal{S}^y$  be an image set labeled y and the corresponding set of concept sets, respectively, and  $\mathcal{S} \in \mathcal{S}^y$ , then a concept subset  $\mathcal{P} \subset \mathcal{S}$  is *primitives* of  $\mathcal{S}$ , if  $\forall \mathcal{S}' \in \mathcal{S}^y$ ,  $\mathcal{P} \subset \mathcal{S}'$ .

Our goal is to **identify a unified primitive representation**  $s^p$ , which is regarded as the linear combination of concepts S. The basic idea is that primitives have a higher probability appearing in xs from the same class  $\mathcal{X}^y$  and are likely to carry important information describing this class. Thus, we have the following two questions: 1) How to represent the selected primitives  $s^p$  from S? and 2) How to minimize the distances among  $s^p$ s extracted from the images in the same class?

**Primitive selection** To answer the first question, we propose a learnable attention-based primitive selection mechanism to aggregate K slots. We use a linear module with layer norm and a tanh activation layer to map slots into a unified similarity space. The similarity to a learnable primitive key  $K^p \in \mathbb{R}^{D_s}$  measures the slot significance. Then this similarity  $w_p$  weights the mapped slots and aggregates them into a single representation  $s^p$  (i.e., primitive representation), which is summarized as follows:

$$\bar{S} = \tanh(\text{Linear}(\text{LN}(S))), \quad w_p = \sigma(\tau_t \bar{S} K^p), \quad s^p = w_p^{\top} \bar{S},$$
 (2)

where  $\tau_t$  is a temperature coefficient controlling the sparsity of slot selection  $w_p$ , which is set to  $100/\sqrt{D_s}$  in practice. A larger  $\tau_t$  indicates a smaller number of slots to be selected to represent this image x.

**Contrastive primitive loss** To answer the second question, we rewrite Definition 2 as follows:

**Theorem 1** (Intra-class consistency). Consider  $S_1, S_2 \in S^y$  and two corresponding largest primitive sets  $\mathcal{P}_1 \subset S_1, \mathcal{P}_2 \subset S_2$  are identical, i.e.,  $\mathcal{P}_1 = \mathcal{P}_2$  and  $||\mathcal{P}_1|| = M$ , where  $||\cdot||$  is the cardinality of set. In other word, consider the pair-wise ordered sets  $\{\mathcal{P}_1^\circ, \mathcal{P}_2^\circ\} = \mathrm{match}(\mathcal{P}_1, \mathcal{P}_2)$ , where  $\mathrm{match}(\cdot, \cdot)$  is a matching algorithm (without loss of generality, Hungarian algorithm (Kuhn, 1955)), then the corresponding matched concepts should be the same:  $\mathcal{P}_1^\circ = \{s_i^1\}_{i=1}^M, \mathcal{P}_2^\circ = \{s_i^2\}_{i=1}^M$  and  $\forall i \in \{1, \ldots, M\}, \sin(s_i^1, s_i^2) = 1$ .

This form of pair-wise primitive similarities from images within the same class motivates the use of label supervision and contrastive learning (Khosla et al., 2020; Chen et al., 2020b). We first collect the normalized similarity  $d_{i,j}^y$  between one-hot label and the softmax similarity  $d_{i,j}^s$  between  $s^p$ . Then, we use a mini-batch clustering loss that a small KL divergence between  $d_{i,j}^y$  and  $d_{i,j}^s$  means a small distance between  $s_i^p$ ,  $s_j^p$  in the same class and a large distance between those in different classes. The primitive loss  $L_p$  is as follows:

$$d_{i,j}^{y} = \frac{\sin(\mathbb{I}_{i}, \mathbb{I}_{j})}{\sum_{\boldsymbol{x}_{k} \in B} \sin(\mathbb{I}_{i}, \mathbb{I}_{k})}, \quad d_{i,j}^{s} = \frac{\exp(\tau_{p} \sin(\boldsymbol{s}_{i}^{\boldsymbol{p}}, \boldsymbol{s}_{j}^{\boldsymbol{p}}))}{\sum_{\boldsymbol{x}_{k} \in B} \exp(\tau_{p} \sin(\boldsymbol{s}_{i}^{\boldsymbol{p}}, \boldsymbol{s}_{k}^{\boldsymbol{p}}))}, \quad L_{p} = \sum_{x_{i}, x_{j} \in B} d_{i,j}^{y} \log \frac{d_{i,j}^{y}}{d_{i,j}^{s}},$$
(3)

where  $\mathbb{I}_i$  is the one-hot label for sample  $x_i$ , and  $\tau_p$  is a temperature coefficient controlling the strength of primitive loss. The learned slot visualizations in section K (including CGQA, COBJ, ImageNet-R, CIFAR-100) demonstrate that meaningful concepts (represented by primitives, third column "Sum") remain stable across tasks for the same images. We attribute this robustness to "concept rehearsal": although class labels change, many visual concepts are shared and recur across tasks, helping stabilize the primitive selection weights. Section K also visualizes the pair-wise primitive similarities, showing that concept relationships are preserved across images of the same class and shared concepts remain consistent even when images are from different tasks.

By jointly minimizing  $L_{re}, L_p$ , the learned slot attention module equips the abilities of extracting concepts, identifying primitives, and achieving intra-class primitive consistency. Specifically, we group these losses as  $L_{slot} = L_{re} + \alpha L_p$ , where  $\alpha$  is a coefficient to balance the impact of  $L_p$ .

#### 4.2 METHOD-AGNOSTIC PRIMITIVE-LOGIT KNOWLEDGE DISTILLATION

The learned primitive  $s^p$  equips a superior property of aggregating important class-relevant concepts. Such understanding can be a self-supervision to regularize the output of the continual learner, i.e., the distribution of logits. Thus, the model gives predictions based on the exact extracted concepts. For example, a *chihuahua* image should have relatively higher logits on other *dog* classes than logits on *cat* classes because they share similar concepts such as *dog body shapes*. Specifically, we design a primitive-logit alignment loss to contrastively distill the learned primitive statistics to logit statistics, i.e., minimizing the KL divergence between softmax logit similarity  $d^l$  and previously learned primitive similarity  $d^s$ , as follows:

$$d_{i,j}^{s} = \frac{\operatorname{sim}_{+}(\boldsymbol{s}_{i}^{\boldsymbol{p}}, \boldsymbol{s}_{j}^{\boldsymbol{p}})}{\sum_{\boldsymbol{x}_{k} \in B} \operatorname{sim}_{+}(\boldsymbol{s}_{i}^{\boldsymbol{p}}, \boldsymbol{s}_{k}^{\boldsymbol{p}})}, \quad d_{i,j}^{l} = \frac{\exp(\tau_{a} \operatorname{sim}(\boldsymbol{l}_{i}, \boldsymbol{l}_{j}))}{\sum_{\boldsymbol{x}_{k} \in B} \exp(\tau_{a} \operatorname{sim}(\boldsymbol{l}_{i}, \boldsymbol{l}_{k}))}, \quad L_{a} = \sum_{\boldsymbol{x}_{i}, \boldsymbol{x}_{j} \in B} d_{i,j}^{s} \log \frac{d_{i,j}^{s}}{d_{i,j}^{l}},$$

$$(4)$$

where  $l_i = h_t(H_i)$  is the logits of  $x_i$ ,  $sim_+(\cdot, \cdot)$  is cosine similarity with min-max normalization, and  $\tau_a$  is a temperature coefficient controlling the loss strength. We employ min-max normalization

(instead of softmax) to sharpen slot supervision. Note that  $L_a$  is method-agnostic as long as the CL method has an FM backbone to support extracting semantic features. Finally with the cross-entropy task loss  $L_{ce}$ , the training loss is as  $L_{tr} = L_{ce} + \beta L_a$ , where  $\beta$  is a coefficient to balance the impact of  $L_a$ .

# 5 EXPERIMENTS

In the experiment part, we highlight the research question we will answer: *How and why does our CompSLOT benefit a large range of continual learning with foundation models?* To answer this, we compare algorithms with and without CompSLOT and perform ablation studies in section 5.2. We analyze the influences of hyperparameters in section H, investigate different backbones in section J, and visualize the slot extraction to analyze how CompSLOT enhances CL performance in section K.

#### 5.1 EXPERIMENTAL SETTINGS

**Baselines** To verify the universality of the proposed CompSLOT, we adopt a wide range of SOTA continual learners with foundation models, including: 1) **prompt-based methods**: CPrompt (Gao et al., 2024); 2) **representation-based methods**: ADAM+adapter (Zhou et al., 2025), Ran-PAC (McDonnell et al., 2023), EASE (Zhou et al., 2024b); 3) **Model-mixture-based methods**: CoFiMA (Marouf et al., 2024), FOSTER\* (Wang et al., 2022a), DER\* (Yan et al., 2021), MEMO\* (Zhou et al., 2023). Methods with a "\*" postfix indicate that they adopt a rehearsal process. Algorithms are implemented using the PILOT (Sun et al., 2025) platform with default hyperparameters. Methods with CompSLOT are denoted with a postfix "†". We also compare recent concept bottleneck models for continual learning, including CLG-CBM (Yu et al., 2025), and another concept knowledge plugin, SACK (Kundargi et al., 2025) with CLIP (Radford et al., 2019) integrated with CODA-Prompt (Smith et al., 2023).

**Benchmarks** We conduct experiments on compositional datasets, including CGQA and COBJ (Liao et al., 2024), and commonly used datasets, including ImageNet-R (Hendrycks et al., 2021). The former classification datasets contain a sufficient number of combinations of concepts, allowing for visual analysis and evaluating the compositionality. When comparing with other concept-based methods, we conduct experiments on CUB200 (Welinder et al., 2010) and CIFAR100 (Krizhevsky et al., 2009). We choose different continual task settings to evaluate different compositionality levels. Specifically, we denote "**F-S tasks**" as that the first task contains **F** classes and the following tasks contain **S** classes. For example, "50-10 tasks" means splitting 100 classes into six tasks with sequence of class numbers [50, 10, 10, 10, 10, 10, 10]. In the main context, we report 10-10 tasks results for CGQA. For results on other benchmarks, please refer to section I. All the experiments are conducted on a single Tesla V100 GPU and we analyze the computational cost in section M.

**Metrics** For continual training stage, we report the average accuracy of all tasks after training the last task  $\mathbf{A}\mathbf{A} = \frac{1}{T}\sum_{t=1}^T \mathbb{E}_{(x_{te},y)\in\mathcal{D}_{te}^t}[\Delta(\operatorname{pred}(x_{te}|P_T),y)]$ , the average cumulative accuracy for each task  $\mathbf{C}\mathbf{A} = \frac{1}{T}\sum_{t=1}^T \frac{1}{T-t+1}\sum_{u=t}^T \mathbb{E}_{(x_{te},y)\in\mathcal{D}_{te}^u}[\Delta(\operatorname{pred}(x_{te}|P_t),y)]$ , and average forgetting for each task  $\mathbf{F}\mathbf{F} = \frac{1}{T}\sum_{t=1}^T \mathbb{E}_{(x_{te},y)\in\mathcal{D}_{te}^t}[\Delta(\operatorname{pred}(x_{te}|P_t),y)] - \mathbf{A}\mathbf{A}$ , where  $\mathcal{D}_{te}^t$  is the testing dataset for task t and  $\Delta(\cdot,\cdot)$  is the equal function. After training on all continual tasks, specifically for CGQA and COBJ, we perform CFST on five compositional test suites including **sys**, **pro**, **sub**, **non**, **noc**, which contain novel recombinations, more combinations, shifting attributes, seen combinations, novel concepts of testing samples, respectively. We generate 300 few-shot tasks for each test suite. For clarity, we calculate the Harmonic mean (i.e.,  $\mathbf{H}\mathbf{n} = 3/(1/\operatorname{sys}+1/\operatorname{pro}+1/\operatorname{sub})$ ,  $\mathbf{H}\mathbf{r} = 2/(1/\operatorname{non}+1/\operatorname{noc})$ , as suggested in Liao et al. (2024). Then we report Hn and the ratio of Hn and Hr (i.e.,  $\mathbf{R} = \operatorname{Hn}/\operatorname{Hr}$ ). For detailed results on each compositional test suite, please refer to section G. Larger Hn and R indicate that the extracted features have better compositional generalization performance.

#### 5.2 RESULTS

**Overall results** We report the statistical results in Table 1. Across all baselines, CompSLOT consistently enhances performance, with the most significant improvement observed in ADAM+adapter (absolute gain: +7.550 in AA). Notably, CA and FF demonstrate consistent superiority over other

Table 1: Main result on CGQA. Methods with CompSLOT are denoted with a postfix " $\dagger$ ". Methods rehearse old samples are denoted with a postfix "\*". We report results over 3 trials with (mean  $\pm$  95% confidence interval).

	Continual			CFST	
Methods	AA (%) ↑	CA (%) ↑	FF (%) ↓	Hn (%) ↑	R↑
CPrompt	46.753±0.570	60.179±1.695	15.670±0.950	78.063±0.817	0.964
CPrompt †	48.537±0.427	61.483±1.645	$18.315 \pm 1.111$	79.091±1.086	0.969
ADAM + adapter	41.930±1.141	53.983±0.444	13.800±0.187	68.649±0.259	0.932
ADAM + adapter †	49.480±1.201	$60.989 \pm 0.641$	$12.896 \pm 0.379$	$74.335 \pm 0.572$	0.958
RanPAC	65.810±0.802	$75.504 \pm 0.318$	10.515±0.176	$78.868 \pm 0.918$	1.016
RanPAC †	$66.753 \pm 0.867$	$76.584 \pm 0.603$	$10.219 \pm 0.281$	$79.815 \pm 0.829$	1.032
EASE	47.657±1.494	$59.475 \pm 2.574$	$18.215 \pm 0.107$	$79.713 \pm 0.449$	0.996
EASE †	49.323±1.165	62.603±1.252	$22.470\pm2.472$	82.887±0.320	1.001
CoFiMA	65.107±0.508	73.227±1.047	15.248±0.542	86.711±0.483	1.011
CoFiMA †	66.170±0.578	$74.322 \pm 0.463$	$14.204 \pm 0.880$	$88.297 \pm 0.278$	1.017
FOSTER*	60.863±0.271	68.800±0.496	2.441±0.122	89.791±0.086	1.087
FOSTER* †	66.290±1.451	$71.828 \pm 2.619$	$6.470 \pm 5.770$	$89.910 \pm 0.710$	1.154
DER*	52.003±1.019	62.675±1.695	40.122±0.907	90.119±0.510	1.080
DER*†	54.900±1.093	$66.020{\pm}1.049$	$38.941 \pm 0.995$	$88.986 \pm 0.129$	1.096
MEMO*	56.553±1.804	66.462±0.702	$9.289 \pm 0.326$	82.425±1.282	1.029
MEMO* †	58.653±1.449	$68.037 {\pm} 1.459$	$8.944 {\pm} 0.268$	$84.003 \pm 1.451$	1.050

methods (except CPrompt and FOSTER, because the original methods do not perform well on the finished tasks, thus, forget less), indicates that our CompSLOT not only mitigates catastrophic forgetting of old tasks but also preserves strong forward adaptation to novel tasks. This robustness is primarily attributed to CompSLOT's improved compositional generalization (manifested by higher Hn and R scores), confirming its ability to learn latent conceptual units and dynamically compose them for robust classification across diverse methodological frameworks.

**Learning curve** Figure 3a shows the learning curves of all methods on the 10-10 tasks from CGQA. We observe that concept learning significantly improves continual learning performance across the entire training process, demonstrating its ability to stabilize learning and mitigate forgetting.

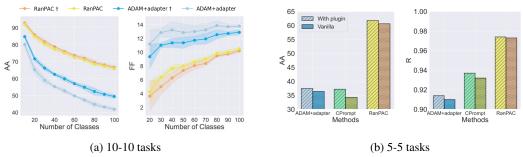


Figure 3: Learning curves and histograms of methods with and without CompSLOT on CGQA a) 10-10 tasks and b) 5-5 tasks. Slot is the case directly using the primitive representation and a cosine similarity classifier for the continual tasks.

**Long task sequence** Figure 3b presents the comparative performance analysis across a challenging long-task sequence of 5-5 CGQA tasks. The results reveal that CompSLOT consistently enhances both continual learning accuracy and compositional generalization performance, even when the slot attention module globally shared across all continual tasks. This finding underscores the remarkable robustness of employing slot attention mechanisms for boosting concept learning in CL scenarios. Notably, the stable improvement suggests that CompSLOT effectively captures transferable

compositional knowledge, enabling better adaptation across sequential tasks without task-specific customization.

**Ablation studies** To evaluate the contribution of each proposed component, we conduct comprehensive ablation experiments, with results presented in Table 2. First, to rule out the possibility that performance gains stem solely from **increased model capacity**, we expand the hidden representation dimensions (denoted as "+param") in RanPAC and CPrompt (see section E for details) to match the parameter count of RanPAC † and CPrompt †, respectively. We further perform the following controlled experiments: 1) **Primitive loss ablation**  $(L_n)$ : We remove the primitive loss term and replace the primitive selection mechanism with a simple slot averaging strategy (avg). 2) **Slot-selection** function ablation: We substitute the softmax operation in Equation 2 with alternative weighting methods, including: averaging (avg), sigmoid (sig), sign quantization (sign), and cosine similarity (cos). Across both methods, disabling

Table 2: Ablation results on CGQA.

Methods	$L_p$	$L_a$	AA (%) ↑	R↑
	<b>X</b> +param	Х	65.080	1.010
	Xavg	✓	58.220	0.969
	✓avg	✓	65.870	1.003
RanPAC	✓sig	/	65.950	1.020
	✓sign	✓	65.140	1.006
	√cos	/	63.910	0.989
	√soft	✓	66.753	1.032
	<b>X</b> +param	X	46.300	0.969
	Xavg	✓	40.230	0.952
	✓avg	1	47.690	0.958
CPrompt	✓sig	/	48.080	0.961
_	✓sign	✓	47.780	0.966
	<b>√</b> cos	✓	47.410	0.964
	✓soft	1	48.537	0.969

 $L_p$  or altering the slot-selection mechanism leads to significant degradations in AA and R scores, demonstrating the critical importance of each component. 1) The primitive loss  $L_p$  ensures intra-class consistency, which is vital for reliable primitive selection and, consequently, improved concept-level class understanding. On the other hand, using all slots indiscriminately allows less relevant concepts (e.g., background) to dilute class-relevant ones, leading to confusion. 2) The softmax-based weighting (as formulated in Equation 2) provides a selection with a convex combination of slots in one image to ensure the primitive representations of images are within an appropriate range, which makes the training robust. A more comprehensive ablation study can be found in section L.

# Comparing with other concept learning methods This paragraph compares CompSLOT with other concept learning methods, i.e., CLG-CBM (Yu et al., 2025) and SACK (Kundargi et al., 2025). We conduct experiments on 10-10 tasks CUB200 and CIFAR100 to show the superiority of CompSLOT with RanPAC. The results are shown in Table 3 with the top perfor-

Table 3: Comparison results on 10-10 tasks CUB200 and CIFAR100.

Datasets	SACK	CLG-CBM	CompSLOT
CUB200	71.78	85.40	88.38
CIFAR100	87.26	84.49	89.57

mance mentioned in their original papers. CompSLOT shows the best AA on coarse-grained and fine-grained benchmarks, because of benefiting from slot attention to extract concept information and the plug-and-play property that can be applied to alternative CL algorithms. Most importantly, CompSLOT fully utilizes the capability of the CL backbone and does not need extra interpretable tools, like ChatGPT.

#### 6 Conclusion

This work propose **CompSLOT**, a framework introducing **concept learning** into the continual learning paradigm for foundation models. The proposed **primitive selection mechanism** effectively extracts class-relevant concepts while maintaining robustness across extended task sequences. Meanwhile, the **primitive-logit knowledge distillation** mechanism enforces concept-based sample similarity regularization, enabling lightweight adaptation to diverse CL methods with foundation models. Experimental results confirm that the performance improvements stem from enhanced **compositional generalization**, offering a novel **concept-level perspective** for the continual learning community. A limitation of our current approach is that concept learning must precede providing conceptual self-supervision to the CL task. Future work will explore end-to-end integration of our mechanism into the continual learning pipeline. We hope this research inspires further advancements in developing resilient and interpretable vision models.

# 7 ETHICS STATEMENT

We hereby affirm our strict adherence to the ICLR Code of Ethics. We have carefully considered the ethical implications of our research throughout the entire process of study design, data collection, experimentation, and manuscript preparation, and we confirm that our work does not violate any of the principles outlined in the ICLR Code of Ethics. All datasets used, including ImageNet-R, CIFAR100, CUB200, CGQA, COBJ, were sourced in compliance with relevant usage guidelines, ensuring no violation of privacy. No personally identifiable information was used, and no experiments were conducted that could raise privacy or security concerns. We are committed to maintaining transparency and integrity throughout the research process.

#### 8 REPRODUCIBILITY STATEMENT

We are committed to ensuring the reproducibility of our research presented in this paper. To facilitate the replication of our results and the verification of our findings, we have provided comprehensive implementation details in section E. Additionally, the datasets we used, are publicly available, ensuring consistent and reproducible evaluation results.

#### REFERENCES

- Alessandro Achille, Tom Eccles, Loic Matthey, Chris Burgess, Nicholas Watters, Alexander Lerchner, and Irina Higgins. Life-long disentangled representation learning with cross-domain latent homologies. *Advances in Neural Information Processing Systems*, 31, 2018.
- Yuval Atzmon, Felix Kreuk, Uri Shalit, and Gal Chechik. A causal view of compositional zero-shot recognition. *Advances in Neural Information Processing Systems*, 33:1462–1473, 2020.
- Magdalena Biesialska, Katarzyna Biesialska, and Marta R Costa-Jussa. Continual lifelong learning in natural language processing: A survey. *arXiv preprint arXiv:2012.09823*, 2020.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Michael Chang, Tom Griffiths, and Sergey Levine. Object representations as fixed points: Training iterative refinement algorithms with implicit differentiation. *Advances in Neural Information Processing Systems*, 35:32694–32708, 2022.
- Hung-Jen Chen, An-Chieh Cheng, Da-Cheng Juan, Wei Wei, and Min Sun. Mitigating forgetting in online continual learning via instance-aware parameterization. *Advances in Neural Information Processing Systems*, 33:17466–17477, 2020a.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PmLR, 2020b.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder—decoder for statistical machine translation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, October 2014.
- Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9285–9295, 2022.
- Qiankun Gao, Chen Zhao, Yifan Sun, Teng Xi, Gang Zhang, Bernard Ghanem, and Jian Zhang. A unified continual learning framework with general parameter-efficient tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 11483–11493, 2023.

- Zhanxin Gao, Jun Cen, and Xiaobin Chang. Consistent prompting for rehearsal-free continual learning.
   In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 28463–28473, 2024.
  - Alexander L Gaunt, Marc Brockschmidt, Nate Kushman, and Daniel Tarlow. Differentiable programs with neural libraries. In *International Conference on Machine Learning*, pp. 1213–1222. PMLR, 2017.
  - Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. On the binding problem in artificial neural networks. *arXiv* preprint arXiv:2012.05208, 2020.
  - Stephen T Grossberg. Studies of mind and brain: Neural principles of learning, perception, development, cognition, and motor control, volume 70. Springer Science & Business Media, 2012.
  - Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 8340–8349, 2021.
  - Michael Hersche, Geethan Karunaratne, Giovanni Cherubini, Luca Benini, Abu Sebastian, and Abbas Rahimi. Constrained few-shot class-incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9057–9067, 2022.
  - Heinke Hihn and Daniel A Braun. Hierarchically structured task-agnostic continual learning. *Machine Learning*, 112(2):655–686, 2023.
  - Dieuwke Hupkes, Verna Dankers, Mathijs Mul, and Elia Bruni. Compositionality decomposed: How do neural networks generalise? *Journal of Artificial Intelligence Research*, 67:757–795, 2020.
  - Baoxiong Jia, Yu Liu, and Siyuan Huang. Improving object-centric learning with query optimization. In *The Eleventh International Conference on Learning Representations*, 2023.
  - Jindong Jiang, Fei Deng, Gautam Singh, and Sungjin Ahn. Object-centric slot diffusion. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
  - Ioannis Kakogeorgiou, Spyros Gidaris, Konstantinos Karantzalos, and Nikos Komodakis. Spot: Self-training with patch-order permutation for object-centric learning with autoregressive transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 22776–22786, June 2024.
  - Karthikeya Ramesh Kaushik and Andrea E Martin. Modelling compositionality and structure dependence in natural language. *arXiv preprint arXiv:2012.02038*, 2020.
  - Daniel Keysers, Nathanael Schärli, Nathan Scales, Hylke Buisman, Daniel Furrer, Sergii Kashubin, Nikola Momchev, Danila Sinopalnikov, Lukasz Stafiniak, Tibor Tihon, Dmitry Tsarkov, Xiao Wang, Marc van Zee, and Olivier Bousquet. Measuring compositional generalization: A comprehensive method on realistic data. In *International Conference on Learning Representations*, 2020.
  - Muhammad Gul Zain Ali Khan, Muhammad Ferjad Naeem, Luc Van Gool, Alain Pagani, Didier Stricker, and Muhammad Zeshan Afzal. Learning attention propagation for compositional zeroshot learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3828–3837, 2023.
  - Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
  - Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
  - 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.

- Avinash Kori, Francesco Locatello, Fabio De Sousa Ribeiro, Francesca Toni, and Ben Glocker. Grounded object-centric learning. In *The Twelfth International Conference on Learning Representations*, 2023.
  - Jonathan Krause, Hailin Jin, Jianchao Yang, and Li Fei-Fei. Fine-grained recognition without part annotations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5546–5555, 2015.
  - Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Technical report*, 2009.
  - Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
  - Shivanand Kundargi, Kowshik Thopalli, and Tejas Gokhale. Sequentially acquiring concept knowledge to guide continual learning. In *Second Workshop on Visual Concepts, CVPR*, 2025.
  - Songning Lai, Mingqian Liao, Zhangyi Hu, Jiayu Yang, Wenshuo Chen, Hongru Xiao, Jianheng Tang, Haicheng Liao, and Yutao Yue. Learning new concepts, remembering the old: Continual learning for multimodal concept bottleneck models. *arXiv*, 2024.
  - Brenden Lake and Marco Baroni. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *International conference on machine learning*, pp. 2873–2882. PMLR, 2018.
  - Minh Le, Huy Nguyen, Trang Nguyen, Trang Pham, Linh Ngo, Nhat Ho, et al. Mixture of experts meets prompt-based continual learning. *Advances in Neural Information Processing Systems*, 37: 119025–119062, 2024.
  - 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, pp. 3925–3934. PMLR, 2019.
  - Yuxiao Li, Eric J Michaud, David D Baek, Joshua Engels, Xiaoqing Sun, and Max Tegmark. The geometry of concepts: Sparse autoencoder feature structure. *Entropy*, 27(4):344, 2025.
  - Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
  - Yan-Shuo Liang and Wu-Jun Li. Inflora: Interference-free low-rank adaptation for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 23638–23647, 2024.
  - Weiduo Liao, Ying Wei, Mingchen Jiang, Qingfu Zhang, and Hisao Ishibuchi. Does continual learning meet compositionality? new benchmarks and an evaluation framework. *Advances in Neural Information Processing Systems*, 36, 2024.
  - Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *Advances in neural information processing systems*, 33:11525–11538, 2020.
  - David Lopez-Paz and Marc'Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
  - Hongyin Luo, Zhiyuan Liu, Huanbo Luan, and Maosong Sun. Online learning of interpretable word embeddings. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1687–1692, 2015.
  - Arun Mallya and Svetlana Lazebnik. PackNet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
  - Imad Eddine Marouf, Subhankar Roy, Enzo Tartaglione, and Stéphane Lathuilière. Weighted ensemble models are strong continual learners. In *European Conference on Computer Vision*, pp. 306–324. Springer, 2024.

- Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pp. 109–165. Elsevier, 1989.
  - Mark D McDonnell, Dong Gong, Amin Parvaneh, Ehsan Abbasnejad, and Anton Van den Hengel. Ranpac: Random projections and pre-trained models for continual learning. *Advances in Neural Information Processing Systems*, 36:12022–12053, 2023.
  - Jorge A Mendez and ERIC EATON. Lifelong learning of compositional structures. In *International Conference on Learning Representations*, 2021.
  - Brian Murphy, Partha Talukdar, and Tom Mitchell. Learning effective and interpretable semantic models using non-negative sparse embedding. In *Proceedings of COLING 2012*, pp. 1933–1950, 2012.
  - Tomáš Musil and David Mareček. Independent components of word embeddings represent semantic features. *arXiv preprint arXiv:2212.09580*, 2022.
  - Tushar Nagarajan and Kristen Grauman. Attributes as operators: Factorizing unseen attribute-object compositions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.
  - Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
  - Oleksiy Ostapenko, Pau Rodriguez, Massimo Caccia, and Laurent Charlin. Continual learning via local module composition. *Advances in Neural Information Processing Systems*, 34:30298–30312, 2021.
  - Kiho Park, Yo Joong Choe, Yibo Jiang, and Victor Veitch. The geometry of categorical and hierarchical concepts in large language models. *arXiv preprint arXiv:2406.01506*, 2024.
  - Benliu Qiu, Hongliang Li, Haitao Wen, Heqian Qiu, Lanxiao Wang, Fanman Meng, Qingbo Wu, and Lili Pan. CafeBoost: Causal feature boost to eliminate task-induced bias for class incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 16016–16025, June 2023.
  - Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
  - Ramesh Rahul and Chaudhari Pratik. Model Zoo: A growing brain that learns continually. In *International Conference on Learning Representations*, 2022.
  - Jathushan Rajasegaran, Munawar Hayat, Salman H Khan, Fahad Shahbaz Khan, and Ling Shao. Random path selection for continual learning. *Advances in Neural Information Processing Systems*, 32, 2019.
  - Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pp. 8821–8831, 2021.
  - Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollar, and Christoph Feichtenhofer. SAM 2: Segment anything in images and videos. In *The Thirteenth International Conference on Learning Representations*, 2025.
  - Mark B Ring. Child: A first step towards continual learning. *Machine Learning*, 28(1):77–104, 1997.

- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. *Advances in neural information processing systems*, 32, 2019.
  - Frank Ruis, Gertjan Burghouts, and Doina Bucur. Independent prototype propagation for zero-shot compositionality. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 10641–10653. Curran Associates, Inc., 2021.
  - Paul Ruvolo and Eric Eaton. ELLA: An efficient lifelong learning algorithm. In *International conference on machine learning*, pp. 507–515. PMLR, 2013.
  - Gautam Singh, Fei Deng, and Sungjin Ahn. Illiterate DALL-e learns to compose. In *International Conference on Learning Representations*, 2022.
  - James Seale Smith, Leonid Karlinsky, Vyshnavi Gutta, Paola Cascante-Bonilla, Donghyun Kim, Assaf Arbelle, Rameswar Panda, Rogerio Feris, and Zsolt Kira. Coda-prompt: Continual decomposed attention-based prompting for rehearsal-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11909–11919, June 2023.
  - Andrew Stange, Robert Lo, Abishek Sridhar, and Kousik Rajesh. Exploring the role of the bottleneck in slot-based models through covariance regularization. *arXiv* preprint arXiv:2306.02577, 2023.
  - Fei Sun, Jiafeng Guo, Yanyan Lan, Jun Xu, and Xueqi Cheng. Sparse word embeddings using 11 regularized online learning. In *Twenty-Fifth International Joint Conference on Artificial Intelligence*, 2016.
  - Hai-Long Sun, Da-Wei Zhou, De-Chuan Zhan, and Han-Jia Ye. Pilot: A pre-trained model-based continual learning toolbox. *SCIENCE CHINA Information Sciences*, 68(4):147101, 2025.
  - Qing Sun, Fan Lyu, Fanhua Shang, Wei Feng, and Liang Wan. Exploring example influence in continual learning. *Advances in Neural Information Processing Systems*, 35:27075–27086, 2022.
  - Zhicheng Sun, Yadong Mu, and Gang Hua. Regularizing second-order influences for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 20166–20175, June 2023.
  - Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 398–414, 2022a.
  - Liyuan Wang, Jingyi Xie, Xingxing Zhang, Mingyi Huang, Hang Su, and Jun Zhu. Hierarchical decomposition of prompt-based continual learning: Rethinking obscured sub-optimality. *Advances in Neural Information Processing Systems*, 36, 2024.
  - Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. In *European Conference on Computer Vision*, pp. 631–648. Springer, 2022b.
  - 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 (CVPR)*, pp. 139–149, 2022c.
  - P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona. Caltech-UCSD Birds 200. Technical Report CNS-TR-2010-001, California Institute of Technology, 2010.
  - Ross Wightman. Pytorch image models, 2019.
  - Tailin Wu, Megan Tjandrasuwita, Zhengxuan Wu, Xuelin Yang, Kevin Liu, Rok Sosič, and Jure Leskovec. ZeroC: A neuro-symbolic model for zero-shot concept recognition and acquisition at inference time. *arXiv* preprint arXiv:2206.15049, 2022.

36:50932–50958, 2023.

756

757

758

	,	
Hi	roaki Yamagiwa, Momose Oyama, and Hidetoshi Shimodaira. Discovering universal geometry i	in
	embeddings with ICA. In Proceedings of the 2023 Conference on Empirical Methods in Nature	
	Language Processing, pp. 4647–4675, 2023.	
-	55.05 Studge 1. 1000550118, Pp. 10.17 10.10, 20.201	
Sh	ipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class	SS
	ncremental learning. In Proceedings of the IEEE/CVF conference on computer vision and patter	
	recognition, pp. 3014–3023, 2021.	
Sir	Han Yang, Tuomas Oikarinen, and Tsui-Wei Weng. Concept-driven continual learning. Transa	c-
	tions on Machine Learning Research, 2024. ISSN 2835-8856.	
	tao Yang, Jie Zhou, Xuanwen Ding, Tianyu Huai, Shunyu Liu, Qin Chen, Yuan Xie, and Lian	
	He. Recent advances of foundation language models-based continual learning: A survey. ACI	И
	Computing Surveys, 57(5):1–38, 2025.	
	Yu, Haoyu Han, Zhe Tao, Hantao Yao, and Changsheng Xu. Language guided concept bottlened	
	models for interpretable continual learning. In Proceedings of the Computer Vision and Patter	n
	Recognition Conference, pp. 14976–14986, 2025.	
_		
	ngwei Zhang, Liyuan Wang, Guoliang Kang, Ling Chen, and Yunchao Wei. Slca: Slow learns	
	with classifier alignment for continual learning on a pre-trained model. In <i>Proceedings of the</i>	ıe
	EEE/CVF International Conference on Computer Vision, pp. 19148–19158, 2023.	
D.,	Weight of Weight of the New Andrew Annal of 602 and a few Tenner	1
	-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Toward	IS
	memory-efficient class-incremental learning. In <i>ICLR</i> , 2023.	
Da	-Wei Zhou, Hai-Long Sun, Jingyi Ning, Han-Jia Ye, and De-Chuan Zhan. Continual learning wit	th
	ore-trained models: A survey. In <i>Proceedings of the Thirty-Third International Joint Conference</i>	
	on Artificial Intelligence (IJCAI), 2024a.	·e
	m Artificial Intelligence (IJCAI), 2024a.	
Da	-Wei Zhou, Hai-Long Sun, Han-Jia Ye, and De-Chuan Zhan. Expandable subspace ensemble for	or
	pre-trained model-based class-incremental learning. In <i>Proceedings of the IEEE/CVF Conference</i>	
	on Computer Vision and Pattern Recognition, pp. 23554–23564, 2024b.	-
	, , , , , , , , , , , , , , , , , , , ,	
Da	-Wei Zhou, Qi-Wei Wang, Zhi-Hong Qi, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Clas	s-
	ncremental learning: A survey. IEEE Transactions on Pattern Analysis and Machine Intelligence	e,
	2024c.	
_		
	-Wei Zhou, Zi-Wen Cai, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Revisiting class-increment	
	earning with pre-trained models: Generalizability and adaptivity are all you need. <i>International</i>	лl
	Journal of Computer Vision, 133(3):1012–1032, 2025.	
<b>.</b>	dans 7 and Champhone 7 hours Heisbook 7 hours Without 1 and D. 1	_ 4
	kiong Zou, Shanghang Zhang, Haichen Zhou, Yuhua Li, and Ruixuan Li. Compositional few-sho	Эt
	class-incremental learning. arXiv preprint arXiv:2405.17022, 2024.	
A	APPENDIX	
C'	ONTENTS	
	71117110	
1	Introduction	1
2	Related Works	3
_	INCIACU WOINS	J
3	Preliminaries	4
	15	
	1 1	

Ziyi Wu, Jingyu Hu, Wuyue Lu, Igor Gilitschenski, and Animesh Garg. Slotdiffusion: Object-centric

generative modeling with diffusion models. Advances in Neural Information Processing Systems,

810	4	Methods	4
811 812		4.1 Concept Learning	4
813			6
814		4.2 Method-agnostic Primitive-Logit Knowledge Distillation	6
815	_		_
816	5	Experiments	7
817 818		5.1 Experimental Settings	7
819		5.2 Results	7
820			
821	6	Conclusion	9
822			
823	7	<b>Ethics Statement</b>	10
824	•	2	
825 826	8	Reproducibility Statement	10
827	0	Reproducibility Statement	10
828			
829	A	Appendix	15
830			
831	В	Discussions	17
832 833			
834	$\mathbf{C}$	Additional Related Works	17
835			
836	D	Theorem	17
837		D.1 Proof of Theorem 1	18
838			
839		D.2 Proof of Theorem 2	18
840 841			
842	$\mathbf{E}$	Hyperparameters and Experimental Settings	19
843			
844	F	Pseudo Code	20
845			
846	G	Detail CFST Results	21
847 848			
849	Н	Influences of Hyperparameters	21
850		• • •	
851	I	Results on Other Benchmarks	24
852	-	results on other benefinarias	
853	J	Results on Other Backbones	24
854 855	J	Results off Other Backbolles	4
856	17	<b>T</b> 70 <b>10</b> (1)	25
857	K	Visualization	25
858			
859	L	Additional Ablation Studies	30
860			
861	M	Algorithm Efficiency Analysis	30
862 863			
500	N	Use of Large Language Models	31

# **B** DISCUSSIONS

Classification Bias from a Concept-Combination Perspective The root cause of compromised stability and plasticity often lies in sub-optimal classifier design, particularly when classifiers develop reliance on inaccurate or incomplete feature representations due to concept biases in training data. To illustrate, consider a scenario where the current vision task  $T_1$  contains two specific classes (a human standing by a tree and a human inside a boat) alongside other classes that lack human-related concepts. In such cases, the learned classifier might develop an over-reliance on distinguishing these two classes based solely on tree and boat concepts while neglecting the more critical human attribute. This limited conceptual understanding creates significant generalization problems when encountering unseen concept combinations. For instance, during task  $T_2$ , a novel image labeled a pig inside a boat would likely receive disproportionately high logits for the human inside a boat class due to the classifier's inability to properly disentangle object-class relationships from spatial-contextual cues. Conversely, a human inside a boat image might similarly activate the pig inside a boat class predictions. This conceptual entanglement manifests as catastrophic forgetting in  $T_1$  (as evidenced by diminished global accuracies post-training on  $T_2$ ) and severely hampers plasticity for  $T_2$  through incorrect plastic responses to novel concept combinations.

Whether concept sharing is a common phenomenon in the real-world? In real-world scenarios, concept sharing is quite common, like, in fine-grained classification cases such as CUB200, and in images with massive objects such as COBJ. This phenomenon is also discussed in other works. For example, Welinder et al. (2010) claims that fine-grained bird classes share some basic parts, and Krause et al. (2015) claims that fine-grained categories share similar shapes. In the experimental results, CompoSLOT consistently brings significant improvements to continual learning algorithms on these real-world cases. In contrast, datasets like CIFAR, which have relatively little concept sharing, are uncommon in complicated real-world scenarios.

#### C ADDITIONAL RELATED WORKS

**Continual learning from scratch** To mitigate forgetting previously learned vision classification tasks and achieve fast adaptation to new ones, the continual learning community has developed three main families of approaches that do not utilize a pre-trained foundation model:

- 1. Rehearsal-based methods (Achille et al., 2018; Rolnick et al., 2019; Rahul & Pratik, 2022; Hersche et al., 2022; Sun et al., 2022; Qiu et al., 2023; Sun et al., 2023) store memory-efficient samples or features from past tasks for reviewing knowledge. However, such buffers can cause significant memory overload as the number of tasks increases.
- 2. Regularization-based methods (Lopez-Paz & Ranzato, 2017; Li & Hoiem, 2017; Kirkpatrick et al., 2017; Achille et al., 2018; Hersche et al., 2022) constrain gradient updates to preserve important knowledge from old tasks, but this can limit the adaptation capability to new tasks.
- 3. Architecture-based methods (Mallya & Lazebnik, 2018; Douillard et al., 2022; Ring, 1997; Ruvolo & Eaton, 2013; Gaunt et al., 2017; Li et al., 2019; Rajasegaran et al., 2019; Chen et al., 2020a; Mendez & EATON, 2021; Ostapenko et al., 2021; Rahul & Pratik, 2022; Hihn & Braun, 2023) aim to create new modules for upcoming tasks, making the determination of module composition crucial for different tasks.

# D THEOREM

Firstly, we define *concepts* as the ground truth slot decomposition of an image. Since slot attention exhibits permutation equivalence w.r.t. the order of the slots (and masks) (Locatello et al., 2020), we regard  $\{\mathcal{S}, \mathcal{A}\}$  as the corresponding set representations of  $\{S, A\}$ , where  $\mathcal{S} = \{s_i\}_{i=1}^K$  and  $\mathcal{A} = \{a_i\}_{i=1}^K$ , with  $s_i \in \mathbb{R}^{D_s}$  and  $a_i \in \mathbb{R}^N$  being the *i*-th row of S and S, respectively.

**Definition 3** (Concept & Disentanglement, equivalent to Def. 1). Let x be an image, then  $\{S, A\}$  is a *disentangled* decomposition of x (a.k.a., *concepts* and corresponding *attention regions*), if 1)  $a_i, a_j \in A, a_i \in \mathbb{R}^N_+, a_i \perp a_j$ , and 2) S satisfies  $\arg\min_{S} \sum_{s_i, s_j \in S} |\sin(s_i, s_j)|$ , where  $s_i \in S$ 

 $\mathbb{R}^{D_s}$ , and  $|\cdot|$  is the absolute value,  $\perp$  is orthogonal symbol, and  $\mathrm{sim}(\cdot,\cdot)$  is an arbitrary similarity score function, e.g., cosine similarity.

Remark 2 (**Requirement 1**: Disentanglement). The competitive spatial attention and the limited capability of a slot naturally achieve the orthogonality of  $\mathcal{A}$ . In practice,  $\mathcal{S}$  is encouraged to be orthogonal (slots bind to different concepts in x) but not ideal since there are some semantically similar concepts, e.g., grass and leaves. Such a disentanglement structure is also mentioned in Park et al. (2024); Li et al. (2025).

**Definition 4** (Primitives, equivalent to Def. 2). Let  $\mathcal{X}^y, \mathcal{S}^y$  be an image set labeled y and the corresponding set of concept sets, and  $\mathcal{S} \in \mathcal{S}^y$ , then a concept subset  $\mathcal{P} \subset \mathcal{S}$  is *primitives* of  $\mathcal{S}$ , if  $\forall \mathcal{S}' \in \mathcal{S}^y, \mathcal{P} \subset \mathcal{S}'$ .

Remark 3. Although  $\mathcal{P}$  is defined at the image level, we can also say that it is unambiguously y's primitive set, denoted  $\mathcal{P}^y$ . In general,  $\mathcal{P} \neq \mathcal{S}$ , because there are always image-specific concepts in the image, e.g., background.

**Theorem 2** (**Requirement 2**: Intra-class consistency, equivalent to Theorem. 1). Consider  $S_1, S_2 \in S^y$  and two corresponding **largest** primitive sets  $\mathcal{P}_1 \subset S_1, \mathcal{P}_2 \subset S_2$  are identical, i.e.,  $\mathcal{P}_1 = \mathcal{P}_2$  and  $||\mathcal{P}_1|| = M$ , where  $||\cdot||$  is the cardinality of set. In other word, consider the pair-wise ordered sets  $\{\mathcal{P}_1^{\circ}, \mathcal{P}_2^{\circ}\} = \mathrm{match}(\mathcal{P}_1, \mathcal{P}_2)$ , where  $\mathrm{match}(\cdot, \cdot)$  is a matching algorithm (without loss of generality, Hungarian algorithm (Kuhn, 1955)), then the corresponding matched concepts should be the same:  $\mathcal{P}_1^{\circ} = \{s_1^1\}_{i=1}^M, \mathcal{P}_2^{\circ} = \{s_i^2\}_{i=1}^M$  and  $\forall i \in \{1, \ldots, M\}, \sin(s_1^1, s_i^2) = 1$ .

**Theorem 3** (Requirement 3: Inter-class concept sharing). *If there is a shared primitive subset between*  $y_1, y_2$ , all images in  $\mathcal{X}^{y_1}, \mathcal{X}^{y_2}$  should contain this subset. If  $\exists \mathcal{P} \subset \mathcal{P}^{y_1}, \mathcal{P} \subset \mathcal{P}^{y_2}, ||\mathcal{P}|| = M > 0$ , then  $\forall \mathcal{S}_1 \in \mathcal{S}^{y_1}, \mathcal{S}_2 \in \mathcal{S}^{y_2}, \mathcal{P} \subset \mathcal{S}_1, \mathcal{P} \subset \mathcal{S}_2$ .

Remark 4 (**Requirement 4**: Inter-task consistency). After trained on future tasks, the concept sets of the same x should be not changed. In T CL tasks,  $\forall x \in \mathcal{X}^u, 1 \leq u < T$ , consider the extracted concept sets  $\{\mathcal{S}^t\}_{t=u}^T$  after task  $t \in \{u, \dots, T\}$ , then  $\forall t_1, t_2 \in \{u, \dots, T\}, \mathcal{S}^{t_1} = \mathcal{S}^{t_2}$ .

CompoSLOT implicitly encourages **Requirement 1** via slot attention's soft-clustering and supports **Requirements 3–4** empirically (Figure 1) and the visualization experiments in section K. **Requirement 2** is enforced through a primitive loss, described in Equation 3, ensuring slot stability across class instances.

#### D.1 PROOF OF THEOREM 1

*Proof.* By the definition of  $\mathcal{P}_1$ ,  $\mathcal{P}_2$  as the *primitive* sets of  $\mathcal{S}_1$ ,  $\mathcal{S}_2$ , respectively, and that  $\mathcal{S}_1$ ,  $\mathcal{S}_2 \in \mathcal{S}^y$ , without loss of generality,  $\mathcal{P}_2$  is also a primitive set of  $\mathcal{S}_1$ . Thus,  $\mathcal{P}_2 \subset \mathcal{S}_1$ . Assume, for the sake of contradiction, that there exists a concept s, such that  $s \in \mathcal{P}_2$  and  $s \notin \mathcal{P}_1$ , i.e.,  $\mathcal{P}_1 \neq \mathcal{P}_2$ . Since  $\mathcal{P}_1$  is the **largest** *primitive* set of  $\mathcal{S}_1$ , we must have  $\mathcal{P}_2 \subseteq \mathcal{P}_1$  and  $\forall \mathcal{P} \subseteq \mathcal{P}_1, s \notin \mathcal{P}$ . This contradicts our initial assumption that  $s \in \mathcal{P}_2$ .

Therefore, the theorem holds.

Remark 5. The matching algorithm facilitates concept alignment across different sets, thereby enabling the computation of our proposed evaluation metrics in section H as well as supporting the visualizations presented in section K. However, this alignment process introduces significant computational overhead that renders it impractical for integration within our distillation framework. To address this limitation, we propose an attention-based primitive selection mechanism (detailed in section 4.1) that ensures permutation invariance to concept ordering in the extracted primitives, effectively eliminating the need for explicit concept matching. This design choice maintains computational efficiency while preserving the critical semantic relationships required for reliable evaluation and visualization.

#### D.2 PROOF OF THEOREM 2

*Proof.* Assume, for the sake of contradiction, that there exists  $\mathcal{P}', \mathcal{S}'$  and  $\mathcal{P}' \subset \mathcal{P}^{y_1}, \mathcal{P}' \subset \mathcal{P}^{y_2}, ||\mathcal{P}'|| = M > 0, \mathcal{S}' \in \mathcal{S}^{y_1}$  (or  $\mathcal{S}^{y_2}$ ), such that  $\mathcal{P}' \not\subset \mathcal{S}'$ . By the definition of  $\mathcal{P}^{y_1}$  as the primitive set for all  $\mathcal{S} \in \mathcal{S}^{y_1}$ , thus  $\mathcal{P}' \subset \mathcal{S}'$ . This contradicts our initial assumption that  $\mathcal{P}' \not\subset \mathcal{S}'$ .

Table 4: Detail hyperparameters for **concept learning stage** in our main experiments.

Hyper-parameters	Value
Optimizer	Adam
LR scheduler	Cosine
LR (1-st task)	1e-4
LR (others)	1e-5
LR (min)	1e-8
Batch size	256
Weight decay	0
Epoch	50
$D_s$	128
K	10
Slot refinemnt iterations $N_s$	5
Slot decoder hidden embedding dim	Linear with ReLU (128 $\rightarrow$ 256 $\rightarrow$ 256 $\rightarrow$ 768)
$ au_t$	100
$\alpha$	10
$ au_p$	10

Table 5: Detail hyperparameters for concept knowledge distillation stage in our main experiments.

Methods	l B	$ au_a$
CD	1.0	
CPrompt	10	0.05
ADAM + adapter	10	0.5
RanPAC	15	0.5
EASE	10	0.1
CoFiMA	1	0.001
FOSTER	2	0.05
DER	7	0.01
MEMO	0.05	0.1

Therefore, the theorem holds.

# E HYPERPARAMETERS AND EXPERIMENTAL SETTINGS

The hyperparameter settings for the concept learning stage are summarized in Table 4, with key values tuned through validation. For the concept knowledge distillation phase, we maintain fairness in comparison by adopting the platform-default hyperparameters from the PILOT framework for both standard CL baselines and CompoSLOT-enhanced variants, with additional parameters introduced in section 4.2 detailed in Table 5. All configurations employ an 80-20 train-validation split using a randomly sampled validation set. To ensure consistent model capacity across methods, all algorithms utilize the ViT-B/16 backbone pretrained on ImageNet as the shared feature extractor. The backbone parameters are sourced from the Python timm (Wightman, 2019) package.

For ablation studies specifically examining CompoSLOT's impact, we appropriately scale model capacities through expanded hidden representations: RanPAC: Increased feedforward layer width (ffn\_num) from 64 to 256; CPrompt: Extended prompt length (prompt\_len) from 50 to 65 tokens. These adjustments ensure fair comparison by matching representational capacity when introducing our architectural modifications, enabling more reliable evaluation of CompoSLOT's actual contribution beyond simple capacity increases.

# F PSEUDO CODE

1026

1027 1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040 1041

1062 1063 1064 In the main paper, we propose a two-stage procedure, including concept learning (aiming to extract concept-level representation by performing slot representation training and primitive selection) and concept knowledge distillation (aiming to distill sample-wise concept-based similarity into logits). We summarize the training framework of CompoSLOT in Algorithm 1. Specifically, we perform concept learning in Lines 4-9. The slot attention and primitive selection module are initialized at first. For each batch of samples in task t, we perform Algorithm 2 and use the obtained primitive loss and reconstruction loss to train slot attention and primitive selection modules in Line 6. After E epochs of training, we perform concept knowledge distillation in Lines 11-18. We calculate pair-wise primitive similarity and obtain primitive-logit alignment loss with Equation 4 in Line 15. We detail the slot representation learning in Algorithm 2. Specifically, we first obtain semantic patch features in Line 3. Then, we use slot attention module to decompose it into a set of slots in Line 4. Next, we reconstruct the patch feature and obtain the reconstruction loss in Lines 6-8. After that, we calculate the primitives in Lines 10-12 and obtain primitive loss with Equation 3 in Lines 14-15.

### Algorithm 1 Continual Learning Framework

```
1042
          1: Input: # tasks T, tasks \mathcal{D}^1, \ldots, \mathcal{D}^T, # epochs E, candidate CL method CL(\cdot | \theta_f, \theta_h).
1043
          2: Initialize slot attention and primitive selection module.
1044
          3: for t from 1 to T do
1045
                 /* Concept Learning */
          4:
1046
          5:
                 for i from 1 to E do
1047
          6:
                    Sample a batch of images (x, y) \sim \mathcal{D}^t.
1048
          7:
                    Perform Algorithm 2 to obtain primitives s^p, contrastive primitive loss L_p, and reconstruc-
1049
                    tion loss L_{re}.
1050
          8:
                    L_{slot} = L_{re} + \alpha L_p.
1051
          9:
                    Backward loss and update.
          10:
1052
                 end for
                 /* Concept Knowledge Distillation */
         11:
1053
         12:
                 for i from 1 to E do
1054
         13:
                    Sample a batch of images (x, y) \sim \mathcal{D}^t.
1055
                    Perform Algorithm 2 to obtain primitives s^p without collecting gradients.
         14:
1056
         15:
                    Perform CL method to obtain logits CL(x|\theta_f,\theta_h) and task loss L_{ce}.
1057
         16:
                    Calculate primitive-logit alignment loss L_a.
1058
         17:
                    L_{tr} = L_{ce} + \beta L_a.
         18:
                    Backward loss and update.
         19:
                 end for
1061
         20: end for
```

#### Algorithm 2 Slot Representation Learning

```
1065
          1: Input: Image batch \{x_i\}_{i=1}^B, CL backbone \theta_f, # slots K, slot dimension D_s, # epochs E,...
          2: Output: Primitive s^p, contrastive primitive loss L_p, reconstruction loss L_{re}.
1067
          3: Obtain semantic patch features E = f(x_i | \theta_f)[1:].
1068
          4: Obtain a set of K slots and the corresponding attentions \{S, A\}.
1069
          5: /* Reconstruction Loss */
1070
          6: Add position embedding for each patch: S_n' = S \oplus pos_n.
1071
          7: Decode S' and re-construct using A: E = A^{\top} d(S' | \theta_d).
1072
          8: L_{re} = || E - \tilde{E} ||_2.
          9: /* Primitive Selection */
1074
         10: Obtain Mapped slots \bar{S} = \tanh(\text{Linear}(\text{LN}(S))).
1075
         11: Obtain weights for each slot w_p = \sigma(\tau_p SK^p).
         12: Obtain primitive s^p = w_p^\top \bar{S}.
1077
         13: /* Contrastive Primitive Loss */
1078
         14: Obtain normalized similarity d_{i,j}^y and softmax primitive similarity d_{i,j}^s for image sample x_i, x_j.
1079
         15: Obtain primitive loss L_p.
```

## G DETAIL CFST RESULTS

The statistical analysis of each compositional test suite for the 10-10 tasks CGQA benchmark is presented in Table 6. All accuracy metrics are reported with their corresponding ±95% confidence intervals to quantify statistical significance. The key metrics include:

- Hn: Harmonic mean of compositional testing metrics (systematicity sys, productivity pro, substitutivity sub);
- Hr: Harmonic mean of reference testing metrics (Non-novel **non**, Not compositional **noc**);
- Ha: Harmonic mean across all test types;
- R=Hn/Hr: Ratio measuring compositional generalization improvement.

The results consistently demonstrate superior performance in both R and Hn (except for DER †), confirming CompoSLOT's ability to enhance compositional generalization, particularly for systematicity and productivity properties. This aligns with our hypothesis that the slot plugin improves compositional reasoning capabilities. However, as previously reported in Liao et al. (2024) for ViT-based architectures, we observe no significant improvement in substitutivity, suggesting inherent limitations of ViT feature extractors in dealing with attribute shifting (e.g., color).

# H INFLUENCES OF HYPERPARAMETERS

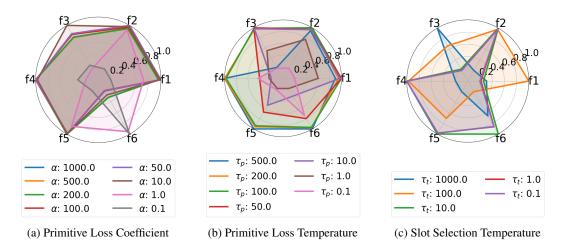


Figure 4: Radars of different hyperparameters in slot representation learning.

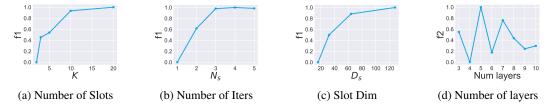


Figure 5: Line charts of different hyperparameters in slot attention architecture.

In this section, we investigate the effect of the introduced hyperparameters in the slot module w.r.t. the slot extraction performance and in the primitive-logit alignment loss. Without loss of generality, we report the model performance after training on the second task of 10-10 tasks CGQA in this section.

**Metrics** We learn slot representation S, attention mask A, and primitive representation  $s^t$  as intermediate products of the forwarding process. Thus, it is necessary to design quantitative metrics to represent the performance of the learned slot as follows:

Table 6: Detail CFST results. We report the average with  $\pm$  95% confidence interval.

Methods	sys	pro	sub	Hn
CPrompt	73.933±1.552	75.367±1.014	<b>85.967</b> ± <b>0.858</b>	78.063±0.817
CPrompt †	75.133±1.835	<b>78.133</b> ± <b>0.971</b>	84.600±0.514	<b>79.091</b> ± <b>1.086</b>
ADAM + adapter	63.400±0.244	68.667±0.838	74.833±0.107	68.649±0.259
ADAM + adapter †	68.533±0.962	<b>75.033</b> ± <b>0.533</b>	<b>80.400</b> ± <b>0.092</b>	<b>74.335</b> ± <b>0.572</b>
RanPAC	74.867±0.912	78.567±0.509	<b>83.667</b> ± <b>1.536</b>	78.868±0.918
RanPAC †	75.833±1.764	<b>80.600</b> ± <b>0.800</b>	83.433±1.783	<b>79.815</b> ± <b>0.829</b>
EASE	74.900±0.423	80.567±0.629	84.233±0.282	79.713±0.449
EASE †	78.267±0.509	<b>84.633</b> ± <b>0.509</b>	<b>86.200</b> ± <b>0.480</b>	<b>82.887</b> ± <b>0.320</b>
CoFiMA	83.100±1.135	86.767±0.267	90.600±0.606	86.711±0.483
CoFiMA †	84.467±0.324	<b>88.967</b> ± <b>0.373</b>	<b>91.767</b> ± <b>0.141</b>	<b>88.297</b> ± <b>0.278</b>
FOSTER	86.900±0.514	91.400±0.489	<b>91.233±0.971</b>	89.791±0.086
FOSTER †	87.600±0.606	<b>91.733</b> ± <b>0.979</b>	90.500±0.733	<b>89.910</b> ± <b>0.710</b>
DER	87.700±0.160	<b>91.733±0.838</b>	<b>91.033±0.828</b>	<b>90.119±0.510</b>
DER †	86.567±0.509	90.300±0.666	90.200±0.320	88.986±0.129
MEMO	78.233±2.189	82.500±1.201	87.033±0.541	82.425±1.282
MEMO †	79.733±1.248	<b>85.133</b> ± <b>1.816</b>	<b>87.533</b> ± <b>1.432</b>	<b>84.003</b> ± <b>1.451</b>
Methods	non	noc	Hr	R
Methods  CPrompt CPrompt †	non	noc	Hr	R
	76.400±0.973	86.033±0.437	80.926±0.360	0.964
	77.167±0.681	<b>86.533</b> ± <b>0.601</b>	<b>81.580</b> ± <b>0.407</b>	<b>0.969</b>
CPrompt	76.400±0.973	86.033±0.437	80.926±0.360	0.964
CPrompt CPrompt † ADAM + adapter	76.400±0.973 77.167±0.681 66.167±0.930	86.033±0.437 <b>86.533</b> ± <b>0.601</b> 82.867±0.615	80.926±0.360 <b>81.580</b> ± <b>0.407</b> 73.580±0.809	0.964 <b>0.969</b> 0.932
CPrompt		86.033±0.437	80.926±0.360	0.964
CPrompt †		86.533±0.601	<b>81.580</b> ± <b>0.407</b>	<b>0.969</b>
ADAM + adapter		82.867±0.615	73.580±0.809	0.932
ADAM + adapter †		84.967±0.192	<b>77.516</b> ± <b>0.323</b>	<b>0.958</b>
RanPAC		80.033±0.833	<b>77.574</b> ± <b>0.813</b>	1.016
CPrompt CPrompt †  ADAM + adapter ADAM + adapter †  RanPAC RanPAC †  EASE	$egin{array}{c c} 76.400 \pm 0.973 \\ 77.167 \pm 0.681 \\ \hline 66.167 \pm 0.930 \\ 71.267 \pm 0.417 \\ \hline 75.267 \pm 1.063 \\ 75.600 \pm 0.606 \\ \hline 76.400 \pm 0.666 \\ \hline \end{array}$	86.033±0.437 <b>86.533</b> ± <b>0.601</b> 82.867±0.615 <b>84.967</b> ± <b>0.192</b> <b>80.033</b> ± <b>0.833</b> 79.133±1.593 83.967±0.141	$80.926\pm0.360$ $81.580\pm0.407$ $73.580\pm0.809$ $77.516\pm0.323$ $77.574\pm0.813$ $77.314\pm0.440$ $80.004\pm0.420$	0.964 <b>0.969</b> 0.932 <b>0.958</b> 1.016 <b>1.032</b> 0.996
CPrompt CPrompt †  ADAM + adapter ADAM + adapter †  RanPAC RanPAC †  EASE EASE †  CoFiMA		$86.033\pm0.437$ $86.533\pm0.601$ $82.867\pm0.615$ $84.967\pm0.192$ $80.033\pm0.833$ $79.133\pm1.593$ $83.967\pm0.141$ $85.867\pm0.541$ $88.233\pm0.509$	$80.926\pm0.360$ $81.580\pm0.407$ $73.580\pm0.809$ $77.516\pm0.323$ $77.574\pm0.813$ $77.314\pm0.440$ $80.004\pm0.420$ $82.775\pm0.255$ $85.729\pm0.353$	0.964 0.969 0.932 0.958 1.016 1.032 0.996 1.001 1.011
CPrompt CPrompt †  ADAM + adapter ADAM + adapter †  RanPAC RanPAC †  EASE EASE †  CoFiMA CoFiMA †  FOSTER	76.400±0.973 77.167±0.681 66.167±0.930 71.267±0.417 75.267±1.063 75.600±0.606 76.400±0.666 79.900±0.185 83.367±0.594 85.600±0.733	86.033±0.437 86.533±0.601 82.867±0.615 84.967±0.192 80.033±0.833 79.133±1.593 83.967±0.141 85.867±0.541 88.233±0.509 89.233±0.385 76.433±0.557	$80.926\pm0.360$ $81.580\pm0.407$ $73.580\pm0.809$ $77.516\pm0.323$ $77.574\pm0.813$ $77.314\pm0.440$ $80.004\pm0.420$ $82.775\pm0.255$ $85.729\pm0.353$ $87.378\pm0.544$ $82.592\pm0.285$	0.964 <b>0.969</b> 0.932 <b>0.958</b> 1.016 <b>1.032</b> 0.996 <b>1.001</b> 1.011 <b>1.017</b>

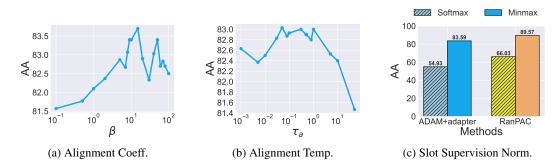


Figure 6: Line charts of different hyperparameters in primitive-logit knowledge distillation.

- Primitive-label matching score:  $\mathbf{f1} = -\text{MAE}(d^s, d^y)$ , where  $d^s$  and  $d^y$  are described in Equation 3.
- **Primitive-concept matching score:**  $\mathbf{f2} = -\text{MAE}(d^s, d^c)$ , where  $d^c$  is similar with  $d^y$  but the one-hot label is replaced with the multi-hot concept label. Note that the concept label is only used to analyze the performance of the learned slots and is never seen during training.
- Task-wise matched attention mask mIoU:  $\mathbf{f3} = \mathrm{Mean}_t\{\mathrm{IoU}(\mathcal{A}_{t-1}^\circ, \mathcal{A}_t^\circ)\}$ , where  $\mathrm{IoU}(\cdot, \cdot)$  is the intersection over union metric and  $\mathcal{A}_{t-1}^\circ, \mathcal{A}_t^\circ$  are matched attention sets (by Hungarian algorithm) extracted from the same image by the learners trained after task t-1 and t, respectively.
- Task-wise weighted attention mask mIoU:  $\mathbf{f4} = \text{Mean}_t\{\text{IoU}(w_{p,t-1} \top A_{t-1}, w_{p,t}^\top A_t)\}.$
- Task-wise matched slot matching score:  $\mathbf{f5} = -\text{MAE}(S_{t-1}^{\circ}, S_t^{\circ})$ .
- Task-wise primitive matching score:  $\mathbf{f6} = -\text{Mean}_x\{\text{MAE}(s_x^{t-1}, s_x^t)\}.$

For clarity, the matching scores are normalized to [0,1] to align with the range of mIoU. A large value of any metric above indicates a better performance according to the corresponding assessment.

**Slot representation learning** First, fixing  $\tau_p=100$ ,  $\tau_t=100$ , we vary the coefficient  $\alpha$  as shown in Figure 4a. While smaller  $\alpha$  values (e.g., 0.1) achieve marginally better f6 scores (indicating greater primitive stability across tasks), they significantly degrade other critical metrics, particularly f1 and f2. This trade-off suggests that excessively stable primitives may fail to adequately capture diverse label semantics necessary for effective primitive-logit alignment.

Next, we examine the temperature parameter  $\tau_p$  by fixing  $\alpha=10, \tau_t=100$  (Figure 4b). The radar chart demonstrates that  $\tau_p=100$  provides optimal balance across all metrics, confirming our hypothesis that moderate temperature settings enable better concept generalization while preventing over-regularization.

Finally, we analyze the task temperature  $\tau_t$  with fixed  $\alpha=10, \tau_p=100$  (Figure 4c). While no single  $\tau_t$  value dominates across all metrics, we observe that  $\tau_t=100$  achieves the highest f1 score. Section K provides  $w_p$  visualizations showing that larger  $\tau_t$  values produce sharper slot selection distributions for primitive construction, which benefits concept representation but may reduce flexibility in extreme cases.

**Slot attention architecture** Figure 5a examines the impact of increasing the number of slots (K). While higher K values initially improve slot performance by enabling representation of more concepts, we observe diminishing returns beyond K=10. This saturation occurs due to two factors: (1) the limited number of visually discriminable concepts per image, and (2) the finite capacity of the slot attention mechanism. Redundant slots tend to converge to similar representations, creating a performance plateau. Our slot mask visualizations in section K confirm this phenomenon, showing that excessive slots merely replicate existing patterns rather than capturing novel information.

Figure 5b investigates the effect of refinement iterations  $(N_s)$  in the slot attention module. While increasing  $N_s$  enhances slot discriminability by promoting greater inter-slot differences, we find that three iterations  $(N_s = 3)$  achieve optimal performance. Further increases do not meaningfully

improve results, suggesting that three iterations strike an effective balance between refinement and computational efficiency.

Figure 5c explores the relationship between slot dimensionality (capability) and performance. We observe that larger slot dimensions consistently improve f1 scores, indicating better concept representation. However, this comes at the cost of increased computational overhead, necessitating careful trade-off considerations for practical applications.

Finally, Figure 5d examines the impact of decoder architecture by varying MLP layer depth. Contrary to expectations, deeper decoders fail to improve extraction performance, suggesting that the current decoder architecture has sufficient capacity for the task.

**Primitive-Logit knowledge distillation** We apply our learned slot attention mechanism to compute concept-based sample-wise similarities on RanPAC, systematically evaluating key hyperparameters in our primitive-logit knowledge distillation framework.

Figure 6a demonstrates that increasing the coefficient  $\beta$  for  $L_a$  consistently improves CL accuracy (AA). This indicates that stronger self-supervision from concept-based similarity effectively enhances the model's ability to preserve task-specific knowledge while adapting to new tasks.

Figure 6b highlights the critical importance of properly tuning the temperature parameter  $\tau_a$ . We observe a performance plateau when  $\tau_a$  is within an optimal range (approximately [0.1, 1.0]). Values beyond this range exhibit clear trade-offs. This is because (1) Large  $\tau_a(>1.0)$  causes excessive emphasis on sample-wise differences, undermining concept sharing; (2) Small  $\tau_a(<0.1)$  produces overly smooth logit similarities, degrading classification performance.

Regarding normalization strategies on primitives (Equation 5 min-max *vs* Equation 3 softmax), Figure 6c shows that min-max normalization outperforms softmax normalization. This advantage stems from min-max normalization's ability to provide sharper supervision through its linear scaling properties, and maintain better sensitivity to subtle concept differences between samples.

#### I RESULTS ON OTHER BENCHMARKS

**COBJ** The results in Table 7 clearly demonstrate that incorporating CompoSLOT into CL methods with FMs leads to significant performance improvements across various metrics. Specifically, CompoSLOT enhances compositional generalization ability, as evidenced by higher Hn and improved R (most significant gain of Hn for ADAM + adapter from 57.793 to 61.581), which in turn drives better overall CL performance (AA for ADAM + adapter improves from 45.75 to 50.15). CompoSLOT's ability to strengthen compositional generalization appears to be the key factor behind these gains, enabling the model to better handle complex concepts and retain knowledge more effectively across tasks.

**ImageNet-R** It can be seen that CompoSLOT can generally improve the performance of CL methods with FMs in Table 8. The improvement is likely due to the observation that the learned slot attention can discover hidden concept sharing between images, as evidenced by the visualization analysis in section K. Rehearsal methods (e.g., FOSTER\* and MEMO\*) achieve better performance in terms of AA and CA, comparing with rehearsal-free methods. This is because rehearsal methods can access old samples, thus, CompoSLOT's primitive-logit alignment loss can provide more pairwise contrastive self-supervision on concept sharing, which enhances the model's compositional generalization performance.

#### J RESULTS ON OTHER BACKBONES

This section is to answer: **Do better vision foundation models contribute to better concept learning and continual learning performance?** We investigate the effect of model scaling via increasing the size and depth of the ViT architecture (e.g., ViT-L16 vs ViT-B16), and the effect of pretraining strategy via leveraging greater pretraining objectives, such as DINO (Oquab et al., 2024) (e.g., ViT-B16-DINO) and SAM (Kirillov et al., 2023) (e.g., ViT-B16-SAM), which have been shown to enhance semantic understanding, especially on segmentation and concept-rich tasks. We conduct

Table 7: Main result on 10-10 tasks COBJ. We report the average accuracy after training the last task (AA), the cumulative average accuracy for each task (CA), and the final forgetting (FF). For CFST, we report the Harmonic mean of compositional testing (Hn) and the ratio of Hn and reference testing (R). Methods with CompoSLOT are denoted with a postfix "†". Methods rehearse old samples are denoted with a postfix "\*". We report results over 3 trials with (mean  $\pm$  95% confidence interval).

Methods	AA (%) ↑	Continual CA (%) ↑	FF (%) ↓	CFST Hn (%) ↑	R↑
CPrompt	42.015±0.118	51.172±9.718	22.575±6.479	58.961±0.409	0.878
CPrompt †	45.520±0.421	<b>52.565</b> ± <b>0.931</b>	<b>19.575</b> ± <b>1.029</b>	<b>59.880</b> ± <b>2.032</b>	<b>0.880</b>
ADAM + adapter	45.750±0.346	52.800±6.121	12.175±1.836	57.793±1.388	0.914
ADAM + adapter †	50.150±0.249	57.767±5.461	11.050±1.802	<b>61.581</b> ± <b>1.399</b>	<b>0.938</b>
RanPAC	59.285±2.377	66.203±4.186	<b>7.450</b> ± <b>0.624</b> 7.875±0.104	60.909±3.240	0.882
RanPAC †	61.950±0.527	67.367±4.075		<b>62.317</b> ± <b>2.447</b>	<b>0.889</b>
CoFiMA	57.330±0.139	<b>64.252</b> ± <b>5.763</b>	17.375±0.035	<b>66.998</b> ± <b>2.112</b> 66.232±2.497	0.890
CoFiMA †	57.435±0.101	63.462±0.599	<b>16.650</b> ± <b>0.207</b>		<b>0.898</b>
FOSTER*	47.800±0.542	53.741±0.290	<b>10.575</b> ± <b>0.759</b>	62.750±0.337	0.852
FOSTER* †	50.980±0.225	<b>59.735</b> ± <b>0.556</b>	14.525±0.240	63.695±0.312	<b>0.908</b>
DER*	55.815±0.714	64.905±3.342	<b>23.650</b> ± <b>2.425</b>	68.558±0.189	0.844
DER* †	56.813±1.808	66.393±3.904	25.800±4.534	68.586±0.441	<b>0.872</b>

Table 8: Main result on 20-20 tasks ImageNet-R. We report the average accuracy after training the last task (AA), the cumulative average accuracy for each task (CA), and the final forgetting (FF). Methods with CompoSLOT are denoted with a postfix " $\dagger$ ". Methods rehearse old samples are denoted with a postfix "\*". The data for methods with citations is reported from the original paper. We report results over 3 trials with (mean  $\pm$  95% confidence interval).

Methods	AA (%) ↑	CA (%) ↑	FF (%) ↓
CPrompt (Gao et al., 2024)	74.790±0.280	<b>81.460±0.930</b> 79.964±1.078	7.340±0.650
CPrompt †	75.225±0.270		<b>6.989</b> ± <b>1.126</b>
RanPAC	78.375±0.062	82.519±0.839	<b>4.856</b> ± <b>0.367</b> 5.294±0.039
RanPAC †	78.550±0.346	<b>82.900</b> ± <b>0.747</b>	
CoFiMA	80.025±0.146	83.927±1.421	7.614±0.142
CoFiMA †	80.250±0.016	<b>84.118</b> ± <b>1.017</b>	<b>7.022</b> ± <b>0.005</b>
FOSTER*	76.001±0.243	80.974±1.083	<b>2.259</b> ± <b>0.526</b>
FOSTER* †	78.950±0.201	<b>82.392</b> ± <b>1.308</b>	2.608±0.720
MEMO*	64.200±1.109	72.118±0.074	<b>4.967</b> ± <b>0.074</b>
MEMO* †	65.200±0.249	<b>72.995</b> ± <b>1.251</b>	5.344±0.256

experiments along the two key dimensions above and report the results in Table 9. The results show that ViT-L16 with larger model sizes demonstrates stronger representation modeling capabilities compared to ViT-B16, thus further boosting the significance of our CompoSLOT. ViT-B16-DINO and ViT-B16-SAM with greater pre-training objectives exhibit better compositionality in decomposing concepts and continual learning performance, as reflected by higher Hn values.

#### K VISUALIZATION

This section investigates how the CompoSLOT framework enhances continual learning performance by first demonstrating through slot attention mask visualizations across various benchmarks that CompoSLOT successfully identifies important concepts (primitives) in an unsupervised manner, and

Table 9: Varying backbone on 10-10 tasks CGQA. We report the average accuracy after training the last task (AA), the cumulative average accuracy for each task (CA), and the final forgetting (FF). The candidate CL algorithm is RanPAC. Methods with CompoSLOT are denoted with a postfix "†"

Backbone	AA (%) ↑	CA (%) ↑	FF (%) ↓	Hn↑
ViT-B16	65.81	75.50	10.51	78.86
ViT-B16 †	66.75	76.58	10.21	79.81
ViT-B16-DINO †	66.58	76.62	10.24	80.39
ViT-B16-SAM †	67.30	77.76	9.67	81.22
ViT-L16 †	67.11	77.54	9.85	80.82

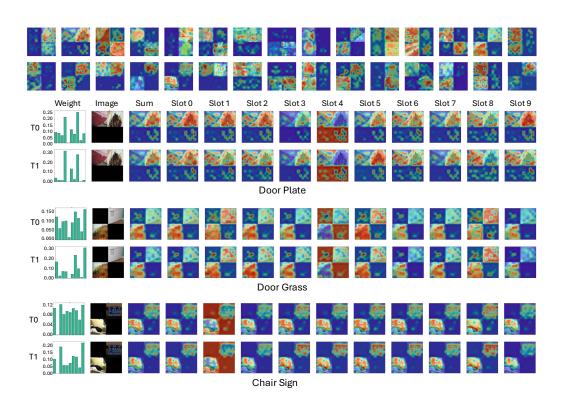


Figure 7: Visualization of learned slots on 30 randomly sampled images (3 images for each class in the first task of the 10-10 tasks) on CGQA. **Top row**: Primitives (weighted-sum of slot masks weighted by  $w_p$ ) for 30 images. **Bottom 3 rows**: Three examples of images from classes (Door Plate), (Door Grass), and (Chair Sign) after being trained on the first task (T0) and on the second task (T1). **From left to right**:  $w_p$ , origin image, primitive (weighted-sum of slot masks), and 10 slot masks, respectively. **Takeaway**: CompoSLOT successfully extracts primitives without any concept label.

then by presenting similarity matrix visualizations of ground truth concepts/primitives/features/logits for specific algorithms to illustrate the regularization effects that improve model compositional generalization and stability during continual learning. We attribute this robustness to "concept rehearsal": although class labels change, many visual concepts are shared and recur across tasks, helping stabilize the primitive selection weights. This is also discussed in Lai et al. (2024).

**Concept learning** We evaluate CompoSLOT on CGQA, COBJ, ImageNet-R, and CIFAR100 benchmarks by randomly selecting three representative images from each class. The extracted slot masks are visualized in Figures 7, 8, 9, and 10, respectively.

On CGQA, the weighted slot masks (using weights  $w_p$ ) effectively localize class-relevant concepts in each image. For instance, in the *Door Plate* class, slot 7 consistently captures the *Plate* concept while

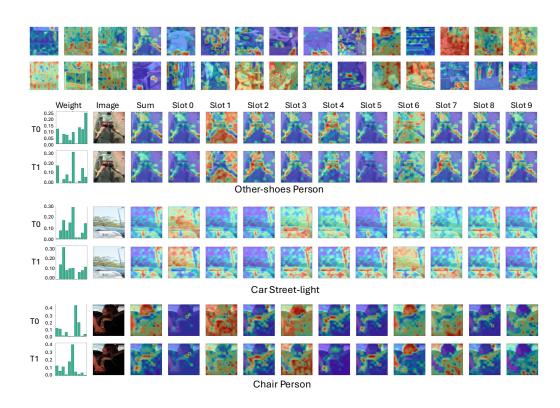


Figure 8: Visualization of learned slots on 30 randomly sampled images (3 images for each class in the first task of the 10-10 tasks) on COBJ. **Top row**: Primitives (weighted-sum of slot masks weighted by  $w_p$ ) for 30 images. **Bottom 3 rows**: Three examples of images from classes (Door Plate), (Door Grass), and (Chair Sign) after being trained on the first task (T0) and on the second task (T1). **From left to right**:  $w_p$ , origin image, primitive (weighted-sum of slot masks), and 10 slot masks, respectively. **Takeaway**: CompoSLOT successfully extracts primitives without any concept label.

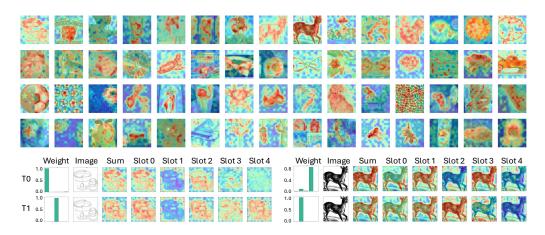


Figure 9: Visualization of learned slots on 60 randomly sampled images (3 images for each class in the first task of the 20-20 tasks) on ImageNet-R. **Top row**: Primitives (weighted-sum of slot masks weighted by  $w_p$ ) for 60 images. **Bottom row**: Two examples of images after being trained on the first task (T0) and on the second task (T1). **From left to right**:  $w_p$ , origin image, primitive (weighted-sum of slot masks), and 5 slot masks, respectively. **Takeaway**: CompoSLOT successfully extracts primitives without any concept label, and the concept sharing is rare between classes.

slot 8 focuses on the *Door*, demonstrating precise concept disentanglement. Notably, the learned

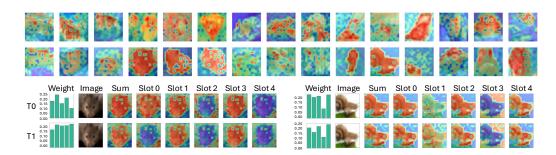


Figure 10: Visualization of learned slots on 30 randomly sampled images (3 images for each class in the first task of the 10-10 tasks) on CIFAR100. **Top row**: Primitives (weighted-sum of slot masks weighted by  $w_p$ ) for 30 images. **Bottom row**: Two examples of images after being trained on the first task (T0) and on the second task (T1). **From left to right**:  $w_p$ , origin image, primitive (weighted-sum of slot masks), and 5 slot masks, respectively. **Takeaway**: CompoSLOT successfully extracts primitives without any concept label, and the concept sharing is rare between classes.

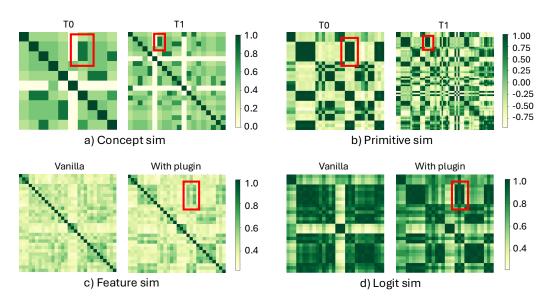


Figure 11: Visualization of a) concept; b) primitive; c) feature; d) logit cosine similarity matrices on sampled images (three images for each class in the first task T0 and second task T1 of the 10-10 tasks) on COBJ. a) **Left**: Multi-hot concept cosine similarity matrix of 30 images for T0; **right**: Multi-hot concept cosine similarity of 60 images (from the first-2 tasks T0 and T1). b) The primitive cosine similarity of the corresponding images. We use the learned pair-wise primitive similarity to mimic the statistics of the pair-wise concept similarity and regularize logits. c) **Left**: The learned feature cosine similarity matrix of 30 images in T0 for ADAM + adapter; **right**: The learned feature cosine similarity matrix of 30 images in T0 for ADAM + adater  $\dagger$ . d) The logit cosine similarity of the corresponding images as in c). **Takeaway**: The learned primitive successfully mimics concept statistics without concept supervision, and our  $L_a$  successfully distills pair-wise primitive similarity into logits and affects the feature representations (as demonstrated with the regions marked with red box), while ADAM + adater does not capture this concept sharing statistic.

primitives maintain visual consistency across tasks, that the primitive representation after task T0 closely resembles that after T1, confirming the stability of CompoSLOT. This phenomenon was similarly observed in Figure 1 of the main paper's introduction.

The more challenging COBJ benchmark presents similar results. For an image in the *Other-shoes Person* class, slot 5 accurately identifies the *Other-shoes* concept while slot 7 correctly localizes the *Person*, even in this complex compositional setting.

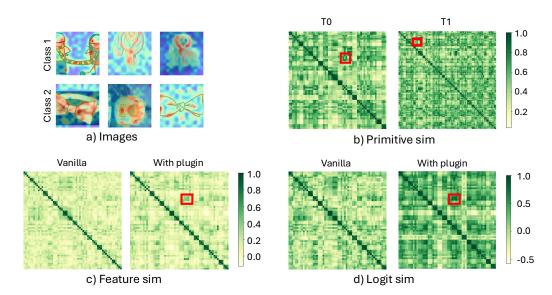


Figure 12: Visualization of a) images related to red box; b) primitive; c) feature; d) logit cosine similarity matrices on sampled images (three images for each class in the first task T0 and second task T1 of the 20-20 tasks) on ImageNet-R. a) Six images from two classes in T0 which corresponding to the red box. b) The primitive cosine similarity of the corresponding images. c) **Left**: The learned feature cosine similarity matrix of 60 images in T0 for FOSTER; right: The learned feature cosine similarity matrix of 60 images in T0 for FOSTER †. d) The logit cosine similarity of the corresponding images as in c). **Takeaway**: The learned primitives show that CompoSLOT discovers hidden relationships based on concept as demonstrated with the regions marked with red box, while FOSTER does not capture this concept sharing statistic.

When examining ImageNet-R and CIFAR100 with K=5 slots, we observe that the primary concept corresponding to each class label is reliably identified, and the representations maintain discriminative power while preserving semantic consistency. However, the concept sharing is visually rare between classes, as demonstrated by the distinct slot activation patterns for different classes.

**Primitive-logit alignment** We conduct in-depth visualization analysis to understand the performance improvement of CompoSLOT on COBJ, using ADAM + adapter as a representative example. We visualize 30 images for T0 and 60 images for T1 (10 old classes and 10 new classes). Figure 11 presents the cosine similarity matrix visualizations including: (a) Ground truth multi-hot concepts; (b) Extracted primitives; (c) Feature representations; (d) Final logits. The red boxes highlight two pairs of classes with concept sharing: (*Other-shoes Person*) and (*Other-shoes Person Sneaker*), as well as (*Person Sneaker*) and (*Other-shoes Person Sneaker*). CompoSLOT successfully captures these shared concepts in the primitive representations (Figure 11b) and effectively distills them into the final logits (Figure 11d). Notably, this alignment process also induces regularization at the feature level, as evidenced by the more coherent feature representations shown in Figure 11c.

We further validate CompoSLOT on ImageNet-R, a standard CL benchmark without ground truth concept labels. Figure 12 shows the case performing CompoSLOT on FOSTER. Our slot attention mechanism identifies shared concepts across six images (highlighted in red boxes), particularly revealing a consistent "Fabric" concept (Figure 12a). This automatic discovery of hidden relationships demonstrates CompoSLOT's ability to generalize concept learning across different benchmarks.

The consistent performance improvements reported in Sections 5 and I validate that CompoSLOT effectively captures meaningful semantic relationships, leading to better generalization and compositional learning capabilities.

Table 10: Additional ablation results on CGQA.

1568	3
1569	9
1570	)
1571	1
1572	2

1609
1610
1611
1612
1613
1614
1615

Methods	$ L_p $	$L_a$	AA (%) ↑	FF (%) ↓
SimpleCIL SimpleCIL	X	X	<b>36.16</b> 24.71	<b>13.9</b> 22.93
RanPAC RanPAC	X ✓	X	<b>65.81</b> 41.59	<b>10.51</b> 11.87

With plugin

Vanilla

With plugin

Vanilla

With plugin

Vanilla

With plugin

Methods

Figure 13: Visualization of the parameter numbers for methods with and without the slot module. Note that the data are collected according to the default implementation in the PILOT (Sun et al., 2025) platform and after the training of the last 10-way CGQA continual task. **Takeaway**: CompoSLOT requires a ViT backbone that is already in any model-based continual learner with a foundation model, thus, it is light-weight and free to be applied.

# L ADDITIONAL ABLATION STUDIES

To clearly substantiate the contribution of slot attention in combination with primitive selection, we conduct an ablation study where we remove knowledge distillation and instead directly use the learned primitive representations with a cosine similarity classifier for continual tasks, as in SimpleCIL (Zhou et al., 2025). We also integrate this strategy into RanPAC and the results are shown in Table 10. This naive approach suffers from severe forgetting, confirming that primitive representations are insufficient for long-term retention when learning new tasks. In contrast, our alignment loss distills pair-wise relational information (i.e., a compact, low-dimensional encoding of concept combinations) rather than high-dimensional raw representations. This enables methods equipped with CompoSLOT to maintain stable performance while accumulating higher accuracies over time, demonstrating the efficacy of CompoSLOT in mitigating catastrophic forgetting.

# M ALGORITHM EFFICIENCY ANALYSIS

**Parameter overhead** We evaluate the parameter overhead introduced by our slot attention module. As this module requires a pretrained ViT as its semantic feature extractor, which is a standard component in all continual learning of foundation models frameworks, the additional trainable parameters are negligible compared to the total model size, as illustrated in Figure 13. This makes our CompoSLOT computationally efficient while delivering significant performance benefits.

**Computation overhead** In Table 11, we study the computational overhead introduced by the slot attention mechanism and primitive extraction. As an example, we choose FOSTER as a representative baseline, since it achieves nearly top performance among others. We compare three cases: 1) Continual training of just our slot module plugin, including both slot attention and primitive selection

Table 11: Computational overhead (h) on CGQA.

Slot module	FOSTER	FOSTER †
5.5	9.1	10.1

components, without applying it to other continual learning algorithms; 2) Full continual training of FOSTER; 3) Full continual learning of FOSTER with a pretrained slot module plugin (FOSTER  $\dagger$ ). We highlight that the slot module can be learned offline as a reusable component which only associated with the benchmark and is independent of algorithms. Once trained, it serves as a pretrained plugin that can be directly loaded for any continual learning algorithm with minimal additional overhead. It only requires adding alignment loss  $L_a$  for logit regularization and spending an additional 10% of total training time for FOSTER from 9.1h to 10.1h. This design is particularly beneficial when running multiple continual learning algorithms on the same data distribution.

Importantly, we conduct an ablation study (Section 5), where we deliberately increase the parameter count of baseline CL methods to match our CompoSLOT-enhanced models. The results demonstrate that the performance gains stem not from increased capacity, but from CompoSLOT's improved compositional generalization capabilities. This confirms that CompoSLOT provides genuine algorithmic advantages rather than simply benefiting from more parameters.

## N USE OF LARGE LANGUAGE MODELS

In the process of preparing this paper, we employed LLMs to polish the writing of the paper. The assistance provided by LLMs was mainly focused on improving the clarity, coherence, and overall quality of the language used in the manuscript. We input sections of the paper into the LLM and requested it to suggest rephrasings, correct grammar and spelling errors, and enhance the readability of the text. It is important to note that LLMs did not play a significant role in the research ideation. The core ideas, research questions, experimental designs, and methodological choices were independently conceived and developed by the human authors. The LLM was not involved in formulating the hypotheses, determining the research direction, or making decisions regarding the data collection and analysis methods.