

---

# PROMPT-BASED CONTINUAL COMPOSITIONAL ZERO-SHOT LEARNING

---

A PREPRINT

**Sauda Maryam, Sara Nadeem, Mohsen Ali**  
Intelligent Machines Lab  
Information Technology University  
{msds22025, phdcs21001, mohsen.ali}@itu.edu.pk

**Faisal Z. Qureshi**  
Visual Computing Lab  
Ontario Tech University  
faisal.qureshi@ontariotechu.ca

December 17, 2025

## ABSTRACT

We tackle continual adaptation of vision–language models to new attributes, objects, and their compositions in Compositional Zero-Shot Learning (CZSL), while preventing forgetting prior knowledge. Unlike classical continual learning where classes are disjoint, CCZSL is more complex as attributes and objects may reoccur across sessions while compositions remain unique. Built on a frozen VLM backbone, we propose the first Prompt-based Continual Compositional Zero-Shot Learning (PromptC-CZSL) framework that retains prior knowledge through recency-weighted multi-teacher distillation. It employs session-aware compositional prompts to fuse multimodal features for new compositions, while attribute and object prompts are learned through session-agnostic fusion to maintain global semantic consistency, which is further stabilized by a Cosine Anchor Loss (CAL) to preserve prior knowledge. To enhance adaptation in the current session, an Orthogonal Projection Loss (OPL) ensures that new attribute and object embeddings remain distinct from previous ones, preventing overlap, while an Intra-Session Diversity Loss (IDL) promotes variation among current-session embeddings for richer, more discriminative representations. We also introduce a comprehensive protocol that jointly measures catastrophic forgetting and compositional generalization. Extensive experiments on UT-Zappos and C-GQA benchmarks demonstrate that PromptCCZSL achieves substantial improvements over prior VLM-based and non-VLM baselines, setting a new benchmark for CCZSL in closed-world setting.

**Keywords** CCZSL · Continual Learning · Compositional Learning · Zero-Shot Continual Compositional Learning

## 1 Introduction

For reliable scene understanding, models must recognize not only *what* an object is but also *how* it appears (e.g., *wet cat*, *broken glass*), and continuously adapt to new compositions without retraining from scratch. Conventional Compositional Zero-Shot Learning (CZSL) frameworks tackle this problem by decomposing visual concepts into attributes and objects, and recombining seen primitives to recognize unseen compositions. However, CZSL assumes a fixed vocabulary of primitives known during training, limiting its applicability in real-world settings where new objects or attributes appear incrementally. Extending CZSL to a continual compositional zero-shot learning (CCZSL) scenario introduces additional complexity: unlike standard continual learning where classes are disjoint, CCZSL must preserve shared primitives (attributes or objects) across sessions while learning novel compositions—making it particularly prone to *catastrophic forgetting* as representations drift when new attributes, objects, and their compositions are introduced (Figure 1). Zhang *et al.* Zhang et al. [2024a] recently proposed the first CCZSL framework by extending standard CZSL to a continual-learning setup through learnable object and attribute embeddings and session-specific super-primitives. While effective in modeling contextual variations across sessions, these embeddings remain confined to the limited semantics of the training data, hindering their ability to adapt to newly introduced primitives.

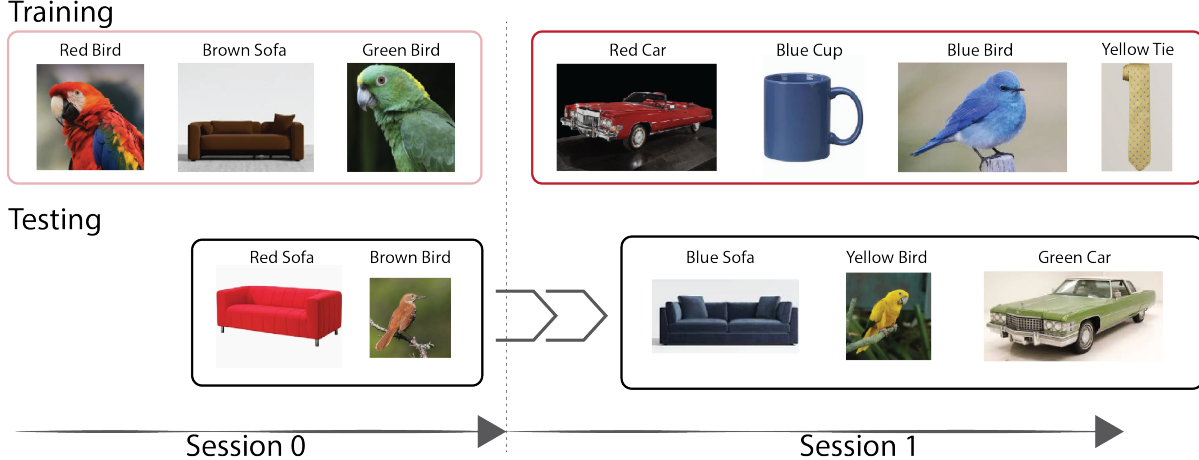


Figure 1: Illustration of the Continual Composition Zero-Shot Learning (CC-ZSL) setup. In Session 1, the model learns compositions such as Red Bird, Brown Sofa, and Green Bird, and is tasked to generalize to unseen combinations like Red Sofa and Brown Bird. In Session 2, without access to prior data, the model learns new compositions (Red Car, Blue Cup, Blue Bird, Yellow Tie) that introduce new attributes (Blue, Yellow) and objects (Car, Cup, Tie). A continual learning approach enables the model to retain earlier concepts (Sofa, Brown, Green) while integrating new ones, allowing correct recognition of both new compositions (Blue Sofa, Yellow Bird, Green Car) and previous unseen ones (Red Sofa, Brown Bird). Images are sourced via Google Search and are used for educational/non-commercial purposes.

Recent advances in vision–language models (VLMs) such as CLIP Radford et al. [2021] have demonstrated strong generalization across diverse visual concepts, enabling compositional reasoning through large-scale contrastive pre-training. Building upon this, several works Lu et al. [2023a], Nayak et al. [2023a], Bao et al. [2024a] have explored fine-tuning CLIP for compositional understanding, employing decomposed embeddings and multimodal prompting mechanisms. However, these approaches are designed under a static compositional setting, where both the attribute and object vocabularies remain fixed. Consequently, they fail to accommodate continual introduction of new primitives, resulting in degraded semantic alignment and representation drift when applied to evolving environments.

To address these challenges, we bridge vision–language pretraining and continual compositional learning, proposing a scalable framework that adapts to new attribute–object compositions while retaining prior knowledge. We introduce Prompt-based Continual Compositional Zero-Shot Learning (Prompt CCZSL), the first framework enabling continual compositional learning within VLMs. Prompt CCZSL leverages a frozen VLM backbone and a shared soft-prompt bank containing learnable embeddings for the attribute and object vocabularies across sessions. Each session introduces new attributes and objects, and the model learns new primitives while preserving the geometric and semantic consistency of the embedding space to maintain attributes, objects, and their compositions from prior sessions. To further mitigate forgetting, we introduce a multi-teacher knowledge distillation strategy that aggregates soft logits from all prior session models weighted by session recency, allowing the model to retain historical knowledge while adapting to new primitives. Additionally, a Cosine Anchor Alignment Loss (CAL) enforces directional consistency between attribute and object embeddings across sessions, providing semantic anchoring for continual updates. To enhance representation quality, we employ orthogonality regularization to ensure separability between session embeddings and intra-session diversification to enrich local representation diversity jointly maintaining a disentangled, compositional embedding space across time. We also propose a comprehensive CCZSL evaluation protocol that quantifies both catastrophic forgetting and compositional generalization using the current session’s zero-shot test set as well as the union of all previous sessions up to the current one. Extensive experiments on UT-Zappos Yu and Grauman [2014] and C-GQA Hudson and Manning [2019] datasets demonstrate that Prompt CCZSL achieves state-of-the-art performance, significantly improving both continual adaptation and compositional generalization compared to existing prompt-based and non-prompt-based baselines. Our main contributions are summarized as follows:

- We present the first Prompt-based Continual Compositional Zero-Shot Learning (Prompt CCZSL) framework, integrating continual learning principles into vision–language models to enable incremental learning of new attribute–object primitives.
- We propose a session-aware compositional multi-modal fusion that stabilizes updates across sessions, allowing the model to adapt to new attributes, objects, and their compositions while preserving the structure of previously learned embeddings.

- We introduce a continual adaptation strategy as a multi-teacher knowledge distillation mechanism that aggregates knowledge from prior sessions with recency weighting, and a Cosine Anchor Alignment Loss that enforces semantic consistency across sessions.
- We employ orthogonality and intra-session diversification to maintain representation separability and promote enriched session-specific semantics in the prompt space.
- We design a comprehensive CCZSL evaluation protocol to measure both compositional generalization and catastrophic forgetting.
- Our modules are plug-and-play compatible with state-of-the-art CZSL methods under continual settings, achieving significant gains on strong CCZSL baselines.
- We introduce two setting of CCZSL, namely *Constrained-CCZSL* where unseen composition in initial session will remain unseen in rest, and *Realistic-CCZSL* a more realistic scenario where composition unseen in one of the session might become seen in next. We proposed results under Constrained-CCZSL.

## 2 Related Work

### 2.1 Compositional Learning

Compositional Zero-Shot Learning (CZSL) aims to recognize unseen attribute–object compositions by leveraging knowledge from seen pairs. Early work learned joint embeddings for images and compositions Misra et al. [2017], modeled attributes as transformations applied to object representations Nagarajan and Grauman [2018], or imposed algebraic constraints (e.g., symmetry/group structures) to stabilize composition mappings Li et al. [2020a]. Later methods improved disentanglement via invariant or contrastive objectives Zhang et al. [2022a], Li et al. [2022a], and recent visual-encoder models enhanced object-conditioned attribute reasoning through attention-based disentanglers and conditional attribute modeling Hao et al. [2023], Wang et al. [2023]. **VLM/CLIP-based CZSL.** With the rise of CLIP-style vision–language models (VLMs), CZSL increasingly adopts parameter-efficient adaptation by freezing backbones and tuning lightweight text-side modules. CSP Nayak et al. [2023b] learns attribute/object tokens within prompts, while DFSP Lu et al. [2023b] introduces decomposed cross-modal fusion for tighter image–text coupling. Troika Huang et al. [2024] employs multi-path prompts with a cross-modal traction step, PLID Bao et al. [2024b] leverages LLM-generated language priors for better generalization, and hierarchical prompt learning Wang and Deng [2023] reduces semantic gaps. Despite strong compositional reasoning, these VLM-based approaches remain static and unable to accommodate continual updates to the attribute–object vocabulary.

### 2.2 Continual Learning

Continual learning (CL) addresses sequential task acquisition while mitigating catastrophic forgetting. Classical CL methods fall into three families: (i) regularization-based approaches constrain updates to parameters important for past tasks (EWC, SI) Kirkpatrick et al. [2017], Zenke et al. [2017]; (ii) replay-based strategies rehearse stored or synthesized exemplars Rebuffi et al. [2017], Castro et al. [2018]; and (iii) parameter-isolation methods allocate task-specific modules or masks to reduce interference Mallya and Lazebnik [2018], Serra et al. [2018]. Knowledge distillation (KD) has become central in deep CL, aligning logits between current and frozen models to retain prior knowledge (LwF) Li and Hoiem [2016]. Extensions combine logit and feature-level alignment for stronger class-incremental stability Castro et al. [2018], Dhar et al. [2019]. **Feature-based distillation.**

Beyond logit alignment, feature-based distillation preserves the teacher’s internal representations Heo et al. [2019], Beyer et al. [2022]. Relational KD Park et al. [2019] maintains pairwise sample geometry, while  $L_2$  and mean-squared error (MSE) directly align hidden features Beyer et al. [2022], Wang et al. [2021], Chen et al. [2022]. Recent studies highlight directional alignment—matching feature orientations while allowing magnitude flexibility Wang et al. [2024].

**Transformer-based distillation.** In Vision Transformers (ViTs), *where* and *how* to distill matters. DeiT Touvron et al. [2021] introduces a distillation token that interacts with teacher features via attention. ViTKD Yang et al. [2024] observes that shallow layers encode transferable low-level structure, while deeper layers specialize to the current task.

**VLMs with parameter-efficient tuning (PEFT).** For large VLMs such as CLIP, freezing pretrained backbones while tuning lightweight adapters, LoRA modules, or soft prompts preserves zero-shot generalization and minimizes forgetting. LADA Luo et al. [2025] attaches label-specific adapters to the frozen image encoder and uses feature distillation to protect prior semantics, yielding strong continual performance without gradient flow through the backbone.

**Prompt-based continual learning.** Prompt-based CL extends PEFT to continual scenarios. L2P Wang et al. [2022a] retrieves task-relevant prompts from a learnable pool; DualPrompt Wang et al. [2022b] separates shared and expert prompts; CODA-Prompt Smith et al. [2023] dynamically composes decomposed components; and CPrompt Gao et al.

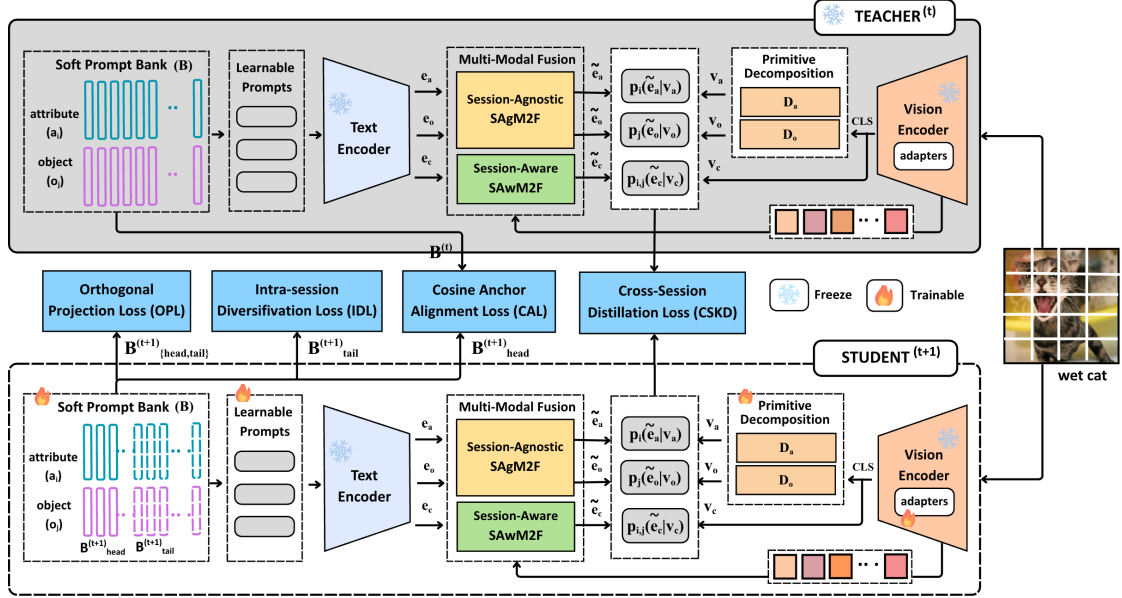


Figure 2: Overview of the proposed Prompt-based Continual Compositional Zero-Shot Learning (PCCZSL) framework. At each training stage  $S^{(t)}$ , the frozen teacher model  $S^{(t-1)}$  transfers prior knowledge via Cross-Session Knowledge Distillation (CSKD), preserving attribute–object relationships from earlier sessions. The Cosine Anchor Loss (CAL) maintains directional alignment of attribute and object prompts across sessions, stabilizing the shared semantic space. Session-Aware Multi-Modal Fusion (SAwM2F) refines session specific prompts through cross-attention with visual features, while the Session-Agnostic (SAgM2F) branch updates all attribute and object prompts for global consistency. Meanwhile, the Orthogonal Projection Loss (OPL) enforces separation between old and new primitives, and the Intra-Session Diversification Loss (IDL) enhances variability among new prompts. Together, these components form a unified continual learning pipeline that preserves past knowledge, integrates new semantics, and mitigates catastrophic forgetting.

[2024] enforces prompt–classifier consistency. As larger ViTs improve plasticity, continual distillation (CDL) methods such as KDP Zhang et al. [2024b] insert learnable distillation prompts into frozen backbones, outperforming standard KD. PromptKD Li et al. [2024] distills precomputed teacher text features via KL divergence on unlabeled data, and dual-teacher setups Zheng et al. [2024] enable joint learning from a frozen foundation (e.g., CLIP) and a task-tuned teacher.

### 2.3 Continual Compositional Zero-Shot Learning

Recently Zhang *et al.* Zhang et al. [2024a] introduces the continual compositional zero-shot learning (CCZSL) framework, incrementally adding attributes and objects across sessions and using dual knowledge distillation to retain prior compositions.

## 3 Method

### 3.1 Problem Setup

We begin by formalizing CZSL. Let  $\mathcal{A}$  denotes the set of attributes and  $\mathcal{O}$  denotes the set of objects then  $\mathcal{C} = \mathcal{A} \times \mathcal{O} = \{(a, o) \mid a \in \mathcal{A}, o \in \mathcal{O}\}$  denotes the set of all possible compositions. Let  $\mathcal{C}_{\text{seen}}, \mathcal{C}_{\text{unseen}} \subset \mathcal{C}$  be two disjoint subsets of  $\mathcal{C}$  consisting of *seen* and *un-seen* compositions,  $\mathcal{C}_{\text{seen}} \cap \mathcal{C}_{\text{unseen}} = \phi$ . Training set consists of images labeled with  $\mathcal{C}_{\text{seen}}$  compositions;  $\mathcal{T}_s = \{(x, c) \mid x \in X, c \in \mathcal{C}_{\text{seen}}\}$ , testing set consists of images from both  $\mathcal{C}_{\text{seen}}$  and  $\mathcal{C}_{\text{unseen}}$  compositions. We follow the *closed-world* setup Zhang et al. [2024a] where testing label space is restricted to only  $\mathcal{C}_{\text{test}} = \mathcal{C}_{\text{seen}} \cup \mathcal{C}_{\text{unseen}}$  and not have to consider all possible compositions (thus avoiding “invalid” compositions).

**Continual-CZSL Setup:** In Continual-CZSL (CCZSL) the attributes, objects, and their composition are partitioned into  $T + 1$  sessions:  $\mathcal{S}^{(s)} = (\mathcal{A}^{(s)}, \mathcal{O}^{(s)}, \mathcal{C}^{(s)})$ , where  $s \in [0, T]$  and  $\mathcal{C}^{(s)} = \mathcal{A}^{(s)} \times \mathcal{O}^{(s)}$ . We impose constraint that  $\cup_{s=0}^T \mathcal{A}^{(s)} = \mathcal{A}$ ,  $\cup_{s=0}^T \mathcal{O}^{(s)} = \mathcal{O}$ , and  $\cup_{s=0}^T \mathcal{C}^{(s)} = \mathcal{C}$ .

For session  $t$ , the model is trained on  $\mathcal{A}^{(t)}$ ,  $\mathcal{O}^{(t)}$ , and  $\mathcal{C}_{\text{seen}}^{(t)}$ . During training at session  $t$ , it may access models trained on earlier sessions  $0 \leq s < t$ . Evaluation is performed on  $\mathcal{C}_{\text{test}}^{(t)} = \bigcup_{0 \leq s \leq t} \{\mathcal{C}_{\text{seen}}^{(s)} \cup \mathcal{C}_{\text{unseen}}^{(s)}\}$ , capturing performance on concepts introduced at session  $t$  as well as retention of prior knowledge.

The CCZSL comes with two settings, in one such that  $\forall_{i \neq j} \mathcal{C}^{(i)} \cap \mathcal{C}^{(j)} = \phi$  basically saying composition that is unseen in initial session will remains unseen in rest. We call this, *Constrained-CCZSL*. A more realistic scenario is when composition unseen in one of the session might become seen in next, we call this *Realistic-CCZSL*, here only  $\forall_{i \neq j} \mathcal{C}_{\text{seen}}^{(i)} \cap \mathcal{C}_{\text{seen}}^{(j)} = \phi$ . In our current work we are exploring only Constrained-CCZSL. In what follows we drop session superscript when unambiguous.

### 3.2 PromptCCZSL Framework

Our proposed PromptCCZSL follows Huang et al. [2024] to use VLM to capture compositional information, by learning prompts for primitives that could be anchored across sessions. Figure 2 illustrates the proposed **PromptCCZSL** framework.

#### 3.2.1 Preliminaries:

Following Huang et al. [2024], Liu et al. [2024], our model is built on pre-trained CLIP backbones and augmented with learnable prompt embeddings for compositional reasoning.

**Primitive Concepts Learning:** Representations for the primitives and their compositions are obtained by performing soft prompts tuning using the frozen pretrained-CLIP’s text-encoder. We maintain a session-shared *Soft Prompt Bank* containing one embedding per primitive concept. To initialize a concept prompt, we tokenize its text with CLIP, discard special tokens e.g., [SOS] / [EOS], and use the mean token embedding as the soft prompt vector, which is then fine-tuned. Collectively, these vectors create dense prompt bank for attributes and objects:  $\mathbf{B}_a \in \mathbb{R}^{|\mathcal{A}| \times d}$  and  $\mathbf{B}_o \in \mathbb{R}^{|\mathcal{O}| \times d}$ . In order to learn the context associated with each primitive, we concatenate them with the prefixes  $[b_1^a, \dots, b_m^a]$ ,  $[b_1^o, \dots, b_m^o]$  and  $[b_1^c, \dots, b_m^c]$ , where  $m$  donates context length. Then fully learnable prompts are  $\hat{b}_a^i = [b_1^a, \dots, b_m^a, \mathbf{B}_a[i]]$ ,  $\hat{b}_o^j = [b_1^o, \dots, b_m^o, \mathbf{B}_o[j]]$  and  $\hat{b}_c^{i,j} = [b_1^c, \dots, b_m^c, \mathbf{B}_a[i], \mathbf{B}_o[j]]$ .

The prompt vectors are mapped to textual embeddings by pre-trained, frozen, text encoder  $f_{\text{text}}$ :

$$\{e_a^i, e_o^j, e_c^{i,j}\} = f_{\text{text}}(\hat{b}_a^i, \hat{b}_o^j, \hat{b}_c^{i,j}). \quad (1)$$

**Image Representation Learning:** In parallel, the vision encoder (pretrained from CLIP) uses adapter-enhanced layers to extract image features, enabling efficient fine-tuning without perturbing pretrained representations.

Given an image  $\mathbf{x}$ , the encoder outputs visual patch tokens and a global [CLS] token:  $\mathbf{V} = [v_1, \dots, v_L] \in \mathbb{R}^{L \times d}$  and  $v_{\text{cls}} \in \mathbb{R}^d$ . A *Primitive Decomposition Module* factorizes these features into attribute, object, and composition embeddings. Let  $\text{Linear}(\cdot)$  denote a token-wise affine map  $\mathbb{R}^d \rightarrow \mathbb{R}^d$  applied to each row of  $\mathbf{V}$ . We compute  $v_* = D_*(\mathbf{V})$ , where  $*$   $\in \{a, o\}$  and  $v_c = v_{\text{cls}}$ . Each  $D_*$  is a feed-forward head consisting of a fully connected layer, batch normalization, ReLU, dropout, and a token-pooling step to produce a single vector in  $\mathbb{R}^d$ . The composition embedding  $v_c$  is taken as the adapter-tuned [CLS] representation, which captures joint attribute–object semantics. Next visual information is used to construct context-aware (fused) textual embeddings as follows.

**Session-agnostic Fusion Module (SAGM2F)** refines attributes and objects textual embeddings by computing cross-attention (CA) with visual patches tokens ( $\mathbf{V}$ ):  $\tilde{e}_a^i = \text{MLP}(\text{CA}(e_a^i W_q, \mathbf{V} W_k, \mathbf{V} W_v))$  and  $\tilde{e}_o^j = \text{MLP}(\text{CA}(e_o^j W_q, \mathbf{V} W_k, \mathbf{V} W_v))$ , where  $W_q$ ,  $W_k$  and  $W_v$  are learnable matrices converting inputs to query, key and value respectively. Since primitives might occur across sessions, session-agnostic module SAGM2F refines embeddings for all attributes and objects encountered up to and including the current session.

**Session-aware Fusion Module (SAWM2F):** Unlike attributes and objects, in our setting compositions are not repeating across the sessions. Therefore, learnable composition prompts are partitioned as follows:  $[e_{c,\text{head}} \parallel e_{c,\text{tail}}] = e_c$ , where  $e_{c,\text{head}}$  includes compositions up to but not including the current session and  $e_{c,\text{tail}}$  lists new compositions for the current session. Here only  $e_{c,\text{tail}}$  are refined by computing cross-attention with image patches:  $\hat{e}_{c,\text{tail}} = \text{CA}(e_{c,\text{tail}} W_q, \mathbf{V} W_k, \mathbf{V} W_v)$ .  $\hat{e}_{c,\text{tail}}$  is passed through an MLP to get context-aware compositions  $\tilde{e}_{c,\text{tail}}$ . Meanwhile  $e_{c,\text{head}}$  are passed through MLP as well. By refining only the tail component while keeping the head fixed, SAWM2F enables the model to learn new attribute–object compositions adaptively without overwriting previously learned semantics, ensuring stable continual learning.

Finally, three classification branches predict the conditional probabilities as  $P(\tilde{e}_a|v_a)$ ,  $P(\tilde{e}_o|v_o)$ , and  $P(\tilde{e}_c|v_c)$ . Given visual embeddings  $(v_a, v_o, v_c)$  and fused textual embeddings  $(\tilde{e}_a^i, \tilde{e}_o^j, \tilde{e}_c^{i,j})$ , alignment is modeled as the probability of

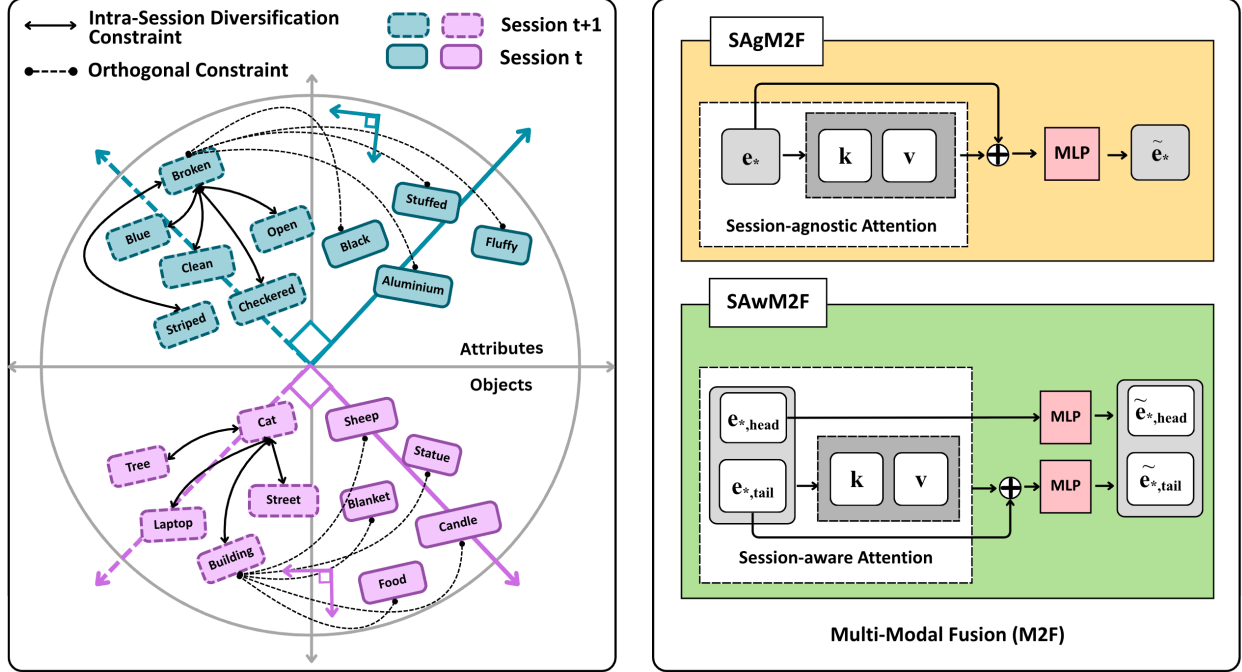


Figure 3: Overview of IDL w OPL and SAwM2F module.

matching each visual factor with its textual counterpart:  $p_i(\tilde{e}_a|v_a)$ ,  $p_j(\tilde{e}_o|v_o)$ , and  $p_{ij}(\tilde{e}_c|v_c)$ , with logits computed as temperature-scaled cosine similarities:  $z = \exp(\beta) \langle v, \tilde{e} \rangle$ , where  $\beta$  is the text encoder’s `logit_scale`. This is applied to each branch:  $z_a = z(v_a, \tilde{e}_a)$ ,  $z_o = z(v_o, \tilde{e}_o)$ , and  $z_c = z(v_c, \tilde{e}_c)$ . Logits are converted to temperature-scaled softmax probabilities:  $p(\tilde{e}_*|v_*) = \text{softmax}(\mathbf{z}_*/\tau)$ , where  $\mathbf{z}_a = \{z_a^i\}_{i=1}^{|\mathcal{A}|}$ ,  $\mathbf{z}_o = \{z_o^j\}_{j=1}^{|\mathcal{O}|}$ , and  $\mathbf{z}_c = \{z_c^{ij}\}_{(i,j) \in \mathcal{C}}$ , with  $\tau$  the temperature.

### 3.3 Training Flow

Each continual learning session  $\mathcal{S}^{(t)}$  incrementally extends the model’s compositional space while preserving knowledge from all preceding sessions  $\{\mathcal{S}^{(1)}, \dots, \mathcal{S}^{(t-1)}\}$ . At the beginning of session  $t$ , the student model parameters are initialized from the frozen teacher model from  $\mathcal{S}^{(t-1)}$ , ensuring continuity within the learned embedding space. The attribute, object, and composition classifiers are expanded to accommodate new compositions in  $\mathcal{C}^{(t)}$ . New primitives are regularized by the *Session-Aware Orthogonal Projection Loss (OPL)* to remain orthogonal to prior subspaces and by the *Intra-Session Diversification Loss (IDL)* to encourage diversity within the current session. During training, only the composition prompts corresponding to the new attribute–object pairs in  $\mathcal{C}^{(t)}$  are refined through the *Session-Aware Multi-Modal Fusion* module, where student prompt tokens cross-attend to visual patches. Attribute and object prompt banks are updated via the *Session-Agnostic Multi-Modal Fusion* mechanism, ensuring global semantic consistency across sessions.

Each branch is optimized via:

$$\mathcal{L}_* = -\frac{1}{|\mathcal{X}|} \sum_{x \in \mathcal{X}} \log p(\tilde{e}_*|v_*), \quad * \in \{a, o, c\},$$

and the overall objective combines the branches as the base compositional loss:

$$\mathcal{L}_{\text{CE}} = \alpha_a \mathcal{L}_a + \alpha_o \mathcal{L}_o + \alpha_c \mathcal{L}_c \quad (2)$$

Previously learned knowledge is preserved via Cross-Session Knowledge Distillation (CSKD), where frozen teachers provide targets for overlapping labels. After convergence, the student from  $\mathcal{S}^{(t)}$  is frozen and stored as the teacher for the next session. It anchors prior prompts, keeps new embeddings orthogonal, and transfers knowledge across sessions, enabling continual growth of the compositional space without catastrophic forgetting.

**Cross-Session Knowledge Distillation (CSKD)** We use multi-teacher knowledge distillation to retain knowledge across sessions. During session  $s$ , all models trained in sessions  $t < s$  are frozen and serve as teacher supervisors. We distill only on the logits corresponding to attributes, objects, and compositions seen in earlier sessions. This enforces compatibility with prior primitives while allowing the student to learn new attribute-object pairs. Each teacher is weighted by recency, with later sessions contributing more. For a given teacher (trained at session)  $t$ , say  $\mathbf{z}_a^{(t)} = \{z_a^i | a \in \mathcal{A}^{(t)} \cap \mathcal{A}^{(s)}\}$  and  $\mathbf{z}_a^{(s)} = \{z_a^i | a \in \mathcal{A}^{(t)} \cap \mathcal{A}^{(s)}\}$  denote the teacher logits and student logits for overlapping attributes then the cross-session distillation loss (for attributes) is:

$$\mathcal{L}_{\text{CSKD},a}^{(t)} = \tau^2 \text{KL} \left[ \text{softmax} \left( \frac{\mathbf{z}_a^{(t)}}{\tau} \right) \parallel \text{softmax} \left( \frac{\mathbf{z}_a^{(s)}}{\tau} \right) \right].$$

We can similarly compute distillation loss over shared objects and composition  $\mathcal{L}_{\text{CSKD},o}^{(t)}$  and  $\mathcal{L}_{\text{CSKD},c}^{(t)}$ . The multi-teacher distillation loss is then

$$\mathcal{L}_{\text{CSKD}} = \sum_{t=1}^{s-1} \pi_t \sum_{* \in (a,o,c)} \lambda_* \mathcal{L}_{\text{CSKD},*}^{(t)},$$

where  $\pi_t$  are monotonically decreasing recency weights and  $\lambda_*$  controls the relative importance of attributes, objects, and composition losses.

**Cosine Anchor Alignment Loss (CAL)** CAL constrains embeddings of overlapping attributes/objects to remain aligned with their anchors from the prior session, stabilizing the prompt bank across sessions. Let  $\mathbf{B}_a^{(s)}$  and  $\mathbf{B}_o^{(s)}$  refer to attribute and object prompt banks at current session  $s$  and  $\mathbf{B}_a^{(t)}$  and  $\mathbf{B}_o^{(t)}$  denote corresponding teacher prompt bank.

$$\begin{aligned} \mathcal{L}_{\text{CAL}}^{(t)} = & \sum_{i \in \mathcal{A}^{(s)} \cap \mathcal{A}^{(t)}} \left( 1 - \cos \left( \mathbf{B}_a^{(s)}[i], \mathbf{B}_a^{(t)}[i] \right) \right) \\ & + \sum_{j \in \mathcal{O}^{(s)} \cap \mathcal{O}^{(t)}} \left( 1 - \cos \left( \mathbf{B}_o^{(s)}[j], \mathbf{B}_o^{(t)}[j] \right) \right). \end{aligned}$$

Here,  $\cos(\cdot, \cdot)$  denotes cosine similarity between normalized vectors. This loss preserves directional consistency of recurring attributes and objects and is computed over previous models as before  $\mathcal{L}_{\text{CAL}} = \sum_{t=1}^{s-1} \pi_t \mathcal{L}_{\text{CAL}}^{(t)}$

**Orthogonal Projection Loss (OPL)** At session  $s$ ,  $\mathbf{B}_a$  represents attribute prompt bank. It list all attributes that the model have encountered thus far. We partition it into  $\mathbf{B}_{a,\text{head}}$  and  $\mathbf{B}_{a,\text{tail}}$  such that  $\mathbf{B}_{a,\text{tail}}$  contains only new attributes encountered in this session. OPL loss enforces orthogonality between rows of  $\mathbf{B}_{a,\text{head}}$  and rows of  $\mathbf{B}_{a,\text{tail}}$  by minimizing the average cosine similarity between the two set of vectors (Figure 3 (left)). We can do the same for object prompt bank. Putting it all together, we get

$$\mathcal{L}_{\text{OPL},a} = \frac{1}{|\mathbf{B}_{a,\text{head}}| |\mathbf{B}_{a,\text{tail}}|} \sum_i \sum_j \langle \mathbf{B}_{a,\text{head}}[i], \mathbf{B}_{a,\text{tail}}[j] \rangle,$$

where  $i$  and  $j$  index over  $\mathbf{B}_{a,\text{head}}$  and  $\mathbf{B}_{a,\text{tail}}$ . We similarly compute  $\mathcal{L}_{\text{OPL},o}$  using object prompt bank. The total OPL loss is  $\mathcal{L}_{\text{OPL}} = \mathcal{L}_{\text{OPL},a} + \mathcal{L}_{\text{OPL},o}$ . While CAL stabilizes recurring primitives, it does not explicitly guard against interference from newly introduced prompts. New attributes or objects may drift into directions already occupied in the embedding space, degrading compositional generalization. To address this, we propose the Orthogonal Projection Loss (OPL), which encourages the prompt subspace of the current session to remain orthogonal to that of prior sessions.

**Intra-Session Diversification Loss (IDL)** While OPL enforces independence between head and tail prompt bank, IDL encourages each new attribute and object to occupy a distinct semantic direction, let  $\mathbf{B}_{a,\text{tail}}^{(s)}$  and  $\mathbf{B}_{o,\text{tail}}^{(s)}$  denote the subsets of attribute and object prompts introduced in the current session. The diversification loss penalizes the mean absolute cosine similarity within attributes and objects:

$$\begin{aligned} \mathcal{L}_{\text{IDL}} = & \frac{1}{|\mathbf{B}_{a,\text{tail}}^{(s)}|^2 - |\mathbf{B}_{a,\text{tail}}^{(s)}|} \sum_{i \neq j} \langle \mathbf{B}_{a,\text{tail}}^{(s)}[i], \mathbf{B}_{a,\text{tail}}^{(s)}[j] \rangle \\ & + \frac{1}{|\mathbf{B}_{o,\text{tail}}^{(s)}|^2 - |\mathbf{B}_{o,\text{tail}}^{(s)}|} \sum_{k \neq l} \langle \mathbf{B}_{o,\text{tail}}^{(s)}[k], \mathbf{B}_{o,\text{tail}}^{(s)}[l] \rangle \end{aligned}$$

Here  $i$  and  $j$  index over  $\mathbf{B}_{a,\text{tail}}^{(s)}$  and  $k$  and  $l$  index over  $\mathbf{B}_{o,\text{tail}}^{(s)}$ .

By minimizing intra-session cosine correlations, IDL promotes decorrelated and diverse prompt embeddings for newly added primitives, reducing representational overlap and enhancing compositional expressiveness.

Table 1: Session splits of UT-Zappos and C-GQA datasets.

Dataset	Session	Attr	Obj	Train	Val	Test
UT-Zappos	0	8	6	24	7	9
	1	4	3	27	10	14
	2	4	3	32	13	13
C-GQA	0	233	363	2392	958	730
	1	35	58	491	168	172
	2	32	67	772	358	265
	3	36	64	836	366	275
	4	39	62	562	225	166
	5	38	60	539	217	203

 Table 2: Performance comparisons with state-of-the-art CZSL methods on the UT-Zappos dataset. We report AUC for each session, the average AUC across sessions, and the final improvement over prior methods. Results are shown for Zero-shot Evaluation (ZSEval) on current-session unseen compositions and Continual Zero-shot Evaluation (CZSEval) on accumulated zero-shot test sets. *Final* is the gain of PromptCCZSL (continual) from comparisons.

Method		UT-Zappos Yu and Grauman [2014] (Session Number)				
Name	Venue	0	1	2	Avg	Final
CCZSL–AoP	ECCV’18	42.21	20.34	16.25	26.27	+29.6
CCZSL–SymNet	CVPR’20	41.93	13.58	9.46	21.66	+34.2
CCZSL–VisProdNN	NeurIPS’21	43.59	17.65	2.78	21.34	+34.5
CCZSL–SCEN	CVPR’22	44.00	15.61	8.74	22.78	+33.1
CCZSL–CANet	CVPR’23	45.48	19.89	5.05	23.47	+32.5
Zhang <i>et al.</i>	IJCAI’24	47.70	24.73	18.96	30.46	+25.4
<b>PromptCCZSL (zero-shot)</b>	–	60.42	34.82	32.89	42.71	–
<b>PromptCCZSL (continual)</b>	–	60.42	40.10	26.50	55.86	–

### 3.4 Overall Training Objective.

The overall training objective is defined as:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{ce}}\mathcal{L}_{\text{CE}} + \lambda_{\text{kd}}\mathcal{L}_{\text{CSKD}} + \lambda_{\text{cal}}\mathcal{L}_{\text{CAL}} + \lambda_{\text{opl}}\mathcal{L}_{\text{OPL}} + \lambda_{\text{idl}}\mathcal{L}_{\text{IDL}}, \quad (3)$$

where  $\lambda_{\text{ce}}, \lambda_{\text{kd}}, \lambda_{\text{cal}}, \lambda_{\text{opl}},$  and  $\lambda_{\text{idl}}$  are scalar weights controlling the contribution of each term.

### 3.5 Inference

At test time, each image is evaluated against the union of all seen  $\bigcup_{k=0}^T \mathcal{C}_k$  and unseen compositions, activating all composition embeddings learned across sessions. Final predictions are obtained by scoring each image against every composition as

$$\hat{c} = \arg \max_{c \in \mathcal{C}} \left( \lambda_c p(\tilde{e}_c | v_c) + \lambda_a p(\tilde{e}_a | v_a) \cdot \lambda_o p(\tilde{e}_o | v_o) \right) \quad (4)$$

where  $\lambda_c, \lambda_a,$  and  $\lambda_o$  control the relative influence of each factor.

## 4 Experiments

### 4.1 Experimental Setting

#### 4.1.1 Datasets.

We evaluate our approach on two widely used compositional zero-shot learning (CZSL) datasets: UT-Zappos Yu and Grauman [2014] and C-GQA Hudson and Manning [2019] splits following the continual CZSL (CCZSL) protocol Zhang et al. [2024a]. UT-Zappos is a fine-grained footwear dataset with 50,025 images. C-GQA is a large-scale natural image dataset with 39,298 images. Table 1 provides a summary of the split sessions.



Table 3: AUC comparison of Troika and CSP using KD, Oracle and Prompt-CCZSL-Troika, partial Prompt-CCZSL-CSP variants on UT-Zappos on Continual zero-shot test sets.

Prompt-Based: Continual Compositional Eval.																		
Method	Session 0 Model (Session 0 Data)						Session 1 Model (Session 0-1 Data)						Session 2 Model (Session 0-1-2 Data)					
	AUC	Attr	Obj	Best-S	Best-U	HM	AUC	Attr	Obj	Best-S	Best-U	HM	AUC	Attr	Obj	Best-S	Best-U	HM
CCZSL-CSP Nayak et al. [2023b]	48.46	61.57	87.95	78.76	80.6	57.69	19.19	51.14	42.83	48.41	48.65	36.38	3.2	42.62	32.67	15.05	31.78	14.05
CCZSL-Troika Radford et al. [2021]	60.42	48.55	86.91	83.33	81.62	68.50	15.39	47.98	43.28	48.41	44.15	31.11	0.16	37.47	30.30	0.64	29.67	0.87
PromptCCZSL-CSP*	48.46	61.57	87.95	78.76	80.6	57.69	19.51	50.85	43.08	29.23	48.29	36.54	3.37	42.76	32.84	15.35	31.31	14.54
PromptCCZSL-Troika	60.42	48.55	86.91	83.33	81.62	68.50	40.10	56.09	71.02	73.97	64.14	52.56	26.81	54.08	58.17	49.36	66.10	43.45
Continual Oracle-CSP	48.46	61.57	87.95	78.76	80.6	57.69	35.64	56.45	77.92	73.14	64.97	47.79	32.6	52.0	72.0	63.5	66.4	46.2
Continual Oracle-Troika	60.42	48.55	86.91	83.33	81.62	68.50	44.22	46.94	74.93	76.8	65.3	57.69	33.2	49.6	68.2	59.8	70.2	47.4

Table 4: Performance comparisons with state-of-the-art CZSL methods on C-GQA Hudson and Manning [2019]. We report AUC for each session, the average AUC across sessions, and the final improvement over prior methods. Results are shown under Continual Zero-shot Evaluation (CZSEval) on the accumulated zero-shot test sets.

Method		C-GQA Hudson and Manning [2019] (Session Number)								
Name	Venue	0	1	2	3	4	5	Avg	Final	
CCZSL–AoP	ECCV’18	2.52	1.32	0.97	0.50	0.34	0.27	0.99	<b>+12.21</b>	
CCZSL–SymNet	CVPR’20	3.42	2.28	1.15	0.65	0.63	0.52	1.44	<b>+11.76</b>	
CCZSL–VisProdNN	NeurIPS’21	4.18	0.40	1.19	0.44	0.24	0.15	1.10	<b>+12.1</b>	
CCZSL–SCEN	CVPR’22	3.43	0.64	0.75	0.26	0.11	0.11	0.88	<b>+12.32</b>	
CCZSL–CANet	CVPR’23	5.17	2.49	1.90	1.17	1.27	1.00	2.17	<b>+11.03</b>	
Zhang <i>et al.</i> Zhang et al. [2024a]	IJCAI’24	5.07	3.89	3.88	2.74	2.32	1.83	3.29	<b>+9.91</b>	
<b>PromptCCZSL(continual)</b>	–	<b>16.58</b>	<b>14.53</b>	<b>13.66</b>	<b>11.4</b>	<b>9.83</b>	–	13.2	–	

#### 4.1.2 Evaluation Metrics.

We follow the standard CZSL evaluation protocol Naeem et al. [2021], Purushwalkam et al. [2019], Zhang et al. [2022b], Chao et al. [2017] in the continual setting of Zhang et al. [2024a], reporting the harmonic mean ( $H$ ) of top-1 accuracy on seen ( $S$ ) and unseen ( $U$ ) compositions, attribute accuracy ( $AttrAcc$ ) and object accuracy ( $ObjAcc$ ), as well as session-wise Area Under the Curve ( $AUC$ ), with overall performance summarized by the average AUC across sessions.

#### 4.1.3 Implementation Details.

Our model is implemented in PyTorch with a frozen CLIP ViT-L/14 backbone Paszke et al. [2019] trained for 15 epochs per session on an NVIDIA V100 32GB GPU using the Adam optimizer with learning rate  $3.12 \times 10^{-6}$ . Session 0 is trained using Eq. 2, while later sessions initialize with previous session’s weights and optimize using Eq. 3. Additional hyperparameter details are provided in the supplementary material.

#### 4.1.4 Baselines.

We compare our method against Zhang *et al.* Zhang et al. [2024a], the first method that studies CZSL within a continual learning setup. Additionally, we incorporate the proposed PromptCCZSL learning framework in Troika Huang et al. [2024], a CLIP-based SOTA scheme for CZSL. We refer to this as PromptCCZSL-Troika. We also compare our scheme against CSP Nayak et al. [2023b] method that is adapted for continual learning setup. Lastly, we set up Oracle models for both Troika and CSP, trained on all the data up to session  $s$ .

### 4.2 Results And Analysis

#### 4.2.1 Comparison with the Baseline.

Our PromptCCZSL framework consistently outperforms the baseline Zhang *et al.* Zhang et al. [2024a], achieving higher AUC across all sessions. On UT-Zappos, it exhibits reduced performance degradation from Session 0 to Session 2, demonstrating improved resistance to catastrophic forgetting. On C-GQA, performance remains stable even as the attribute-object space scales to hundreds of primitives, confirming the framework’s scalability in large compositional domains. Specifically, Tables 2 and 4 show that PromptCCZSL attains an average AUC of 55.86% on the UT-Zappos dataset, representing a +25.4% absolute improvement over Zhang *et al.* (30.46%). On the C-GQA dataset, PromptCCZSL achieves an average AUC of 13.2%, a +9.91% absolute improvement over the baseline (3.29%). Compared to earlier non-CLIP methods (AoP Misra and Gupta [2018], SymNet Li et al. [2020b], VisProdNN Saini et al. [2021], SCEN Li et al. [2022b], CANet Liu et al. [2023]), CLIP-based baselines already show superior compositional

Table 5: Comparison of AUC, composition accuracy (Comp), attribute accuracy (Attr), object accuracy (Obj), and harmonic mean (HM) for PromptCCZSL and the baseline CCZSL-Troika on C-GQA Hudson and Manning [2019]. Results are reported under continual evaluation over all zero-shot test sets accumulated up to the current session. Our method consistently outperforms the baseline.

Prompt-Based: Continual Compositional Eval. (Constrained CCZSL)															
Method	Session 0 - S0 Data					Session 1 - S01 Data					Session 2 - S012 Data				
	AUC	HM	Comp	Attr	Obj	AUC	HM	Comp	Attr	Obj	AUC	HM	Comp	Attr	Obj
CCZSL-Troika	16.58	34.57	34.98	45.59	62.32	2.16	11.19	9.4	20.23	27.58	3.82	16.51	15.02	25.47	41.17
PromptCCZSL-Troika	16.58	34.57	34.98	45.59	62.32	14.53	32.02	31.54	41.54	57.24	13.66	30.94	32.0	43.89	58.43

Table 6: Ablation results on UT-zappos under continual zero-shot evaluation, showing the effect of each module when added individually. We report performance for Session 1, Session 2 with single-teacher, and Session 2 with dual-teacher.

Prompt-Based: Continual Compositional Eval.									
Method	Session 1 Model			Session 2 w/Single-Teacher KD <sup>(t-1)</sup>			Session 2 w/Dual-Teacher KD		
	Session 01 Data			Session 012 Data			Session 012 Data		
	AUC	AttrAcc	ObjAcc	AUC	AttrAcc	ObjAcc	AUC	AttrAcc	ObjAcc
+ $\mathcal{L}_{CSKD}$	15.4	49.0	43.2	18.09	42.35	54.29	1.01	37.51	28.86
+SAwM2F	34.5	54.2	65.4	17.74	45.33	52.51	25.06	53.23	54.53
+ $\mathcal{L}_{CAL}$	36.9	53.9	70.0	15.2	43.0	53.5	28.69	53.5	56.45
+ $\mathcal{L}_{CAL_{dual}}$	—	—	—	—	—	—	28.92	54.21	56.52
+ $\mathcal{L}_{OPL}$	39.34	54.2	70.02	17.34	47.15	51.65	28.15	53.12	57.1
+ $\mathcal{L}_{IDL}$	40.1	56.3	70.98	22.21	48.7	55.53	26.5	52.88	57.41

reasoning, and PromptCCZSL further delivers consistent gains across all continual stages. These results highlight the effectiveness of combining multi-teacher, recency-weighted knowledge distillation with our Cosine Anchor, Orthogonal Projection, and Intra-Session Diversification losses.

#### 4.2.2 Comparison with Troika Baselines on UT-Zappos and CGQA

Table 3 and 5 reports the comparison between our proposed PromptCCZSL framework, the original *Troika* model, and its continual variant trained using a simple knowledge distillation strategy (*KD on Troika*) on the UT-Zappos and CGQA dataset. All models are trained on identical dataset splits following the same continual compositional setup to ensure fairness. The *Oracle Troika* represents the upper bound obtained by training the model independently on each session with full access to all data, and thus serves as a non-continual reference for the achievable maximum performance. When the same model is trained continually using only a basic KD objective (*KD on Troika*), performance degrades rapidly over sessions—AUC drops from 61.9% in Session 0 to 0.16% in Session 2, and harmonic mean (H) falls to 0.87%. This demonstrates that simple distillation is insufficient to mitigate catastrophic forgetting in compositional settings. In contrast, our proposed PromptCCZSL maintains stable performance across all sessions, achieving an AUC of 60.42% in Session 0, 40.10% in Session 1 and 26.50% in Session 2, with a final harmonic mean of 43.45%. You can also review Table 5 to see the performance improvements on the C-GQA dataset. The baseline (CCZSL-Troika) exhibits a consistent decline in performance as training progresses. In contrast, PromptCCZSL consistently outperforms the baseline and achieves the highest AUC across sessions. These results confirm that integrating multi-teacher, recency-weighted distillation with semantic anchoring and prompt-space regularization allows PromptCCZSL to effectively preserve prior attribute-object knowledge while adapting to new compositions in a continual manner.

#### 4.2.3 Qualitative Results

In Fig. 4, we present qualitative results on unseen attribute-object compositions from the UT-Zappos benchmark. The images in (a-c) are sampled from the Session-0 unseen split ( $\mathcal{C}^{(0)}_{unseen}$ ), while (d-e) correspond to the Session-1 unseen split ( $\mathcal{C}^{(1)}_{unseen}$ ). For each image, we report predictions from the Session-0 ( $\mathcal{S}^{(0)}$ ), Session-1 ( $\mathcal{S}^{(1)}$ ), and Session-2 ( $\mathcal{S}^{(2)}$ ) PromptCCZSL models. Columns (a-b) show cases where the ( $\mathcal{S}^{(0)}$ ) model correctly predicts the ( $\mathcal{C}^{(0)}_{unseen}$ ). In (a),  $\mathcal{S}^{(0)}, \mathcal{S}^{(1)}, \mathcal{S}^{(2)}$  retain the correct prediction, illustrating successful preservation of early-session knowledge. In (b), the  $\mathcal{S}^{(0)}$  and  $\mathcal{S}^{(1)}$  predict results correctly, whereas the  $\mathcal{S}^{(2)}$  forgets the earlier composition, demonstrating where catastrophic forgetting occurs. Column (c) presents a case where  $\mathcal{S}^{(1)}$  fails but  $\mathcal{S}^{(2)}$  recovers the correct prediction, enabled by our dual-distillation and dual cosine anchor alignment, which reinforce earlier attribute-object relationships. Column (d-e) illustrate examples from the  $\mathcal{C}^{(1)}_{unseen}$ . In (d), both the  $\mathcal{S}^{(1)}$  and  $\mathcal{S}^{(2)}$  models predict correctly, showing strong generalization on newly introduced unseen compositions. In (e), although the  $\mathcal{S}^{(1)}$  model predicts correctly, the  $\mathcal{S}^{(2)}$  model forgets this composition, reflecting a challenging case of interference from later-session updates.

### 4.3 Ablation Analysis

We conduct extensive ablations on UT-Zappos to evaluate the effect of different architectural and training components in our framework (Table 6).

**Session-Aware Multi-Modal Fusion (SAwM2F).** Without SAwM2F, the model forgets earlier compositions (AUC drops to 15.4 in Session 1 and 1.01 in Session 2). Updating only session specific tail compositions restores performance (34.5 and 25.06), and freezing the head drops performance because it breaks alignment, while updating both head and tail keeps meanings consistent and improves continual generalization.

**Cosine Anchor Alignment Loss (CAL).** Adding the cosine loss boosts AUC from 34.5 to 36.9 in Session 1 and from 25.06 to 28.92 in Session 2. A moderate weight ( $\lambda_{cal} = 0.05\text{--}0.1$ ) gives the best balance, as higher values overfit to old sessions and lower values cause forgetting. We therefore adjust  $\lambda_{cal}$  based on overlapping attributes and objects to keep stability and flexibility balanced. Other regularizers (e.g., Gram-matrix or diversity losses) performed worse, showing that cosine anchoring is the most stable option. **Orthogonality between Prompt Spaces (OPL).** OPL helps the model remember earlier sessions and separate new compositions, improving AUC to 39.3 in Session 1 and strengthening retention in Session 2 when combined with cosine anchoring. Moderate OPL leads to smoother updates, while excessive regularization ( $\lambda_{opl} > 0.1$ ) restricts learning new data. **Intra-Session Diversity Loss (IDL).** Adding IDL increases prompt diversity within each session, improving AUC to 40.1 in Session 1 compared to OPL alone. With a small weight ( $\lambda_{idl} = 0.005$ ), IDL adds variety to the representations without hurting stability, leading to smoother and richer learning within each session. **Single vs. Dual Distillation and Anchoring.** Using multiple teachers and cosine anchoring instead of a single teacher improves performance on earlier sessions, raising the overall AUC from 15.2 to 28.92. Further details are provided in the supplementary material.






	$\mathcal{C}_{\text{unseen}}^{(0)}$			$\mathcal{C}_{\text{unseen}}^{(1)}$	
Ground-truth	Leather__Boots.Mid-Calf	Suede__Shoes.Boat.Shoes	Synthetic__Boots.Mid-Calf	Canvas__Shoes.Loafers	Faux.L Leather__Boots.Ankle
UT-Zappos					
$\mathcal{S}^{(0)}$	Leather__Boots.Mid-Calf	Suede__Shoes.Boat.Shoes	Synthetic__Boots.Mid-Calf	-----	-----
$\mathcal{S}^{(1)}$	Leather__Boots.Mid-Calf	Suede__Shoes.Boat.Shoes	<b>Faux.Fur__Boots.Mid-Calf</b>	Canvas__Shoes.Loafers	Faux.L Leather__Boots.Ankle
$\mathcal{S}^{(2)}$	Leather__Boots.Mid-Calf	Suede__Shoes.Loafers	Synthetic__Boots.Mid-Calf	Canvas__Shoes.Loafers	<b>Leather__Boots.Ankle</b>
	(a)	(b)	(c)	(d)	(e)

Figure 4: **Qualitative results on the UT-Zappos** Yu and Grauman [2014] benchmark. (a–c) show the PromptCCZSL model’s ( $\mathcal{S}^{(0)}$ ,  $\mathcal{S}^{(1)}$ ,  $\mathcal{S}^{(2)}$ ) performance on Session-0 unseen zero-shot compositions ( $\mathcal{C}^{(0)}$ unseen), while (d–e) show performance on Session-1 unseen compositions ( $\mathcal{C}^{(1)}$ unseen). (c) highlights that dual distillation helps in improving generalization on previous session results. PromptCCZSL effectively retains knowledge from previous sessions as training progresses. Correct predictions are shown in bold, and incorrect predictions are shown in red.

## 5 Limitations

Our experiments were conducted using an NVIDIA V100 with 32 GB memory. As the number of compositions increases progressively across CCZSL sessions, the computational and memory requirements also grow substantially. This highlights a key CCZSL limitation: the cumulative expansion of compositions can make later sessions prohibitively expensive to train on mid-range GPU hardware.

## 6 Conclusion

We introduce a prompt-based framework designed to support continual compositional zero-shot learning within vision–language models. By constructing a shared soft-prompt bank and enforcing session-consistent alignment through

cosine anchoring and multi-teacher distillation, our method brings structured semantic knowledge into the training pipeline while maintaining flexibility across sessions. These components collectively stabilize primitive representations, mitigate drift, and enable the model to acquire new compositions without catastrophic forgetting issue on earlier ones. Extensive experiments on UT-Zappos demonstrate that our approach achieves SOTA results under closed world setting.

## References

- Yang Zhang, Songhe Feng, and Jiazheng Yuan. Continual compositional zero-shot learning. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 1724–1732, 2024a.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.
- Xiaocheng Lu, Song Guo, Ziming Liu, and Jingcai Guo. Decomposed soft prompt guided fusion enhancing for compositional zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23560–23569, 2023a.
- Nihal V Nayak, Peilin Yu, and Stephen Bach. Learning to compose soft prompts for compositional zero-shot learning. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Wentao Bao, Lichang Chen, Heng Huang, and Yu Kong. Prompting language-informed distribution for compositional zero-shot learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 107–123. Springer, 2024a. URL <https://arxiv.org/abs/2407.01779>.
- Aron Yu and Kristen Grauman. Fine-grained visual comparisons with local learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 192–199, 2014. URL [https://openaccess.thecvf.com/content\\_cvpr\\_2014/papers/Yu\\_Fine\\_Grained\\_Visual\\_Comparisons\\_2014\\_CVPR\\_paper.pdf](https://openaccess.thecvf.com/content_cvpr_2014/papers/Yu_Fine_Grained_Visual_Comparisons_2014_CVPR_paper.pdf).
- Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6700–6709, 2019. URL [https://openaccess.thecvf.com/content\\_CVPR\\_2019/papers/Hudson\\_GQA\\_A\\_New\\_Dataset\\_for\\_Real\\_World\\_Visual\\_Reasoning\\_and\\_Compositional\\_CVPR\\_2019\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2019/papers/Hudson_GQA_A_New_Dataset_for_Real_World_Visual_Reasoning_and_Compositional_CVPR_2019_paper.pdf).
- I. Misra, A. Gupta, and M. Hebert. From red wine to red tomato: Composition with context. In *CVPR*, 2017. [https://openaccess.thecvf.com/content\\_cvpr\\_2017/papers/Misra\\_From\\_Red\\_Wine\\_CVPR\\_2017\\_paper.pdf](https://openaccess.thecvf.com/content_cvpr_2017/papers/Misra_From_Red_Wine_CVPR_2017_paper.pdf).
- T. Nagarajan and K. Grauman. Attributes as operators: Factorizing unseen attribute–object compositions. In *ECCV*, 2018. [https://openaccess.thecvf.com/content\\_ECCV\\_2018/papers/Tushar\\_Nagarajan\\_Attributes\\_as\\_Operators\\_ECCV\\_2018\\_paper.pdf](https://openaccess.thecvf.com/content_ECCV_2018/papers/Tushar_Nagarajan_Attributes_as_Operators_ECCV_2018_paper.pdf).
- Y.-L. Li, X. Xu, L. Nie, T.-S. Chua, and Z. Zhang. Symmetry and group in attribute–object compositions. In *CVPR*, 2020a. [https://openaccess.thecvf.com/content\\_CVPR\\_2020/papers/Li\\_Symmetry\\_and\\_Group\\_in\\_Attribute-Object\\_Compositions\\_CVPR\\_2020\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2020/papers/Li_Symmetry_and_Group_in_Attribute-Object_Compositions_CVPR_2020_paper.pdf).
- T. Zhang, J. Chen, L. Wang, et al. Learning invariant visual representations for compositional zero-shot learning. In *ECCV*, 2022a. [https://www.ecva.net/papers/eccv\\_2022/papers\\_ECCV/papers/136840335.pdf](https://www.ecva.net/papers/eccv_2022/papers_ECCV/papers/136840335.pdf).
- X. Li, J. Zhou, L. Liu, et al. Siamese contrastive embedding network for compositional zero-shot learning. In *CVPR*, 2022a. [https://openaccess.thecvf.com/content/CVPR2022/papers/Li\\_Siamese\\_Contrastive\\_Embedding\\_Network\\_for\\_Compositional\\_Zero-Shot\\_Learning\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2022/papers/Li_Siamese_Contrastive_Embedding_Network_for_Compositional_Zero-Shot_Learning_CVPR_2022_paper.pdf).
- Y. Hao, K. Han, and K.-Y. K. Wong. Learning attention as disentangler for compositional zero-shot learning. In *CVPR*, 2023. <https://arxiv.org/abs/2303.15111>.
- Z. Wang, J. Chen, J. Chen, et al. Learning conditional attributes for compositional zero-shot learning (canet). In *CVPR*, 2023. <https://jingchenchen.github.io/files/papers/2023/CVPR-CANET.pdf>.
- N. V. Nayak, P. Yu, and S. Bach. Learning to compose soft prompts for compositional zero-shot learning (csp). In *ICLR*, 2023b. <https://openreview.net/pdf?id=S8-A2FXnIh>.
- Y. Lu, Z. Wu, S. Zhang, et al. Decomposed soft prompt guided fusion enhancing for compositional zero-shot learning (dfsp). In *CVPR*, 2023b. [https://openaccess.thecvf.com/content/CVPR2023/papers/Lu\\_Decomposed\\_Soft\\_Prompt\\_Guided\\_Fusion\\_Enhancing\\_for\\_Compositional\\_Zero-Shot\\_Learning\\_CVPR\\_2023\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2023/papers/Lu_Decomposed_Soft_Prompt_Guided_Fusion_Enhancing_for_Compositional_Zero-Shot_Learning_CVPR_2023_paper.pdf).

- S. Huang, B. Gong, Y. Feng, M. Zhang, Y. Lv, and D. Wang. Troika: Multi-path cross-modal traction for compositional zero-shot learning. In *CVPR*, 2024. [https://openaccess.thecvf.com/content/CVPR2024/papers/Huang\\_Troika\\_Multi-Path\\_Cross-Modal\\_Traction\\_for\\_Compositional\\_Zero-Shot\\_Learning\\_CVPR\\_2024\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2024/papers/Huang_Troika_Multi-Path_Cross-Modal_Traction_for_Compositional_Zero-Shot_Learning_CVPR_2024_paper.pdf).
- W. Bao, L. Chen, H. Huang, and Y. Kong. Prompting language-informed distribution for compositional zero-shot learning (plid). In *ECCV*, 2024b. <https://arxiv.org/pdf/2305.14428>.
- Y. Wang and J. Deng. Hierarchical prompt learning for compositional zero-shot learning. In *IJCAI*, 2023. <https://www.ijcai.org/proceedings/2023/0163.pdf>.
- J. Kirkpatrick, R. Pascanu, N. Rabinowitz, and et al. Overcoming catastrophic forgetting in neural networks. In *PNAS*, 2017. <https://www.pnas.org/doi/10.1073/pnas.1611835114>.
- F. Zenke, B. Poole, and S. Ganguli. Continual learning through synaptic intelligence. In *ICML*, 2017. <https://proceedings.mlr.press/v70/zenke17a.html>.
- S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017. [https://openaccess.thecvf.com/content\\_cvpr\\_2017/papers/Rebuffi\\_iCaRL\\_Incremental\\_Classifier\\_CVPR\\_2017\\_paper.pdf](https://openaccess.thecvf.com/content_cvpr_2017/papers/Rebuffi_iCaRL_Incremental_Classifier_CVPR_2017_paper.pdf).
- F. M. Castro, M. J. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari. End-to-end incremental learning. In *ECCV*, 2018. [https://openaccess.thecvf.com/content\\_ECCV\\_2018/papers/Francisco\\_M.\\_Castro\\_End-to-End\\_Incremental\\_Learning\\_ECCV\\_2018\\_paper.pdf](https://openaccess.thecvf.com/content_ECCV_2018/papers/Francisco_M._Castro_End-to-End_Incremental_Learning_ECCV_2018_paper.pdf).
- A. Mallya and S. Lazebnik. Packnet: Adding multiple tasks to a single network via layerwise pruning. In *CVPR*, 2018. [https://openaccess.thecvf.com/content\\_cvpr\\_2018/papers/Mallya\\_PackNet\\_Adding\\_Multiple\\_CVPR\\_2018\\_paper.pdf](https://openaccess.thecvf.com/content_cvpr_2018/papers/Mallya_PackNet_Adding_Multiple_CVPR_2018_paper.pdf).
- J. Serra, D. Suris, M. Miron, and A. Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *ICML*, 2018. <https://proceedings.mlr.press/v80/serra18a.html>.
- Z. Li and D. Hoiem. Learning without forgetting. In *ECCV*, 2016. <https://arxiv.org/abs/1606.09282>.
- P. Dhar, R. V. Singh, K.-C. Peng, Z. Wu, and R. Chellappa. Learning without memorizing. In *CVPR*, 2019. [https://openaccess.thecvf.com/content\\_CVPR\\_2019/papers/Dhar\\_Learning\\_Without\\_Memorizing\\_CVPR\\_2019\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2019/papers/Dhar_Learning_Without_Memorizing_CVPR_2019_paper.pdf).
- B. Heo, M. Lee, S. Yun, and J. Y. Choi. A comprehensive overhaul of feature distillation. In *ICCV*, 2019. [https://openaccess.thecvf.com/content\\_ICCV\\_2019/papers/Heo\\_A\\_Comprehensive\\_Overhaul\\_of\\_Feature\\_Distillation\\_ICCV\\_2019\\_paper.pdf](https://openaccess.thecvf.com/content_ICCV_2019/papers/Heo_A_Comprehensive_Overhaul_of_Feature_Distillation_ICCV_2019_paper.pdf).
- L. Beyer, X. Zhai, A. Kolesnikov, and N. Houlsby. Knowledge distillation: A good teacher is patient and consistent. In *CVPR*, 2022. [https://openaccess.thecvf.com/content/CVPR2022/papers/Beyer\\_Knowledge\\_Distillation\\_A\\_Good\\_Teacher\\_Is\\_Patient\\_and\\_Consistent\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2022/papers/Beyer_Knowledge_Distillation_A_Good_Teacher_Is_Patient_and_Consistent_CVPR_2022_paper.pdf).
- W. Park, D. Kim, Y. Lu, and M. Cho. Relational knowledge distillation. In *CVPR*, 2019. [https://openaccess.thecvf.com/content\\_CVPR\\_2019/papers/Park\\_Relational\\_Knowledge\\_Distillation\\_CVPR\\_2019\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2019/papers/Park_Relational_Knowledge_Distillation_CVPR_2019_paper.pdf).
- G. H. Wang, Y. Ge, and J. Wu. Distilling knowledge by mimicking features. *IEEE TIP*, 2021. <https://arxiv.org/pdf/2011.01424>.
- D. Chen, J. Mei, C. Wang, Y. Feng, and C. Chen. Knowledge distillation with the reused teacher classifier. In *CVPR*, 2022. [https://openaccess.thecvf.com/content/CVPR2022/papers/Chen\\_Knowledge\\_Distillation\\_With\\_the\\_Reused\\_Teacher\\_Classifier\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2022/papers/Chen_Knowledge_Distillation_With_the_Reused_Teacher_Classifier_CVPR_2022_paper.pdf).
- Y. Wang, X. Zhang, et al. Improving kd via regularizing feature direction and norm. In *ECCV*, 2024. [https://www.ecva.net/papers/eccv\\_2024/papers\\_ECCV/papers/03432.pdf](https://www.ecva.net/papers/eccv_2024/papers_ECCV/papers/03432.pdf).
- H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou. Training data-efficient image transformers (deit). In *ICML*, 2021. <https://arxiv.org/pdf/2012.12877>.
- C. Yang et al. Vitkd: Feature-based knowledge distillation for vision transformers. In *CVPRW*, 2024. [https://openaccess.thecvf.com/content/CVPR2024W/PBDL/papers/Yang\\_ViTKD\\_Feature-based\\_Knowledge\\_Distillation\\_for\\_Vision\\_Transformers\\_CVPRW\\_2024\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2024W/PBDL/papers/Yang_ViTKD_Feature-based_Knowledge_Distillation_for_Vision_Transformers_CVPRW_2024_paper.pdf).
- M. Luo et al. Label-specific adapter (lada) for clip in continual learning, 2025. <https://icml.cc/virtual/2025/poster/43751>.
- J. Wang et al. Learning to prompt for continual learning (l2p). In *CVPR*, 2022a. [https://openaccess.thecvf.com/content/CVPR2022/papers/Wang\\_Learning\\_To\\_Prompt\\_for\\_Continual\\_Learning\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2022/papers/Wang_Learning_To_Prompt_for_Continual_Learning_CVPR_2022_paper.pdf).

- J. Wang et al. Dualprompt: Complementary prompting for continual learning. In *ECCV*, 2022b. [https://www.ecva.net/papers/eccv\\_2022/papers\\_ECCV/papers/136860617.pdf](https://www.ecva.net/papers/eccv_2022/papers_ECCV/papers/136860617.pdf).
- J. S. Smith et al. Coda-prompt: Continual decomposed attention-based prompting. In *CVPR*, 2023. [https://openaccess.thecvf.com/content/CVPR2023/papers/Smith\\_CODA-Prompt\\_COntinual\\_Decomposed\\_Attention-Based\\_Prompting\\_for\\_Rehearsal-Free\\_Continual\\_Learning\\_CVPR\\_2023\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2023/papers/Smith_CODA-Prompt_COntinual_Decomposed_Attention-Based_Prompting_for_Rehearsal-Free_Continual_Learning_CVPR_2023_paper.pdf).
- X. Gao et al. Cprompt: Classifier/prompt consistency for continual prompting. In *CVPR*, 2024. <https://arxiv.org/pdf/2403.08568>.
- Y. Zhang et al. Kdp: Knowledge distillation based on prompts for cdl. In *arXiv*, 2024b. <https://arxiv.org/pdf/2407.13911>.
- Q. Li et al. Promptkd: Unsupervised prompt distillation for vlms. In *CVPR*, 2024. [https://openaccess.thecvf.com/content/CVPR2024/papers/Li\\_PromptKD\\_Unsupervised\\_Prompt\\_Distillation\\_for\\_Vision-Language\\_Models\\_CVPR\\_2024\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2024/papers/Li_PromptKD_Unsupervised_Prompt_Distillation_for_Vision-Language_Models_CVPR_2024_paper.pdf).
- M. Zheng et al. Adapt without forgetting: Distill proximity from dual teachers. In *ECCV*, 2024. <https://arxiv.org/abs/2403.09296>.
- W. Liu, Y. Ding, H. Li, X. Li, and X. Zhu. Continual compositional zero-shot learning with super-primitives and dual distillation. In *Proceedings of the 33rd International Joint Conference on Artificial Intelligence (IJCAI)*, 2024. <https://arxiv.org/pdf/2402.17251>.
- Muhammad Ferjad Naeem, Yongqin Xian, Federico Tombari, and Zeynep Akata. Learning graph embeddings for compositional zero-shot learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 953–962, 2021. URL <https://arxiv.org/abs/2102.01987>.
- Senthil Purushwalkam, Maximilian Nickel, Marc’Aurelio Ranzato, and Abhinav Gupta. Task-driven modular networks for zero-shot compositional learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3593–3602, 2019. URL <https://arxiv.org/abs/1905.05908>.
- Tian Zhang, Kongming Liang, Ruoyi Du, Xian Sun, Zhanyu Ma, and Jun Guo. Learning invariant visual representations for compositional zero-shot learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 55–72. Springer, 2022b. URL <https://arxiv.org/abs/2206.00415>.
- Wei-Lun Chao, Soravit Changpinyo, Boqing Gong, and Fei Sha. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 895–903, 2017. URL <https://arxiv.org/abs/1605.04253>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, pages 8024–8035, 2019. URL [https://papers.nips.cc/paper\\_files/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html](https://papers.nips.cc/paper_files/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html).
- I. Misra and A. Gupta. Attributes as operators: Factorizing unseen attribute-object compositions. In *ECCV*, 2018. <https://arxiv.org/pdf/1803.09851>.
- Y. Li, C. Xu, X. Mao, and et al. Symmetry and group in attribute-object compositions. In *CVPR*, 2020b. [https://openaccess.thecvf.com/content\\_CVPR\\_2020/papers/Li\\_Symmetry\\_and\\_Group\\_in\\_Attribute-Object\\_Compositions\\_CVPR\\_2020\\_paper.pdf](https://openaccess.thecvf.com/content_CVPR_2020/papers/Li_Symmetry_and_Group_in_Attribute-Object_Compositions_CVPR_2020_paper.pdf).
- N. Saini, V. Verma, and P. Rai. Disentangling visual product representations for compositional zero-shot learning. In *NeurIPS*, 2021. <https://neurips.cc/virtual/2021/35448>.
- Y. Li, C. Xu, X. Mao, and et al. Siamese contrastive embedding network for compositional zero-shot learning. In *CVPR*, 2022b. [https://openaccess.thecvf.com/content/CVPR2022/papers/Li\\_Siamese\\_Contrastive\\_Embedding\\_Network\\_for\\_Compositional\\_Zero-Shot\\_Learning\\_CVPR\\_2022\\_paper.pdf](https://openaccess.thecvf.com/content/CVPR2022/papers/Li_Siamese_Contrastive_Embedding_Network_for_Compositional_Zero-Shot_Learning_CVPR_2022_paper.pdf).
- W. Liu, F. Shen, Z. Lin, and et al. Canet: Compositional attribute network for zero-shot learning. In *CVPR*, 2023. <https://arxiv.org/pdf/2305.17940>.
- Zhiqi Kang, Liyuan Wang, Xingxing Zhang, and Karteek Alahari. Advancing prompt-based methods for replay-independent general continual learning. In *International Conference on Learning Representations (ICLR)*, 2025.
- Ketan Rajshekhar Shahapure and Charles Nicholas. Cluster quality analysis using silhouette score. In *2020 IEEE 7th international conference on data science and advanced analytics (DSAA)*, pages 747–748. IEEE, 2020.

## Supplementary Material

This document is supplementary material for our work on *Continual Composition Zero-shot Learning (CCZSL)*. Here, we present following sections (a) Evaluation settings for catastrophic forgetting in CCZSL and results thereof, (b) Results on C-GQA dataset, (c) Ablation on hyper-parameters, (d) Qualitative analysis, and finally (e) limitations and future research direction.

## 7 Definitions

### 7.1 Evaluation Settings

CCZSL method should be evaluated both for its ability to adapt to new compositions and its ability to retain previously acquired zero-shot knowledge, thereby mitigating catastrophic forgetting. In main paper we have provided two settings, (i) Zero-shot evaluation (ZSEval), where model trained on session  $t$  in continual setting is only tested on unseen (zero-shot) compositions in session  $t$ , (ii) Continual zero-shot evaluation (CZSEval), model trained on session  $t$  in continual setting is tested on unseen (zero-shot) compositions in session  $t$  and all the previous sessions. Below we present evaluation setting for the catastrophic forgetting and for composition accuracy.

**Catastrophic Forgetting in CCZSL (CFZSEval):** Here objective is to measure catastrophic forgetting on unseen compositions, that is whether model is still able to accurately prediction unseen composition which it was able to do in previous session (see Figures and ). Following Kang et al. [2025] we define the catastrophic forgetting measure for unseen (zero-shot) compositions as Let  $\Xi_i$  be test-samples labeled with unseen (zero-shot) compositions in session  $i$ . Let model  $M_t$  denote the model trained on session  $t$  and  $AUC_i^t$  be AUC computed on  $\Xi_i$  using  $M_t$ , where  $0 \leq i \leq t$ . Then average forgetting will be

$$\mathcal{F}_{AUC}^i = \frac{1}{(T-i)} \sum_{t=i}^{T-1} \text{abs}(AUC_i^t - AUC_i^i) \quad (5)$$

$$\mathcal{F}_{AUC} = \frac{1}{T} \sum_{i=0}^{T-1} \mathcal{F}_{AUC}^i \quad (6)$$

Larger  $\mathcal{F}_{AUC}$  means more forgetting, therefore lower is better. In the Constrained-CCZSL setting, an unseen composition in an earlier session remains unseen in all the next sessions. Note that, objects and attributes might reoccur in next sessions, that's why we evaluate catastrophic forgetting using  $\mathcal{F}_{AUC}$  (Tab. 7).

Empirically, we also show in Fig. 5 and 6 under the Constrained-CCZSL setting on UT-Zappos Hudson and Manning [2019], the baseline shows increased forgetting as sessions progress, whereas PromptCCZSL–Troika preserves substantially higher performance on early-session zero-shot compositions.

**Composition Accuracy.** We used Composition Accuracy (CompAcc) alongside Attribute Accuracy (AttrAcc) and Object Accuracy (ObjAcc). AttrAcc and ObjAcc assess how well the model retains and adapts knowledge of individual attributes and objects, whereas CompAcc evaluates performance on both seen and unseen attribute–object pairs. For a zero-shot composition  $(a, o) \in \mathcal{C}_{\text{unseen}}$ , a prediction is correct only if  $\hat{a} = a$  and  $\hat{o} = o$ . CompAcc is computed as the proportion of samples satisfying this condition. While AttrAcc and ObjAcc capture single-factor recognition, CompAcc is essential in CCZSL because it directly measures compositional generalization to unseen attribute–object combinations.

## 8 More Experiment Results

**subsectionAblation Results on Hyperparameters** **Cross Session Knowledge Distillation (CSKD)** In Tab.8 we find a trade-off between balancing Cross Entropy (CE) and Cross-Session Knowledge Distillation (CSKD). Given a higher CSKD weight helps the student retain the previous session knowledge. In our experiments, KD = 0.65 and CE = 0.35 gives the best AUC, outperforming the other settings.

**Orthogonality between Prompt Spaces (OPL).** OPL contributes to help in preventing new attributes and objects from overlapping with those learned in earlier sessions. With a very small weight (0.001), the new features become somewhat orthogonal compared to the previous setup, but noticeable overlap still remains in the t-SNE plots. Increasing

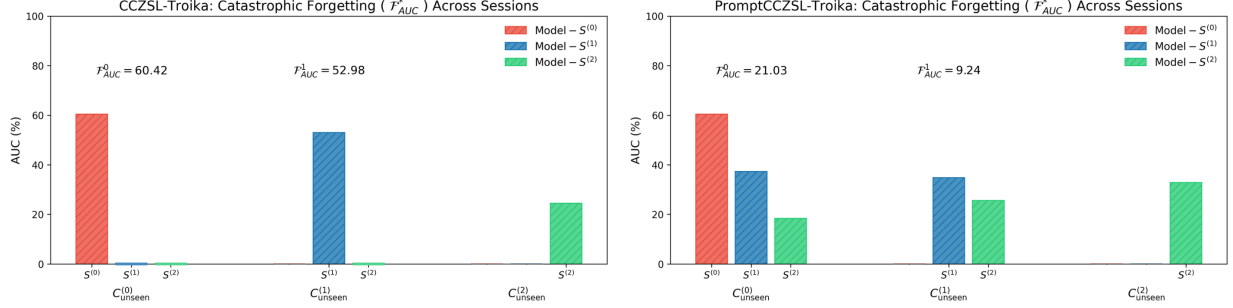


Figure 5: **Catastrophic Forgetting Problem in Continual Composition Zero-shot Learning (CCZSL) through the lens of AUC.** We use AUC computed across training sessions to capture catastrophic forgetting for **(top)** CCZSL-Troika (baseline) and **(bottom)** PromptCCZSL-Troika (our model) on UT-Zappos benchmark. Here  $S^{(0)}$  denote the model trained using  $C_{seen}^{(0)}$ ,  $S^{(1)}$  denote the model trained using  $C_{seen}^{(0)} \rightarrow C_{seen}^{(1)}$ , and  $S^{(2)}$  denote the model trained using  $C_{seen}^{(0)} \rightarrow C_{seen}^{(1)} \rightarrow C_{seen}^{(2)}$ .  $S^{(0)}$  is evaluated on  $C_{unseen}^{(0)}$ ,  $S^{(1)}$  is evaluated on  $C_{unseen}^{(0)}$  and  $C_{unseen}^{(1)}$ , and  $S^{(2)}$  is evaluated on  $C_{unseen}^{(0)}$ ,  $C_{unseen}^{(1)}$ , and  $C_{unseen}^{(2)}$ . Top plot suggests that CCZSL-Troika suffers from severe catastrophic forgetting. When the model is trained on new data, it loses all ability to perform compositional inference on unseen data from previous session as measured using AUC, e.g., AUC values for both  $S^{(1)}$  and  $S^{(2)}$  on  $C_{unseen}^{(0)}$  is close to zero. Similarly, only  $S^{(1)}$  posts good AUC scores on  $C_{unseen}^{(1)}$  and only  $S^{(2)}$  posts good scores on  $C_{unseen}^{(2)}$ . Bottom plots shows a different trend. For our model, both  $S^{(1)}$  and  $S^{(2)}$  are able to carry out compositional inference on  $C_{unseen}^{(0)}$ . This suggests that as the model is trained on new data, it retains the ability to perform inference on unseen data from previous sessions. Similarly,  $S^{(2)}$  is still able to achieve good AUC on  $C_{unseen}^{(1)}$ . These plots suggest that the proposed model gracefully handles catastrophic forgetting.

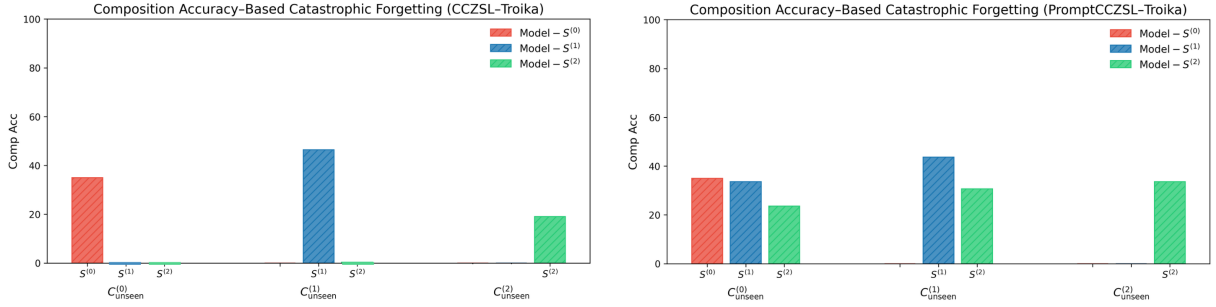


Figure 6: **Catastrophic Forgetting Problem in Continual Composition Zero-shot Learning (CCZSL) through the lens of Composition Accuracy.** We use composition accuracy computed across training sessions to capture catastrophic forgetting for **(top)** CCZSL-Troika (baseline) and **(bottom)** PromptCCZSL-Troika (our model) on UT-Zappos benchmark. Here  $S^{(0)}$  denote the model trained using  $C_{seen}^{(0)}$ ,  $S^{(1)}$  denote the model trained using  $C_{seen}^{(0)} \rightarrow C_{seen}^{(1)}$ , and  $S^{(2)}$  denote the model trained using  $C_{seen}^{(0)} \rightarrow C_{seen}^{(1)} \rightarrow C_{seen}^{(2)}$ .  $S^{(0)}$  is evaluated on  $C_{unseen}^{(0)}$ ,  $S^{(1)}$  is evaluated on  $C_{unseen}^{(0)}$  and  $C_{unseen}^{(1)}$ , and  $S^{(2)}$  is evaluated on  $C_{unseen}^{(0)}$ ,  $C_{unseen}^{(1)}$ , and  $C_{unseen}^{(2)}$ . Top plot suggests that CCZSL-Troika suffers from severe catastrophic forgetting. When the model is trained on new data, it loses all ability to perform compositional inference on unseen data from previous session as measured using composition accuracy, e.g., composition accuracy values for both  $S^{(1)}$  and  $S^{(2)}$  on  $C_{unseen}^{(0)}$  is close to zero. Similarly, only  $S^{(1)}$  posts good composition accuracy scores on  $C_{unseen}^{(1)}$  and only  $S^{(2)}$  posts good scores on  $C_{unseen}^{(2)}$ . Bottom plots shows a different trend. For our model, both  $S^{(1)}$  and  $S^{(2)}$  are able to carry out compositional inference on  $C_{unseen}^{(0)}$ . This suggests that as the model is trained on new data, it retains the ability to perform inference on unseen data from previous sessions. Similarly,  $S^{(2)}$  is still able to achieve good AUC on  $C_{unseen}^{(1)}$ . These plots suggest that the proposed model gracefully handles catastrophic forgetting.

the OPL weight especially up to 0.05 gives the best AUC and produces the most clearly separated, non-overlapping compositions, as shown in Tab. 9.

**Intra-Session Diversity Loss (IDL).** We first make the new session embeddings orthogonal to the previous ones, and then add a small diversity loss so the intra-session attribute and object embeddings do not overlap. Applying IDL on top of different Orthogonal Projection Loss (OPL) weights shows clear AUC improvements, as the model learns current-session features without mixing them together. The best separation and highest AUC occur with Intra-session Diversification Loss (IDL) = 0.005 and OPL = 0.05, as shown in Tab. 10 and tsn plot (Fig. 7).



Table 7: Performance comparison between (baseline) CCZSL–Troika and (ours) PromptCCZSL–Troika on the UT-Zappos dataset Yu and Grauman [2014] under the constrained continual CZSL setting. The baseline (CCZSL–Troika) exhibits strong catastrophic forgetting, with accuracy on early sessions (Session 0 and Session 1) collapsing as later sessions are learned. In contrast, PromptCCZSL–Troika maintains substantially higher AUC and attribute, object accuracy on earlier sessions, yielding much lower forgetting scores  $\mathcal{F}_{AUC}$  (from Eq. 5). These results highlight the effectiveness of our prompt-based continual learning strategy in preserving prior knowledge while maintaining competitive generalization across sessions.

CCZSL – Troika												
Model \ Data	Session 0 ( $\mathcal{C}_{\text{unseen}}^{(0)}$ )				Session 1 ( $\mathcal{C}_{\text{unseen}}^{(1)}$ )				Session 2 ( $\mathcal{C}_{\text{unseen}}^{(2)}$ )			
	AUC	Attr	Obj	HM	AUC	Attr	Obj	HM	AUC	Attr	Obj	HM
Session 0 ( $\mathcal{S}^{(0)}$ )	60.42	48.55	86.91	68.50	–	–	–	–	–	–	–	–
Session 1 ( $\mathcal{S}^{(1)}$ )	0.0	46.12	17.04	0.0	52.98	60.15	82.73	59.24	–	–	–	–
Session 2 ( $\mathcal{S}^{(2)}$ )	0.0	37.47	25.35	0.0	0.0	38.29	8.2	0.0	24.49	36.15	77.8	38.0
$\mathcal{F}_{AUC}$	60.42	–	–	–	52.98	–	–	–	–	–	–	–

PromptCCZSL – Troika												
Model \ Data	Session 0 ( $\mathcal{C}_{\text{unseen}}^{(0)}$ )				Session 1 ( $\mathcal{C}_{\text{unseen}}^{(1)}$ )				Session 2 ( $\mathcal{C}_{\text{unseen}}^{(2)}$ )			
	AUC	Attr	Obj	HM	AUC	Attr	Obj	HM	AUC	Attr	Obj	HM
Session 0 ( $\mathcal{S}^{(0)}$ )	60.42	48.55	86.91	68.50	–	–	–	–	–	–	–	–
Session 1 ( $\mathcal{S}^{(1)}$ )	37.35	55.61	6.69	52.39	34.82	57.34	77.42	47.06	–	–	–	–
Session 2 ( $\mathcal{S}^{(2)}$ )	18.36	54.64	52.08	34.68	25.58	52.97	55.46	41.11	32.89	45.78	78.59	48.08
$\mathcal{F}_{AUC}$	21.03	–	–	–	9.24	–	–	–	–	–	–	–

**Session-Aware Multi-Modal Fusion (SAwM2F) with CAL** Tab. 11 showed how updating or freezing the head–tail textual embeddings affects the Session-Aware module and found that enhancing both parts with cosine anchor alignment loss (CAL) gives a very low AUC (23.69). Freezing the head slightly improves AUC but causes the model to forget earlier sessions. Our Session-Aware Module, where only the tail (current-session embedding) is updated, achieves a much higher AUC of 47.15, showing that session specific adaptation works best when previous-session embeddings remain intact.

Table 8: Ablation of CSKD weightage ( $\lambda_{kd}$ ) and compositional loss weightage ( $\lambda_{ce}$ ) on UT-Zappos.

Session 1 Model		Continual – S01 Data			
$\lambda_{kd}$	$\lambda_{ce}$	AUC	CompAcc	AttrAcc	ObjAcc
0.5	1	23.17	23.38	39.40	61.67
0.65	1	22.6	22.06	39.17	61.11
1.0	1	23.56	23.17	38.05	59.82
0.65–0.35	0.35–0.65	23.5	23.17	38.05	59.82

Table 9: Ablation of OPL weights ( $\lambda_{opl}$ ) on UT-Zappos50K Yu and Grauman [2014], evaluated using the validation set.

Session 1 Model				
$\lambda_{opl}$	continual – S01 Data			
	AUC	CompAcc	AttrAcc	ObjAcc
0.01	43.74	36.4	49.86	74.76
0.05	51.22	41.4	53.11	77.33
decay <sub>[0.01]</sub> ↓	44.07	37.66	50.03	75.32

## 9 Qualitative Analysis

### 9.1 t-SNE Visualization

To better understand how each module in our PromptCCZSL–Troika framework influences the representation space, we compute the silhouette score Shahapure and Nicholas [2020] which measures how well samples are clustered with respect to their semantic group (higher is better and indicates stronger separation). We visualize the textual embeddings for Session-1 using t-SNE under two configurations. In Fig. 7 (a), we apply Cross-Session Knowledge Distillation

Table 10: Ablation of OPL weights ( $\lambda_{opl}$ ) on UT-Zappos50K Yu and Grauman [2014] with IDL weights ( $\lambda_{idl}=0.005$ ), evaluated using the validation set.

Session 1 Model				
$\lambda_{opl}, \lambda_{idl}=0.005$	continual – S01 Data			
	AUC	CompAcc	AttrAcc	ObjAcc
0.009	39.7	36.3	47.17	77.11
0.01	39.44	35.25	48.18	74.43
0.05	46.26	36.37	48.63	75.94
decay <sub>[0.01]</sub> ↓	35.7	34.30	46.28	76.83

Table 11: Ablation: Impact of freezing vs. updating cross-session textual embeddings (head/tail) in SAwM2F through cross-attention (CA) with visual tokens, evaluated on the validation set.

Session 1 Model					
SAwM2F (CA)		continual – S01 Data			
head	tail	AUC	CompAcc	AttrAcc	ObjAcc
✓	✓	23.69	22.89	37.55	60.04
✗	✓	47.15	38.33	49.58	77.11
freeze	✓	24.26	27.53	42.87	63.07

(CSKD), Cosine Anchor Alignment Loss (CAL), and Session-Aware Multi-Modal Fusion (SAwM2F). While the model retains several Session-0 composition clusters, noticeable overlap remains between attribute, object, and composition embeddings, reflected by the lower silhouette scores. In Fig. 7 (b), we additionally incorporate Orthogonal Projection Loss (OPL) and Intra-session Diversification Loss (IDL). Orthogonalizing the current session embeddings against prior-session representations and enforcing intra-session diversity produces significantly cleaner, non-overlapping clusters with a consistent increase in silhouette score across attributes, objects, and compositions. These results demonstrate that OPL and IDL are crucial for preventing embedding collapse, improving cluster separability, and preserving previous-session knowledge during continual compositional learning.

## 10 Limitations

Our experiments were conducted using an NVIDIA V100 with 32 GB memory. As the number of compositions increases progressively across CCZSL sessions, the computational and memory requirements also grow substantially. This highlights a key CCZSL limitation: the cumulative expansion of compositions can make later sessions prohibitively expensive to train on mid-range GPU hardware.

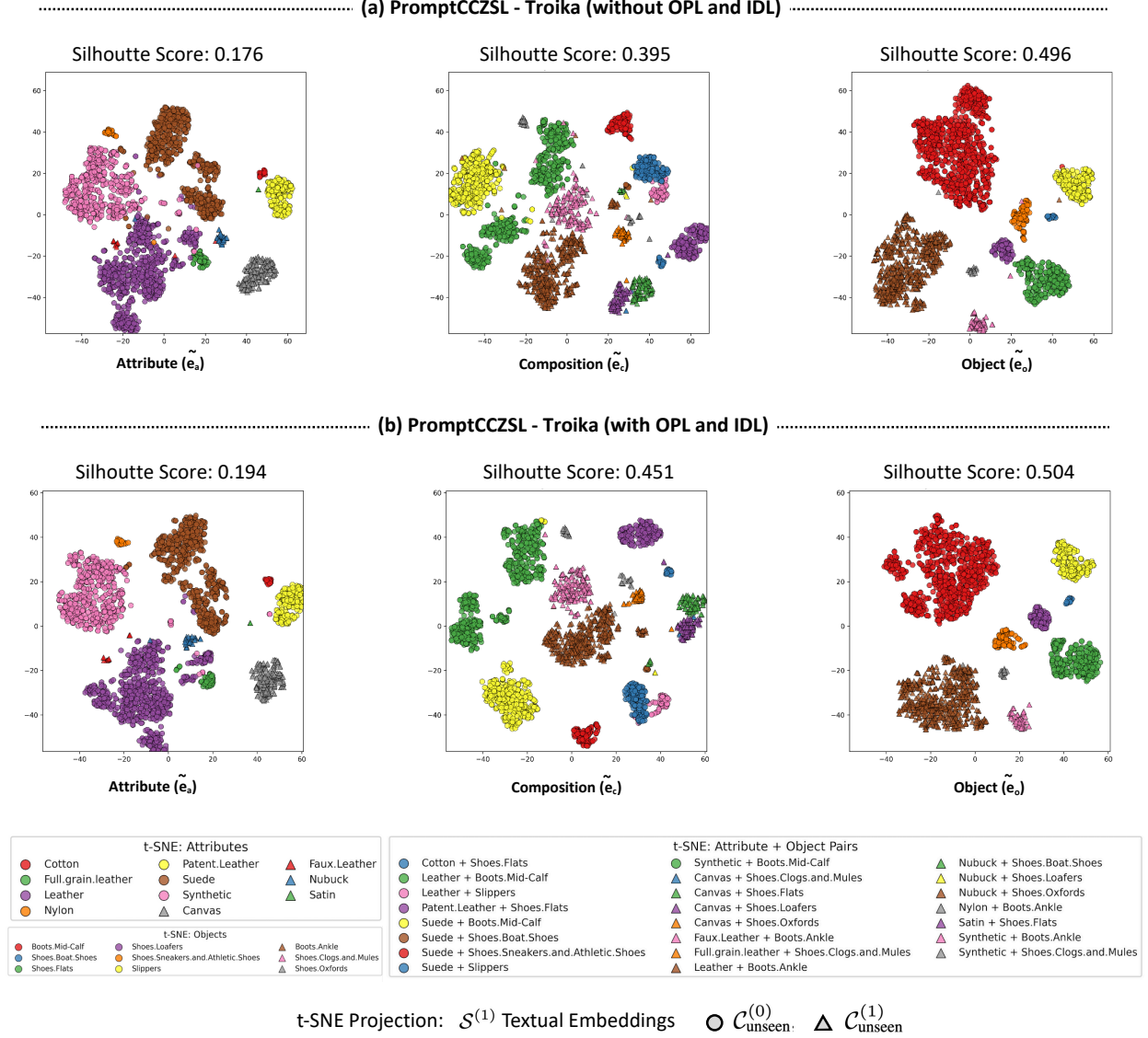


Figure 7: **t-SNE Projections of Textual Embeddings** from the Session-1 model on the UT-Zappos Yu and Grauman [2014] Session-01 Unseen dataset for PromptCCZSL-Troika. (a) The variant integrated Cross-Session Knowledge Distillation (CSKD), Cosine Anchor Alignment Loss (CAL), and Session-Aware Multi-Modal Fusion (SAwM2F) modules exhibits overlapping attribute, object, and composition clusters, yielding lower silhouette scores Shahapure and Nicholas [2020]. (b) Incorporating Orthogonal Projection Loss (OPL) and Intra-session Diversification Loss (IDL) produces more clearly separated clusters, increases silhouette scores, and mitigates catastrophic forgetting by enforcing orthogonality between new-session embeddings and previously learned representations. Average silhouette scores summarizes overall clustering quality (ranging from  $-1$  to  $1$ , with higher values indicating better cluster separation).