This appendix provides additional details for the ICLR 2026, titled "One Patch, One Text: Sparse Alignment for Closing CLIP's Modality Gap for Compositional Zero-Shot Learning". It is orgnized as follows:

- §A Detailed Experiment Settings.
- §B More Ablation Experiments.
- §C More Qualitative Experiments.
- §D Pseudo-code.
- §E Statement for Using Large Language Models.

### A DETAILED EXPERIMENT SETTINGS

**Detailed Dataset Split Statistics.** We conduct experiments on three widely-used datasets: UT-Zappos, MIT-States, and C-GQA. UT-Zappos is a fine-grained dataset composed of 50,025 shoes images with 16 attributes (e.g., Cotton, Nylon), 12 objects (e.g., Shoes.Heels, Boots.Ankle) and 116 compositions. MIT-States contains 53,753 natural images with 115 attributes (e.g., Ancient, Broken), 245 objects (e.g., Computer, Tree) and 1962 compositions. C-GQA is the most extensive dataset containing 39,298 images with 453 attributes, 870 objects and more than 9,500 compositions. Following the standard split, we divide the compositions into *train / validation / test* splits. The detailed splits are shown in Tab. 7.  $|C_s|$  indicates the number of seen compositions,  $|C_u|$  is the number of unseen compositions,  $|C_u|$  is the number of samples in the corresponding splits.

Table 7: Detail of data split statistics.

D-44		Compo	sitions	Tı	rain	Val		Test		
Dataset	$ \mathcal{A} $	$ \mathcal{O} $	$ \mathcal{A} \times  \mathcal{O} $	$ \mathcal{C}_s $	$ \mathcal{X} $	$ \mathcal{C}_s / \mathcal{C}_u $	$ \mathcal{X} $	$ \mathcal{C}_s / \mathcal{C}_u $	$ \mathcal{X} $	
UT-Zappos	16	12	192	83	22998	15 / 15	3214	18 / 18	2914	
MIT-States	115	245	28175	1262	30338	300 / 300	10420	400 / 400	12995	
C-GQA	413	674	278362	5592	26920	1252 / 1040	7280	888 / 923	5098	

**Detailed Evaluation Metrics.** Following the generalized CZSL evaluation protocol (Naeem et al., 2021; Liu et al., 2023), our method is evaluated on both seen and unseen compositions. We report the four widely used metrics for a comprehensive evaluation. Seen Accuracy (S) and Unseen Accuracy (U) are computed to evaluate the best classification performance on seen and unseen compositions. Using Seen Accuracy as *x*-axis and Unseen Accuracy as *y*-axis, we calibrate a bias scalar (Naeem et al., 2021) on Unseen Accuracy and obtain a seen-unseen accuracy curve. Then, we compute and report the Area Under the Curve (AUC). Meanwhile, we compute the best Harmonic Mean (HM) between Seen Accuracy and Unseen Accuracy at a specific bias scalar.

More Implementation Details. For network initialization, we load the weights of CLIP (Radford et al., 2021) and tune the image encoder with LoRA (Zanella & Ben Ayed, 2024). The Sparse Alignment suppresses semantically irrelevant regions to achieve information balance in image-text pairs. The overall pipeline of *Sparse Alignment* is illustrated in Fig. 7. Visual Adaptive condensation module is implemented with K blocks composed of multi-head attention and feed-forward network. The number of blocks K is set to 3, 3 and 1 for UT-Zappos, MIT-States and C-GQA, respectively. The Dynamically Updated Memory Bank does not introduce additional parameters, as the retrieval and prediction processes are calculated on the condensed visual representations without transformation. The coefficient  $\alpha$  for distillation loss in Eq. 8 is set to 0.5, 0.9 and 0.5 for three datasets. The coefficient  $\beta$  in Eq. 12 is set to 0.3, 0.7 and 0.7. The coefficient  $\gamma$  in Eq. 12 is set to 0.5, 0.4 and 0.1. For the number of stored samples in Dynamically Updated Memory Bank is set to 16, 24 and 16. We train our model for 15, 10 and 15 epochs with Adam Optimizer (Kingma & Ba, 2014). The learning rates are initialized at 2e-4, 5e-5 and 5e-4, where the learning rate is scheduled by the StepLR (PyTorch, 2025). During training, we set batch size to 64, 64 and 16 for three datasets. All the experiments are conducted on a single NVIDIA RTX 3090 GPU. More ablation experiments on hyper-parameters is presented in Sec. B.

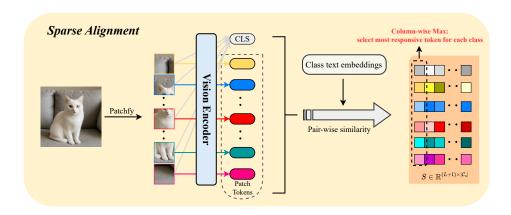


Figure 7: Pipeline of sparse alignment.

### B MORE ABLATION EXPERIMENTS

**More Comparison with SOTA Methods.** Due to space limitations, we report here a more comprehensive comparison of experiments, in which we additionally include more impressive CLIP-based methods.

Method	UT-Zappos					MIT-	State	s	C-GQA			
	S	U	HM	AUC	S	U	HM	AUC	S	U	HM	AUC
CLIP <sub>[ICML'21]</sub> (Radford et al., 2021)	15.8	49.1	15.6	5.0	30.2	46.0	26.1	11.0	7.5	25.0	8.6	1.4
CoOp <sub>[IJCV'22]</sub> (Zhou et al., 2022)	52.1	49.3	34.6	18.8	34.4	47.6	29.8	13.5	20.5	26.8	17.1	4.4
PCVL[Arxiv'22] (Xu et al., 2022)	64.4	64.0	46.1	32.2	48.5	47.2	35.3	18.3	-	-	-	-
HPL <sub>[IJCAI'23]</sub> (Wang et al., 2023a)	63.0	68.8	48.2	35.0	47.5	50.6	37.3	20.2	30.8	28.4	22.4	7.2
CSP <sub>[ICLR'23]</sub> (Nayak et al., 2023)	64.2	66.2	46.6	33.0	46.6	49.9	36.3	19.4	28.8	26.8	20.6	6.2
DFSP <sub>[CVPR'23]</sub> (Lu et al., 2023)	66.7	71.7	47.2	36.0	46.9	52.0	37.3	20.6	38.2	32.0	27.1	10.5
PLID <sub>[ECCV'24]</sub> (Bao et al., 2023)	67.3	68.8	52.4	38.7	49.7	52.4	39.0	22.1	38.8	33.0	27.9	11.0
CDS <sub>[CVPR'24]</sub> (Li et al., 2024)	63.9	74.8	52.7	39.5	50.3	52.9	39.2	22.4	38.3	34.2	28.1	11.1
Troika <sub>[CVPR'24]</sub> (Huang et al., 2024)	66.8	73.8	54.6	41.7	49.0	53.0	39.3	22.1	41.0	35.7	29.4	12.4
CAILA <sub>[WACV'24]</sub> (Zheng et al., 2024)	67.8	74.0	57.0	44.1	51.0	53.9	39.9	23.4	43.9	38.5	32.7	14.8
RAPR <sub>[AAAI'24]</sub> (Jing et al., 2024)	69.4	72.8	56.5	44.5	50.0	53.3	39.2	22.5	45.6	36.0	32.0	14.4
LogiCzsl <sub>[CVPR'25]</sub> (Wu et al., 2025)	69.6	74.9	57.8	45.8	50.8	53.9	40.5	23.4	44.4	39.4	33.3	15.3
ClusPro <sub>[ICLR'25]</sub> (Qu et al., 2025)	70.7	76.0	58.5	46.6	52.1	54.0	40.7	23.8	44.3	37.8	32.8	14.9
SAC	73.3	76.8	62.0	50.0	53.2	53.0	40.8	24.0	45.8	39.5	34.8	16.2

Table 8: The experimental results on closed-world settings.

More Ablation Study on Hyper-Parameteres. We further study the impact of hyper-parameters on performance, including weight coefficient  $\alpha$  in distillation loss Eq. 7, weight coefficient  $\beta$ ,  $\gamma$  in inference Eq. 12 and number of blocks K in VAC.

Influence of Loss Coefficient Weight  $\alpha$ . First, we conduct experiments on  $\alpha$  to investigate the impact of the distillation loss in Eq. 8 on the Visual Adaptive Condensation module and the results are reported in Fig. 8. According to the analysis, we set the  $\alpha$  as 0.5, 0.9 and 0.5 for UT-Zappos, MIT-States and C-GQA, respectively. As we can see, the distillation loss provides a clear performance gain for the VAC module. Ablating this loss (e.g., setting the weight to 0) reduces the VAC objective to a standard classification loss, resulting in notably poorer performance.

Influence of Inference Weight of  $\beta$  and  $\gamma$ . Then, we conduct experiments on inference weight  $\beta$  and  $\gamma$  in Eq. 12 and the results are reported in Fig. 9 and Fig. 10, respectively. We observe that the optimal parameter settings differ across benchmarks. We hypothesize that this arises from varying dataset characteristics, including differences in object or attribute contamination from surrounding regions. Consequently, adjusting the contribution of our modules yields different levels of perfor-

Table 9: The experimental results on open-world settings.

Method	UT-Zappos					MIT-States				C-GQA			
1,100110u	S	U	HM	AUC	S	U	HM	AUC	S	U	НМ	AUC	
CLIP <sub>[ICML'21]</sub> (Radford et al., 2021)	15.7	20.6	11.2	2.2	30.1	14.3	12.8	3.0	7.5	4.6	4.0	0.3	
CoOp <sub>[IJCV'22]</sub> (Zhou et al., 2022)	52.1	31.5	28.9	13.2	34.6	9.3	12.3	2.8	21.0	4.6	5.5	0.7	
PCVL[Arxiv'22] (Xu et al., 2022)	64.6	44.0	37.1	21.6	48.5	16.0	17.7	6.1	-	-	-	-	
HPL <sub>[IJCAI'23]</sub> (Wang et al., 2023a)	63.4	48.1	40.2	24.6	46.4	18.9	19.8	6.9	30.1	5.8	7.5	1.4	
CSP <sub>[ICLR'23]</sub> (Nayak et al., 2023)	64.1	44.1	38.9	22.7	46.3	15.7	17.4	5.7	28.7	5.2	6.9	1.2	
DFSP <sub>[CVPR'23]</sub> (Lu et al., 2023)	66.8	60.0	44.0	30.3	47.5	18.5	19.3	6.8	38.3	7.2	10.4	2.4	
PLID <sub>[ECCV'24]</sub> (Bao et al., 2023)	67.6	55.5	46.6	30.8	49.1	18.7	20.0	7.3	39.1	7.5	10.6	2.5	
CDS <sub>[CVPR'24]</sub> (Li et al., 2024)	64.7	61.3	48.2	32.3	49.4	21.8	22.1	8.5	37.6	8.2	11.6	2.7	
Troika <sub>[CVPR'24]</sub> (Huang et al., 2024)	66.4	61.2	47.8	33.0	48.8	18.7	20.1	7.2	40.8	7.9	10.9	2.7	
CAILA <sub>[WACV'24]</sub> (Zheng et al., 2024)	67.8	59.7	49.4	32.8	51.0	20.2	21.6	8.2	43.9	8.0	11.5	3.1	
RAPR <sub>[AAAI'24]</sub> (Jing et al., 2024)	69.4	59.4	47.9	33.3	49.9	20.1	21.8	8.2	45.5	11.2	14.6	4.4	
LogiCzsl <sub>[CVPR'25]</sub> (Wu et al., 2025)	69.6	63.7	50.8	36.2	50.7	21.4	22.4	8.7	43.7	9.3	12.6	3.4	
ClusPro <sub>[ICLR'25]</sub> (Qu et al., 2025)	71.0	66.2	54.1	39.5	51.2	22.1	23.0	9.3	41.6	8.3	11.6	3.0	
SAC	72.9	66.7	54.8	42.3	52.9	21.2	23.1	9.4	45.5	11.5	15.3	4.6	

mance gain. Therefore, based on our experimental results, we set  $\beta$  as 0.3, 0.7 and 0.7, and set  $\gamma$  as 0.5, 0.4 and 0.1 for UT-Zappos, MIT-States and C-GQA, respectively.

**Influence of Number of Blocks in VAC.** In addition, we report the ablation study for K, number of blocks in *Visual Adaptive Condensation* module. After a comprehensive evaluation, we ultimately set K as 3, 3 and 1 for UT-zappos, MIT-States and C-GQA, respectively. The detailed performance are reported in Tab. 10.

Influence of Number of Stored Samples in Memory Bank. We empirically investigate the impact of the number of stored samples per composition in the memory bank. As illustrated in Tab. 11, we observe that a small memory size leads to suboptimal and unstable performance due to limited sample diversity, and the performance becomes consistent as the memory size increases. However, continually increasing the memory size (e.g., by initializing new slots as zero vectors) may dilute retrieval weights in Eq. 10. Based on this analysis, we set the number of samples to 16, 24, and 16 for datasets UT-Zappos, MIT-States, and C-GQA, respectively, and set the temperature  $\tau_{mb}$  as 0.1 in Eq. 10 to sharpen the weight of effective samples.

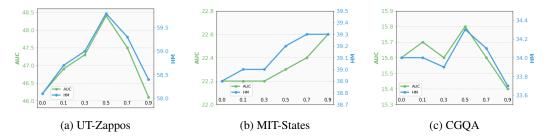


Figure 8: Impact of  $\alpha$  across three datasets.

#### C More Qualitative Experiments

More Qualitative Results. Here, we report more qualitative results in UT-Zappos, MIT-States and C-GQA datasets. As shown in Fig. C, our method can predict accurate results where the baseline makes mistakes. For example, baseline is struggle to distinguish similar objects, e.g., "countertop" and "drawer", "box" and "cooler". Meanwhile, without filtering redundant information, baseline is misled by extraneous visual content, e.g., baseline focuses on object "iron fence", not "calm water". These results demonstrate the effectiveness of our method: by suppressing redundant information, our method is able to make more accurate predictions.

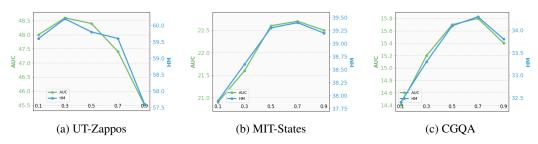


Figure 9: Impact of  $\beta$  across three datasets.

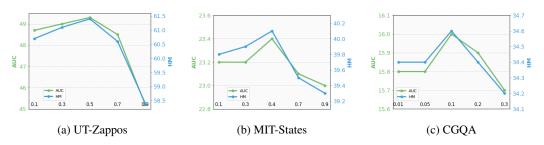


Figure 10: Impact of  $\gamma$  across three datasets.

**More Visualization Results.** As shown in Fig. 12, we present more visualization results of *Visual Adaptive Condensation* (VAC) module in C-GQA dataset. We can observe that our proposed VAC is capable of excavating critical visual information without disturbing by redundant visual cues, such as, "bear" in "green leaf", "wall" in "mess fence" and "cat" in "gray seat", where the main objects are more salient and occupy greater space. These results demonstrate the effectiveness of our proposed VAC.

Table 10: Impact of K in VAC across three datasets.

(a) UT-Zappos	(b) MIT-States	(c) C-GQA					
S U HM AUC	S U HM AUC	S U HM AUC					
K=1   67.6 74.1 57.0 43.7	K=1   50.3 52.1 38.6 22.0	K=1   45.6 38.6 34.1 15.7					
K=3 71.2 76.2 58.9 46.7	K=3 50.6 52.1 39.3 22.4	<i>K</i> =2 45.2 <b>39.1</b> 33.9 15.7					
K=5   69.5 76.0 57.7 45.5	K=5 50.6 51.9 39.1 22.2	K=3 45.7 37.8 33.6 15.4					

# D PSEUDO-CODE

**Training Scheme for SAC.** In this section, we provide a detailed training scheme for our proposed SAC framework, which can be divided into three stages. **Stage I: Sparse Alignment**, we conduct sparse alignment between textual representations and patch visual representations. Leveraging this information-balanced training data, we optimize LoRA (Zanella & Ben Ayed, 2024) for the visual encoder in CLIP. **Stage II: Visual Adaptive Condensation**, with the reduced visual information in the above alignment, the module is guided to adaptively excavate critical visual information within the image, which preserves potential discarded yet valuable information in stage I. **Stage III: Dynamically Updated Memory Bank**, we first initialize memory bank through training data and dynamicall update the memory bank during inference.

# E STATEMENT FOR USING LARGE LANGUAGE MODELS.

In this section, we illustrate the contributions of author contributions and LLMs tools.

Table 11: Impact of N in memory bank across three datasets.

(a) UT-Zappos						(b) N	⁄IIT-S	tates			(c) C-GQA					
	S	U	HM	AUC		S	U	HM	AUC		S	U	HM	AUC		
N=2	72.3	76.2	61.1	48.6	N=4	50.5	52.6	38.7	22.2	N=2	43.9	38.8	33.5	15.2		
N=4	72.6	76.2	61.1	48.8	N=8	51.4	52.6	39.3	22.6	N=4	44.7	38.8	33.6	15.4		
N=8	72.3	76.2	60.9	48.8	N=16	51.5	52.6	39.3	22.6	N=8	44.9	38.8	33.9	15.5		
N=16	73.1	76.2	61.0	49.2	N=24	51.9	52.6	39.2	22.7	N=16	45.1	38.8	34.0	15.6		
N=24	72.9	76.2	60.9	49.0	N=24	51.9	52.6	39.0	22.7	N=24	45.0	38.8	34.0	15.6		



Figure 11: More qualitative results of our method on three datasets.

**Core Contributions** (by the authors): Conception of the Sparse Alignment idea, design the architecture and loss function for modules, design and implementation of all experiments, data analysis, and interpretation of all results.

Assistance from LLMs: In the final stages of manuscript preparation, we used AI tools (ChatGPT/GPT-4 and DeepSeek) for specific, non-intellectual tasks to improve presentation quality. Their use was strictly limited to: 1) Language Polishing: Identifying and correcting typographical, grammatical, and spelling errors. 2) Syntax and Style: Rephrasing sentences for improved readability and academic tone, without altering technical meaning. 3) LaTeX Code Debugging: Ensuring consistency in reference formatting, figure/table labels, and other LaTeX conventions.

The models did not contribute to the scientific ideas, experimental design, or conclusions of this work. The authors reviewed and edited all AI-suggested changes and assume full responsibility for the published content.

979

993 994

995 996

997

998

999

1000

1001

1002

1003

1004

1007

1008

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1023

1024

1025

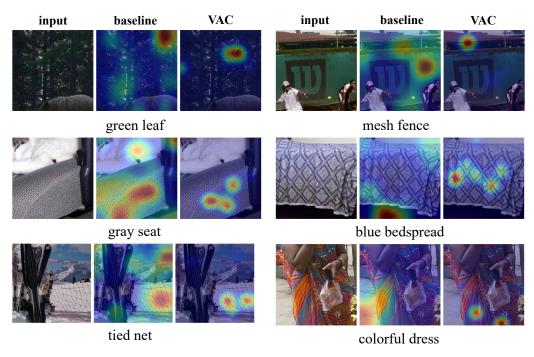


Figure 12: More visualization results of VAC module in C-GQA dataset.

# Algorithm 1 Training Scheme for SAC.

```
Input: training data \mathcal{D}_{tr}, visual encoder of CLIP \phi_{vis}, textual encoder of CLIP \psi_{txt},
learnable soft prompts \theta_t = [\theta_a, \theta_o, \theta_c], visual adaptive condensation module \theta_{vac},
LoRA weight \theta_{LoRA}, memory bank B.
```

**Output:** optimized: LoRA weight  $\theta_{LoRA}$ , learnable soft prompts  $\theta_t = [\theta_a, \theta_o, \theta_c]$ , visual adaptive condensation module  $\theta_{vac}$ ; updated memory bank **B**.

- 1: Stage I:, randomly initialize parameters  $\theta_{LoRA}$ ; load pre-trained parameters visual encoder of CLIP  $\phi_{\text{vis}}$ , textual encoder of CLIP $\psi_{\text{txt}}$ , learnable soft prompts  $[\theta_a, \theta_o, \theta_c]$ .
- 2: while not converged do
- batch of training data  $(\mathcal{X}_b, \mathcal{Y}_b)$ 3:
- 4: conducting sparse alignment by visual reduction in Eq. 2
- calculating basic learning objective  $\mathcal{L}_{base}$  in Eq. 4
- 6: optimize parameters  $\theta$  ( $\theta_{LoRA}$ ,  $\theta_t$ ) =  $\theta - \nabla_{\theta}(\mathcal{L}_{base}(\mathcal{X}_b, \mathcal{Y}_b; \theta))$
- 7: end while
- 8: **Stage II:** randomly initialize parameters  $\theta_{VAC}$ .
- 9: while not converged do
- batch of training data  $(\mathcal{X}_b, \mathcal{Y}_b)$ 10:
- 11: condense visual information within image into  $v_q$
- 12: calculation prediction  $p_{vac}$  of VAC by Eq. 5 and  $p_{sa}$  of SA by Eq. 3
- 13:
- calculating learning objective  $\mathcal{L}_{base}^{vac}$  in Eq. 6 and  $\mathcal{L}_{kl}$  in Eq. 7 optimize parameters  $\boldsymbol{\theta}\left(\boldsymbol{\theta}_{vac}\right) = \boldsymbol{\theta} \nabla_{\boldsymbol{\theta}}((1-\alpha) \cdot \mathcal{L}_{base}^{vac}(\mathcal{X}_b, \mathcal{Y}_b; \boldsymbol{\theta}) + \alpha \cdot \mathcal{L}_{kl}(\mathcal{X}_b, \mathcal{Y}_b; \boldsymbol{\theta})$ 14:
- 15: end while
- 16: **Stage III:** initialize stored samples for seen compositions in memory bank **B** by Eq. 9.
- 17: **for** batch of testing data  $\mathcal{X}_b$  **do**
- calculating predictions  $p_{sa}$ ,  $p_{vac}$  and  $p_{bank}$  from three modules by Eq. 3, Eq. 5 and Eq. 10, respectively
- 19: obtain final prediction by Eq. 12
  - 20: utilizing  $p_{vac}$  to update memory bank by Eq. 9
- 21: **end for** 
  - 22: calculating the results of each evaluation metric with the final prediction
  - 23: **return** optimized LoRA weight  $\theta_{LoRA}$ , learnable soft prompts  $\theta_t = |\theta_a, \theta_o, \theta_c|$ , visual adaptive condensation module  $\theta_{vac}$ ; updated memory bank **B**.