

VILA: Learning Image Aesthetics from User Comments with Vision-Language Pretraining

CVPR 2023

● Introduction

- Humanlabeled rating scores simplifies the process of aesthetic perception.

User comments provide more comprehensive information and human are better at expressing aesthetic preferences through natural language.



we propose learning image aesthetics from **user comments**, and exploring vision-language pretraining methods to learn **multimodal aesthetic representations**

● CLIP-IQA



Look

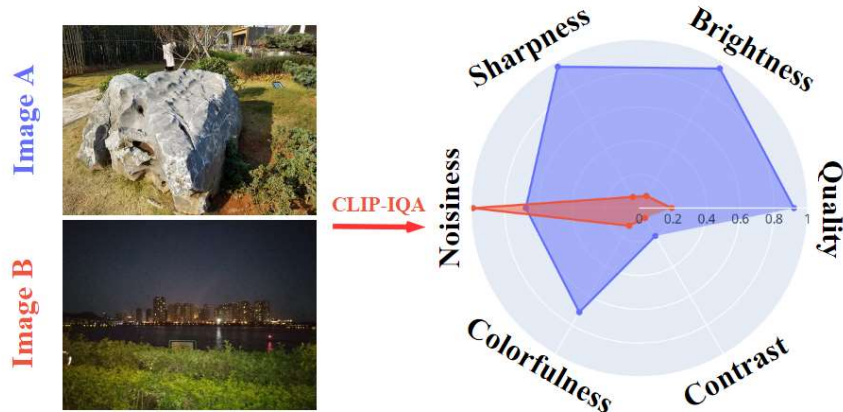
Brightness, Colorfulness, Noise

Feel

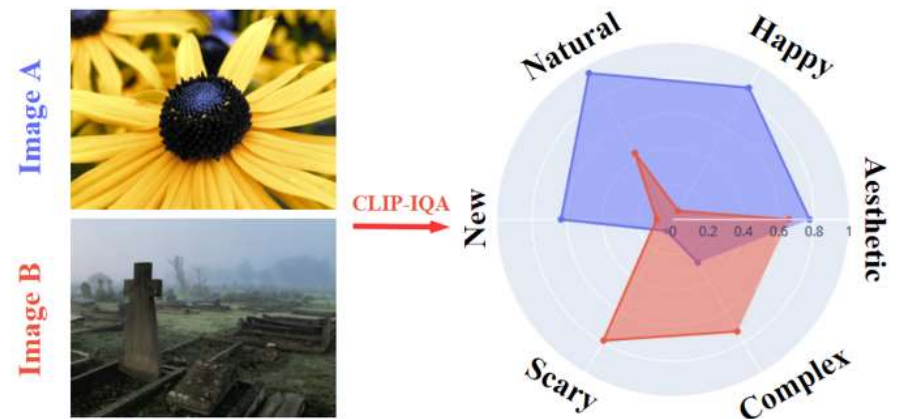
Emotion, Aesthetic

quality perception

abstract perception



(a) CLIP for quality perception

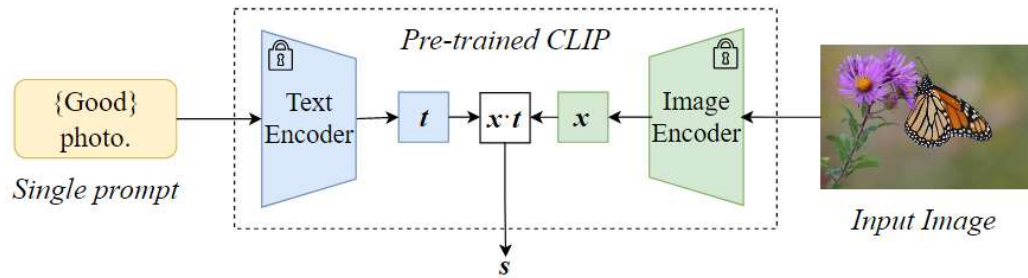


(b) CLIP for abstract perception

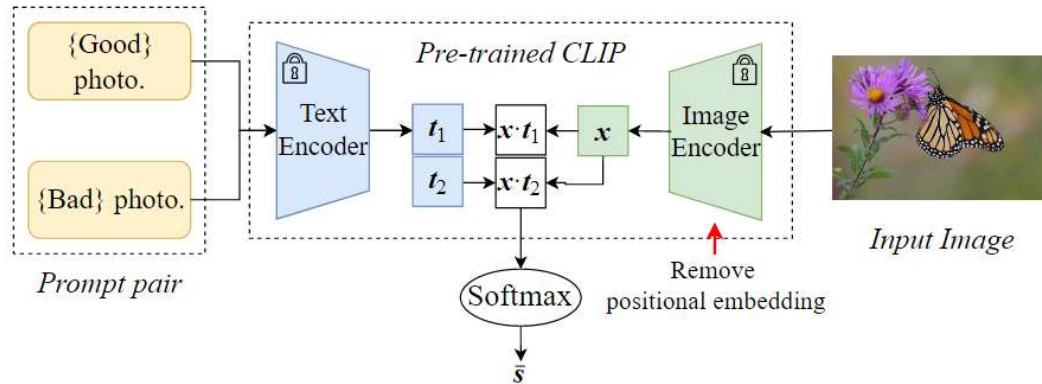
Wang J, Chan K C K, Loy C C. Exploring clip for assessing the look and feel of images[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2023, 37(2): 2555-2563.

CLIP-IQA

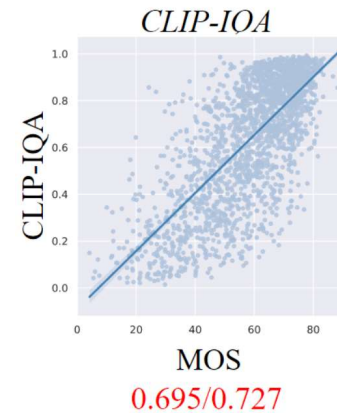
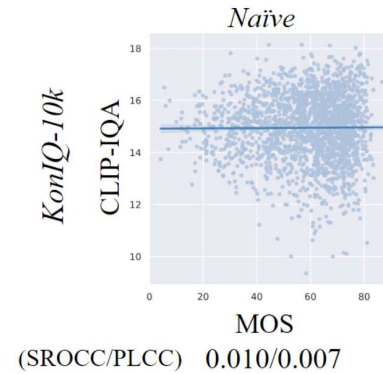
(a) Naïve way to adopt CLIP for visual perception



(b) Framework for the proposed CLIP-IQA



(c) CLIP-IQA score vs. MOS score



➤ Antonym prompt pairing

$$s = \frac{x \odot t}{\|x\| \cdot \|t\|},$$

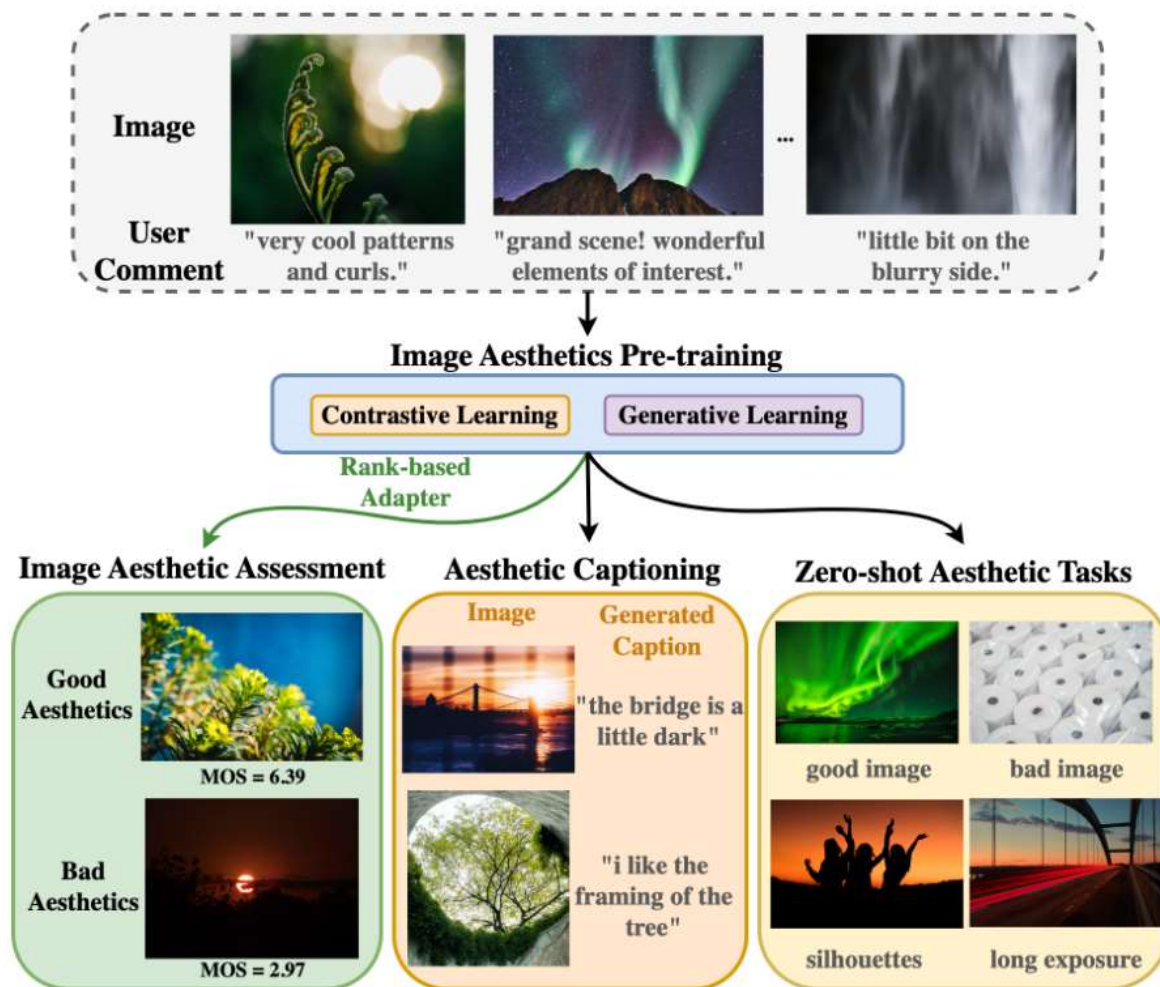


$$s_i = \frac{x \odot t_i}{\|x\| \cdot \|t_i\|}, \quad i \in \{1, 2\},$$

$$\bar{s} = \frac{e^{s_1}}{e^{s_1} + e^{s_2}}.$$

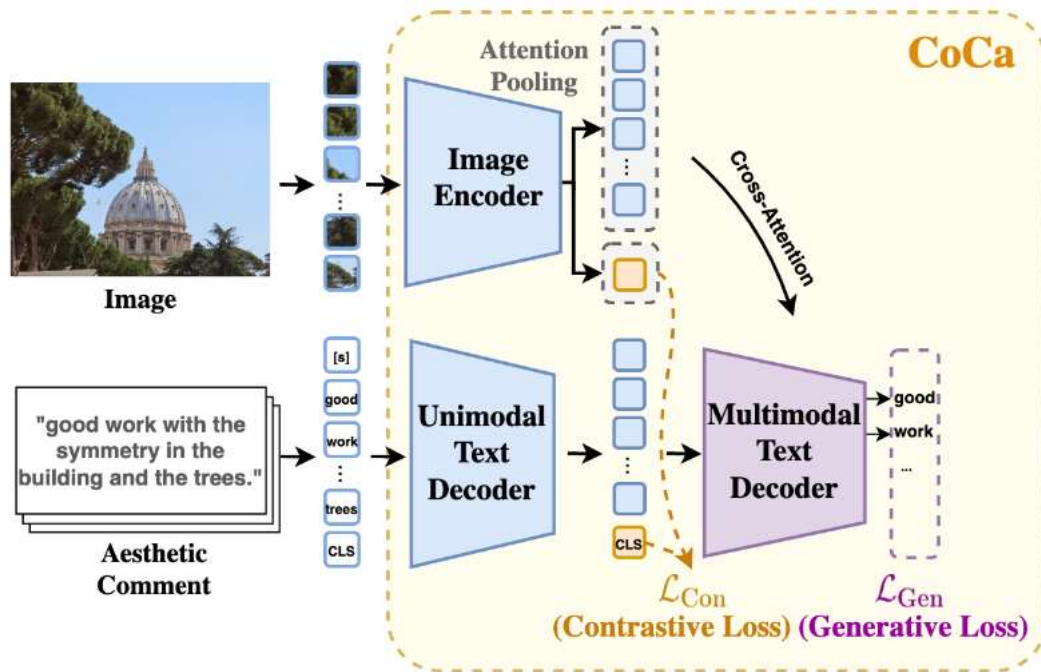
Wang J, Chan K C K, Loy C C. Exploring clip for assessing the look and feel of images[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2023, 37(2): 2555-2563.

● VILA

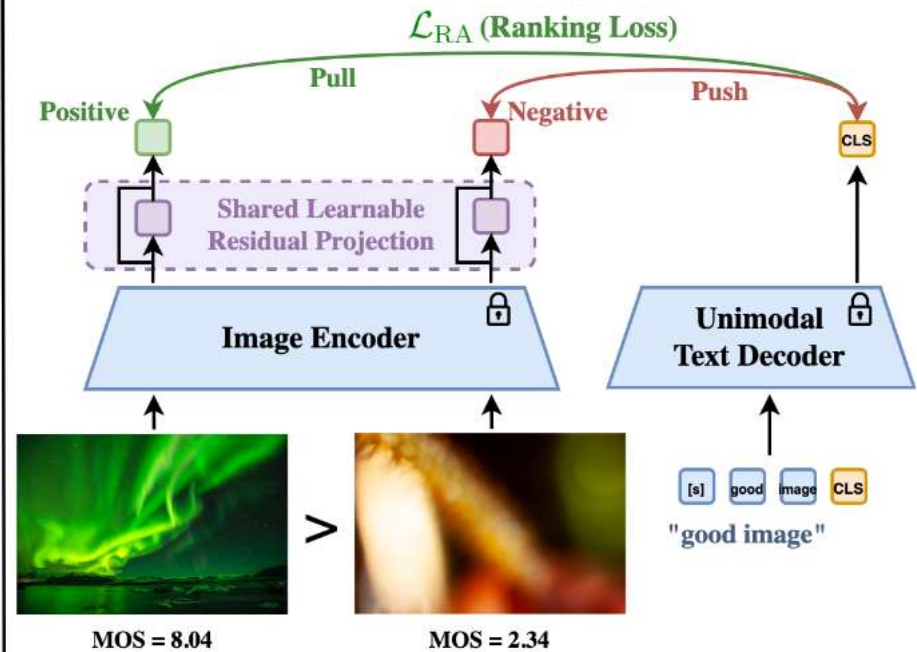


● Method

(1) VILA-P: Vision-Language Aesthetics Pretraining

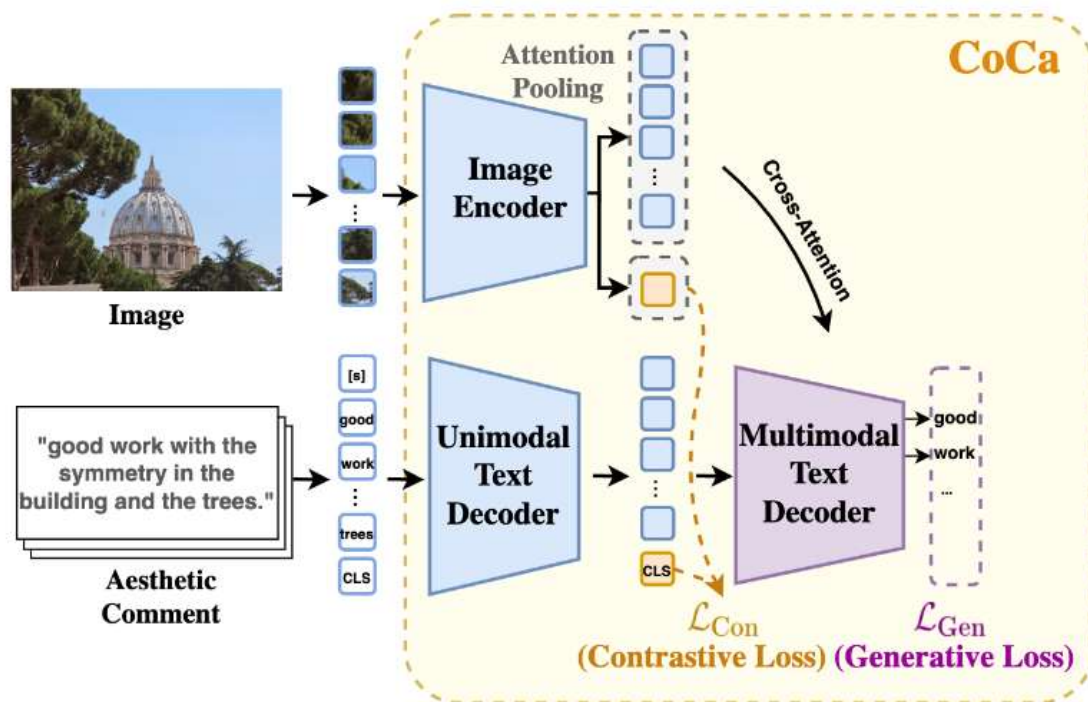


(2) VILA-R: Rank-based Adapter for IAA



● VILA-P

(1) VILA-P: Vision-Language Aesthetics Pretraining



Preliminary of CoCa

Contrastive Learning Objective

$$\mathcal{L}_{\text{Con}}^{i2t} = -\frac{1}{N} \left(\sum_i \log \frac{\exp(\mathbf{x}_i^\top \mathbf{y}_i) / \tau}{\sum_{j=1}^N \exp(\mathbf{x}_i^\top \mathbf{y}_j / \tau)} \right)$$

$$\mathcal{L}_{\text{Con}}^{t2i} = -\frac{1}{N} \left(\sum_i \log \frac{\exp(\mathbf{y}_i^\top \mathbf{x}_i) / \tau}{\sum_{j=1}^N \exp(\mathbf{y}_i^\top \mathbf{x}_j / \tau)} \right)$$

$$\mathcal{L}_{\text{Con}} = \mathcal{L}_{\text{Con}}^{i2t} + \mathcal{L}_{\text{Con}}^{t2i}$$

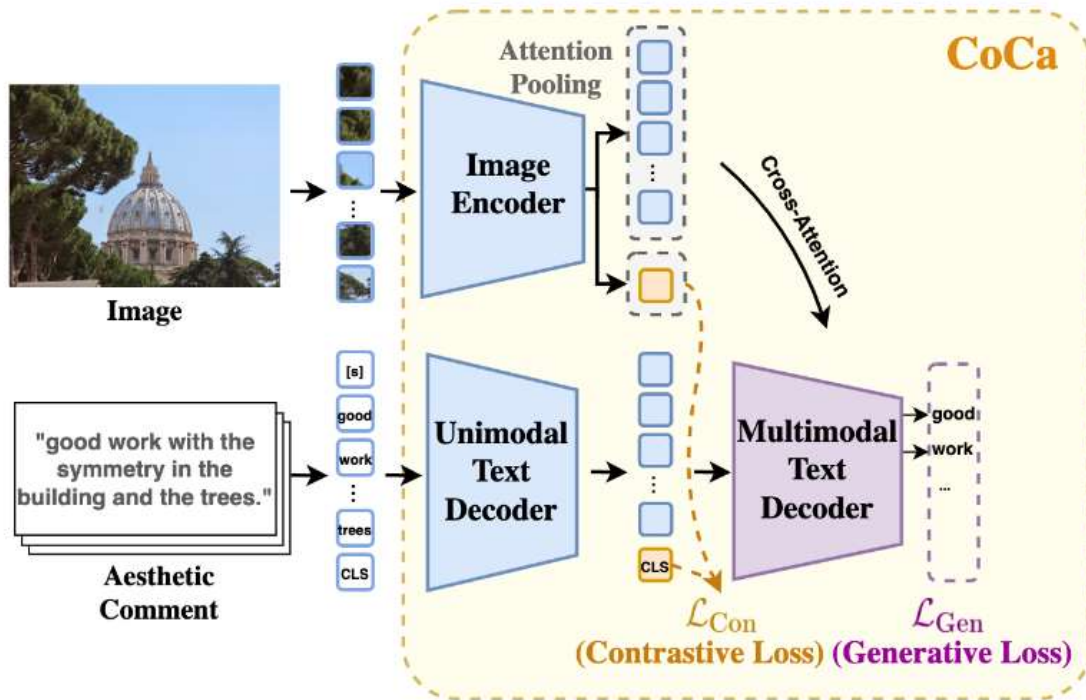
Generative Learning Objective

$$\mathcal{L}_{\text{Gen}} = -\sum_{t=1}^L \log P(w_t | w_{<t}, V).$$

$$\mathcal{L} = \alpha \mathcal{L}_{\text{Con}} + \beta \mathcal{L}_{\text{Gen}}.$$

● VILA-P

(1) VILA-P: Vision-Language Aesthetics Pretraining



Two-stage pretraining approach

➤ General Pretraining

LAION-5B-English dataset

650M images – text pairs

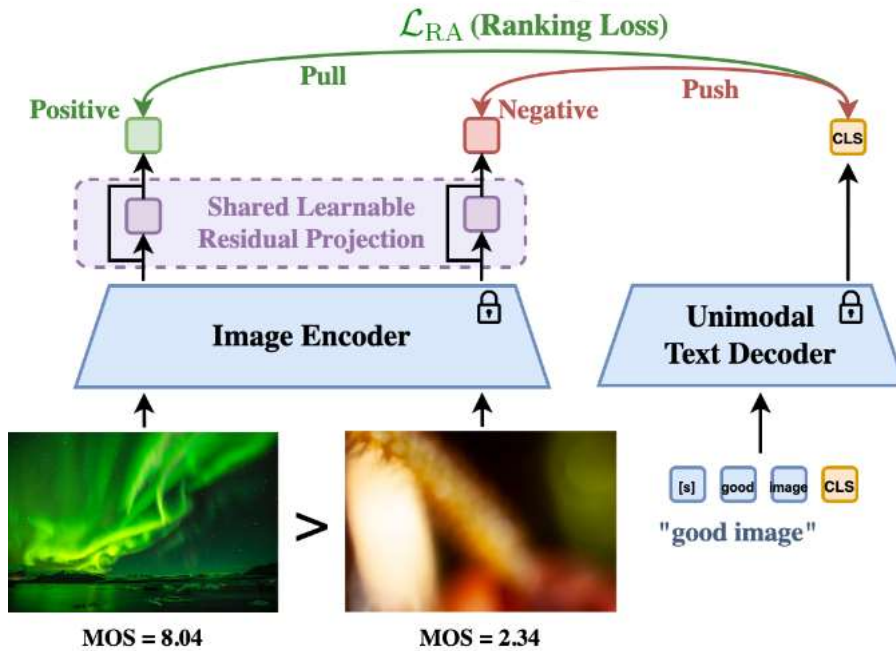
➤ Aesthetic Pretraining

AVA-Captions dataset

230k images + 1.5M captions

● VILA-R

(2) VILA-R: Rank-based Adapter for IAA



IAA Formulation

$$v = E(I, \theta_{frozen}),$$

$$r = F(v, \gamma),$$

Rank-based Adapter

$$\tilde{v} = \text{normalize}(v^\top H + v),$$

$$r = \tilde{v}^\top w_p$$

triplet ranking loss:

$$\mathcal{L}_{RA} = \frac{1}{P} \sum_{i,j, i \neq j, l_i > l_j} \max(0, m - \tilde{v}_i^\top w_p + \tilde{v}_j^\top w_p)$$

● Experiments

Method	SRCC	PLCC
Kong <i>et al.</i> [24]	0.558	-
NIMA (Inception-v2) [42]	0.612	0.636
AFDC + SPP [2]	0.649	0.671
MaxViT [45]	0.708	0.745
AMP [31]	0.709	-
Zeng <i>et al.</i> (resnet101) [55]	0.719	0.720
MUSIQ [19]	0.726	0.738
Hentschel <i>et al.</i> [13]	0.731	0.741
Niu <i>et al.</i> [33]	0.734	0.740
MLSP (Pool-3FC) [15]	0.756	0.757
TANet [12]	0.758	0.765
GAT _{×3} -GATP [11]	0.762	0.764
Zero-shot Learning		
VILA-P (single prompt)	0.605	0.617
VILA-P (ensemble prompts)	0.657	0.663
VILA-R	0.774	0.774

Table 1. Results on AVA dataset. **Blue** and **black** numbers in bold represent the best and second best respectively. First group shows baselines, second group shows ZSL results using our model from Sec. 3, final line shows our result combining Sec. 3 and Sec. 4.

	Prompts	
	p_g	p_b
Single Prompt	“good image”	“bad image”
Ensemble of Prompts	“good image”	“bad image”
	“good lighting”	“bad lighting”
	“good content”	“bad content”
	“good background”	“bad background”
	“good foreground”	“bad foreground”
	“good composition”	“bad composition”

Table 11. Text prompts used in ZSL for IAA.

● Experiments

➤ Effects of image-text pretraining rank-based adapter

	ZSL Ens. Prompts			w/ Our Adapter		
General Pretraining	✓		✓	✓		✓
Aesthetic Pretraining		✓	✓		✓	✓
SRCC	0.228	0.265	0.657	0.746	0.566	0.774
PLCC	0.228	0.276	0.663	0.750	0.575	0.774

Table 2. Effects of image-text pretraining on AVA. Different pre-training schema are employed for each column and two settings are reported: 1) ZSL using an ensemble of prompts; 2) further finetuned using our proposed rank-based adapter.

Method	SRCC	PLCC
VILA-P w/ L2 Loss	0.757	0.756
VILA-P w/ EMD Loss [42]	0.759	0.759
VILA-R w/o Text Anchor	0.763	0.764
VILA-R w/o Residual	0.766	0.766
VILA-R (Ours)	0.774	0.774
VILA-R Finetune Image Encoder	0.780	0.780

Table 3. Ablation for the proposed rank-based adapter (Sec. 4) on AVA. First two groups use frozen pretrained image encoder.

Experiments

Zero-shot Aesthetic Style Classification

Method	mAP (%)
Murray <i>et al.</i> [32]	53.9
Karayev <i>et al.</i> [18]	58.1
Lu <i>et al.</i> [29]	64.1
MNet [41]	65.5
Sal-RGB [9]	71.8
Zero-shot Learning	
General Pretraining (single prompt)	29.3
General Pretraining (ensemble prompts)	32.6
VILA-P (single prompt)	62.3
VILA-P (ensemble prompts)	69.0

AVA Comments Generation

Method	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE	CIDEr
CWS [10]	0.535	0.282	0.150	0.074	0.254	0.059
Yeo <i>et al.</i> [52]	0.464	0.238	0.122	0.063	0.262	0.051
VILA	0.503	0.288	0.170	0.113	0.262	0.076

Table 5. Results on AVA-Captions dataset.



Figure 5. Aesthetic comments generated by VILA.



Thanks!