

# CLAP: Isolating Content from Style through Contrastive Learning with Augmented Prompts

Yichao Cai (✉), Yuhang Liu, Zhen Zhang, and Javen Qinfeng Shi

Australian Institute for Machine Learning (AIML), The University of Adelaide, SA 5000, Australia

## Motivation

**CLIP** - Trained on massive image-text data with a symmetric InfoNCE loss.

### Pros

- ✓ Mitigated modality gap
- ✓ High zero-shot performance
- ✓ Great generalization ability

### Cons

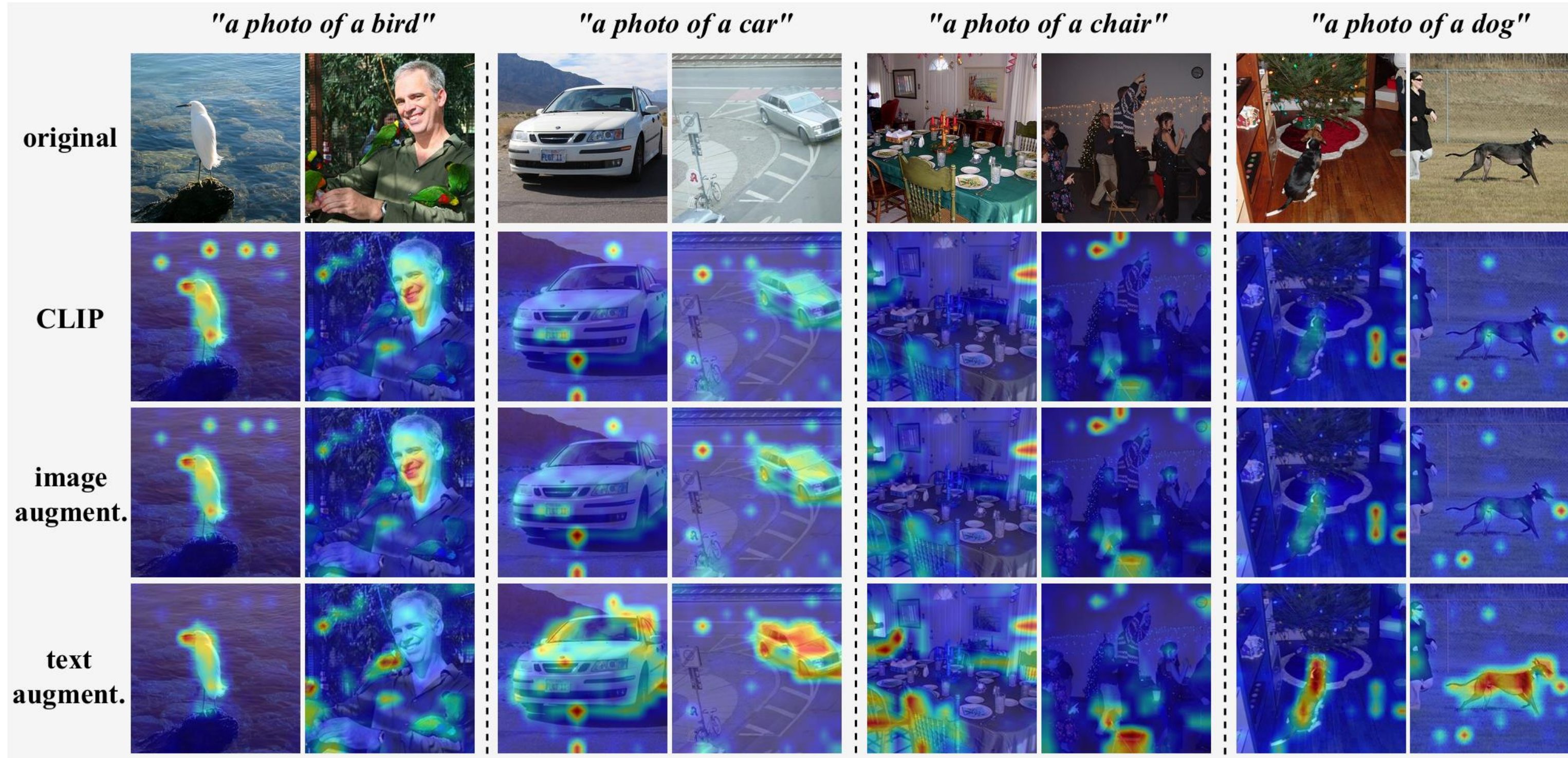
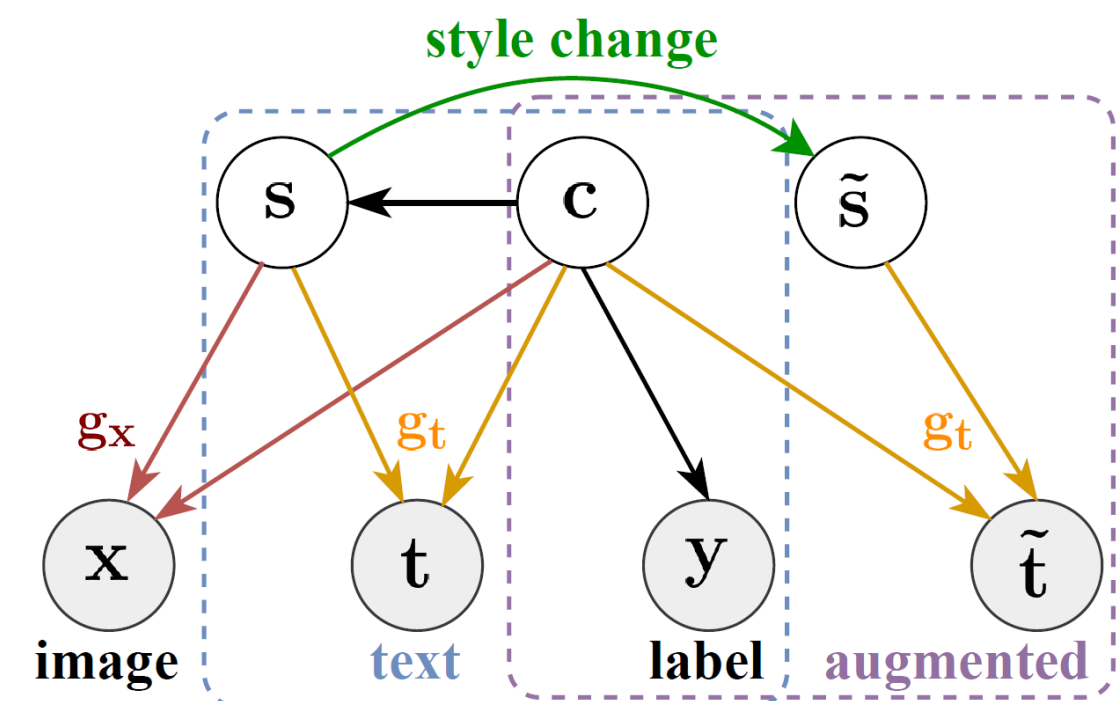
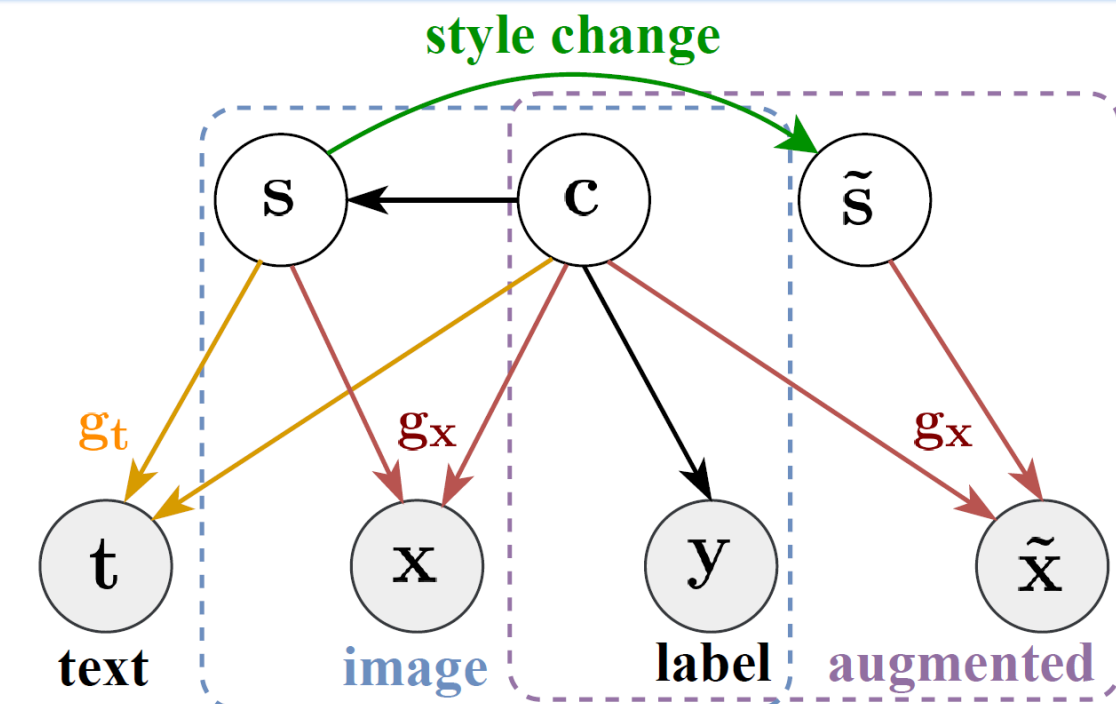
- Zero-shot sensitive to prompts
- Few-shot performance degrades
- Vulnerable to adversarial attacks

**Intuition** - There exists spurious correlations in CLIP features, i.e., style-related information is erroneously used to predict class labels.

## Contribution

- We propose a contrastive learning method to disentangle content and style in pretrained CLIP-like models.
- Our disentangled network, trained on either the image or text encoder, can be seamlessly applied to both modalities.
- Leveraging the high semantic structure of text data, we introduce **CLAP** (Contrastive Learning with Augmented Prompts) to isolate content features within the pre-trained CLIP feature space.

## Causal Generative Models of Vision-Language Data

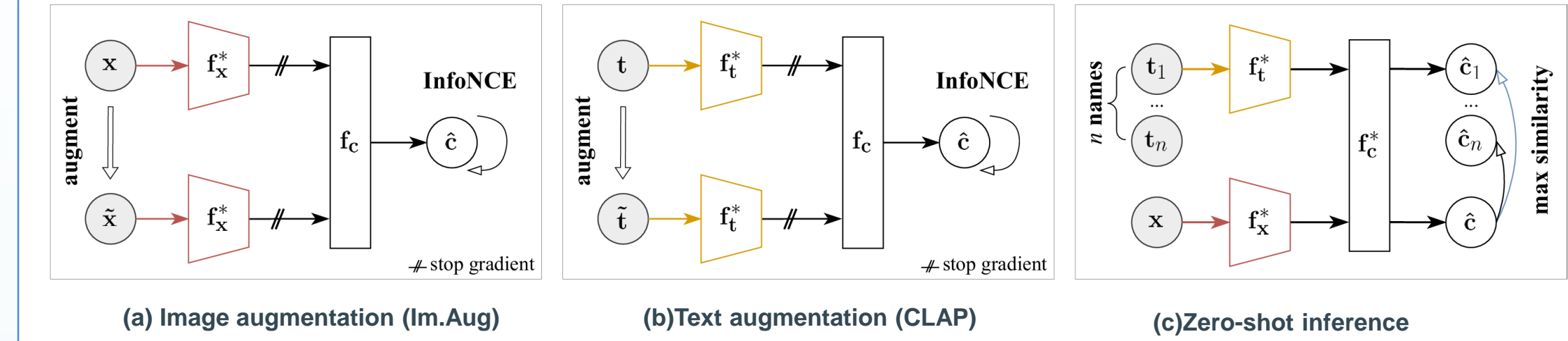


**CAM Visualization:** Comparative visualization of features from CLIP, Image augmentation (Im.Aug), and Text augmentation (CLAP) for zero-shot inference.

**Insight** – Image  $x$  and text  $t$ , derived from a unified latent space with content  $c$  and style  $s$ , follow distinct deterministic processes,  $g_x$  and  $g_t$ . The class label  $y$  is determined solely by the latent content. (a) Soft interventions on style variables generate augmented images  $\tilde{x}$ . (b) Similar interventions produce augmented text  $\tilde{t}$  due to the shared latent space.

## Isolating Content through Data Augmentation

### Framework



**Innovation** - Contrastive learning with data augmentation in one modality benefits both, with text data being more amendable for style changes due to its semantic structure. The trained adapting network can be seamlessly applied in both modality for zero-shot inference.

### Contrastive Learning with Augmented Images (Im.Aug)

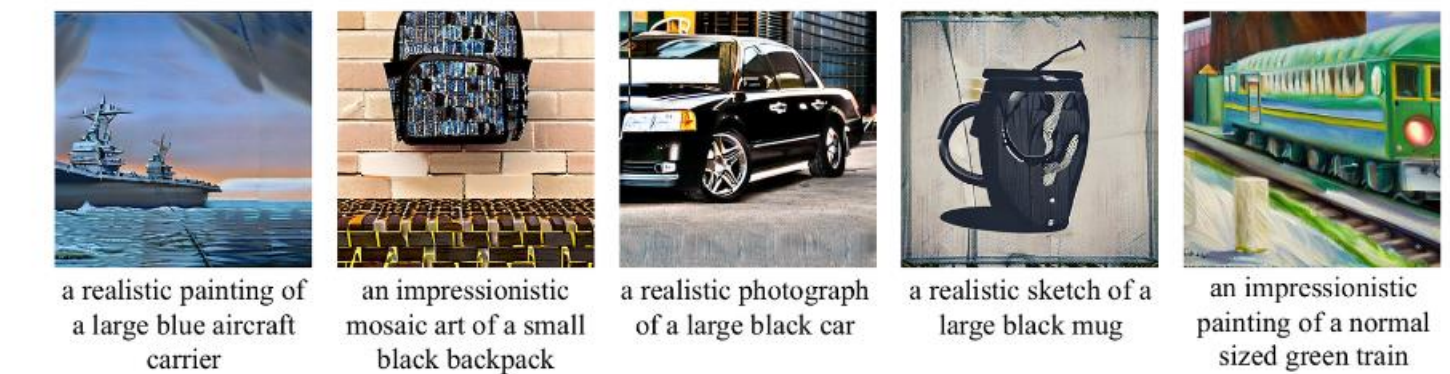
#### Template Prompts

"an [Art Style] [Image Type] of an [Object Size] [Object Color] [Class]"

Art Style	Image Type	Object Size	Object Color
realistic, impressionistic	painting, cartoon, infographic, sketch, photograph, clipart, mosaic art, sculpture	large, small, normal sized	yellow, green, black, blue, multicolored, orange, red, white, brown, purple

Text-to-image diffusion model

#### Synthetic Images



#### Image Augmentation

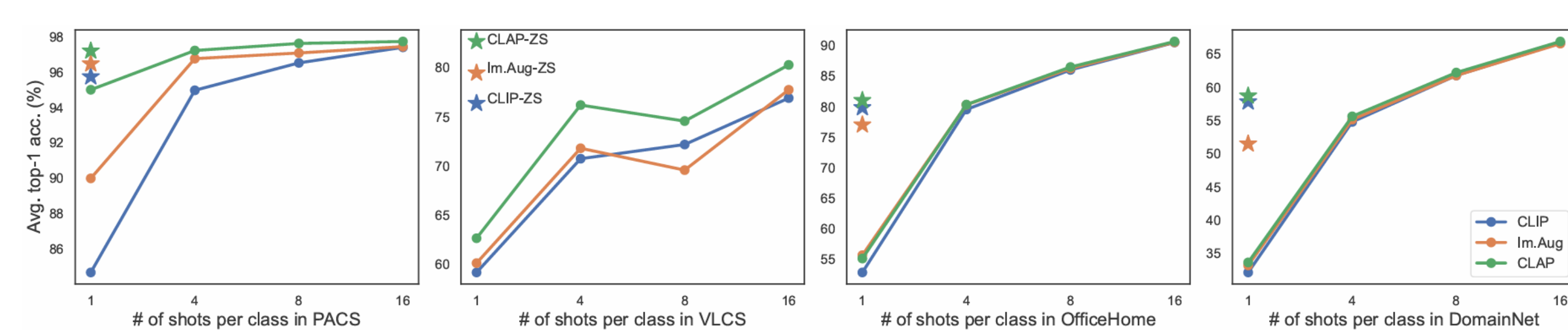
Random Cropping + Color distortion

## Experimental Results

### Zero-shot performance

Prompt Method	Zero-shot performance, avg. top-1 acc. (%) (↑)				
	PACS	VLCS	Off.Home	Dom.Net	Overall
ZS(C)					
CLIP	95.7	76.4	79.8	57.8	77.4
Im.Aug	96.5	79.5	77.0	51.5	76.1
CLAP	<b>97.2</b>	<b>82.6</b>	<b>81.0</b>	<b>58.7</b>	<b>79.9</b>
ZS(CP)					
CLIP	95.2	82.0	79.5	57.0	78.4
Im.Aug	96.3	82.9	75.8	50.7	76.4
CLAP	<b>97.3</b>	<b>83.4</b>	<b>80.5</b>	<b>58.0</b>	<b>79.8</b>
ZS(PC)					
CLIP	96.1	82.4	82.5	57.7	79.7
Im.Aug	96.5	83.0	78.6	51.6	77.4
CLAP	<b>97.2</b>	<b>83.4</b>	<b>83.0</b>	<b>59.0</b>	<b>80.6</b>
ZS(NC)					
CLIP	90.8	68.3	71.5	51.0	70.4
Im.Aug	94.8	73.1	67.5	44.0	69.9
CLAP	<b>97.2</b>	<b>81.0</b>	<b>73.5</b>	<b>52.6</b>	<b>76.1</b>

### Few-shot performance



### Adversarial robustness

Setting Method		Avg. top-1 acc. (%) under adversarial attacks(↑)											
		FGSM				PGD-20				CW-20			
		PACS	VLCS	O.H.	D.N.	PACS	VLCS	O.H.	D.N.	PACS	VLCS	O.H.	D.N.
ZS(C)	CLIP	86.8	65.6	57.9	22.5	29.1	2.0	10.1	10.7	27.4	1.5	7.4	7.6
	Im.Aug	88.0	69.6	55.1	37.9	<b>31.3</b>	2.1	10.4	9.0	29.4	1.7	7.0	5.8
	CLAP	<b>88.7</b>	<b>71.9</b>	<b>58.5</b>	<b>44.2</b>	30.8	<b>3.2</b>	<b>10.6</b>	<b>11.2</b>	<b>29.8</b>	<b>2.3</b>	<b>8.1</b>	<b>8.0</b>
1-shot	CLIP	66.7	45.2	34.3	22.5	34.8	16.0	5.6	11.3	18.9	0.7	4.5	3.2
	Im.Aug	79.4	47.1	<b>37.1</b>	23.5	55.2	16.1	<b>8.5</b>	<b>12.5</b>	23.2	0.9	<b>5.1</b>	3.4
	CLAP	<b>89.6</b>	<b>52.2</b>	<b>37.1</b>	<b>23.9</b>	<b>73.4</b>	<b>21.2</b>	7.4	<b>12.5</b>	<b>27.0</b>	<b>1.1</b>	5.0	<b>3.5</b>

- CLAP successfully refines pretrained CLIP features by effectively isolating content from style, as demonstrated by improved zero-shot and few-shot performance, increased robustness against adversarial attacks, and enhanced qualitative visualizations.

- Exploring the more efficient data augmentation techniques, including the combination of augmentations from both modalities could be valuable, given their complementary strengths.



Scan me!

### Contrastive Learning with Augmented Prompts (CLAP)

#### Template Prompts

"an [Art Style] [Image Type] of an [Object Size] [Object Color] [Class]"

e.g., "a realistic painting of a large red car"

#### Text augmentation

Randomly apply the combination of these techniques:

Object Size Deletion (OSD)	Object Color Deletion (OCD)	Image Type Deletion (ITD)	Art Style Deletion (ASD)	Swapping Prompt Order (SPO)	Inserting Gaussian Noise (IGN)
a realistic painting of a red car	a realistic painting of a large car	a realistic of a large red car	a painting of a large red car	a large red car in a realistic painting	[noise] a realistic painting of a large red car

Metric Method		Performance variance, avg. top-1 acc. (%) (↓)				
		PACS	VLCS	Off.Home	Dom.Net	Overall
$R$	CLIP	0.9	6.1	3.1	<b>0.8</b>	2.7
	Im.Aug	<b>0.1</b>	3.6	2.8	0.9	1.9
	CLAP	<b>0.1</b>	<b>0.8</b>	<b>2.5</b>	1.0	<b>1.1</b>
$\delta$	CLIP	0.4	2.8	1.4	<b>0.4</b>	1.2
	Im.Aug	0.1	1.7	1.2	<b>0.4</b>	0.8
	CLAP	<b>0.0</b>	<b>0.4</b>	<b>1.1</b>	<b>0.4</b>	<b>0.5</b>
$\Delta_{(NC)}$	CLIP	4.9	8.1	8.3	6.8	7.0
	Im.Aug	1.6	6.4	9.5	7.5	6.3
	CLAP	<b>0.0</b>	<b>1.6</b>	<b>7.5</b>	<b>6.1</b>	<b>3.8</b>

Inference prompts: ZS(C) – "[class]" ZS(CP) – "a [class] in a photo" ZS(PC) – "a photo of a [class]" ZS(NC) – "[Gaussian noise][class]"