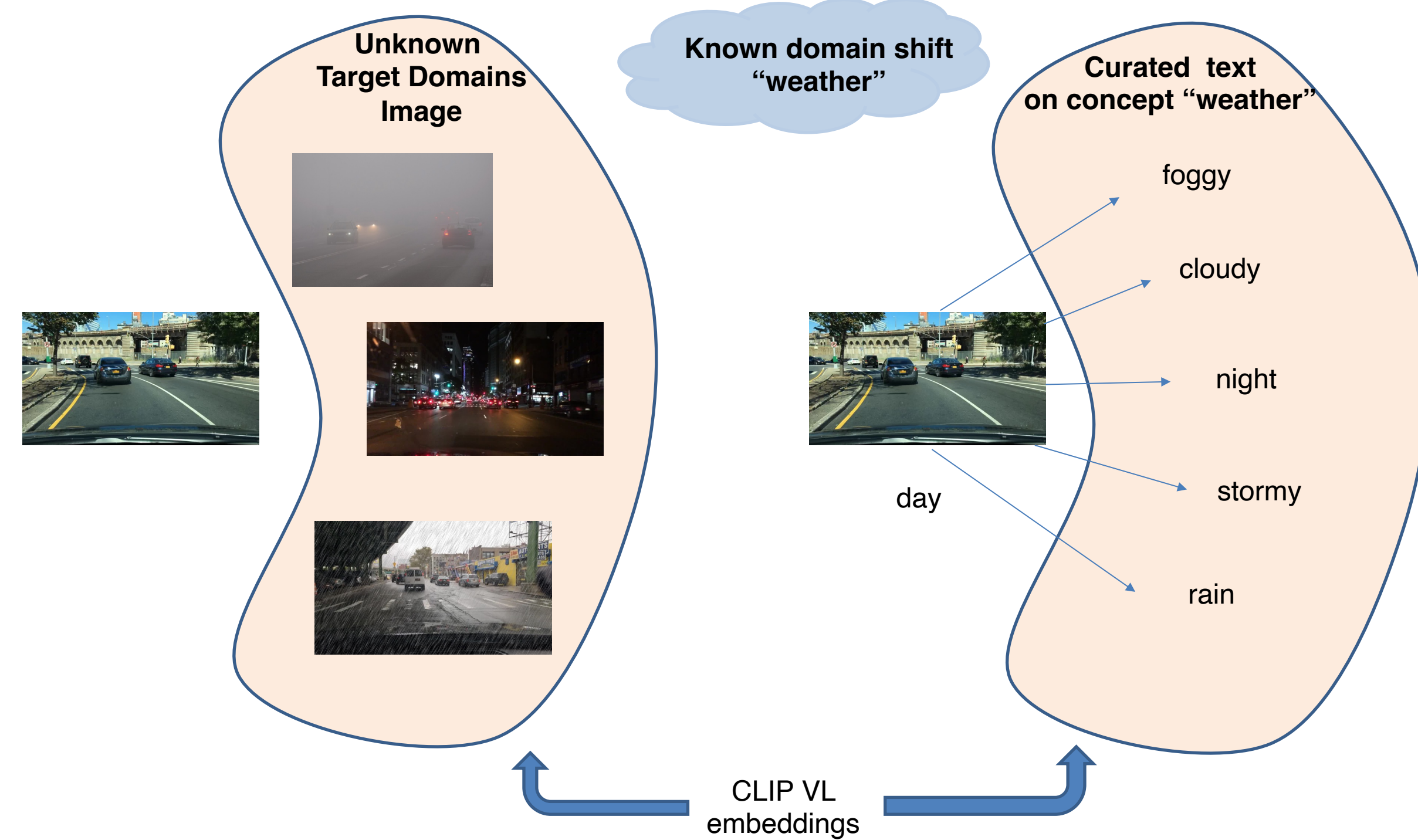


Introduction

- ❖ SDG aims for generalization to new domains
- ❖ One source domain is available
- ❖ Domain shifts can be broadly defined via text prompts
- ❖ Vision-Language models like CLIP [1] can be helpful

Text Prompts are used to estimate target domain features

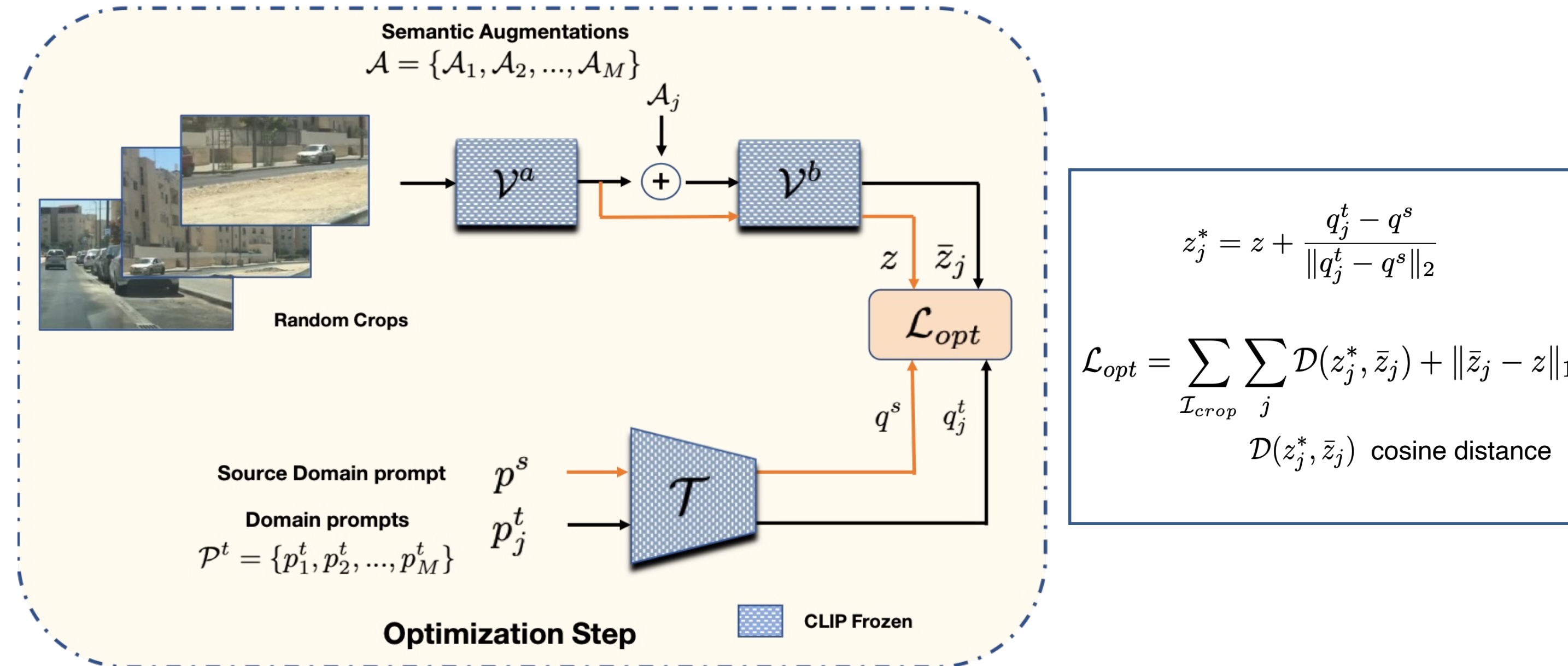


Contribution

- ❖ Leveraging vision - text aligned embeddings
- ❖ Use textual domain prompts to generate semantic augmentations
- ❖ Train-time only augmentation

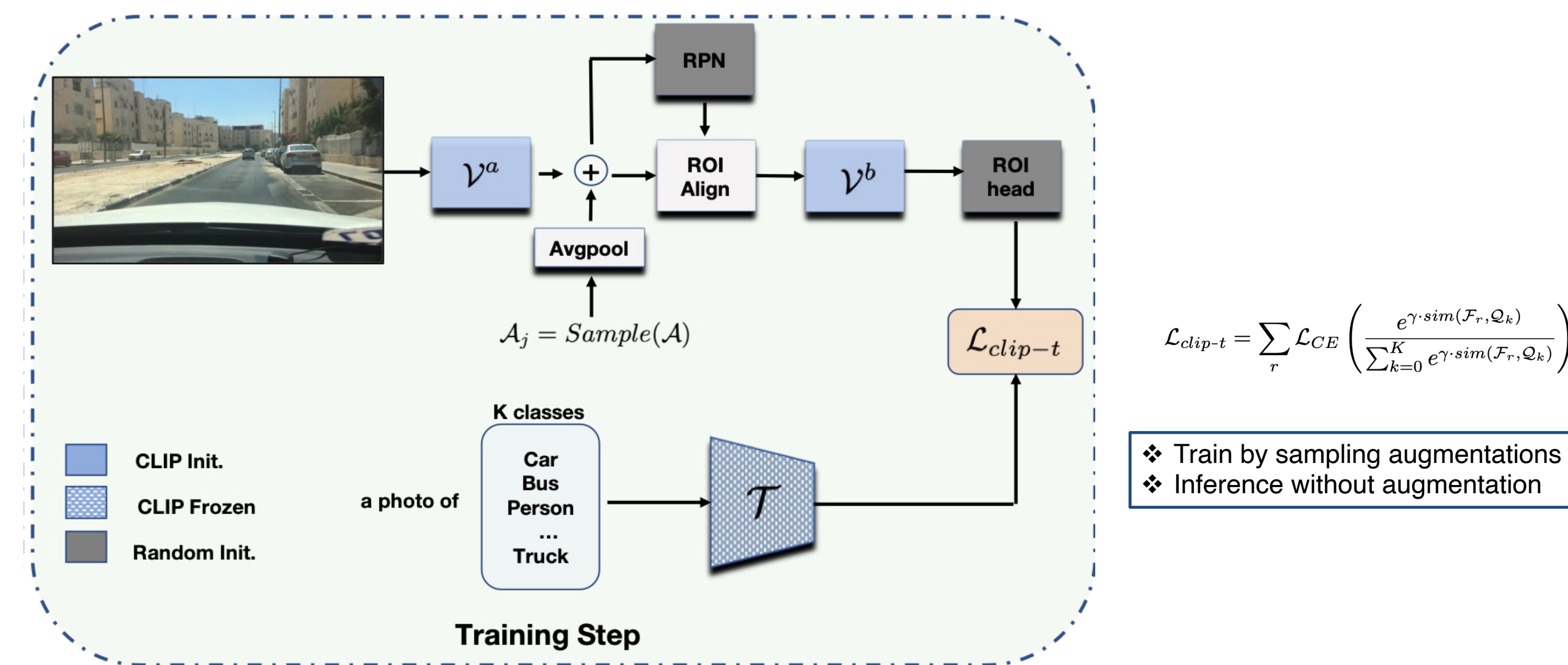
Method

Step 1.



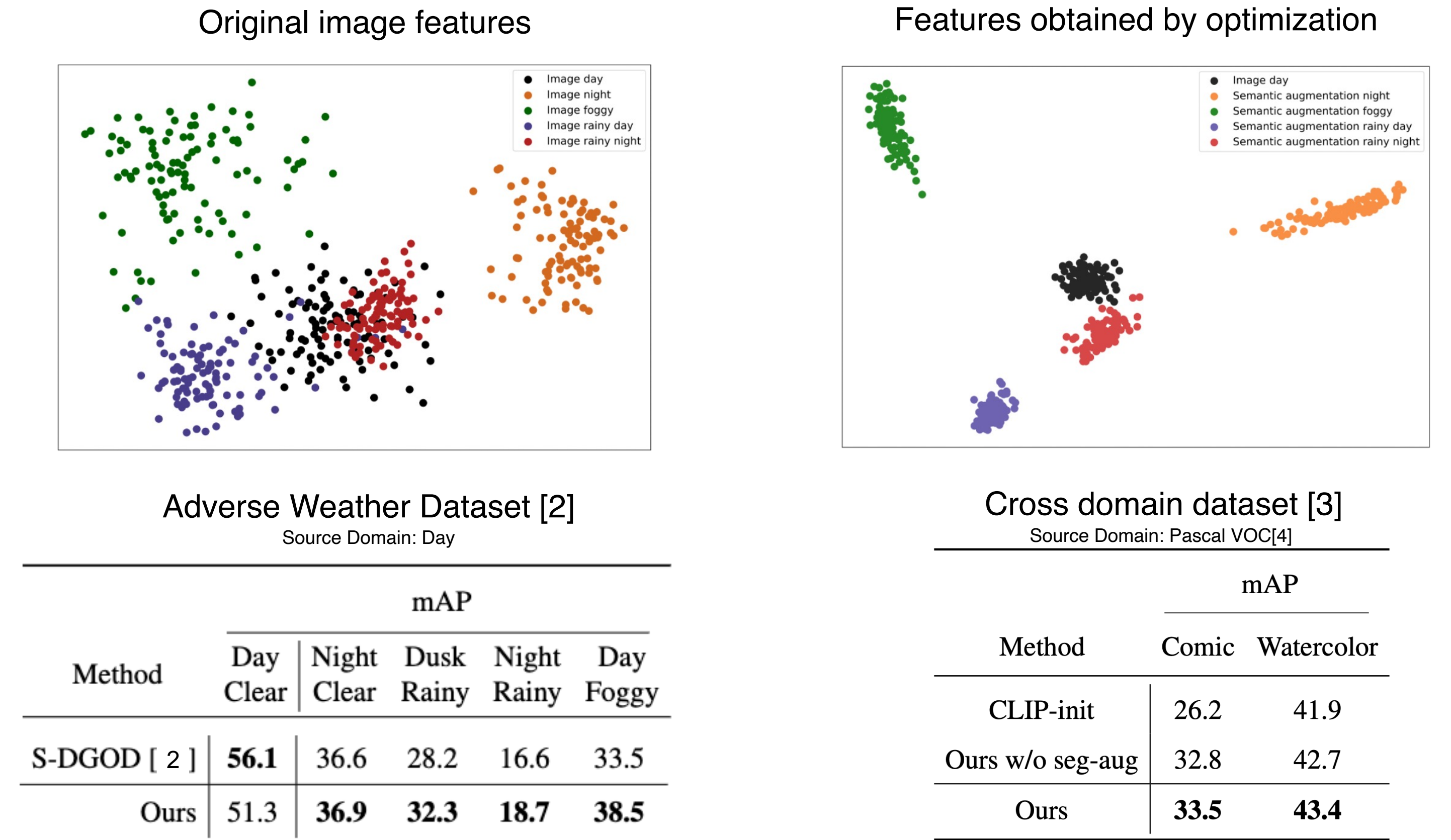
- ❖ We curate a set of textual prompts \mathcal{P}^t which can represent new domains
- ❖ For *weather* dataset, we can create prompts \mathcal{P}^t as an image taken on a {weather} {time of the day.}
- ❖ Source domain prompt p^s for the same is an image taken during the day.

Step 2.



- ❖ Train by sampling augmentations
- ❖ Inference without augmentation

Results



Ablations

Model Component	mAP							
	Source		Target					
CLIP init	\mathcal{L}_{clip-t}	Attn. Pool	Sem. Aug	Day Clear	Night Clear	Dusk Rainy	Night Rainy	Day Foggy
				48.1	34.4	26.0	12.4	32.0
✓				51.2	37.0	31.0	15.7	37.5
✓	✓			50.7	36.0	31.3	16.3	36.9
✓	✓	✓		51.0	35.9	31.3	16.7	37.7
✓	✓	✓	✓	51.3	36.9	32.3	18.7	38.5

Aug. Type	mAP				
	Day Clear	Night Clear	Dusk Rainy	Night Rainy	Day Foggy
no-aug.	51.0	35.9	31.3	16.7	37.7
random	51.2	36.0	30.4	15.3	37.3
clip-random	51.5	36.4	30.2	15.9	37.9
Ours w/ seg.aug	51.3	36.9	32.3	18.7	38.5

Conclusion

- ❖ Textual description of underlying domain shift can be helpful
- ❖ Vision-Text embeddings helps in augmenting missing target image features

