

# Motivation:

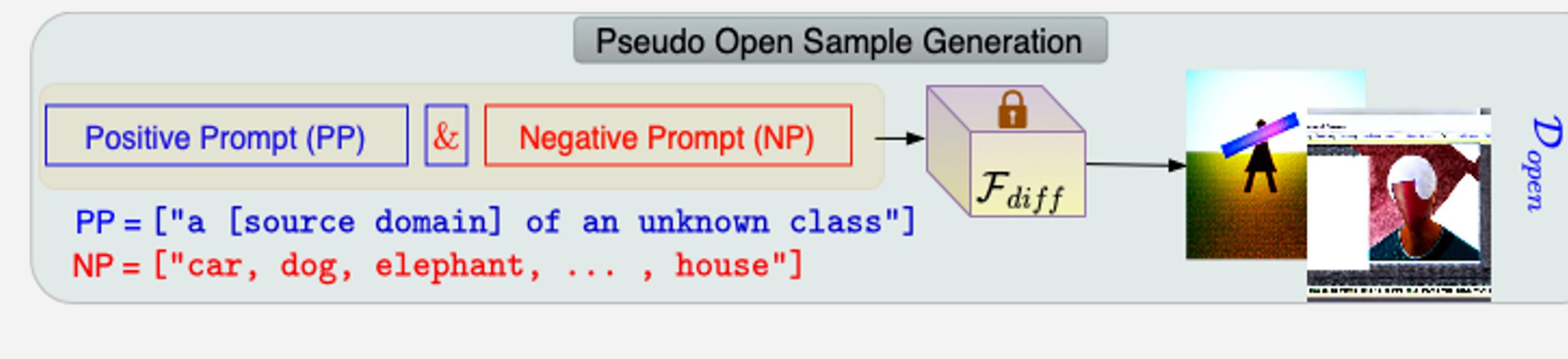
- CNN based generalized models can effectively learn distribution shifts, but often misclassify novel classes during inference, especially when these classes were absent during training.
  - Even zero-shot transfer in CLIP and generalized transfer in existing vision language models (VLMs) encounter challenges in detecting completely unlabeled outlier or **unknown** classes, as they rely on class names in the textual input prompt.

# Our Proposal:

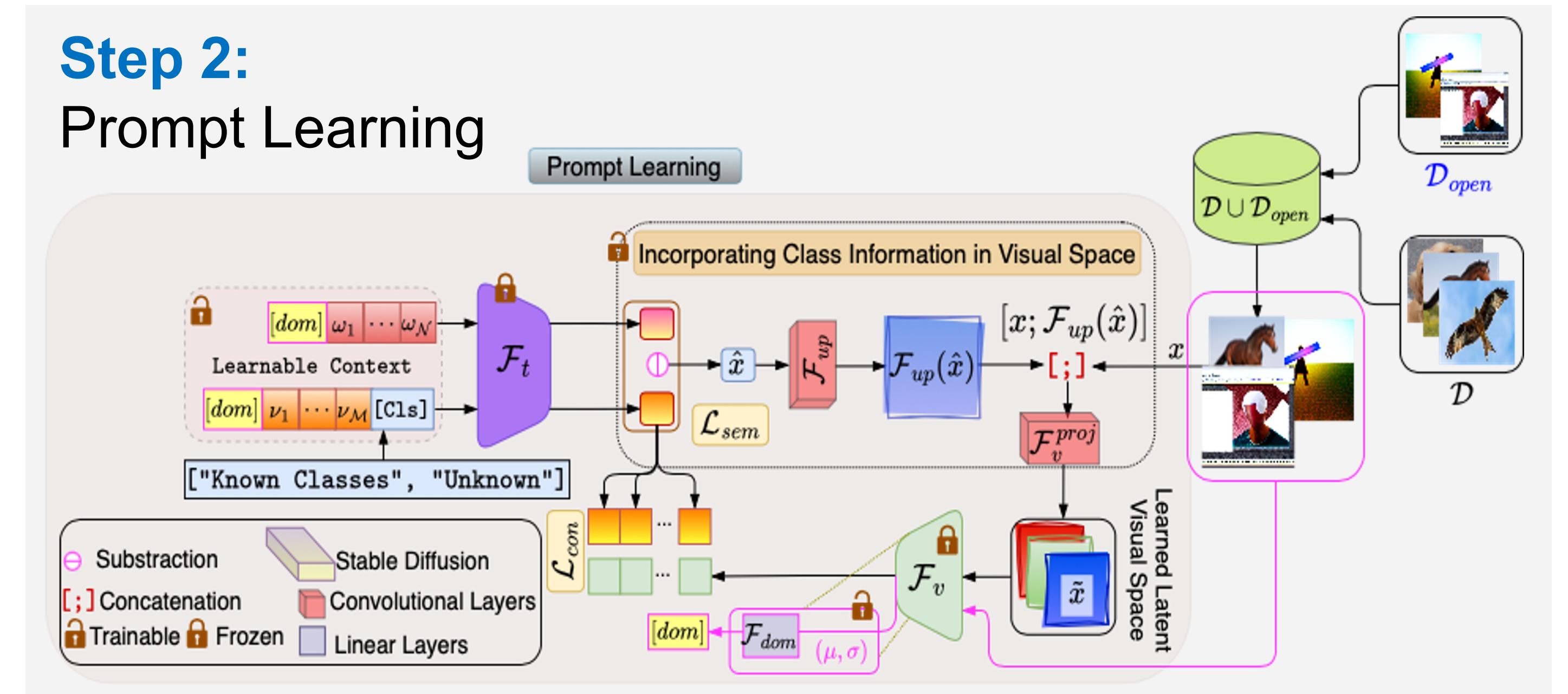
- Leveraging the pre-trained stable diffusion model for the generation of the pseudo-open training samples (**unknown**) using positive and negative prompts.
  - Introducing a novel domain-agnostic prompt learning strategy to generalize the domain-specific style information in the text space.
  - Enhance visual semantics in a learnable latent space by maximizing the similarities between the CLIP's textual and visual features.

# Step 1:

# Pseudo Open Sample Generation



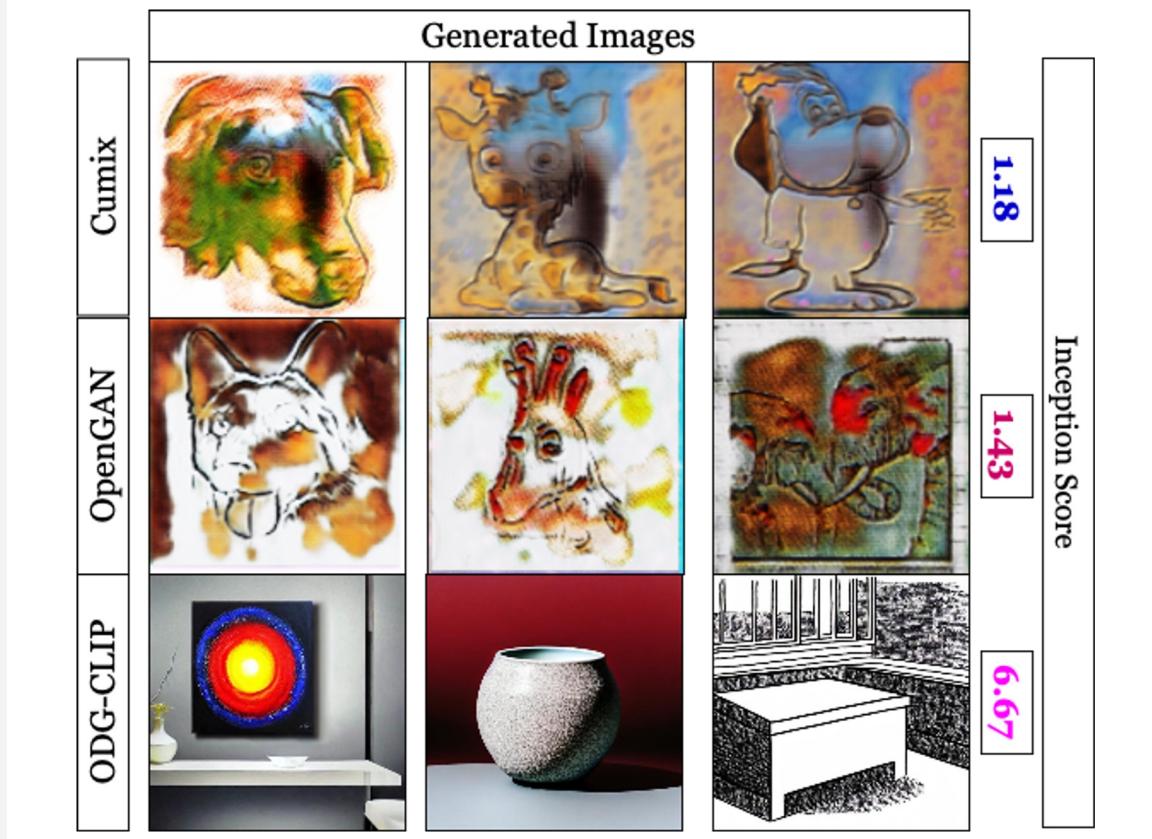
# Step 2: Prompt Learning



# **Proposed Methodology:**

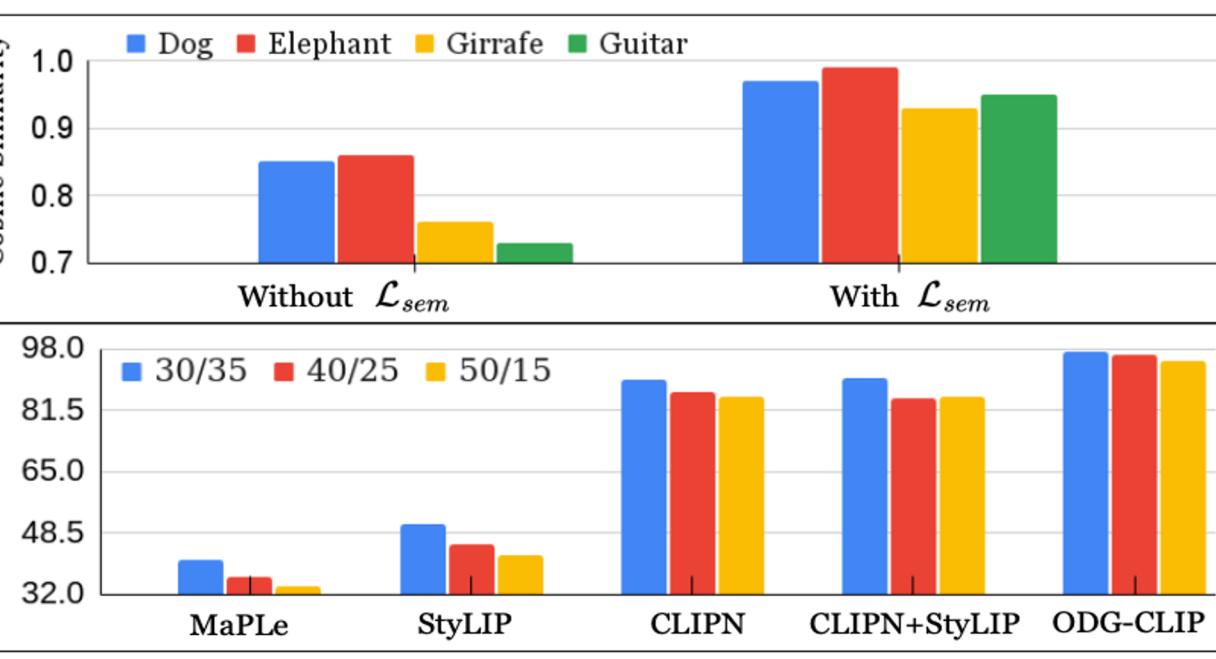
- Aim to improve the quality and semantic versatility of pseudo-open samples by leveraging the pretrained Stable Diffusion v1-5 model. This entails utilizing **unknown** or unseen categories as positive prompts, while known classes serve as negative prompts (indicating what we do not want to generate).
  - Our classification strategy encompasses C+1 classes, where the (C+1)-th index is designated for the novel unknown-class.
  - We advocate for adaptive prompt learning across all classes, enabling the capture of domain-specific distributions and overarching semantic contents through distinct token sets.
  - We refine visual-textual contrastive learning for ODG-CLIP by enhancing the discriminability of visual embeddings. This is achieved by establishing a latent visual space, guided by the prompts we've developed.
  - Also, we propose a cross-domain semantic consistency loss to cultivate a robust class-wise correlation in the derived latent visual representations across images from different domains but sharing identical class labels.

# Experiments and Results:



Methods	PACS		O.H.		M.Data		M.DNet	
	Acc	H	Acc	H	Acc	H	Acc	H
Only PP	92.45	92.16	93.72	90.83	78.20	81.57	91.30	89.78
PP + NP	<b>99.53</b>	<b>99.70</b>	<b>98.32</b>	<b>96.08</b>	<b>84.60</b>	<b>90.00</b>	<b>95.68</b>	<b>94.48</b>

Methods	PACS		O.H.		M.Data		M.DNet	
	Acc	H	Acc	H	Acc	H	Acc	H
w/o $\hat{x}$ and $\mathcal{L}_{sem}$	90.47	88.34	92.21	87.00	73.56	75.73	87.24	83.51
w/o $\mathcal{L}_{sem}$ , with $\hat{x}$	94.21	92.56	95.67	91.56	80.34	85.32	91.24	90.88
Manual $\hat{x}$	93.54	92.82	95.31	91.22	78.53	79.26	90.65	86.52
<b>Full (ours)</b>	<b>99.53</b>	<b>99.70</b>	<b>98.32</b>	<b>96.08</b>	<b>84.60</b>	<b>90.00</b>	<b>95.68</b>	<b>94.48</b>



# Summary:

- Open Domain Generalization (ODG) addresses both known and novel classes across out-of-distribution domains. During inference, the target domain may include **completely unlabeled** familiar categories and entirely **new categories**, resembling real-world recognition scenarios.
  - Our proposed ODG-CLIP harness the semantic capabilities of the vision-language pretrained model CLIP, to solve the ODG task, just using a text prompt of **unknown** class.

