

# Do-Prompt: Causal Interventions Meet Variational Prompt Bottlenecks

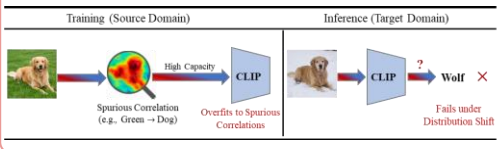
Xueting Chen, Jun-Jie Huang, Yan Yan, Long Lan, Yuhua Tang, Wenjing Yang  
<sup>1</sup> National University of Defense Technology <sup>2</sup> Xiamen University



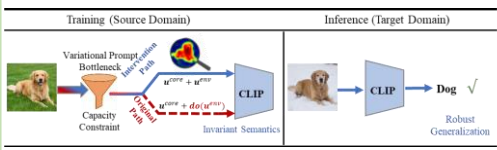
## 1. Motivation

Why standard prompt tuning fails under shift?

(a) Standard Prompt Tuning (Baseline)



(b) Do-Prompt (Ours)



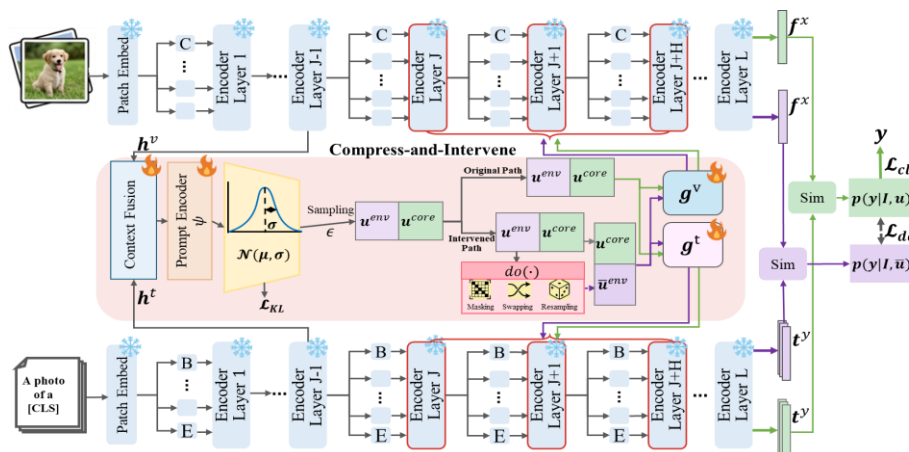
Why Do-Prompt?

- Prompts can become a high-capacity shortcut pathway that encodes environment-specific cues.
- This hurts transfer under distribution shift, especially in multi-modal prompt learning.
- Do-Prompt constrains prompt capacity and intervenes on environment-related prompt content.

### Our Contributions

- Variational prompt bottleneck for stochastic, information-constrained prompts.
- Prompt-level causal interventions on the environment branch (mask / swap / resample).
- Plug-and-play integration with existing multi-modal prompt tuning pipelines.

## 2. Method: Do-Prompt Framework



Key Formulations

$$q_\phi(\mathbf{u}_i | \mathbf{h}) = \mathcal{N}(\mathbf{u}_i(\mathbf{h}), \text{diag}(\sigma_i^2(\mathbf{h})))$$

$$\mathbf{u}_i = [\mathbf{u}_i^{\text{core}}; \mathbf{u}_i^{\text{env}}]$$

$$\mathcal{L}_{\text{KL}} = \sum_{j=1}^{J+H-1} \text{KL}(q_\phi(\mathbf{u}_i | \mathbf{h}) \| p(\mathbf{u}_i))$$

$$\mathcal{L}_{\text{do}} = \mathbb{E}_k [\text{KL}(p_\theta(y | I; \{\mathbf{u}_i\}) \| p_\theta(y | I; \{\mathbf{u}_i^{(k)}\}))]$$

What Do-Prompt Does

- VIB limits channel capacity.
- do-interventions perturb only environment-related prompt components.
- Consistency regularization encourages invariant semantics.

## 3. Main Results

### A. Base to novel generalization results

Method	Average		ImageNet		Cats101		OxfordPets	
	Base	Novel	Base	Novel	Base	Novel	Base	Novel
CoOp (Zhou et al., 2022a)	82.69	63.22	70.83	76.47	69.88	71.92	98.00	89.81
CoCoOp (Zhou et al., 2022b)	80.47	71.69	75.83	75.98	70.43	73.10	97.96	93.81
APP (Choi et al., 2024)	83.0	65.8	72.61	69.9	63.2	66.4	95.2	91.0
PromptSRC++ (Khattak et al., 2023b)	84.93	74.49	78.61	76.77	67.8	72.01	98.07	94.03
MaPLE* (Khattak et al., 2023a)	82.03	75.03	78.37	74.96	66.97	70.74	97.83	94.87
<b>MaPLE* + Do-Prompt (Ours)</b>	<b>83.82</b>	<b>76.38</b>	<b>80.12</b>	<b>76.06</b>	<b>68.80</b>	<b>72.25</b>	<b>98.35</b>	<b>94.87</b>
Improvements $\Delta$	+1.79	+1.65	+1.75	+1.10	+1.83	+1.51	+0.52	0.0
MMRL* (Guo & Gu, 2025)	85.54	76.52	80.59	77.55	67.43	72.14	98.93	94.60
<b>MMRL* + Do-Prompt (Ours)</b>	<b>86.38</b>	<b>79.16</b>	<b>82.56</b>	<b>77.90</b>	<b>68.83</b>	<b>73.09</b>	<b>99.10</b>	<b>94.27</b>
Improvements $\Delta$	<b>+0.84</b>	<b>+2.64</b>	<b>+1.97</b>	<b>+0.35</b>	<b>+1.40</b>	<b>+0.95</b>	<b>+0.17</b>	<b>-0.33</b>

### B. Base to novel generalization results

Source	Target											
	ImageNet	Cats101	OxfordPets	Sun397	StanfordCars	Flow101	Food101	FCV-Cameras	SUN397	DTD	ImageSAT	UZF100
MaPLE* (Khattak et al., 2023a)	67.96	93.17	90.20	65.97	71.07	86.33	23.23	67.23	47.20	45.70	66.27	65.63
<b>MaPLE* + Do-Prompt (Ours)</b>	<b>68.66</b>	<b>93.80</b>	<b>90.73</b>	<b>66.30</b>	<b>74.57</b>	<b>86.40</b>	<b>25.37</b>	<b>67.70</b>	<b>48.17</b>	<b>54.90</b>	<b>70.87</b>	<b>67.19</b>
Improvements $\Delta$	<b>+0.70</b>	<b>+0.63</b>	<b>+0.53</b>	<b>+0.33</b>	<b>+3.50</b>	<b>+0.07</b>	<b>+2.14</b>	<b>+0.47</b>	<b>+0.97</b>	<b>+9.20</b>	<b>+4.60</b>	<b>+1.56</b>
MMRL* (Guo & Gu, 2025)	70.13	94.30	91.00	66.20	71.53	86.20	26.07	67.50	46.93	49.83	69.13	66.87
<b>MMRL* + Do-Prompt (Ours)</b>	<b>71.00</b>	<b>94.60</b>	<b>91.97</b>	<b>66.63</b>	<b>75.57</b>	<b>86.43</b>	<b>28.33</b>	<b>67.90</b>	<b>48.27</b>	<b>58.97</b>	<b>72.50</b>	<b>68.84</b>
Improvements $\Delta$	<b>+0.87</b>	<b>+0.30</b>	<b>+0.97</b>	<b>+0.43</b>	<b>+4.04</b>	<b>+0.23</b>	<b>+2.26</b>	<b>+0.40</b>	<b>+1.34</b>	<b>+9.14</b>	<b>+3.37</b>	<b>+1.97</b>

### C. Base to novel generalization results

Method	Source	Target			
	ImageNet	ImageNetV2	ImageNet-S	ImageNet-A	ImageNet-R
MaPLE* (Khattak et al., 2023a)	67.96	61.57	47.70	48.80	75.33
<b>MaPLE* + Do-Prompt (Ours)</b>	<b>68.66</b>	<b>62.80</b>	<b>49.20</b>	<b>50.30</b>	<b>76.53</b>
Improvements $\Delta$	<b>+0.70</b>	<b>+1.23</b>	<b>+1.50</b>	<b>+1.50</b>	<b>+1.20</b>
MMRL* (Guo & Gu, 2025)	70.13	62.20	47.80	48.90	75.03
<b>MMRL* + Do-Prompt (Ours)</b>	<b>71.00</b>	<b>63.80</b>	<b>49.70</b>	<b>50.80</b>	<b>76.83</b>
Improvements $\Delta$	<b>+0.87</b>	<b>+1.60</b>	<b>+1.90</b>	<b>+1.90</b>	<b>+1.80</b>

## 4. Ablation & Analysis

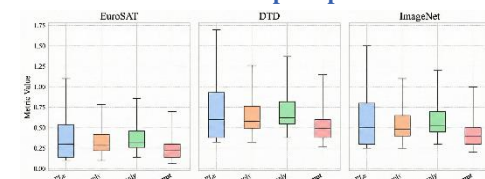
### A. Ablation on key components

Variant	VIB	do-int.	HM
MaPLE* baseline			78.37
+ VIB only	✓		78.93
+ do-intervention only		✓	79.02
<b>MaPLE* + Do-Prompt</b>	✓	✓	<b>80.12</b>

### B. Intervention operators

Intervention	Base	Novel	Avg HM
None (MaPLE)	82.03	75.03	78.37
Mask only	83.01	75.97	79.01
Swap only	82.78	75.75	78.94
Resample only	83.06	75.82	79.15
<b>All (Do-Prompt)</b>	<b>83.82</b>	<b>76.68</b>	<b>80.12</b>

### C. Prediction invariance under prompt interventions



### D. Overhead of Do-Prompt

Method	Train Time / Iter	Inference Time / Img	Params ( $\Delta$ )
MaPLE (Khattak et al., 2023a)	1.00x	1.00x	-
<b>MaPLE + Do-Prompt</b>	1.8x	1.02x	<b>+4.0M</b>
MMRL (Guo & Gu, 2025)	1.00x	1.00x	-
<b>MMRL + Do-Prompt</b>	1.9x	1.03x	<b>+7.0M</b>