





























Table 6: Description of key symbols in this paper

Symbol	Description
$S^i = x_i^k, y_i^k \_{k=1}^{n_i}$	is the definition of source domain in DG problem
$P_{data}^{S^i}$	is the joint distribution concerning the data and the label space of source domain
$S^{N+1} = x_t^k, y_t^k \_{k=1}^{n_t}$	is the definition of target domain in DG problem
$P_{data}^t$	is the joint distribution concerning the data and the label space of target domain
$g()$	is the text encoder in CLIP model
$t$	is the prompt given to the text encoder in our method
$L_f$	is the hidden layer features obtained through linear mapping of the learnable prompt.
$J$	is the number of layers fine-tuned using our method in the image encoder.
$scale_i, bias_i$	is the parameters applied to the image encoder layer $i$ ( $i \in \{1, \dots, J\}$ )
$F_i$	is the output of the image encoder layer $i$ ( $i \in \{1, \dots, J\}$ )

Table 7: Comparison of our proposed method with CLIP- based state-of-the-art methods for upper bound on DomainNet and Office-Home. CLIP Liner is a model obtained by training an additional linear classifier on top of CLIP

Method	CLIP Liner	CoOp	CoCoOp	Im-Tuning
upper bound on DomainNet	76.2	79.3	81.0	<b>88.6</b>
upper bound on Office-Home	86.2	88.1	89.2	<b>94.5</b>

Table 8: Investigations on initialization of our method. We report the average top-1 classification performance for domain generalization on DomainNet and Office-Home.

Initialization	DomainNet	Office-Home
random initialization	75.2	85.6
"a photo of a"	75.2	85.5

Figure 6: t-SNE plots of image features in prompt tuning method CoOp, and our Im-Tuning on Office-Home dataset. Im-Tuning shows better separability.

Table 9: Ablation analysis of the number of layers in Im-Tuning using ViT-B/16 backbone (In%)

Baselines	DomainNet
context length=4,transformer layer 1-3	74.3
context length=4,transformer layer 1-6	74.8
context length=4,transformer layer 1-9	<b>75.2</b>
context length=4,transformer layer 1-12	74.9
context length=16,transformer layer 1-3	74.2
context length=16,transformer layer 1-6	74.6
context length=16,transformer layer 1-9	<b>75.0</b>
context length=16,transformer layer 1-12	74.8