
On the Role of Mechanistic Interpretability for Vision-Language Prompt Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent advances in mechanistic interpretability of vision-language models (VLMs)
2 such as CLIP propose using sparse autoencoders (SAEs) to discover monoseman-
3 tic, human-understandable features that explain CLIP’s internal representations.
4 Existing work using SAEs to probe VLMs primarily focuses on post-hoc inter-
5 pretability analysis. We posit that SAE-based interpretability methods are not
6 just probing tools, but can also serve as meaningful training guides for adapt-
7 ing VLMs to downstream tasks. To this end, we propose IPL (Interpretability-
8 Guided Prompt Learning), which leverages SAE decoders to extract interpretable
9 concept directions, composes them into prompt tokens via a learnable attention
10 selector, and injects the resulting tokens into both the vision and text encoder
11 layers of CLIP for adaptation. We further study how prompt tokens obtained
12 by probing vision-only, text-only, and unified concept directions from respective
13 interpretability methods affect performance on downstream tasks. We perform
14 extensive experiments across downstream settings such as base-to-novel general-
15 ization, domain generalization, cross-dataset transfer, and few-shot learning. While
16 IPL using vision-only and text-only concept directions obtains decent gains, IPL
17 with unified concept directions achieves the strongest results, outperforming prior
18 prompt-learning methods over 15 datasets across all downstream settings. Code:
19 <https://anonymous.4open.science/r/IPL-34EF/>

20 1 Introduction

21 Mechanistic interpretability (MI) research on vision-language models (VLMs) like CLIP [30] has
22 made rapid progress in the past years. Sparse autoencoders (SAEs) [23, 32, 2] have emerged as a
23 prominent tool in the MI toolkit and they have been shown to be able to decompose intermediate
24 representations of VLMs into human-understandable concept directions: dictionaries of features that
25 fire selectively on coherent semantic categories such as object parts, textures, or scene types, with
26 explicit per-input activation magnitudes that quantify each concept’s presence in any given image.
27 These tools have given researchers a better picture of what CLIP has learned. So far, however, this
28 growing toolbox has been used almost exclusively for post-hoc purposes like *interpreting* model
29 behavior after the model has been trained. Our central research question is *can interpretability tools*
30 *act as a training guide for adapting VLM’s like CLIP to different downstream settings?*

31 We posit that this is exactly where mechanistic interpretability can be practically useful, and we
32 ground the question in a concrete use case, i.e., **prompt learning**. Prompt learning has emerged as
33 the dominant parameter-efficient strategy for adapting CLIP to downstream tasks, with approaches
34 like CoOp [45], MaPLe [20], and MMRL [15] introducing learnable prompt tokens at various
35 layers of CLIP. While effective, these prompts are opaque continuous vectors with no human-
36 interpretable semantics. Recent attempts to make prompts interpretable like CPL [42], ArGue [35],
37 IntCoOp [14], and IPO [10], do so by importing concepts from *outside* CLIP: querying language

38 models or vision-language models for natural-language attributes and hard-coding them into the
39 prompt. The interpretability is thus outsourced to *external* tools; it does not engage with the rich
40 *internal* concept structure that mechanistic interpretability methods [23, 2, 32] have shown already
41 exists *inside* CLIP itself.

42 This gap motivates our central proposal. We treat SAE-derived concept dictionaries not as analytical
43 artifacts/probes but as a semantically grounded concept directions from which prompt tokens can
44 be composed. Rather than asking a large language model (LLM) to provide semantic attributes for
45 prompt tokens, we leverage CLIP’s own interpretable features (the same features that MI researchers
46 have already identified, labeled, and inspected) to guide the prompt-learning process. To this end, we
47 propose **IPL** (Interpretability-Guided Prompt Learning), a prompt-learning framework to adapt to a
48 dataset, that extracts a concept dictionary from a frozen SAE for that dataset, composes prompt tokens
49 as learnable attention-weighted combinations of those concepts, and projects the resulting tokens
50 into both vision and text encoder layers. A projection-preservation regularizer keeps prompt-token
51 geometry aligned with the underlying concept directions, so that the learned prompts remain anchored
52 to the interpretable concept directions they were composed from. The framework is agnostic to the
53 choice of SAE based probing methods: vision-only, text-only, and unified concept dictionaries can all
54 be plugged in.

55 Overall, our contributions and key findings can be summarized as follows:

- 56 • We introduce IPL, a prompt-learning framework that composes prompt tokens via attention over
57 interpretable concept directions and preserves their geometry through a projection preservation
58 regularizer. The framework is agnostic to the source SAE, supporting vision-only, text-only, and
59 unified concept dictionaries.
- 60 • Through this framework, we show that mechanistic interpretability tools can be use as a training
61 guide for vision-language prompt learning.
- 62 • We systematically study how the modality of interpretability signals (vision, text, or unified) affects
63 downstream performance, finding that unified signals from VL-SAE [32] produce the strong gains.
- 64 • Across four standard prompt-learning downstream task settings and 15 evaluation datasets, IPL
65 outperforms prior interpretable prompt-learning methods by up to +2.49% in harmonic mean on
66 base-to-novel generalization while remaining competitive with strong black-box prompt-learning
67 baselines.

68 2 Related Works

69 **Mechanistic interpretability of vision-language models.** Mechanistic interpretability aims to de-
70 compose neural computations into human-understandable circuits and features. Sparse autoencoders
71 (SAEs) [7] are a promising tool for this purpose: by reconstructing intermediate activations with
72 a sparse overcomplete dictionary, SAE decoder directions often capture monosemantic concepts.
73 PatchSAE [23] applies SAEs to CLIP vision-transformer residual streams to reveal patch-level visual
74 concepts. SpLiCE [2] decomposes CLIP image embeddings into sparse non-negative combinations of
75 text-aligned concept directions. VL-SAE [32] learns a unified latent concept space for image and text
76 inputs with modality-specific decoders. So far, these methods have mainly been used for post-hoc
77 interpretation. In contrast, IPL uses SAE-derived concept directions as a structured basis for prompt
78 learning, bridging mechanistic interpretability with downstream adaptation.

79 **Prompt learning for vision-language models.** CLIP [30], ALIGN [19], and related contrastive
80 vision-language models achieve strong zero-shot recognition using hand-crafted text prompts. Prompt
81 learning adapts these models efficiently to downstream tasks while preserving pretrained capabilities.
82 CoOp [45] learns continuous text prompts, and CoCoOp [44] makes them image-conditional to
83 reduce base-class overfitting. KgCoOp [40] regularizes prompts toward hand-crafted templates, while
84 ProDA [26] models prompt distributions for diverse class semantics. Multimodal methods such
85 as MaPLe [20] learn coupled vision-text prompts, while PromptSRC [4], MMA [39], TCP [41],
86 2SFS [11], and SkipT [37] improve generalization through regularization or architectural changes.
87 More recently, MMRL [15] injects multimodal representation tokens into intermediate layers, and
88 VaMP [5] learns sample-conditioned text prompts via variational inference. Despite their effectiveness,
89 these methods learn opaque continuous prompt vectors without human-interpretable semantics.

90 **Interpretable prompt learning.** A separate line of work improves prompt interpretability by using
91 human-readable concepts as prompt content. CPL [42] caches class-relevant concept descriptions

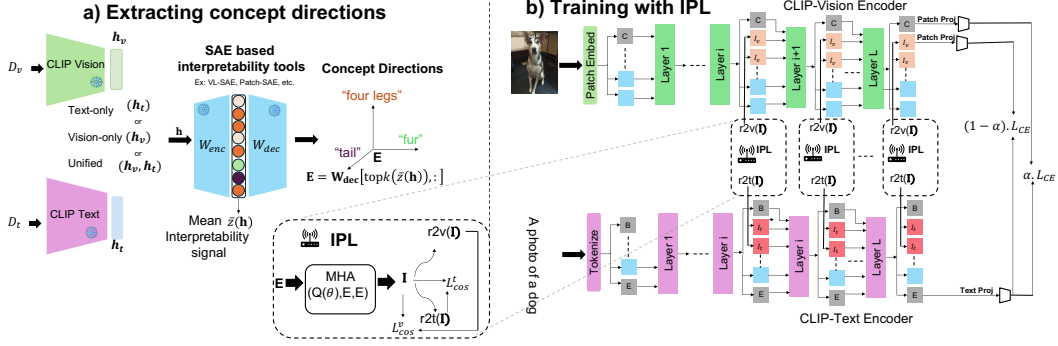


Figure 1: a) We extract concept directions from the training set of dataset to be adapted. We use these concept directions to construct concept attention weighted prompt tokens. b) IPL prompt tokens are projected at different layers in both text and vision encoder of CLIP. A projection regularization is applied to ensure concept geometry after projection.

92 as context, while ArGue and ArGue-N [35] use GPT-3 to generate visual attributes per class (e.g.,
 93 a long tail,” black paws”) and concatenate selected attributes into prompts. IntCoOp [14] extracts
 94 attribute labels with BLIP-2 and aligns them with class embeddings during prompt tuning. IPO [10]
 95 optimizes prompts in discrete token space so learned tokens remain natural-language phrases. Overall,
 96 these methods rely on *externally generated* concepts or attributes, so interpretability comes from
 97 the prompt input form rather than CLIP’s internal computation. In contrast, IPL derives concept
 98 directions from *inside* CLIP using mechanistic interpretability methods, explicitly connecting prompt
 99 learning to CLIP’s internal representations. Please refer to appendix Section E for detailed related
 100 works.

101 3 Method.

102 Our approach **IPL** (Interpretability-Guided Prompt Learning) leverages mechanistic interpretability
 103 to guide prompt learning for adapting VLMs such as CLIP [30] to downstream tasks like base to
 104 novel generalization, domain generalization, cross-dataset evaluation and few-shot learning. We first
 105 introduce preliminaries elucidating zero-shot classification, prompt learning with CLIP, and sparse
 106 autoencoder-based interpretability probing in Sections 3.1.1 to 3.1.3 respectively. Then, we show
 107 how our approach leverages CLIP-based mechanistic interpretability methods like VL-SAE [32],
 108 PatchSAE [23] etc, to obtain vision, language and unified interpretability signals in Section 3.2.1, and
 109 explain how we use the interpretability signals to guide construction of prompt tokens in Sections 3.2.2
 110 and 3.2.3. In Section 3.2.4, we use the interpretability signal guided prompt tokens as learnable
 111 representations and describe the training and inference process.

112 3.1 Preliminaries

113 3.1.1 Revisiting CLIP zero-shot classification.

114 CLIP consists of paired image encoder f_v and text encoder f_t trained contrastively to align image and
 115 text embeddings in a shared latent space. **Image encoding:** Given an input image $\mathbf{X} \in \mathbb{R}^{H \times W \times 3}$,
 116 the visual transformer backbone [9] first patchifies \mathbf{X} into \mathcal{M} patch tokens, prepends a [CLS] class
 117 token, and adds positional embeddings to obtain the input sequence $\mathbf{X}^{(0)} \in \mathbb{R}^{\mathcal{M} \times d_v}$, with $d_v = 768$
 118 for ViT-B/16. This sequence is processed by \mathcal{L} transformer layers as

$$[CLS_i, X_i] = f_v([CLS_{i-1}, X_{i-1}]), \quad \forall i \in \{1, \dots, L\} \quad (1)$$

119 The final image embeddings are obtained by applying layer normalization and a patch projection
 120 layer P_v to the [CLS] token of the last layer thus projecting the output of [CLS] token into shared
 121 vision-language latent space.

$$feat_v = P_v(CLS_L) \quad feat_v \in \mathbb{R}^{d_{vi}} \quad (2)$$

122 **Text encoding:** A textual prompt like “a photo of a [class]” is tokenized and embedded, pro-
 123 ducing \mathcal{G} token embeddings $w_0 = [w_0^1, w_0^2, \dots, w_0^{\mathcal{G}}] \in \mathbb{R}^{\mathcal{G} \times d_t}$. After processing through \mathcal{L} text

124 transformer layers $[w_i] = f_t([w_{i-1}]); \forall i \in \{1, \dots, L\}$, the text embedding features are obtained by
 125 and passing through a final projection layer P_t

$$feat_t = P_t(w_L^G) \quad feat_t \in \mathbb{R}^{d_{vt}} \quad (3)$$

126 **Zero-shot classification** For an image \mathbf{X} and text prompts with class labels $\mathcal{Y} \in \{1, 2, \dots, C\}$. the
 127 prediction is decided by highest cosine similarity, using temperature coefficient τ :

$$\hat{y} = \arg \max_y \frac{\exp(\text{sim}(feat_v, feat_t)/\tau)}{\sum_{i=1}^C \exp(\text{sim}(feat_v, feat_t^i)/\tau)} \quad (4)$$

128 3.1.2 Understanding prompt learning for efficient adaptation.

129 Hand-engineered prompts like ‘‘a photo of a [class]’’ that are used as inputs to CLIP are
 130 often sub-optimal for downstream tasks on complex and fine-grained datasets as they may not
 131 be aligned with the semantic context of different domains. One way out is to fine-tune CLIP on
 132 prompts engineered for each domain but it can lead to catastrophic forgetting and is computationally
 133 expensive. The alternative, prompt learning, introduces continuous learnable tokens either in the
 134 input [20, 45, 44, 10, 42] or intermediate layers [15, 39, 5] for adapting CLIP to downstream tasks
 135 without hand-engineering the prompts and keeping the CLIP’s encoder frozen. Existing approaches
 136 differ in the layers and encoders in which these prompt tokens are inserted. **Language prompt**
 137 **learning** introduces randomly initialized learnable prompt tokens $I_h^t \in \mathbb{R}^{N \times d_t}$, consisting of N
 138 tokens, into each H consecutive transformer layers $h \in \{J, J+H-1\}$ of the text encoder f_t , leaving
 139 the rest of the layers unchanged. Similarly, **vision prompt learning** introduces N learnable prompt
 140 tokens $I_h^v \in \mathbb{R}^{N \times d_v}$. In **Multimodal prompt learning** prompt tokens are added to both text and
 141 vision branches. MaPLe [20] passes language prompt tokens to a learnable coupling function to
 142 obtain vision prompts tokens, MMRL [15] proposes a learnable and shared representation space with
 143 tokens initialized by sampling from gaussian distribution and uses a linear projection layer to get
 144 vision and text branch prompt tokens, and VaMP [5] injects sample specific text prompts obtained by
 145 variation inference conditioned on image features, and also simultaneously injects vision prompts.

146 3.1.3 Sparse-Autoencoder based mechanistic interpretability

147 Recent mechanistic interpretability studies [23, 32, 7] have shown that sparse autoencoders (SAE’s)
 148 find highly interpretable features in their latent space. We will first elucidate the working mechanism
 149 and terminologies for SAE’s and then define what constitutes a ‘‘concept’’ and ‘‘interpretability
 150 signal’’. **Sparse Autoencoder:** Given an input image \mathbf{X} to CLIP image encoder, $\mathbf{h} \in \mathbb{R}^{d_v}$ is CLIP
 151 vision encoder layer’s ($l \in L$) representation. These are given as an input to SAE with encoder
 152 $\mathbf{W}_{\text{enc}} \in \mathbb{R}^{d_v \times K}$ produces concept activations

$$\mathbf{z}(\mathbf{h}) = \sigma(\mathbf{h}\mathbf{W}_{\text{enc}} + \mathbf{b}_{\text{enc}}) \in \mathbb{R}^K, \quad (5)$$

153 and decoder $\mathbf{W}_{\text{dec}} \in \mathbb{R}^{K \times d_v}$ reconstructs $\hat{\mathbf{h}}$ from \mathbf{z}

$$\hat{\mathbf{h}}(\mathbf{z}(\mathbf{h})) = \mathbf{W}_{\text{dec}}\mathbf{z}(\mathbf{h}) + \mathbf{b}_{\text{dec}}, \quad (6)$$

154 Here, σ is a non-negative activation (ReLU or TopK) inducing sparsity, \mathbf{b}_{enc} and \mathbf{b}_{dec} are the bias
 155 terms. The SAE is trained to minimize $\|\mathbf{h} - \hat{\mathbf{h}}(\mathbf{z}(\mathbf{h}))\|^2 + \lambda\|\mathbf{z}(\mathbf{h})\|$ on a huge corpus of representations
 156 (VL-SAE [32] is trained with CC3M dataset). K is dictionary size (typically $K \gg d_v$. For instance,
 157 in PatchSAE[23] K is 64 times the dimension of d_v). Each row $\mathbf{W}_{\text{dec}}[i, :]$ is interpreted as concept
 158 i ’s direction in the SAE representation space. Each entry in \mathbf{z}_i measures concept i ’s presence in \mathbf{h} .

159 **Definition 1** (Concept). *Each SAE decoder row can be viewed as a candidate concept feature with a*
 160 *unit-norm direction.*

$$\mathbf{v}_i = \frac{\mathbf{W}_{\text{dec}}[i, :]}{\|\mathbf{W}_{\text{dec}}[i, :]\|} \in \mathbb{R}^{d_v}, \quad i \in \{1, \dots, K\}, \quad (7)$$

161 *Consider an image \mathbf{X} of a dog passed through CLIP, yielding a residual-stream represen-*
 162 *tation $\mathbf{h} \in \mathbb{R}^{d_v}$. A sparse autoencoder trained on CLIP’s residuals might yield a con-*
 163 *cept dictionary containing candidate concept directions that can be empirically annotated as*
 164 *\mathbf{v}_{fur} , $\mathbf{v}_{\text{four-legged}}$, \mathbf{v}_{tail} , \mathbf{v}_{sky} , $\mathbf{v}_{\text{indoor scene}}$, \dots , based on the training images on which the correspond-*
 165 *ing SAE latent activates strongly. Thus, each candidate concept feature is paired with an evidence set*
 166 *of training images supporting its semantic interpretation.*

167 **Definition 2 (Interpretability Signal).** Given a representation $\mathbf{h} \in \mathbb{R}^{d_v}$ and a candidate concept
 168 feature direction \mathbf{v}_i , the **interpretability signal** for \mathbf{v}_i on \mathbf{h} is defined as the corresponding SAE latent
 169 activation

$$z_i(\mathbf{h}) = \sigma(\mathbf{W}_{\text{enc}}[:, i]^\top \mathbf{h} + b_{\text{enc}, i}) \in \mathbb{R}_{\geq 0}, \quad (8)$$

170 i.e., the i -th entry of $\mathbf{z}(\mathbf{h})$ in Equation (5). The signal $z_i(\mathbf{h})$ serves as an input-dependent activation
 171 score for the candidate concept feature \mathbf{v}_i . Interpretability signals are therefore input-dependent: the
 172 same candidate concept feature can produce different signal magnitudes for different representations.

173 **Vision, text, and unified interpretability signals.** The interpretability methods we use yield three
 174 modalities of signal, distinguished by which CLIP encoder produces the representation \mathbf{h} that drives
 175 concept activation. **Vision interpretability signals:** are obtained by running images through CLIP’s
 176 vision encoder and reading concept activations from a vision-side SAE. Methods like PatchSAE [23],
 177 VL-SAE-vision [32] operate in this mode. **Text interpretability signals:** are obtained by running
 178 per-class textual prompts through CLIP’s text encoder and reading activations from a text-side SAE.
 179 Methods like VL-SAE-text [32] operate in this mode. **Unified interpretability signals:** constitute a
 180 unified concept set obtained by leveraging both image and text embeddings from CLIP. VL-SAE-
 181 unified [32] instantiates this construction.

182 3.2 Interpretability Guided Prompt Learning

183 3.2.1 Extracting concepts from interpretability methods

184 Given a frozen sparse autoencoder provided by any CLIP [30] based mechanistic interpretability
 185 method (we use pretrained SAE provided by VL-SAE [32] and PatchSAE [23]), we construct an
 186 interpretable concept dictionary $\mathbf{E} \in \mathbb{R}^{k \times d_v}$, $k \subset K$, where each row of \mathbf{E} is a concept direction
 187 (Definition 1) corresponding to a human-meaningful concept in CLIP’s residual-stream space. The
 188 construction of \mathbf{E} has three steps.

189 **Computing mean interpretability signal.** Let $\mathcal{D}_{\text{base}}$ denote the base-class training set, comprising
 190 images and their associated class prompts. We forward each sample through the CLIP encoder then
 191 pass the resulting representation \mathbf{h} through the SAE encoder to obtain per-sample interpretability
 192 signals $z(\mathbf{h})$ over the SAE’s full dictionary of K concepts (Definition 2). So for \mathcal{S} samples in $\mathcal{D}_{\text{base}}$,
 193 we then compute its mean signal \bar{z}_i across all samples in $\mathcal{D}_{\text{base}}$, yielding a per-concept mean signal:
 194 $\bar{\mathbf{z}}(\mathbf{h}) = \sum_{i=1}^{\mathcal{S}} (z_i(\mathbf{h}))$

195 We refer to $\bar{\mathbf{z}}(\mathbf{h})$ as a *vision signal* when \mathbf{h} is produced by the image encoder, a *text signal* when \mathbf{h} is
 196 produced by the text encoder, and a *unified signal* when the two are concatenated to obtain a unified
 197 concept activation.

198 **Top- k selection.** We retain the indices \mathcal{I} of the top- k , where $k \subset K$ obtained by mean signal
 199 magnitude $|\bar{\mathbf{z}}(\mathbf{h})|$. The resulting concept dictionary is

$$\mathbf{E} = \mathbf{W}_{\text{dec}}[\mathcal{I}, :] \in \mathbb{R}^{k \times d_v}, \quad (9)$$

200 We instantiate three modality variants : $\mathbf{E}^{\text{vision}}$ is constructed using the SAE decoder weights ranked
 201 by mean of image interpretability signals, \mathbf{E}^{text} is constructed using the SAE decoder representations
 202 ranked by mean text interpretability signals on per-class textual prompts, and $\mathbf{E}^{\text{unified}}$ is constructed
 203 using the SAE decoder representations ranked by unified interpretability signals like VL-SAE [32]
 204 encoder is trained to produce unified interpretability signals in the shared latent space from both
 205 vision and text inputs.

206 3.2.2 Interpretability signal guided prompt tokens

207 We introduce N learnable representation queries $\mathbf{Q} = [\mathbf{q}_1; \dots; \mathbf{q}_M] \in \mathbb{R}^{M \times d_v}$, trained jointly with
 208 the rest of the prompt-learning objective. Unlike prior works like [15, 5] which uses Gaussian-
 209 initialized free-form representations, our queries select content from the frozen interpretable dictionary
 210 \mathbf{E} via attention. For each rep-injection layer $\ell \in \mathcal{L}$ transformer layers, we instantiate a multi-head
 211 attention selector $\text{MHA}^{(\ell)}$ that attends from \mathbf{Q} to \mathbf{E} :

$$\mathbf{I}^{(\ell)} = \text{MHA}^{(\ell)}(\mathbf{Q}, \mathbf{E}, \mathbf{E}) \in \mathbb{R}^{N \times d_v}, \quad (10)$$

212 where \mathbf{Q} provides the queries and \mathbf{E} provides both keys and values. Each row of $\mathbf{I}^{(\ell)}$ is a layer-specific
 213 prompt token expressed as an attention-weighted combination of interpretable concepts.

214 **3.2.3 Projecting prompt tokens into encoder layers**

215 The prompt tokens $\mathbf{I}^{(\ell)}$ live in CLIP’s vision-encoder residual space (d_v) and must be projected
 216 into the form expected at each injection layer in both the vision and text encoders. We employ two
 217 per-layer projections:

$$\mathbf{I}_v^{(\ell)} = r2v^{(\ell)}(\mathbf{I}^{(\ell)}) \in \mathbb{R}^{N \times d_v} \quad \text{and} \quad \mathbf{I}_t^{(\ell)} = r2t^{(\ell)}(\mathbf{I}^{(\ell)}) \in \mathbb{R}^{N \times d_t}, \quad (11)$$

218 where $r2v^{(\ell)}$ is a linear map preserving residual width and $r2t^{(\ell)}$ projects to the text encoder width
 219 d_t . The projected tokens are concatenated to the existing token sequence at layer ℓ in their respective
 220 encoders as defined in Section 3.1.2

221 **3.2.4 Learning with projected prompt tokens**

222 **Training.** For an image-class pair (\mathbf{x}, y) , IPL produces two complementary classification logits. The
 223 first reflects the standard image-text similarity using the [CLS] embedding.

$$\mathbf{I}_{\text{cls}}^{(c)}(\mathbf{x}) = \frac{\exp(\text{sim}(\text{feat}_v, \text{feat}_t^{(c)})/\tau)}{\sum_{i=1}^C \exp(\text{sim}(\text{feat}_v, \text{feat}_t^{(i)})/\tau)}, \quad (12)$$

224 where feat_v and $\text{feat}_t^{(c)}$ denote [CLS] token projection and CLIP text projection features computed
 225 with prompt tokens injected at layers \mathcal{L} , and τ is CLIP’s learned temperature coefficient.

226 The second logit derives image features from the prompt tokens themselves: we take the mean-pooled
 227 prompt-token output at the final injection layer, denoted $\bar{\mathbf{I}}(\mathbf{x}) = \text{mean}(\mathbf{I}(\mathbf{x}))$, and use it in place of
 228 the CLS-derived image feature:

$$\mathbf{I}_{\text{rep}}^{(c)}(\mathbf{x}) = \frac{\exp(\text{sim}(\bar{\mathbf{I}}(\mathbf{x}), \text{feat}_t^{(c)})/\tau)}{\sum_{i=1}^C \exp(\text{sim}(\bar{\mathbf{I}}(\mathbf{x}), \text{feat}_t^{(i)})/\tau)}. \quad (13)$$

229 **Projection-preservation regularization.** A naive linear projection performed to obtain prompt
 230 tokens (as described in Section 3.2.3) can rotate concept structure arbitrarily, decoupling the prompt
 231 tokens’ geometry from \mathbf{I} ’s interpretable structure. To encourage the projections to preserve concept
 232 geometry, we introduce a projection-preservation loss:

$$\mathcal{L}_{\text{proj}} = \frac{1}{|\mathcal{L}|} \sum_{\ell \in \mathcal{L}} \left(\underbrace{1 - \cos\left(r2v^{(\ell)}(\mathbf{I}_v^{(\ell)}), \mathbf{I}_v^{(\ell)}\right)}_{\text{vision: cosine alignment}(\mathcal{L}_{\text{cos}}^v)} + \underbrace{1 - \cos\left(r2t^{(\ell)}(\mathbf{I}_t^{(\ell)}), \mathbf{I}_t^{(\ell)}\right)}_{\text{text: cosine alignment}(\mathcal{L}_{\text{cos}}^t)} \right), \quad (14)$$

233 where $\cos(\cdot, \cdot)$ averages cosine similarity row-wise. The vision term encourages $r2v^{(\ell)}$ to act near-
 234 identically on $\mathbf{I}_v^{(\ell)}$, while the text term preserves the pairwise similarity structure of $\mathbf{I}_t^{(\ell)}$ under $r2t^{(\ell)}$
 235 even though dimensionality changes. Together, these constrain the projections to keep prompt-token
 236 geometry similar to concept directions.

237 Our training objective is the Lagrangian with α balancing the two cross-entropy objectives, and λ_{proj}
 238 is the projection-preservation loss weight tuned per dataset:

$$\mathcal{L} = \alpha \mathcal{L}_{\text{CE}}(\mathbf{I}_{\text{cls}}, y) + (1 - \alpha) \mathcal{L}_{\text{CE}}(\mathbf{I}_{\text{rep}}, y) + \lambda_{\text{proj}} \mathcal{L}_{\text{proj}}. \quad (15)$$

239 **The trainable parameters are:** queries \mathbf{Q} in layer-specific selectors $\{\text{MHA}^{(\ell)}\}$, and projections
 240 $\{r2v^{(\ell)}, r2t^{(\ell)}\}$. CLIP’s encoders, \mathbf{E} , and the interpretability methods’ parameters remain frozen
 241 throughout training.

242 **Inference.** For base-class evaluation, we use the fused logits $\alpha \mathbf{I}_{\text{cls}} + (1 - \alpha) \mathbf{I}_{\text{rep}}$ to leverage both the
 243 standard and rep-token classification signals. For novel-class evaluation, we use only \mathbf{I}_{cls} , since rep
 244 tokens are trained on base classes and may bias predictions toward seen categories. This asymmetry
 245 preserves the zero-shot transfer capability of the underlying CLIP alignment for unseen classes.

246 **4 Experiments**

247 We evaluate IPL on four standard downstream tasks for prompt learning: base-to-novel generalization,
 248 few-shot learning, domain generalization, and cross-dataset transfer.

Table 1: Base-to-novel generalization accuracies (B: Base, N: Novel, H: Harmonic Mean) compared to state-of-the-art interpretable and blackbox methods. Bold indicates best overall and underline indicates best interpretable method. **Note: the average reported in this table is over all 11 datasets, please refer to Section A for detailed results with on other datasets.**

| Method | Average | | | ImageNet | | | StanfordCars | | | SUN397 | | | UCF101 | | | |
|----------------------|---------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|-------|
| | B | N | H | B | N | H | B | N | H | B | N | H | B | N | H | |
| <i>Interpretable</i> | CPL [42] | 84.38 | 78.03 | 81.08 | 78.74 | 72.03 | 75.24 | 79.31 | 76.65 | 77.96 | 81.88 | 79.65 | 80.75 | 86.73 | 80.17 | 83.32 |
| | IPO [10] | 79.92 | 80.51 | 80.21 | 77.83 | 72.45 | 75.04 | 73.42 | 75.71 | 74.55 | 81.25 | 80.92 | 81.08 | 85.32 | 80.92 | 83.06 |
| | IntCoOp [†] [14] | 83.20 | 78.72 | 80.75 | 75.99 | 72.67 | 74.29 | 77.04 | 76.32 | 76.67 | 81.63 | 79.33 | 80.46 | 86.76 | 79.42 | 82.92 |
| | ArGue-N [35] | 83.77 | 78.74 | 81.18 | 76.95 | 71.86 | 74.32 | 75.06 | 74.18 | 74.62 | 81.89 | 80.48 | 81.18 | 86.00 | 79.43 | 82.58 |
| | IPL (Vision) | 85.80 | 77.15 | 81.25 | 78.30 | 70.80 | 74.36 | 81.10 | 75.20 | 78.04 | 83.00 | 79.20 | 81.06 | 88.60 | 79.90 | 84.03 |
| | IPL (Text) | 85.75 | 76.59 | 80.91 | 78.40 | 70.70 | 74.35 | 81.40 | 75.50 | 78.34 | 83.10 | 79.40 | 81.21 | 88.80 | 80.00 | 84.17 |
| IPL (Unified) | <u>86.26</u> | <u>81.23</u> | <u>83.67</u> | <u>79.30</u> | <u>74.50</u> | <u>76.83</u> | <u>81.90</u> | <u>80.40</u> | <u>81.14</u> | 83.00 | <u>82.10</u> | <u>82.55</u> | <u>89.30</u> | <u>85.80</u> | <u>87.52</u> | |
| <i>Blackbox</i> | CLIP [30] | 69.34 | 74.22 | 71.70 | 72.43 | 68.14 | 70.22 | 63.37 | 74.89 | 68.65 | 69.36 | 75.35 | 72.23 | 70.53 | 77.50 | 73.85 |
| | CoOp [45] | 82.69 | 63.22 | 71.66 | 76.47 | 67.88 | 71.92 | 78.12 | 60.40 | 68.13 | 80.60 | 65.89 | 72.51 | 84.69 | 56.05 | 67.46 |
| | CoCoOp [44] | 80.47 | 71.69 | 75.83 | 75.98 | 70.43 | 73.10 | 70.49 | 73.59 | 72.01 | 79.74 | 76.86 | 78.27 | 82.33 | 73.45 | 77.64 |
| | ProDA [26] | 81.56 | 72.30 | 76.65 | 75.40 | 70.23 | 72.72 | 74.70 | 71.20 | 72.91 | 78.67 | 76.93 | 77.79 | 85.23 | 71.97 | 78.04 |
| | KgCoOp [40] | 80.73 | 73.60 | 77.00 | 75.83 | 69.96 | 72.78 | 71.76 | 75.04 | 73.36 | 80.29 | 76.53 | 78.36 | 82.89 | 76.67 | 79.65 |
| | MaPLe [20] | 82.28 | 75.14 | 78.55 | 76.66 | 70.54 | 73.47 | 72.94 | 74.00 | 73.47 | 80.82 | 78.70 | 79.75 | 83.00 | 78.66 | 80.77 |
| | PromptSRC [4] | 84.26 | 76.10 | 79.97 | 77.60 | 70.73 | 74.01 | 78.27 | 74.97 | 76.58 | 82.67 | 78.47 | 80.52 | 87.10 | 78.80 | 82.74 |
| | TCP [41] | 84.13 | 75.36 | 79.51 | 77.27 | 69.87 | 73.38 | 80.80 | 74.13 | 77.32 | 82.63 | 78.20 | 80.35 | 87.13 | 80.77 | 83.83 |
| | MMA [39] | 83.20 | 76.80 | 79.87 | 77.31 | 71.00 | 74.02 | 78.50 | 73.10 | 75.70 | 82.27 | 78.57 | 80.38 | 86.23 | 80.03 | 82.20 |
| | 2SFS [11] | 85.55 | 75.48 | 80.20 | 77.71 | 70.99 | 74.20 | 82.50 | 74.80 | 78.46 | 82.59 | 78.91 | 80.70 | 87.85 | 78.19 | 82.74 |
| | SkipT [37] | 85.04 | 77.53 | 81.11 | 77.73 | 70.40 | 73.89 | 82.93 | 72.50 | 77.37 | 82.40 | 79.03 | 80.68 | 87.30 | 82.47 | 84.81 |
| | MMRL [15] | 85.68 | 77.16 | 81.20 | 77.90 | 71.30 | 74.45 | 81.30 | 75.07 | 78.06 | 83.20 | 79.30 | 81.20 | 88.10 | 80.07 | 83.89 |
| | VaMP [5] | 86.45 | 78.67 | 82.37 | 78.98 | 73.45 | 76.11 | 83.78 | 80.14 | 81.91 | 83.37 | 78.95 | 81.09 | 88.52 | 78.99 | 83.48 |

249 **Additional analyses in the appendix.** Due to space constraints, several supporting analyses are
250 deferred to the appendix: (i) loss-component ablations (Section C.1), (ii) the effect of prompt-
251 token injection layers (Section C.2), (iii) the effect of concept-dictionary size (Section C.3) k , (iv)
252 sensitivity to the loss weights α and λ_{proj} (Section D), (v) the effect of the number of prompt tokens
253 (Section D.2), (vi) the effect of the layer range used for prompt-token injection (Section C.2), (vii)
254 a computational complexity analysis (Section C.4), and (viii) IPL’s performance across different
255 SAE-based interpretability tools (Section D.1).

256 4.1 Experimental Settings

257 **Datasets.** We evaluate on the 15 standard image-classification datasets used by the prompt-learning
258 literature: For Base to novel generalization, Few-shot Learning, and cross-dataset evaluation we use
259 11 datasets ImageNet [8], Caltech101 [22], OxfordPets [29], StanfordCars [21], OxfordFlowers [28],
260 Food101 [3], FGVCAircraft [27], SUN397 [38], DTD [6], EuroSAT [16], and UCF101 [33]. For
261 domain generalization we use ImageNetV2 [31], ImageNet-Sketch [36], ImageNet-A [18], and
262 ImageNet-R [17] as target domains.

263 **Interpretability Methods.** We instantiate IPL with three interpretability tool variants. VL-SAE [32]
264 provides three modalities: VL-SAE-vision (vision decoder), VL-SAE-text (text decoder), and VL-
265 SAE-unified (concatenation of both decoders, see Section 3.2.1).

266 **Hyperparameter details.** All IPL runs use $N = 5$ representation tokens, $\alpha = 0.7$, soft injection
267 mode, and inject prompt tokens at layers $\mathcal{L}_{\text{rep}} = \{6, 7, 8, 9, 10, 11, 12\}$ in both vision and text
268 encoders. We train for 10 epochs with a learning rate of 0.0035 (cosine schedule, 1-epoch warmup),
269 batch size 4, and SGD optimizer following the convention of MMRL [15]. All our experiments
270 use CLIP ViT-B/16 as the frozen vision-language backbone. Image and text encoder weights from
271 OpenAI’s released checkpoints are used throughout, with **no fine-tuning of CLIP at any point**. All
272 experiments use 16-shot training and results reported over 3 seeds.

273 **Baselines.** We compare IPL against two categories of baselines. black-box prompt-learning: which
274 inject randomly initialized prompt tokens into CLIP encoder layer and optimize over them like
275 zero-shot CLIP, CoOp [45], CoCoOp [44], MaPLe [20], PromptSRC [4], MMA [39], TCP [41],
276 2SFS [11], SkipT [37], MMRL [15], and VaMP [5]. interpretable prompt-learning: which inject the
277 LLM queries attributes into CLIP encoder layers and optimize over them like CPL [42], ArGue-N [35],
278 IntCoOp [14], and IPO [10]. Numbers are reported from original papers.

Base to novel generalization In this experiment, we split each dataset equally into base and novel classes across 11 datasets. During training, only the base classes are used, while evaluation is performed on both base and novel classes. This setting allows us to examine how interpretability-signal-guided prompts affect VLM generalization to both seen and unseen classes.

In Table 1, we report the performance of IPL using vision-only interpretability signals (IPL(vision)), text-only interpretability signals (IPL (text)), and unified interpretability signals (IPL (unified)). We observe that IPL (vision) and IPL (text) achieve gains over previous interpretable prompt-learning methods on the base and novel classes for a few datasets, such as StanfordCars[21], UCF101[33], and Caltech101[22]. On average, they also improve base-class performance. However, their performance comparatively drops on several other datasets, particularly on novel classes. This is likely because interpretability signals derived from a single modality may not be sufficient to produce well-aligned logits that improve generalization.

In contrast, IPL with unified signals achieves the strongest average results on both base and novel classes. Notably, IPL (unified) also outperforms the black-box state-of-the-art method VaMP on several datasets, including ImageNet[8], Caltech101[22], OxfordPets[29], Food101[3], and UCF101[33].

Domain generalization In this experiment, we evaluate whether interpretability-guided prompt learning enables VLMs to generalize to distribution shifts during inference. We use ImageNet as the source dataset for training and evaluate the learned model on target ImageNet variants, including ImageNet-V2[31], ImageNet-Sketch[36], ImageNet-A[18], and ImageNet-R[17]. Table 2 shows that IPL (vision) and IPL (text) achieve strong source-domain accuracy on ImageNet, but show slightly weaker generalization on ImageNet-V2, ImageNet-Sketch, and ImageNet-R compared to previous interpretability-based state-of-the-art methods. IPL (unified) achieves strong results, outperforming prior interpretability-based state-of-the-art methods across all target domains and on the average target-domain accuracy.

Few-shot learning In this experiment, we evaluate IPL’s ability to adapt under different few-shot regimes using all the classes of ImageNet dataset. From Figure 2 We observe that IPL (unified) consistently outperforms previous baselines across different numbers of shots. Since only a limited number of prior methods report few-shot learning results, we compare IPL against the baselines for which such evaluations were reported. Results on other datasets in Section B.

Cross-dataset evaluation In this experiment, we evaluate the transferability of IPL to different target datasets during evaluation. Following the setting of CoCoOp [44], we train the source model on ImageNet and evaluate its performance on 10 other datasets at inference time in the few-shot setting as shown in Table 3. Similar to trends observed in base-to-novel generalization and domain generalization, IPL (unified) achieves the better overall performance across the target datasets.

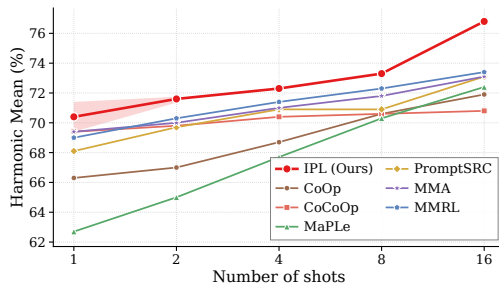


Figure 2: Few-shot performance on ImageNet under the base-to-novel protocol. IPL consistently improves harmonic mean over base and novel classes across shot counts, with the largest gain at 16 shots. Shaded bands denote \pm std over 3 seeds.

Table 2: Domain generalization. All methods are trained on ImageNet (16-shot, ViT-B/16) and evaluated on four ImageNet variants. Results are top-1 accuracy (%). **Best** and **second-best** per column. The leftmost column groups methods into *Interpretable* (concept/attribute-based prompt learning) and *Blackbox* (other prompt learning).

| Method | Source | | Target | | | | Avg. |
|---------------------------|------------------|------------------|------------------|------------------|------------------|------------------|------|
| | ImageNet | -V2 | -Sketch | -A | -R | | |
| IntCoOp ^l [14] | 71.85 | 65.21 | 49.20 | 51.55 | 76.88 | 60.71 | |
| ArGue-N ⁱ [35] | 71.84 | 65.02 | 49.25 | 51.47 | 76.96 | 60.68 | |
| IPL (Vision) | 73.63 \pm 0.05 | 64.93 \pm 0.09 | 48.10 \pm 0.08 | 51.90 \pm 0.14 | 75.90 \pm 0.37 | 60.20 \pm 0.05 | |
| IPL (Text) | 73.57 \pm 0.05 | 64.77 \pm 0.12 | 48.07 \pm 0.05 | 51.73 \pm 0.49 | 75.67 \pm 0.25 | 60.06 \pm 0.15 | |
| IPL (Unified) | 74.50 \pm 0.08 | 65.83 \pm 0.12 | 50.17 \pm 0.17 | 52.13 \pm 0.21 | 77.87 \pm 0.40 | 61.50 \pm 0.04 | |
| CLIP [30] | 66.73 | 60.83 | 46.15 | 47.77 | 73.96 | 57.18 | |
| CoOp [45] | 71.51 | 64.20 | 47.99 | 49.71 | 75.21 | 59.28 | |
| CoCoOp [44] | 71.02 | 64.07 | 48.75 | 50.63 | 76.18 | 59.91 | |
| MaPLe [20] | 70.72 | 64.07 | 49.15 | 50.90 | 76.98 | 60.26 | |
| PromptSRC [4] | 71.27 | 64.35 | 49.55 | 50.90 | 77.80 | 60.65 | |
| MMRL [15] | 72.03 | 64.47 | 49.17 | 51.20 | 77.53 | 60.59 | |
| VaMP [5] | 72.83 | 64.96 | 49.69 | 51.97 | 78.01 | 61.16 | |

Blackbox CLIP [30], CoOp [45], CoCoOp [44], MaPLe [20], PromptSRC [4], MMRL [15], VaMP [5]

Interpretable IntCoOp^l [14], ArGue-Nⁱ [35], IPL (Vision), IPL (Text), IPL (Unified)

Table 3: Cross-dataset transfer. All methods are trained on ImageNet (16-shot, ViT-B/16) and evaluated zero-shot on 10 downstream datasets. Top-1 accuracy (%) averaged over 3 seeds. The methods are grouped into *Interpretable* and *Blackbox*.

| Method | Source | | | | | Target | | | | | Avg. | |
|----------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| | ImageNet | Caltech101 | OxfordPets | StanfordCars | Flowers102 | Food101 | Aircraft | SUN397 | DTD | EuroSAT | | UCF101 |
| <i>Interpretable</i> | | | | | | | | | | | | |
| CPL [42] | 73.53 | 95.52 | 91.64 | 66.17 | 73.35 | 87.68 | 27.36 | 68.24 | 48.96 | 51.25 | 70.52 | 68.07 |
| IPO [10] | 72.15 | 94.34 | 90.96 | 66.10 | 72.75 | 86.75 | 25.14 | 67.97 | 47.01 | 48.56 | 69.23 | 67.36 |
| IPL (Vision) | 73.70\pm0.05 | 94.80\pm0.61 | 91.60\pm0.57 | 65.20\pm0.50 | 72.20\pm0.31 | 85.20\pm0.29 | 25.30\pm0.26 | 66.60\pm0.24 | 44.90\pm1.63 | 53.80\pm3.06 | 68.60\pm0.73 | 66.82\pm0.33 |
| IPL (Text) | 73.60\pm0.05 | 94.90\pm0.42 | 91.30\pm0.37 | 65.00\pm0.45 | 72.80\pm0.83 | 84.90\pm0.12 | 26.00\pm0.59 | 66.60\pm0.31 | 45.00\pm0.37 | 54.50\pm4.21 | 68.40\pm0.46 | 66.94\pm0.45 |
| IPL (Unified) | 73.70\pm0.08 | 95.90\pm0.39 | 91.30\pm0.33 | 66.80\pm0.79 | 73.40\pm0.08 | 87.60\pm0.17 | 27.40\pm0.90 | 68.80\pm0.43 | 45.90\pm0.43 | 54.50\pm4.21 | 70.70\pm0.66 | 68.23\pm0.60 |
| <i>Blackbox</i> | | | | | | | | | | | | |
| CoOp [45] | 71.51 | 93.70 | 89.14 | 64.51 | 68.71 | 85.30 | 18.47 | 64.15 | 41.92 | 46.39 | 66.55 | 63.88 |
| CoCoOp [44] | 71.02 | 94.43 | 90.14 | 65.32 | 71.88 | 86.06 | 22.94 | 67.36 | 45.73 | 45.37 | 68.21 | 65.74 |
| MaPLc [20] | 70.72 | 93.53 | 90.49 | 65.57 | 72.23 | 86.20 | 24.74 | 67.01 | 46.49 | 48.06 | 68.69 | 66.30 |
| PromptSRC [4] | 71.27 | 93.60 | 90.25 | 65.70 | 70.25 | 86.15 | 23.90 | 67.10 | 46.87 | 45.50 | 68.75 | 65.81 |
| TCP [41] | 71.40 | 93.97 | 91.25 | 64.69 | 71.21 | 86.69 | 23.45 | 67.15 | 44.35 | 51.45 | 68.73 | 66.29 |
| MMA [39] | 71.00 | 93.80 | 90.30 | 66.13 | 72.07 | 86.12 | 25.33 | 68.17 | 46.57 | 49.24 | 68.32 | 66.61 |
| MMRL [15] | 72.03 | 94.67 | 91.43 | 66.10 | 72.77 | 86.40 | 26.30 | 67.57 | 45.90 | 53.10 | 68.27 | 67.25 |
| VaMP [5] | 72.83 | 94.96 | 91.79 | 66.10 | 73.18 | 86.97 | 26.76 | 68.04 | 46.82 | 53.82 | 68.93 | 67.74 |

5 Qualitative Analysis

333 **Qualitative visualization.** Figure 3 traces the interpretability chain that IPL exploits for prompt construction, illustrated on an Oxford Flowers dataset. We use VL-SAE [32] as the interpretability tool, which provides a pretrained shared latent space along with human-readable concept labels for each decoder direction. Initially we rank concepts by mean SAE activation $\bar{z}_i(\mathbf{h})$ over the dataset’s base classes (Definition 2) yields the unit-norm directions $\mathbf{v}_i \in \mathbf{E}^{\text{unified}}$ that populate IPL’s concept dictionary, which are visually meaningful and dataset-relevant (*pointed floret tips*, *pinched bud*). A learned prompt token $\mathbf{I}^{(\ell)}$ at injection layer ℓ then composes these directions via attention, with cosine contributions revealing that $\mathbf{I}^{(\ell)}$ re-weights concepts beyond their raw activation ranking, indicating that the multi-head attention selector learns task-discriminative combinations rather than just copying $\bar{z}(\mathbf{h})$. Comparing Grad-CAM saliency, vanilla CLIP attends to semantically irrelevant regions, while IPL concentrates on the flower head, suggesting that interpretability-guided prompts induce more semantically grounded attention.

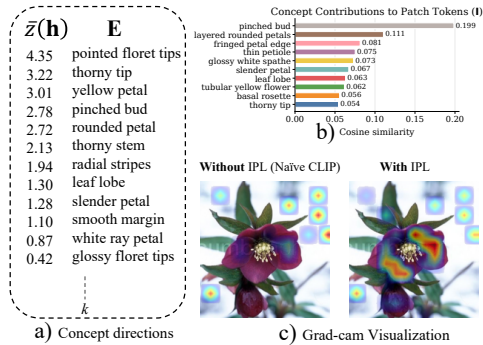


Figure 3: (a) top-ranked concept directions in \mathbf{E} with mean activations $\bar{z}(\mathbf{h})$ on oxford flowers dataset, (b) contributions of these concepts to a learned prompt token \mathbf{I} in terms of cosine similarity, and (c) Grad-CAM heatmaps showing that IPL focuses attention on the semantic features flower head while naive CLIP does not.

6 Conclusions

359 We presented IPL, a prompt-learning framework that leverages SAE based mechanistic interpretability
360 tools as active training guides. By extracting concept directions from frozen sparse autoencoders,
361 composing them into prompt tokens via learnable attention selector, and preserving their geometry
362 through a projection regularizer, IPL grounds prompt learning in CLIP’s own interpretable structure
363 rather than externally generated attributes. Across base-to-novel generalization, domain generaliza-
364 tion, cross-dataset transfer, and few-shot learning, IPL with unified concept directions performs better
365 in average accuracy compared to most of the previous state-of-the-art methods. As a limitation we
366 acknowledge that IPL’s effectiveness depends on the quality of the underlying interpretability tool.
367 For detailed analysis of our other limitations please refer to appendix Section C.4.

368 **References**

- 369 [1] Aurélien Bellet, Amaury Habrard, and Marc Sebban. A survey on metric learning for feature
370 vectors and structured data. *arXiv preprint arXiv:1306.6709*, 2013.
- 371 [2] Usha Bhalla, Alex Oesterling, Suraj Srinivas, Flavio P Calmon, and Himabindu Lakkaraju.
372 Interpreting clip with sparse linear concept embeddings (splice). *Advances in Neural Information*
373 *Processing Systems*, 37:84298–84328, 2024.
- 374 [3] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative
375 components with random forests. In *European conference on computer vision*, pages 446–461.
376 Springer, 2014.
- 377 [4] Guangyi Chen, Weiran Yao, Xiangchen Song, Xinyue Li, Yongming Rao, and Kun Zhang.
378 Plot: Prompt learning with optimal transport for vision-language models. *arXiv preprint*
379 *arXiv:2210.01253*, 2022.
- 380 [5] Silin Cheng and Kai Han. Vamp: Variational multi-modal prompt learning for vision-language
381 models. *arXiv preprint arXiv:2511.22664*, 2025.
- 382 [6] Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi.
383 Describing textures in the wild. In *Proceedings of the IEEE conference on computer vision and*
384 *pattern recognition*, pages 3606–3613, 2014.
- 385 [7] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. Sparse autoen-
386 coders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*,
387 2023.
- 388 [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
389 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
390 *recognition*, pages 248–255. Ieee, 2009.
- 391 [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
392 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
393 An image is worth 16x16 words: Transformers for image recognition at scale. In *International*
394 *Conference on Learning Representations*.
- 395 [10] Yingjun Du, Wenfang Sun, and Cees G Snoek. Ipo: Interpretable prompt optimization for vision-
396 language models. *Advances in Neural Information Processing Systems*, 37:126725–126766,
397 2024.
- 398 [11] Matteo Farina, Massimiliano Mancini, Giovanni Iacca, and Elisa Ricci. Rethinking few-shot
399 adaptation of vision-language models in two stages. In *Proceedings of the IEEE/CVF Conference*
400 *on Computer Vision and Pattern Recognition*, pages 29989–29998, 2025.
- 401 [12] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adap-
402 tation of deep networks. In *International conference on machine learning*, pages 1126–1135.
403 PMLR, 2017.
- 404 [13] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
405 Laviolette, Mario March, and Victor Lempitsky. Domain-adversarial training of neural networks.
406 *Journal of machine learning research*, 17(59):1–35, 2016.
- 407 [14] Soumya Suvra Ghosal, Samyadeep Basu, Soheil Feizi, and Dinesh Manocha. Intcoop:
408 Interpretability-aware vision-language prompt tuning. In *Proceedings of the 2024 Confer-*
409 *ence on Empirical Methods in Natural Language Processing*, pages 19584–19601, 2024.
- 410 [15] Yuncheng Guo and Xiaodong Gu. Mmrl: Multi-modal representation learning for vision-
411 language models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*,
412 pages 25015–25025, 2025.
- 413 [16] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel
414 dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal*
415 *of Selected Topics in Applied Earth Observations and Remote Sensing*, 12(7):2217–2226, 2019.

- 416 [17] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo,
417 Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness:
418 A critical analysis of out-of-distribution generalization. In *Proceedings of the IEEE/CVF*
419 *international conference on computer vision*, pages 8340–8349, 2021.
- 420 [18] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural
421 adversarial examples. In *Proceedings of the IEEE/CVF conference on computer vision and*
422 *pattern recognition*, pages 15262–15271, 2021.
- 423 [19] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan
424 Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning
425 with noisy text supervision. In *International conference on machine learning*, pages 4904–4916.
426 PMLR, 2021.
- 427 [20] Muhammad Uzair Khattak, Hanoona Rasheed, Muhammad Maaz, Salman Khan, and Fa-
428 had Shahbaz Khan. Maple: Multi-modal prompt learning. In *Proceedings of the IEEE/CVF*
429 *conference on computer vision and pattern recognition*, pages 19113–19122, 2023.
- 430 [21] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for
431 fine-grained categorization. In *4th International IEEE Workshop on 3D Representation and*
432 *Recognition (3dRR-13)*, Sydney, Australia, 2013.
- 433 [22] Fei-Fei Li, Marco Andreeto, Marc’Aurelio Ranzato, and Pietro Perona. Caltech 101, Apr 2022.
- 434 [23] Hyesu Lim, Jinho Choi, Jaegul Choo, and Steffen Schneider. Sparse autoencoders reveal
435 selective remapping of visual concepts during adaptation. *arXiv preprint arXiv:2412.05276*,
436 2024.
- 437 [24] Hong Liu, Jianmin Wang, and Mingsheng Long. Cycle self-training for domain adaptation.
438 *Advances in Neural Information Processing Systems*, 34:22968–22981, 2021.
- 439 [25] Wang Lu, Jindong Wang, Han Yu, Lei Huang, Xiang Zhang, Yiqiang Chen, and Xing Xie.
440 Fixed: Frustratingly easy domain generalization with mixup. In *Conference on Parsimony and*
441 *Learning*, pages 159–178. PMLR, 2024.
- 442 [26] Yuning Lu, Jianzhuang Liu, Yonggang Zhang, Yajing Liu, and Xinmei Tian. Prompt distri-
443 bution learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
444 *recognition*, pages 5206–5215, 2022.
- 445 [27] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-
446 grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 447 [28] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large
448 number of classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*,
449 Dec 2008.
- 450 [29] Omkar M. Parkhi, Andrea Vedaldi, Andrew Zisserman, and C. V. Jawahar. Cats and dogs. In
451 *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- 452 [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
453 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
454 models from natural language supervision. In *International conference on machine learning*,
455 pages 8748–8763. PmLR, 2021.
- 456 [31] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet
457 classifiers generalize to imagenet? In *International conference on machine learning*, pages
458 5389–5400. PMLR, 2019.
- 459 [32] Shufan Shen, Junshu Sun, Qingming Huang, and Shuhui Wang. Vl-sae: Interpreting and
460 enhancing vision-language alignment with a unified concept set. In *The Thirty-ninth Annual*
461 *Conference on Neural Information Processing Systems*.
- 462 [33] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human
463 actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.

- 464 [34] Baochen Sun, Jiashi Feng, and Kate Saenko. Correlation alignment for unsupervised domain
465 adaptation. In *Domain adaptation in computer vision applications*, pages 153–171. Springer,
466 2017.
- 467 [35] Xinyu Tian, Shu Zou, Zhaoyuan Yang, and Jing Zhang. Argue: Attribute-guided prompt tuning
468 for vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
469 and Pattern Recognition*, pages 28578–28587, 2024.
- 470 [36] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global repre-
471 sentations by penalizing local predictive power. *Advances in neural information processing
472 systems*, 32, 2019.
- 473 [37] Shihan Wu, Ji Zhang, Pengpeng Zeng, Lianli Gao, Jingkuan Song, and Heng Tao Shen. Skip
474 tuning: Pre-trained vision-language models are effective and efficient adapters themselves. In
475 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 14723–14732,
476 2025.
- 477 [38] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database:
478 Large-scale scene recognition from abbey to zoo. In *2010 IEEE computer society conference
479 on computer vision and pattern recognition*, pages 3485–3492. IEEE, 2010.
- 480 [39] Lingxiao Yang, Ru-Yuan Zhang, Yanchen Wang, and Xiaohua Xie. Mma: Multi-modal adapter
481 for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision
482 and pattern recognition*, pages 23826–23837, 2024.
- 483 [40] Hantao Yao, Rui Zhang, and Changsheng Xu. Visual-language prompt tuning with knowledge-
484 guided context optimization. In *Proceedings of the IEEE/CVF conference on computer vision
485 and pattern recognition*, pages 6757–6767, 2023.
- 486 [41] Hantao Yao, Rui Zhang, and Changsheng Xu. Tcp: Textual-based class-aware prompt tuning
487 for visual-language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision
488 and Pattern Recognition*, pages 23438–23448, 2024.
- 489 [42] Yi Zhang, Ce Zhang, Ke Yu, Yushun Tang, and Zhihai He. Concept-guided prompt learning for
490 generalization in vision-language models. In *Proceedings of the AAAI Conference on Artificial
491 Intelligence*, volume 38, pages 7377–7386, 2024.
- 492 [43] Han Zhao, Remi Tachet Des Combes, Kun Zhang, and Geoffrey Gordon. On learning invariant
493 representations for domain adaptation. In *International conference on machine learning*, pages
494 7523–7532. PMLR, 2019.
- 495 [44] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning
496 for vision-language models. In *Proceedings of the IEEE/CVF conference on computer vision
497 and pattern recognition*, pages 16816–16825, 2022.
- 498 [45] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for
499 vision-language models. *International journal of computer vision*, 130(9):2337–2348, 2022.

501 A Base to novel generalization results

Table 4: Comparison with state-of-the-art methods on base-to-novel generalization across 11 datasets. Methods are grouped into *Interpretable* and *Blackbox*. Our method (IPL) is reported using vision-only, text-only and unified interpretability signals using VL-SAE[32]. Bold indicates best per column and underline indicates second-best per column. Delta rows show absolute accuracy-point change over ArGue-N, with positive gains in green and negative changes in red. [†]IntCoOp does not report DTD; its average in the source paper is over 10 datasets.

| Method | Average | | | Caltech101 | | | OxfordPets | | | Flowers102 | | |
|---------------------------|--------------|-------------------|--------------|-------------------|-------------------|--------------|--------------|-------------------|--------------|--------------|--------------|--------------|
| | B | N | H | B | N | H | B | N | H | B | N | H |
| CPL [42] | 84.38 | 78.03 | 81.08 | 98.35 | 95.13 | 96.71 | 95.86 | 98.21 | 97.02 | 98.07 | 80.43 | 88.38 |
| IPO [10] | 79.92 | 80.51 | 80.21 | 97.32 | 95.23 | 96.26 | 95.21 | 98.23 | 96.70 | 96.78 | 78.32 | 86.58 |
| IntCoOp [†] [14] | 83.20 | 78.72 | 80.75 | 97.80 | 94.76 | 96.25 | 95.92 | 98.20 | 97.04 | 97.82 | 75.54 | 85.24 |
| ArGue-N [35] | 83.77 | 78.74 | 81.18 | 98.63 | 94.70 | 96.63 | 96.23 | 98.59 | 97.40 | 98.62 | 77.96 | 87.08 |
| IPL (Vision) | 85.80±0.42 | 77.15±1.09 | 81.25 | <u>99.10±0.19</u> | 95.50±0.75 | 97.27 | 95.90±0.24 | 97.70±0.60 | 96.79 | 98.70±0.09 | 76.60±0.45 | 86.26 |
| Δ vs ArGue-N | +2.03 | -1.59 | +0.07 | +0.47 | +0.80 | +0.64 | -0.33 | -0.89 | -0.61 | +0.08 | -1.36 | -0.82 |
| IPL (Text) | 85.75±0.36 | 76.59±0.83 | 80.91 | 99.30±0.21 | 94.50±0.38 | 96.84 | 96.30±0.33 | 97.10±0.59 | 96.70 | 98.60±0.16 | 77.10±0.57 | 86.53 |
| Δ vs ArGue-N | +1.98 | -2.15 | -0.27 | +0.67 | -0.20 | +0.21 | +0.07 | -1.49 | -0.70 | -0.02 | -0.86 | -0.55 |
| IPL (Unified) | 86.26±0.40 | 81.23±1.05 | 83.67 | 99.50±0.24 | 96.40±0.75 | 97.93 | 96.60±0.62 | 98.60±0.29 | 97.59 | 98.80±0.37 | 81.80±0.38 | 89.50 |
| Δ vs ArGue-N | +2.49 | +2.49 | +2.49 | +0.87 | +1.70 | +1.30 | +0.37 | +0.01 | +0.19 | +0.18 | +3.84 | +2.42 |
| CLIP [30] | 69.34 | 74.22 | 71.70 | 96.84 | 94.00 | 95.40 | 91.17 | 97.26 | 94.12 | 72.08 | 77.80 | 74.83 |
| CoOp [45] | 82.69 | 63.22 | 71.66 | 98.00 | 89.81 | 93.73 | 93.67 | 95.29 | 94.47 | 97.60 | 59.67 | 74.06 |
| CoCoOp [44] | 80.47 | 71.69 | 75.83 | 97.96 | 93.81 | 95.84 | 95.20 | 97.69 | 96.43 | 94.87 | 71.75 | 81.71 |
| ProDA [26] | 81.56 | 72.30 | 76.65 | 98.27 | 93.23 | 95.68 | 95.43 | 97.83 | 96.62 | 97.70 | 68.68 | 80.66 |
| KgCoOp [40] | 80.73 | 73.60 | 77.00 | 97.72 | 94.39 | 96.03 | 94.65 | 97.76 | 96.18 | 95.00 | 74.73 | 83.65 |
| MaPLe [20] | 82.28 | 75.14 | 78.55 | 97.74 | 94.36 | 96.02 | 95.43 | 97.76 | 96.58 | 95.92 | 72.46 | 82.56 |
| PromptSRC [4] | 84.26 | 76.10 | 79.97 | 98.10 | 94.03 | 96.02 | 95.33 | 97.30 | 96.30 | 98.07 | 76.50 | 85.95 |
| TCP [41] | 84.13 | 75.36 | 79.51 | 98.23 | 94.67 | 96.42 | 94.67 | 97.20 | 95.92 | 97.73 | 75.57 | 85.23 |
| MMA [39] | 83.20 | 76.80 | 79.87 | 98.40 | 94.00 | 96.15 | 95.40 | 98.07 | 96.72 | 97.77 | 75.93 | 85.48 |
| 2SFS [11] | 85.55 | 75.48 | 80.20 | 98.71 | 94.43 | 96.52 | 95.32 | 97.82 | 96.55 | 98.29 | 76.17 | 85.83 |
| SkipT [37] | 85.04 | 77.53 | 81.11 | 98.50 | 95.33 | 96.89 | 95.70 | 97.87 | 96.77 | 98.57 | 75.80 | 85.70 |
| MMRL [15] | 85.68 | 77.16 | 81.20 | 98.97 | 94.50 | 96.68 | 95.90 | 97.60 | 96.74 | 98.97 | 77.27 | 86.78 |
| VaMP [5] | 86.45 | 78.67 | <u>82.37</u> | 98.95 | <u>95.96</u> | <u>97.43</u> | 96.95 | 95.24 | 96.08 | <u>98.96</u> | 83.97 | 90.85 |

| Method | Food101 | | | FGVCAircraft | | | DTD | | | EuroSAT | | |
|---------------------------|--------------|-------------------|--------------|--------------|--------------|--------------|-------------------|--------------|--------------|-------------------|---------------|--------------|
| | B | N | H | B | N | H | B | N | H | B | N | H |
| CPL [42] | 91.92 | 93.87 | 92.88 | 42.27 | 38.85 | 40.49 | 80.92 | 62.27 | 70.38 | 94.18 | 81.05 | 87.12 |
| IPO [10] | 90.92 | 93.08 | 91.99 | 41.21 | 41.42 | 41.31 | 82.14 | 66.81 | 73.69 | 94.25 | 80.11 | 86.61 |
| IntCoOp [†] [14] | 91.45 | 91.99 | 91.72 | 38.55 | 35.90 | 37.17 | — | — | — | 95.26 | 78.01 | 85.77 |
| ArGue-N [35] | 91.42 | 92.40 | 91.91 | 41.29 | 38.80 | 40.01 | 80.33 | 67.03 | 73.08 | 95.10 | 90.68 | 92.84 |
| IPL (Vision) | 90.70±0.12 | 91.50±0.05 | 91.10 | 46.30±0.66 | 38.20±1.02 | 41.86 | 86.10±0.61 | 64.50±1.55 | 73.75 | 96.00±0.53 | 79.60±6.25 | 87.03 |
| Δ vs ArGue-N | -0.72 | -0.90 | -0.81 | +5.01 | -0.60 | +1.85 | +5.77 | -2.53 | +0.67 | +0.90 | -11.08 | -5.81 |
| IPL (Text) | 90.90±0.22 | 91.50±0.00 | 91.20 | 45.90±1.25 | 37.90±0.62 | 41.52 | 85.00±0.39 | 64.70±0.95 | 73.47 | 95.60±0.42 | 74.10±4.78 | 83.49 |
| Δ vs ArGue-N | -0.52 | -0.90 | -0.71 | +4.61 | -0.90 | +1.51 | +4.67 | -2.33 | +0.39 | +0.50 | -16.58 | -9.35 |
| IPL (Unified) | 90.60±0.05 | 97.70±0.05 | 94.02 | 46.50±0.57 | 42.20±0.70 | 44.25 | 86.30±0.29 | 66.50±1.87 | <u>75.12</u> | 97.10±0.48 | 87.50±5.34 | <u>92.05</u> |
| Δ vs ArGue-N | -0.82 | +5.30 | +2.11 | +5.21 | +3.40 | +4.24 | +5.97 | -0.53 | +2.04 | +2.00 | -3.18 | -0.79 |
| CLIP [30] | 90.10 | 91.22 | 90.66 | 27.19 | 36.29 | 31.09 | 53.24 | 59.90 | 56.37 | 56.48 | 64.05 | 60.03 |
| CoOp [45] | 88.33 | 82.26 | 85.19 | 40.44 | 22.30 | 28.75 | 79.44 | 41.18 | 54.24 | 92.19 | 54.74 | 68.69 |
| CoCoOp [44] | 90.70 | 91.29 | 90.99 | 33.41 | 23.71 | 27.74 | 77.01 | 56.00 | 64.85 | 87.49 | 60.04 | 71.21 |
| ProDA [26] | 90.30 | 88.57 | 89.43 | 36.90 | 34.13 | 35.46 | 80.67 | 56.48 | 66.44 | 83.90 | 66.00 | 73.88 |
| KgCoOp [40] | 90.50 | 91.70 | 91.09 | 36.21 | 33.55 | 34.83 | 77.55 | 54.99 | 64.35 | 85.64 | 64.34 | 73.48 |
| MaPLe [20] | 90.71 | 92.05 | 91.38 | 37.44 | 35.61 | 36.50 | 80.36 | 59.18 | 68.16 | 94.07 | 73.23 | 82.35 |
| PromptSRC [4] | 90.67 | 91.53 | 91.10 | 42.73 | 37.87 | 40.15 | 83.37 | 62.97 | 71.75 | 92.90 | 73.90 | 82.32 |
| TCP [41] | 90.57 | 91.37 | 90.97 | 41.97 | 34.43 | 37.83 | 82.77 | 58.07 | 68.25 | 91.63 | 74.73 | 82.32 |
| MMA [39] | 90.13 | 91.30 | 90.71 | 40.57 | 36.33 | 38.33 | 83.20 | 65.63 | 73.38 | 85.46 | 82.34 | 83.87 |
| 2SFS [11] | 90.71 | 91.34 | 91.02 | 47.48 | 35.51 | 40.63 | 84.60 | 65.01 | 73.52 | 96.91 | 67.09 | 79.29 |
| SkipT [37] | 90.67 | 92.03 | 91.34 | 45.37 | 37.13 | 40.84 | 83.77 | 67.23 | 74.59 | 92.47 | 83.00 | 87.48 |
| MMRL [15] | 90.57 | 91.50 | 91.03 | 46.30 | 37.03 | 41.15 | 85.67 | 65.00 | 73.82 | 95.60 | 80.17 | 87.21 |
| VaMP [5] | 92.77 | 93.16 | 92.96 | 46.77 | 41.13 | <u>43.76</u> | <u>86.14</u> | 67.20 | 75.50 | 95.78 | 77.21 | 85.49 |

502 Table 1 reports IPL’s base-to-novel results on the remaining datasets, complementing the subset
503 shown in the main paper (Table 4). The trends mirror those in the main results: IPL with unified
504 concept directions consistently achieves the strongest harmonic mean, while the vision-only and
505 text-only variants remain competitive with prior interpretable prompt-learning methods.

506 B Results on few shot learning other datasets.

507 In Table 5, we study performance of few-shot learning with (1, 2, 4, 8, and 16) shots for different
508 datasets. We observe that IPL performs better than previous blackbox based prompt learning methods.
509 Note that interpretability based prompt learning methods do not provide results for few-shot learning.

Table 5: Few-shot classification accuracies across selected datasets under 1, 2, 4, 8, and 16 shots.

| Method | Average | | | | | Caltech101 | | | | | OxfordPets | | | | | Flowers102 | | | | |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 |
| IPL (Vision) | 69.51 | 72.54 | 77.06 | 81.15 | 84.62 | 93.93 | 94.63 | 95.57 | 96.37 | 96.63 | 88.10 | 91.33 | 92.70 | 93.20 | 93.77 | 76.50 | 83.23 | 88.90 | 94.57 | 96.93 |
| IPL (Text) | 69.47 | 72.82 | 76.15 | 80.75 | 84.95 | 93.97 | 94.43 | 95.77 | 96.13 | 96.63 | 88.30 | 90.97 | 93.03 | 92.53 | 93.80 | 76.40 | 83.53 | 89.97 | 94.50 | 96.97 |
| IPL (Unified) | 74.03 | 77.29 | 80.47 | 82.54 | 87.07 | 93.73 | 94.53 | 96.80 | 96.53 | 97.73 | 90.60 | 91.83 | 92.83 | 93.13 | 94.80 | 86.37 | 91.33 | 95.77 | 96.50 | 98.90 |
| Linear probe CLIP [30] | 48.72 | 60.44 | 70.51 | 76.41 | 80.51 | 79.88 | 89.01 | 92.05 | 93.41 | 95.43 | 44.06 | 58.37 | 71.17 | 78.36 | 85.34 | 69.74 | 85.07 | 92.02 | 96.10 | 97.37 |
| CoOp [45] | 67.34 | 71.37 | 74.85 | 77.87 | 80.99 | 92.60 | 93.07 | 94.40 | 94.37 | 95.57 | 90.37 | 89.80 | 92.57 | 91.27 | 91.87 | 77.53 | 87.33 | 92.17 | 94.97 | 97.07 |
| CoCoOp [44] | 65.63 | 66.21 | 71.21 | 73.35 | 75.88 | 93.83 | 94.82 | 94.98 | 95.04 | 95.16 | 91.27 | 92.64 | 92.81 | 93.45 | 93.34 | 72.08 | 75.79 | 78.40 | 84.30 | 87.84 |
| MaPLE [20] | 70.88 | 74.28 | 77.31 | 80.49 | 83.32 | 92.57 | 93.97 | 94.43 | 95.20 | 96.00 | 89.10 | 90.87 | 91.90 | 92.57 | 92.83 | 83.30 | 88.93 | 92.67 | 95.80 | 97.00 |
| PromptSRC [4] | 73.36 | 76.42 | 79.72 | 82.04 | 84.40 | 93.67 | 94.53 | 95.27 | 95.67 | 96.07 | 92.00 | 92.50 | 93.43 | 93.50 | 93.67 | 85.93 | 91.17 | 93.87 | 96.27 | 97.60 |
| MMA [39] | 69.55 | 72.48 | 77.51 | 80.91 | 84.38 | 92.90 | 94.00 | 94.33 | 95.37 | 96.33 | 91.23 | 91.97 | 92.23 | 92.77 | 93.23 | 83.60 | 90.30 | 93.00 | 95.97 | 97.97 |
| MMRL [15] | 73.83 | 77.37 | 80.71 | 82.95 | 86.09 | 94.17 | 94.83 | 96.03 | 96.27 | 97.13 | 90.87 | 91.57 | 92.57 | 93.03 | 93.83 | 85.97 | 91.20 | 94.60 | 96.60 | 98.40 |

| Method | Food101 | | | | | FGVCAircraft | | | | | DTD | | | | | EuroSAT | | | | |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 | 1 | 2 | 4 | 8 | 16 |
| IPL (Vision) | 85.13 | 85.47 | 86.23 | 86.97 | 87.50 | 29.10 | 34.40 | 37.47 | 44.17 | 52.70 | 52.43 | 56.93 | 64.93 | 69.58 | 73.87 | 61.40 | 61.77 | 73.63 | 83.20 | 90.93 |
| IPL (Text) | 85.07 | 85.47 | 86.27 | 87.03 | 87.43 | 28.80 | 34.43 | 37.40 | 44.27 | 54.33 | 50.07 | 56.90 | 64.63 | 69.87 | 74.23 | 63.67 | 64.03 | 65.97 | 80.93 | 91.27 |
| IPL (Unified) | 85.10 | 85.40 | 86.17 | 86.87 | 87.47 | 29.17 | 35.30 | 39.37 | 46.93 | 59.60 | 56.50 | 60.03 | 66.07 | 69.60 | 75.03 | 76.73 | 82.60 | 86.30 | 88.23 | 95.97 |
| Linear probe CLIP [30] | 43.96 | 61.51 | 73.19 | 79.79 | 82.90 | 19.61 | 26.41 | 32.33 | 39.35 | 45.36 | 34.59 | 40.76 | 55.71 | 63.46 | 69.96 | 49.23 | 61.98 | 77.09 | 84.43 | 87.21 |
| CoOp [45] | 84.33 | 84.40 | 84.47 | 82.67 | 84.20 | 21.37 | 26.20 | 30.83 | 39.00 | 43.40 | 50.23 | 53.60 | 58.70 | 64.77 | 69.87 | 54.93 | 65.17 | 70.80 | 78.07 | 84.93 |
| CoCoOp [44] | 85.65 | 86.22 | 86.88 | 86.97 | 87.25 | 12.68 | 15.06 | 24.79 | 26.61 | 31.21 | 48.54 | 52.17 | 55.04 | 58.89 | 63.04 | 55.33 | 46.74 | 65.56 | 68.21 | 73.32 |
| MaPLE [20] | 80.50 | 81.47 | 81.77 | 83.60 | 85.33 | 26.73 | 30.90 | 34.87 | 42.00 | 48.40 | 52.13 | 55.50 | 61.00 | 66.50 | 71.33 | 71.80 | 78.30 | 84.50 | 87.73 | 92.33 |
| PromptSRC [4] | 84.87 | 85.70 | 86.17 | 86.90 | 87.50 | 27.67 | 31.70 | 37.47 | 43.27 | 50.83 | 56.23 | 59.97 | 65.53 | 69.87 | 72.73 | 73.13 | 79.37 | 86.30 | 88.80 | 92.43 |
| MMA [39] | 83.03 | 82.50 | 82.13 | 83.00 | 84.57 | 28.73 | 31.90 | 37.57 | 44.83 | 52.70 | 52.27 | 56.90 | 63.93 | 67.97 | 73.47 | 55.07 | 59.80 | 79.40 | 86.47 | 92.37 |
| MMRL [15] | 84.87 | 85.53 | 85.77 | 86.33 | 87.03 | 28.53 | 34.23 | 40.47 | 48.07 | 57.60 | 56.37 | 61.37 | 67.87 | 71.60 | 75.30 | 76.00 | 82.87 | 87.67 | 88.73 | 93.37 |

510 C Ablation Studies

511 C.1 Losses.

512 Table 6 studies the effect of each loss component
 513 in IPL using the unified variant on ImageNet[8].
 514 The results show that the projection-preservation
 515 regularization losses, given by the vision- and
 516 text-side cosine losses $\mathcal{L}_{\text{cos}}^v$ and $\mathcal{L}_{\text{cos}}^t$, are impor-
 517 tant for improving generalization. The cross-
 518 entropy losses provide task-level supervision,
 519 with $\mathcal{L}_{\text{CE-rep}}$ further strengthening representa-
 520 tion learning. The best performance is obtained
 521 when all losses are used together, indicating that
 522 task discrimination and preservation of inter-
 523 pretable concept directions are complementary
 524 for balancing base-class adaptation and novel-class generalization.

Table 6: Loss-component ablation on ImageNet for IPL (unified). \checkmark indicates the loss is active during training.

| \mathcal{L}_{CE} | $\mathcal{L}_{\text{CE-rep}}$ | $\mathcal{L}_{\text{cos}}^v$ | $\mathcal{L}_{\text{cos}}^t$ | Base | Novel | HM |
|---------------------------|-------------------------------|------------------------------|------------------------------|------------------------------------|------------------------------------|--------------|
| \checkmark | \times | \times | \times | 72.43 \pm 0.90 | 68.70 \pm 0.10 | 70.52 |
| \checkmark | \checkmark | \times | \times | 77.93 \pm 0.29 | 69.40 \pm 0.44 | 73.42 |
| \checkmark | \times | \checkmark | \times | 72.27 \pm 0.71 | 68.77 \pm 0.29 | 70.47 |
| \checkmark | \times | \times | \checkmark | 72.47 \pm 0.67 | 70.30 \pm 0.26 | 71.37 |
| \checkmark | \checkmark | \checkmark | \times | 77.93 \pm 0.15 | 69.27 \pm 0.15 | 73.34 |
| \checkmark | \checkmark | \times | \checkmark | 78.27 \pm 0.25 | 70.60 \pm 0.20 | 74.24 |
| \checkmark | \times | \checkmark | \checkmark | 72.50 \pm 0.66 | 70.30 \pm 0.26 | 71.38 |
| \checkmark | \checkmark | \checkmark | \checkmark | 79.30 \pm 0.20 | 74.50 \pm 0.10 | 76.83 |

525 C.2 Layers for prompt token insertion.

526 Figure 4 analyzes the effect of inserting IPL
 527 prompt tokens at different CLIP encoder layers on
 528 ImageNet. We observe that mid-to-late layer inser-
 529 tion performs best, with layers 6-12 achieving the
 530 strongest harmonic mean. Layers 3-12 show a simi-
 531 lar trend, suggesting that interpretability-guided
 532 prompts are most effective when injected into se-
 533 mantic representations. In contrast, early-layer
 534 insertion leads to weaker generalization, likely
 535 because early layers encode lower-level visual fea-
 536 tures less suited for concept-guided adaptation.

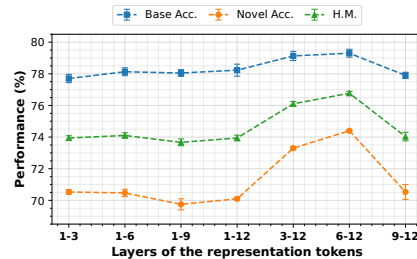


Figure 4: Effect of representation-token insertion layers on ImageNet base-to-novel generalization.

537 C.3 Different sizes of $k \in K$ for \mathbf{E} concepts.

538 In this ablation, we analyze the impact of number
 539 of concept directions \mathbf{E} used to construct prompt tokens. From Figure 5, we observe that $E =$
 540 256 achieves comparatively better harmonic mean, indicating that a compact set of high-ranked
 541 interpretable concept directions is sufficient for IPL’s prompt token construction. performance
 542 comparatively degrades with larger \mathbf{E} , suggesting that lower ranked concepts add noise rather than
 543 useful structure. We use $E = 256$ as the default in all main experiments.

544 **C.4 Limitations analysis and future scope**

545 IPL introduces additional compute relative to existing prompt-learning baselines. As reported in Table 7, IPL requires approximately 23.1 GMAC per forward pass compared to 20.9 GMAC for MMRL [15] and 20.5 GMAC for ArGue [35], a $\sim 10\text{--}13\%$ increase, attributable entirely to the per-layer multi-head attention selectors and $r2v/r2t$ projection layers introduced for interpretability-guided prompt construction. We note that this overhead is a deliberate consequence of IPL’s per-layer composition: the selector $\text{MHA}^{(\ell)}$ and projections are instantiated independently at each rep-injection layer ℓ , which provides per-layer flexibility in how concepts are composed at different depths but multiplicatively scales prompt-side FLOPs with $|\mathcal{L}_{\text{rep}}|$. Several straightforward extensions can mitigate this cost: (i) sharing MHA and projection parameters across injection layers, which would reduce prompt-side FLOPs by roughly $7\times$ in our default configuration at the cost of less per-layer specialization; (ii) injecting at fewer layers, which our ablation (Section C.2) shows incurs only a small accuracy drop; or (iii) lower-rank parameterizations of the projection layers. A systematic study of these efficiency-accuracy trade-offs is left to future work.

Table 7: Computational cost on CLIP ViT-B/16 (single forward pass). *Prompt* columns report FLOPs introduced by each method’s prompt-token machinery. Δ is the percentage increase in total GMAC over ArGue.

| Method | Vision (GMAC) | Text (GMAC) | Prompt (GMAC) | Total (GMAC) | Δ (%) |
|----------------------|---------------|-------------|---------------|--------------|--------------|
| ArGue [35] | 17.563 | 2.980 | 0.000 | 20.543 | — |
| MMRL [15] | 17.833 | 3.096 | 0.023 | 20.951 | +1.99 |
| IPL (Vision) | 17.833 | 3.096 | 2.203 | 23.132 | +12.60 |
| IPL (Text) | 17.833 | 3.096 | 2.203 | 23.132 | +12.60 |
| IPL (Unified) | 17.833 | 3.096 | 2.203 | 23.132 | +12.60 |

567 **D Hyperparameter details**

Table 8: Ablation analysis of α and λ on base-to-new generalization.

| α | Base | New | HM |
|----------|------------------------------------|------------------------------------|--------------|
| 0.3 | 76.30 \pm 0.20 | 73.96 \pm 0.16 | 75.11 |
| 0.5 | 77.63 \pm 0.30 | 73.84 \pm 0.08 | 75.69 |
| 0.7 | 79.30 \pm 0.29 | 74.50 \pm 0.06 | 76.83 |
| 0.9 | 77.41 \pm 0.35 | 73.29 \pm 0.10 | 75.29 |
| 1.0 | 76.73 \pm 0.25 | 69.78 \pm 0.15 | 73.09 |

| λ | α | Base | New | H |
|-----------|----------|--------------|--------------|--------------|
| 0.01 | 0.7 | 77.70 | 71.10 | 74.25 |
| 0.1 | 0.7 | 77.70 | 73.03 | 75.29 |
| 0.2 | 0.7 | 77.67 | 71.07 | 74.22 |
| 0.5 | 0.7 | 77.70 | 73.00 | 75.28 |
| 3.0 | 0.7 | 77.60 | 74.50 | 76.02 |
| 4.0 | 0.7 | 78.67 | 71.90 | 75.13 |
| 5.0 | 0.7 | 77.67 | 73.10 | 75.32 |
| 6.0 | 0.7 | 79.30 | 70.83 | 74.83 |
| 7.0 | 0.7 | 77.53 | 71.83 | 74.57 |

568 IPL is trained for 10 epochs with batch size 4, SGD, and a cosine schedule peaking at lr 0.0035 after a one-epoch warmup. We use $N=5$ rep tokens injected at layers $\{6, \dots, 12\}$, $\alpha=0.7$, and a concept budget k of 256. Results are averaged over 3 seeds at 16 shots. From Table 8, we observe that IPL is most effective at $\alpha = 0.7$, which balances the standard and prompt-token cross-entropy objectives, and at $\lambda_{\text{proj}} = 3.0$, where the projection-preservation loss provides sufficient geometric regularization without overconstraining the projection layers.

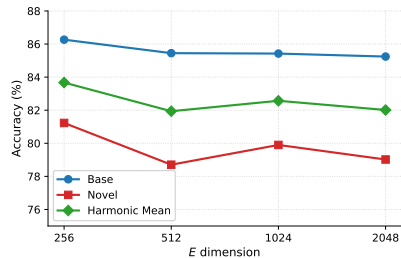


Figure 5: Effect of concept direction size \mathbf{E} averaged across 11 datasets for base to novel generalization.

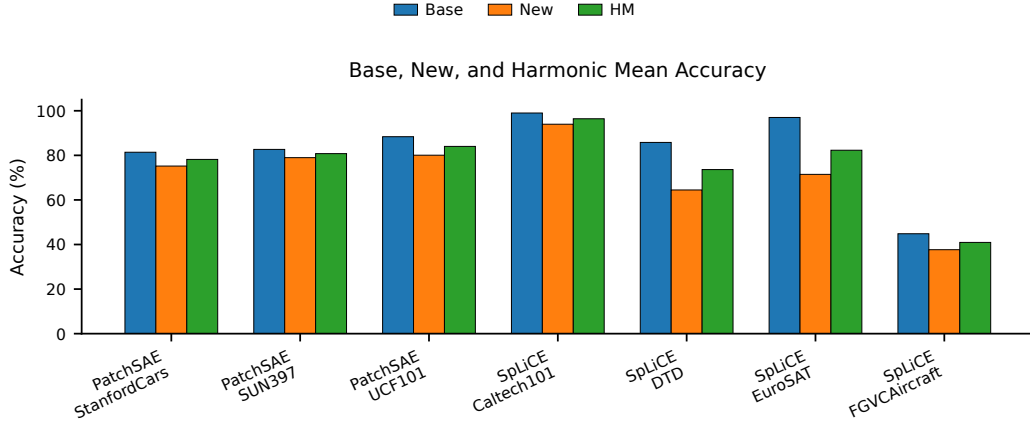


Figure 6: Performace of other interpretability tools like PatchSAE [23] and SpLiCE [2] on base to novel generalization for different datasets.

580 D.1 Results

581 on other Interpretability Methods

582 In Figure 6, we study how different interpretability tools like PatchSAE[23] and splice [2] perform
 583 on base to novel generalization for different dataset. We observe that these tools perform on par with
 584 VL-SAE (vision).

585 D.2 Representation-token count.

586 Figure 7 shows that performance improves rapidly as the number of representation tokens increases
 587 from 1 to 5 and then saturates. We therefore use $N_{\text{rep}} = 5$ throughout the main experiments, since
 588 larger values provide only marginal gains while increasing prompt length.

589 E Detailed literature survey

590 Distribution shifts and dealing with novelty.

591 Adapting models under distribution shift has a
 592 long history in classical machine learning, predat-
 593 ing vision-language models. *Domain adaptation*
 594 (DA) methods minimize feature-distribution diver-
 595 gence between a labeled source and unlabeled tar-
 596 get domain via discrepancy measures [13], adver-
 597 sarial alignment [34], or self-training on pseudo-
 598 labels [24]. *Domain generalization* (DG) extends
 599 this to the harder setting where the target domain
 600 is unseen at training time, with approaches based
 601 on domain-invariant feature learning [43], meta-
 602 learning across source domains [12], and feature
 603 augmentation or mixup [25]. *Few-shot learning*
 604 similarly has a rich literature on metric learning [1]
 605 and meta-learning [12] for adaptation under lim-
 606 ited supervision.

607 These classical formulations target a fundamen-
 608 tally different setting from ours: they assume full access to model weights and adapt task-specific
 609 networks, often trained from scratch, by reshaping their feature space to bridge source-target gaps.
 610 In contrast, prompt-learning for vision-language models, the regime IPL operates in and adapts a
 611 single *frozen* foundation model (CLIP) by learning a small set of continuous prompt tokens, without
 612 modifying any backbone parameters. The challenge is therefore not to learn a domain-invariant

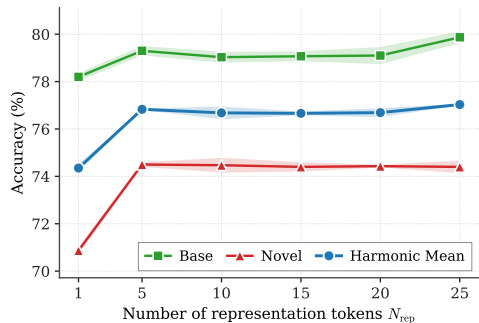


Figure 7: Effect of the number of representation tokens N_{rep} on ImageNet base-to-novel generalization. Performance saturates at $N_{\text{rep}} = 5$, which we use in all main experiments. Shaded bands denote \pm std over 3 seeds.

613 representation from data, but to steer a pretrained representation that is already broadly capable toward
614 task-, domain-, or class-specific behavior without overwriting its zero-shot generalization. Recent
615 work formalizes the various distributional axes along which this steering must succeed: base-to-novel
616 generalization [44, 20, 4] (unseen classes within a dataset), domain generalization [31, 36, 18, 17]
617 (input shifts at fixed label space), cross-dataset transfer [44] (zero-shot transfer to disjoint label
618 spaces), and few-shot learning under limited supervision. Existing prompt-learning methods often
619 optimize design choices for a single protocol [40, 4, 41, 35, 14]. IPL treats all four protocols as
620 instances of a single underlying problem, adapting prompts under limited supervision while preserv-
621 ing CLIP’s pretrained generalization, and addresses them with one mechanism: anchoring prompt
622 tokens in interpretable concept directions discovered inside CLIP itself. This unified framing avoids
623 per-protocol architectural changes, and is empirically borne out by IPL’s consistent gains across all
624 four evaluation regimes (Section 4.2).