

Shuvendu Roy, Ali Etemad
Queen's University, Kingston, Canada

Introduction

We present SelfPrompt, a novel semi-supervised prompt-tuning approach for tuning vision-language models (VLMs) in a semi-supervised learning setup. Existing methods for tuning VLMs in semi-supervised setups struggle with efficiently using the limited label set budget, accumulating noisy pseudo-labels, and properly utilizing unlabeled data. SelfPrompt addresses these challenges by introducing (a) a weakly-supervised sampling technique that selects a diverse and representative labelled set, (b) a cluster-guided pseudo-labelling method that improves pseudo-label accuracy, and (c) a confidence-aware semi-supervised learning module that maximizes the utilization of unlabeled data by learning from high- and low-confidence pseudo-labels differently.

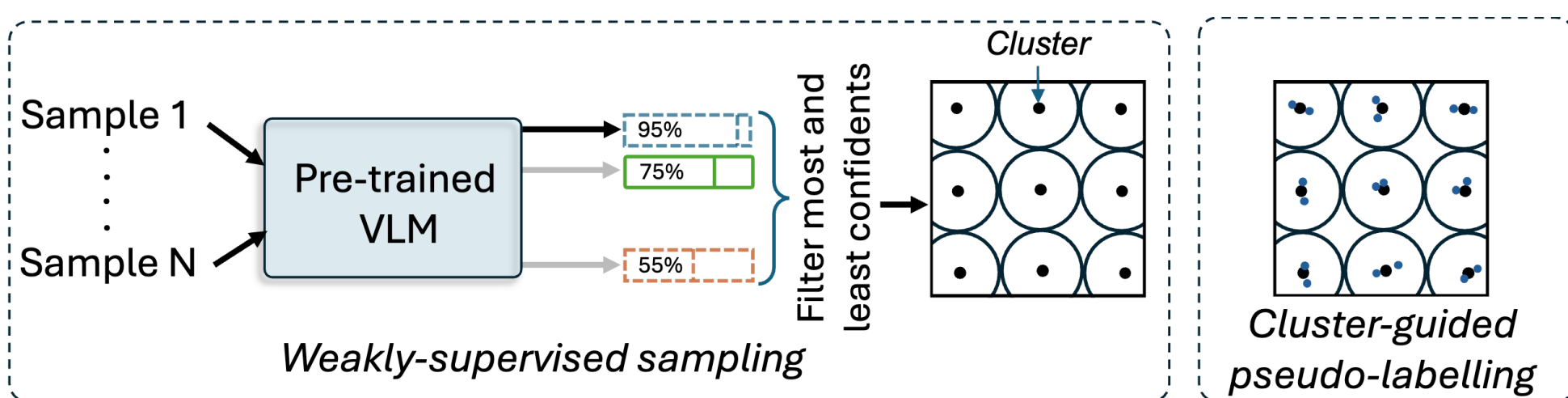


Figure 1: (left) A visual illustration of the weakly-supervised sampling module. Using predictions from the pre-trained VLM, the least and most confident samples, which are not representative of the downstream data, are filtered out. The remaining feature space is then clustered into a number of clusters equal to the labelling budget to ensure maximum diversity among the selected samples. (right) Cluster-guided pseudo-labelling assigns the same class label to samples near the cluster centers as the pseudo-label.

Method

Weakly-supervised sampling

To overcome the limitations of random selection, we introduce a weakly-supervised sampling module that selects the most diverse and representative N samples from the unlabelled set. This module operates through a two-step protocol:

Step 1: Filtering with weak supervision. We leverage the zero-shot predictions of the pre-trained VLM as weak supervision to filter the unlabelled set U . Specifically, we remove samples with both the highest and lowest confidence predictions by the VLM. To this end, we divide the sorted samples into q quantiles, $\{Q_1, Q_2, \dots, Q_q\}$, and select $\mathcal{D}_{\text{filtered}} = \bigcup_{k=2}^{q-1} Q_k$.

Step 2: Diversity Sampling. We select N diverse samples from the filtered dataset with a cluster-based sampling technique. To this end, we apply k -means clustering to group the samples into N clusters and select one sample per cluster closest to the cluster center.

Cluster-guided pseudo-labelling

To improve the pseudo-label quality, especially at the beginning of the training, we propose a novel clustering-guided pseudo-labelling approach that does not utilize the zero-shot prediction from the VLM as the pseudo-label. Specifically, for each cluster \mathcal{C}_j , we pick the p samples closest to the cluster centers to form a pseudo-label set $\mathcal{P}_j = \{x_j^1, x_j^2, \dots, x_j^p\}$.

Confidence-aware semi-supervised learning

To make the best use of the unlabelled data, we propose a confidence-aware semi-supervised module that learns from the high-confident samples in a supervised learning setup, while learning from the low-confident samples in a weakly-supervised setting. We first predict the output distribution for each sample in the unlabelled set $\mathbf{p}_i = f(x_i) \in \mathbb{R}^C$. Then we incorporate the t most confident samples-per-class into our pseudo-label set as:

$$\mathcal{X}^+ = \mathcal{X}_P \cup \left(\bigcup_{c=1}^C \text{top}_t(\{x_i | \arg \max(\mathbf{p}_i) = c\}) \right)$$

Finally, we learn from the labelled set, pseudo-labeled set, and weakly labelled set, together as follow:

$$\mathcal{L}_{\text{final}} = \frac{1}{|\mathcal{X}_L|} \sum_{(x,y) \in \mathcal{X}_L} \ell(f(x), y) + \frac{1}{|\mathcal{X}^+|} \sum_{(x,y) \in \mathcal{X}^+} \ell(f(x), y) + \frac{\lambda}{|\mathcal{X}_{\text{weak}}|} \sum_{(x,s) \in \mathcal{X}_{\text{weak}}} \ell_w(f(x), s)$$

Here, ℓ_w is a partial label learning loss defined as:

$$\ell_w(f(x), \mathbf{s}) = - \sum_{c \in C} \mathbf{s}^c \log(p(c|x))$$

Algorithm 1 SelfPrompt

- 1: **Input:** Unlabelled set U with M samples, label budget N , pre-trained VLM (θ, ϕ) , learnable prompt P , number of sessions S , hyper-parameters t , number of clusters N , pseudo-labels per cluster p
- 2: // Filtering with weak supervision
- 3: **for each** $x_i \in U$ **do**
- 4: Compute the class probability distribution using Eq. 1 as: $\mathbf{p}_i = [p_i^1, p_i^2, \dots, p_i^C]$
- 5: Compute confidence scores, $c_i = \max_{1 \leq c \leq C} (p_i^c)$
- 6: **end for**
- 7: Sort samples in descending order of c_i and divide into q quantiles, $\{Q_1, Q_2, \dots, Q_q\}$.
- 8: Remove samples from first and last quantiles to get $\mathcal{D}_{\text{filtered}} = \bigcup_{k=2}^{q-1} Q_k$
- 9: // Diversity Sampling
- 10: Extract embeddings $\mathbf{z}_i = \theta(x_i)$ for $x_i \in \mathcal{D}_{\text{filtered}}$.
- 11: Perform k -means clustering on $\{\mathbf{z}_i\}$ to form N clusters $\{\mathcal{C}_1, \mathcal{C}_2, \dots, \mathcal{C}_N\}$.
- 12: Select a sample from each cluster j , $x_j^* = \arg \min_{x_i \in \mathcal{C}_j} \|\mathbf{z}_i - \mu_j\|^2$ for $j = 1, \dots, N$, where μ_j is the cluster center.
- 13: Form labelled set $\mathcal{X}_L = \{(x_1^*, y_1), (x_2^*, y_2), \dots, (x_N^*, y_N)\}$.
- 14: // Cluster-guided pseudo-labelling
- 15: **for** $j = 1$ to N **do**
- 16: Select p additional samples per cluster \mathcal{C}_j nearest to the cluster center x_j^* .
- 17: Assign cluster label to selected samples: $\mathcal{P}_j = \{(x_{jk}, y_j)\}_{k=1}^p$
- 18: **end for**
- 19: Create pseudo-label set: $\mathcal{X}_P = \bigcup_j \mathcal{P}_j$
- 20: // Confidence-aware semi-supervised learning
- 21: **for** $s = 1$ to S **do**
- 22: **if** $s == 1$ **then continue**
- 23: **end if**
- 24: Predict probability distribution, $\mathbf{p}_i = f(x_i) \in \mathbb{R}^C$, for $x_i \in U$
- 25: Update \mathcal{X}^+ as, $\mathcal{X}^+ = \mathcal{X}_P \cup (\bigcup_{c=1}^C \text{top}_t(\{x_i | \arg \max(\mathbf{p}_i) = c\}))$
- 26: Form weakly-labelled set: $\mathcal{X}_{\text{weak}} = \{(x_i, s_i) | x_i \in U \setminus \mathcal{X}^+\}$
- 27: Train VLM using loss:

$$\mathcal{L}_{\text{final}} = \frac{1}{|\mathcal{X}_L|} \sum_{(x,y) \in \mathcal{X}_L} \ell(f(x), y) + \frac{1}{|\mathcal{X}^+|} \sum_{(x,y) \in \mathcal{X}^+} \ell(f(x), y) + \frac{\lambda}{|\mathcal{X}_{\text{weak}}|} \sum_{(x,s) \in \mathcal{X}_{\text{weak}}} \ell_w(f(x), s)$$

28: **end for**

Experiments

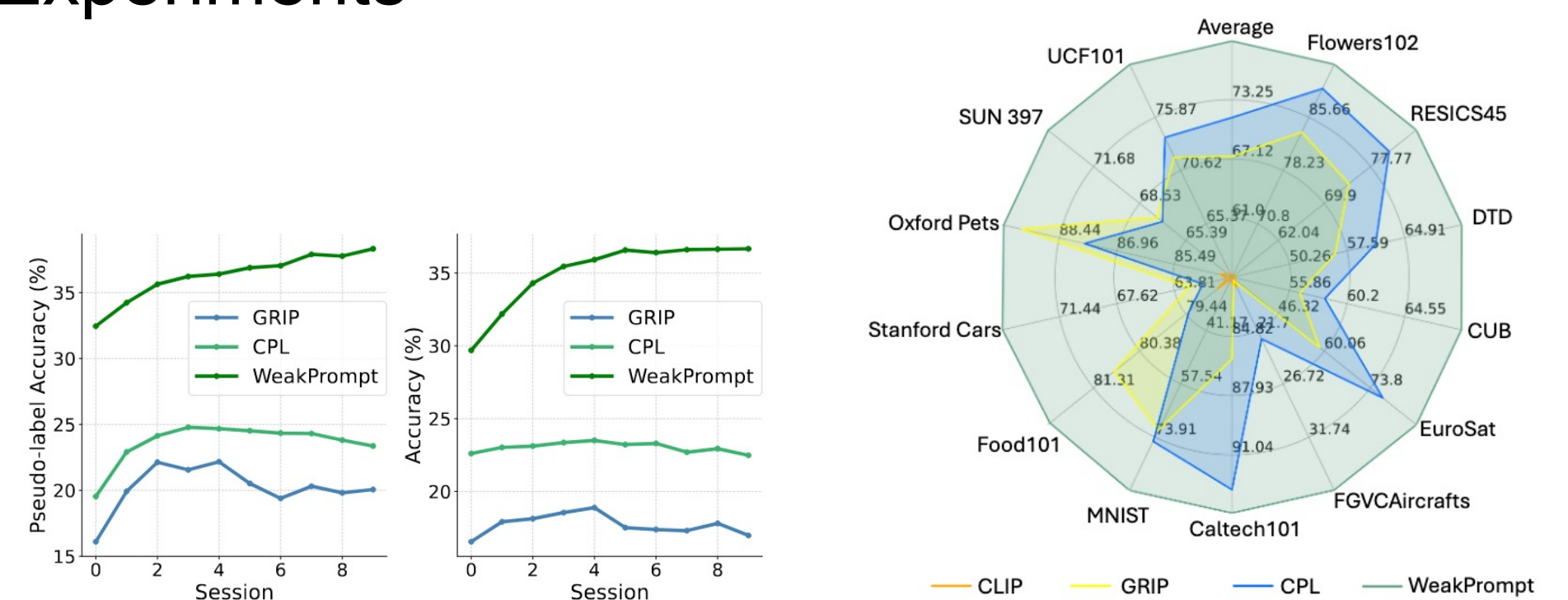


Figure 2: (left) Pseudo-label accuracy; (right) Figure 3: Performance comparison to prior works on semi-supervised tuning of VLMs.

Table 1: Comparison results of top-1 test accuracy (%) on 13 benchmarks on semi-supervised learning with textual prompt strategy.

Methods	Average	Flowers102	RESISC45	DTD	CUB	EuroSAT	FGVCAircraft
Zero-shot CLIP	55.17	63.67 _{0.00}	54.48 _{0.00}	43.24 _{0.00}	51.82 _{0.00}	32.88 _{0.00}	17.58 _{0.00}
CoOp	62.28	75.96 _{0.74}	68.13 _{0.55}	37.10 _{5.45}	55.29 _{0.59}	62.05 _{1.64}	20.02 _{0.77}
GRIP	67.40	83.60 _{0.48}	74.11 _{0.68}	56.07 _{0.79}	56.65 _{0.33}	58.66 _{2.64}	16.98 _{0.20}
CPL	71.41	89.66 _{0.36}	80.98 _{0.11}	61.21 _{0.56}	58.53 _{0.24}	77.51 _{0.80}	22.48 _{0.63}
SelfPrompt	79.33	93.04_{0.33}	85.58_{0.18}	72.18_{0.78}	68.84_{0.16}	87.49_{0.12}	36.71_{0.70}
Δ	$\uparrow 7.92$	$\uparrow 3.38$	$\uparrow 4.60$	$\uparrow 10.97$	$\uparrow 12.31$	$\uparrow 9.98$	$\uparrow 14.23$

	Caltech101	MNIST	Food101	StanfordCars	OxfordPets	SUN397	UCF101
Zero-shot CLIP	82.01 _{0.00}	25.10 _{0.00}	78.81 _{0.00}	60.29 _{0.00}	84.32 _{0.00}	62.54 _{0.00}	60.42 _{0.00}
CoOp	84.69 _{1.43}	58.22 _{1.98}	76.23 _{1.45}	58.23 _{2.45}	82.34 _{1.44}	62.19 _{1.78}	69.19 _{1.03}
GRIP	85.99 _{1.06}	71.78 _{2.59}	80.89 _{1.14}	62.83 _{1.42}	89.40 _{0.33}	67.34 _{0.98}	71.94 _{0.95}
CPL	92.87 _{1.14}	75.18 _{4.40}	79.38 _{1.05}	61.93 _{1.30}	87.79 _{1.31}	66.98 _{0.65}	73.88 _{1.32}
SelfPrompt	94.10_{0.92}	90.23_{0.36}	82.19_{0.17}	75.21_{0.33}	89.86_{0.48}	74.77_{0.18}	81.07_{0.44}
Δ	$\uparrow 1.23$	$\uparrow 15.05$	$\uparrow 2.81$	$\uparrow 13.28$	$\uparrow 2.07$	$\uparrow 7.79$	$\uparrow 7.19$

Table 2: Comparison with existing SOTA on base-to-novel generalization in a 2-shot training setup.

VIT-B/16	ul. Base	Novel	HM
CLIP	✗ 69.3	74.2	71.7
Co-CoOp	✗ 71.9	73.4	72.6
MaPLe	✗ 74.9	73.3	74.0
PromptSRC	✗ 78.1	74.7	76.3
CoPrompt	✗ 74.2	72.4	73.1
PromptKD	✓ 79.7	76.8	78.1
SelfPrompt	✓ 85.6	80.8	83.0
Δ	$\uparrow 5.9$	$\uparrow 4.0$	$\uparrow 4.9$

Table 3: Ablation Study

W.S.S.	C.G.P.	C.A.SSL	Accuracy
✓	✓	✓	79.33
✗	✓	✓	76.12
✓	✗	✓	74.39
✓	✓	✗	78.01
✗	✗	✓	73.49
✓	✓	✗	75.67
✓	✗	✗	73.08
✗	✗	✗	71.41

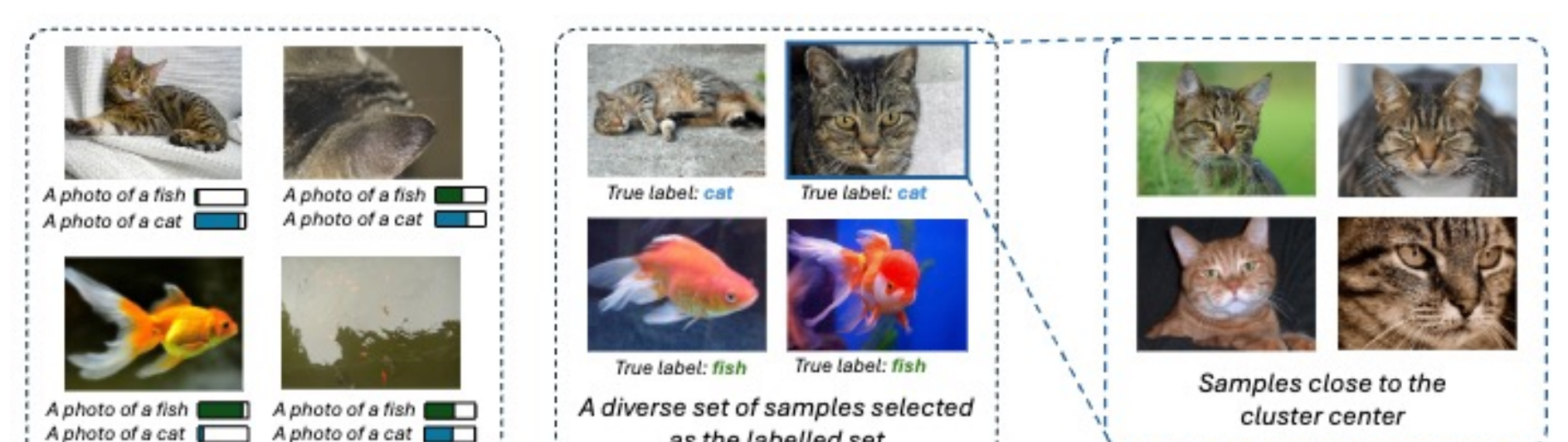


Figure 4: Qualitative analysis of weakly-supervised sampling and cluster-guided pseudo-labelling with two classes (fish and cat). (left) Illustrations of the most confident samples, which provide minimal information gain, alongside the least confident samples, which are less representative of their respective classes. (middle) Examples of selected samples demonstrating high semantic diversity. (right) Samples close to the cluster centers exhibit high visual and semantic similarity.