

# Vision-Language Pseudo-Labels for Single-Positive Multi-Label Learning

Xin Xing<sup>1</sup> Zhexiao Xiong<sup>2</sup> Abby Stylianou<sup>3</sup> Srikumar Sastry<sup>2</sup>  
Liyu Gong<sup>4</sup> Nathan Jacobs<sup>2</sup>

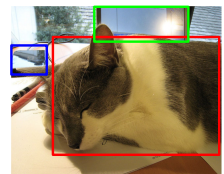
<sup>1</sup> University of Nebraska Omaha <sup>2</sup> Washington University in St. Louis <sup>3</sup> Saint Louis University <sup>4</sup> Oracle Inc

## Abstract

We study a limited label problem and present a novel approach to Single-Positive Multi-label Learning. In the multi-label learning setting, a model learns to predict multiple labels or categories for a single input image. This contrasts with standard multi-class image classification, where the task is to predict a single label from many possible labels for an image. Single-Positive Multi-label Learning specifically considers learning to predict multiple labels when there is only one annotation per image in the training data. Multi-label learning is a more natural task than single-label learning because real-world data often involves instances belonging to multiple categories simultaneously; however, most computer vision datasets contain single labels due to the inherent complexity and cost of collecting multiple high-quality annotations per image. We propose a novel approach called Vision-Language Pseudo-Labeling, which uses a vision-language model, CLIP, to suggest strong positive and negative pseudo-labels. The experiment performance shows the effectiveness of the proposed model. Our code and data will be made publicly available at <https://github.com/mvrl/VLPL>.

## 1. Introduction

Most image classification approaches focus on performing multi-class classification: given an input image, predict which of many possible labels is the most appropriate. In Figure 1, a standard image classification model would likely predict the label ‘cat.’ Most images, however, have more than just one appropriate class. For example, Figure 1 shows a cat, a cell phone, and a laptop – all of which would be appropriate labels for the image. Predicting multiple labels for an input image falls in the domain of multi-label learning. One of the largest challenges for multi-label learning is that most common computer vision datasets only provide a single annotation, even though most images contain multiple objects or classes. In [26], the authors found that the ImageNet dataset contains 1.22 classes per image on average, even though the dataset only includes a single



	cat	dog	cellphone	laptop	boat	person
(a)	✓	✗	✓	✓	✗	✗
(b)	✓	✗	?	✓	?	?
(c)	✓	?	?	?	?	?

Figure 1. This figure shows different levels of available annotation for multi-label learning tasks: (a) full annotation (we know all ground truth positive and negative labels), (b) partial annotation (we know partial ground truth labels, and the rest labels are unknown), and (c) single positive annotation (we only know one positive ground truth label, and the rest labels are all unknown).

label per image. Collecting all possible labels for an image is time consuming, costly, and error-prone, especially when an image has a large number of classes and some of them may only be visible in a very small part of the image. This problem domain – where the training data contains only a single label, but the task is predicting multiple labels – is called Single-Positive Multi-Label Learning (SPML).

There are a variety of different works that focus on the SPML task [2, 8, 24]. These works mainly concentrate on pseudo-labeling approaches and novel loss definitions that use these labels. Pseudo-labeling uses different types of weak supervision to identify potential positive labels for an image. These labels may come from pre-trained multi-class classification backbones, or from label-to-label association, which focuses on leveraging known or inferred relationships and dependencies between different labels. Novel losses explore how to utilize these pseudo-labels in multi-label training.

In this paper, we introduce a novel approach called Vision-Language Pseudo-Label (VLPL) for SPML. Prior pseudo-labeling work has largely focused on setting a score threshold for extracted features [4], or incorporating uncertainty in the pre-trained features in the pseudo-labeling [21]. More recently, researchers have considered using Vision-Language Models like CLIP in the pseudo-labeling process. In DualCoOp([22]), the authors use a fixed CLIP model and learn positive and negative prompt contexts per

image as pseudo-labels, which are then fed into an asymmetric loss [20] for limited-annotation (but not single positive) multi-label classification. In [3], the authors focus on incorporating label-to-label correspondence priors using a structured prior derived from a CLIP model and a Semantic Correspondence Prompt network that keys on label-to-label correspondences.

Our VLPL approach is most similar to DualCoOp. We, however, show that only using positive pseudo-labels extracted based on CLIP image-text similarity, and using an Entropy Maximization loss, can achieve SOTA performance in the SPML setting, on Pascal VOC, MS-COCO, and CUB-Birds with a significantly simpler approach.

We demonstrate the superior performance of the VLPL model compared to baseline models, evaluate the influence of hyperparameters on the VLPL model’s performance, and explore the proposed model’s performance under different scenarios. The contributions of our study are as follows:

- Inspired by the VLM application, we proposed a novel model called vision-language pseudo-label (VLPL) for SPML, which aims to produce an accurate and robust pseudo-label to boost the model performance.
- We conducted experiments on four benchmark multi-label datasets, i.e., PASCAL VOC [5], MS-COCO [13], NUS-WIDE [1], and CUB [25]. Our initial experiments achieve superior performance over the baseline models, proving the effectiveness of the proposed method. Besides, by further exploration of the backbone, our method achieves new state-of-the-art results, pushing the performance boundary improvement to  $mAP = 93.37$ ,  $mAP = 84.65$ ,  $mAP = 57.12$ , and  $mAP = 26.04$  over the four benchmarks, respectively.
- To further investigate the impact on the VLPL model’s performance, we conduct more experiments to systematically evaluate our model. We examine how varying hyperparameters affected the effectiveness of the model, discuss the positive-negative imbalance in our study, and visualize the final prediction probabilities. For more details, please refer to our experiment section.

## 2. Related Work

### 2.1. Loss-function Focused Methods

For SPML tasks, much of the work focuses on developing novel loss functions to train models. Assume Negative (AN) Loss [2] is a simple method that assumes all the unknown labels are negative, inevitably introducing some number of false negatives in the implementation. Though the performance of AN is unsatisfactory, AN is still a widely used baseline for comparison. Entropy-Maximization (EM) loss [27] leverages the idea of acknowledging unknown labels and aims to maximize the entropy of predicted probabilities for unannotated labels. The Weak Assume Nega-

tive (WAN) loss [2, 15] is an updated AN method, wherein the negative labels are weighted by a ‘weak’ coefficient to reduce the impact of false negatives. Regularized Online Label Estimation (ROLE) [2] mirrors the expectation-maximization algorithm in jointly training the image classifier and concurrently estimating potential labels online. Large Loss (LL) [10] proposes to overcome the memorization effect, which the model first learns the representation of clean labels, and then starts memorizing noisy labels. Our study does not primarily focus on the loss function – we use the EM loss, which is suitable for our model as our label prediction also includes a number of unknown labels.

### 2.2. Pseudo Label Focused Methods

Pseudo-labeling is another popular method to overcome the problem of limited annotations per image. Given the imbalance between positive and negative labels (where the number of negative labels significantly outnumbers positive labels), an intuitive approach is to sample a portion of the unknown labels and assign them as negative “pseudo-labels” [12]. Clustering methods like [24] leverage a distance metric to facilitate weakly-supervised or self-supervised learning for pseudo-labeling. Asymmetric Pseudo-Labeling (APL) [27] assigns positive and negative pseudo-labels with asymmetric tolerance. This approach is often employed in conjunction with the previously mentioned EM loss. Unlike the prior pseudo-labeling methods that lack robust reference, our strategy involves leveraging an aligned vision-language embedding space to predict positive and negative labels. Our method demonstrates robustness and introduces fewer noisy labels during implementation, thereby enhancing the quality and effectiveness of the labeling process.

### 2.3. Vision-Language Model

Vision-Language Model (VLM) [9, 17] is a multi-modality model, using the image and its corresponding text as supervision signal to help us better understand vision-language correlation. The most commonly used VLM model is named Contrastive Language-Image Pre-Training (CLIP) [17]. Since CLIP is a well-trained vision-language model, it has impressive potential as the tool for different downstream tasks including multi-label learning. In DualCoOp [22], the authors leverage CLIP for multi-label recognition tasks in limited-annotation domains (although not single-label). They highlight the benefit of learning the relationship between different category names in the multi-label recognition task and observe that the aligned image and textual CLIP spaces can be used for this purpose. Concretely, they learn a positive and negative “prompt context” – a sequence of embedding vectors – for possible target category names. These prompt contexts can then be used as classifiers by computing the similarity between local features in

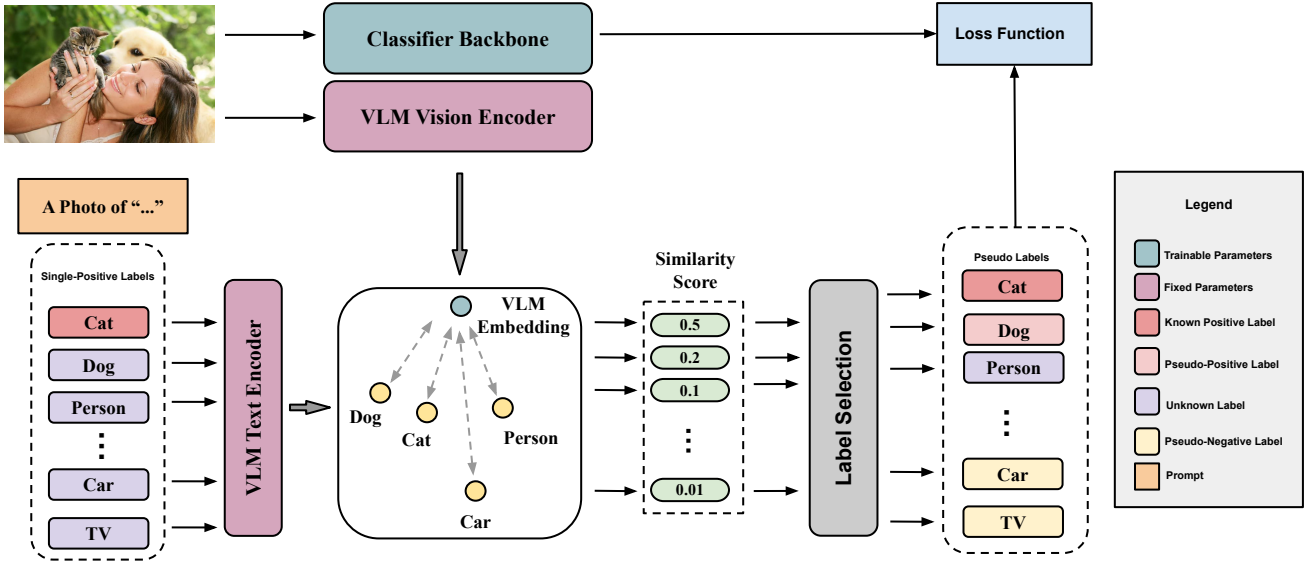


Figure 2. The architecture of the proposed model. The image encoder is a trainable model for feature extraction, such as ResNet and Vision Transformer, meanwhile, the VLM Vision and Text encoders are fixed for vision-language embedding space. By the vision-language embedding operation, instead of the original single-positive label offering, we can access a brand new pseudo-label as the reference for the final prediction.

an image and each of the context vectors, and assigning a positive or negative label for each category (at each location) based on whichever context has the higher similarity. SCPNet [3] proposes to explore structured semantic prior information to better understand label-to-label associations in images. The authors use a CLIP model to extract an object association matrix as the prior information to achieve better performance. MKT [7] is proposed for the zero-shot multi-label learning by applying the knowledge transfer model with CLIP model’s initial weights. In our study, we propose a simpler approach to leveraging CLIP features, called Vision-Language Pseudo-Label (VLPL), that works in the single-label annotation domain and only requires the selection of positive pseudo-labels.

### 3. Approach

In this section, we describe the problem definition, detail the architecture of the proposed model, illustrate the vision language pseudo-labeling method, and describe our loss function.

#### 3.1. Problem Definition

In the context of multi-label learning (MLL), we are given a dataset  $D = \{X_i, Y_i\}_{i=1}^N$  consisting of  $N$  training samples. Each sample  $X_i$  is an input image, and its corresponding label vector  $Y_i \in \{-1, 1\}^L$  has a length  $L$ . Within this label vector,  $Y_{il} = 1$  designates a positive label that is relevant to  $X_i$ , whereas  $Y_{il} = -1$  signifies a negative label that

is irrelevant to  $X_i$ . Our study focuses on Single-Positive Multi-Label learning (SPML), where there exists only one positive label and all others are unknown. To indicate this, we modify the label vector annotation as  $Y_i \in \{-1, \emptyset, 1\}^L$ . The symbol  $\emptyset$  in  $Y_{il}$  indicates that the association of the  $l$ -th label with the input image  $X_i$  remains undetermined.

We have the annotation  $\sum_{l=1}^L 1_{y_{il}} = 1$ , where  $1_{[\cdot]}$  represents an indicator function. This implies that each input image has only one observed positive label, and the rest remain unknown. The main objective of the SPML study is to learn a mapping function  $f : X \rightarrow Y$  from the dataset  $D$ . The ground truth label is  $Y = \{-1, 1\}$ , while the observed label is limited to  $Y' = \{\emptyset, 1\}$ , making SPML a challenging task within the realm of MLL, operating under limited supervision.

#### 3.2. Architecture

In Figure 2, we present the architecture of our proposed model. This model is designed with two branches: the upper branch is comprised of a trainable feature extraction component, with the flexibility to use any image encoder as its backbone. For further details about the specific image encoder utilized in our experiments, please refer to Section 4. The lower branch includes the fixed VLM, which we use for pseudo-label prediction. We initialize a Vision Encoder and Text Encoder with pre-trained CLIP weights, keeping these weights fixed throughout the process. These encoders produce 768-dimensional output em-

beddings. To generate embeddings for each possible image label, we compute the CLIP text embedding based on the prompt “A photo of  $X$ ”. As the text encoder remains static, these embeddings only have to be computed once. During the model’s operation, the input image is fed into both the trainable image encoder and the fixed CLIP vision encoder, resulting in visual embeddings. This CLIP visual embedding can be compared with the text embeddings generated from text prompts in the form of “A Photo of ...”. To determine pseudo-label assignments, we compute the cosine similarity between the visual and label embeddings from the text prompts. Based on this similarity measure, we select pseudo-labels for each image. Subsequently, we utilize these pseudo-labels, in conjunction with the single positive label, as the final reference for training the entire model, enabling it to make accurate predictions in tasks.

### 3.3. Vision Language Pseudo-Labeling

---

#### Algorithm 1 Vision Language Pseudo-Labeling

---

**Input:** The visual embedding  $e_I$  and the query label embedding  $e_L^i$

**Parameter:** positive pseudo-label threshold  $\theta$ , negative pseudo-label partial  $\delta\%$

**Output:** Pseudo-Labeling  $Y$  of input image

- 1: Compute the image and the query labels similarity as equation 1
  - 2: Let  $i = 0$ .
  - 3: **while**  $i < L - 1$  **do**
  - 4:   **if**  $p_i > \theta_p$  **then**
  - 5:      $Y_i = 1(\text{Positive} - \text{label})$
  - 6:   **else if**  $p_i < \theta_n$  **then**
  - 7:      $Y_i = -1(\text{Negative} - \text{label})$
  - 8:   **else**
  - 9:      $Y_i = \emptyset(\text{Unknow} - \text{label})$
  - 10:   **end if**
  - 11:    $i = i + 1$
  - 12: **end while**
  - 13: **return**  $Y$
- 

Our pseudo-labeling method benefits from the “open-world” capabilities of the large VLM, enabling the use of rich, free-form text with a long list vocabulary of visual categories. In zero-shot learning applications, these VLMs allow us to obtain visual embeddings and potential label embeddings with ease. The cosine similarity between these embeddings can then be used to perform image classification. We adopt this methodology for multi-label learning in our work. As shown in Figure 2, the VLPL employs a visual encoder  $E_V : R^{w \times h \times 3} \rightarrow R^d$  and a text encoder  $E_L : R^{m \times d_c} \rightarrow R^d$  to extract the image and text embeddings, respectively. The visual inputs are 3-channel

images of shape  $w \times h$ , while the text inputs are prompts consisting of  $m$  words, each of which is embedded into a  $d_c$ -dimensional vector. Both the image and text inputs are mapped into a  $d$ -dimensional latent space. We obtain a visual embedding vector  $e_I \in R^d$  and  $n$  label embedding vector  $e_L^1, e_L^2, \dots, e_L^n \in R^d$ , where  $n$  denotes the label number of the category space. We compute the dot product between  $e_I$  and each of the  $e_L^i$ , resulting in an  $n$ -dimensional vector, where the  $i$ -th element means the similarity between the image and the  $i$ -th label query. This similarity vector can be used to predict the label of the image, and the probability  $p_i$  of the appearance of the  $i$ -th label on the image is computed by the temperature softmax function

$$p_i = \frac{\exp(\langle e_I, e_L^i \rangle / \tau)}{\sum_{j=1}^L \exp(\langle e_I, e_L^j \rangle / \tau)} \quad (1)$$

where  $\langle \cdot, \cdot \rangle$  denotes the dot product and  $\tau$  is a temperature scalar of the softmax function.

As shown in Algorithm 1, we use the equation 1 to measure the similarity between the image and query labels. We determine the pseudo-label in three formations: positive, negative, and unknown using one threshold, namely, a positive threshold  $\theta$  and a negative label percentage coefficient  $\delta\%$ . For the input image  $I$  and the  $i$ -th query label, if the measurement similarity  $p_i$  is over than  $\theta$ , we set  $Y_i$  as positive. In terms of the large label space of multi-label learning, a given image has a few positive labels, and the rest are negative. We can rank the similarity vector  $P = [p_1, \dots, p_L]$  and set  $\delta\%$  of the smallest similarity values as negative labels. Afterward, the rest are set unknown. By following this simple algorithm, we can generate a new pseudo-label vector given the input image.

### 3.4. Loss Function

In the SPML domain, there is only one ground truth positive label, and the rest are unknown. How these labels are utilized by the loss function plays a crucial role in the model training. The Assuming-Negative (AN) Loss, where the unknown labels are assumed to be negative, is commonly used as the baseline loss for SPML tasks, but leads to the generation of many false negatives during the model training. We instead use the more recently proposed Entropy-Maximization (EM) loss, which acknowledges the un-annotated labels as unknown, rather than negative, and seeks to maximize the entropy of predicted probabilities for the unknown labels.

$$\begin{aligned} Loss_{EM}(x^{(n)}, y^{(n)}) = & -\frac{1}{L} \sum_{l=1}^L [1_{[y_l^n=1]} \log(f_l(x^{(n)})) \\ & + 1_{[y_l^n=\emptyset]} \alpha H(f_l(x^{(n)}))] \end{aligned} \quad (2)$$



$$H(f_l(x^{(n)})) = -[f_l(x^{(n)})\log(f_l(x^{(n)})) + (1 - f_l(x^{(n)}))\log(1 - f_l(x^{(n)}))] \quad (3)$$

where  $H(f(x^{(n)}))$  is the entropy loss of the unknown labels.

In our model, the VLPL will generate pseudo-labels for positive, negative, and unknown categories. We use the EM loss strategy to acknowledge the unknown labels and integrate the pseudo-label loss of our model. Our loss is the following:

$$\begin{aligned} Loss(x^{(n)}, y^{(n)}) = & -\frac{1}{L}\sum_{l=1}^L [1_{[y_l^n=1]}\log(f_l(x^{(n)})) \\ & + 1_{[y_l^n=\emptyset]}\alpha H(f_l(x^{(n)})) \quad (4) \\ & + 1_{[y_l^n=\hat{1}]} \beta S(f_l(x^{(n)})) + 1_{[y_l^n=-\hat{1}]} \gamma S(f_l(x^{(n)}))] \end{aligned}$$

$$S(f(x^{(n)})) = (1 - \rho)\log(1 - f_l(x^{(n)})) - \rho\log(f_l(x^{(n)})) \quad (5)$$

where,  $1_{[y_l^n=\hat{1}]}$  denotes the pseudo positive-label,  $1_{[y_l^n=-\hat{1}]}$  denotes the pseudo negative-label,  $S(f(x^{(n)}))$  is the pseudo-label loss with labeling smooth  $\rho$ ,  $\alpha$ ,  $\beta$ ,  $\gamma$  are the coefficients of each loss section.

However, in experimentation (discussed in Section 4.5), we found that the best model performance is achieved when we only use pseudo-positive labels, while keep the rest labels unknown, rather than including pseudo-negative labels. Therefore our final loss function is:

$$\begin{aligned} Loss(x^{(n)}, y^{(n)}) = & -\frac{1}{L}\sum_{l=1}^L [1_{[y_l^n=1]}\log(f_l(x^{(n)})) \\ & + 1_{[y_l^n=\emptyset]}\alpha H(f_l(x^{(n)})) + 1_{[y_l^n=\hat{1}]} \beta S(f_l(x^{(n)})) \quad (6) \end{aligned}$$

## 4. Evaluation

### 4.1. Implementation Details

Our models are implemented using PyTorch. We train the model for 10 epochs, using the Adam [11] optimizer. The batch size is 8. The learning rate is determined by grid search in the range of the  $\{1e - 3, 1e - 4, 1e - 5\}$ , and we find  $1e - 5$  yields the best performance. For data augmentation, we use horizontal flipping for the training dataset with 50% probability. Images are resized to  $448 \times 448$  for both our proposed model and all baseline models. We follow the conventions of previous works in multi-label classification [2, 27] for model evaluation and report the mean average precision (mAP).

### 4.2. Dataset

Since there are currently no datasets explicitly designed for Single Positive Multi-label Learning (SPML), we use [2]

and others’ adaptation of existing large-scale multi-label datasets to simulate a “single positive” scenario. This method allows us to retain all ground-truth labels for performance evaluation and training phenomena analysis. After setting aside 20% of the training images for validation, one random positive label is kept for each training image, treating all other labels as un-annotated. This operation is performed once for each dataset. It’s important to note that the validation and test sets remain fully labeled. We use four well-known datasets in our experiments, namely: PASCAL VOC (VOC) [5], MS-COCO (COCO) [13], NUS-WIDE(NUS) [1], and CUB-200-2011(CUB) [25].

**PASCAL Visual Object Classes Challenge (VOC2007)** [5] is a widely used dataset for multi-label recognition. It contains 5,011 images in the training/validation set, and 4,952 images as the test set. There are 20 possible classes, with an average of 2.5 categories per image.

**Microsoft COCO** [13] (MS-COCO) is another widely used benchmark for multi-label image recognition. It contains 82,801 training images and 40,504 validation images. There are 80 categorized objects in this dataset, with an average of 2.9 object labels per image. Since this data set lacks of test set, the validation images are often used for evaluation in the literature.

**NUS-WIDE** [1] is a real-world web image dataset. Originally, the dataset contains 269,648 Flickr images with 81 manually annotated visual concepts. However, due to some Flickr image downloading links expiring, it is impractical to evaluate our model using the original dataset. We use the curated dataset provided by the authors of AckUnknown [27], which has 150,000 training images and 60,260 test images, to conduct our experiment and provide a fair comparison to recent work.

**CUB** [25] is a dataset for fine-grained visual categorization task. It contains 11,788 images of 200 subcategories belonging to birds, with 5,994 training images and 5,794 testing images. In the experiment, the model will predict 312 attributes of each bird images.

### 4.3. Baseline Performance

Most of the current state-of-the-art models use a ResNet50 backbone. Because of this, we first explore the performance of our proposed model using a ResNet50. Our experimental setting – using a single-positive label adaptation of the datasets and reporting the mAP evaluation metric, evaluated on the model that achieves the highest accuracy on a withheld validation set – is the same as the previous methods [2, 10, 27]. Table 1 reports the performance results of the different models. The proposed model, VLPL, achieves  $mAP = 89.10$  on the VOC dataset,  $mAP = 71.45$  on the COCO dataset,  $mAP = 49.55$  on the NUS-WIDE dataset, and  $mAP = 24.02$  on the CUB dataset. Compared with

Ann. Labels	Methods	VOC	COCO	NUS	CUB
<i>Full Annotation</i>					
All P. & All N	BCE Loss	89.42	76.78	52.08	30.90
<i>Limited Annotation</i>					
1 P. & All N	BCE Loss	87.60	71.39	46.45	20.65
1 P. & 0 N.	AN Loss	85.89	64.92	42.27	18.31
	DW [2]	86.98	67.59	45.71	19.15
	L1R [2]	85.97	64.44	42.15	17.59
	L2R [2]	85.96	64.41	42.72	17.71
	LS [16]	87.90	67.15	43.77	16.26
	N-LS [2]	88.12	67.15	43.86	16.82
	EntMin [6]	53.16	32.52	19.38	13.08
	Focal Loss [14]	87.59	69.79	47.00	19.80
	ASL [20]	87.76	68.78	46.93	18.81
	ROLE [2]	87.77	67.04	41.63	13.66
	ROLE+LI [2]	88.26	69.12	45.98	14.86
	EM [27]	89.09	70.70	47.15	20.85
	EM+APL [27]	89.19	70.87	47.59	21.84
	LL-R [10]	<b>89.2</b>	71.0	47.4	19.5
	LL-Ct [10]	89.0	70.5	48.0	20.4
LL-Cp [10]	88.4	70.7	48.3	20.1	
DualCoOp [22]	83.6	69.2	42.8	–	
1 P. & 0 N.	VLPL(Ours)	89.10	<b>71.45</b>	<b>49.55</b>	<b>24.02</b>

Table 1. Results of the different models with the same experimental settings as the [2]. Using the same input image size setting  $448 \times 448$ , our model outperforms the limited-annotation baseline models for the benchmark COCO, NUS-WIDE, and CUB.

the SOTA baselines under the same experimental setting, VLPL demonstrated superior performance across all the benchmarks, indicating the effectiveness of our proposed model. While we focus on the limited annotation setting, we also compare our performance to a model trained using full ground truth annotations. Our performance using limited annotations is competitive with the model using the full set of ground truth annotations.

#### 4.4. Backbone Experimentation

In the previous sections, we prove the effectiveness of our proposed model using a ResNet50 model and explore the model performance with different conditions. The current state-of-the-art methods for SPML concentrate more on the methodology but ignore the power of the classifier backbone architecture. It’s widely known, however, that for computer vision tasks, an optimal network architecture with suitable pretraining initialization can substantially improve the model’s performance [18]. To this end, we explored various network architectures and pretraining initializations. In our implementation, we select ConvNeXt-XL and ViT-L as the backbone. For the pretraining initialization, we choose the ConvNeXt-XL model pretrained under ImageNet1k and ImageNet22k, and the ViT-L model initialized with CLIP

weights. Table 2 presents the different pretrained backbone performances. We observe that all the chosen large network architectures outperform the ResNet50 backbone, indicating that model architecture scaling can be beneficial for performance improvement. While the model scale has a significant impact on performance, the performance gap between different pretraining initializations is not significant.

While the ResNet50 baseline trained using VLPL already outperformed other baseline methods, the best larger models achieve even more impressive performance improvement. Specifically, compared with ResNet50 backbone, we achieved a 5.7% increase in mAP (up to 94.16) on the VOC dataset, a 18.5% increase (up to 84.65) on the COCO dataset, a 15.3% increase (up to 57.12) on the NUS-WIDE dataset, and an 8.4% increase (up to 26.04) on the CUB dataset. Considering the model performance across all four benchmarks, the ConvNeXt-XL model pretrained with ImageNet-22k achieves the best performance on average.

#### 4.5. Positive Labels vs. Negative Labels

In this section, we explore the influence of pseudo-positive and pseudo-negative labels on the performance of the VLPL model. Multi-label learning datasets are typically charac-

Backbone	Pretrained	VOC	COCO	NUS	CUB
ResNet-50	ImageNet1k	89.10	71.45	49.55	24.02
ConvNeXt-XL	ImageNet1k	93.31	83.37	56.11	25.49
ConvNeXt-XL	ImageNet22k	93.37	<b>84.65</b>	<b>57.12</b>	<b>26.04</b>
ViT-L	CLIP	<b>94.16</b>	80.55	52.53	12.82

Table 2. Results of the different network architectures and pretraining initializations for the proposed model. Compared with ResNet50 backbone, all of the chosen large network architectures show better performance. Further, it shows the ConvNeXt-XL with ImageNet22k pretraining initialization weights achieves comparable performance across the different datasets.

$\delta\%$	10%	20%	30%	40%
Pos. + Neg.	88.51	88.43	88.33	88.13
Pos. only	89.10			

Table 3. The results on different pseudo-negative label percentages over the total label vector (from 10% to 40%). We conducted the experiments on benchmark PASCAL VOC.

terized by an imbalance between positive and negative labels, as noted in previous studies [19]. For instance, the Pascal VOC dataset has 20 potential labels, yet the average image contains only 2.5 positive labels. This imbalance creates a challenging environment for model learning. In the context of our method, for a predicted label vector, we can predict pseudo-negative labels by ranking the similarity scores between visual and label embeddings. We then select the labels with the least  $\delta\%$  similarity scores as pseudo-negative labels. The rationale behind this is that the labels with the least similarity scores are likely to be the ones that are most irrelevant or ‘negative’ to the given image. To understand the impact of the number of pseudo-negative labels on model performance, we conducted experiments using varying percentages of pseudo-negative labels, ranging from 10% to 40% of the total label count on the benchmark Pascal VOC dataset.

Table 3 shows the results of model performance on different pseudo-negative labels. We observe a decrease in model performance as the number of pseudo-negative labels increases. Our hypothesis for the cause of this is that the imbalance between positive and negative labels negatively impacts model performance. By comparing models that use both pseudo-positive and pseudo-negative labels against models that only utilize pseudo-positive labels it appears that the latter approach provides better performance. In light of these findings, our practical implementation uses pseudo-positive labeling exclusively, while treating the remaining labels as unknown.

#### 4.6. Label Smoothing

Label smoothing (LS) is a method to overcome overfitting and mitigate the label noise for multi-class classifiers [16, 23]. In our implementation, we adopt LS for

	VOC	COCO	NUS	CUB
w/o LS	88.69	70.96	49.13	23.71
w LS	89.10	71.34	49.55	24.02

Table 4. The results on model performance with and without label smooth (LS) on the pseudo-labeling. We conducted the experiments on four benchmarks. The results indicate that LS contributes more to the model performance.

our pseudo-labeling and set the  $\epsilon = 0.9$ . Therefore, the  $i$ -th pseudo positive-label category entropy loss function is  $loss_i = -[\epsilon(\log(f(x))) + (1-\epsilon)(\log(1-f(x)))]$ . As shown in Table 4, we conduct the experiments with and without LS over four benchmarks. The results indicate applying LS improves the model performance compared to the model without it. We set the label smoothing as the default setting in our whole experiment.

#### 4.7. Temperature Hyperparameter & Threshold

The temperature scalar  $\tau$  of equation 1 and pseudo-labeling threshold  $\theta$  are crucial hyperparameters, as they directly influence the pseudo-labeling processing. We conducted a hyperparameter search experiment on a range of values for these two hyperparameters. Considering the different benchmark datasets under different label query numbers, especially for the CUB dataset. We set  $\tau = [0.01, 0.03, 0.05, 0.07, 0.09]$  and  $\theta = [0.1, 0.2, 0.3]$  for VOC, COCO, and NUS-WIDE datasets. Meanwhile, we set  $\tau = [0.01, 0.03, 0.05, 0.07, 0.09]$  and  $\theta = [0.01, 0.03, 0.05]$  for CUB to identify the optimal settings. Our results, shown in Figure 3, indicate that the model performance is sensitive to the  $\tau$  and  $\theta$ . To visualize the performance of the different hyperparameters, we set each threshold  $\theta$  as the condition and plot the different temperature value  $\tau$  over the same  $\theta$ . In terms of the performance over different benchmarks,  $\tau = 0.03$  and  $\theta = 0.3$  yield superior performance for VOC ( $mAP = 89.10$ ) dataset,  $\tau = 0.01$  and  $\theta = 0.3$  yield superior performance for COCO ( $mAP = 71.45$ ) dataset,  $\tau = 0.03$  and  $\theta = 0.1$  yield superior performance for NUS-WIDE ( $mAP = 49.55$ ), and  $\tau = 0.03$  and  $\theta = 0.01$  yield superior performance for CUB ( $mAP = 24.02$ ).

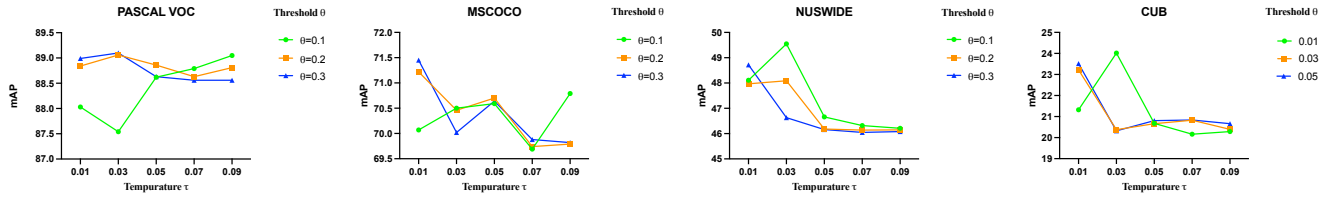


Figure 3. The hyperparameter search of the temperature scalar  $\tau$ , and pseudo-labeling threshold  $\theta$  across the four benchmark datasets: PASCAL VOC, COCO, NUS-WIDE, and CUB.

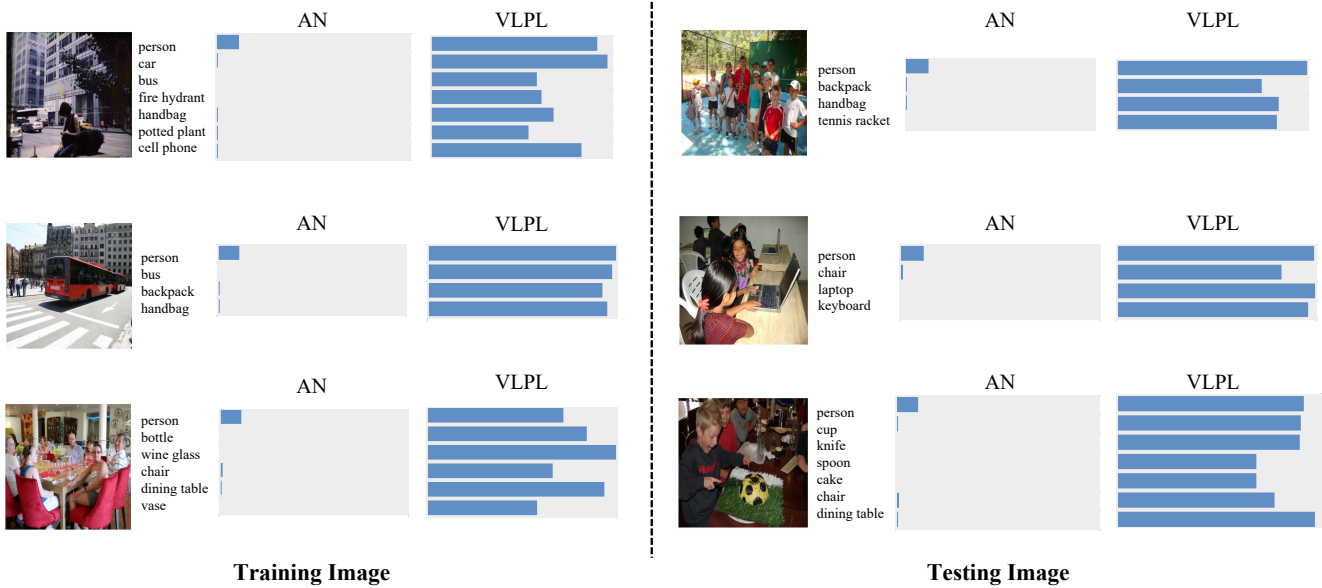


Figure 4. The visualization results on the training and testing images of COCO dataset. The blue bar is the prediction probability of each positive label. Compared with the baseline method AN, our VLPL shows superior performance for the final label prediction.

#### 4.8. Results Visualization

To demonstrate the efficacy of our approach, we provide visualizations in Figure 4, highlighting the label prediction probability generated by our method compared with those from the baseline model (referred to as AN). Our proposed model demonstrates a high level of confidence in accurately identifying positive labels, compared with the predictions offered by the baseline model. For the purposes of our visualization, we have intentionally selected instances from both the training and testing datasets. This choice allows us to showcase how our VLPL method consistently excels in performance across different stages of model training and testing. The visual evidence serves as a testament to the advantages of employing VLPL, both during the training phase and in the application to unseen data.

#### 5. Discussion

We introduced VLPL, an innovative yet simple approach to single-positive multi-label learning. VLPL leverages a

large-scale vision-language model and utilizes the aligned visual and textual embedding similarities to generate pseudo-labels. Our method consists of simple components and results in significant performance improvements across several popular datasets when compared to existing, more complex approaches. We carried out a comprehensive set of experiments and ablations to better understand the impact of various factors within the VLPL framework and explore how to maximize accuracy.

SPML is an extreme challenge of the weakly-supervised multi-label classification task. Meanwhile, pseudo-labeling is one of the most effective methods for SPML task. Previous pseudo-labeling methods concentrate on the single modality (visual embedding) and lacking a stable label reference. However, CLIP, a model jointly mapping visual and label embeddings into the same space, makes it possible to refer the different visual embeddings to the corresponding labels. What's more, as a foundation model, CLIP is well-pretrained with a large-scale dataset, making it capable of extracting the discriminative features.



## References

- [1] Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, pages 1–9, 2009. 2, 5
- [2] Elijah Cole, Oisín Mac Aodha, Titouan Lorieul, Pietro Perona, Dan Morris, and Nebojsa Jojic. Multi-label learning from single positive labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 933–942, 2021. 1, 2, 5, 6
- [3] Zixuan Ding, Ao Wang, Hui Chen, Qiang Zhang, Pengzhang Liu, Yongjun Bao, Weipeng Yan, and Jungong Han. Exploring structured semantic prior for multi label recognition with incomplete labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3398–3407, 2023. 2, 3
- [4] Thibaut Durand, Nazanin Mehrasa, and Greg Mori. Learning a deep convnet for multi-label classification with partial labels. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 647–657, 2019. 1
- [5] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. <http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html>. 2, 5
- [6] Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. *Advances in neural information processing systems*, 17, 2004. 6
- [7] Sunan He, Taian Guo, Tao Dai, Ruizhi Qiao, Xiujun Shu, Bo Ren, and Shu-Tao Xia. Open-vocabulary multi-label classification via multi-modal knowledge transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 808–816, 2023. 3
- [8] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5070–5079, 2019. 1
- [9] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*, pages 4904–4916. PMLR, 2021. 2
- [10] Youngwook Kim, Jae Myung Kim, Zeynep Akata, and Jungwoo Lee. Large loss matters in weakly supervised multi-label classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14156–14165, 2022. 2, 5, 6
- [11] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [12] Dong-Hyun Lee. The simple and efficient semi-supervised learning method for deep neural networks. In *ICML 2013 Workshop: Challenges in Representation Learning*, pages 1–6, 2013. 2
- [13] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 740–755. Springer, 2014. 2, 5
- [14] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 6
- [15] Oisín Mac Aodha, Elijah Cole, and Pietro Perona. Presence-only geographical priors for fine-grained image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9596–9606, 2019. 2
- [16] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? *Advances in neural information processing systems*, 32, 2019. 6, 7
- [17] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [18] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. *arXiv preprint arXiv:2104.10972*, 2021. 6
- [19] Tal Ridnik, Emanuel Ben-Baruch, Nadav Zamir, Asaf Noy, Itamar Friedman, Matan Protter, and Lihi Zelnik-Manor. Asymmetric loss for multi-label classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 82–91, 2021. 7
- [20] Tal Ridnik, Emanuel Ben-Baruch, Nadav Zamir, Asaf Noy, Itamar Friedman, Matan Protter, and Lihi Zelnik-Manor. Asymmetric loss for multi-label classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 82–91, 2021. 2, 6
- [21] Mamshad Nayeem Rizve, Kevin Duarte, Yogesh S Rawat, and Mubarak Shah. In defense of pseudo-labeling: An uncertainty-aware pseudo-label selection framework for semi-supervised learning. *arXiv preprint arXiv:2101.06329*, 2021. 1
- [22] Ximeng Sun, Ping Hu, and Kate Saenko. Dualcoop: Fast adaptation to multi-label recognition with limited annotations. *Advances in Neural Information Processing Systems*, 35:30569–30582, 2022. 1, 2, 6
- [23] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. 7
- [24] Fariborz Taherkhani, Ali Dabouei, Sobhan Soleymani, Jeremy Dawson, and Nasser M Nasrabadi. Self-supervised wasserstein pseudo-labeling for semi-supervised image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12267–12277, 2021. 1, 2

- [25] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset, 2011. [2](#), [5](#)
- [26] Sangdoon Yun, Seong Joon Oh, Byeongho Heo, Dongyoon Han, Junsuk Choe, and Sanghyuk Chun. Re-labeling imagenet: from single to multi-labels, from global to localized labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2340–2350, 2021. [1](#)
- [27] Donghao Zhou, Pengfei Chen, Qiong Wang, Guangyong Chen, and Pheng-Ann Heng. Acknowledging the unknown for multi-label learning with single positive labels. In *European Conference on Computer Vision*, pages 423–440. Springer, 2022. [2](#), [5](#), [6](#)