

More Reliable Pseudo-labels, Better Performance: A Generalized Approach to Single Positive Multi-label Learning

Luong Tran
FPT Software AI Center
luongtk@fpt.com

Thieu Vo
National University of Singapore
thieuvo@nus.edu.sg

Anh Nguyen
University of Liverpool
anh.nguyen@liverpool.ac.uk

Sang Dinh
Hanoi University of Science and Technology
sangdv@soict.hust.edu.vn

Van Nguyen
FPT Software AI Center
vanth19@fpt.com

Abstract

Multi-label learning is a challenging computer vision task that requires assigning multiple categories to each image. However, fully annotating large-scale datasets is often impractical due to high costs and effort, motivating the study of learning from partially annotated data. In the extreme case of Single Positive Multi-Label Learning (SPML), each image is provided with only one positive label, while all other labels remain unannotated. Traditional SPML methods that treat missing labels as unknown or negative tend to yield inaccuracies and false negatives, and integrating various pseudo-labeling strategies can introduce additional noise. To address these challenges, we propose the Generalized Pseudo-Label Robust Loss (GPR Loss), a novel loss function that effectively learns from diverse pseudo-labels while mitigating noise. Complementing this, we introduce a simple yet effective Dynamic Augmented Multi-focus Pseudo-labeling (DAMP) technique. Together, these contributions form the Adaptive and Efficient Vision-Language Pseudo-Labeling (AEVLP) framework. Extensive experiments on four benchmark datasets demonstrate that our framework significantly advances multi-label classification, achieving state-of-the-art results.

1. Introduction

Most visual classification studies focus on multi-class settings which aim to identify the most suitable label for a given input image from a set of possible options. However, in real-world scenarios, an image often contains multiple objects or attributes that can be considered as their labels. Therefore, multi-label learning [20, 31] has emerged as a solution to develop predictive models capable of assigning multiple labels to an unseen image. In early stages, multi-

label learning required annotation of all relevant classes for each training instance; however, such an exhaustive annotation scheme is often costly and impractical [7].

Due to the cost and impracticality of exhaustive annotation [7], researchers have explored multi-label learning with missing labels [25, 26, 29, 33]. In its extreme setting, known as Single-Positive Multi-Label Learning (SPML) [5], each image is annotated with only one confirmed positive label, leaving the remaining labels unknown. This approach significantly reduces labeling cost and demonstrates that multi-label classifiers can be learned with minimal supervision.

A common strategy in SPML is to treat missing labels as negative (the Assume Negative method [5]), and many existing methods [2, 19, 28, 34] further rely on the model’s internal knowledge to generate pseudo-labels for fully supervised learning. However, this approach often produces inherently noisy pseudo-labels, leading to repeated errors during training, an issue that prior studies have not effectively addressed. To tackle this challenge, we propose the Generalized Pseudo-label Robust Loss (GPR Loss), a novel loss function designed to learn effectively from pseudo-labels generated by an external knowledge-based method while robustly mitigating the noise associated with uncertain labels.

Recently, the emergence of Vision-Language Models (VLMs) has enhanced performance in various recognition tasks [30], including multi-label learning. One approach involves fine-tuning these models with additional network layers [8], while another leverages model outputs with pseudo-labels in an unsupervised, annotation-free manner [1]. Recent work [27] generates a fixed pseudo-label vector for each image using a pretrained large CLIP model. However, relying on a single fixed pseudo-label vector prevents the model from re-evaluating and correctly associating an image with categories that were initially misclassified as

false negatives during training.

Based on these observations, to verify the strength of GPR Loss, we introduce a simple yet effective dynamic pseudo-labeling process that encourages CLIP to attend not only to the global image but also to various local areas, which are randomly selected and augmented during training. This simple pseudo-labeling approach allows the pseudo-label vector for each image to vary from epoch to epoch, reducing the chance of missing correct associations.

In summary, the main contributions of this work are twofold. *First*, we propose the Generalized Pseudo-label Robust Loss (GPR Loss) to effectively learn from pseudo-labels generated by an external knowledge-based method while mitigating the noise from uncertain labels. *Second*, we introduce the Dynamic Augmented Multi-focus Pseudo-labeling (DAMP) technique to generate reliable pseudo-labels in a stable manner, which, together with GPR Loss, forms the AEVLP framework. Through extensive experiments, to the best of our knowledge, we demonstrate that this framework achieves state-of-the-art (SOTA) results for the SPML problem compared to existing methods based on loss design and pseudo-labeling.

2. Related Works

In SPML, two strategies are employed to handle missing labels. The first strategy treats missing labels as unknown variables to be predicted [23]. The second strategy assumes that all missing labels are negative, transforming SPML into a fully supervised multi-label learning problem [5, 13, 14, 28, 34]. The Assume Negative assumption remains popular in multi-label learning with missing labels including SPML [19] and often serves as a baseline in benchmark experiments.

However, this naive assumption oversimplifies the problem and can lead to numerous false negative labels [32]. To address this issue, Cole et al. [5] propose the online estimation of labels during training. Furthermore, Kim et al. [13] show that reducing large losses, which are potentially introduced by false negatives, can significantly improve performance. Meanwhile, the authors in [14] modify the final stage of the CNN architecture to increase the attribution scores of output logits. In [34], entropy maximization (EM) and asymmetric pseudo-labeling (APL) are combined to improve robustness. Additionally, Chen et al. [2] adapt class- and instance-wise loss concepts to SPML within a comprehensive framework. In [19, 28], the approach introduces a label enhancement process, enhancing model performance with minimal positive labels per instance. While our study continues to rely on the Assume Negative assumption and addresses the resulting false negatives, we also focus on learning from pseudo-labels and mitigating the noise they may introduce.

To address the challenges caused by incomplete label-

ing in multi-label learning, the ability of the CLIP model [22] for zero-shot classification has emerged as an effective solution. Recent work by [27] forces the model to learn from pseudo-labels generated by CLIP. Another approach [8] finetunes CLIP model with additional layers of a Graph Convolutional Network (GCN) [16]. Additionally, in [1], the approach is to generate soft pseudo-labels based on global-local image-text similarity aggregation. Our approach also leverages CLIP for pseudo-labeling; however, unlike previous methods, it employs dynamic pseudo-labeling throughout the training process, re-assigning each image a *new* set of pseudo-labels at every epoch. This helps mitigate pseudo-labeling errors while expanding the range of correct pseudo-labels.

3. Problem Statement

In the original multi-label learning problem, given a dataset $\mathcal{D} = \{(x_n, y_n)\}_{n=1}^N$, where $x_n \in \mathcal{X}$ (input space) and $y_n \in \mathcal{Y} = \{0, 1\}^C$ (label space), and a validation dataset $\mathcal{D}^{\text{val}} = \{(x_n^{\text{val}}, y_n^{\text{val}})\}_{n=1}^{N'}$, the goal is to learn a mapping $f: \mathcal{X} \rightarrow \mathcal{Y}$ that predicts multiple labels for each x_n . We denote $p_n = f(x_n) = \sigma(s_n)$, where s_n represents the model’s logits and σ is the sigmoid activation function. The loss function is defined as follows:

$$\mathcal{L} = -\frac{1}{NC} \sum_{n=1}^N \sum_{i=1}^C y_{n,i} \log p_{n,i} + (1 - y_{n,i}) \log(1 - p_{n,i})$$

To address the SPML problem, we apply the Assume Negative assumption [5], which treats unobserved labels as negative, transforming the problem into the original multi-label learning setup. The original dataset \mathcal{D} is replaced with a pseudo-multi-label dataset $\hat{\mathcal{D}} = \{(x_n, \hat{y}_n)\}_{n=1}^N$, where $\hat{y}_n \in \hat{\mathcal{Y}} = \{0, 1\}^C$, ensuring $\sum_{i=1}^C \hat{y}_{n,i} = 1$. Each $\hat{y}_{n,i} = 1$ if x_n is relevant to class i (positive label) and 0 if the label is unknown or unannotated.

In SPML, based on pseudo-labeling, an instance can be represented as a triplet (x_n, \hat{y}_n, l_n) , where x_n is the input, \hat{y}_n is the label vector based on the Assume Negative assumption, and $l_n \in \{-1, 0, 1\}^C$ represents pseudo-labels generated by an arbitrary method \mathcal{M} to approximate additional labels. For class i , $l_{n,i} = 1$ indicates a positive pseudo-label, $l_{n,i} = -1$ a negative pseudo-label, and $l_{n,i} = 0$ means the label is undefined by \mathcal{M} .

4. Generalized Pseudo-Label Robust Loss

4.1. Background

To address the noise introduced in the Assume Negative assumption, the authors in [2] propose the Generalized Robust Loss (GR Loss) as follows:

$$\mathcal{L}^{GR} = \frac{1}{NC} \sum_{n=1}^N \sum_{i=1}^C v^{old}(p_{n,i}; \alpha) \cdot \mathcal{L}_{n,i}^{old} \quad (1)$$

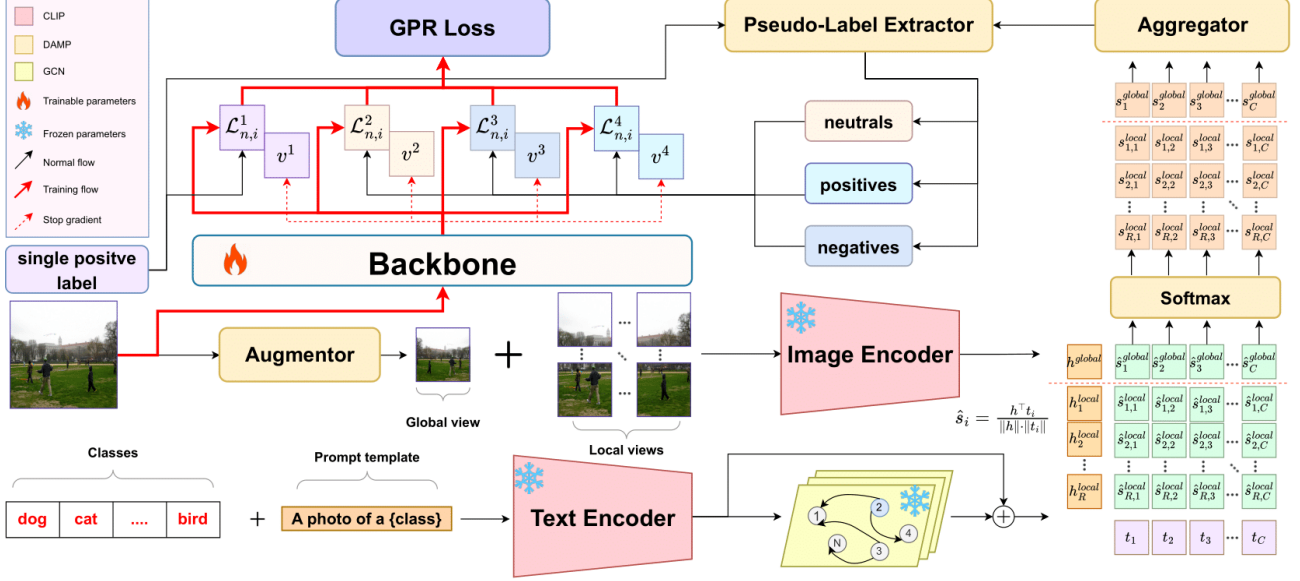


Figure 1. An overview of the proposed method - Adaptive and Efficient Vision-Language Pseudo-Labeling Framework (AEVLP)

where the probability of class i and instance n is denoted as $p_{n,i}$; $v^{old}(p_{n,i}; \alpha)$ and $\mathcal{L}_{n,i}^{old}$ are the class-and-instance-specific weight and loss, respectively. Let $\mathbb{1}_{[\cdot]}$ is the indicator function, the weight term $v^{old}(p_{n,i}; \alpha)$ is expressed as:

$$v^{old}(p_{n,i}; \alpha) = \mathbb{1}_{[\hat{y}_{n,i}=1]} v^1(p_{n,i}; \alpha) + \mathbb{1}_{[\hat{y}_{n,i}=0]} v^2(p_{n,i}; \alpha),$$

where $v^1(p; \alpha) = 1$ and $v^2(p; \alpha) = \exp\left(-\frac{(p-\mu)^2}{2\sigma^2}\right)$. Here, $\alpha = [\sigma, \mu]$ is updated linearly over training epochs, following [2], as detailed in the supplementary materials. The loss term $\mathcal{L}_{n,i}^{old}$ is defined as:

$$\mathcal{L}_{n,i}^{old} = \mathbb{1}_{[\hat{y}_{n,i}=1]} \mathcal{L}_{n,i}^1 + \mathbb{1}_{[\hat{y}_{n,i}=0]} \mathcal{L}_{n,i}^2,$$

where $\mathcal{L}_{n,i}^1 = -\log p_{n,i}$ and the term $\mathcal{L}_{n,i}^2$ is computed as:

$$\mathcal{L}_{n,i}^2 = (1 - \hat{k}(p_{n,i}; \beta)) \frac{1 - (1 - p_{n,i})^{q_1}}{q_1} + \hat{k}(p_{n,i}; \beta) \frac{1 - p_{n,i}^{q_2}}{q_2}.$$

Here, $\hat{k}(p_{n,i}; \beta)$ is the estimate of the label being a false negative during training with a parameter β , while q_1 and q_2 , treated as hyperparameters, balance the trade-off between MAE and BCE loss, as described in [2] and illustrated in the supplementary materials. Note that both $\hat{k}(p_{n,i}; \beta)$ and $v(p_{n,i}; \alpha)$ are implemented with gradients stopped with respect to $p_{n,i}$.

4.2. Proposed Method

The above method, GR Loss in Sec. 4.1, only works with the naive Assume Negative assumption and cannot handle external pseudo-labels generated by a reliable method \mathcal{M} .

Inspired by [2, 27], we propose the Generalized Pseudo-Label Robust Loss (GPR Loss), as follows:

$$\mathcal{L}^{GPR} = \frac{1}{NC} \sum_{n=1}^N \sum_{i=1}^C v^{new}(p_{n,i}; \alpha) \cdot \mathcal{L}_{n,i}^{new} + \eta \mathcal{R}. \quad (2)$$

We explain the terms in Eq. (2) below. The loss term $\mathcal{L}_{n,i}^{new}$ and the weight term $v^{new}(p_{n,i}; \alpha)$ are expressed as follows:

$$\begin{aligned} \mathcal{L}_{n,i}^{new} &= \mathbb{1}_{[\hat{y}_{n,i}=1]} \mathcal{L}_{n,i}^1 + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=0]} \mathcal{L}_{n,i}^2 \\ &\quad + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=-1]} \mathcal{L}_{n,i}^3 + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=1]} \mathcal{L}_{n,i}^4, \\ v^{new}(p_{n,i}; \alpha) &= \mathbb{1}_{[\hat{y}_{n,i}=1]} v^1(p_{n,i}; \alpha) \\ &\quad + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=0]} v^2(p_{n,i}; \alpha) \\ &\quad + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=-1]} v^3(p_{n,i}; \alpha) \\ &\quad + \mathbb{1}_{[\hat{y}_{n,i}=0 \wedge l_{n,i}=1]} v^4(p_{n,i}; \alpha). \end{aligned}$$

The new loss terms $\mathcal{L}_{n,i}^3$ and $\mathcal{L}_{n,i}^4$ are given by:

$$\begin{aligned} \mathcal{L}_{n,i}^3 &= -\log(1 - p_{n,i}), \\ \mathcal{L}_{n,i}^4 &= -(1 - q_3) \log(1 - p_{n,i}) - q_3 \log(p_{n,i}). \end{aligned} \quad (3)$$

To provide a clearer understanding of our loss function's design, we outline how it handles various cases based on label status as follows.

Confirmed Positives ($\hat{y}_{n,i} = 1$): The loss term $\mathcal{L}_{n,i}^1$ is retained to ensure that the model assigns high probability to the true positive classes.

Undefined Pseudo-labels ($\hat{y}_{n,i} = 0 \wedge l_{n,i} = 0$): The term $\mathcal{L}_{n,i}^2$ (inherited from GR Loss) is used when no additional information from pseudo labeling is provided, and

the label is naively assumed to be negative due to the high negative rate.

Negative Pseudo-labels ($\hat{y}_{n,i} = 0 \wedge l_{n,i} = -1$): The loss term $\mathcal{L}_{n,i}^3$, unlike $\mathcal{L}_{n,i}^2$ which includes the control of $\hat{k}(p_{n,i})$, directly penalizes the model for assigning a high probability to classes that are considered negative by the pseudo-labels. Moreover, by employing the same weight function $v^3(p; \alpha) = \exp\left(-\frac{(p-\mu)^2}{2\sigma^2}\right)$, which is also used in $\mathcal{L}_{n,i}^2$ for undefined pseudo-labels, the method balances the contributions between positive and negative cases.

Positive Pseudo-labels ($\hat{y}_{n,i} = 0 \wedge l_{n,i} = 1$): The loss term $\mathcal{L}_{n,i}^4$ is a softened version of $\mathcal{L}_{n,i}^1$ through the coefficient q_3 . This label smoothing helps to temper the model’s confidence, reducing overfitting to noisy pseudo-labels while still promoting the learning of positive associations when a pseudo-label is provided. The weight function

$$v^4(p; \alpha) = \min\left(\max\left(1 - \exp\left(-\frac{(p-\mu)^2}{2\sigma^2}\right), \lambda_1\right), \lambda_2\right)$$

constrains the influence of this term within a desirable range, controlled by λ_1 and λ_2 .

For stable learning when introducing additional pseudo-labels, a regularization term \mathcal{R} is used to restrict the number of positive predictions with a coefficient η , as in [5]:

$$\mathcal{R} = \left(\frac{\hat{m} - m}{C}\right)^2, \quad (4)$$

where m is the expected number of positive labels per image that can either be estimated from the available data or treated as a hyperparameter, such as:

$$m \approx E_{pos} = \mathbf{E}_{(x,y) \sim p_{data}(x,y)} \sum_{i=1}^C \mathbb{1}_{[y_i=1]}, \quad (5)$$

and \hat{m} represents the average sum of probabilities per image, computed as:

$$\hat{m} = \frac{1}{N} \sum_{n=1}^N \sum_{i=1}^C p_{n,i}. \quad (6)$$

Our loss function \mathcal{L}^{GPR} is a generalization of \mathcal{L}^{GR} as shown in the following theorem.

Theorem 4.1 (Our GPR Loss generalizes GR Loss). *Given a training dataset $\{(x_n, \hat{y}_n, l_n)\}_{n=1}^N$ and a validation set $\{(x_n^{val}, y_n^{val})\}_{n=1}^{N'}$ as described in Sec. 3, define:*

$$\mathbf{C}(\mathcal{M}) = \exp\left(\sum_{n=1}^N \sum_{i=1}^C \mathbb{1}_{[\hat{y}_{n,i}=1]} \log P(l_{n,i} = 1 \mid x_n, \mathcal{M})\right),$$

and:

$$m' = \frac{1}{N'} \sum_{n=1}^{N'} \sum_{i=1}^C \mathbb{1}_{[y_{n,i}^{val}=1]}.$$

Then \mathcal{L}^{GPR} tends to \mathcal{L}^{GR} when $\max(\mathbf{C}(\mathcal{M}), |m' - \hat{m}|)$ tends to zero. Here, \hat{m} is defined in Eq. (6).

Remark 1. *Intuitively, the above theorem says that our loss function \mathcal{L}^{GPR} will become the loss function \mathcal{L}^{GR} when $\max(\mathbf{C}(\mathcal{M}), |m' - \hat{m}|)$ tends to zero. The quantity $\mathbf{C}(\mathcal{M})$ represents the confidence of method \mathcal{M} within the range $(0, e]$ when generating pseudo-labels based on the given single positive label ($\hat{y}_{n,i} = 1$). Lower values of $\mathbf{C}(\mathcal{M})$, particularly those below ϵ , indicate a significant divergence between the pseudo-label probability distribution and the ground truth distribution, leading to $l_{n,i} = 0$ for all n, i .*

Remark 2. *When N' is sufficiently large, $m' \rightarrow E_{pos}$, with E_{pos} defined in Eq. (5), then $\hat{m} \rightarrow E_{pos}$. As a result, the regularization term in Eq. (2) vanishes, i.e., $\mathcal{R} \rightarrow 0$. This demonstrates that the GPR Loss naturally calibrates itself when the pseudo-label confidence and the expected positive label statistics align.*

The proof of this theorem can be found in the supplementary materials.

5. DAMP

In this section, we introduce a simple yet effective pseudo-labeling method that integrates with GPR Loss via a Dynamic Augmented Multi-focus Pseudo-labeling (DAMP) approach to form the framework AEVLP, as shown in Fig. 1, for SPML. Further implementation details are provided in the supplementary materials.

5.1. CLIP Inference and Strengthening with Noise

As introduced in [22], given an image input x and the i -th class from a set of C classes, the corresponding visual embedding and textual embedding are $h = E_v(x) \in \mathbb{R}^K$ and $t_i = E_t(\mathcal{P}_i) \in \mathbb{R}^K$, respectively. Here, E_v and E_t are the image and text encoders of the CLIP model with dimension K , and \mathcal{P}_i is a predefined prompt for class i , such as “a photo of a {class}”. Let $\hat{s}_i = \text{Cos}(h, t_i)$, $s_i = \text{Softmax}\left(\frac{\hat{s}_i}{\tau}\right)$, and denote $S = \{s_1, s_2, \dots, s_C\}$ as the probability distribution for x across C classes.

Several works, including [12, 17, 35], have studied enhancing model performance during fine-tuning by adding controlled noise to model embeddings. Additionally, in [3, 8], label-to-label relationships are presented by GCN. Inspired by this, we propose adding controlled label-to-label correspondence noise to the text embeddings of the CLIP model. Concretely, we redefine the text embedding and update the GCN as: $t_i = G(E_t(\mathcal{P}_i)) + E_t(\mathcal{P}_i)$ and $H_{l+1} = \text{LeakyReLU}(A^* H_l W_l)$ for $l \in [0, L)$, with $H_0 = \{E_t(\mathcal{P}_i) \mid 1 \leq i \leq C\}$. The graph G remains frozen during training, and the weights W_l are initialized from a uniform distribution. The adjacency matrix A^* is derived from the cosine similarity scores between the text embeddings of the classes produced by the CLIP text encoder.

5.2. Dynamic Augmented Multi-focus Pseudo-labeling

Augmentor. Let I be an image. It is partitioned into overlapping patches $\{P_z\}_{z=1}^R$ by dividing it into a grid and slightly enlarging each patch by a random ratio to ensure overlap. We then process the image and its patches using a transformation pipeline $T(\cdot)$, which includes standard pre-processing and weak data augmentation techniques, to generate various views for CLIP: the global view $x^{global} = T(I)$ and the local views $x_z^{local} = T(P_z)$. Following Sec. 5.1, these views yield the probability distributions $S^{global} = \{s_1^{global}, s_2^{global}, \dots, s_C^{global}\}$ and $S_z^{local} = \{s_{z,1}^{local}, s_{z,2}^{local}, \dots, s_{z,C}^{local}\}$.

Local threshold based on single positives. Let \hat{c} be the given single positive label, according to the SPML setting in Sec. 3, for the image I . The local threshold ζ^{local} , which defines the patches to be trusted, is adjusted based on $\zeta^{local} = \min(s_{\hat{c}}^{global}, \nu)$, where ν is the general local threshold, set as a hyperparameter. In some cases, if ν is set too high to recognize hard positives, we should consider the scores above $s_{\hat{c}}^{global}$, as these can be meaningful, since \hat{c} is one of the true labels of the global view.

Aggregator. From the local distributions $\{S_z^{local}\}_{z=1}^R$, we aggregate a unified local distribution S^{agg} following [1]. For each class c , we compute $\omega_c = \max_{z=1, \dots, R} s_{z,c}^{local}$ and $\psi_c = \min_{z=1, \dots, R} s_{z,c}^{local}$, and define the aggregation score as $s_c^{agg} = \mathbb{1}_{[\omega_c \geq \zeta^{local}]} \omega_c + \mathbb{1}_{[\omega_c < \zeta^{local}]} \psi_c$; this yields the soft aggregation vector $S^{agg} = \{s_1^{agg}, s_2^{agg}, \dots, s_C^{agg}\}$ for each input image.

Positive pseudo-labels. To extract reliable positive pseudo-labels, we integrate both global and aggregated local similarities into $S^{final} = \frac{1}{2}(S^{global} + S^{agg})$. Let $Q' = \{l'_1, l'_2, \dots, l'_C\}$ be the pseudo labels of the image I . We convert the soft similarity scores S^{final} into hard pseudo-labels as follows:

$$l'_c = \begin{cases} 1, & s_c^{final} \in \text{TopK}(S^{final}, k) \ \& \ s_c^{final} \geq \zeta^{global} \\ 0, & \text{otherwise,} \end{cases}$$

where ζ^{global} is the global threshold for high-confidence positive pseudo-labels, set as a hyperparameter, and k limits the number of positive pseudo-labels.

Negative pseudo-labels. To identify potential negative pseudo-labels we compute average similarity scores as:

$$S^{avg} = \frac{1}{2} \left(S^{global} + \frac{1}{R} \sum_{z=1}^R S_z^{local} \right).$$

We use S^{avg} to refine Q' by assigning negative pseudo-labels, producing the final pseudo-labels $Q = \{l_1, l_2, \dots, l_C\}$. A class c is designated as a negative pseudo-label if its score s_c^{avg} falls within the lowest $\Delta_{neg}\%$ of values in S^{avg} . Assuming a potential negative pseudo-label has low scores in both image I and every patch P_z according to the VLM, we define the assignment as:

$$l_c = \begin{cases} -1, & s_c^{avg} \leq \theta_{\Delta_{neg}}(S^{avg}) \\ l'_c, & \text{otherwise} \end{cases}$$

where $\theta_{\Delta_{neg}}(S^{avg})$ denotes the Δ_{neg} -th-percentile of S^{avg} , serving as the threshold to identify the lowest $\Delta_{neg}\%$ of values in S^{avg} as negative pseudo-labels.

6. Experiments

6.1. Setup

Datasets. In this study, our method is evaluated through environmental experiments similar to those in [2, 5, 13, 14, 27] across four standard benchmark datasets: PASCAL VOC 2012 (VOC) [10], MS-COCO 2014 (COCO) [18], NUS-WIDE (NUS) [4], and CUB-200-2011 (CUB) [24]. From these fully labeled multi-label datasets, we simulate single positive training data by randomly retaining one positive label per training example, using the same seed as in [5]. Twenty percent of the training set for each dataset is set aside for validation purposes. Both the validation and test sets are fully labeled.

Implementation details. For a fair comparison, we follow the standard SPML implementation described in [2, 5, 14, 27], using the ResNet-50 architecture [11], pre-trained on ImageNet [6]. Each image is resized to 448×448 and augmented with random horizontal flipping. For the VLPL method in [27], we use CLIP ViT-L/16, while our approach uses CLIP ViT-B/16, both adopted from [22] to balance pseudo-labeling quality and training time. We train the models for 10 epochs on CUB and NUS, and 8 epochs on COCO and VOC, using the Adam optimizer [15]. The batch size is set to 8 for CUB and VOC, and 16 for NUS and COCO. A learning rate of 1×10^{-5} is used across all experiments. We evaluate performance using the mean Average Precision (mAP) metric, conducting the final evaluation on the test set with the model that achieves the highest performance on the withheld validation set, in line with previous methods [2, 5, 13, 14].

Baselines. For a comprehensive understanding, we compare our proposed method, AEVLP, against several state-of-the-art methods within the SPML setting. Specifically, we evaluated against the following methods: Assume Negative [5], ROLE [5], EM [34], EM + APL [34], BoostLU + LL-R

Table 1. Experimental results on various benchmarks with SPML setting according to [5]. Our results are averaged over three runs as suggested in [14].

Method	VOC	COCO	NUS	CUB
Full-label (BCE)	89.42	76.78	52.08	30.90
Assume Negative (CVPR’21) [5]	85.89	64.92	42.27	18.31
ROLE (CVPR’21) [5]	87.77	67.04	41.63	13.66
EM (ECCV’22) [34]	89.09	70.70	47.15	20.85
EM + APL (ECCV’22) [34]	89.19	70.87	47.59	21.84
BoostLU + LL-R (CVPR’23) [14]	89.29	72.89	49.59	19.8
SMILE (NeurIPS’22) [28]	86.31	63.33	43.61	18.61
MIME (ICML’23) [19]	89.20	72.92	48.74	21.89
GR Loss (IJCAI’24) [2]	89.83	73.17	49.08	21.64
VLPL (CVPRW’24) [27]	89.10	71.45	49.55	24.02
AEVLP (Ours)	90.46	73.54	50.70	24.89

Table 2. Comparison of mAP scores when applying GPR Loss to different pseudo-labeling strategies.

Method	VOC	COCO	NUSWIDE	CUB	Average
BCE + Random	86.84	67.46	41.23	10.69	51.56
GPR + Random	87.33	70.00	43.60	14.98	53.98 (+ 2.42)
VLPL	89.1	71.45	49.55	24.02	58.53
GPR + VLPL	89.52	72.83	49.93	24.22	59.13 (+ 0.6)
LL-Ct	89.00	70.50	48.00	20.40	56.98
GPR + LL-Ct	89.5	71.66	47.72	20.41	57.32 (+ 0.34)
BCE + DAMP	88.72	71.89	48.60	24.01	58.31
GPR + DAMP	90.46	73.54	50.7	24.89	59.90 (+ 1.59)

[14], SMILE [28], MINE [19], GR Loss [2], VLPL [27]. All methods above, including our AEVLP, were evaluated on the same benchmark datasets: VOC, COCO, NUS, and CUB for a fair comparison. The results were averaged over three runs, as recommended in previous studies [14].

6.2. Results and Discussion

Table 1 shows experimental results on four benchmark datasets: VOC, COCO, NUS, and CUB, using the previously described settings. Our method, AEVLP, achieves state-of-the-art performance across all datasets, demonstrating both effectiveness and robustness. Notably, it surpasses fully labeled multi-label classification on the VOC dataset. Although AEVLP can be seen as a fusion of GR Loss and VLPL, it not only combines their strengths but also surpasses both methods. VLPL, in particular, shows limitations in 3 out of 4 benchmark datasets compared with traditional state-of-the-art SPML methods. In contrast, AEVLP demonstrates superior performance, underscoring the effectiveness of introducing different focuses of the same input image to a VLM for pseudo-label generation. Notably, VLPL uses the CLIP model with a ViT Large backbone (428M parameters) [9], with 14×14 patches and a 336-pixel resolution, whereas AEVLP achieves better results while using the CLIP model with a smaller ViT Base backbone

Table 3. DAMP performance on several benchmark datasets.

Metrics	VOC	COCO	NUS	CUB
Average Precision	65.13	84.79	37.09	19.07
Accumulated Precision	60.33	81.59	34.49	18.29
Average Recall	48.29	24.68	8.84	11.37
Accumulated Recall	51.1	26.84	10.06	13.88

(149.62M parameters) with 16×16 patches and a 224×224 input resolution (ViT-B/16). As can be seen from Table 1, traditional SPML methods [2, 5, 19, 28] based on pseudo-labeling approach struggle with datasets containing a large number of labels, such as CUB (312 labels). AEVLP, in contrast, shows a relative improvement of more than 15% over GR Loss method [2] and 3.62% over the next best-performing method, VLPL [27], extending GR Loss’s capability to work with pseudo-labels. Our proposed GPR Loss, described in Sec. 4.2, can cope with pseudo-labels generated by various mechanisms, as long as they provide an acceptable level of confidence.

6.3. Analysis

GPR Loss robustness. GPR Loss is designed to work with various pseudo-labeling strategies, not only by effectively utilizing pseudo-labels but also by demonstrating robustness to noise. To validate this, we apply GPR Loss to pseudo-labels generated by different methods, including DAMP (AEVLP), Random Pseudo-labels, VLPL [27], and LL-Ct [13]. Random Pseudo-labels are obtained by randomly selecting unknown labels (assumed negative) as pseudo-positive labels, based on the average number per image in the dataset. VLPL provides pseudo-labels generated from the CLIP model, while LL-Ct converts assumed negative labels, particularly those with high loss values, into pseudo-positive labels. For the baseline of DAMP and Random Pseudo-labels, we use BCE for training; for the

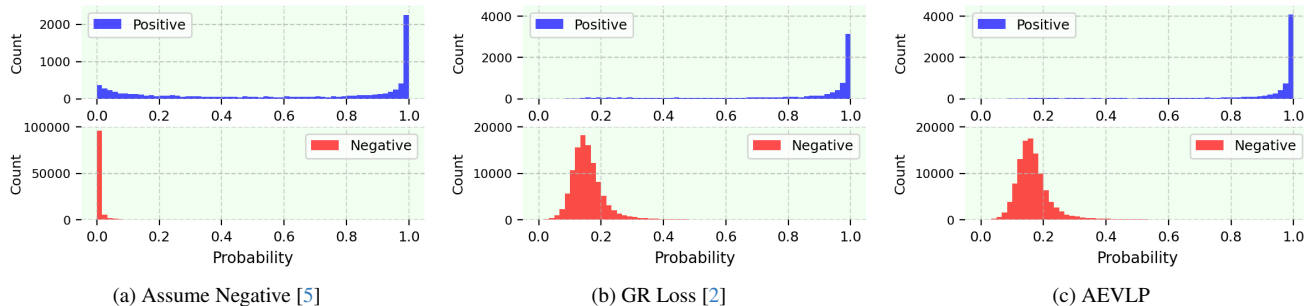


Figure 2. Distribution of output probabilities for the positive and negative classes of the VOC test set, as predicted by the ResNet-50 backbone classifier trained in the SPML setting.

Table 4. Ablation study on the main components of the AEVLP framework: noise addition $G(\cdot)$, augmentation $T(\cdot)$, overlapping ratio r , loss re-weighting for positive pseudo-labels $v^4(p; \alpha)$, regularization term \mathcal{R} , and negative pseudo-labels loss $\mathcal{L}_{n,i}^3$. Note that \times for $v^4(p; \alpha)$ indicates it is set to a constant, and for r , it indicates $r = 0$.

AEVLP						Datasets			
DAMP			GPR Loss			VOC	COCO	NUS	CUB
$G(\cdot)$	$T(\cdot)$	r	$v^4(p; \alpha)$	\mathcal{R}	$\mathcal{L}_{n,i}^3$				
\times	\checkmark	\times	\times	\times	\times	90.28	73.07	50.00	24.14
\checkmark	\checkmark	\times	\times	\times	\times	90.21	73.37	50.32	24.27
\checkmark	\times	\checkmark	\times	\times	\times	90.13	73.32	50.46	24.15
\checkmark	\checkmark	\checkmark	\times	\times	\times	90.22	73.14	50.57	24.24
\checkmark	\checkmark	\checkmark	\checkmark	\times	\times	90.31	73.40	50.37	24.73
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\times	90.38	73.51	50.64	24.83
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	90.46	73.54	50.70	24.89

other pseudo-labeling methods, we replace their original loss function with GPR Loss. As shown in Tab. 2, incorporating GPR Loss consistently enhances mAP scores across all datasets for each pseudo-labeling strategy. This highlights GPR Loss’s effectiveness in mitigating the inherent noise in pseudo-labels, regardless of the generation method.

DAMP performance. We apply pseudo-labeling to the training data of the VOC, COCO, NUS, and CUB datasets using the DAMP technique, evaluating over the same number of epochs used for training the main results in Tab. 1. Since pseudo-labels of a given image may vary across epochs, we report the results as average precision and recall across all epochs to assess the supervision quality that DAMP provides for training the backbone. For strict evaluation and to address label imbalance, we focus only on the missing positive labels in the training data according to the SPML setting, ignoring predictions for the confirmed single positive and true negative labels. To validate coverage of missing positive labels, we accumulate predictions across epochs: each new positive pseudo-label prediction is recorded, and precision and recall are computed on these accumulated positive labels. As illustrated in Tab. 3, DAMP, a strategy of randomization, enhances positive label cov-

erage, with accumulated recall exceeding average recall, while maintaining a stable quality of pseudo-labels. Notably, precision across epochs varies only slightly from the average precision (see supplementary materials). In contrast, using fixed pseudo-labels results in lower coverage, and a straightforward accumulation approach introduces more noise.

Probability distribution analysis. Figure 2 plots the distribution of the classifier output probabilities on the VOC test set from three different methods: Assume Negative, GR Loss, and AEVLP (ours). The ResNet-50 backbone is used in all three methods. Compared to existing method in Tab. 1, the Assume Negative method represents a naive approach, while GR Loss shows robustness across all datasets. As shown in Fig. 2a, the Assume Negative method performs the worst: although it maintains high confidence for negative labels (close to 0), it includes numerous low-confidence predictions for positive labels, which can lead to significant errors. For both GR Loss in Fig. 2b and AEVLP in Fig. 2c, the probability distributions of negative labels have a similar shape, concentrated around 0.2. However, AEVLP outperforms GR Loss, with its distribution of positive labels shifted closer to 1, even when learning from pseudo-labels

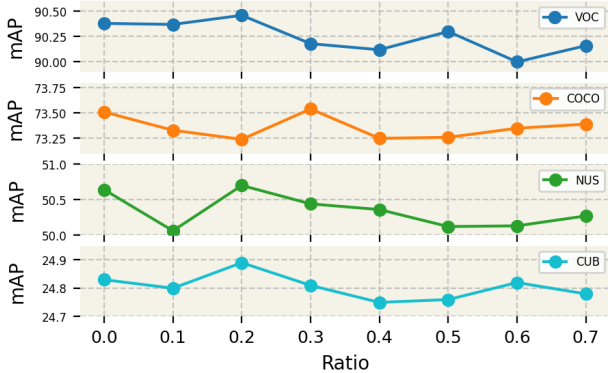


Figure 3. The performance of the model in learning from varying ratios of negative pseudo-labels across different datasets using the AEVLP method.

that may introduce noise.

Mining negative pseudo-labels. Fig. 3 shows the impact of varying the negative pseudo-label ratio $\Delta_{neg}\%$ on model performance. At approximately 20% to 30% negative pseudo-labels, the model achieves its best mAP performance across all datasets. This suggests that introducing a moderate proportion of negative pseudo-labels helps improve diversity in the learning process without significantly increasing noise. As the negative pseudo-label ratio increases further, the mAP slightly decreases but remains relatively stable. This is likely due to the influence of GPR Loss, which helps mitigate the negative impact of potentially inaccurate pseudo-labels in situations where there is an imbalance between positive and negative labels.

6.4. Ablation Study

Components of AEVLP. In this section, we analyze the impact of each component in the AEVLP framework, as shown in Tab. 4. The ablation results in Tab. 4 shows that each component contributes to the overall performance improvement. Some observations are relevant here to describe the relationships between closely linked components that complement each other. The combination of noise addition $G(\cdot)$ with either augmentation $T(\cdot)$ or overlapping r provides a noticeable boost, as these pairs introduce valuable pseudo-labels that may be difficult to capture with the standard input alone. The positive pseudo-label re-weighting $v^4(p; \alpha)$ and regularization term \mathcal{R} help stabilize learning from uncertain pseudo-labels. Additionally, the negative pseudo-label loss $\mathcal{L}_{n,i}^3$ encourages a stronger focus on potential negative labels, thereby enhancing discrimination.

Impact of the grid size. We evaluate grid sizes g (from 2 to 6) on four datasets: VOC (20 labels), COCO (80 labels), NUS (81 labels), and CUB (312 labels). The results (see

Table 5. mAP scores for various grid sizes g on VOC and CUB datasets.

Grid size	2	3	4	5	6
VOC	90.13	90.08	90.46	90.21	90.21
COCO	73.07	73.13	73.54	73.48	73.32
NUS	50.25	50.19	50.7	50.35	50.15
CUB	24.37	24.38	24.48	24.89	24.67

Table 6. AEVLP performance on various backbone architectures. Note that 1K and 22K refer to pretrained models on ImageNet-1K and ImageNet-22K, respectively.

Backbone	ResNet-50	ConvNeXt-L-1K	ConvNeXt-L-22K	ViT-L/14
VOC	90.46	92.92	93.11	93.10
COCO	73.54	76.14	82.69	81.39
NUS	50.7	50.9	54.74	54.98
CUB	24.89	25.26	24.87	23.36

Tab. 5) show that a grid size of 4 yields the best performance for VOC, COCO, and NUS, suggesting that a moderate grid size effectively balances spatial detail and computational efficiency. In contrast, CUB achieves its highest mAP with a grid size of 5, indicating that datasets with finer-grained labels benefit from slightly larger grid sizes to capture richer object details. While increasing grid size reduces the number of labels per patch, enhancing label separability, excessively small patches may fail to capture entire object instances and incur higher computational costs. These findings underscore the need to balance label reduction, object capture, and efficiency when selecting the optimal grid size.

Backbone architecture. We evaluate several backbone architectures, as shown in Tab. 6. Specifically, we replace the baseline backbone with ConvNeXt-L [21], using two different pretrained versions (ImageNet-1K and ImageNet-22K), as well as ViT-L/14 [22]. The results demonstrate a notable improvement in performance with these architectures. This analysis confirms AEVLP’s flexibility, showing it can effectively incorporate advanced models to achieve substantial performance gains.

7. Conclusion

In this paper, we propose AEVLP, a novel framework for SPML. First, our approach introduces a new GPR loss function that effectively complements the pseudo-labels. We then present a dynamic augmented multi-focus pseudo-labeling strategy designed to overcome the limitations of fixed pseudo-labeling approaches, where the pseudo-labels for each image remain constant throughout training. Experimental results demonstrate that AEVLP achieves state-of-the-art performance on four benchmark datasets, reducing reliance on fully annotated data while addressing the shortcomings of traditional multi-label learning methods.

References

- [1] Rabab Abdelfattah, Qing Guo, Xiaoguang Li, Xiaofeng Wang, and Song Wang. Cdul: Clip-driven unsupervised learning for multi-label image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1348–1357, 2023. 1, 2, 5
- [2] Yanxi Chen, Chunxiao Li, Xinyang Dai, Jinhuan Li, Weiyu Sun, Yiming Wang, Renyuan Zhang, Tinghe Zhang, and Bo Wang. Boosting single positive multi-label classification with generalized robust loss. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24*, pages 3825–3833. International Joint Conferences on Artificial Intelligence Organization, 2024. Main Track. 1, 2, 3, 5, 6, 7
- [3] Zhao-Min Chen, Xiu-Shen Wei, Peng Wang, and Yanwen Guo. Multi-label image recognition with graph convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5177–5186, 2019. 4
- [4] 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. 5
- [5] Elijah Cole, Oisín Mac Aodha, Titouan Llorca, Pietro Perona, Dan Morris, and Nebojsa Jojic. Multi-label learning from single positive labels. In *CVPR-21*, 2021. 1, 2, 4, 5, 6, 7
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 5
- [7] Jia Deng, Olga Russakovsky, Jonathan Krause, Michael S Bernstein, Alex Berg, and Li Fei-Fei. Scalable multi-label annotation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pages 3099–3102, 2014. 1
- [8] 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. 1, 2, 4
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 6
- [10] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88:303–338, 2010. 5
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [12] Neel Jain, Ping-yeh Chiang, Yuxin Wen, John Kirchenbauer, Hong-Min Chu, Gowthami Somepalli, Brian R. Bartoldson, Bhavya Kailkhura, Avi Schwarzschild, Aniruddha Saha, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Neftune: Noisy embeddings improve instruction finetuning. *CoRR*, abs/2310.05914, 2023. 4
- [13] 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
- [14] Youngwook Kim, Jae Myung Kim, Jieun Jeong, Cordelia Schmid, Zeynep Akata, and Jungwoo Lee. Bridging the gap between model explanations in partially annotated multi-label classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3408–3417, 2023. 2, 5, 6
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 5
- [16] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. 2
- [17] Kezhi Kong, Guohao Li, Mucong Ding, Zuxuan Wu, Chen Zhu, Bernard Ghanem, Gavin Taylor, and Tom Goldstein. Robust optimization as data augmentation for large-scale graphs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 60–69, 2022. 4
- [18] 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. 5
- [19] Biao Liu, Ning Xu, Jiaqi Lv, and Xin Geng. Revisiting pseudo-label for single-positive multi-label learning. In *International Conference on Machine Learning*, pages 22249–22265. PMLR, 2023. 1, 2, 6
- [20] Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W Tsang. The emerging trends of multi-label learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(11): 7955–7974, 2021. 1
- [21] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11976–11986, 2022. 8
- [22] 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, 4, 5, 8
- [23] Reshma Rastogi and Sayed Mortaza. Multi-label classification with missing labels using label correlation and robust

- structural learning. *Knowledge-Based Systems*, 229:107336, 2021. [2](#)
- [24] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011. [5](#)
- [25] Baoyuan Wu, Zhilei Liu, Shangfei Wang, Bao-Gang Hu, and Qiang Ji. Multi-label learning with missing labels. In *2014 22nd International Conference on Pattern Recognition*, pages 1964–1968, 2014. [1](#)
- [26] Baoyuan Wu, Siwei Lyu, and Bernard Ghanem. Ml-mg: Multi-label learning with missing labels using a mixed graph. In *Proceedings of the IEEE international conference on computer vision*, pages 4157–4165, 2015. [1](#)
- [27] Xin Xing, Zhexiao Xiong, Abby Stylianou, Srikumar Sastry, Liyu Gong, and Nathan Jacobs. Vision-language pseudo-labels for single-positive multi-label learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7799–7808, 2024. [1](#), [2](#), [3](#), [5](#), [6](#)
- [28] Ning Xu, Congyu Qiao, Jiaqi Lv, Xin Geng, and Min-Ling Zhang. One positive label is sufficient: Single-positive multi-label learning with label enhancement. *Advances in Neural Information Processing Systems*, 35:21765–21776, 2022. [1](#), [2](#), [6](#)
- [29] Hsiang-Fu Yu, Prateek Jain, Purushottam Kar, and Inderjit Dhillon. Large-scale multi-label learning with missing labels. In *International conference on machine learning*, pages 593–601. PMLR, 2014. [1](#)
- [30] Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024. [1](#)
- [31] Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering*, 26(8):1819–1837, 2013. [1](#)
- [32] Wenqiao Zhang, Changshuo Liu, Lingze Zeng, Bengchin Ooi, Siliang Tang, and Yueting Zhuang. Learning in imperfect environment: Multi-label classification with long-tailed distribution and partial labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1423–1432, 2023. [2](#)
- [33] Youcai Zhang, Yuhao Cheng, Xinyu Huang, Fei Wen, Rui Feng, Yaqian Li, and Yandong Guo. Simple and robust loss design for multi-label learning with missing labels. *arXiv preprint arXiv:2112.07368*, 2021. [1](#)
- [34] 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. [1](#), [2](#), [5](#), [6](#)
- [35] Chen Zhu, Yu Cheng, Zhe Gan, Siqi Sun, Tom Goldstein, and Jingjing Liu. Freelb: Enhanced adversarial training for natural language understanding. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. [4](#)