

PiCO+: Contrastive Label Disambiguation for Robust Partial Label Learning

Haobo Wang, Ruixuan Xiao, Yixuan Li, Lei Feng, Gang Niu, Gang Chen, Junbo Zhao

Abstract—Partial label learning (PLL) is an important problem that allows each training example to be labeled with a coarse candidate set, which well suits many real-world data annotation scenarios with label ambiguity. Despite the promise, the performance of PLL often lags behind the supervised counterpart. In this work, we bridge the gap by addressing two key research challenges in PLL—representation learning and label disambiguation—in one coherent framework. Specifically, our proposed framework PiCO consists of a contrastive learning module along with a novel class prototype-based label disambiguation algorithm. PiCO produces closely aligned representations for examples from the same classes and facilitates label disambiguation. Theoretically, we show that these two components are mutually beneficial, and can be rigorously justified from an expectation-maximization (EM) algorithm perspective. Moreover, we study a challenging yet practical *noisy partial label learning* setup, where the ground-truth may not be included in the candidate set. To remedy this problem, we present an extension PiCO+ that performs distance-based clean sample selection and learns robust classifiers by a semi-supervised contrastive learning algorithm. Extensive experiments demonstrate that our proposed methods significantly outperform the current state-of-the-art approaches in standard and noisy PLL tasks and even achieve comparable results to fully supervised learning.

Index Terms—Partial Label Learning, Contrastive Learning, Prototype-based Disambiguation, Noisy Label Learning.

1 INTRODUCTION

THE training of modern deep neural networks typically requires massive labeled data, which imposes formidable obstacles in data collection. Of a particular challenge, data annotation in the real-world can naturally be subject to inherent label ambiguity and noise. For example, as shown in Figure 1, identifying an Alaskan Malamute from a Siberian Husky can be difficult for a human annotator. The issue of labeling ambiguity is prevalent yet often overlooked in many applications, such as web mining [1] and automatic image annotation [2]. This gives rise to the importance of *partial label learning* (PLL) [3], [4], where each training example is equipped with a set of candidate labels instead of the exact ground-truth label. This stands in contrast to its supervised counterpart where one label must be chosen as the “gold”. Arguably, the PLL problem is deemed more common and practical in various situations due to its relatively lower cost to annotations.

Despite the promise, a core challenge in PLL is label disambiguation, i.e., identifying the ground-truth label from the candidate label set. Existing methods typically require a



A dog image x_i with candidate labels Y_i :
Husky, Malamute,
Samoyed

Fig. 1. An input image with three candidate labels, where the ground-truth is Malamute.

good feature representation [5], [6], [7], and operate under the assumption that data points closer in the feature space are more likely to share the same ground-truth label. However, the reliance on representations has led to a non-trivial dilemma—the inherent label uncertainty can undesirably manifest in the representation learning process—the quality of which may, in turn, prevent effective label disambiguation. To date, few efforts have been made to resolve this.

This paper bridges the gap by reconciling the intrinsic tension between the two highly dependent problems—representation learning and label disambiguation—in one coherent and synergistic framework. Our framework, **P**artial label learning with **C**ontrastive label disambiguation (dubbed **PiCO**), produces closely aligned representations for examples from the same classes and facilitates label disambiguation. Specifically, PiCO encapsulates two key components. First, we leverage contrastive learning (CL) [8] to partial label learning, which is unexplored in previous PLL literature. To mitigate the key challenge of constructing positive pairs, we employ the classifier’s output and generate pseudo positive pairs for contrastive comparison

- Haobo Wang, Ruixuan Xiao, Gang Chen and Junbo Zhao are with Key Lab of Intelligent Computing Based Big Data of Zhejiang Province, Zhejiang University, Hangzhou, China, 310027. E-mail: {wanghaobo, xiaoruixuan, cg, j.zhao}@zju.edu.cn.
- Yixuan Li is with Department of Computer Sciences, University of WisconsinMadison, Madison, WI, United States, 53706. E-mail: sharonli@cs.wisc.edu.
- Lei Feng is with the College of Computer Science, Chongqing University, Chongqing, China, 400044. E-mail: lfeng@cqu.edu.cn.
- Gang Niu is with RIKEN Center for Advanced Intelligence Project, Tokyo, Japan. E-mail: gang.niu.ml@gmail.com.
- Junbo Zhao is the corresponding author.

(Section 3.1). Second, based on the learned embeddings, we propose a novel prototype-based label disambiguation strategy (Section 3.2). Key to our method, we gradually update the pseudo target for classification, based on the closest class prototype. By alternating the two steps above, PiCO converges to a solution with a highly distinguishable representation for accurate classification. Empirically, PiCO establishes *state-of-the-art* performance on three benchmark datasets, outperforming the baselines by a significant margin (Section 5) and obtains results that are competitive with *fully supervised learning*.

Theoretically, we demonstrate that our contrastive representation learning and prototype-based label disambiguation are mutually beneficial, and can be rigorously interpreted from an Expectation-Maximization (EM) algorithm perspective (Section 6). First, the refined pseudo labeling improves contrastive learning by selecting pseudo positive examples accurately. This can be analogous to the E-step, where we utilize the classifier’s output to assign each data example to one label-specific cluster. Second, better contrastive performance in turn improves the quality of representations and thus the effectiveness of label disambiguation. This can be reasoned from an M-step perspective, where the contrastive loss partially maximizes the likelihood by clustering similar data examples. Finally, the training data will be mapped to a mixture of von Mises-Fisher distributions on the unit hypersphere, which facilitates label disambiguation by using the component-specific label.

We have presented preliminary results of this work in [9]. While the earlier version focuses on the standard PLL setup, it holds a strong assumption that the gold labels are ensured to be included in the candidate sets. However, lacking domain knowledge, it is likely the annotators take the wrong set of labels as the candidates and dismiss the true one. Lv *et al.* [10] formalize this problem as *noisy partial label learning* (noisy PLL). But, they focus on analyzing the theoretical robustness of existing PLL methods instead of providing new solutions. Experimentally, we find that current best-performing PLL algorithms, including PiCO, display degenerated performance in the noisy PLL setup (Section 5.3), e.g. the accuracy of PRODEN drops -10.62% on CIFAR-10 with 20% wrong candidate sets. The main reason is that the label disambiguation procedure of current PLL methods is restricted to the candidate labels, which causes severe overfitting on wrong labels.

In this work, we propose PiCO+, an extension of PiCO, to tackle the noisy PLL problem. PiCO+ additionally incorporates two mechanisms. First, we present a novel distance-based clean sample detection technique that chooses near-prototype examples as clean. Second, to handle the remaining noisy examples, we develop a semi-supervised contrastive learning framework that generalizes the two PiCO modules by (i)-*contrastive learning*: construct the positive set by noisy prediction and nearest neighbor embeddings; (ii)-*label disambiguation*: guess the prototype-based soft target for noisy samples. Through extensive experiments, PiCO+ exhibits significant improvement to state-of-the-art PLL approaches on noisy PLL datasets.

Our **main contributions** are summarized as follows:

1 (Methodology). To the best of our knowledge, our paper pioneers the exploration of contrastive learning for

partial label learning and proposes a novel framework termed PiCO. As an integral part of our algorithm, we also introduce a new prototype-based label disambiguation mechanism, that leverages the contrastively learned embeddings.

2 (Practicality). Additionally, we propose PiCO+, an extension of PiCO, that targets to mitigate noisy candidate label sets. It further integrates a distance-based clean sample detection mechanism along with a semi-supervised contrastive learning module. We believe our work makes a serious attempt at improving the practicality of PLL in open-world environments.

3 (Experiments). Empirically, our proposed PiCO framework establishes the *state-of-the-art* performance on three PLL tasks. Moreover, we make the first attempt to conduct experiments on fine-grained classification datasets, where we show classification performance improvement by up to **9.61%** compared with the best baseline on the CUB-200 dataset. We also evaluate our new PiCO+ framework on the noisy variants of these datasets, where PiCO+ outperforms the best baseline by up to **12.52%** on the noisy PLL version of the CIFAR-10 dataset.

4 (Theory). We theoretically interpret our framework from the expectation-maximization perspective. Our derivation is also generalizable to other CL methods and shows the *alignment* property in CL [11] mathematically equals the M-step in center-based clustering algorithms.

2 BACKGROUND

The problem of *partial label learning* (PLL) is defined using the following setup. Let \mathcal{X} be the input space, and $\mathcal{Y} = \{1, 2, \dots, C\}$ be the output label space. We consider a training dataset $\mathcal{D} = \{(\mathbf{x}_i, Y_i)\}_{i=1}^n$, where each tuple comprises of an image $\mathbf{x}_i \in \mathcal{X}$ and a *candidate label set* $Y_i \subset \mathcal{Y}$. Identical to the supervised learning setup, the goal of PLL is to obtain a functional mapping that predicts the one true label associated with the input. Yet differently, the PLL setup bears significantly more uncertainty in the label space. A basic assumption of PLL is that the ground-truth label y_i is concealed in its candidate set, i.e., $y_i \in Y_i$, and is invisible to the learner. For this reason, the learning process can suffer from inherent ambiguity, compared with the supervised learning task with explicit ground-truth.

The key challenge of PLL is to identify the ground-truth label from the candidate label set. During training, we assign each image \mathbf{x}_i a normalized vector $\mathbf{s}_i \in [0, 1]^C$ as the *pseudo target*, whose entries denote the probability of labels being the ground-truth. The total probability mass of 1 is allocated among candidate labels in Y_i . Note that \mathbf{s}_i will be updated during the training procedure. Ideally, \mathbf{s}_i should put more probability mass on the (unknown) ground-truth label y_i over the course of training. We train a classifier $f : \mathcal{X} \rightarrow [0, 1]^C$ using cross-entropy loss, with \mathbf{s}_i being the target prediction. The per-sample loss is given by:

$$\begin{aligned} \mathcal{L}_{\text{cls}}(f; \mathbf{x}_i, Y_i) &= \sum_{j=1}^C -s_{i,j} \log(f^j(\mathbf{x}_i)) \\ \text{s.t. } \sum_{j \in Y_i} s_{i,j} &= 1 \text{ and } s_{i,j} = 0, \forall j \notin Y_i, \end{aligned} \quad (1)$$

where j denotes the indices of labels. $s_{i,j}$ denotes the j -th pseudo target of \mathbf{x}_i . Here f is the softmax output of

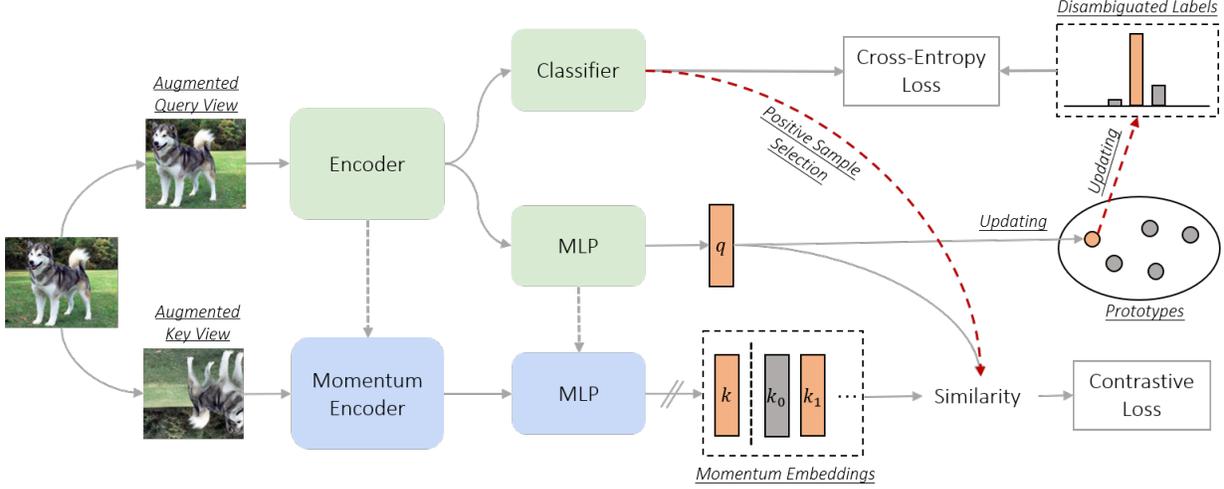


Fig. 2. Illustration of PiCO. The classifier’s output is used to determine the positive peers for contrastive learning. The contrastive prototypes are then used to gradually update the pseudo target. The momentum embeddings are maintained by a queue structure. ‘//’ means stop gradient.

the networks and we denote f^j as its j -th entry. In the remainder of this paper, we omit the sample index i when the context is clear. We proceed by describing our proposed framework.

3 THE PiCO FRAMEWORK

In this section, we describe our novel Partial label learning with **C**ontrastive label disambiguation (PiCO) framework in detail. In a nutshell, PiCO comprises two key components tackling the representation quality (Section 3.1) and label ambiguity respectively (Section 3.2). The two components systematically work as a whole and reciprocate each other. We further rigorously provide a theoretical interpretation of PiCO from an EM perspective in Section 6.

3.1 Contrastive Representation Learning for PLL

The uncertainty in the label space posits a unique obstacle for learning effective representations. In PiCO, we couple the classification loss in Eq. (1) with a contrastive term that facilitates a clustering effect in the embedding space. While contrastive learning has been extensively studied in recent literature, it remains untapped in the domain of PLL. The main challenge lies in the construction of a positive sample set. In conventional supervised CL frameworks, the positive sample pairs can be easily drawn according to the ground-truth labels [8]. However, this is not straightforward in the setting of PLL.

Training Objective. To begin with, we describe the standard contrastive loss term. We adopt the most popular setups by closely following MoCo [12] and SupCon [8]. Given each sample (x, Y) , we generate two views—a query view and a key view—by way of randomized data augmentation $\text{Aug}(x)$. The two images are then fed into a query network $g(\cdot)$ and a key network $g'(\cdot)$, yielding a pair of L_2 -normalized embeddings $q = g(\text{Aug}_q(x))$ and $k = g'(\text{Aug}_k(x))$. In implementations, the query network shares the same convolutional blocks as the classifier, followed by a prediction head (see Figure 2). Following

MoCo, the key network uses a momentum update with the query network. We additionally maintain a queue storing the most current key embeddings k , and we update the queue chronologically. To this end, we have the following contrastive embedding pool:

$$A = B_q \cup B_k \cup \text{queue}, \quad (2)$$

where B_q and B_k are vectorial embeddings corresponding to the query and key views of the current mini-batch. Given an example x , the per-sample contrastive loss is defined by contrasting its query embedding with the remainder of the pool A ,

$$\mathcal{L}_{\text{cont}}(g; x, \tau, A) = -\frac{1}{|P(x)|} \sum_{k_+ \in P(x)} \log \frac{\exp(q^\top k_+ / \tau)}{\sum_{k' \in A(x)} \exp(q^\top k' / \tau)}, \quad (3)$$

where $P(x)$ is the positive set and $A(x) = A \setminus \{q\}$. $\tau \geq 0$ is the temperature.

Positive Set Selection. As mentioned earlier, the crucial challenge is how to construct the positive set $P(x)$. We propose utilizing the predicted label $\tilde{y} = \arg \max_{j \in Y} f^j(\text{Aug}_q(x))$ from the classifier. Note that we restrict the predicted label to be in the candidate set Y . The positive examples are then selected as follows,

$$P(x) = \{k' | k' \in A(x), \tilde{y}' = \tilde{y}\}. \quad (4)$$

where \tilde{y}' is the predicted label for the corresponding training example of k' . For computational efficiency, we also maintain a label queue to store past predictions. In other words, we define the positive set of x to be those examples carrying the same *approximated* label prediction \tilde{y} . Despite its simplicity, we show that our selection strategy can be theoretically justified (Section 6) and also lead to superior empirical results (Section 5). Note that more sophisticated selection strategies can be explored, for which we discuss in Appendix B.3. Putting it all together, we jointly train the classifier as well as the contrastive network. The overall loss function is:

$$\mathcal{L}_{\text{pico}} = \mathcal{L}_{\text{cls}} + \lambda \mathcal{L}_{\text{cont}}. \quad (5)$$

Still, our goal of learning high-quality representation by CL relies on accurate classifier prediction for positive set selection, which remains unsolved in the presence of label ambiguity. To this end, we further propose a novel label disambiguation mechanism based on contrastive embeddings and show that these two components are mutually beneficial.

3.2 Prototype-based Label Disambiguation

As we mentioned (and later theoretically prove in Section 6), the contrastive loss poses a clustering effect in the embedding space. As a collaborative algorithm, we introduce our novel prototype-based label disambiguation strategy. Importantly, we keep a *prototype* embedding vector μ_c corresponding to each class $c \in \{1, 2, \dots, C\}$, which can be deemed as a set of representative embedding vectors. Categorically, a naive version of the pseudo target assignment is to find the nearest prototype of the current embedding vector. Notably this primitive resembles a clustering step. We additionally soften this hard label assignment version by using a moving-average style formula. To this end, we may posit intuitively that the employment of the prototype builds a connection with the clustering effect in the embedding space brought by the contrastive term (Section 3.1). We provide a more rigorous justification in Section 6.

Pseudo Target Updating. We propose a softened and moving-average style strategy to update the pseudo targets. Specifically, we first initialize the pseudo targets with a uniform distribution, $s_j = \frac{1}{|Y|} \mathbb{I}(j \in Y)$. We then iteratively update it by the following moving-average mechanism,

$$s = \phi s + (1 - \phi)z, \quad z_c = \begin{cases} 1 & \text{if } c = \arg \max_{j \in Y} \mathbf{q}^\top \mu_j, \\ 0 & \text{else} \end{cases} \quad (6)$$

where $\phi \in (0, 1)$ is a positive constant, and μ_j is a *prototype* corresponding to the j -th class. The intuition is that fitting uniform pseudo targets results in a good initialization for the classifier since the contrastive embeddings are less distinguishable at the beginning. The moving-average style strategy then smoothly updates the pseudo targets towards the correct ones, and meanwhile ensures stable dynamics of training; see Appendix B.1. With more rigorous validation provided later in Section 6, we provide an explanation for the prototype as follows: (i)-for a given input x , the closest prototype is indicative of its ground-truth class label. At each step, s has the tendency to slightly move toward the one-hot distribution defined by z based on Eq. (6); (ii)-if an example consistently points to one prototype, the pseudo target s can converge (almost) to a one-hot vector with the least ambiguity.

Prototype Updating. The most canonical way to update the prototype embeddings is to compute it in every iteration of training. However, this would extract a heavy computational toll and in turn cause unbearable training latency. As a result, we update the class-conditional prototype vector similarly in a moving-average style:

$$\mu_c = \text{Normalize}(\gamma \mu_c + (1 - \gamma) \mathbf{q}), \quad (7)$$

$$\text{if } c = \arg \max_{j \in Y} f^j(\text{Aug}_q(\mathbf{x})),$$

where the momentum prototype μ_c of class c is defined by the moving-average of the normalized query embeddings

Algorithm 1: Pseudo-code of PiCO (one epoch).

```

1 Input: Training dataset  $\mathcal{D}$ , classifier  $f$ , query network  $g$ ,
   key network  $g'$ , momentum queue, uniform
   pseudo-labels  $s_i$  associated with  $\mathbf{x}_i$  in  $\mathcal{D}$ , class
   prototypes  $\mu_j$  ( $1 \leq j \leq C$ ).
2 for  $iter = 1, 2, \dots$ , do
3   sample a mini-batch  $B$  from  $\mathcal{D}$ 
4   // query and key embeddings generation
5    $B_q = \{\mathbf{q}_i = g(\text{Aug}_q(\mathbf{x}_i)) | \mathbf{x}_i \in B\}$ 
6    $B_k = \{\mathbf{k}_i = g'(\text{Aug}_k(\mathbf{x}_i)) | \mathbf{x}_i \in B\}$ 
7    $A = B_q \cup B_k \cup \text{queue}$ 
8   for  $\mathbf{x}_i \in B$  do
9     // classifier prediction
10     $\tilde{y}_i = \arg \max_{j \in Y_i} f^j(\text{Aug}_q(\mathbf{x}_i))$ 
11    // momentum prototype updating
12     $\mu_c = \text{Normalize}(\gamma \mu_c + (1 - \gamma) \mathbf{q}_i)$ , if  $\tilde{y}_i = c$ 
13    // positive set generation
14     $P(\mathbf{x}_i) = \{\mathbf{k}' | \mathbf{k}' \in A(\mathbf{x}_i), \tilde{y}' = \tilde{y}_i\}$ 
15  end
16  // prototype-based label disambiguation
17  for  $\mathbf{q}_i \in B_q$  do
18     $\mathbf{z}_i = \text{OneHot}(\arg \max_{j \in Y_i} \mathbf{q}_i^\top \mu_j)$ 
19     $\mathbf{s}_i = \phi \mathbf{s}_i + (1 - \phi) \mathbf{z}_i$ 
20  end
21  // network updating
22  minimize loss  $\mathcal{L}_{\text{pico}} = \mathcal{L}_{\text{cls}} + \lambda \mathcal{L}_{\text{cont}}$ 
23  // update the key network and momentum
24  queue
25  momentum update  $g'$  by using  $g$ 
26  enqueue  $B_k$  and classifier predictions and dequeue
27 end

```

\mathbf{q} whose predicted class conforms to c . γ is a tunable hyperparameter.

3.3 Synergy between Contrastive Learning and Label Disambiguation

While seemingly separated from each other, the two key components of PiCO work in a collaborative fashion. First, as the contrastive term favorably manifests a clustering effect in the embedding space, the label disambiguation module further leverages via setting more precise prototypes. Second, a set of well-polished label disambiguation results may, in turn, reciprocate the positive set construction which serves as a crucial part in the contrastive learning stage. The entire training process converges when the two components perform satisfactorily. We further rigorously draw a resemblance of PiCO with a classical EM-style clustering algorithm in Section 6. Our experiments, particularly the ablation study displayed in Section 5.2.2, further justify the mutual dependency of the synergy between the two components. The pseudo-code of our complete algorithm is shown in Algorithm 1.

4 PICO+ FOR NOISY PARTIAL LABELS

In this section, we investigate a more practical setup called *noisy partial label learning* [10], where the true label potentially lies outside the candidate set, i.e., $y_i \notin Y_i$. Empirically, we observe that the current PLL methods, including PiCO, exhibit a significant performance drop on noisy PLL datasets (Section 5.3). One main reason is that these methods

mostly rely on a within-set pseudo-label updating step, but the noisy candidate sets in noisy PLL can mislead them towards overfitting on a wrong label.

To this end, we propose PiCO+, an extension of PiCO, which learns robust classifiers from noisy partial labels. First, we introduce a distance-based sample selection mechanism that detects *clean examples* whose candidate sets contain the ground-truth labels. Then, we develop a semi-supervised contrastive PLL framework to handle data with noisy candidates. In what follows, we elaborate on our novel PiCO+ framework in detail.

4.1 Distance-based Clean Example Detection

To remedy overfitting on noisy candidates, we would like to select those reliable candidates, which contain the true labels, to run the PiCO method. In the noisy label learning regime, a widely-adopted strategy is to employ the small-loss selection criterion [13], which is based on the observation that noisy examples typically demonstrate a large loss. However, in the noisy PLL problem, the loss values are less informative since examples with more candidate labels can also incur a larger loss, even if the candidate sets are reliable.

To address this problem, we propose a distance-based selection mechanism as follows,

$$\mathcal{D}_{\text{clean}} = \{(\mathbf{x}_i, Y_i) | \mathbf{q}_i^\top \boldsymbol{\mu}_{\tilde{y}_i} > \kappa_\delta\}, \quad (8)$$

where $\tilde{y}_i = \arg \max_{j \in Y_i} f^j(\text{Aug}_q(\mathbf{x}_i))$ is the classifier prediction. κ_δ is the $(100 - \delta)$ -percentile of the cosine similarity between the query embedding and the \tilde{y}_i -th prototype. For example, when $\delta = 60$, it means 60% of examples are above the threshold. Our motivation is that the clustering effect of contrastive learning makes the clean examples dominate the prototype calculation and thus they are distributed close to at least one prototype inside the candidate set. On the other hand, the noisy candidates mostly deviate from all candidate prototypes as their true labels are not contained in their candidate sets.

4.2 Semi-supervised Contrastive Learning

While leveraging the clean examples to run a PiCO model is straightforward, the low data utility restricts it from better performance. Consequently, we regard noisy examples as unlabeled ones and develop a semi-supervised learning framework to learn from the two sets. On clean datasets, we assume the true labels are included in the candidate and run our PiCO method. It should be particularly noted that we also restrict the positive set construction and prototype updating procedures to merely employ clean examples. But, the pseudo-target updating is performed over all data points.

On the noisy dataset, which is denoted by $\mathcal{D}_{\text{noisy}} = \mathcal{D} \setminus \mathcal{D}_{\text{clean}}$, we follow our design patterns in PiCO to synergistically train the contrastive branch as well as the classifier. It is achieved by the following components:

Neighbor-Augmented Contrastive Learning. Recall that the crucial step for designing contrastive loss is to construct the positive set, which is challenging for noisy samples as their candidate sets are unreliable. To this end, we first propose a label-driven construction approach. By

regarding noisy samples as unlabeled data, it is intuitive to treat all labels as candidates. This gives rise to the following noisy positive set,

$$P_{\text{noisy}}(\mathbf{x}) = \{\mathbf{k}' | \mathbf{k}' \in A(\mathbf{x}), \hat{y}' = \hat{y}\},$$

$$\text{where } \hat{y} = \begin{cases} \arg \max_{1 \leq j \leq L} f^j(\text{Aug}_q(\mathbf{x})) & \text{if } \mathbf{x} \in \mathcal{D}_{\text{noisy}}, \\ \arg \max_{j \in Y} f^j(\text{Aug}_q(\mathbf{x})) & \text{else.} \end{cases} \quad (9)$$

That is, we choose the within-set classifier prediction for clean examples and choose the full-label prediction for the remaining. We apply this noisy positive set to all data in \mathcal{D} and calculate a noisy contrastive loss $\mathcal{L}_{\text{n-cont}}$ by Eq. (3). This objective serves as our ultimate goal of semi-supervised training that recovers the PiCO algorithm on clean PLL data and noisy unlabeled data.

Nevertheless, as we are less knowledgeable of noisy examples, our estimation on $P_{\text{noisy}}(\mathbf{x})$ can be inaccurate and assigns wrong cluster centers. To this end, we incorporate a data-driven technique to regularize the noisy contrastive loss. In specific, we collect the nearest neighbors of noisy samples to be positive peers,

$$P_{\text{knn}}(\mathbf{x}) = \{\mathbf{k}' | \mathbf{k}' \in A(\mathbf{x}) \cap \mathcal{N}_k(\mathbf{x})\}, \quad (10)$$

where $\mathcal{N}_k(\mathbf{x})$ is the embedding set of \mathbf{x} 's k -nearest neighbors in the embedding space. Then, we calculate the k NN-based contrastive loss \mathcal{L}_{knn} on noisy examples. Note that the original CL objective naturally encourages the examples to be locally smooth by aligning augmented copies. Our neighbor-augmented loss further enhances this effect to ensure the examples in a local region share the same label. By then, the noisy examples are aligned to their local neighbors, which promotes their clustering effect towards the right labels.

Prototype-based Label Guessing. Similarly, we would also like to identify the true labels on noisy examples to enhance the classifier training. Although the noisy examples are treated as unlabeled, it is not appropriate to directly set their labels to uniform labels $\frac{1}{L}$ as in PiCO, since the data separation procedure dynamically changes during training. To this end, we leverage the class prototypes to guess their pseudo-targets s' by,

$$s'_j = \frac{\exp(\mathbf{q}^\top \boldsymbol{\mu}_j / \tau)}{\sum_{t=1}^L \exp(\mathbf{q}^\top \boldsymbol{\mu}_t / \tau)}, \quad \forall 1 \leq j \leq L. \quad (11)$$

We calculate the cross-entropy loss on $\mathcal{D}_{\text{noisy}}$ as defined in Eq. (1), which is term as $\mathcal{L}_{\text{n-cls}}$.

Mixup Training. Recently, the mixup regularization technique has widely adopted to improve the robustness of weakly-supervised learning algorithms [14], [15]. Therefore, we also incorporate it into PiCO+ for boosted performance. Formally, given a pair of images \mathbf{x}_i and \mathbf{x}_j , we create a virtual training example by linearly interpolating both,

$$\begin{aligned} \mathbf{x}^m &= \sigma \text{Aug}_q(\mathbf{x}_i) + (1 - \sigma) \text{Aug}_q(\mathbf{x}_j), \\ \mathbf{s}^m &= \sigma \hat{\mathbf{s}}_i + (1 - \sigma) \hat{\mathbf{s}}_j, \end{aligned} \quad (12)$$

where $\sigma \sim \text{Beta}(\zeta, \zeta)$ and ζ is a hyperparameter. Here, we take the pseudo-target of PiCO on clean examples, and the guessed label on noisy examples, i.e., $\hat{\mathbf{s}} = \mathbf{s}$ if $\mathbf{x} \in \mathcal{D}_{\text{clean}}$ else $\hat{\mathbf{s}} = \mathbf{s}'$. We define the mixup loss \mathcal{L}_{mix} as the cross-entropy on \mathbf{x}^m and \mathbf{s}^m .

TABLE 1

Accuracy comparisons on standard PLL datasets. Bold indicates superior results. Notably, PiCO+ achieves comparable results to the fully supervised learning (less than 1% in accuracy with ≈ 1 false candidate).

Dataset	Method	$q = 0.1$	$q = 0.3$	$q = 0.5$
CIFAR-10	PiCO+ (ours)	95.99 \pm 0.03%	95.73 \pm 0.10%	95.33 \pm 0.06%
	PiCO (ours)	94.39 \pm 0.18%	94.18 \pm 0.12%	93.58 \pm 0.06%
	LWS	90.30 \pm 0.60%	88.99 \pm 1.43%	86.16 \pm 0.85%
	PRODEN	90.24 \pm 0.32%	89.38 \pm 0.31%	87.78 \pm 0.07%
	CC	82.30 \pm 0.21%	79.08 \pm 0.07%	74.05 \pm 0.35%
	MSE	79.97 \pm 0.45%	75.64 \pm 0.28%	67.09 \pm 0.66%
	EXP	79.23 \pm 0.10%	75.79 \pm 0.21%	70.34 \pm 1.32%
Dataset	Method	$q = 0.01$	$q = 0.05$	$q = 0.1$
CIFAR-100	PiCO+ (ours)	76.29 \pm 0.42%	76.17 \pm 0.18%	75.55 \pm 0.21%
	PiCO (ours)	73.09 \pm 0.34%	72.74 \pm 0.30%	69.91 \pm 0.24%
	LWS	65.78 \pm 0.02%	59.56 \pm 0.33%	53.53 \pm 0.08%
	PRODEN	62.60 \pm 0.02%	60.73 \pm 0.03%	56.80 \pm 0.29%
	CC	49.76 \pm 0.45%	47.62 \pm 0.08%	35.72 \pm 0.47%
	MSE	49.17 \pm 0.05%	46.02 \pm 1.82%	43.81 \pm 0.49%
	EXP	44.45 \pm 1.50%	41.05 \pm 1.40%	29.27 \pm 2.81%

Finally, we aggregate the above losses together,

$$\mathcal{L}_{\text{pico+}} = \mathcal{L}_{\text{mix}} + \alpha \mathcal{L}_{\text{clean}} + \beta (\mathcal{L}_{\text{n-cont}} + \mathcal{L}_{\text{knn}} + \mathcal{L}_{\text{n-cls}}), \quad (13)$$

where $\mathcal{L}_{\text{clean}}$ is the PiCO loss on clean examples. Note that over-trusting the guessed labels and positive sets on noisy samples may cause confirmation bias [16] and makes the model overfit on wrong labels. Empirically, we set a larger α and a smaller β (e.g. $\alpha = 2, \beta = 0.1$), and thus, the clean samples dominate the learning procedure to ensure favorable noisy example detection ability.

5 EXPERIMENTS

In this section, we present our main results and part of ablation results to show the effectiveness of PiCO and PiCO+ methods. We refer readers to Appendix B for more experimental results and analysis. Code and data are available at: <https://github.com/hbzju/PiCO>.

5.1 Setup

Datasets. First, we evaluate PiCO on two commonly used benchmarks — CIFAR-10 and CIFAR-100 [17]. Adopting the identical experimental settings in previous work [18], [19], we generate conventional partially labeled datasets by flipping negative labels $\bar{y} \neq y$ to false positive labels with a probability $q = P(\bar{y} \in Y | \bar{y} \neq y)$. In other words, all $C - 1$ negative labels have a uniform probability to be false positive and we aggregate the flipped ones with the ground-truth to form the candidate set. We consider $q \in \{0.1, 0.3, 0.5\}$ for CIFAR-10 and $q \in \{0.01, 0.05, 0.1\}$ for CIFAR-100. In Section 5.4, we further evaluate our method on fine-grained classification tasks, where label disambiguation can be more challenging.

For the noisy PLL task, we introduce a noise controlling parameter $\eta = 1 - P(y \in Y | y)$ that controls the probability of the ground-truth label not being selected as a candidate. As it is possible that some instances are not assigned any

candidate and we simply re-generate the candidate set until it has at least one candidate label. We select $\eta \in \{0.1, 0.2\}$ for both datasets; results with stronger noisy ratios are shown in Section 5.3.3.

Baselines. We choose the five best-performed partial label learning algorithms to date: 1) LWS [19] weights the risk function by means of a trade-off between losses on candidate labels and the remaining; 2) PRODEN [18] iteratively updates the latent label distribution in a self-training style; 3) CC [20] is a classifier-consistent method that assumes set-level uniform data generation process; 4) MSE and EXP [21] are two simple baselines that adopt mean square error and exponential loss as the risks. For the noisy PLL task, we additionally incorporate one baseline method GCE [10], [22] that generalizes the cross-entropy loss via negative Box-Cox transformation. The hyperparameters are tuned according to the original methods. For all experiments, we report the mean and standard deviation based on 5 independent runs (with different random seeds).

Implementation Details. Following the standard experimental setup in PLL [19], [20], we split a clean validation set (10% of training data) from the training set to select the hyperparameters. After that, we transform the validation set back to its original PLL form and incorporate it into the training set to accomplish the final model training. We use an 18-layer ResNet as the backbone for feature extraction. Most of experimental setups for the contrastive network follow previous works [8], [12]. The projection head of the contrastive network is a 2-layer MLP that outputs 128-dimensional embeddings. We use two data augmentation modules SimAugment [8] and RandAugment [23] for query and key data augmentation respectively. Empirically, we find that even weak augmentation for key embeddings also leads to good results. The size of the queue that stores key embeddings is fixed to be 8192. The momentum coefficients are set as 0.999 for contrastive network updating and $\gamma = 0.99$ for prototype calculation. For pseudo target

TABLE 2
Ablation study of PiCO on standard partial label learning datasets CIFAR-10 ($q = 0.5$) and CIFAR-100 ($q = 0.05$).

Ablation	$\mathcal{L}_{\text{cont}}$	Label Disambiguation	CIFAR-10	CIFAR-100
PiCO	✓	Ours	93.58	72.74
PiCO w/o Disambiguation	✓	Uniform Pseudo Target	84.50	64.11
PiCO w/o $\mathcal{L}_{\text{cont}}$	✗	Uniform Pseudo Target	76.46	56.87
PiCO with $\phi = 0$	✓	Soft Prototype Probs	91.60	71.07
PiCO with $\phi = 0$	✓	One-hot Prototype	91.41	70.10
PiCO	✓	MA Soft Prototype Probs	81.67	63.75
Fully-Supervised	✓	N/A	94.91	73.56

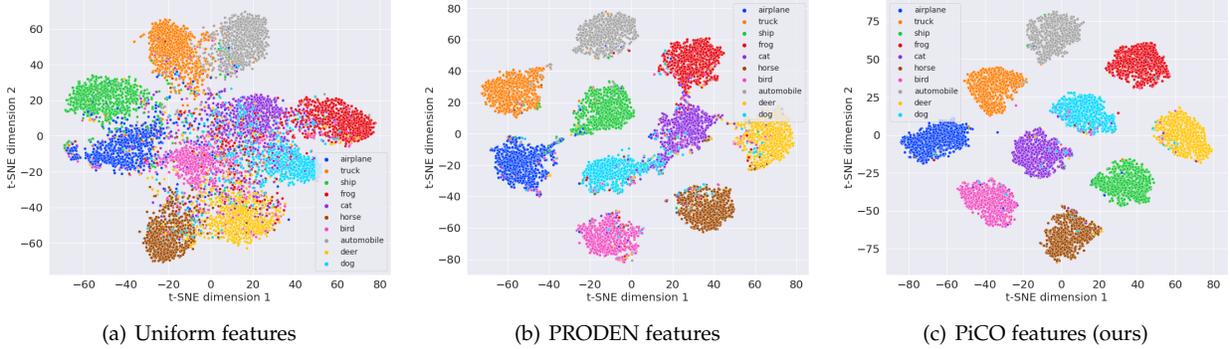


Fig. 3. T-SNE visualization of the image representation on CIFAR-10 ($q = 0.5$). Different colors represent the corresponding classes.

updating, we linearly ramp down ϕ from 0.95 to 0.8. The temperature parameter is set as $\tau = 0.07$. The loss weighting factor is set as $\lambda = 0.5$. The model is trained by a standard SGD optimizer with a momentum of 0.9 and the batch size is 256. We train the model for 800 epochs with cosine learning rate scheduling. We also empirically find that classifier warm up leads to better performance when there are many candidates. Hence, we disable contrastive learning in the first 100 epoch for CIFAR-100 with $q = 0.1$ and 1 epoch for the remaining experiments.

For PiCO+, we basically follow the original PiCO method. The clean sample selection ratio parameter δ is set as 0.8/0.6 for noisy ratio 0.1/0.2, respectively. For neighbor augmentation, we set $k = 16$ for CIFAR-10 and a smaller $k = 5$ for CIFAR-100. In the beginning, the embeddings can be less meaningful and thus, we enable k NN augmentation after the first 100 epochs. We fix the shape parameter of the Beta distribution to $\varsigma = 4$ for mixup training. We set the loss weighting factor $\alpha = 2$ and $\beta = 0.1$. Similar to the standard PLL setup, we warm up the model by fitting uniform targets for 5 and 50 epochs on CIFAR-10 and CIFAR-100 datasets respectively.

5.2 Experimental Results on standard PLL

5.2.1 Main Results

PiCO achieves SOTA results on standard PLL task. As shown in Table 1, PiCO significantly outperforms all the rivals by a significant margin on all datasets. Specifically, on CIFAR-10 dataset, we improve upon the best baseline by **4.09%**, **4.80%**, and **5.80%** where q is set to 0.1, 0.3, 0.5 respectively. Moreover, PiCO consistently achieves superior results as the size of the candidate set increases, while the baselines demonstrate a significant performance drop.

Besides, it is worth pointing out that previous works [18], [19] are typically evaluated on datasets with a small label space ($C = 10$). We challenge this by showing additional results on CIFAR-100. When $q = 0.1$, all the baselines fail to obtain satisfactory performance, whereas PiCO remains competitive. Moreover, we observe that PiCO achieves results that are comparable to the *fully supervised contrastive learning model* (in Table 2), showing that disambiguation is sufficiently accomplished. The comparison highlights the superiority of our label disambiguation strategy. Lastly, we evaluated the PiCO+ method in the standard PLL setup, which equals a PiCO model with mixup training. It can be shown that PiCO+ further improves PiCO by a notable margin, validating the robustness of the PiCO+ method.

PiCO learns more distinguishable representations. We visualize the image representation produced by the feature encoder using t-SNE [24] in Figure 3. Different colors represent different ground-truth class labels. We use the CIFAR-10 dataset with $q = 0.5$. We contrast the t-SNE embeddings of three approaches: (a) a model trained with uniform pseudo targets, i.e., $s_j = 1/|Y|$ ($j \in Y$), (b) the best baseline PRODEN, and (c) our method PiCO. We can observe that the representation of the uniform model is indistinguishable since its supervision signals suffer from high uncertainty. The features of PRODEN are improved, yet with some class overlapping (e.g., blue and purple). In contrast, PiCO produces well-separated clusters and more distinguishable representations, which validates its effectiveness in learning high-quality representation.

5.2.2 Ablation Studies of PiCO

Effect of $\mathcal{L}_{\text{cont}}$ and label disambiguation. We ablate the contributions of two key components of PiCO: contrastive

TABLE 3
Accuracy comparisons on noisy PLL datasets. Bold indicates superior results.

Dataset	Method	$q = 0.3$		$q = 0.5$	
		$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$
CIFAR-10	PiCO+ (ours)	95.11 ± 0.13%	93.98 ± 0.39%	94.45 ± 0.27%	92.59 ± 0.22%
	PiCO (ours)	89.47 ± 0.37%	84.13 ± 0.53%	87.79 ± 0.09%	80.07 ± 0.60%
	LWS	84.51 ± 0.13%	77.98 ± 0.09%	71.02 ± 7.27%	61.96 ± 3.22%
	PRODEN	84.56 ± 0.16%	79.35 ± 0.12%	81.97 ± 0.59%	77.15 ± 0.08%
	CC	72.16 ± 0.93%	68.42 ± 0.37%	65.61 ± 0.43%	51.82 ± 4.13%
	MSE	53.77 ± 1.44%	49.73 ± 2.94%	46.56 ± 0.18%	39.80 ± 2.82%
	EXP	75.81 ± 0.09%	69.97 ± 0.39%	64.26 ± 1.02%	54.93 ± 1.11%
	GCE	74.32 ± 1.04%	69.90 ± 0.80%	70.38 ± 7.63%	50.57 ± 0.95%
Dataset	Method	$q = 0.05$		$q = 0.1$	
		$\eta = 0.1$	$\eta = 0.2$	$\eta = 0.1$	$\eta = 0.2$
CIFAR-100	PiCO+ (ours)	74.68 ± 0.15%	72.98 ± 0.22%	67.58 ± 1.05%	62.24 ± 0.97%
	PiCO (ours)	66.29 ± 0.10%	59.81 ± 0.24%	66.15 ± 0.03%	45.32 ± 0.89%
	LWS	52.20 ± 1.47%	42.31 ± 1.05%	20.54 ± 4.77%	17.76 ± 4.47%
	PRODEN	53.40 ± 0.61%	46.11 ± 0.38%	47.34 ± 1.39%	38.03 ± 1.79%
	CC	42.06 ± 0.67%	37.90 ± 3.27%	32.11 ± 3.95%	22.28 ± 6.18%
	MSE	31.06 ± 2.46%	27.36 ± 0.40%	25.86 ± 1.87%	22.98 ± 1.74%
	EXP	23.98 ± 4.25%	22.37 ± 5.45%	23.78 ± 4.59%	22.27 ± 3.38%
	GCE	35.85 ± 1.37%	31.65 ± 0.71%	27.79 ± 2.80%	24.21 ± 1.67%

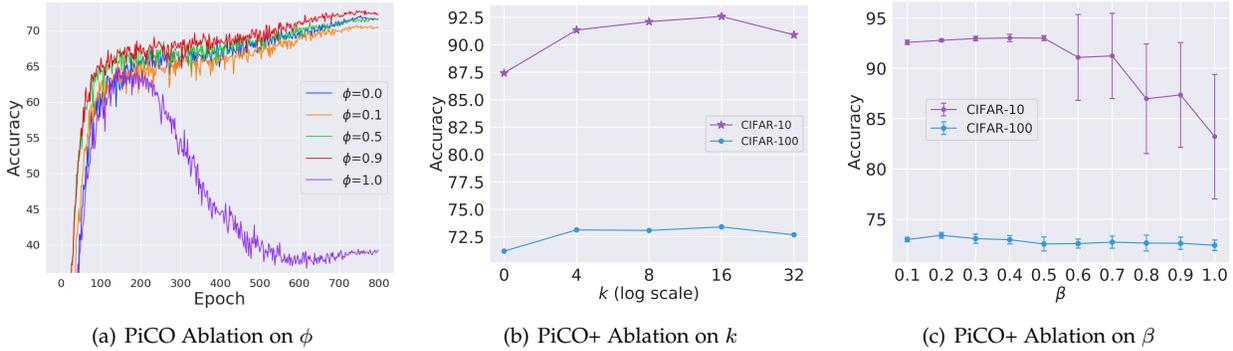


Fig. 4. (a) Performance of PiCO with varying moving-average factor ϕ on CIFAR-100 ($q = 0.05$). (c) Performance of PiCO+ with varying neighbor number k on CIFAR-10 ($q = 0.5, \eta = 0.2$) and CIFAR-100 ($q = 0.05, \eta = 0.2$). (b) Performance of PiCO+ with varying semi-supervised loss weighting factor β on CIFAR-10 ($q = 0.5, \eta = 0.2$).

learning and prototype-based label disambiguation. In particular, we compare PiCO with two variants: 1) *PiCO w/o disambiguation* which keeps the pseudo target as uniform $1/|Y|$; and 2) *PiCO w/o \mathcal{L}_{cont}* which further removes the contrastive learning and only trains a classifier with uniform pseudo targets. From Table 2, we can observe that variant 1 substantially outperforms variant 2 (e.g., +8.04% on CIFAR-10), which signifies the importance of contrastive learning for producing better representations. Moreover, with label disambiguation, PiCO obtains results close to *fully supervised setting*, which verifies the ability of PiCO in identifying the ground-truth.

Different disambiguation strategy. Based on the contrastive prototypes, various strategies can be used to disambiguate the labels, which corresponds to the E-step in our theoretical analysis. We choose the following variants:

- 1) *One-hot Prototype* always assigns a one-hot pseudo target $s = z$ by using the nearest prototype ($\phi = 0$);
- 2) *Soft Prototype Probs* follows [25] and uses a soft class probability $s_i = \frac{\exp(\mathbf{q}^T \mu_i / \tau)}{\sum_{j \in Y} \exp(\mathbf{q}^T \mu_j / \tau)}$ as the pseudo target ($\phi = 0$);
- 3) *MA Soft Prototype Probs* gradually updates pseudo target from uniform by using the soft probabilities in a moving-average style. From Table 2, we can see that directly using either soft or hard prototype-based label assignment leads to competitive results. This corroborates our theoretical analysis in Section 6, since center-based class probability estimation is common in clustering algorithms. However, *MA Soft Prototype Probs* displays degenerated performance, suggesting soft label assignment is less reliable in identifying the ground-truth. Finally, PiCO outperforms the best variant by $\approx 2\%$ in accuracy on both datasets, showing the superiority of our label disambiguation strategy.

TABLE 4

Ablation study of PiCO on noisy partial label learning datasets CIFAR-10 ($q = 0.5, \eta = 0.2$) and CIFAR-100 ($q = 0.05, \eta = 0.2$).

Ablation	\mathcal{L}_{n-cls}	\mathcal{L}_{n-cont}	kNN	Mixup	CIFAR-10	CIFAR-100
PiCO+	✓	✓	✓	All Data	92.59	72.98
PiCO+ with Small Loss	✓	✓	✓	All Data	91.22	72.54
PiCO+ with Only Clean	✗	✗	✗	Only Clean	87.63	70.72
PiCO+ w/o \mathcal{L}_{n-cls}	✗	✓	✓	All Data	91.14	71.98
PiCO+ w/o \mathcal{L}_{n-cont}	✓	✗	✓	All Data	90.87	72.89
PiCO+ w/o kNN	✓	✓	✗	All Data	87.43	71.19
PiCO+ w/o Mixup	✓	✓	✓	No Mixup	89.67	68.51
Fully-Supervised [†]	-	-	-	All Data	95.96 [‡]	76.36

[†] These supervised results are evaluated with mixup like PiCO+ and thus are different from Table 2.

[‡] It is slightly smaller than PiCO+ in Table 1 ($q = 0.1$) because of randomized running, but they have no statistically significant difference.

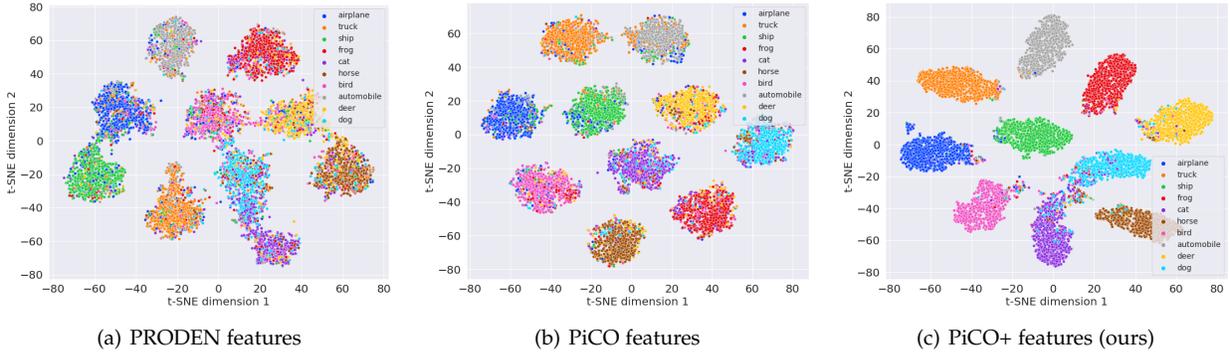


Fig. 5. T-SNE visualization of the image representation on noisy PLL version of CIFAR-10 ($q = 0.5, \eta = 0.2$). Different colors represent the corresponding classes.

Effect of moving-average factor ϕ . We then explore the effect of pseudo target updating factor ϕ on PiCO performance. Figure 4 (a) shows the learning curves of PiCO on CIFAR-100 ($q = 0.05$). We can see that the best result is achieved at $\phi = 0.9$ and the performance drops when ϕ takes a smaller value, particularly on the early stage. When $\phi = 0$, PiCO obtains a competitive result but is much lower than $\phi = 0.9$. This confirms that trusting the uniform pseudo targets at the early stage is crucial in obtaining superior performance. At the other extreme value $\phi = 1$, uniform pseudo targets are used, and PiCO demonstrates a degenerated performance and severe overfitting phenomena. In general, PiCO performs well when $\phi \approx 0.9$.

5.3 Main Empirical Results on Noisy PLL

5.3.1 Main Results

PiCO+ achieves SOTA results on noisy PLL task. In Table 3, we compare PiCO+ with competitive PLL methods on CIFAR datasets, where PiCO+ significantly outperforms baselines. In specific, on CIFAR-10 dataset with $q = 0.5$, PiCO+ improves upon the best competitor by **6.66%** and **12.52%** when η is set to 0.1, 0.2 respectively. Notably, under the noisy PLL setup, even when only 10% examples have wrong candidate sets, the baseline algorithms (including PiCO) exhibit severe performance degradation. This is further aggravated on CIFAR-100 with a larger label space, while PiCO+ consistently retains its great robustness.

PiCO+ learns compact and distinguishable features. In Figure 5, we visualize the feature representations of PiCO+ on the noisy PLL CIFAR-10 dataset with $q = 0.5, \eta = 0.2$. It can be shown that PiCO generates compact representations even with noisy candidate sets, which further supports the clustering effect of contrastive learning. However, both PiCO and PRODEN exhibit severe overfitting on wrong labels. Instead, the features of our PiCO+ is both compact and distinguishable.

5.3.2 Ablation Studies of PiCO+

Effect of sample selection. We first study the effectiveness of our distance-based selection mechanism by comparing PiCO+ with two variants: 1) *PiCO+ with Small Loss* selects clean examples by sorting cross-entropy losses; 2) *PiCO+ with Only Clean* employs the clean samples to run a simple PiCO method. As reported in Table 4, PiCO+ with only clean data exhibits better performance than the vanilla PiCO method (e.g. +7.56%) on CIFAR-10, indicates the selected examples enjoy high purity. Moreover, PiCO+ with Small Loss underperforms PiCO+, which verifies that the loss values are indeed less informative in the presence of candidate labels and that our strategy is a better alternative.

Effect of semi-supervised contrastive training. Next, we explore the effect of each component in SSL training. We compare PiCO+ with three variants: 1) *PiCO+ w/o \mathcal{L}_{n-cls}* removes the label guessing technique; 2) *PiCO+ w/o \mathcal{L}_{n-cont}* removes the label-driven contrastive loss; 3) *PiCO+ w/o kNN* disables the neighbor-augmented contrastive loss;

TABLE 5
Accuracy comparisons on noisy PLL datasets with more noisy samples. Bold indicates superior results.

Method	CIFAR-10 ($q = 0.5$)		CIFAR-100 ($q = 0.05$)		
	$\eta = 0.3$	$\eta = 0.4$	$\eta = 0.3$	$\eta = 0.4$	$\eta = 0.5$
PiCO+ (ours)	90.12 \pm 0.51%	76.09 \pm 3.62%	70.46 \pm 0.51%	66.41 \pm 0.58%	60.50 \pm 0.99%
PiCO (ours)	64.79 \pm 2.08%	34.59 \pm 7.26%	52.18 \pm 0.52%	44.17 \pm 0.08%	35.51 \pm 1.14%
PRODEN	50.32 \pm 1.07%	29.68 \pm 13.29%	39.19 \pm 0.20%	33.64 \pm 0.82%	26.91 \pm 0.83%

TABLE 6
Accuracy comparisons on fine-grained classification datasets with standard PLL labels.

Method	CUB-200 ($q = 0.05$)	CIFAR-100-H ($q = 0.5$)
PiCO+	72.05 \pm 0.80%	75.38 \pm 0.52%
PiCO	72.17 \pm 0.72%	72.04 \pm 0.31%
LWS	39.74 \pm 0.47%	57.25 \pm 0.02%
PRODEN	62.56 \pm 0.10%	60.89 \pm 0.03%
CC	55.61 \pm 0.02%	42.60 \pm 0.11%
MSE	22.07 \pm 2.36%	39.52 \pm 0.28%
EXP	9.44 \pm 2.32%	35.08 \pm 1.71%

TABLE 7
Accuracy comparisons on fine-grained classification datasets with noisy PLL labels.

Method	CUB-200 ($q = 0.05, \eta = 0.2$)	CIFAR-100-H ($q = 0.5, \eta = 0.2$)
PiCO+	60.65 \pm 0.79%	68.31 \pm 0.47%
PiCO	53.05 \pm 2.03%	59.81 \pm 0.25%
LWS	18.65 \pm 2.15%	22.18 \pm 6.12%
PRODEN	44.74 \pm 2.47%	48.03 \pm 0.47%
CC	26.98 \pm 1.16%	34.57 \pm 0.99%
MSE	20.92 \pm 1.20%	35.20 \pm 1.03%
EXP	2.81 \pm 14.46%	20.80 \pm 4.62%
GCE	5.13 \pm 38.65%	33.21 \pm 2.03%

4) *PiCO+* w/o *Mixup* disables the mixup training. From Table 4, we can see that all the components of our SSL framework contribute to the performance improvements. In particular, \mathcal{L}_{n-cont} has a stronger positive effect on CIFAR-10 and the label guessing component brings extra performance improvements for both datasets. The mixup and neighbor augmentation are the most crucial to the final performance. We note that the mixup training technique also improves the fully-supervised model (e.g. +1.05% on CIFAR-10). But the performance improvement is more substantial on the noisy PLL tasks (e.g. +2.92% on CIFAR-10), indicating its robustness on noisy data. Lastly, Figure 4 (b) shows the influence of neighbor number k , where *PiCO+* works well in a wide range of k values. Nevertheless, a too large k may collect many noisy positive peers and slightly drops the performance.

Effect of loss weighting factor β . Figure 4 reports the performance of *PiCO+* with varying β values. On the CIFAR-10 dataset, we observe a severe performance degradation with β being larger. The variance becomes increasingly larger as well. Similar trends can also be observed on CIFAR-100, though the results are much stabler. It suggests that the usage of noisy examples should be careful as they may result in confirmation bias.

5.3.3 The Robustness of *PiCO+* with Severe Noise

Finally, we conduct experiments on noisy PLL datasets that contain much more severe noise to show the robustness of our *PiCO+* method. In particular, we choose $\eta \in \{0.3, 0.4\}$ and $\eta \in \{0.3, 0.4, 0.5\}$ for CIFAR-10 and CIFAR-100 respectively. Accordingly, we adjust the selection ratio to $\delta = 0.5, 0.4$ when $\eta = 0.4, 0.5$, without changing other setups. Table 5 compares *PiCO+* with two most competitive baselines *PiCO* and *PRODEN*, where *PiCO+* obtains very

impressive performance. For example, the gaps between *PiCO+* and the best baseline are **41.50%** and **24.99%** on CIFAR-10 with $\eta = 0.4$ and CIFAR-100 with $\eta = 0.5$. We conclude that *PiCO+* is indeed much more robust than existing PLL algorithms.

5.4 Further Extension: Fine-Grained Partial Label Learning

Recall the dog example highlighted in Section 2, where semantically similar classes are more likely to cause label ambiguity. It begs the question of whether *PiCO* is effective on the challenging fine-grained image classification tasks. To verify this, we conduct experiments on two datasets: 1) CUB-200 dataset [26] contains 200 bird species; 2) CIFAR-100 with hierarchical labels (CIFAR-100-H), where we generate candidate labels that belong to the same superclass¹. We set $q = 0.05$ for CUB-200 and $q = 0.5$ for CIFAR-100 with hierarchical labels. In Table 6, we compare *PiCO* with baselines, where *PiCO* outperforms the best method *PRODEN* by a large margin (+9.61% on CUB-200 and +11.15% on CIFAR-100-H). In addition, we test the performance of *PiCO+* on both standard and noisy PLL versions of fine-grained datasets. For the noisy version, we set the number of neighbors by $k = 3$ for CUB-200 and $\eta = 0.2$ for both. The results are listed in Table 6 and 7, where *PiCO+* achieves substantially better performance than all the baselines. Our results validate the effectiveness of our framework, even in the presence of strong label ambiguity.

1. CIFAR-100 dataset consists of 20 superclasses, with 5 classes in each superclass.

6 WHY PiCO IMPROVES PARTIAL LABEL LEARNING?

In this section, we provide theoretical justification on why the contrastive prototypes help disambiguate the ground-truth label. We show that the *alignment* property in contrastive learning [11] intrinsically minimizes the intraclass covariance in the embedding space, which coincides with the objective of classical clustering algorithms. It motivates us to interpret PiCO through the lens of the expectation-maximization algorithm. To see this, we consider an ideal setting: in each training step, all data examples are accessible and the augmentation copies are also included in the training set, i.e., $A = \mathcal{D}$. Then, the contrastive loss is calculated as,

$$\begin{aligned} \tilde{\mathcal{L}}_{\text{cont}}(g; \tau, \mathcal{D}) &= \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} \left\{ -\frac{1}{|P(\mathbf{x})|} \sum_{\mathbf{k}_+ \in P(\mathbf{x})} \log \frac{\exp(\mathbf{q}^\top \mathbf{k}_+ / \tau)}{\sum_{\mathbf{k}' \in A(\mathbf{x})} \exp(\mathbf{q}^\top \mathbf{k}' / \tau)} \right\} \\ &= \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} \underbrace{\left\{ -\frac{1}{|P(\mathbf{x})|} \sum_{\mathbf{k}_+ \in P(\mathbf{x})} (\mathbf{q}^\top \mathbf{k}_+ / \tau) \right\}}_{(a)} \\ &\quad + \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} \underbrace{\left\{ \log \sum_{\mathbf{k}' \in A(\mathbf{x})} \exp(\mathbf{q}^\top \mathbf{k}' / \tau) \right\}}_{(b)}. \end{aligned} \quad (14)$$

We focus on analyzing the first term (a), which is often dubbed as the *alignment* term [11]. The main functionality of this term is to optimize the tightness of the clusters in the embedding space. In this work, we connect it with classical clustering algorithms. We first split the dataset to C subsets $S_j \in \mathcal{D}_C$ ($1 \leq j \leq C$), where each subset contains examples possessing the same predicted labels. In effect, our selection strategy in Eq. (4) constructs the positive set by selecting examples from the same subset. Therefore, we have,

$$\begin{aligned} (a) &= \frac{1}{n} \sum_{\mathbf{x} \in \mathcal{D}} \frac{1}{|P(\mathbf{x})|} \sum_{\mathbf{k}_+ \in P(\mathbf{x})} (|\mathbf{q} - \mathbf{k}_+|^2 - 2) / (2\tau) \\ &\approx \frac{1}{2\tau n} \sum_{S_j \in \mathcal{D}_C} \frac{1}{|S_j|} \sum_{\mathbf{x}, \mathbf{x}' \in S_j} \|g(\mathbf{x}) - g(\mathbf{x}')\|^2 + K \\ &= \frac{1}{\tau n} \sum_{S_j \in \mathcal{D}_C} \sum_{\mathbf{x} \in S_j} \|g(\mathbf{x}) - \boldsymbol{\mu}_j\|^2 + K, \end{aligned} \quad (15)$$

where K is a constant and $\boldsymbol{\mu}_j$ is the mean center of S_j . Here we approximate $\frac{1}{|S_j|} \approx \frac{1}{|S_j|-1} = \frac{1}{|P(\mathbf{x})|}$ since n is usually large. We omitted the augmentation operation for simplicity. The *uniformity* term (b) can benefit information-preserving, and has been analyzed in [11].

We are now ready to interpret the PiCO algorithm as an expectation-maximization algorithm that maximizes the likelihood of a generative model. At the E-step, the classifier assigns each data example to one specific cluster. At the M-step, the contrastive loss concentrates the embeddings to their cluster mean direction, which is achieved by minimizing Eq. (15). Finally, the training data will be mapped to a mixture of von Mises-Fisher distributions on the unit hypersphere.

EM Perspective. Recall that the candidate label set is a noisy version of the ground-truth. To estimate the likelihood

$P(Y_i, \mathbf{x}_i)$, we need to establish the relationship between the candidate and the ground-truth label. Following [5], we make a mild assumption,

Assumption 1. All labels y_i in the candidate label set have the same probability of generating Y_i , but no label outside of Y_i can generate Y_i , i.e. $P(Y_i|y_i) = h(Y_i)$ if $y_i \in Y_i$ else 0. Here $h(\cdot)$ is some function making it a valid probability distribution.

Then, we show that the PiCO implicitly maximizes the likelihood as follows,

E-Step. First, we introduce some distributions over all examples and the candidates $\pi_i^j \geq 0$ ($1 \leq i \leq n, 1 \leq j \leq C$) such that $\pi_i^j = 0$ if $j \notin Y_i$ and $\sum_{j \in Y_i} \pi_i^j = 1$. Let θ be the parameters of g . Our goal is to maximize the likelihood below,

$$\begin{aligned} &\arg \max_{\theta} \sum_{i=1}^n \log P(Y_i, \mathbf{x}_i | \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^n \log \sum_{y_i \in Y_i} P(\mathbf{x}_i, y_i | \theta) + \sum_{i=1}^n \log(h(Y_i)) \\ &= \arg \max_{\theta} \sum_{i=1}^n \log \sum_{y_i \in Y_i} \pi_i^{y_i} \frac{P(\mathbf{x}_i, y_i | \theta)}{\pi_i^{y_i}} \\ &\geq \arg \max_{\theta} \sum_{i=1}^n \sum_{y_i \in Y_i} \pi_i^{y_i} \log \frac{P(\mathbf{x}_i, y_i | \theta)}{\pi_i^{y_i}}. \end{aligned} \quad (16)$$

The last step of the derivation uses Jensen's inequality. By using the fact that $\log(\cdot)$ function is concave, the inequality holds with equality when $\frac{P(\mathbf{x}_i, y_i | \theta)}{\pi_i^{y_i}}$ is some constant. Therefore,

$$\pi_i^{y_i} = \frac{P(\mathbf{x}_i, y_i | \theta)}{\sum_{y_i \in Y_i} P(\mathbf{x}_i, y_i | \theta)} = \frac{P(\mathbf{x}_i, y_i | \theta)}{P(\mathbf{x}_i | \theta)} = P(y_i | \mathbf{x}_i, \theta), \quad (17)$$

which is the posterior class probability. In PiCO, it is estimated by using the classifier's output.

To estimate $P(y_i | \mathbf{x}_i, \theta)$, classical unsupervised clustering methods intuitively assign the data examples to the cluster centers, e.g. k-means. As in the supervised learning setting, we can directly use the ground-truth. However, under the setting of PLL, the supervision signals are situated between the supervised and unsupervised setups. Based on empirical findings, the candidate labels are more reliable for posterior estimation at the beginning; yet alongside the training process, the prototypes tend to become more trustful. This empirical observation has motivated us to update the pseudo targets in a moving-average style. Thereby, we have a good initialization in estimating class posterior, and it will be smoothly refined during the training procedure. This is verified in our empirical studies; see Section 5.2.2 and Appendix B.1. Finally, we take one-hot prediction $\tilde{y}_i = \arg \max_{j \in Y} f^j(\mathbf{x}_i)$ since each example inherently belongs to exactly one label and hence, we have $\pi_i^j = \mathbb{I}(\tilde{y}_i = j)$.

M-Step. At this step, we aim at maximizing the likelihood under the assumption that the posterior class probability is known. We show that under mild assumptions, minimizing Eq. (15) also maximizes a lower bound of likelihood in Eq. (16).

Theorem 1. Assume data from the same class in the contrastive output space follow a d -variate von Mises-Fisher (vMF) distribution whose probabilistic density is given by $f(\mathbf{x} | \boldsymbol{\mu}_i, \kappa) =$

$c_d(\kappa)e^{\kappa\bar{\mu}_i^\top g(x)}$, where $\bar{\mu}_i = \mu_i/\|\mu_i\|$ is the mean direction, κ is the concentration parameter, and $c_d(\kappa)$ is the normalization factor. We further assume a uniform class prior $P(y_i = j) = 1/C$. Let $n_j = |S_j|$. Then, optimizing Eq. (15) and Eq. (16) equal to maximize R_1 and R_2 below, respectively.

$$R_1 = \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\mu_j\|^2 \leq \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\mu_j\| = R_2. \quad (18)$$

The proof can be found in the Appendix A. Theorem 1 indicates that minimizing Eq. (15) also maximizes a lower bound of the likelihood in Eq. (16). The lower bound is tight when $\|\mu_j\|$ is close to 1, which in effect means a strong intraclass concentration on the hypersphere. Intuitively, when the hypothesis space is rich enough, it is possible to achieve a low intraclass covariance in the Euclidean space, resulting in a large norm of the mean vector $\|\mu_j\|$. Then, normalized embeddings in the hypersphere also have an intraclass concentration in a strong sense, because a large $\|\mu_j\|$ also results in a large κ [27]. Regarding the visualized representation in Figure 3, we note that PiCO is indeed able to learn compact clusters. Therefore, we have that minimizing the contrastive loss also partially maximizes the likelihood defined in Eq. (16).

7 RELATED WORKS

Partial Label Learning (PLL) allows each training example to be annotated with a candidate label set, in which the ground-truth is guaranteed to be included. The most intuitive solution is average-based methods [3], [4], [28], which treat all candidates equally. However, the key and obvious drawback is that the predictions can be severely misled by false positive labels. To disambiguate the ground-truth from the candidates, identification-based methods [29], which regard the ground-truth as a latent variable, have recently attracted increasing attention; representative approaches include maximum margin-based methods [30], [31], graph-based methods [6], [7], [32], [33], and clustering-based approaches [5]. Recently, self-training methods [18], [19], [20] have achieved state-of-the-art results on various benchmark datasets, which disambiguate the candidate label sets by means of the model outputs themselves. But, few efforts have been made to learn high-quality representations to reciprocate label disambiguation.

Contrastive Learning (CL) [12], [34] is a framework that learns discriminative representations through the use of instance similarity/dissimilarity. A plethora of works has explored the effectiveness of CL in unsupervised representation learning [12], [34], [35]. Recently, [8] propose supervised contrastive learning (SCL), an approach that aggregates data from the same class as the positive set and obtains improved performance on various supervised learning tasks. The success of SCL has motivated a series of works to apply CL to a number of weakly supervised learning tasks, including noisy label learning [25], [36], semi-supervised learning [37], [38], etc. Despite promising empirical results, however, these works, lack theoretical understanding. [11] theoretically show that the CL favors *alignment* and *uniformity*, and thoroughly analyzed the properties of uniformity. But, to date, the terminology *alignment* remains confusing; we show it inherently maps data points to a mixture of vMF distributions.

Noisy Label Learning (NLL) [39] aims at mitigating overfitting on mislabeled samples. One popular strategy is to design robust risk functions, including but does not limit to robust cross-entropy losses [22], [40], [41], sample re-weighting [42], [43], [44] and noise transition matrix-based loss correction [45], [46], [47]. Another active line of research relies on selecting clean samples from noisy ones [48]. Most of them adopt the small-loss selection criterion [13], [49] which is motivated by the fact that deep neural networks tend to memorize easy patterns first [50]. Based on that, the state-of-the-art NLL algorithms [15], [25], [51] regard the unchosen samples as unlabeled and incorporate semi-supervised learning (SSL) for boosted performance. Inspired by these works, PiCO+ incorporates a new distance-based selection criterion and extends the contrastive learning framework to facilitate SSL training. The most related one to our work is [10] which theoretically analyzes the robustness of average-based loss functions for the noisy PLL task. But, [10] does not provide a new empirically strong solution. Instead, our PiCO+ framework establishes promising results against label noise that makes the PLL problem more practical for open-world applications.

8 CONCLUSION

In this work, we propose a novel partial label learning framework PiCO. The key idea is to identify the ground-truth from the candidate set by using contrastively learned embedding prototypes. Empirically, we conducted extensive experiments and show that PiCO establishes state-of-the-art performance. Our results are competitive with the fully supervised setting, where the ground-truth label is given explicitly. Theoretical analysis shows that PiCO can be interpreted from an EM-algorithm perspective. Additionally, we extend the PiCO framework to PiCO+ which is able to learn robust classifiers from noisy partial labels. Applications of multi-class classification with ambiguous labeling can benefit from our method, and we anticipate further research in PLL to extend this framework to tasks beyond image classification. We hope our work will draw more attention from the community toward a broader view of using contrastive prototypes for partial label learning.

ACKNOWLEDGMENTS

HW, RX, GC and JZ are supported by the Key R&D Program of Zhejiang Province (Grant No. 2020C01024). JZ would also like to thank the Fundamental Research Funds for the Central Universities. YL is supported by Wisconsin Alumni Research Foundation (WARF), Facebook Research Award, and funding from Google Research. LF is supported by the National Natural Science Foundation of China (Grant No. 62106028).

REFERENCES

- [1] J. Luo and F. Orabona, "Learning from candidate labeling sets," in *NeurIPS*. Curran Associates, Inc., 2010, pp. 1504–1512.
- [2] C. Chen, V. M. Patel, and R. Chellappa, "Learning from ambiguously labeled face images," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 7, pp. 1653–1667, 2018.
- [3] E. Hüllermeier and J. Beringer, "Learning from ambiguously labeled examples," *Intell. Data Anal.*, vol. 10, no. 5, pp. 419–439, 2006.

- [4] T. Cour, B. Sapp, and B. Taskar, "Learning from partial labels," *J. Mach. Learn. Res.*, vol. 12, pp. 1501–1536, 2011.
- [5] L. Liu and T. G. Dietterich, "A conditional multinomial mixture model for superset label learning," in *NeurIPS*, 2012, pp. 557–565.
- [6] M. Zhang, B. Zhou, and X. Liu, "Partial label learning via feature-aware disambiguation," in *KDD*. ACM, 2016, pp. 1335–1344.
- [7] G. Lyu, S. Feng, T. Wang, C. Lang, and Y. Li, "GM-PLL: graph matching based partial label learning," *IEEE Trans. Knowl. Data Eng.*, vol. 33, no. 2, pp. 521–535, 2021.
- [8] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," in *NeurIPS*, 2020.
- [9] H. Wang, R. Xiao, Y. Li, L. Feng, G. Niu, G. Chen, and J. Zhao, "Pico: Contrastive label disambiguation for partial label learning," in *ICLR*. OpenReview.net, 2022.
- [10] J. Lv, L. Feng, M. Xu, B. An, G. Niu, X. Geng, and M. Sugiyama, "On the robustness of average losses for partial-label learning," *CoRR*, vol. abs/2106.06152, 2021.
- [11] T. Wang and P. Isola, "Understanding contrastive representation learning through alignment and uniformity on the hypersphere," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 9929–9939.
- [12] K. He, H. Fan, Y. Wu, S. Xie, and R. B. Girshick, "Momentum contrast for unsupervised visual representation learning," in *CVPR*. IEEE, 2020, pp. 9726–9735.
- [13] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. W. Tsang, and M. Sugiyama, "Co-teaching: Robust training of deep neural networks with extremely noisy labels," in *NeurIPS*, 2018, pp. 8536–8546.
- [14] D. Berthelot, N. Carlini, I. J. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," in *NeurIPS*, 2019, pp. 5050–5060.
- [15] J. Li, R. Socher, and S. C. H. Hoi, "Dividemix: Learning with noisy labels as semi-supervised learning," in *ICLR*. OpenReview.net, 2020.
- [16] E. Arazo, D. Ortego, P. Albert, N. E. O'Connor, and K. McGuinness, "Pseudo-labeling and confirmation bias in deep semi-supervised learning," in *IJCNN*. IEEE, 2020, pp. 1–8.
- [17] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [18] J. Lv, M. Xu, L. Feng, G. Niu, X. Geng, and M. Sugiyama, "Progressive identification of true labels for partial-label learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 6500–6510.
- [19] H. Wen, J. Cui, H. Hang, J. Liu, Y. Wang, and Z. Lin, "Leveraged weighted loss for partial label learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 139. PMLR, 2021, pp. 11 091–11 100.
- [20] L. Feng, J. Lv, B. Han, M. Xu, G. Niu, X. Geng, B. An, and M. Sugiyama, "Provably consistent partial-label learning," in *NeurIPS*, 2020.
- [21] L. Feng, T. Kaneko, B. Han, G. Niu, B. An, and M. Sugiyama, "Learning with multiple complementary labels," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 3072–3081.
- [22] Z. Zhang and M. R. Sabuncu, "Generalized cross entropy loss for training deep neural networks with noisy labels," in *NeurIPS*, 2018, pp. 8792–8802.
- [23] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "RandAugment: Practical data augmentation with no separate search," *CoRR*, vol. abs/1909.13719, 2019.
- [24] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [25] J. Li, C. Xiong, and S. C. H. Hoi, "Mopro: Webly supervised learning with momentum prototypes," in *ICLR*. OpenReview.net, 2021.
- [26] P. Welinder, S. Branson, T. Mita, C. Wah, F. Schroff, S. Belongie, and P. Perona, "Caltech-UCSD Birds 200," California Institute of Technology, Tech. Rep. CNS-TR-2010-001, 2010.
- [27] A. Banerjee, I. S. Dhillon, J. Ghosh, and S. Sra, "Clustering on the unit hypersphere using von mises-fisher distributions," *J. Mach. Learn. Res.*, vol. 6, pp. 1345–1382, 2005.
- [28] M. Zhang and F. Yu, "Solving the partial label learning problem: An instance-based approach," in *IJCAI*. AAAI Press, 2015, pp. 4048–4054.
- [29] R. Jin and Z. Ghahramani, "Learning with multiple labels," in *NeurIPS*. MIT Press, 2002, pp. 897–904.
- [30] N. Nguyen and R. Caruana, "Classification with partial labels," in *SIGKDD*. ACM, 2008, pp. 551–559.
- [31] H. Wang, Y. Qiang, C. Chen, W. Liu, T. Hu, Z. Li, and G. Chen, "Online partial label learning," in *ECML PKDD*, ser. Lecture Notes in Computer Science, vol. 12458. Springer, 2020, pp. 455–470.
- [32] D. Wang, L. Li, and M. Zhang, "Adaptive graph guided disambiguation for partial label learning," in *SIGKDD*. ACM, 2019, pp. 83–91.
- [33] N. Xu, J. Lv, and X. Geng, "Partial label learning via label enhancement," in *AAAI*. AAAI Press, 2019, pp. 5557–5564.
- [34] A. van den Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," *CoRR*, vol. abs/1807.03748, 2018.
- [35] T. Chen, S. Kornblith, M. Norouzi, and G. E. Hinton, "A simple framework for contrastive learning of visual representations," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 119. PMLR, 2020, pp. 1597–1607.
- [36] Z. Wu, T. Wei, J. Jiang, C. Mao, M. Tang, and Y. Li, "NGC: A unified framework for learning with open-world noisy data," *CoRR*, vol. abs/2108.11035, 2021.
- [37] J. Li, C. Xiong, and S. C. H. Hoi, "Comatch: Semi-supervised learning with contrastive graph regularization," *CoRR*, vol. abs/2011.11183, 2020.
- [38] Y. Zhang, X. Zhang, R. C. Qiu, J. Li, H. Xu, and Q. Tian, "Semi-supervised contrastive learning with similarity co-calibration," *CoRR*, vol. abs/2105.07387, 2021.
- [39] B. Han, Q. Yao, T. Liu, G. Niu, I. W. Tsang, J. T. Kwok, and M. Sugiyama, "A survey of label-noise representation learning: Past, present and future," *CoRR*, vol. abs/2011.04406, 2020.
- [40] Y. Wang, X. Ma, Z. Chen, Y. Luo, J. Yi, and J. Bailey, "Symmetric cross entropy for robust learning with noisy labels," in *ICCV*. IEEE, 2019, pp. 322–330.
- [41] L. Feng, S. Shu, Z. Lin, F. Lv, L. Li, and B. An, "Can cross entropy loss be robust to label noise?" in *IJCAI*. ijcai.org, 2020, pp. 2206–2212.
- [42] L. Jiang, Z. Zhou, T. Leung, L. Li, and L. Fei-Fei, "Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 80. PMLR, 2018, pp. 2309–2318.
- [43] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 80. PMLR, 2018, pp. 4331–4340.
- [44] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," in *NeurIPS*, 2019, pp. 1917–1928.
- [45] N. Natarajan, I. S. Dhillon, P. Ravikumar, and A. Tewari, "Learning with noisy labels," in *NeurIPS*, 2013.
- [46] G. Patrini, A. Rozza, A. K. Menon, R. Nock, and L. Qu, "Making deep neural networks robust to label noise: A loss correction approach," in *CVPR*. IEEE Computer Society, 2017, pp. 2233–2241.
- [47] Y. Yao, T. Liu, B. Han, M. Gong, J. Deng, G. Niu, and M. Sugiyama, "Dual T: reducing estimation error for transition matrix in label-noise learning," in *NeurIPS*, 2020.
- [48] X. Xia, T. Liu, B. Han, M. Gong, J. Yu, G. Niu, and M. Sugiyama, "Sample selection with uncertainty of losses for learning with noisy labels," in *ICLR*. OpenReview.net, 2022.
- [49] H. Wei, L. Feng, X. Chen, and B. An, "Combating noisy labels by agreement: A joint training method with co-regularization," in *CVPR*. Computer Vision Foundation / IEEE, 2020, pp. 13 723–13 732.
- [50] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," in *ICLR*. OpenReview.net, 2017.
- [51] D. T. Nguyen, C. K. Mummadi, T. Ngo, T. H. P. Nguyen, L. Beggel, and T. Brox, "SELF: learning to filter noisy labels with self-ensembling," in *ICLR*. OpenReview.net, 2020.
- [52] L. Bottou and Y. Bengio, "Convergence properties of the k-means algorithms," in *NIPS*, G. Tesauro, D. S. Touretzky, and T. K. Leen, Eds. MIT Press, 1994, pp. 585–592.
- [53] J. Li, P. Zhou, C. Xiong, and S. C. H. Hoi, "Prototypical contrastive learning of unsupervised representations," in *ICLR*. OpenReview.net, 2021.
- [54] N. Saunshi, O. Plevrakis, S. Arora, M. Khodak, and H. Khandeparkar, "A theoretical analysis of contrastive unsupervised representation learning," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 97. PMLR, 2019, pp. 5628–5637.

- [55] M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *NeurIPS*, 2020.
- [56] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in *ICLR*. OpenReview.net, 2017.
- [57] K. Sohn, D. Berthelot, N. Carlini, Z. Zhang, H. Zhang, C. Raffel, E. D. Cubuk, A. Kurakin, and C. Li, "Fixmatch: Simplifying semi-supervised learning with consistency and confidence," in *NeurIPS*, 2020.
- [58] J. Snell, K. Swersky, and R. S. Zemel, "Prototypical networks for few-shot learning," in *NeurIPS*, 2017, pp. 4077–4087.
- [59] W. Xu, Y. Xian, J. Wang, B. Schiele, and Z. Akata, "Attribute prototype network for zero-shot learning," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- [60] S. Ö. Arik and T. Pfister, "Attention-based prototypical learning towards interpretable, confident and robust deep neural networks," *CoRR*, vol. abs/1902.06292, 2019.
- [61] T. Han, J. Gao, Y. Yuan, and Q. Wang, "Unsupervised semantic aggregation and deformable template matching for semi-supervised learning," in *NeurIPS 2020*, 2020.

APPENDIX A THEORETICAL ANALYSIS

Derivation of Eq. (15). We provide the derivation of the second equality in Eq. (15). It suffices to show that LHS = $\frac{1}{2n_j} \sum_{\mathbf{x}, \mathbf{x}' \in S_j} \|\mathbf{g}(\mathbf{x}) - \mathbf{g}(\mathbf{x}')\|^2 = \sum_{\mathbf{x} \in S_j} \|\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_j\|^2 =$ RHS. We have,

$$\begin{aligned}
 \text{LHS} &= \frac{1}{2n_j} \sum_{\mathbf{x} \in S_j} \sum_{\mathbf{x}' \in S_j} (\|\mathbf{g}(\mathbf{x})\|^2 - 2\mathbf{g}(\mathbf{x})^\top \mathbf{g}(\mathbf{x}') + \|\mathbf{g}(\mathbf{x}')\|^2) \\
 &= \frac{1}{n_j} \sum_{\mathbf{x} \in S_j} (n_j - \mathbf{g}(\mathbf{x})^\top (\sum_{\mathbf{x}' \in S_j} \mathbf{g}(\mathbf{x}')))) \\
 &= \frac{1}{n_j} \sum_{\mathbf{x} \in S_j} (n_j - \mathbf{g}(\mathbf{x})^\top (n_j \boldsymbol{\mu}_j)) \\
 &= n_j - (\sum_{\mathbf{x} \in S_j} \mathbf{g}(\mathbf{x}))^\top \boldsymbol{\mu}_j = n_j(1 - \|\boldsymbol{\mu}_j\|^2).
 \end{aligned} \tag{19}$$

On the other hand,

$$\begin{aligned}
 \text{RHS} &= \sum_{\mathbf{x} \in S_j} (\|\mathbf{g}(\mathbf{x})\|^2 - 2\mathbf{g}(\mathbf{x})^\top \boldsymbol{\mu}_j + \|\boldsymbol{\mu}_j\|^2) \\
 &= (n_j - 2(\sum_{\mathbf{x} \in S_j} \mathbf{g}(\mathbf{x}))^\top \boldsymbol{\mu}_j + n_j \|\boldsymbol{\mu}_j\|^2) \\
 &= n_j(1 - \|\boldsymbol{\mu}_j\|^2) = \text{LHS}.
 \end{aligned} \tag{20}$$

Proof of Theorem 1. By regarding π_i^j as constants w.r.t

θ , we can get the following derivation from Eq. (16),

$$\begin{aligned}
 &\arg \max_{\theta} \sum_{i=1}^n \sum_{y_i \in Y_i} \pi_i^{y_i} \log \frac{P(\mathbf{x}_i, y_i | \theta)}{\pi_i^{y_i}} \\
 &= \arg \max_{\theta} \sum_{i=1}^n \sum_{y_i \in Y_i} \pi_i^{y_i} \log P(\mathbf{x}_i | y_i, \theta) P(y_i) \\
 &= \arg \max_{\theta} \sum_{i=1}^n \sum_{y_i \in Y_i} \mathbb{I}(\tilde{y}_i = y_i) \log P(\mathbf{x}_i | y_i, \theta) \\
 &= \arg \max_{\theta} \sum_{S_j \in \mathcal{D}_C} \sum_{\mathbf{x} \in S_j} \log P(\mathbf{x} | y = j, \theta) \\
 &= \arg \max_{\theta} \sum_{S_j \in \mathcal{D}_C} \sum_{\mathbf{x} \in S_j} (\kappa \bar{\boldsymbol{\mu}}_j^\top \mathbf{g}(\mathbf{x}) + \log(c_d(\kappa))) \\
 &= \arg \max_{\theta} \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\boldsymbol{\mu}_j\|
 \end{aligned} \tag{21}$$

where $n_j = |S_j|$. Here we ignore the constant factor $-\sum_{i=1}^n \sum_{y_i \in Y_i} \pi_i^{y_i} \log \pi_i^{y_i}$ w.r.t. θ in the first equality. In the last equality, we use the fact that $\boldsymbol{\mu}_j = \frac{1}{n_j} \sum_{\mathbf{x} \in S_j} \mathbf{g}(\mathbf{x})$ and $\bar{\boldsymbol{\mu}}_j$ is the unit directional vector of $\boldsymbol{\mu}_j$. From Eq. (15), we have that,

$$\begin{aligned}
 &\arg \min_{\theta} \sum_{S_j \in \mathcal{D}_C} \sum_{\mathbf{x} \in S_j} \|\mathbf{g}(\mathbf{x}) - \boldsymbol{\mu}_j\|^2 \\
 &= \arg \min_{\theta} \sum_{S_j \in \mathcal{D}_C} \sum_{\mathbf{x} \in S_j} (\|\mathbf{g}(\mathbf{x})\|^2 - 2\mathbf{g}(\mathbf{x})^\top \boldsymbol{\mu}_j + \|\boldsymbol{\mu}_j\|^2) \\
 &= \arg \min_{\theta} \sum_{S_j \in \mathcal{D}_C} (n_j - n_j \|\boldsymbol{\mu}_j\|^2) \\
 &= \arg \max_{\theta} \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\boldsymbol{\mu}_j\|^2.
 \end{aligned} \tag{22}$$

Note that the contrastive embeddings are distributed on the hypersphere S^{d-1} and thus $\|\boldsymbol{\mu}_j\| \in [0, 1]$. It can be directly derived that,

$$R_1 = \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\boldsymbol{\mu}_j\|^2 \leq \sum_{S_j \in \mathcal{D}_C} \frac{n_j}{n} \|\boldsymbol{\mu}_j\| = R_2. \tag{23}$$

Therefore, maximizing the intraclass covariance in Eq. (15) is equivalent to maximizing a lower bound of the likelihood in Eq. (16). It can also be shown that $(R_2)^2 \leq R_1$ followed by the convexity of the squared function. Since $\arg \max R_2 = \arg \max (R_2)^2$, we have that the contrastive loss also maximizes an upper bound of R_2 .

Unfortunately, there is no guarantee that the lower bound is tight without further assumptions. To see this, assume that we have two classes $y \in \{1, 2\}$ with equal-sized samples and their mean vectors have the norm of $\|\boldsymbol{\mu}_1\| = 0.5$ and $\|\boldsymbol{\mu}_2\| = 1$. We have that $R_1 = 0.625$ and $R_2 = 0.75$, which demonstrates a large discrepancy. It is interesting to see that when the norm of the mean vectors are the same, i.e. $\|\boldsymbol{\mu}_j\| = \|\boldsymbol{\mu}_k\|$ for all $1 \leq j \leq k \leq C$, we have $(R_2)^2 = R_1$ by the Jensen's inequality, making the upper bound tight. But, it is not a trivial condition.

To obtain a tight lower bound, what we need is a rich enough hypothesis space to achieve a low intraclass covariance in Eq. (15), and hence a large R_1 . We show that it inherently produces compact vMF distributions. To see this, it should be noted that the concentration parameter κ

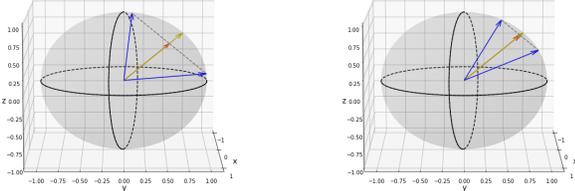


Fig. 6. A large norm of Euclidean’s mean vector also leads to a strong concentration of unit vectors to its mean direction.

of a vMF distribution is given by the inverse of the ratio of Bessel functions of mean vector μ_j . Though it is not possible to obtain an analytic solution of κ , we have the following well-known approximation [27],

$$\begin{aligned} \kappa &\approx \frac{d-1}{2(1-\|\mu_j\|)}, && \text{valid for large } \|\mu_j\|, \\ \kappa &\approx d\|\mu_j\| \left(1 + \frac{d}{d+2}\|\mu_j\|^2 + \frac{d^2(d+8)}{(d+2)^2(d+4)}\|\mu_j\|^4 \right), && \text{valid for small } \|\mu_j\|. \end{aligned} \tag{24}$$

The above approximations show that a larger norm of Euclidean’s mean μ_j typically leads to a stronger concentration on the hypersphere. By Eq. (22), we know that contrastive loss encourages a large norm of μ_j , and thus also tightly clusters the embeddings on the hypersphere; see Figure 6.

We further note that we do not include a k-means process in our PiCO method. PiCO is related to center-based clustering algorithms in our theoretical analysis. Since we restrict the gold label to be included in the candidate set, we believe that this piece of information could largely help avoid the bad optimum problem that occurs in a pure unsupervised setup. For the convergence properties of our PiCO algorithm, we did not empirically find any issues with PiCO converging to a (perhaps locally) optimal point. However, we want to refer the readers to the proof of k-means clustering in [52] for the convergence analysis.

Discussion. Regarding our empirical results where PiCO does indeed learn compact representations for each class, we can conclude that PiCO implicitly clusters data points in the contrastive embeddings space as a mixture of vMF distributions. In each iteration, our algorithm can be viewed as alternating the two steps until convergence, though different in detail. First, it is intractable to handle the whole training dataset, and thus we accelerate via a MoCo-style dictionary and MA-updated prototypes. Second, the contrastive loss also encourages the *uniformity* of the embeddings to maximally preserve information, which serves as a regularizer and typically leads to better representation. Finally, we use two copies of data examples such that data are also aligned in its local region. Moreover, our theoretical result also answers why merely taking the pseudo-labels to select positive examples also leads to superior results, since the selected positive set will be refined as the training procedure proceeds.

Our theoretical results are also generalizable to other contrastive learning methods. For example, the classical unsupervised contrastive learning [12] actually assumes n -prototypes to cluster all locally augmented data; prototypical contrastive learning [53] directly assigns the data

TABLE 8
Performance of PiCO with varying γ on CIFAR-10 ($q = 0.5$) and CIFAR-100 ($q = 0.05$).

Dataset	$\gamma = 0.1$	0.5	0.9	0.99	0.999
CIFAR-10	93.61	93.51	93.52	93.58	93.66
CIFAR-100	72.87	73.09	72.54	72.74	67.33

TABLE 9
Training accuracy of pseudo targets on CIFAR-10 and CIFAR-100.

Dataset	CIFAR-10			CIFAR-100		
q	0.1	0.3	0.5	0.01	0.05	0.1
Accuracy	98.28	98.26	96.79	99.06	96.27	90.58

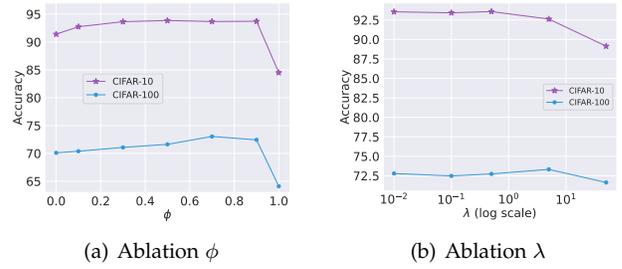


Fig. 7. More ablation results on CIFAR-10 ($q = 0.5$) and CIFAR-100 ($q = 0.05$). (a) Performance of PiCO with varying ϕ . (b) Performance of PiCO with varying λ .

examples to one cluster to get the posterior, since there are no supervision signals; supervised contrastive learning [8] chooses the known ground-truth as the posterior. In our problem, we follow two extreme settings to progressively obtain an accurate posterior estimation. Finally, it is also noteworthy that the objective in Eq. (15) has a close connection to the intra-class deviation [54], minimizing which is proven to be beneficial in obtaining tighter generalization error bound on downstream tasks. It should further be noted that our work differs from existing clustering-based CL methods [53], [55], which explicitly involves clustering to aggregate the embeddings; instead, our results are derived from the loss itself.

APPENDIX B MORE EXPERIMENTAL RESULTS

B.1 Ablation Studies

Moving-average updating factor ϕ . We first present more ablation results about the effect of pseudo target updating factor ϕ on PiCO performance. Figure 7 (a) shows the results on two datasets CIFAR-10 ($q = 0.5$) and CIFAR-100 ($q = 0.05$). The overall trends on both datasets follow our arguments in Section 5.2.2. Specifically, performance on CIFAR-100 achieves the best result when $\phi = 0.7$, and slight drops when $\phi = 0.9$. Therefore, in practice, we may achieve a better result with careful fine-tuning on ϕ value. In contrast, PiCO works well in a wide range of ϕ values on CIFAR-10. The reason might be that CIFAR-10 is a simpler

TABLE 10
Accuracy comparisons on fine-grained datasets. Bold indicates superior results.

Dataset	Method	$q = 0.1$	$q = 0.5$	$q = 0.8$
CIFAR-100-H	PiCO (ours)	73.41 \pm 0.27%	72.04 \pm 0.31%	66.17 \pm 0.23%
	LWS	62.41 \pm 0.03%	57.25 \pm 0.02%	20.64 \pm 0.48%
	PRODEN	62.91 \pm 0.01%	60.89 \pm 0.03%	43.64 \pm 1.82%
	CC	50.40 \pm 0.20%	42.60 \pm 0.11%	37.80 \pm 0.09%
	MSE	46.05 \pm 0.17%	39.52 \pm 0.28%	15.18 \pm 0.73%
	EXP	45.73 \pm 0.22%	35.08 \pm 1.71%	22.31 \pm 0.39%
Dataset	Method	$q = 0.01$	$q = 0.05$	$q = 0.1$
CUB-200	PiCO (ours)	74.14 \pm 0.24%	72.17 \pm 0.72%	62.02 \pm 1.16%
	LWS	73.74 \pm 0.23%	39.74 \pm 0.47%	12.30 \pm 0.77%
	PRODEN	72.34 \pm 0.04%	62.56 \pm 0.10%	35.89 \pm 0.05%
	CC	56.63 \pm 0.01%	55.61 \pm 0.02%	17.01 \pm 1.44%
	MSE	61.12 \pm 0.51%	22.07 \pm 2.36%	11.40 \pm 2.42%
	EXP	55.62 \pm 2.25%	9.44 \pm 2.32%	7.3 \pm 0.99%
	Fully Supervised		76.02 \pm 0.19%	

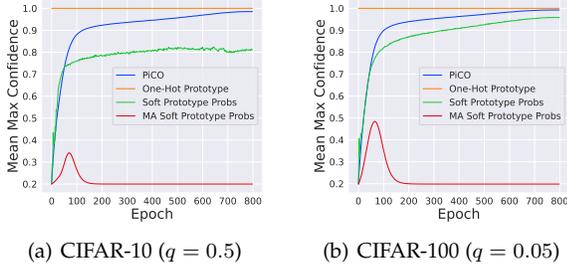


Fig. 8. The mean max confidence curves of different label disambiguation strategies.

version of CIFAR-100, and thus the prototypes can be high-quality quickly. But, setting ϕ to either 0 or 1 leads to a worse result, which has been discussed earlier.

Loss weight λ . Figure 7 (b) reports the performance of PiCO with varying λ values that trade-off the classification and contrastive losses. λ is selected from $\{0.01, 0.1, 0.5, 5, 50\}$. We can observe that on CIFAR-10, the performance is stable, but on CIFAR-100, the best performance is obtained at $\lambda = 5$. When $\lambda = 50$, PiCO shows inferior results on both two datasets. In general, a relatively small λ (< 10) usually leads to good performance than a much larger value. When λ is large, the contrastive network tends to fit noisy labels at the early stage of training.

Prototype updating factor γ . Then, we show the effect of γ that controls the speed of prototype updating and the results are listed in Table 8. On the CIFAR-10 dataset, the performance is stable with varying γ . But, on CIFAR-100, it can be seen that too large λ leads to a significant performance drop, which may be caused by insufficient label disambiguation.

The disambiguation ability of PiCO. Next, we evaluate the disambiguation ability of the proposed PiCO. To see this, we calculate the max confidence to represent the uncertainty of an example, which has been widely used in recent works [56]. If one example is uncertain about its ground-truth, then it typically associates with low max confidence $\max_j s_j$. To represent the uncertainty of the whole training dataset, we calculate the mean max confidence (MMC) score. In Figure 8, we plot the MMC scores of different label disambiguation

TABLE 11
Performance of PiCO+ with varying α on CIFAR-10 ($q = 0.5, \eta = 0.2$) and CIFAR-100 ($q = 0.05, \eta = 0.2$).

Dataset	$\alpha = 0.1$	0.5	1	2	10
CIFAR-10	67.41	76.16	92.77	92.59	90.92
CIFAR-100	75.46	74.83	73.93	72.98	69.75

strategies in different training epochs. First, we can observe the MMC score of PiCO smoothly increases and finally achieves near 1 results, which means most of the labels are well disambiguated. In contrast, the Soft Prototype Probs strategy oscillates at the beginning, and then also increases to a certain value, which means that directly adopting the soft class probability also helps disambiguation. But, it is worth noting that it ends with a smaller MMC score compared with PiCO. The reason might be that the cosine distances to non-ground-truth prototypes are still at a scale. Hence, the Soft Prototype Probs strategy always holds a certain degree of ambiguity. Finally, we can see that the MA Soft Prototype Probs strategy fails to achieve great disambiguation ability. Compared with the non-moving-average version, it fails to get rid of severe label ambiguity and will finally converge to uniform pseudo targets again.

Furthermore, we evaluated the accuracy of the pseudo targets over training examples. From Table 9, we can find that the pseudo targets achieve high training accuracy. Combined with the fact that the mean max confidence score of the pseudo target is close to 1, the training examples finally become near supervised ones. Thus, the proposed PiCO is able to achieve near-supervised learning performance, especially when the label ambiguity is not too high. These results verify that PiCO has a strong label disambiguation ability to handle the PLL problem.

Clean loss weight α for PiCO+. Lastly, we show the effect of clean loss weight α for PiCO+ in Table 11. With a too large α , the performance of PiCO+ typically drops since the regularization effect of mixup training and semi-supervised learning becomes weak. When α is too small, we observe that PiCO+ exhibits degenerated performance on CIFAR-10. The reason is that the clean examples are dominated by unreliable samples, resulting in severe confirmation bias

TABLE 12
Accuracy comparisons with different data generation processes on CIFAR-10.

Method	Case (1)	Case (2)
PiCO	94.49 ± 0.08%	94.11 ± 0.25%
LWS	90.78 ± 0.01%	68.37 ± 0.04%
PRODEN	90.53 ± 0.01%	87.02 ± 0.02%
CC	75.81 ± 0.13%	66.51 ± 0.11%
MSE	68.11 ± 0.23%	39.49 ± 0.41%
EXP	71.62 ± 0.79%	48.87 ± 2.32%

and wrong selection.

B.2 Additional Results on Fine-Grained Classification

In the sequel, we show the full setups and experimental results on fine-grained classification datasets. In particular, on CUB-200, we set the length of the momentum queue as 4192. For CUB-200, we set the input image resolution as 224×224 and select $q \in \{0.01, 0.05, 0.1\}$. When $q = 0.05/0.1$, we warm up for 20/100 epochs and train the model for 200/300 epochs, respectively. For CIFAR-100-H, we select $q \in \{0.1, 0.5, 0.8\}$ and warm up the model for 100 epochs when $q = 0.8$. Other hyperparameters are the same as our default setting. The baselines are also fine-tuned to achieve their best results.

From Table 10, we can observe that PiCO significantly outperforms all the baselines on all the datasets. Moreover, as the size of candidate sets grows larger, PiCO consistently leads by an even wider margin. For example, on CIFAR-100-H, compared with the best baseline, performance improvement reaches **9.50%**, **11.15%** and **22.53%** in accuracy when $q = 0.1, 0.5, 0.8$, respectively. The comparison emerges the dominance of our label disambiguation strategy among semantically similar classes.

B.3 Strategies for Positive Selection

While our positive set selection strategy is simple and effective, one may still explore more complicated strategies to boost performance. We have empirically tested two strategies: 1) **Filter-based**: we set a filter $\frac{|Y_i \cap Y_j|}{|Y_i \cup Y_j|} \leq \rho$ ($\rho = 0.5$) to remove example pairs who have dissimilar candidate sets at the early stage. 2) **Threshold-based**: we set a threshold $\max f^j(\text{Aug}_q(\mathbf{x})) \leq \delta$ ($\delta = 0.95$) to remove those uncertain examples at the end of training, which has been widely used in semi-supervised learning [57]. Our basic principle is that contrastive learning is robust to noisy negative pairs and thus, we can flip those less reliable positive pairs to negative. Unfortunately, we did not observe statistically significant improvement to our vanilla strategy in experiments. These negative results suggest that the proposed PiCO has a strong error correction ability, which corroborates our theoretical analysis.

B.4 The Influence of Data Generation

In practice, some labels may be more analogous to the true label than others, which makes their probability of label flipping q larger than others. In other words, the data generation procedure is non-uniform. In Section 5.4, we have

Algorithm 2: Pseudo-code of PiCO+.

```

1 Input: Training dataset  $\mathcal{D}$ , selection ratio  $\delta$ , number of
   nearest neighbors  $k$ , beta distribution shape parameter
    $\varsigma$ , loss weighting factors  $\alpha, \beta$ .
2 warm up by running PiCO on  $\mathcal{D}$ 
3 for  $epoch = 1, 2, \dots$ , do
4   split a clean set by  $\mathcal{D}_{\text{clean}} = \{(\mathbf{x}_i, Y_i) | \mathbf{q}_i^\top \boldsymbol{\mu}_{\hat{y}_i} > \kappa_\delta\}$ 
5   set  $\mathcal{D}_{\text{noisy}} = \mathcal{D} \setminus \mathcal{D}_{\text{clean}}$ 
6   for  $iter = 1, 2, \dots$ , do
7     sample a mini-batch  $B$  from  $\mathcal{D}$ 
8      $B_{\text{clean}} = B \cap \mathcal{D}_{\text{clean}}, B_{\text{noisy}} = B \cap \mathcal{D}_{\text{noisy}}$ 
9     run PiCO on  $B_{\text{clean}}$  to get loss  $\mathcal{L}_{\text{clean}}$ 
10    for  $\mathbf{x}_i \in B$  do
11      // label-driven positive set
12       $P_{\text{noisy}}(\mathbf{x}_i) = \{\mathbf{k}' | \mathbf{k}' \in A(\mathbf{x}_i), \hat{y}'_i = \hat{y}_i\}$ 
13    end
14    for  $\mathbf{x}_i \in B_{\text{noisy}}$  do
15      // neighbors-based positive set
16       $P_{\text{knn}}(\mathbf{x}_i) = \{\mathbf{k}' | \mathbf{k}' \in A(\mathbf{x}_i) \cap \mathcal{N}_k(\mathbf{x}_i)\}$ 
17      // label guessing
18       $s'_{ij} = \frac{\exp(\mathbf{q}_i^\top \boldsymbol{\mu}_j / \tau)}{\sum_{t=1}^L \exp(\mathbf{q}_i^\top \boldsymbol{\mu}_t / \tau)}, \forall 1 \leq j \leq L$ 
19    end
20    for  $(\mathbf{x}_i, (\mathbf{x}_j) \in B$  do
21       $\sigma \sim \text{Beta}(\varsigma, \varsigma)$ 
22      // mixup samples and pseudo-targets
23       $\mathbf{x}_{ij}^m = \sigma \text{Aug}_q(\mathbf{x}_i) + (1 - \sigma) \text{Aug}_q(\mathbf{x}_j)$ 
24       $\mathbf{s}_{ij}^m = \sigma \hat{\mathbf{s}}_i + (1 - \sigma) \hat{\mathbf{s}}_j$ 
25    end
26  calculate  $\mathcal{L}_{\text{n-cont}}, \mathcal{L}_{\text{knn}}, \mathcal{L}_{\text{n-cls}}, \mathcal{L}_{\text{mix}}$ 
27  minimize loss
28   $\mathcal{L}_{\text{pico+}} = \mathcal{L}_{\text{mix}} + \alpha \mathcal{L}_{\text{clean}} + \beta (\mathcal{L}_{\text{n-cont}} + \mathcal{L}_{\text{knn}} + \mathcal{L}_{\text{n-cls}})$ 
29 end

```

shown one such case on CIFAR-100-H, where semantically similar labels have a larger probability of being a false positive. Moreover, we follow [19] to conduct empirical comparisons on data with alternative generation processes. In particular, we test two commonly used cases on CIFAR-10 with the following flipping matrix, respectively:

$$(1) = \begin{bmatrix} 1 & 0.5 & 0 & \dots & 0 \\ 0 & 1 & 0.5 & \dots & 0 \\ \vdots & & \dots & \ddots & \vdots \\ 0.5 & 0 & 0 & \dots & 1 \end{bmatrix},$$

$$(2) = \begin{bmatrix} 1 & 0.9 & 0.7 & 0.5 & 0.3 & 0.1 & 0 & \dots & 0 \\ 0 & 1 & 0.9 & 0.7 & 0.5 & 0.3 & 0.1 & \dots & 0 \\ \vdots & & & \dots & & & & \ddots & \vdots \\ 0.9 & 0.7 & 0.5 & 0.3 & 0.1 & 0 & 0 & \dots & 1 \end{bmatrix}$$

where each entry denotes the probability of a label being a candidate. As shown in Table 12, PiCO outperforms other baselines in both cases. It is worth noting that in Case (2), each ground-truth label has a maximum probability of 0.9 of being coupled with the same false positive label. In such a challenging setup, PiCO still achieves promising results that are competitive with the supervised performance, which further verifies its strong disambiguation ability.

B.5 The Influence of Prototype Calculation

There are several ways to calculate the prototypes and hence, we further test a variant of PiCO that re-computes

TABLE 13
Training time (min/epoch) and accuracy of different prototype calculation methods.

Dataset	Method	Time	Accuracy
CIFAR-10 ($q = 0.5$)	PiCO	0.94	93.58
	Re-Compute	1.39	93.55
CIFAR-100 ($q = 0.05$)	PiCO	0.96	72.74
	Re-Compute	1.40	72.35

the prototypes by averaging embeddings of all training examples at the end of each epoch. We train the models using one Quadro P5000 GPU respectively and evaluate the average training time per epoch. From Table 13, we can observe that the Re-Compute variant achieves competitive results, but is much slower than PiCO.

APPENDIX C PSEUDO-CODE OF PICO+

We summarize the pseudo-code of our PiCO+ method in Algorithm 2.

APPENDIX D THE LITERATURE OF PROTOTYPE LEARNING

Prototype learning (PL) aims to learn a metric space where examples are enforced to be closer to its class prototype. PL is typically more robust in handling few-shot learning [58], zero-shot learning [59], and out-of-distribution samples [60]. Recently, PL has demonstrated promising results in weakly-supervised learning, such as semi-supervised learning [61], noisy-label learning [25], etc. For example, USADTM [61] shows that informative class prototypes usually lead to better pseudo-labels for semi-supervised learning than classical pseudo-labeling algorithms [57] which reuse classifier outputs. Motivated by this, we also employ contrastive prototypes for label disambiguation.