# CLImage: Human-Annotated Datasets for Complementary-Label Learning

Hsiu-Hsuan Wang, Tan-Ha Mai, Nai-Xuan Ye, Wei-I Lin, Hsuan-Tien Lin National Taiwan University {b09902033, d10922024, b09902008, r10922076, htlin}@csie.ntu.edu.tw

#### Abstract

Complementary-label learning (CLL) is a weakly-supervised learning paradigm 1 2 that aims to train a multi-class classifier using only complementary labels, which indicate classes to which an instance does not belong. Despite numerous algo-3 rithmic proposals for CLL, their practical applicability remains unverified for two 4 reasons. Firstly, these algorithms often rely on assumptions about the generation of 5 complementary labels, and it is not clear how far the assumptions are from reality. 6 7 Secondly, their evaluation has been limited to synthetic datasets. To gain insights into the real-world performance of CLL algorithms, we developed a protocol to 8 collect complementary labels from human annotators. Our efforts resulted in the 9 creation of four datasets: CLCIFAR10, CLCIFAR20, CLMicroImageNet10, and 10 CLMicroImageNet20, derived from well-known classification datasets CIFAR10, 11 CIFAR100, and TinyImageNet200. These datasets represent the very first real-12 world CLL datasets. Through extensive benchmark experiments, we discovered 13 a notable decrease in performance when transitioning from synthetic datasets to 14 15 real-world datasets. We investigated the key factors contributing to the decrease with a thorough dataset-level ablation study. Our analyses highlight annotation 16 noise as the most influential factor in the real-world datasets. In addition, we 17 discover that the biased-nature of human-annotated complementary labels and the 18 difficulty to validate with only complementary labels are two outstanding barriers 19 to practical CLL. These findings suggest that the community focus more research 20 efforts on developing CLL algorithms and validation schemes that are robust to 21 noisy and biased complementary-label distributions. 22

# 23 **1** Introduction

Ordinary multi-class classification methods rely heavily on high-quality labels to train effective
classifiers. However, such labels can be expensive and time-consuming to collect in many real-world
applications. To address this challenge, researchers have turned their attention towards weaklysupervised learning, which aims to learn from incomplete, inexact, or inaccurate data sources [20, 28].
This learning paradigm includes but is not limited to noisy-label learning [5], partial-label learning [2],
positive-unlabeled learning [3], and complementary-label learning [8].

In this work, we focus on complementary-label learning (CLL). This learning problem involves
 training a multi-class classifier using only complementary labels, which indicate the classes that a
 data instance does not belong to. Although several algorithms have been proposed to learn from

33 complementary labels, they were only benchmarked on synthetic datasets with some idealistic

assumptions on complementary-label generation [1, 8, 9, 16, 21]. Thus, it remains unclear how well
 these algorithms perform in practical scenarios.

In particular, current CLL algorithms heavily rely on the *uniform assumption* for generating comple-36 mentary labels [8], which specifies that complementary labels are generated by uniformly sampling 37 from the set of all possible complementary labels. To alleviate the restrictiveness of the uniform 38 assumption, Yu et al. [27] considered a more general *class-conditional assumption*, where the dis-39 tribution of the complementary labels only depends on its ordinary labels. These assumptions have 40 been used in many subsequent works to generate the synthetic complementary datasets for examining 41 CLL algorithms [1, 9, 16, 21, 25]. Although these assumptions simplify the design and analysis of 42 CLL algorithms, it remains unknown whether these assumptions hold true in practice and whether 43 violation of these assumptions will significantly affect the performance of CLL algorithms. In 44 45 addition to the uniform or class-conditional assumptions, most existing studies implicitly assumes that the complementary labels are noise-free. That is, they do not mistakenly represent the ordinary 46 47 labels. While some studies claim to be more robust to noisy complementary labels [14], they were only tested on synthetic scenarios. It remains unclear how noisy the real-world datasets are, and how 48 such noise affects the performance of current CLL algorithms. 49

To understand how much the real-world scenario differs from the assumptions, we started by collecting 50 the datasets CLCIFAR10 and CLCIFAR20, which are derived from the famous CIFAR datasets 51 for ordinary multi-class classification [12]. Since their release in 2023, the datasets [22] have 52 been utilized by several emerging CLL studies [15, 23, 24, 26], demonstrating their instantaneous 53 impact. We continue to extend the collection and form two additional human-annotated datasets, 54 CLMicroImageNet10 and CLMicroImageNet20, which are derived from TinyImageNet200 [13, 19]. 55 The extension verifies that our observations on CIFAR-derived datasets hold true for other image 56 datasets. For all four datasets, we analyze the collected complementary labels, including their noise 57 rates and non-uniform nature. Then, we perform benchmark experiments with diverse state-of-the-art 58 CLL algorithms and conduct dataset-level ablation study on the assumptions of complementary-label 59 generation using the collected datasets. Our studies reveal annotation noise as the most influential 60 factor in the real-world datasets, and confirm that the non-uniform nature of human-annotated 61 complementary labels cause certain CLL algorithms more susceptible to overfitting. These findings 62 immediately suggest that the community focus more research efforts on developing CLL algorithms 63 64 that are robust to noisy and non-uniform complementary-label distributions. In addition, we used the collected datasets to demonstrate that existing complementary-label-only validation schemes are 65 not mature yet, suggesting the community a novel research direction for making CLL practical. Our 66 contributions are summarized as follows: 67

- We designed a collection protocol of complementary labels (CLs) for images, and verified that the protocol collects reasonable human-annotated CLs across different datasets.
- We released CLImage, the collected set of four real-world CL datasets to support the continuous research of the community, publicly released at https://github.com/ntucllab/
   CLImage\_Dataset.
- We analyzed the collected datasets with extensive benchmarking experiments, which provides novel and valuable insights for the community.

# 75 2 Preliminaries on CLL

#### 76 2.1 Complementary-label learning

In ordinary multi-class classification, a dataset  $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  that is *i.i.d.* sampled from an unknown distribution is given to the learning algorithm. For each  $i, \mathbf{x}_i \in \mathbb{R}^M$  represents the *M*-dimension feature of the *i*-th instance and  $y_i \in [K] = \{1, 2, ..., K\}$  represents the class  $\mathbf{x}_i$ belongs to. The goal of the learning algorithm is to learn a classifier from *D* that can predict the labels of unseen instances correctly. The classifier is typically parameterized by a scoring function  $\mathbf{g}: \mathbb{R}^M \to \mathbb{R}^K$ , and the prediction is made by  $\arg \max_{k \in [K]} \mathbf{g}(\mathbf{x})_k$  given an instance  $\mathbf{x}$ , where  $\mathbf{g}(\mathbf{x})_k$  denotes the k-th output of  $\mathbf{g}(\mathbf{x})$ . In contrast to ordinary multi-class classification, CLL shares the same goal of learning a classifier but trains with different labels. In CLL, the ordinary label  $y_i$  is not accessible to the learning algorithm. Instead, a complementary label  $\bar{y}_i$  is provided, which is a class that the instance  $\mathbf{x}_i$  does *not* belong to. The goal of CLL is to learn a classifier that is able to predict the correct labels of unseen instances from a complementary-label dataset  $\bar{D} = \{(\mathbf{x}_i, \bar{y}_i)\}_{i=1}^n$ .

#### 88 2.2 Common assumptions on CLL

Researchers have made some additional assumptions on the generation process of complementary 89 labels to facilitate the analysis and design of CLL algorithms. One common assumption is the 90 *class-conditional assumption* [27]. It assumes that the distribution of a complementary label only 91 depends on its ordinary label and is independent of the underlying example's feature, i.e.,  $P(\bar{y}_i \mid i)$ 92  $\mathbf{x}_i, y_i = P(\bar{y}_i \mid y_i)$  for each *i*. One special case of the class-conditional assumption is the *uniform* 93 assumption, which further specifies that the complementary labels are generated uniformly. That is, 94  $P(\bar{y}_i = k | y_i = j) = \frac{1}{K-1}$  for all  $k \in [K] \setminus \{j\}$  [8, 9, 14]. 95 For convenience, a  $K \times K$  matrix T, called *transition matrix*, is often used to represent how the 96 complementary labels are generated under the class-conditional assumption.  $T_{i,k}$  is defined to be 97 the probability of obtaining a complementary label k if the underlying ordinary label is j, i.e., 98  $T_{j,k} = P(\bar{y} = k \mid y = j)$  for each  $j, k \in [K]$ . The diagonals of T hold the conditional probabilities 99 that a complementary label mistakenly represents the ordinary label. That is, they indicate the noise 100 level of the complementary labels. When T contains all zeros on its diagonals, the CLL scenario is 101

called *noiseless*. For instance, the uniform and noiseless assumption can be represented by  $T_{j,j} = 0$ for each  $j \in [K]$  and  $T_{j,k} = \frac{1}{K-1}$  for each  $k \neq j$ . Class-conditional CLL scenarios based on any other non-uniform T are often called *biased*.

#### 105 2.3 A brief overview of CLL algorithms

The pioneering work by Ishida et al. [8] studied how to learn from complementary labels under 106 the *uniform assumption* by converting the risk estimator in ordinary multi-class classification to an 107 unbiased risk estimator (URE) in CLL [8]. URE is then found to be prone to overfitting because 108 of negative empirical risks, and is upgraded with two tricks, non-negative risk estimator (URE-NN) 109 and gradient accent (URE-GA) [9]. The surrogate complementary loss (SCL) algorithm mitigates 110 the overfitting issue of **URE** by a different loss design that decreases the variance of the empirical 111 estimation. However, these algorithms either rely on the uniform assumption in design or are only 112 tested on the synthetic datasets that obeys the uniform assumption. 113

To make CLL one step closer to practice, researchers have explored algorithms to go beyond the 114 uniform (and thus noiseless) assumption. Yu et al. [27] utilized the forward-correction loss (FWD) 115 116 to accommodate biased complementary label generation by adapting techniques from noisy label learning [18] to change the loss. Additionally, Gao and Zhang [6] proposed the L-W algorithm based 117 on discriminatively modeling the distribution of complementary labels through a weighting function, 118 further improving the performance in bias scenario. Furthermore, Ishiguro et al. [10] designed robust 119 loss functions for learning from noisy CLs, including MAE and WMAE, by applying the gradient 120 ascent technique [9] to handle noisy scenarios. 121

Besides CLL algorithms, a crucial component for making CLL practical is model validation. In ordinary-label learning, this can be done by naively calculating the classification accuracy on a validation dataset. In CLL, this scheme can be intractable if there are not enough ordinary labels. One generic way of model validation is based on the result of Ishida et al. [9] by calculating the unbiased risk estimator of the zero-one loss, i.e.,

$$\hat{R}_{01}(\mathbf{g}) = \frac{1}{N} \sum_{i=1}^{N} e_{\overline{y}_i}^{\top}(T^{-1}) \ell_{01}(\mathbf{g}(x_i))$$
(1)

where  $e_{\overline{y}_i}$  denotes the one-hot vector of  $\overline{y}_i$ ,  $\ell_{01}(\mathbf{g}(x_i))$  denotes the K-dimensional vector ( $\ell_{01}(\mathbf{g}(x_i), 1), \ldots, \ell_{01}(\mathbf{g}(x_i)), K)$ )<sup>T</sup>, and  $\ell_{01}(\mathbf{g}(x_i), k) = 0$  if  $\arg \max_{k \in [K]} \mathbf{g}(x_i) = k$  and 1 otherwise, representing the zero-one loss of  $g(x_i)$  if the ordinary label is k. This estimator will be used in the experiments in Section 6. Another validation objective, surrogate complementary esimation loss (SCEL), was proposed by Lin and Lin [14]. SCEL measures the log loss of the complementary probability estimates induced by the probability estimates on the ordinary label space. The formula to calculate SCEL is as follows,

$$\hat{R}_{\text{SCEL}}(\mathbf{g}) = \frac{1}{N} \sum_{i=1}^{N} -\log\left(e_{\overline{y}_i}^{\top} T^{\top} \operatorname{softmax}(\mathbf{g}(x_i))\right).$$
(2)

# **34 3 Construction of the CLImage collection**

In this section, we introduce the four complementary-labeled datasets that we collected, CLCIFAR10,
 CLCIFAR20, CLMicroImageNet10 and CLMicroImageNet20. All datasets are labeled by human
 annotators on Amazon Mechanical Turk (MTurk)<sup>1</sup>.

#### 138 3.1 Datasets and goals

The complementary-labeled datasets are derived from ordinary multi-class classification datasets. CIFAR10, CIFAR100 and TinyImageNet200 [12, 13, 19]. This selection is motivated by the realworld noisy label dataset by Wei et al. [25]. Building upon the CIFAR and TinyImageNet200 datasets allow us to estimate the noise rate and the empirical transition matrix easily, as they already contain nearly noise-free ordinary labels. In addition, many of the state-of-the-art CLL algorithms have been benchmarked on synthetic complementary labels with the CIFAR datasets [4, 11, 17]. Our CLCIFAR counterparts immediately allow a fair comparison to those results with the same network architecture.

In addition to our CLCIFAR extensions, we are the first to introduce (Tiny)ImageNet-derived datasets to the CLL literature. Such datasets serve two purposes. First, it allows us to confirm the validity of our collection protocol and findings beyond CIFAR-derived datasets. Second, ImageNet knowingly contains images of higher complexity than CIFAR and can thus be used to challenge the ability of existing CLL algorithms more realistically.

There is a historical note that is worth sharing with the community: We initially attempted to collect complementary labels based on the 100 classes in CIFAR100. But some preliminary testing soon revealed that state-of-the-art CLL algorithms cannot produce meaningful classifiers for 100 classes even on synthetic complementary labels that are uniformly and noiselessly generated. We thus set our collection goals to be 10-class classification, which is the focus of most current CLL studies, and 20-class classification, which extends the horizon of CLL and matches the 20 super-class structure in CIFAR.

#### **158 3.2** Complementary label collection protocol

To collect only complementary labels from the CIFAR, TinyImageNet datasets, for each image in the training split, we first randomly sample four distinct labels and ask the human annotators to select any of the *incorrect* one from them. To leave room for analyzing the annotators' behavior, each image is labeled by three different annotators. The four labels are re-sampled for each annotator on each image. That is, each annotator possibly receives a different set of four labels to choose from. An algorithmic description of the protocol is as follows. For each image x,

165 1. Uniformly sample four labels without replacement from the label set [K].

- 166 2. Ask the annotator to select any one of the complementary label  $\bar{y}$  from the four sampled 167 labels.
- 168 3. Add the pair  $(\mathbf{x}, \bar{y})$  to the complementary dataset.

<sup>&</sup>lt;sup>1</sup>https://www.mturk.com/

Note that if the annotators always select one of the correct complementary labels uniformly, the 169 empirical transition matrix will also be uniform in expectation. We will inspect the empirical transition 170 matrix in Section 4. The labeling tasks are deployed on MTurk by dividing them into smaller we first 171 divide the total images into smaller human intelligence tasks (HITs). For instance, for constructing 172 the CLCIFAR datasets, we first divide the 50,000 images into five batches of 10,000 images. Then, 173 each batch is further divided into 1,000 HITs with each HIT containing 10 images. Each HIT is 174 deployed to three annotators, who receive 0.03 dollar as the reward by annotating 10 images. To 175 make the labeling task easier and increase clarity, the size of the images are enlarged to  $200 \times 200$ 176 pixels. 177

#### 178 **4 Result analysis**

Next, we closely examine the collected complementary labels. We first analyze the error rates of the
 collected labels, and then verify whether the transition matrix is uniform or not. Finally, we end with
 an analysis on the behavior of the human annotators observed in the label collection protocol.



Figure 1: The label distribution of CLCIFAR10 and CLMicroImageNet10 datasets.

Observation 1: noise rate compared to ordinary label collection We first look at the noise rate of 182 the collected complementary labels. A complementary label is considered to be incorrect if it is actu-183 184 ally the ordinary label. The mean error rate made by the human annotators is 3.93% for CLCIFAR10, 2.80% for CLCIFAR20, 5.19% for CLMicroImageNet10 and 3.21% for CLMicroImageNet20. In 185 theory, we can estimate a random annotator achieves a noise rate of  $\frac{1}{K}$  for complementary label 186 annotation and a noise rate of  $\frac{K-1}{K}$  for ordinary label annotation. If we compare the human annotators 187 to a random annotator, then for CLCIFAR10, human annotators have 60.7% less noisy labels than 188 the random annotator whereas for CIFAR10-N, human anotators have 80% less noisy labels. This 189 demonstrates that human annotators are more competent compared to a random annotator in the 190 ordinary-label annotation. Similarly, human annotators have 44% less noise than a random annotator 191 for CLCIFAR20 and 73.05% less noise for CIFAR100N-coarse. This observation reveals that while 192 the absolute noise rate is lower in annotating complementary labels, it may be more difficult to be 193 competent against random labels than the ordinary label annotation. 194

Observation 2: imbalanced complementary label annotation Next, we analyze the distribution of 195 the collected complementary labels. The frequency of the complementary labels for the CLCIFAR10 196 and CLMicroImageNet10 (CLMIN10) datasets are reported in Figure 1. As we can see in the 197 figure, the annotators exhibit specific biases towards certain labels. For instance, in CLCIFAR10, 198 annotators prefer "airplane" and "automobile," while in CLMIN10, they prefer "pizza" and "torch". In 199 CLCIFAR10, the bias is towards labels in different categories, as vehicles ("airplane," "automobile") 200 versus animals ("cat", "bird"). In contrast, in CLMIN10, the bias is towards items that are easily 201 recognizable ("pizza" and "torch") and against those that are less familiar ("cardigan" or "alp"). 202

**Observation 3: biased transition matrix** Finally, we visualize the empirical transition matrix using the collected CLs in Figure 2. Based on the first two observations, we could imagine that the transition matrix is biased. By inspecting Figure 2, we further discover that the bias in the complementary labels are dependent on the true labels. For instance, in CLCIFAR10, despite we see more annotations on airplane and automobile in aggregate, conditioning on the transportation-related labels ("airplane",



Figure 2: The empirical transition matrices of CLCIFAR10 and CLMicroImageNet10.

"automobile", etc), the distribution of the complementary labels becomes more biased towards other
animal-related labels ("bird", "cat", etc.) Furthermore, this observation holds true on CLMIN10 as
well. Next, we study the impact of the bias and noise on existing CLL algorithms.

We discovered similar patterns in all four human-annotated datasets, validating that our design methodology is practical for collecting real-world CLL image datasets. Due to space limitations, we have included the detailed analysis of CLCIFAR20 and CLMicroImageNet20 in Appendix B.4.

# 214 **5 Experiments**

In this section, we benchmarked several state-of-the-art CLL algorithms on CLImage. A significant 215 performance gap between the models trained on the humanly annotated CLCIFAR, CLMicroImageNet 216 dataset and those trained on the synthetically generated complementary labels (CL) was observed 217 in Section 5.1, which motivates us to analyze the possible reasons for the gap with the following 218 experiments. To do so, we discuss the effect of three factors in the label generating process, feature 219 dependency, noise, and biasedness, in Section 5.2, Section 5.3, and Section 5.4, respectively. From 220 our experiment results, we conclude that noise is the dominant factor affecting the performance of 221 the CLL algorithms on CLCIFAR<sup>2</sup>. 222

#### 223 5.1 Standard benchmark on CLImage

Baseline methods Several state-of-the-art CLL algorithms were selected for this benchmark. Some 224 of them take the transition matrix T as inputs, which we call T-informed methods, including two 225 version of forward correction [27]: FWD-U and FWD-R, two version of unbiased risk estimator 226 with gradient ascent [9]: URE-GA-U and URE-GA-R, and robust loss [10] for learning from noisy 227 CL: CCE, MAE, WMAE, GCE, and  $SL^3$ . We also included some algorithms that assume the 228 transition matrix T to be uniform, called T-agnostic methods, including surrogate complementary 229 loss SCL-NL and SCL-EXP [1], discriminative modeling L-W and its weighted variant (L-UW) [6], 230 and pairwise-comparison (PC) with the sigmoid loss [8]. The details of the algorithms mentioned 231 above are discussed in Appendix D. 232

Implementation details We collected and released three CLs per image to prepare for future studies. However, for this standard benchmark, we chose the first CL from the collected labels for each data instances to form a single CLL dataset, ensuring reproducibility. Then, we trained a ResNet18 [7] model using the baseline methods mentioned above on the single CLL dataset using

<sup>&</sup>lt;sup>2</sup>Due to space and time constraints, we only provide the results and discussion on the CLCIFAR datasets.

<sup>&</sup>lt;sup>3</sup>Due to space limitations, we only provided the results of MAE. The remaining results and discussions related to the robust loss methods can be found in Appendix B.3

Table 1: Standard benchmark results on CLCIFAR/CLMicroImageNet(CLMIN) and uniform-CIFAR/ MicroImageNet(MIN) datasets. Mean accuracy ( $\pm$  standard deviation) on the testing dataset from four trials with different random seeds. Highest accuracy in each column is highlighted in bold.

	uniform-CIFAR10	uniform-CIFAR20	uniform-MIN10	uniform-MIN20	CLCIFAR10	CLCIFAR20	CLMIN10	CLMIN20
FWD-U	64.19±0.57	21.54±0.37	36.30±1.12	12.57±2.94	34.83±0.50	8.03±0.74	$23.85 \pm 2.76$	$6.33 \pm 1.04$
FWD-R	61.32±0.90	$21.50 \pm 0.38$	35.70±1.19	$14.85 \pm 1.75$	$38.13 \pm 0.88$	$20.27 \pm 0.53$	30.15±1.83	$10.60{\scriptstyle\pm0.82}$
URE-GA-U	50.24±1.11	16.67±1.35	35.70±1.97	$11.65 \pm 1.90$	$34.72 \pm 0.40$	$10.49 \pm 0.52$	22.90±2.97	$5.75 \pm 0.43$
URE-GA-R	50.73±1.83	$17.57 \pm 0.61$	$33.65 \pm 1.40$	9.78±3.88	$30.23 \pm 0.70$	$6.17 \pm 0.82$	$13.25 \pm 5.11$	$6.50 \pm 0.35$
SCL-NL	63.76±0.09	$21.37 \pm 1.18$	37.05±1.40	$13.00 \pm 2.80$	$34.77 \pm 0.60$	$8.02 \pm 0.36$ 7.70 $\pm 0.41$ 7.71 $\pm 0.35$	21.80±1.85 24.80±1.14 23.80±2.64	$6.17 \pm 0.49$
SCL-EXP	63.29±1.02	$21.57 \pm 1.13$	$36.55 \pm 1.28$	$12.95 \pm 3.38$	$35.18 \pm 0.67$			5.58±0.13 6.40±0.29
L-W	$54.32 \pm 0.41$	$19.59 \pm 0.99$	33.80±2.66	$12.70 \pm 2.35$	$32.99 \pm 1.01$			
L-UW	57.52±0.59	$20.71 \pm 0.92$	35.10±2.74	$12.12 \pm 3.13$	$34.69 \pm 0.32$	$8.15 \pm 0.30$	$22.40 \pm 1.67$	$6.35 \pm 0.86$
PC-sigmoid	$37.78 \pm 0.80$	$14.48 \pm 0.47$	$29.10 \pm 0.98$	$10.72 \pm 1.38$	$32.15 \pm 0.80$	$12.11 \pm 0.46$	$23.15 \pm 0.46$	$6.90 \pm 1.04$
ROB-MAE	$59.38 \pm 0.63$	$18.17 \pm 1.31$	$31.50 \pm 1.81$	$6.35{\pm}0.86$	$20.23 \pm 1.02$	$5.40{\pm}0.59$	$14.15{\scriptstyle\pm0.68}$	$5.38{\scriptstyle\pm0.33}$
	CIFA	CIFAR10		AR20	MI	N10	MIN20	
standard supervision	82.80	)±0.28	63.80	<b>)</b> ±0.49	68.70	€±1.53	63.90±1.00	

the Adam optimizer for 300 epochs without learning rate scheduling. The weight decay was fixed 237 at  $10^{-4}$  and the batch size was set to 512. The experiments were run with Tesla V100-SXM2. For 238 better generalization, we applied standard data augmentation technique, RandomHorizontalFlip, 239 RandomCrop, and normalization to each image. The learning rate was selected from  $\{10^{-3}, 5 \times 10^{-4}, 5 \times$ 240  $10^{-4}$ ,  $5 \times 10^{-5}$ ,  $10^{-5}$ } using a 10% hold-out validation set. We selected the learning rate with the 241 best classification accuracy on the validation dataset. Note that here we assumed the ordinary labels 242 in the validation dataset are known. We will discuss other validation objectives that rely only on 243 complementary labels in Section 6. As CLL algorithms are prone to overfitting [1, 9], some previous 244 works did not use the model after training for evaluation. Instead, previous works were performed by 245 evaluating the model on the validation dataset and selecting the epoch with the highest validation 246 accuracy. In this work, we also follow the same aforementioned technique to validate testing set. For 247 reference, we also performed the experiments on synthetically-generated CLL dataset, where the CLs 248 were generated uniformly and noiselessly, denoted uniform-CIFAR. 249

**Results and discussion** As we can observe in Table 1, there is a significant performance gap between the humanly annotated dataset, CLCIFAR, and the synthetically generated dataset, uniform-CIFAR. The difference between the two datasets can be divided into three parts: (a) whether the generation of complementary labels depends on the feature, (b) whether there is noise, and (c) whether the complementary labels are generated with bias. A negative answer to those questions simplify the problem of CLL. We can gradually simplify CLCIFAR to uniform-CIFAR by chaining those assumptions as follows <sup>4</sup>:



In the following subsections, we will analyze how these three factors affect the performance of the CLL algorithms.

#### 259 5.2 Feature dependency

In this experiment, we verified whether the performance gap resulted from the feature-dependent generation of practical CLs. Conceivably, even if two images belong to the same class, the distribution on the complementary labels could be different. On the other hand, the distributional difference could also be too small to affect model performance, e.g., if  $P(\bar{y} \mid y, \mathbf{x}) \approx P(\bar{y} \mid y)$  for most  $\mathbf{x}$ . Consequently, we decided to further look into whether this assumption can explain the performance gap. To observe the effects of approximating  $P(\bar{y} \mid y, \mathbf{x})$  with  $P(\bar{y} \mid y)$ , we generated two synthetic

<sup>&</sup>lt;sup>3</sup>Note that FWD-R and URE-GA-R assume the empirical transition matrix  $T_e$  to be provided. The empirical transition matrix is computed from the labels in the training set, so it is slightly different from a uniform transition matrix  $T_u$  in the uniform-CIFAR datasets. As a result, the performances of FWD-R and URE-GA-R do not exactly match those of FWD-U and URE-GA-U, respectively, in the uniform-CIFAR datasets.

<sup>&</sup>lt;sup>4</sup>The "interpolation" between CLCIFAR and uniform-CIFAR does not necessarily have to be this way. For instance, one can remove the biasedness before removing the noise. We chose this order to reflect the advance of CLL algorithms. First, researchers address the uniform case [8], then generalize to the biased case [27], then consider noisy labels [10]. There is no work considering feature-dependent complementary labels yet.

complementary datasets, CLCIFAR10-*iid* and CLCIFAR20-*iid* by i.i.d. sampling CLs from the empirical transition matrix in CLCIFAR10 and CLCIFAR20, respectively. We proceeded to benchmark
 the CLL algorithms on CLCIFAR-*iid* and presented the accuracy difference compared to CLCIFAR
 in Table 2.

**Results and discussion** From Table 2, we observed that the accuracy barely changes on the resampled

271 CLCIFAR-*iid*, suggesting that even if the complementary labels in CLCIFAR could be feature-

dependent, this dependency does not affect the model performance significantly. Hence, there might

<sup>273</sup> be other factors contributing to the performance gap.

Table 2: Mean accuracy difference ( $\pm$  standard deviation) of different CLL algorithms. A plus indicates the performance on is calculated as CLCIFAR-*i.i.d.* accuracy minus CLCIFAR accuracy.

	FWD-U	FWD-R	URE-GA-U	URE-GA-R	SCL-NL	SCL-EXP	L-W	L-UW	PC-sigmoid
CLCIFAR10-iid CLCIFAR20-iid	-1.1±2.17 -0.64±0.39	$\substack{-0.36 \pm 1.15 \\ -3.53 \pm 1.13}$	$-3.03 \pm 1.25$ $-0.37 \pm 0.51$	$\substack{0.74 \pm 0.35 \\ 1.79 \pm 2.34}$	$\substack{-0.67 \pm 1.81 \\ -0.28 \pm 0.61}$	$-1.97{\pm}1.16$ $-0.39{\pm}0.69$	-2.5±0.56 -0.5±1.37	$-3.53 \pm 1.36$ $-0.82 \pm 0.04$	$-2.03\pm2.05$ $-2.24\pm0.52$

#### 274 5.3 Labeling noise

In this experiment, we further investigated the impact of the label noise on the performance gap. Specifically, we measured the accuracy on the noise-removed versions of CLCIFAR datasets, where varying percentages (0%, 25%, 50%, 75%, or 100%) of noisy labels are eliminated.

Results and discussion We present the performance of FWD trained on the noise-removed CLCIFAR10 dataset in the left figure in Figure 3. The results for other algorithms and the noise-removed
CLCIFAR20 dataset can be found in Appendix E. From the figure, we observe a strong positive
correlation between the performance and the proportion of removed noisy labels. When more noisy
labels are removed, the performance gap diminishes and the accuracy approaches that of the ideal
uniform-CLFAR dataset. Therefore, we conclude that the performance gap between the humanly
annotated CLs and the synthetically generated CLs are primarily attributed to the label noise.

#### 285 5.4 Biasedness of complementary labels

To further study the biasedness of CL as a potential factor contributing to the performance gap, we removed the biasedness from the noise-removed CLCIFAR dataset and examined the resulting accuracy. Specifically, we introduced the same level of uniform noise in uniform-CIFAR dataset and reevaluated the performance of FWD algorithms.

**Results and discussion** The striking similarity between the two curves in the right figure in Figure 3 shows that the accuracy is significantly influenced by label noise, while the biasedness of CL has a negligible impact on the results. Furthermore, we observe that the accuracy difference between the results of the last epoch and the best accuracy of validation set (or early-stopping: **ES**) results becomes smaller when the model is trained on the uniformly generated CLs. That is, the *T*-informed methods are more prone to overfitting when there is a bias in the CL generation.

With the experiment results in Section 5.2, 5.3, and 5.4, we can conclude that the performance gap between humanly annotated CL and synthetically generated CL is primarily attributed to label noise. Additionally, the biasedness of CLs may potentially contribute to overfitting, while the featuredependent CLs do not detrimentally affect performance empirically. It is worth noting that in the last row of Table 1, the MAE methods that can learn from noisy CL fails to generalize well in the practical dataset. These results suggest that more research on learning with noisy complementary labels can potentially make CLL more realistic.

# **303 6 Validation Objectives**

Validation is a crucial component in applying CLL algorithms in practice. With the collection of the real-world datasets, we are now able to estimate the difference between using ordinary labels for validation (the common practice in existing CLL studies, as what we do in Section 5) and using complementary labels for validation.



Figure 3: Accuracy of FWD-U and FWD-R on the noise-removed CLCIFAR10 dataset (Left) and the uniform-CIFAR10 dataset with uniform noise (**Right**) at varying noise rates.

Table 3: The testing accuracy of models evaluated with URE and SCEL.

	CLCIFAR10				CLCIFAR20				CLMIN10			CLMIN20				
	URE	SCEL	valid acc	gap $(\downarrow)$	URE	SCEL	valid acc	gap (↓)	URE	SCEL	valid acc	gap (↓)	URE	SCEL	valid acc	gap (↓)
FWD-U	$33.13 \pm 1.30$	$31.86{\scriptstyle\pm1.52}$	$34.83{\scriptstyle\pm0.50}$	1.70	$6.70 \pm 0.46$	$7.10{\scriptstyle\pm0.48}$	$8.03{\scriptstyle \pm 0.74}$	0.93	20.75±2.12	$20.20{\scriptstyle\pm0.72}$	$23.85{\scriptstyle\pm2.76}$	3.10	4.97±0.72	$4.55{\scriptstyle\pm 0.81}$	$6.33{\pm}1.04$	1.35
FWD-R	33.70±3.38	35.64±1.37	38.13±0.88	2.49	17.35±2.32	18.40±1.56	20.27±0.53	1.86	22.15±4.15	29.15±1.93	30.15±1.83	1.00	8.60±1.32	9.90±1.19	$10.60 \pm 0.82$	0.70
URE-GA-U	30.45±3.58	$33.21 \pm 1.12$	$34.72 \pm 0.40$	1.51	7.03±0.61	$8.71 \pm 0.74$	$10.49 \pm 0.52$	1.79	17.05±3.35	$21.30 \pm 3.01$	22.90±2.97	1.60	4.27±0.80	$5.03 \pm 0.48$	$5.75 \pm 0.43$	0.72
URE-GA-R	27.39±1.89	$28.32 \pm 1.38$	$30.23 \pm 0.70$	1.91	$3.58 \pm 0.47$	$5.42 \pm 0.96$	$6.17 \pm 0.82$	0.75	8.90±1.03	$10.30 \pm 1.53$	13.25±5.11	2.95	5.15±0.62	$5.57 \pm 1.54$	$6.50 \pm 0.35$	0.93
SCL-NL	$33.55 \pm 0.79$	33.70±1.33	$34.77 \pm 0.60$	1.07	6.73±0.51	$7.47 \pm 0.56$	$8.02 \pm 0.36$	0.55	19.55±1.37	22.15±1.76	$21.80 \pm 1.85$	-0.35	4.83±1.12	$5.20 \pm 0.51$	$6.17 \pm 0.49$	0.98
SCL-EXP	$31.30 \pm 2.62$	$33.47 \pm 1.16$	$35.18 \pm 0.67$	1.71	$6.83 \pm 0.23$	$7.03 \pm 0.62$	$7.70 \pm 0.41$	0.66	$18.35 \pm 1.60$	$20.65 \pm 1.39$	$24.80 \pm 1.14$	4.15	5.05±0.56	$4.45{\scriptstyle \pm 0.74}$	$5.58 \pm 0.13$	0.52
L-W	$27.49 \pm 4.30$	$30.32 \pm 2.40$	$32.99 \pm 1.01$	2.67	5.90±0.29	$7.18 \pm 0.31$	$7.71 \pm 0.35$	0.53	19.30±4.66	$18.95 \pm 2.30$	$23.80 \pm 2.64$	4.50	5.97±0.33	$5.55 \pm 0.17$	$6.40 \pm 0.29$	0.43
L-UW	$28.90 \pm 2.01$	29.78±2.69	$34.69 \pm 0.32$	4.91	$6.40 \pm 0.42$	$8.16 \pm 0.30$	$8.15 \pm 0.30$	-0.01	$18.25 \pm 4.31$	$19.80 \pm 1.61$	$22.40 \pm 1.67$	2.60	5.82±0.77	$6.48 \pm 1.03$	$6.35 \pm 0.86$	-0.13
PC-sigmoid	24.83±5.94	$31.48 \pm 1.93$	$32.15 \pm 0.80$	0.67	7.98±2.47	$10.59 \pm 0.87$	$12.11 \pm 0.46$	1.51	12.55±1.31	$17.85 \pm 4.61$	$23.15 \pm 0.46$	5.30	6.40±1.19	5.33±1.28	$6.90 \pm 1.04$	0.50
ROB-MAE	$18.80 \pm 1.64$	$18.75 \pm 0.99$	$20.23 \pm 1.02$	1.43	$4.70 \pm 0.43$	$4.87 \pm 0.32$	$5.40 \pm 0.59$	0.53	$11.80 \pm 2.92$	$14.35 \pm 1.59$	$14.15 \pm 0.68$	-0.20	5.08±0.44	$4.62 \pm 0.66$	$5.38 \pm 0.33$	0.30

Validation objectives As discussed in Section 2, validating the model performance solely with 308 complementary labels poses a non-trivial challenge. To the best of our knowledge, only two existing 309 CLL studies offer some possibility to evaluate a classifier with only complementary labels. They are 310 URE [9] and SCEL [14]. We take these two validation objectives to select the optimal learning rate 311 from  $\{10^{-3}, 5 \times 10^{-4}, 10^{-4}, 5 \times 10^{-5}, 10^{-5}\}$  and provides the accuracy on testing set in Table 3. 312 We compare the result to another validation objective that computes the accuracy on an equal number 313 of *ordinary labels*. Our goal was to determine the gap between using complementary labels and 314 ordinary labels for validation. We selected the best learning rate based on the validation objectives 315 for URE, SCEL, and ordinary-label accuracy, and then report the test performance, as shown in Table 316 3 for real-world datasets and Table 4 in the Appendix for synthetic datasets. 317

Results and discussion Firstly, there appears no clear winner between URE and SCEL, both using 318 only CLs for validation. Validating with the ordinary-label accuracy generally provides stronger 319 performance than URE/SCEL, and the test performance gap between validating with ordinary labels 320 and validating with complementary labels can be as big as nearly 5%. These findings suggest that 321 using purely complementary labels for validation, whether through URE or SCEL, still suffers from a 322 non-negligible performance drop compared to using ordinary validation. That is, the numbers reported 323 in existing studies, which validates with ordinal labels, can be optimistic for practice. Whether this 324 gap can be further reduced remains an open research problem and the community can pay more 325 326 attention on that to make CLL more practical.

# 327 7 Conclusion

In this paper, we devised a protocol to collect complementary labels from human annotators. Utilizing 328 this protocol, we curated four real-world datasets, CLCIFAR10, CLCIFAR20, CLMicroImageNet10, 329 and CLMicroImageNet20 and made them publicly available to the research community. Through 330 our meticulous analysis of these datasets, we confirmed the presence of noise and bias in the human-331 annotated complementary labels, challenging some of the underlying assumptions of existing CLL 332 algorithms. Extensive benchmarking experiments revealed that noise is a critical factor that under-333 mines the effectiveness of most existing CLL algorithms. Furthermore, the biased complementary 334 labels can trigger overfitting, even for algorithms explicitly designed to leverage this bias information. 335 In addition, our study on the validation objective for CLL suggests that validating with only com-336 plementary labels causes significant performance degrading. These findings emphasize the need for 337 the community to dedicate more effort on those issues. The curated datasets pave the way for the 338 community to create more practical and applicable CLL solutions. 339

#### 340 **References**

- [1] Y.-T. Chou, G. Niu, H.-T. Lin, and M. Sugiyama. Unbiased risk estimators can mislead: A case study of learning with complementary labels, 2020.
- [2] T. Cour, B. Sapp, and B. Taskar. Learning from partial labels. *The Journal of Machine Learning Research*, 12:1501–1536, 2011.
- [3] F. Denis. Pac learning from positive statistical queries. In *Algorithmic Learning Theory: 9th International Conference, ALT'98 Otzenhausen, Germany, October 8–10, 1998 Proceedings 9*, pages 112–126. Springer, 1998.
- [4] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
   M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
   image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [5] B. Frénay and M. Verleysen. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems*, 25(5):845–869, 2013.
- [6] Y. Gao and M.-L. Zhang. Discriminative complementary-label learning with weighted loss. In
   *International Conference on Machine Learning*, pages 3587–3597. PMLR, 2021.
- [7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. 2016.
- [8] T. Ishida, G. Niu, W. Hu, and M. Sugiyama. Learning from complementary labels. *Advances in neural information processing systems*, 30, 2017.
- [9] T. Ishida, G. Niu, A. K. Menon, and M. Sugiyama. Complementary-label learning for arbitrary
   losses and models, 2019.
- [10] H. Ishiguro, T. Ishida, and M. Sugiyama. Learning from noisy complementary labels with
   robust loss functions. *IEICE TRANSACTIONS on Information and Systems*, 105(2):364–376,
   2022.
- [11] A. Kolesnikov, L. Beyer, X. Zhai, J. Puigcerver, J. Yung, S. Gelly, and N. Houlsby. Big
   transfer (bit): General visual representation learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages
   491–507. Springer, 2020.
- [12] A. Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [13] Y. Le and X. S. Yang. Tiny imagenet visual recognition challenge. 2015.
- [14] W.-I. Lin and H.-T. Lin. Reduction from complementary-label learning to probability estimates.
   In *Proceedings of the Pacific-Asia Conference on Knowledge Discovery and Data Mining* (*PAKDD*), May 2023.
- W.-Y. Lin. Reduction from complementary-label learning to probability estimates. Master's thesis, 2023.
- [16] S. Liu, Y. Cao, Q. Zhang, L. Feng, and B. An. Consistent complementary-label learning via
   order-preserving losses. In *International Conference on Artificial Intelligence and Statistics*,
   pages 8734–8748. PMLR, 2023.
- [17] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza,
   F. Massa, A. El-Nouby, et al. Dinov2: Learning robust visual features without supervision.
   *arXiv preprint arXiv:2304.07193*, 2023.
- [18] G. Patrini, A. Rozza, A. Menon, R. Nock, and L. Qu. Making deep neural networks robust to
   label noise: a loss correction approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017.
- [19] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy,
   A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition
   challenge. *Int. J. Comput. Vision*, 115(3), 2015.

- [20] M. Sugiyama, H. Bao, T. Ishida, N. Lu, T. Sakai, and G. Niu. *Machine learning from weak supervision: An empirical risk minimization approach.* MIT Press, 2022.
- [21] D.-B. Wang, L. Feng, and M.-L. Zhang. Learning from complementary labels via partial-output consistency regularization. In *IJCAI*, pages 3075–3081, 2021.
- [22] H.-H. Wang, W.-I. Lin, and H.-T. Lin. Clcifar: Cifar-derived benchmark datasets with human
   annotated complementary labels, 2023.
- [23] W. Wang, T. Ishida, Y.-J. Zhang, G. Niu, and M. Sugiyama. Learning with complementary
   labels revisited: The selected-completely-at-random setting is more practical.
- W. Wang, T. Ishida, Y.-J. Zhang, G. Niu, and M. Sugiyama. Learning with complementary labels revisited: A consistent approach via negative-unlabeled learning. *arXiv preprint* arXiv:2311.15502, 2023.
- [25] J. Wei, Z. Zhu, H. Cheng, T. Liu, G. Niu, and Y. Liu. Learning with noisy labels revisited: A
   study using real-world human annotations, 2022.
- [26] Y. You, J. Huang, B. Wang, and Q. Tong. Rethinking one-vs-the-rest loss for instance-dependent
   complementary label learning.
- <sup>402</sup> [27] X. Yu, T. Liu, M. Gong, and D. Tao. Learning with biased complementary labels, 2018.
- [28] Z.-H. Zhou. A brief introduction to weakly supervised learning. *National science review*, 5(1):
   44–53, 2018.

# 405 Checklist

406	1. For all authors
407 408	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
409 410	(b) Did you describe the limitations of your work? [Yes] We describe some potential direction to make CLL algorithms more practical in the conclusion section.
411	(c) Did you discuss any potential negative societal impacts of your work? [Yes]
412 413	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
414	2. If you are including theoretical results
415	(a) Did you state the full set of assumptions of all theoretical results? [N/A]
416	(b) Did you include complete proofs of all theoretical results? [N/A]
417	3. If you ran experiments (e.g. for benchmarks)
418 419 420	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)? [Yes] Please see the link in Section 1 or Appendix J.
421 422	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] Please see the "Implementation Details" paragraph in Section 5.1
423 424 425	(c) Did you report error bars (e.g., with respect to the random seed after running experi- ments multiple times)? [Yes] The standard deviation of four trials is indicated after ± sign in Table 1
426 427	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We mentioned them in Section 5.1.
428	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
429 430	(a) If your work uses existing assets, did you cite the creators? [Yes] Please refer to the footnote 1
431 432	(b) Did you mention the license of the assets? [N/A] The authors of the CIFAR, TinyImageNet datasets require the paper using them cite their paper, which we did.
433 434	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] The reference is at 1
435 436	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A]
437 438	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
439	5. If you used crowdsourcing or conducted research with human subjects
440 441	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [Yes] We provide the screenshots of the mTurk interface at our github repo.
442 443	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
444 445 446	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [Yes] We descirbe how we paid the annonators in Section 3.2.