## A. Algorithm of the proposed method RETINA

We provide the details of our proposed method RETINA in Algorithm 1.

---

**Algorithm 1** The training of RETINA

---

1: **procedure** MULTI-TEACHER STUDENT TRAINING($\{\mathcal{D}^{(m)}\}_{m=1}^{M}, \lambda_u, n_{\text{warm\_up}}, n_e$)
2:    $\triangleright \mathcal{D}^{(m)}$: the noisy label dataset of annotator $m$                                  $\triangleleft$
3:    $\triangleright \lambda_u$: a hyper-parameter that weights the loss of noisy label samples              $\triangleleft$
4:    $\triangleright n_{\text{warm\_up}}$: the number of warmup epochs                             $\triangleleft$
5:    $\triangleright n_{\text{epoch}}$: the number of training epochs                              $\triangleleft$
6:    initialize $M$ model parameters: $\{\theta^{(m)}\}_{m=1}^{M}$
7:    warm-up on noisy datasets: $\theta^{(m)} \leftarrow \text{WARM-UP}(\mathcal{D}^{(m)}, \theta^{(m)}, n_{\text{warm\_up}}), \forall m \in \{1, \ldots, M\}$
8:    **for** epoch $= n_{\text{warm\_up\_epoch}} + 1 : n_{\text{epoch}}$ **do**
9:        **for** $m = 1 : M$ **do**
10:           $f_{\theta^{(m)}} \sim_{\text{w/o}} \mathcal{F}_s(\{f_{\theta^{(m)}}\}_{m=1}^{M}, \{\mathcal{D}^{(m)}\}_{m=1}^{M})$        $\triangleright$ *Sample a student without replacement*
11:           $f_{\theta^{(n)}} \sim \mathcal{F}_t(\{f_{\theta^{(n)}}\}_{n=1}^{M}, \{\mathcal{D}^{(m)}\}_{m=1}^{M}), \text{student} = f_{\theta^{(m)}}, n \neq m)$        $\triangleright$ *Select the teacher*
12:           $\mathcal{D}_{\text{clean}}^{(m)}, \mathcal{D}_{\text{noisy}}^{(m)} \leftarrow \text{SAMPLE-SELECTION}(f_{\theta^{(n)}}, \mathcal{D}^{(m)})$           $\triangleright$ *Eq. (6)*
13:           $\mathsf{L} = \ell_{\text{CLEAN}}(\mathcal{D}_{\text{clean}}^{(m)}, \theta^{(m)}) + \lambda_u \ell_{\text{NOISE}}(\mathcal{D}_{\text{noisy}}^{(m)}, \theta^{(m)}) + \lambda_r \ell_{\text{REG}}$        $\triangleright$ *loss to train the student model*
14:           $\theta^{(m)} \leftarrow \mathsf{SGD}(\mathsf{L}, \theta^{(m)})$        $\triangleright$ *train student model and update model parameters*
15:    **return** $\{\theta^{(m)}\}_{m=1}^{M}$

---

## B. Detailed experiment setting

**Noise ratios in different datasets**　Tab. 3 presents the noise rates (in percentage) for both individual annotators and their aggregated majority vote labels across several datasets evaluated in our experiments:

| № annotator | Noise rates of individual annotators and majority vote labels (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | a1 | a2 | a3 | a4 | a5 | m1∼2 | m1∼3 | m1∼4 | m1∼5 |
| CIFAR100-IDN30 | 30.07 | 30.11 | 30.26 | 29.91 | 30.28 | 25.16 | 13.63 | 5.82 | 2.46 |
| CIFAR100-IDN50 | 50.3 | 49.8 | 50.11 | 49.91 | 49.94 | 45.10 | 33.64 | 23.06 | 15.29 |
| CIFAR100-IDN70 | 70.05 | 70.13 | 70.19 | 70.04 | 69.76 | 65.71 | 59.47 | 51.62 | 44.12 |
| dopanim | 39.38 | 39.54 | 40.95 | 40.12 | - | 38.85 | 31.82 | 29.30 | - |
| Flickr_LDL | 62.32 | 58.94 | 58.10 | 55.23 | 50.11 | 61.09 | 53.18 | 47.87 | 41.99 |
| Chaoyang | 13.70 | 18.68 | 0.5 | - | - | 18.01 | 0.2 | - | - |

Table 3. Noise ratios(%) of annotators and their majority votes labels on three synthesized and real-world datasets. The № annotator starts with 'a' means the individual annotator in the dataset, the number after 'a' is the serial number of individual annotators, while the 'm' in m1 ∼ n(n ∈ 2, 3, 4, 5) indicates the noise rates of majority vote labels derived from individual annotator a1, a2, ..., an. For the dopanim and Chaoyang datasets, '-' means that the corresponding items do not exist since the annotators in distinct data versions are different.

**Existing multi-rater methods**　We compare RETINA with the following state-of-the-art/baseline multi-rater methods: 1) **Ensemble** an algorithm that averages the outputs of several single classifiers, each classifier being trained by a regular classifier using the noisy labels from a single annotator; 2) **Majority Vote** which trains one model with a regular classifier using the majority voting labels from all annotators; 3) **CrowdLayer** [36], an end-to-end algorithm that directly trains deep learning models from the noisy labels of multiple annotators using only backpropagation; 4) **FDS** [37] proposed an EM-based algorithm to predict the aggregated consensus labels, then one classifier is trained based on these labels, 5) **Trace-reg** [39] proposed an approach that simultaneously estimates individual annotator reliability through confusion matrices and learns the true label distribution from noisy annotations by incorporating a trace regularization term, 6) **CrowdLab** [17], a two-step algorithm that firstly estimates consensus labels for data examples by aggregating the individual annotations, then a classifier is trained on these consensus labels; 7) **UnionNet** [44], an end-to-end model that maximizes the likelihood of the union of one-hot encoded vectors of labels provided by all annotators with the help of a parametric transition matrix; 8) **Conal** [10], an end-to-end learning solution with two parallel noise adaptation layers that decompose crowdsourced annotation noise into shared confusions across annotators and annotator-specific confusions, 9) **MaDL** [23], an end-to-end algorithm that jointly trains a ground-truth model and an annotator model by presenting a probabilistic training framework; 10) **CrowdAR** [5], an end-to-end algorithm that models the reliability of annotators and is then further used to construct a soft annotation for training; 11) **GeoCrowdNet** [25], an end-to-end system that learns the label correction mechanism and the neural classifer simultaneously; and 12) **Annot-Mix** [24], an algorithm that maximizes the marginal likelihood of observed noisy class labels during the joint training of a classification and an annotator model, thus separating the noise from the true labels.

# C. Setting of Experiments

| | Backbone | WarmUp Epochs | Epochs | Optimizer | LR Scheduler | Batch Size | Initial_LR |
|---|---|---|---|---|---|---|---|
| RETINA (DivideMix) | | | 300 | | | 128 | 0.02 |
| RETINA (ProMix) | PreAct-ResNet-18 | 30 | 600 | SGD | Cosine Annealing | 256 | 0.05 |
| RETINE (Anne) | | | 300 | | | 128 | 0.02 |

Table 4. Experimental setting of proposed methods on Cifar100-IDN datasets.

| | Backbone | WarmUp Epochs | Epochs | Optimizer | LR Scheduler | Batch Size | Initial_LR |
|---|---|---|---|---|---|---|---|
| RETINA (DivideMix) | | | 50 | | | | |
| RETINA (ProMix) | Pretrained DINO-V2 | 1 | 50 | Adam | Cosine Annealing | 64 | 0.02 |
| RETINA (Anne) | | | 50 | | | | |

Table 5. Experimental setting of proposed methods on dopanim dataset.

| | Backbone | WarmUp Epochs | Epochs | Optimizer | LR Scheduler | Batch Size | Initial_LR |
|---|---|---|---|---|---|---|---|
| RETINA (DivideMix) | | | 100 | | | 64 | |
| RETINA (ProMix) | ResNet-18 | 10 | 100 | SGD | Cosine Annealing | 128 | 0.02 |
| RETINA (Anne) | | | 100 | | | 64 | |

Table 6. Experimental setting of proposed methods on Flickr_LDL and Chaoyang.

# D. Results

We report the accuracy results of our proposed algorithm, RETINA, on the three synthesized CIFAR100-IDN variants datasets and three real-world datasets, dopanim, Flickr_LDL, and Chaoyang. The experimental results of real-world datasets are as shown below. The bold font indicates the highest accuracy.

| № annotators | Test Accuracy (%) | | |
|---|---|---|---|
| | **2** | **3** | **4** |
| Ensemble | $50.74 \pm 0.30$ | $51.68 \pm 0.28$ | $52.81 \pm 0.25$ |
| Majority Vote | $50.27 \pm 0.32$ | $66.98 \pm 0.27$ | $77.87 \pm 0.31$ |
| Crowdlayer [36] | $44.89 \pm 0.23$ | $59.05 \pm 2.14$ | $68.05 \pm 3.93$ |
| FDS [37] | $64.22 \pm 0.35$ | $75.93 \pm 0.23$ | $76.23 \pm 0.21$ |
| Trace-reg [39] | $46.06 \pm 0.55$ | $61.72 \pm 0.63$ | $74.35 \pm 0.55$ |
| CrowdLab [17] | $48.89 \pm 0.23$ | $65.05 \pm 0.64$ | $75.67 \pm 0.39$ |
| Conal [10] | $46.11 \pm 0.39$ | $62.03 \pm 1.10$ | $75.69 \pm 0.14$ |
| UnionNet [44] | $49.82 \pm 0.25$ | $62.77 \pm 0.80$ | $74.17 \pm 0.19$ |
| MaDL [23] | $47.76 \pm 0.61$ | $64.06 \pm 1.33$ | $75.24 \pm 1.08$ |
| CrowdAR [5] | $45.76 \pm 0.14$ | $61.51 \pm 0.76$ | $73.79 \pm 0.31$ |
| GeoCrowdNet [25] | $50.26 \pm 0.27$ | $63.41 \pm 0.55$ | $74.31 \pm 0.68$ |
| Annot-Mix [24] | $51.56 \pm 0.13$ | $67.41 \pm 0.68$ | $78.30 \pm 0.21$ |
| RETINA (DivideMix) | $51.80 \pm 0.21$ | $68.81 \pm 0.17$ | $79.31 \pm 0.14$ |
| RETINA (ProMix) | $53.15 \pm 0.15$ | $70.27 \pm 0.16$ | $80.65 \pm 0.15$ |
| RETINA (ANNE) | $\mathbf{55.46 \pm 0.16}$ | $\mathbf{71.08 \pm 0.13}$ | $\mathbf{81.77 \pm 0.12}$ |

Table 7. Test accuracy (%) on dopanim dataset.

| № annotators | Test Accuracy (%) | | | |
|---|---|---|---|---|
| | **2** | **3** | **4** | **5** |
| Ensemble | $40.37 \pm 0.55$ | $44.41 \pm 0.64$ | $47.17 \pm 0.38$ | $48.02 \pm 0.40$ |
| Majority Vote | $37.39 \pm 0.53$ | $41.08 \pm 0.51$ | $42.22 \pm 0.37$ | $46.88 \pm 0.34$ |
| Crowdlayer [36] | $52.46 \pm 0.39$ | $53.87 \pm 0.55$ | $56.87 \pm 0.79$ | $59.77 \pm 0.90$ |
| FDS [37] | $42.71 \pm 0.51$ | $49.57 \pm 0.43$ | $57.22 \pm 0.51$ | $58.14 \pm 0.46$ |
| Trace-reg [39] | $47.73 \pm 0.43$ | $49.69 \pm 0.59$ | $51.39 \pm 0.33$ | $53.80 \pm 0.36$ |
| CrowdLab [17] | $44.22 \pm 0.76$ | $50.18 \pm 0.43$ | $54.76 \pm 0.31$ | $58.21 \pm 0.20$ |
| Conal [10] | $48.89 \pm 0.59$ | $50.11 \pm 0.22$ | $52.46 \pm 0.66$ | $54.58 \pm 1.02$ |
| UnionNet [44] | $4.41 \pm 1.77$ | $6.31 \pm 0.75$ | $10.52 \pm 0.83$ | $11.60 \pm 2.21$ |
| MaDL [23] | $47.45 \pm 0.28$ | $49.83 \pm 1.64$ | $52.17 \pm 0.82$ | $54.89 \pm 1.18$ |
| CrowdAR [5] | $50.12 \pm 0.53$ | $50.33 \pm 0.91$ | $52.46 \pm 0.90$ | $56.26 \pm 0.11$ |
| GeoCrowdNet [25] | $51.23 \pm 0.46$ | $53.31 \pm 0.66$ | $55.31 \pm 0.31$ | $58.24 \pm 0.22$ |
| Annot-Mix [24] | $50.57 \pm 1.07$ | $53.02 \pm 0.76$ | $55.73 \pm 0.74$ | $58.76 \pm 0.10$ |
| RETINA (DivideMix) | $57.36 \pm 0.44$ | $58.28 \pm 0.41$ | $59.12 \pm 0.20$ | $60.76 \pm 0.11$ |
| RETINA (ProMix) | $58.64 \pm 0.39$ | $59.95 \pm 0.26$ | $60.24 \pm 0.14$ | $61.37 \pm 0.10$ |
| RETINA (ANNE) | $\mathbf{60.05 \pm 0.25}$ | $\mathbf{61.04 \pm 0.26}$ | $\mathbf{61.46 \pm 0.18}$ | $\mathbf{63.45 \pm 0.14}$ |

Table 8. Test accuracy (%) on Flickr_LDL dataset.

| Nº annotators | Test Accuracy (%) | |
| --- | --- | --- |
| | **2** | **3** |
| Ensemble | 82.70± 0.22 | 83.22 ± 0.18 |
| Majority Vote | 75.88 ± 0.31 | 83.09 ± 0.17 |
| Crowdlayer [36] | 74.16 ± 0.41 | 81.39 ± 0.29 |
| FDS [37] | 82.00 ± 0.20 | 80.22 ± 0.19 |
| Trace-reg [39] | 80.69 ± 0.27 | 84.32 ± 0.16 |
| CrowdLab [17] | 81.45 ± 0.24 | 83.26 ± 0.18 |
| Conal [10] | 81.49 ± 0.28 | 75.33 ± 0.35 |
| UnionNet [44] | 80.27 ± 0.30 | 81.81 ± 0.22 |
| MaDL [23] | 80.78 ± 0.23 | 83.37 ± 0.17 |
| CrowdAR [5] | 79.15 ± 0.32 | 81.76 ± 0.20 |
| GeoCrowdNet [25] | 78.82 ± 0.36 | 83.77 ± 0.29 |
| Annot-Mix [24] | 79.15 ± 0.18 | 81.90 ± 0.13 |
| RETINA (DivideMix) | 83.18 ± 0.13 | 84.89 ± 0.09 |
| RETINA (ProMix) | 83.59 ± 0.11 | 85.09 ± 0.08 |
| RETINA (ANNE) | **83.42 ± 0.10** | **85.61 ± 0.07** |

Table 9. Test accuracy (%) on Chaoyang dataset.

# E. Ablation Study

| № annotators | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| **CIFAR100-IDN70** | | | | |
| Biased Selection | $72.65 \pm 0.64$ | $74.56 \pm 0.41$ | $\mathbf{75.38 \pm 0.25}$ | $\mathbf{75.69 \pm 0.24}$ |
| Random Selection | $\mathbf{73.00 \pm 0.65}$ | $\mathbf{74.63 \pm 0.46}$ | $75.03 \pm 0.44$ | $75.41 \pm 0.37$ |
| **dopanim** | | | | |
| Biased Selection | $51.32 \pm 0.27$ | $67.62 \pm 0.46$ | $\mathbf{79.32 \pm 0.22}$ | - |
| Random Selection | $\mathbf{51.80 \pm 0.21}$ | $\mathbf{68.81 \pm 0.17}$ | $79.31 \pm 0.10$ | - |

Table 10. Test accuracy (%) on different teacher-student selection algorithms based on: (top) CIFAR100-IDN70, and (bottom) dopanim.