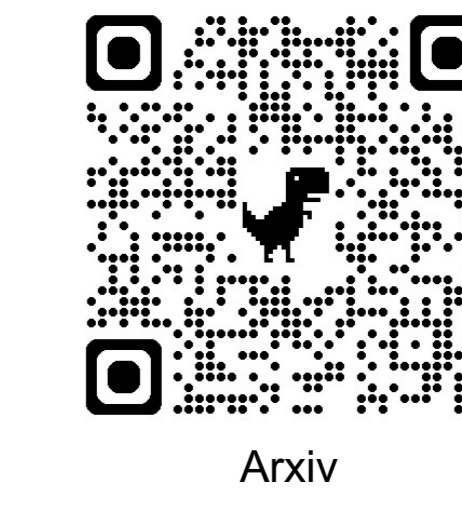
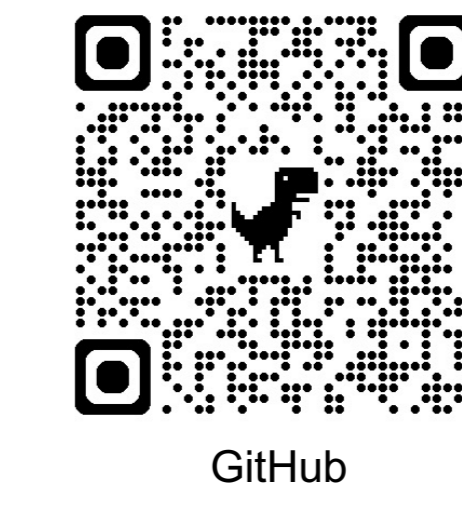


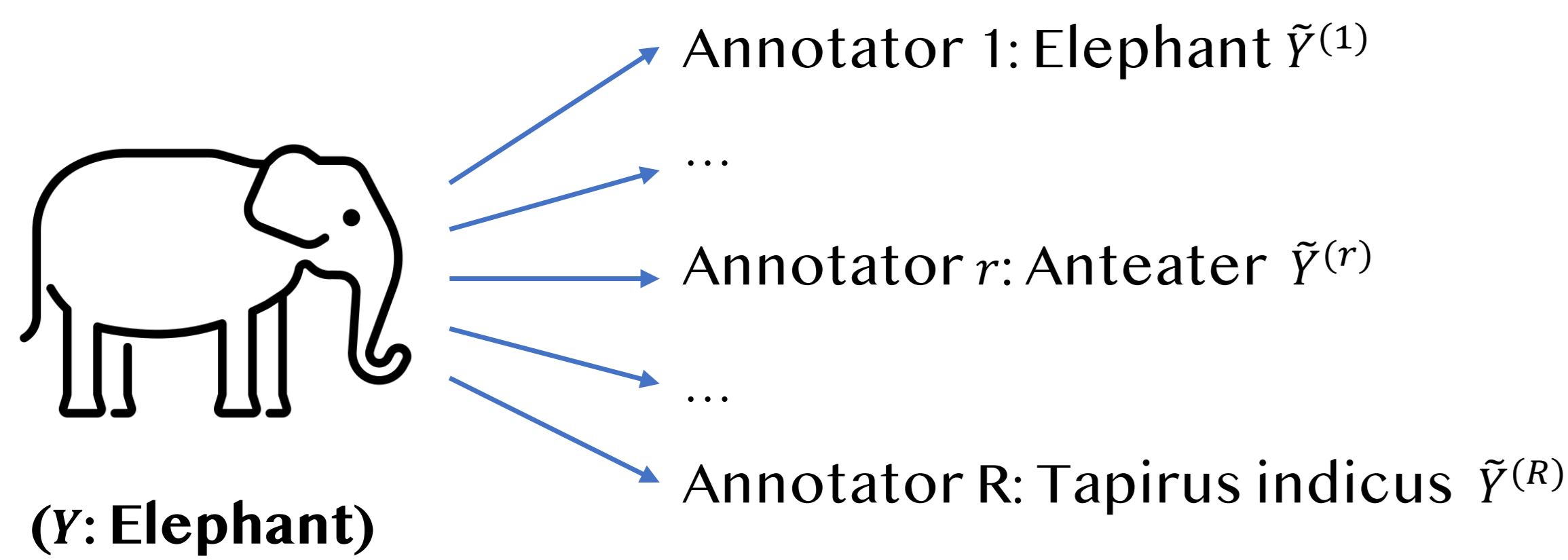
# Learning from Noisy Labels via Conditional Distributionally Robust Optimization

Hui Guo, Grace Y. Yi, Boyu Wang University of Western Ontario



## Problem

- **Learning with noisy labels:** the true label is unobserved; instead, noisy labels are collected
- **Crowdsourcing:** each data item is labeled by multiple annotators with diverse expertise or skills



Data ( $\mathbf{X}$ ) Crowdsourced Noisy Labels ( $\tilde{\mathbf{Y}}$ )

- **Noisy training data:**  $\mathcal{D} = \{\mathbf{X}_i, \tilde{\mathbf{Y}}_i\}_{i=1}^n$ 
  - $\mathbf{X}_i$ : input data
  - $\mathbf{Y}_i$ : unobserved true label
  - $\tilde{\mathbf{Y}}_i$ : noisy labels provided by  $R$  annotators for  $\mathbf{X}_i$
- **Goal:** train a classifier  $\psi$  using the noisy dataset  $\mathcal{D}$  to accurately predict the true label for future instances

## Motivation

- **Existing methods:**
  - **Approximate** the posterior distribution of the true label  $\mathbf{Y}$ , given the observed data  $\mathbf{X}$  and  $\tilde{\mathbf{Y}}$ , denoted  $P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}$
  - **Directly apply**  $P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}$  to either **infer the true labels** or to **weight the loss functions**, **without considering potential misspecification of the associated model**
  - **Our goal:** mitigate the effects of model misspecifications
- **Conditional distributionally robust risk optimization problem (CDRP):**

$$\inf_{\psi \in \Psi} R_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}), \leftarrow \text{robust/worst-case risk}$$

with  $R_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}) \triangleq \mathbb{E}_{\mathbf{x},\tilde{\mathbf{y}}} \left[ \sup_{Q_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}} \in \Gamma_\epsilon(P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}})} \mathbb{E}_{Q_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}} \{ \ell(\psi(\mathbf{X}), \mathbf{Y}) \} \right]$  loss function

ambiguity set centered around the reference probability distribution  $P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}$ :  
 $\Gamma_\epsilon(P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}) = \{Q_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}} \in \mathcal{P}(\mathcal{Y}) : d(Q_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}, P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}) \leq \epsilon\}$  ( $d$ :  $p$ -Wasserstein distance)

  - **Challenge 1:** construct a reliable reference distribution without access to true labels
  - **Challenge 2:** solve the minimax optimization problem

## AdaptCDRP: Theoretical Analysis and Methodology

- **Proposition (Dual problem):**  $R_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}) = \mathbb{E}_{\mathbf{x},\tilde{\mathbf{y}}} \left\{ \inf_{\gamma \geq 0} \left( \gamma \epsilon^p + \mathbb{E}_{P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}} \left[ \sup_{y' \in \mathcal{Y}} \{ \ell(\psi(\mathbf{X}), y') - \gamma c^p(y', \mathbf{Y}) \} \right] \right) \right\}$
- **Relaxed problem:**  $\mathfrak{R}_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}) \triangleq \inf_{\gamma \geq 0} \mathbb{E}_{\mathbf{x},\tilde{\mathbf{y}}} \left( \gamma \epsilon^p + \mathbb{E}_{P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}}} \left[ \sup_{y' \in \mathcal{Y}} \{ \ell(\psi(\mathbf{X}), y') - \gamma c^p(y', \mathbf{Y}) \} \right] \right) \Rightarrow$  Empirical version:  $\hat{\mathfrak{R}}_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}})$ , or  $\hat{\mathfrak{R}}_\epsilon$

Learning from Noisy Labels via Conditional Distributionally Robust True Label Posterior with an Adaptive Lagrange multiplier (AdaptCDRP)

- ❖ Warm up classifiers  $\psi^{(1)}$  and  $\psi^{(2)}$ ; Approximate noise transition probabilities  $\hat{\tau}_j(\tilde{\mathbf{y}})$  for  $P^*(\tilde{\mathbf{Y}} = \tilde{\mathbf{y}} | \mathbf{Y} = j, \mathbf{x})$
- ❖ **Epoch  $t$ :** Update the classifiers with pseudo-empirical distribution

- Update approximated true label posteriors (Bayes' Rule):

$$\hat{P}_j^{(t)}(\mathbf{x}, \tilde{\mathbf{y}}) \triangleq \hat{P}^{(t)}(\mathbf{Y} = j | \tilde{\mathbf{Y}} = \tilde{\mathbf{y}}, \mathbf{X} = \mathbf{x}) \propto \psi_j^{(t)}(\mathbf{x}) \cdot \hat{\tau}_j(\tilde{\mathbf{y}})$$

- For each instance  $\mathbf{x}_i$ :

$$\frac{\hat{P}_{k^*}^{(t)}(\mathbf{x}_i, \tilde{\mathbf{y}}_i)}{\max_{j \neq k^*} \hat{P}_j^{(t)}(\mathbf{x}_i, \tilde{\mathbf{y}}_i)} \geq \mathcal{C} \implies y_i^* = k^* \text{ and collect } (\mathbf{x}_i, \tilde{\mathbf{y}}_i, y_i^*) \text{ into } \mathcal{D}_{t,\ell}^*$$

- Update the pseudo-empirical distribution  $P_{t,\ell}^*$  based on  $\mathcal{D}_{t,\ell}^*$

- Update  $\psi^{(t)}$  by minimizing  $\mathfrak{R}_\epsilon(\psi; P_{\mathbf{y}|\mathbf{x},\tilde{\mathbf{y}}})$

- reference distribution  $P_{t,\ell}^*$ ; Lagrange multiplier  $\gamma_{t-1}^{(t)}$

### Addressing Challenge 1:

- Find the optimal action  $\psi^*$  for each data point  $(\mathbf{x}, \tilde{\mathbf{y}})$
- Construct the reference distribution

**Theorem (Informal).** Let  $P_j \triangleq P(\mathbf{Y} = j | \mathbf{x}, \tilde{\mathbf{y}})$ , and let  $P^{(j)}$  denote the  $j$ -th largest element of  $\{P_k\}_{k=1}^K$ . Minimizing  $\hat{\mathfrak{R}}_\epsilon$  yields

$$\psi^{*(j)} = \frac{1}{k_0} \mathbf{1}(j \leq k_0) \text{ for some } k_0 \in [K],$$

where  $\psi^{*(j)}$  corresponds to the index of  $P^{(j)}$ . In particular, if the difference between  $P^{(1)}$  and  $P^{(2)}$  is large enough, then  $\psi^{*(1)} = 1$  and  $\psi^{*(j)} = 0$  for  $j = 2, \dots, K$ .  $\implies$  **Robust pseudo-label  $y^*$**

- ❖ **Epoch  $t$ :** Update the Lagrange multiplier

- Compute  $\alpha'_i s$  for  $i \in [nK]$  and  $s^*$

- Compute the reference value for  $\gamma$ :  $\gamma_{0,t} = \alpha^{(s^*)}$

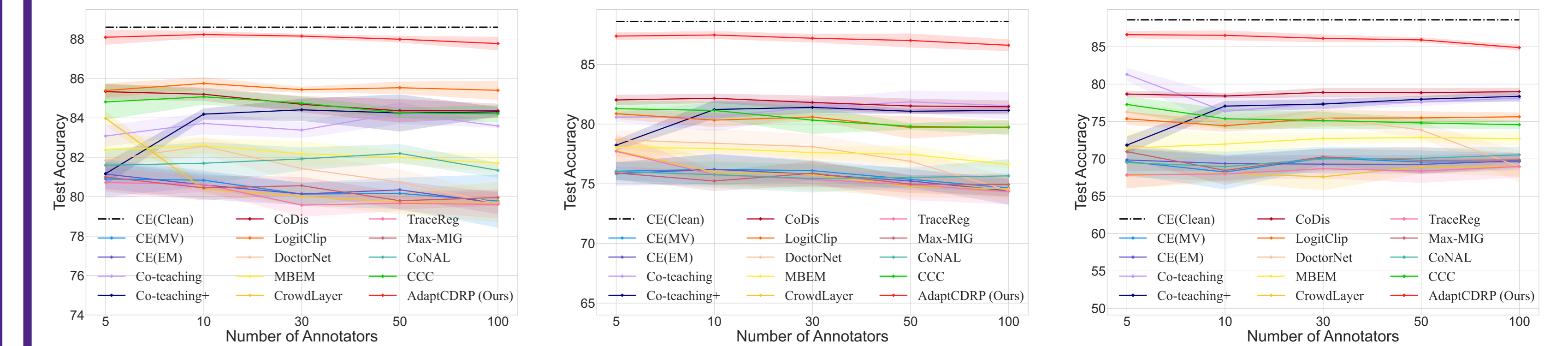
- Update  $\gamma$ :  $\gamma_t^{(t)} = \gamma_{0,t} - \frac{1}{\lambda} \{ \epsilon^p - \mathbb{E}_{P_{t,\ell}^*} c^p(y', \mathbf{Y}) \}$

### Addressing Challenge 2: Derive the optimal Lagrange multiplier $\gamma^*$

**Theorem (Informal).** Consider the loss function of the form  $\ell(\psi(\mathbf{x}), \mathbf{y}) = \sum_{j=1}^K \mathbf{1}(\mathbf{y} = j) \mathcal{T}(\psi(\mathbf{x})_j)$ . Let  $P_{i,j} \triangleq P(\mathbf{Y} = j | \mathbf{x}_i, \tilde{\mathbf{y}}_i)$  and  $\psi_{i,j} \triangleq \psi(\mathbf{x}_i)_j$  for  $i \in [n]$  and  $j \in [K]$ . An integer  $\alpha^* \in [nK + 1]$  can be defined using  $\{P_{i,j}\}$ . Let  $\alpha_{i,j} \triangleq \mathcal{T}(\min_{j \in [K]} \psi_{i,j}) - \mathcal{T}(\psi_{i,j})$ , and sort  $\{\alpha_{i,j}\}$  as  $\alpha^{(1)} \geq \dots \geq \alpha^{(nK)}$ . Then,  $\gamma^* = \alpha^{(s^*)}$ .

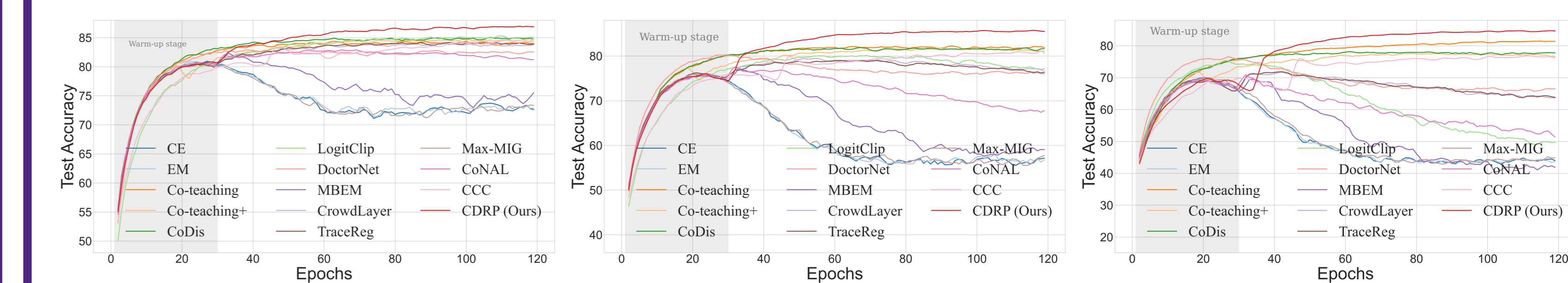
## Empirical Results

Average test accuracy on CIFAR-10 with varying numbers of annotators:



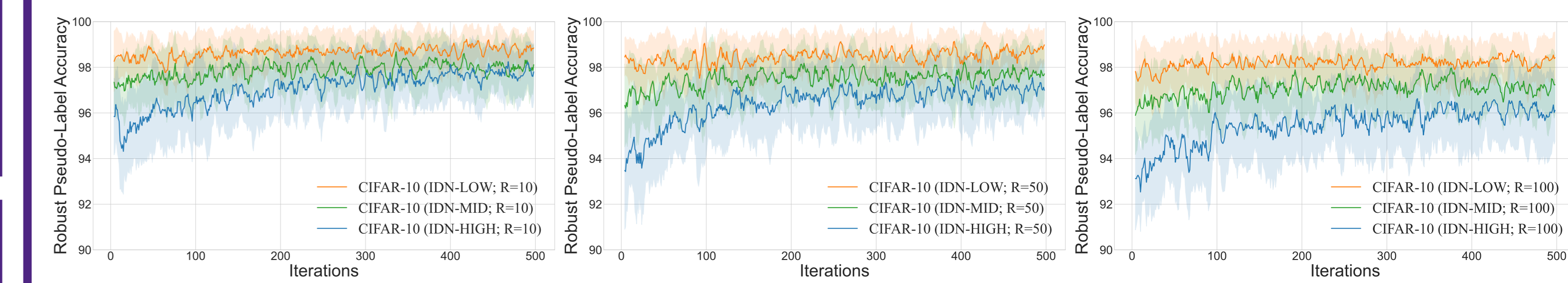
(a) IDN-LOW (b) IDN-MID (c) IDN-HIGH

Training dynamics on CIFAR-10 ( $R = 5$ ):



(d) IDN-LOW (e) IDN-MID (f) IDN-HIGH

Accuracy of robust pseudo-labels with varying numbers of annotators:



(g) IDN-LOW (h) IDN-MID (i) IDN-HIGH

Table 1: Average test accuracies on CIFAR-10 and real datasets with different transition matrix estimation methods.

Method	CIFAR-10			Real datasets	
	IDN-LOW	IDN-MID	IDN-HIGH	Animal10N	LabelMe
TraceReg	80.72±0.79	77.71±1.36	67.86±1.77	80.34±0.66	83.10±0.15
TraceReg+Ours	87.74±0.26	86.76±0.07	85.83±0.37	83.05±0.26	83.80±0.44
GeoCrowdNet (F)	84.73±0.39	81.44±1.00	77.29±1.23	81.07±0.45	84.59±0.19
GeoCrowdNet (F) + Ours	88.06±0.33	87.43±0.29	86.69±0.13	83.12±0.42	86.20±0.48
GeoCrowdNet (W)	83.82±0.53	75.72±1.10	64.64±2.23	80.19±0.33	81.63±1.49
GeoCrowdNet (W) + Ours	87.94±0.35	87.21±0.33	83.48±5.69	82.41±0.04	83.32±0.51
BayesianIDNT	86.46±1.07	85.14±0.96	82.49±2.86	81.22±0.59	83.01±0.32
BayesianIDNT + Ours	87.66±0.85	86.44±0.57	84.38±0.10	83.80±0.44	84.09±0.53

## Contributions

1. We **formulate learning with noisy labels as a CDRO problem** and **develop its dual form** to tackle the challenge of potential model misspecification in estimating the true label posterior using noisy data.
2. We derive an **analytical solution to the dual problem** for each data point, and propose a novel algorithm that **constructs a robust reference distribution** for this problem.
3. By deriving the **optimal Lagrange multiplier for the empirical robust risk**, we develop a one-step update method for the Lagrange multiplier, allowing for a principled balance between robustness and model fitting.