

# CSOT: Curriculum and Structure-Aware Optimal Transport for Learning with Noisy Labels

Wanxing Chang; Ye Shi; and Jingya Wang

37th International Conference on Neural Information Processing Systems (NeurIPS 2023)

**Learning with Noisy Labels, Optimal Transport, Curriculum Learning**

<https://dl.acm.org/doi/10.5555/3666122.3666495>

<https://github.com/changwx/CSOT-for-LNL>

## Abstract

Learning with noisy labels (LNL) poses a significant challenge in training a well-generalized model while avoiding overfitting to corrupted labels. Recent advances have achieved impressive performance by identifying clean labels and correcting corrupted labels for training. However, the current approaches rely heavily on the model's predictions and evaluate each sample independently without considering either the global or local structure of the sample distribution. These limitations typically result in a suboptimal solution for the identification and correction processes, which eventually leads to models overfitting to incorrect labels. In this paper, we propose a novel optimal transport (OT) formulation, called Curriculum and Structure-aware Optimal Transport (CSOT). CSOT concurrently considers the inter- and intra-distribution structure of the samples to construct a robust denoising and relabeling allocator. During the training process, the allocator incrementally assigns reliable labels to a fraction of the samples with the highest confidence. These labels have both global discriminability and local coherence. Notably, CSOT is a new OT formulation with a nonconvex objective function and curriculum constraints, so it is not directly compatible with classical OT solvers. Here, we develop a lightspeed computational method that involves a scaling iteration within a generalized conditional gradient framework to solve CSOT efficiently. Extensive experiments demonstrate the superiority of our method over the current state-of-the-arts in LNL.

## Problem Statement and Research Objectives

Mining large-scale labeled data based on a web search and user tags can provide a cost-effective way to collect labels, but this approach inevitably introduces noisy labels. Since DNNs can so easily overfit to noisy labels, such label noise can significantly degrade performance, giving rise to a challenging task: learning with noisy labels (LNL).

### Learning with noisy labels

- Identifying clean labels:** These methods often *model per-sample loss distributions* using a Beta Mixture Model or a Gaussian Mixture Model, **treating samples with smaller loss as clean ones**.
- Label correction methods:** Typically adopt a pseudo-labeling strategy that **leverages the DNNs predictions to correct the labels**

→ However, these approaches evaluate each sample independently **without considering the correlations among samples**, which leads to a suboptimal identification and correction solution.

### Optimal transport-based pseudo-labeling

- OT-based PL optimizes the **mapping samples to class centroids**, while considering the global structure of the sample distribution in terms of marginal constraints instead of per-sample predictions.

→ However, these approaches only consider the inter-distribution matching of samples and classes but **do not consider the intra-distribution coherence structure of samples**.

- More specifically, the cost matrix in OT relies on pairwise metrics, so two nearby samples could be mapped to two far-away class centroids (Fig. 1).

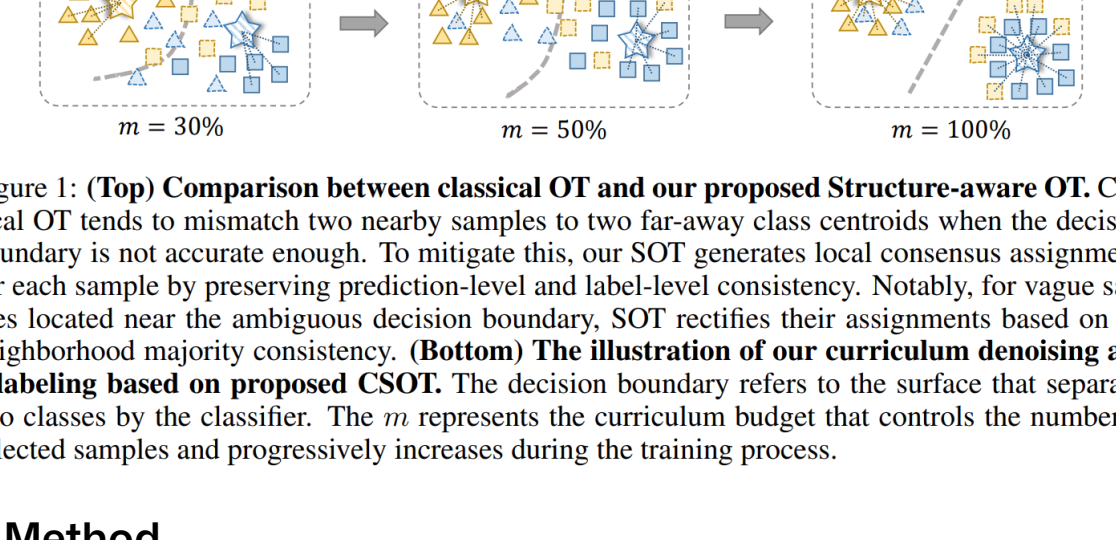


Figure 1: **(Top) Comparison between classical OT and our proposed Structure-aware OT.** Classical OT tends to mismatch two nearby samples to two far-away class centroids when the decision boundary is not accurate enough. To mitigate this, our SOT generates local consensus assignments for each sample by preserving prediction-level and label-level consistency. Notably, for vague samples located near the ambiguous decision boundary, SOT rectifies their assignments based on the neighborhood majority consistency. **(Bottom) The illustration of our curriculum denoising and relabeling based on proposed CSOT.** The decision boundary refers to the surface that separates two classes by the classifier. The  $m$  represents the curriculum budget that controls the number of selected samples and progressively increases during the training process.

## Proposed Method

In this paper, to enhance intra-distribution coherence, we propose a new OT formulation for denoising and relabeling, called **Structure-aware Optimal Transport (SOT)**. This formulation fully considers the intra-distribution structure of the samples and produces robust assignments with both *global discriminability* and *local coherence*.

- Technically speaking, we introduce **local coherent regularized terms** to encourage both prediction- and label-level local consistency in the assignments.
- Furthermore, to avoid generating incorrect labels in the early stages of training or cases with high noise ratios, we devise **Curriculum and Structure-aware Optimal Transport (CSOT)** based on SOT.
  - CSOT constructs a robust denoising and relabeling allocator by relaxing one of the equality constraints to allow only a fraction of the samples with the highest confidence to be selected.
  - These samples are then assigned with reliable pseudo labels.
  - The allocator progressively selects and relabels batches of high-confidence samples based on an increasing budget factor that controls the number of selected samples.
  - Notably, CSOT is a new OT formulation with a nonconvex objective function and curriculum constraints, so it is significantly different from the classical OT formulations.

### 1. Structure-Aware Optimal Transport for Denoising and Relabeling

Our proposed SOT for denoising and relabeling is formulated by **adding two local coherent regularized terms** based on Problem (2). Given a cosine similarity  $\mathbf{S} \in \mathbb{R}^{B \times B}$  among samples in feature space, a one-hot label matrix  $\mathbf{L} \in \mathbb{R}^{B \times C}$  transformed from given noisy labels, SOT is formulated as follows

$$\min_{\mathbf{Q} \in \Pi(\frac{1}{B}\mathbf{1}_B, \frac{1}{C}\mathbf{1}_C)} \langle -\log \mathbf{P}, \mathbf{Q} \rangle + \kappa (\Omega^{\mathbf{P}}(\mathbf{Q}) + \Omega^{\mathbf{L}}(\mathbf{Q})) \quad (3)$$

where the local coherent regularized terms  $\Omega^{\mathbf{P}}$  and  $\Omega^{\mathbf{L}}$  encourages prediction-level and label-level local consistency respectively, and are defined as follows (where  $\odot$  indicates element-wise multiplication)

$$\Omega^{\mathbf{P}}(\mathbf{Q}) = -\sum_{i,j} \mathbf{S}_{ij} \sum_k \mathbf{P}_{ik} \mathbf{P}_{jk} \mathbf{Q}_{ik} \mathbf{Q}_{jk} = -\langle \mathbf{S}, (\mathbf{P} \odot \mathbf{Q})(\mathbf{P} \odot \mathbf{Q})^{\top} \rangle \quad (4)$$

$$\Omega^{\mathbf{L}}(\mathbf{Q}) = -\sum_{i,j} \mathbf{S}_{ij} \sum_k \mathbf{L}_{ik} \mathbf{L}_{jk} \mathbf{Q}_{ik} \mathbf{Q}_{jk} = -\langle \mathbf{S}, (\mathbf{L} \odot \mathbf{Q})(\mathbf{L} \odot \mathbf{Q})^{\top} \rangle \quad (5)$$

To be more specific,  $\Omega^{\mathbf{P}}$  encourages assigning larger weight to  $\mathbf{Q}_{ik}$  and  $\mathbf{Q}_{jk}$  if the  $i$ -th sample is very close to the  $j$ -th sample, and their predictions  $\mathbf{P}_{ik}$  and  $\mathbf{P}_{jk}$  from the  $k$ -th class centroid are simultaneously high.

### 2. Curriculum and Structure-Aware Optimal Transport for Denoising and Relabeling

For the purpose of robust clean label identification and corrupted label correction, we further propose a Curriculum and Structure-aware Optimal Transport (CSOT), which constructs a robust curriculum allocator.

- This curriculum allocator gradually selects a fraction of the samples with high confidence from the noisy training set, controlled by a budget factor, then assigns reliable pseudo labels for them.
- Our proposed CSOT for denoising and relabeling is formulated by **introducing new curriculum constraints based on SOT** in Problem (3).
  - Given curriculum budget factor  $m \in [0, 1]$ , **our CSOT seeks optimal coupling matrix  $\mathbf{Q}$  by minimizing following objective**

$$\min_{\mathbf{Q}} \langle -\log \mathbf{P}, \mathbf{Q} \rangle + \kappa (\Omega^{\mathbf{P}}(\mathbf{Q}) + \Omega^{\mathbf{L}}(\mathbf{Q}))$$

$$\text{s.t. } \mathbf{Q} \in \{\mathbf{Q} \in \mathbb{R}_+^{B \times C} | \mathbf{Q}\mathbf{1}_C \leq \frac{1}{B}\mathbf{1}_B, \mathbf{Q}^{\top}\mathbf{1}_B \leq \frac{m}{C}\mathbf{1}_C\} \quad (6)$$
- Unlike SOT, which enforces an equality constraint on the samples, **CSOT relaxes this constraint and defines the total coupling budget** as  $m \in [0, 1]$ , where  $m$  represents the expected total sum of  $\mathbf{Q}$ .
  - Intuitively speaking,  $m = 0.5$  indicates that top 50% confident samples are selected from all the classes, **avoiding only selecting the same class for all the samples within a mini-batch**.
- In addition, we define the general confident scores of samples as  $\mathcal{W} = \{w_0, w_1, \dots, w_{B-1}\}$ , where  $w_i = \mathbf{Q}_{ii} / (m/C)$ .
  - Since our curriculum allocator assigns weight to only a fraction of samples controlled by  $m$ , we use  $\text{topK}(\mathcal{S}, k)$  operation (return top- $k$  indices of input set  $\mathcal{S}$ ) to identify selected samples denoted as  $\delta_i$ 

$$\delta_i = \begin{cases} 1, & \text{topK}(\mathcal{W}, [mB]) \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

### 3. Training Objectives

- To avoid error accumulation in the early stage of training, we adopt a two-stage training scheme.
  - In the first stage, the model is **supervised by progressively selected clean labels and self-supervised by unselected samples**.
  - In the second stage, the model is **semi-supervised by all denoised labels**.
- Notably, we construct our training objective based on Mixup loss  $\mathcal{L}^{\text{mix}}$  and Label consistency loss  $\mathcal{L}^{\text{lab}}$  same as NCE [45], and a self-supervised loss  $\mathcal{L}^{\text{simiam}}$  proposed in SimSiam [15].

Our two-stage training objective can be constructed as follows

$$\mathcal{L}^{\text{sup}} = \mathcal{L}_{\text{clean}}^{\text{mix}} + \mathcal{L}_{\text{clean}}^{\text{lab}} + \lambda_1 \mathcal{L}_{\text{corrupted}}^{\text{simiam}} \quad (9)$$

$$\mathcal{L}^{\text{semi}} = \mathcal{L}_{\text{clean}}^{\text{mix}} + \mathcal{L}_{\text{clean}}^{\text{lab}} + \lambda_2 \mathcal{L}_{\text{corrupted}}^{\text{lab}} \quad (10)$$

### 4. Lightspeed Computation for CSOT

- Solving Curriculum Optimal Transport
  - Problem (6) can be solved by performing iterative KL projection between  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , namely Dykstra's algorithm
$$\mathcal{C}_1 \stackrel{\text{def}}{=} \{\mathbf{Q} \in \mathbb{R}_+^{|\alpha| \times |\beta|} | \mathbf{Q}\mathbf{1}_{\beta} \geq \alpha\}$$

$$\mathcal{C}_2 \stackrel{\text{def}}{=} \{\mathbf{Q} \in \mathbb{R}_+^{|\alpha| \times |\beta|} | \mathbf{Q}^{\top}\mathbf{1}_{\alpha} \geq \beta\}$$
- Solving Curriculum and Structure-Aware Optimal Transport

#### Algorithm 1 Efficient scaling iteration for entropic regularized Curriculum OT

- Input:** Cost matrix  $\mathbf{C}$ , marginal constraints vectors  $\alpha$  and  $\beta$ , entropic regularization weight  $\varepsilon$
- Initialize:  $\mathbf{K} \leftarrow e^{-\mathbf{C}/\varepsilon}$ ,  $\mathbf{v}^{(0)} \leftarrow \mathbf{1}_{|\beta|}$
- Compute:  $\mathbf{K}_{\alpha} \leftarrow \frac{\mathbf{K}}{\text{diag}(\alpha) \mathbf{1}_{|\alpha| \times |\beta|}}$ ,  $\mathbf{K}_{\beta}^{\top} \leftarrow \frac{\mathbf{K}^{\top}}{\text{diag}(\beta) \mathbf{1}_{|\beta| \times |\alpha|}}$  // Saving computation
- for**  $n = 1, 2, 3, \dots$  **do**
- $\mathbf{u}^{(n)} \leftarrow \min \left( \frac{\mathbf{1}_{|\alpha|}}{\mathbf{K}_{\alpha} \mathbf{v}^{(n-1)}}, \mathbf{1}_{|\alpha|} \right)$
- $\mathbf{v}^{(n)} \leftarrow \frac{\mathbf{1}_{|\beta|}}{\mathbf{K}_{\beta}^{\top} \mathbf{u}^{(n)}}$
- end for**
- Return:**  $\text{diag}(\mathbf{u}^{(n)}) \mathbf{K} \text{diag}(\mathbf{v}^{(n)})$

#### Algorithm 2 Generalized conditional gradient algorithm for entropic regularized CSOT

- Input:** Cost matrix  $\mathbf{C}$ , marginal constraints vectors  $\alpha$  and  $\beta$ , entropic regularization weight  $\varepsilon$ , local coherent regularization weight  $\kappa$ , local coherent regularization function  $\mathbb{R} : \mathbb{R}^{|\alpha| \times |\beta|} \rightarrow \mathbb{R}$ , and its gradient function  $\nabla \mathbb{R} : \mathbb{R}^{|\alpha| \times |\beta|} \rightarrow \mathbb{R}^{|\alpha| \times |\beta|}$
- Initialize:  $\mathbf{Q}^{(0)} \leftarrow \alpha \beta^{\top}$
- for**  $i = 1, 2, 3, \dots$  **do**
- $\mathbf{G}^{(i)} \leftarrow \mathbf{Q}^{(i)} + \kappa \nabla \mathbb{R}(\mathbf{Q}^{(i)})$  // Gradient computation
- $\mathbf{Q}^{(i)} \leftarrow \arg\min_{\mathbf{Q} \in \Pi(\alpha, \beta)} \langle \mathbf{G}^{(i)}, \mathbf{Q} \rangle + \varepsilon \langle \mathbf{Q}, \log \mathbf{Q} \rangle$  // Linearization, solved efficiently by Algorithm 1
- Choose  $\eta^{(i)} \in [0, 1]$  so that it satisfies the Armijo rule // Backtracking line-search
- $\mathbf{Q}^{(i+1)} \leftarrow (1 - \eta^{(i)}) \mathbf{Q}^{(i)} + \eta^{(i)} \mathbf{Q}^{(i)}$  // Update
- end for**
- Return:**  $\mathbf{Q}^{(i)}$

## Evaluation and Results

### 1. Comparison with the State-of-the-Arts

Table 1: **Comparison with state-of-the-art methods in test accuracy (%) on CIFAR-10 and CIFAR-100.** The results are mainly copied from [45, 48]. We present the performance of our CSOT method using the "mean±variance" format, which is obtained from 3 trials with different seeds.

Learning Technique	(g) CSOT repl. $\mathcal{L}^{sup}$ with $\mathcal{L}^{co}$	93.47	81.93	53.45	91.43
	(h) CSOT w/o $\mathcal{L}^{semi}$	95.34	93.04	88.9	94.11
	(i) CSOT repl. correction with CT (0.95)	95.46	90.73	89.09	95.21
	(j) CSOT w/o $\mathcal{L}^{sym}$	95.92	94.17	89.31	95.21
	(k) CSOT w/o $\mathcal{L}^{corrupted}$				
	CSOT	<b>96.20</b>	<b>94.39</b>	<b>90.65</b>	<b>95.50</b>

Accuracy of selected clean labels

Method	Start	End
CSOT	0.95	0.85
CSOT w/o $\mathcal{L}^{semi}$	0.95	0.80
CSOT repl. $\mathcal{L}^{sup}$ with $\mathcal{L}^{co}$	0.45	0.85
CSOT w/o $\mathcal{L}^{sym}$	0.95	0.90
CSOT w/o $\mathcal{L}^{corrupted}$	0.95	0.95

Accuracy of corrected labels

Method	Start	End
CSOT	0.55	0.95
CSOT w/o $\mathcal{L}^{semi}$	0.55	0.90
CSOT repl. $\mathcal{L}^{sup}$ with $\mathcal{L}^{co}$	0.55	0.75
CSOT w/o $\mathcal{L}^{sym}$	0.55	0.90
CSOT w/o $\mathcal{L}^{corrupted}$	0.55	0.90

Recall rate (%)

Method	Recall rate (%)
CSOT	92.6

Table 2: **Comparison with SOTA methods in top-1 / 5 test accuracy (%) on the Webvision and ImageNet ILSVRC12 validation sets.**

Method	Webvision	top-1	top-5	ILSVRC12	top-1	top-5
F-correction [56]	61.12	82.68	57.36	82.36	-	-
Decoupling [52]	62.54	84.74	58.26	82.26	-	-
MentorNet [39]	63.00	81.40	57.80	79.92	-	-
Co-teaching [31]	63.58	85.20	61.48	84.70	-	-
DivideMix [46]	72.32	91.64	75.20	90.84	-	-
ELR [50]	76.26	91.26	68.71	87.84	-	-
ELR+ [50]	77.78	91.68	70.29	89.76	-	-
NCE [45]	79.20	91.80	74.40	91.00	-	-
RRL [48]	77.80	91.30	74.40	90.90	-	-
UncCon [44]	77.60	93.44	75.29	93.72	-	-
MOIT [51]	77.90	91.90	73.80	91.70	-	-
NCE [45]	79.50	93.80	76.30	94.10	-	-
CSOT	79.67	91.95	76.64	91.67	-	-

### 2. Ablation Studies and Analysis

Table 4: **Ablation studies under multiple label noise ratios on CIFAR-10 and CIFAR-100.** "repl." is an abbreviation for "replaced", and  $\mathcal{L}^{\text{ce}}$  represents a cross-entropy loss. GMM refers to the selection of clean labels based on small-loss criterion [46]. CT (confidence thresholding [62]) is a relabeling scheme where we set the CT value to 0.95.

<https://github.com/lijiachang/LNL-NCE>

**Curriculum learning (CL):** attempts to **gradually increase the difficulty of the training samples** allowing the model to learn progressively from easier concepts to more complex ones.

- CL has been applied to various machine learning tasks, including image classification and reinforcement learning.

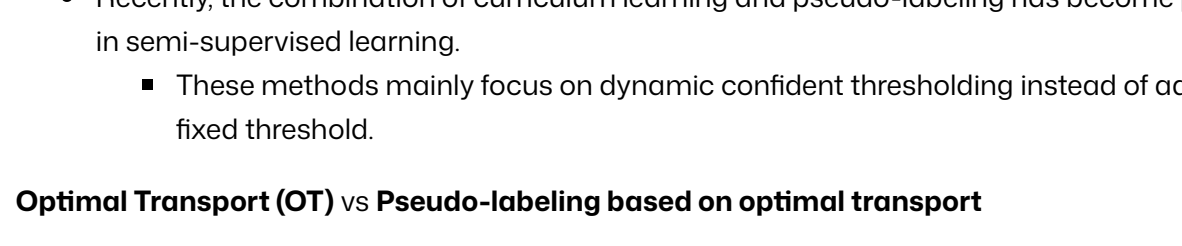


Figure 2: **Performance comparison for clean label identification and corrupted label correction.**

## Notes

- <https://github.com/lijichang/LNL-NCE>
- Curriculum learning (CL):** attempts to **gradually increase the difficulty of the training samples**, allowing the model to learn progressively from easier concepts to more complex ones.
  - CL has been applied to various machine learning tasks, including image classification, and reinforcement learning.
  - Recently, the combination of curriculum learning and pseudo-labeling has become popular in semi-supervised learning.
    - These methods mainly focus on dynamic confident thresholding instead of adopting a fixed threshold.
- Optimal Transport (OT) vs Pseudo-labeling based on optimal transport**

Optimal Transport (OT)	Pseudo-labeling based on optimal transport
$\min_{\mathbf{Q} \in \Pi(\alpha, \beta)} \langle \mathbf{C}, \mathbf{Q} \rangle \quad (1)$ <p>* <math>\langle \cdot, \cdot \rangle</math> : Frobenius dot-product</p> <ul style="list-style-type: none"> <li><math>\alpha, \beta</math> : probability vectors (indicating two distributions)                             <ul style="list-style-type: none"> <li><math> \alpha </math> : the dimension of <math>\alpha</math></li> </ul> </li> <li>Cost matrix : <math>\mathbf{C} \in \mathbb{R}^{ \alpha  \times  \beta }</math></li> <li>Coupling matrix <math>\mathbf{Q}</math> <ul style="list-style-type: none"> <li><math>\mathbf{Q}_{ik}</math> indicates how the mass is moved from the <math>i</math> of the <math>\alpha</math> to the <math>k</math> of the <math>\beta</math>.</li> <li>Constraints:                                     <math display="block">\mathbf{Q}\mathbf{1}_{ \beta } = \alpha, \quad \mathbf{Q}^{\top}\mathbf{1}_{ \alpha } = \beta</math> </li> </ul> </li> </ul>	$\min_{\mathbf{Q} \in \Pi(\frac{1}{B}\mathbf{1}_B, \frac{1}{C}\mathbf{1}_C)} \langle -\log \mathbf{P}, \mathbf{Q} \rangle \quad (2)$ <p>* <math>\mathbf{1}_d</math> : <math>d</math>-dimensional vector of ones</p> <ul style="list-style-type: none"> <li><math>B, C</math> : the batch size of samples, and the number of classes</li> <li>Cost matrix : <math>\mathbf{C} = -\log \mathbf{P}</math> <ul style="list-style-type: none"> <li><math>\mathbf{P} \in \mathbb{R}_+^{B \times C}</math> : classifier softmax predictions</li> </ul> </li> <li>Coupling matrix <math>\mathbf{Q}</math> <ul style="list-style-type: none"> <li><math>\mathbf{Q}_{ik}</math> indicates the probability that the sample <math>i</math> will be mapped to a class <math>k</math></li> <li>Constraints:                                     <math display="block">\mathbf{Q}\mathbf{1}_C = \frac{1}{B}\mathbf{1}_B, \quad \mathbf{Q}^{\top}\mathbf{1}_B = \frac{1}{C}\mathbf{1}_C</math> </li> </ul> </li> </ul>

- Directly optimizing the exact OT problem would be time-consuming, and an entropic regularization term is introduced:

$$\min_{\mathbf{Q} \in \Pi(\alpha, \beta)} \langle \mathbf{C}, \mathbf{Q} \rangle + \epsilon \langle \mathbf{Q}, \log \mathbf{Q} \rangle, \quad \text{where } \epsilon > 0$$