
Characterizing Out-of-Distribution Error via Optimal Transport

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Out-of-distribution (OOD) data poses serious challenges in deployed machine
2 learning models, so methods of predicting a model’s performance on OOD data
3 without labels are important for machine learning safety. While a number of meth-
4 ods have been proposed by prior work, they often underestimate the actual error,
5 sometimes by a large margin, which greatly impacts their applicability to real tasks.
6 In this work, we identify *pseudo label shift*, or the difference between the predicted
7 and true OOD label distributions, as a key indicator to this under-estimation. Based
8 on this observation, we introduce a novel method for estimating model performance
9 by leveraging optimal transport theory, Confidence Optimal Transport (COT), and
10 show that it provably provides more robust error estimates in the presence of
11 pseudo label shift. Additionally, we introduce an empirically-motivated variant of
12 COT, Confidence Optimal Transport with Thresholding (COTT), which applies
13 thresholding to the individual transport costs and further improves the accuracy
14 of COT’s error estimates. We evaluate COT and COTT on a variety of standard
15 benchmarks that induce various types of distribution shift – synthetic, novel sub-
16 population, and natural – and show that our approaches significantly outperform
17 existing state-of-the-art methods with an up to 3x lower prediction error.

18 1 Introduction

19 Machine Learning methods are largely based on the assumption that test samples are drawn from
20 the same distribution as training samples, providing a basis for generalization. However, this i.i.d.
21 assumption is often violated in real-world applications where test samples are found to be out-of-
22 distribution (OOD) – sampled from a different distribution than during training. This may result in a
23 significant negative impact on model performance [14, 32, 24]. A common practice for alleviating
24 this issue is to regularly gauge the model’s performance on a set of labeled data from the current *target*
25 data distribution, and update the model if necessary. When labeled data is not available, however, one
26 needs to predict the model’s performance on the target distribution with unlabeled data, a task known
27 as *OOD performance prediction*.

28 Performance prediction on unlabeled data has previously been shown to be impossible without
29 imposing additional constraints over the unknown target distribution [7, 11, 5, 25], due to the fact that
30 target samples may take any label. Thus, the feasibility of this task is dependent on what assumptions
31 we make regarding the shift between the train and target distributions. Prior works often make the
32 assumption that the conditional density $P(y|x)$ remains fixed in the presence of covariate shift [35].
33 However, this tells us little when x falls outside the support of the train distribution. Despite the
34 theoretical difficulty, prior works have proposed a number of heuristic methods to estimate the
35 performance of a model based on unlabeled target samples. For instance, Average Confidence
36 (AC) [19, 15] estimates error based on the average maximum softmax score for target samples,

37 assuming the model has been calibrated for the train distribution. This method was further improved
 38 with the addition of a learned threshold [12], for which the error is predicted as the fraction of
 39 samples with confidence falling below it. Other approaches estimate model performance based
 40 on a disagreement score computed between the predictions of two models trained over the same
 41 dataset [21, 3]. Some works have found that applying a transformation over target samples and
 42 estimating the effect of the transformation leads to a reliable prediction of model performance in
 43 vision tasks [9, 10]. However, many of these prior methods have been shown to underestimate
 44 the model’s error when it is *miscalibrated* in the target distribution [12, 21]; that is, the predicted
 45 softmax scores differ from the true class likelihoods. We empirically observe that this miscalibration
 46 is strongly positively correlated with *pseudo-label shift*, which is the difference between the predicted
 47 target label distribution $P_T(\tilde{y})$ and the true label distribution $P_T(y)$. Thus, we treat pseudo label shift
 48 as a key indicator of error underestimation.

49 In this work, we propose an approach which provides robust error estimates in the presence of pseudo-
 50 label shift. Our approach, Confidence Optimal Transport (COT), leverages the optimal transport
 51 framework to predict the error of a model as the Wasserstein distance between the predicted target
 52 class probabilities and the true source label distribution. We theoretically derive lower bounds for
 53 COT’s predicted error. This results in a more provably conservative error prediction than AC, which
 54 is crucial for safety in many real-world machine learning applications. In addition, we introduce
 55 a variant of COT, Confidence Optimal Transport with Thresholding (COTT), which introduces a
 56 learned threshold over the optimal transportation costs in line with prior work [12] and empirically
 57 improves upon COT’s performance. We compare our proposed methods to existing state-of-the-art
 58 approaches in extensive empirical experiments over eleven datasets exhibiting distribution shift from
 59 multiple vision and language domains. These distribution shifts include: visual corruptions (e.g.,
 60 blurred image data), novel subpopulation shifts (e.g., novel appearances of a category), and natural
 61 shifts in the wild (e.g., different stain colors for pathology images), all of which are frequently
 62 experienced in the real world. We find that COT and COTT consistently avoid the significant error
 63 underestimations suffered by previous methods. In particular, COTT achieves significantly lower
 64 prediction errors than existing methods for most models and datasets (up to 3x better), establishing
 65 new state-of-the-art results.

66 2 Preliminaries and Motivation

67 In this section, we introduce the problem setup of OOD error prediction and show how a popular
 68 baseline, Average Confidence (AC), can be understood as the Wasserstein Distance (WD) between
 69 a reference label distribution and the softmax output distribution from an Optimal Transport (OT)
 70 perspective. We then demonstrate why the reference label distribution AC uses is problematic by
 71 explaining pseudo-label shift and its correlation with miscalibration. After establishing the relation
 72 and utility of OT to OOD error prediction, we formally introduce our proposed methods in Sec. 3.

73 2.1 OOD Performance Prediction

74 In this work, we address the OOD problem in the domain of classification tasks. Let $\mathcal{X} \subseteq \mathbb{R}^d$ be the
 75 input space of size d , and $\mathcal{Y} = \{1, \dots, k\}$ be the label space, where k is the number of classes. Let
 76 the source distribution over $\mathcal{X} \times \mathcal{Y}$ be $P_S(x, y)$ and the target distribution be $P_T(x, y)$. A classifier
 77 $\vec{f}: \mathcal{X} \rightarrow \Delta^{k-1}$ maps an input to a *confidence vector* (i.e. the output of the softmax layer), where
 78 $\Delta^{k-1} = \{(z_1, \dots, z_k) : z_1 + \dots + z_k = 1, z_i \geq 0\}$ is the k -dimensional unit simplex. Let the training
 79 samples be $\{(x_{\text{train}}^{(i)}, y_{\text{train}}^{(i)})\} \sim P_S(x, y)$, and the validation samples be $\{(x_{\text{val}}^{(i)}, y_{\text{val}}^{(i)})\} \sim P_S(x, y)$. The
 80 validation set is used in some previous methods and will also be used in COT and COTT.

81 The OOD performance prediction problem is formally stated as follows: Given an unlabeled test
 82 set $\{x_T^{(i)}\} \sim P_T(x)$, and a classifier \vec{f} trained over the training samples, predict the accuracy of
 83 \vec{f} over the test set $\alpha = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[y_T^{(i)} - \arg \max_j \vec{f}_j(x_T^{(i)})]$, where $y_T^{(i)}$ is the ground truth label.
 84 Equivalently, we can also predict the error $\epsilon = 1 - \alpha$ with an estimate $\hat{\epsilon}$.

85 In this work, we are further investigating *the distribution of confidence vectors*. While a confidence
 86 vector $\vec{f}(x)$ itself is a distribution of labels in the k -dimensional simplex Δ^{k-1} , the distribution of
 87 confidence vectors $\vec{f}_{\#}P(\vec{c})$ is a distribution of distributions, where $\vec{c} \in \Delta^{k-1}$. $\vec{f}_{\#}P(\vec{c})$ is defined

88 to be the *pushforward* of a covariate distribution $P(x)$ using \vec{f} : $\vec{f}_\# P(\vec{c}) = P(\vec{f}^{-1}(\vec{c}))$. Consider
 89 the following example: Suppose we have a uniform covariate distribution $P(x)$ on $x \in \{A, B, C\}$
 90 and \vec{f} maps A, B to $[0.5, 0.5]^\top$ and C to $[0.7, 0.3]^\top$. Then, $\vec{f}_\# P(\vec{c})$ will assign $\frac{2}{3}$ mass to $[0.5, 0.5]^\top$
 91 and $\frac{1}{3}$ mass to $[0.7, 0.3]^\top$. To facilitate easy comparison between confidence vectors and labels, we
 92 denote $\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}$ to be the one-hot representation of $y \in \mathcal{Y}$, where the only non-zero
 93 element in \vec{y} is the y -th element, i.e. $\vec{y}_j = \mathbb{1}[y = j]$. We denote *the distribution of one-hot labels*
 94 as $P(\vec{y})$. For a covariate distribution $P(x)$, we will refer to the distribution of predicted labels
 95 from classifier \vec{f} as the *pseudo-label distribution* $P_{\text{pseudo}}(\vec{y})$. We reuse the notation to denote the
 96 probability mass of \vec{y} also as which is given by the mass of the inputs that give this prediction, i.e.
 97 $P_{\text{pseudo}}(\vec{y}) = P(\{x \in \mathcal{X} \mid \arg \max_j \vec{f}_j(x) = y\})$.

98 2.2 Wasserstein Distance and Optimal Transport

99 In recent years, optimal transport theory has found numerous applications in the field of machine
 100 learning [2, 4, 1, 26]. Optimal transport aims to move one distribution of mass to another as efficiently
 101 as possible under a given cost function. In the Kantorovich formulation of optimal transport, we
 102 are given two distributions $\mu(x)$ over \mathcal{X} and $\nu(y)$ over \mathcal{Y} and a cost function $c(x, y)$ that tells us
 103 the cost of transporting from location x to location y . Here, we aim to find a transport plan $\pi(x, y)$
 104 that minimizes the total transport cost. The transport plan is a joint distribution with marginal
 105 $\pi(\cdot, \mathcal{Y}) = \mu(\cdot)$ and $\pi(\mathcal{X}, \cdot) = \nu(\cdot)$. The conditional distribution $\frac{\pi(x, y)}{\mu(x)}$ of the transport plan informs
 106 us how much mass is moved from x to y . Let $\Pi(\mu, \nu)$ be the set of all transport plans. More formally,
 107 the Wasserstein Distance is defined as:

$$W(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \int_{\mathcal{X} \times \mathcal{Y}} c(x, y) d\pi(x, y)$$

108 The Wasserstein distance satisfies the definition of a metric (non-negativity, symmetry, and sub-
 109 additivity), inducing a metric space over a space of probability distributions. Unlike other metrics,
 110 such as total variation, the Wasserstein metric induces a weaker topology and provides a robust
 111 framework for comparing probability distributions that respect the underlying space geometry [2].

112 For discrete distributions μ, ν such as the empirical distributions, the Wasserstein distance simplifies
 113 to the following linear programming problem:

$$W(\mu, \nu) = \min_{\mathbf{P} \in \Pi(\mu, \nu)} \langle \mathbf{P}, \mathbf{C} \rangle = \min_{\mathbf{P} \in \Pi(\mu, \nu)} \sum_{i, j} \mathbf{P}_{ij} \mathbf{C}_{ij}$$

114 where $\mathbf{C}, \mathbf{P} \in \mathbb{R}^{m \times n}$ are the cost matrix and the plan matrix respectively and \mathbf{C}_{ij} is the transport
 115 cost from sample i to sample j . When $n = m$, the optimal transport problem reduces to the
 116 optimal matching problem (Proposition 2.1 in Peyré et al. [30]), i.e. the optimal transport plan
 117 $\mathbf{P}^* = \arg \min_{\mathbf{P} \in \Pi(\mu, \nu)} \langle \mathbf{P}, \mathbf{C} \rangle$ is a permutation matrix. Not only does this constraint enable efficient
 118 algorithms like the Hungarian algorithm, but it will also help draw the connection between pseudo-
 119 label shift and target error that is central to our Confidence Optimal Transport method.

120 2.3 Average Confidence as Wasserstein Distance

121 We use Average Confidence, a popular OOD performance prediction method, as the starting point
 122 of our analysis. Average Confidence with Max Confidence (AC-MC) estimates the target accuracy
 123 by taking the empirical mean of maximum confidence of the classifier over all target samples
 124 $x^{(i)} \sim P_T(x)$, i.e. $\frac{1}{n} \sum_{i=1}^n \max_j \vec{f}_j(x^{(i)})$. Its corresponding target error estimate is therefore
 125 $\hat{\epsilon}_{\text{AC}} = 1 - \frac{1}{n} \sum_{i=1}^n \max_j \vec{f}_j(x^{(i)})$. By definition, $|\epsilon - \hat{\epsilon}_{\text{AC}}|$ measures the miscalibration of the model,
 126 i.e. how far away the model's confidence is from its actual accuracy. In the following proposition, we
 127 connect the AC-MC estimates with distances in the Wasserstein Metric Space:

128 **Proposition 1** ($\hat{\epsilon}_{\text{AC}} - W_\infty$ Equivalence). *Let $(\mathcal{P}(\mathbb{R}^k), W_\infty)$ be the metric space of all distributions
 129 over \mathbb{R}^k , where W_∞ is the Wasserstein distance with $c(x, y) = \|x - y\|_\infty$. Then, the estimated error
 130 of AC-MC is given by $\hat{\epsilon}_{\text{AC}} = W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y}))$.*

131 We defer all proofs to the supplementary material. Here, we appeal to pictorial intuition in Figure
 132 1, where δ_0 (Dirac delta over the zero vector) represents the origin. All one-hot label distributions,

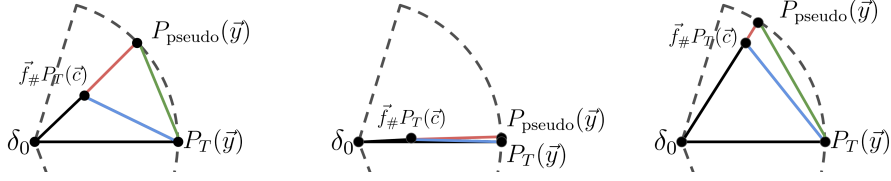


Figure 1: The AC-COT Triangle in the Wasserstein Space ($\mathcal{P}(\mathbb{R}^k), W_{\infty}$). **Red** line: AC-MC error estimate. **Blue** line: (our) COT error estimate (assuming $P_T(\vec{y}) \approx P_S(\vec{y})$). **Green** line: Pseudo-label shift. **Left:** AC-MC predicts the target error as the distance between the *distribution of confidence vectors* and its projection on the unit sphere, the smallest among all label distributions. This makes AC-MC prone to *underestimating* the target error. **Middle:** Mild pseudo-label shift. **Right:** Severe pseudo-label shift.

133 including $P_T(\vec{y})$ and $P_{\text{pseudo}}(\vec{y})$, are represented as points on the (dashed) unit sphere around δ_0 .
 134 A distribution of confidence vectors $\vec{f}_{\#}P(\vec{c})$ is simply a point within the unit ball around δ_0 . The
 135 AC-MC accuracy estimate measures how far the point is to the origin δ_0 . Additionally, we can project
 136 the point to the unit sphere, resulting in the projection point $P_{\text{pseudo}}(\vec{y})$, which is the pseudo-label
 137 distribution. The AC-MC error estimate measures *the length of the projection*.

138 **Corollary 1** ($P_{\text{pseudo}}(\vec{y})$ is closest to $\vec{f}_{\#}P(\vec{c})$). Let $P'(\vec{y}) \in \mathcal{P}(\{0, 1\}^k \cap \Delta^{k-1})$ be a one-hot label
 139 distribution. Then $W_{\infty}(\vec{f}_{\#}P(\vec{c}), P'(\vec{y})) \geq W_{\infty}(\vec{f}_{\#}P(\vec{c}), P_{\text{pseudo}}(\vec{y}))$.

140 This suggests that AC-MC, among all possible reference one-hot label distributions, selects the
 141 closest one and reports the distance to be the predicted error. Thus, we can see that AC-MC is
 142 a very optimistic prediction strategy, so it is not surprising that AC-MC is widely reported to be
 143 over-estimating the performance on real tasks [12, 15, 21, 3].

144 2.4 Pseudo-Label Shift and its Correlation with Miscalibration

145 The pseudo-label shift is defined as $W_{\infty}(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$, i.e. the distance from the pseudo label
 146 distribution $P_{\text{pseudo}}(\vec{y})$ to the ground truth target label distribution $P_T(\vec{y})$, which is also the length of
 147 the green line segment in Figure 1. An important property is $W_{\infty}(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) \leq \epsilon \leq 1$, i.e.
 148 the pseudo-label shift is a lower bound of the true target error of the model (as true target error can be
 149 interpreted as the average transport cost of a suboptimal matching). We can clearly see how a large
 150 pseudo-label shift could potentially destroy the prediction of AC-MC from Figure 1 (right), where
 151 the red line segment is much shorter than the green one which should be a lower bound of the true
 152 error. This lower bound, however, can be loose when it is small, as shown in Figure 1 (middle).

153 Most existing prediction methods heavily rely on the model being well-calibrated, as pointed out by
 154 Jiang et al. [21], Garg et al. [12], yet the underlying difficulty is a lack of precise understanding of
 155 *when* and *by how much* neural networks become miscalibrated. This is evident in the large-scale
 156 empirical study conducted in [28], showing that the variance of the error becomes much larger among
 157 the different types of shift studied, despite some positive correlation between the expected calibration
 158 error and the strength of the distribution shift. In this section, we empirically show that there is
 159 a strong positive correlation between the pseudo-label shift and $|\epsilon - \hat{\epsilon}_{\text{AC}}|$, which we define as the
 160 model’s miscalibration. Because of this correlation, we consider pseudo-label shift as the key signal
 161 to why existing methods have undesirable performance.

162 We evaluate the pseudo label shift and the prediction error of AC-MC $|\epsilon - \hat{\epsilon}_{\text{AC}}|$ on CIFAR-10 and
 163 CIFAR-100 under different distribution shifts, and plot the results in Figure 2 (left). The plots
 164 demonstrate a strong positive correlation between these two quantities, which not only means that
 165 the performance of AC-MC worsens as the pseudo-label shift gets larger, but also implies that the
 166 performance of existing methods [15, 12, 21] that depend on model calibration will drop.

167 Our exploration of miscalibration and pseudo-label shift reveals a previously unexplored tradeoff:
 168 When the pseudo-label shift is small, prior work has given us some reassurance to trust the calibration
 169 of neural networks, which motivated methods such as AC [19] and ATC [12]. More critically, as
 170 miscalibration worsens, the pseudo-label shift becomes a tighter low bound of the error and thus a
 171 more trustworthy estimate of the error.



Figure 2: **Left:** Absolute difference between the model’s Average Confidence and its true error under different levels of pseudo-label shift. This difference measures the degree of miscalibration. We see a strong correlation between the miscalibration and pseudo-label shift on the common corruption benchmarks (CIFAR-10-C, CIFAR-100-C). **Right:** Absolute difference between GDE error estimate and true error under different levels of pseudo-label shift. The strong correlation is also observed.



Figure 3: We further compare the sensitivity of COT and COTT’s estimation error to the degree of pseudo-label shift. Compared to Figure 2, we can clearly see that the correlation between the prediction error and pseudo-label shift weakens significantly. Moreover, COTT is even more robust to pseudo-label shift compared to COT.

172 Motivated by this observation, we leverage the information of the pseudo-label shift to improve the
 173 performance of confidence-based prediction methods for miscalibrated models. The problem, as
 174 mentioned earlier, is that without any information about the target label distribution, it is theoretically
 175 impossible to estimate the pseudo label shift when the test set contains data from unseen domains.
 176 Thus, the only option is to make assumptions on the target label distribution.

177 In our proposed method, we make a natural assumption: the target label distribution is close to the
 178 source label distribution. This assumption aligns with most natural shifts that can be observed in
 179 real-world datasets. This extra assumption allows us to develop COT and COTT, which perform
 180 much better than existing methods in most cases (See Table 1), especially when the pseudo-label
 181 shift is large. Given this assumption, Figure 3 (Left), shows the prediction error of COT versus the
 182 pseudo-label shift. We can see that equipped with the extra information, COT is able to maintain
 183 a low prediction error even under very large pseudo label shift, and the correlation between COT
 184 performance and the pseudo label shift is weak. Since the pseudo label shift is strongly correlated
 185 with miscalibration, this implies that COT is much more robust to miscalibration than existing
 186 miscalibration-sensitive prediction methods.

187 3 Methods

188 In this section, we formally introduce our proposed method – Confidence Optimal Transport – and
 189 propose an additional variation, COTT, that utilizes a threshold over the optimal transportation costs
 190 between two distributions, instead of taking a simple average.

191 3.1 Confidence Optimal Transport

192 Let $\hat{P}_S(\vec{y})$ denote the empirical source label distribution. The predicted error of COT, which is the
 193 length of the blue line segment (assuming $P_T(\vec{c}) = P_S(\vec{y})$) in Figure 1, is given by

$$\hat{e}_{\text{COT}} = W_\infty(\vec{f}_\# \hat{P}_T(\vec{c}), \hat{P}_S(\vec{y})).$$

194 By Corollary 1, $\hat{\epsilon}_{\text{COT}} \geq \hat{\epsilon}_{\text{AC}}$, which means that COT provably predicts a larger error than AC-MC,
 195 which empirically tends to produce more accurate predictions as we find overestimation of the error
 196 far less common than underestimation. As mentioned earlier, in the case where the pseudo-label shift
 197 is large, such as in Figure 1 (Right), AC-MC can have arbitrarily large prediction error, while COT
 198 always has the following guarantee:

199 **Proposition 2** (Calibration independent lower bound of COT). *Under the assumption that $P_T(\vec{y}) =$
 200 $P_S(\vec{y})$, we always have $\hat{\epsilon}_{\text{COT}} \geq 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$.*

201 Thus, we can see that COT is by nature different from existing methods because it is a *miscalibration-*
 202 *robust method*. Its success is dependent on the difference between $P_T(\vec{y})$ and $P_S(\vec{y})$, which is
 203 relatively small in most real-world scenarios, while the dependency on calibration of existing methods
 204 can not always be controlled. The geometric explanation of this is that COT measures a fundamentally
 205 different length in the Wasserstein Space $(\mathcal{P}(\mathbb{R}^k), W_\infty)$, allowing for the above guarantee which
 206 does not exist in previous methods. Moreover, as we will empirically demonstrate in the next section,
 207 large pseudo-label shift is prevalent in real models and datasets. Consequently, COT performs much
 208 better than miscalibration-sensitive methods in most cases.

209 3.2 Thresholding as a Robust Transport Cost Statistic

210 Computationally, COT (or Wasserstein Distance in general) is implemented as a two-step process: 1)
 211 calculating individual transport costs for all samples; 2) returning the mean across all transport costs
 212 as the error estimate. While we have seen that COT has some protection against miscalibration, it
 213 is not completely immune. Another outstanding issue lies in the second step - computing the mean
 214 of all transport costs. In statistics, it is well known that the mean is less robust to outliers than the
 215 median [20]. In the supplemental material, we show that the empirical transport cost distribution is
 216 also heavy-tailed and therefore COT is susceptible to outlier costs. A large fluctuation in the outliers
 217 impacts the mean much more than it does to the median. Therefore, a more robust statistic is desired
 218 to protect COT against these outliers.

219 To search for such a robust statistic that is suitable for our purpose, we turn to Average Thresholded
 220 Confidence (ATC) [12] for inspiration. ATC tremendously improves the performance of AC precisely
 221 by offering such protection against the overconfident outlier predictions (shown in supp material) by
 222 returning the percentile above the threshold rather than the mean. In a similar vein, we safeguard
 223 COT against outlier costs by introducing COT with Thresholding, COTT.

224 Specifically, to compute the threshold $t \in [0, 1]$, we first sample a validation set $\{(x_{\text{val}}^{(i)}, y_{\text{val}}^{(i)})\}_{i=1}^n \sim$
 225 $P_S(x, y)$. Let \hat{P}_{val} denote the empirical validation distribution. We compute the optimal transport plan
 226 $\mathbf{P}_{\text{val}}^* = \arg \min_{\mathbf{P} \in \Pi(\vec{f}_{\#} \hat{P}_{\text{val}}(\vec{c}), \hat{P}_S(\vec{y}))} \langle \mathbf{P}, \mathbf{C} \rangle$, where \mathbf{C} is the cost matrix determined by the L-infinity
 227 distance. Note that $\mathbf{P}_{\text{val}}^* \in \{0, 1\}^{n \times n}$ is a permutation matrix due to the equal number of confidence
 228 vectors and one-hot labels (Proposition 2.1 in [30]). Then the threshold t is set such that the validation
 229 of error the classifier equals the fraction of samples with transport cost higher than t , i.e.

$$\epsilon_{\text{val}} = \frac{1}{n} |\{ \mathbf{C}_{ij} \geq t | \mathbf{P}_{\text{val}ij}^* = 1 \}|$$

230 With the threshold t learned from the validation set, we can compute the COTT target error estimate.
 231 Let $\hat{P}_S(\vec{y})$ denote the empirical source label distribution and $\vec{f}_{\#} \hat{P}_T(\vec{c})$ denote the empirical distribu-
 232 tion of confidence vectors of samples from the target distribution $P_T(x)$. We compute the optimal
 233 transport plan $\mathbf{P}^* = \arg \min_{\mathbf{P} \in \Pi(\vec{f}_{\#} \hat{P}_T(\vec{c}), \hat{P}_S(\vec{y}))} \langle \mathbf{P}, \mathbf{C} \rangle$. Our COTT estimate is then defined as

$$\hat{\epsilon}_{\text{COTT}} = \frac{1}{n} |\{ \mathbf{C}_{ij} \geq t | \mathbf{P}_{ij}^* = 1 \}|$$

234 4 Experiments

235 In this section, we empirically compare COT and COTT with existing methods on various benchmark
 236 datasets. For all experiments, we trained the model on in-distribution data and froze the model after
 237 convergence. To predict the model’s performance on the target domain, we only used unlabeled data
 238 from the target domain. When we have a test set size greater than 10,000, we show the results of

Table 1: Mean Absolute Error (MAE) between the estimated error and ground truth error to compare different methods. The "shift" column denotes the nature of distribution shifts for each dataset. For vision datasets, we reported results for ResNet18 and ResNet50; for language datasets, we reported results for DistilBERT-base-uncased. The results are averaged over 3 random seeds. We highlight the best-performing method. We defer the full table with std to the supplemental material.

Dataset	Shift	Baselines						Ours	
		AC	DoC	IM	GDE	ATC-MC	ATC-NE	COT	COTT
CIFAR10	Natural	5.97	5.38	5.87	5.9	3.38	3.15	5.41	3.33
	Synthetic	9.1	8.53	9.26	8.84	4.2	3.37	2.17	1.7
CIFAR100	Synthetic	10.83	8.76	12.07	11.36	6.8	6.63	2.09	2.59
ImageNet	Natural	8.5	7.43	8.62	5.62	3.57	2.6	3.88	2.41
	Synthetic	10.34	9.28	12.87	6.54	1.59	3.41	3.24	1.42
Entity13	Same	19.63	19.2	17.5	15.37	8.09	7.23	8.47	2.61
	Novel	29.61	29.18	27.22	24.48	14.54	9.49	15.9	5.46
Entity30	Same	16.97	16.21	13.56	13.98	8.19	9.08	5.9	2.46
	Novel	27.57	26.81	23.96	23.4	13.46	8.57	15.11	5.94
Living17	Same	14.84	14.67	11.22	9.94	4.88	5.43	6.25	2.94
	Novel	29.61	29.18	27.22	24.48	14.54	9.49	15.9	5.53
Nonliving26	Same	19.25	18.43	16.6	12.77	11.18	9.69	7.06	3.34
	Novel	31.37	30.54	28.79	23.37	19.93	16.56	17.8	10.46
Camelyon17-WILDS	Natural	9.44	9.44	10.24	5.19	7.73	7.73	7.27	5.71
RxRx1-WILDS	Natural	5.21	8.44	8.09	7.48	6.53	6.86	3.25	5.82
Amazon-WILDS	Natural	2.62	2.35	2.34	17.04	1.63	1.54	2.43	2.01
CivilCom.-WILDS	Natural	1.54	0.96	0.86	8.7	2.3	2.3	1.23	4.68

239 utilizing the batched version of COT and COTT detailed in Section 3.2. For context, solving the OT
 240 problem of size 10,000 only takes around 10s, thus adding only negligible computation overhead.

241 4.1 Datasets and Nature of Shift

242 In our comprehensive evaluation, we consider more than 10 benchmark datasets across multiple
 243 modalities, including vision and language, with a variety of distribution shifts:

244 **Synthetic Shift:** First, we consider distribution shifts caused by common visual corruptions, such
 245 as brightness, defocusing, and blurriness, which are common in real-world settings. We used the
 246 corrupted versions of CIFAR10, CIFAR100, and ImageNet proposed in [18], which includes 19 types
 247 of common visual corruptions across 5 levels of severity.

248 **Novel Subpopulation Shift:** Next, we consider novel subpopulation shifts, where the subpopulations
 249 in the train and test sets differ. For example, models might have only observed golden retrievers for
 250 the dog class but not huskies. In our experiments, we used the BREEDS benchmark [34], which
 251 leveraged the ImageNet[8] class hierarchy to create 4 datasets, Living-17, Nonliving-26, Entity13,
 252 and Entity-30.

253 **Natural Shift:** Finally, we consider non-simulated shifts in which the distribution shifts are induced
 254 through differences in the data collection process, such as ImageNet-V2 and CIFAR10-V2 proposed
 255 in [32]. We also include ImageNet-Sketch [37], which consists of sketched images of the original
 256 ImageNet classes. Additionally, we consider distribution shifts faced in the wild, such as ones
 257 curated in the WILDS benchmark [23]. We consider four WILDS datasets, two for language tasks
 258 (Amazon-WILDS, CivilComments-WILDS) and two for vision tasks (Camelyon17, RxRx1).

259 4.2 Architectures and Evaluations

260 For vision tasks (CIFAR10, CIFAR100, ImageNet, Living17, Nonliving26, Entity13, Entity30,
 261 Camelyon17-WILDS, RxRx1-WILDS), We trained ResNet18 and ResNet50 [17]; for language
 262 tasks (Amazon-WILDS, CivilComments-WILDS), we fine-tuned DistilBERT-base-uncased [33]. We

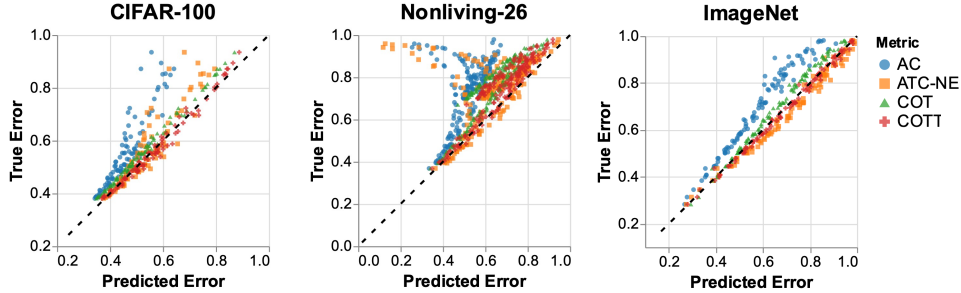


Figure 4: Qualitative results for AC, ATC, COT, and COTT, comparing error estimates vs. ground truth target error. Accurate estimates should be close to $y = x$ (dashed black line). Notably, COT and COTT remain accurate even when the shifts are large. By contrast, AC and ATC often severely underestimate the error, which is particularly evident in the Nonliving-26 dataset.

263 followed training setups from previous works [12] and provided the full details in the supplemental
 264 material. After training, we calibrated models using Temperature Scaling (TS) [16] on the in-
 265 distribution validation data, effectively adjusting the output probabilities of the neural network
 266 to match the actual correctness likelihood. This approach has previously demonstrated that TS
 267 consistently improves error estimation performance for all methods [12]. To evaluate different
 268 methods, we utilized the mean absolute difference between their predicted errors and the true errors,
 269 which are obtained using ground truth labels. We refer to this metric as Mean Absolute Error (MAE).

270 4.3 Results

271 We systematically evaluate our methods against an array of baselines, including *Average Confidence*
 272 (*AC*) [19], *Difference of Confidence* (DoC a.k.a. *DOC-Feat*) [15], *Importance Re-weighting* (IM)
 273 [6], *Generalized Disagreement Equality* (GDE) [21], *Average Thresholded Confidence* (ATC) [12],
 274 *ProjNorm* [38]. See the supplementary material for a review of detailed definitions.

275 In Table 1, we report the MAE results grouped by datasets and nature of shifts. Across all benchmarks,
 276 we observe that COT always obtains lower estimation error than AC, supporting our theoretical
 277 analysis that COT additionally leverages pseudo-label shift to fight miscalibration. On synthetic
 278 shift benchmarks (CIFAR10-Synthetic, CIFAR100-Synthetic, ImageNet-Synthetic BREEDS-same),
 279 COT is $2 - 4\times$ better than AC, drastically reducing the estimation error. On novel subpopulation
 280 shift (BREEDS-novel), COT cuts about half of the AC error. On natural shift, COT also improves
 281 upon AC. On ImageNet natural shift datasets, COT reduces the error from 8.5 to 3.88 compared to
 282 AC. COTT, presents the best overall results, surpassing the best current method (ATC) by a notable
 283 margin. On synthetic shift benchmarks, COTT is $2-3\times$ better than ATC-NE, the stronger version of
 284 ATC that uses a negative entropy score function. On novel subpopulation shift benchmarks, COTT is
 285 at least 4 absolute percent better than ATC-NE. On natural shift benchmarks, however, we observed
 286 mixed results. On ImageNet-Natural, COTT is better than ATC-NE while on CIFAR10-Natural,
 287 COTT is marginally worse. On WILDS benchmarks, no single method dominantly outperforms
 288 others. Note that the WILDS benchmark datasets contain label shift, meaning $P_S(y) \neq P_T(y)$. This
 289 leads COTT to overestimate the error on the CivilComments-WILDS. Nonetheless, COTT has the
 290 smallest worst-case error of 5.82 compared to ATC-NE’s 7.73, demonstrating its robustness even on
 291 distribution shifts faced in the wild.

292 In Figure 4, we use scatterplots to visualize estimations given by different methods, notably AC, ATC,
 293 COT, and COTT, where we plot the predicted error against the true error. Ideally, these scattered
 294 points should demonstrate a strong linear correlation and closely follow the $y = x$ line. While this
 295 is true for COT and COTT, AC and ATC often underestimate, sometimes by a large margin. In
 296 Nonliving-26, we can see data points representing shifts whose ATC predicted errors are around
 297 0.1 while true errors are close to 1. These observations corroborate the necessity of leveraging
 298 pseudo-label shifts to guard against such catastrophic failures of existing confidence-based methods.
 299 We include the scatterplots for all datasets in the supplemental material, where we show that COT
 300 and COTT successfully prevent severe underestimation seen in other methods.

References

- 301
- 302 [1] David Alvarez-Melis and Nicolo Fusi. Geometric dataset distances via optimal transport.
303 *Advances in Neural Information Processing Systems*, 33:21428–21439, 2020.
- 304 [2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial
305 networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- 306 [3] Christina Baek, Yiding Jiang, Aditi Raghunathan, and Zico Kolter. Agreement-on-the-line:
307 Predicting the performance of neural networks under distribution shift, 2023.
- 308 [4] Yogesh Balaji, Rama Chellappa, and Soheil Feizi. Normalized wasserstein for mixture distribu-
309 tions with applications in adversarial learning and domain adaptation. In *Proceedings of the*
310 *IEEE/CVF International Conference on Computer Vision*, pages 6500–6508, 2019.
- 311 [5] Krishnakumar Balasubramanian, Pinar Donmez, and Guy Lebanon. Unsupervised supervised
312 learning ii: Margin-based classification without labels. In *Proceedings of the Fourteenth Inter-*
313 *national Conference on Artificial Intelligence and Statistics*, pages 137–145. JMLR Workshop
314 and Conference Proceedings, 2011.
- 315 [6] Mayee Chen, Karan Goel, Nimit S Sohoni, Fait Poms, Kayvon Fatahalian, and Christopher Ré.
316 Mandoline: Model evaluation under distribution shift. In *International Conference on Machine*
317 *Learning*, pages 1617–1629. PMLR, 2021.
- 318 [7] Shai Ben David, Tyler Lu, Teresa Luu, and Dávid Pál. Impossibility theorems for domain
319 adaptation. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence*
320 *and Statistics*, pages 129–136. JMLR Workshop and Conference Proceedings, 2010.
- 321 [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
322 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
323 *recognition*, pages 248–255. Ieee, 2009.
- 324 [9] Weijian Deng and Liang Zheng. Are labels always necessary for classifier accuracy evaluation?
325 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
326 pages 15069–15078, 2021.
- 327 [10] Weijian Deng, Stephen Gould, and Liang Zheng. What does rotation prediction tell us about
328 classifier accuracy under varying testing environments? In *International Conference on Machine*
329 *Learning*, pages 2579–2589. PMLR, 2021.
- 330 [11] Pinar Donmez, Guy Lebanon, and Krishnakumar Balasubramanian. Unsupervised supervised
331 learning i: Estimating classification and regression errors without labels. *Journal of Machine*
332 *Learning Research*, 11(4), 2010.
- 333 [12] Saurabh Garg, Sivaraman Balakrishnan, Zachary C Lipton, Behnam Neyshabur, and Hanie
334 Sedghi. Leveraging unlabeled data to predict out-of-distribution performance. *arXiv preprint*
335 *arXiv:2201.04234*, 2022.
- 336 [13] Saurabh Garg, Nick Erickson, James Sharpnack, Alex Smola, Sivaraman Balakrishnan, and
337 Zachary Chase Lipton. RLSBench: A large-scale empirical study of domain adaptation under
338 relaxed label shift. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting Methods and*
339 *Applications*, 2022. URL <https://openreview.net/forum?id=kGgutmhd1H>.
- 340 [14] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann,
341 and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias
342 improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018.
- 343 [15] Devin Guillory, Vaishaal Shankar, Sayna Ebrahimi, Trevor Darrell, and Ludwig Schmidt. Pre-
344 dicting with confidence on unseen distributions. In *Proceedings of the IEEE/CVF International*
345 *Conference on Computer Vision*, pages 1134–1144, 2021.
- 346 [16] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
347 networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017.

- 348 [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
349 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
350 pages 770–778, 2016.
- 351 [18] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
352 corruptions and perturbations. In *International Conference on Learning Representations*, 2018.
- 353 [19] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
354 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 355 [20] Peter J Huber. Robust statistics. In *International encyclopedia of statistical science*, pages
356 1248–1251. Springer, 2011.
- 357 [21] Yiding Jiang, Vaishnavh Nagarajan, Christina Baek, and J Zico Kolter. Assessing generalization
358 of sgd via disagreement. *arXiv preprint arXiv:2106.13799*, 2021.
- 359 [22] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
360 *arXiv:1412.6980*, 2014.
- 361 [23] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
362 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee,
363 Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure
364 Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang.
365 WILDS: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine*
366 *Learning (ICML)*, 2021.
- 367 [24] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay
368 Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al.
369 Wilds: A benchmark of in-the-wild distribution shifts. In *International Conference on Machine*
370 *Learning*, pages 5637–5664. PMLR, 2021.
- 371 [25] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift
372 with black box predictors. In *International conference on machine learning*, pages 3122–3130.
373 PMLR, 2018.
- 374 [26] Xinran Liu, Yikun Bai, Yuzhe Lu, Andrea Soltoggio, and Soheil Kolouri. Wasserstein task
375 embedding for measuring task similarities. *arXiv preprint arXiv:2208.11726*, 2022.
- 376 [27] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
377 *arXiv:1711.05101*, 2017.
- 378 [28] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua
379 Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty?
380 evaluating predictive uncertainty under dataset shift. *Advances in neural information processing*
381 *systems*, 32, 2019.
- 382 [29] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
383 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative
384 style, high-performance deep learning library. *Advances in neural information processing*
385 *systems*, 32, 2019.
- 386 [30] Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data
387 science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- 388 [31] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do cifar-10
389 classifiers generalize to cifar-10? *arXiv preprint arXiv:1806.00451*, 2018.
- 390 [32] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet
391 classifiers generalize to imagenet? In *International conference on machine learning*, pages
392 5389–5400. PMLR, 2019.
- 393 [33] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version
394 of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.

- 395 [34] Shibani Santurkar, Dimitris Tsipras, and Aleksander Madry. Breeds: Benchmarks for subpopu-
396 lation shift. *arXiv preprint arXiv:2008.04859*, 2020.
- 397 [35] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the
398 log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- 399 [36] Cédric Villani et al. *Optimal transport: old and new*, volume 338. Springer, 2009.
- 400 [37] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global represen-
401 tations by penalizing local predictive power. In *Advances in Neural Information Processing*
402 *Systems*, pages 10506–10518, 2019.
- 403 [38] Yaodong Yu, Zitong Yang, Alexander Wei, Yi Ma, and Jacob Steinhardt. Predicting out-of-
404 distribution error with the projection norm. *arXiv preprint arXiv:2202.05834*, 2022.

405 **A Deferred Proofs**

406 For readers' convenience, we review the statements of the propositions and corollaries and provide
407 the full proofs below.

408 **A.1 Proof of Proposition 1**

409 **Proposition 1** ($\hat{\epsilon}_{AC}$ - W_∞ Equivalence). *Let $(\mathcal{P}(\mathbb{R}^k), W_\infty)$ be the metric space of all distributions
410 over \mathbb{R}^k , where W_∞ is the Wasserstein distance with $c(x, y) = \|x - y\|_\infty$. Then, the estimated error
411 of AC-MC is given by $\hat{\epsilon}_{AC} = W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y}))$.*

412 *Proof.* We first show the following equality. Let $j^* = \arg \max_j \vec{f}_j(x)$

$$1 - \vec{f}_{j^*}(x) = \|\vec{y} - \vec{f}(x)\|_\infty \quad (1)$$

413 where $\vec{y}_j = \mathbb{1}[j = j^*]$. Let $\vec{f}_{-j^*}(x)$ denote the vector $\vec{f}(x)$ with j^* -th element removed. Since for a
414 confidence vector $\|\vec{f}(x)\|_1 = 1$,

$$1 - \vec{f}_{j^*}(x) = \|\vec{f}_{-j^*}(x)\|_1 \geq \|\vec{f}_{-j^*}(x)\|_\infty$$

415 Therefore, we have obtained the desired Equality 1:

$$\|\vec{y} - \vec{f}(x)\|_\infty = \max\{\|\vec{f}_{-j^*}(x)\|_\infty, 1 - \vec{f}_{j^*}(x)\} = 1 - \vec{f}_{j^*}(x)$$

416 Next, we consider the optimal transport plan between $\vec{f}_\# P(\vec{c})$ and $P_{\text{pseudo}}(\vec{y})$. Namely, we show *all*
417 confidence vectors $\vec{f}(x^{(i)})$ are coupled with their one-hot pseudo-labels $\vec{y}^{(i)}$. This can be observed by
418 the fact that the one-hot pseudo-label is the one-hot label that achieves the lowest L-infinity cost, i.e.

$$\|\vec{f}(x^{(i)}) - \vec{y}^{(i)}\|_\infty \leq \|\vec{f}(x^{(i)}) - \vec{y}'\|_\infty, \forall \vec{y}' \in \{0, 1\}^k \cap \Delta^{k-1}$$

419 Suppose there exist confidence vectors that are not coupled with their one-hot pseudo-labels, then *all*
420 individual costs are suboptimal and the total cost is suboptimal as well, contradicting the assumption
421 that the transport plan is optimal. Therefore,

$$W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y})) = \frac{1}{n} \sum_{i=1}^n \|\vec{f}(x^{(i)}) - \vec{y}^{(i)}\|_\infty \quad (2)$$

422 Combining Equality 1 and Equality 2, we obtain the desired relationship between AC error estimate
423 and W_∞ distance

$$\hat{\epsilon}_{AC} = \frac{1}{n} \sum_{i=1}^n (1 - \max_j \vec{f}_j(x^{(i)})) = \frac{1}{n} \sum_{i=1}^n \|\vec{f}(x^{(i)}) - \vec{y}^{(i)}\|_\infty = W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y}))$$

424 □

425 **A.2 Proof of Corollary 1**

426 **Corollary 1** ($P_{\text{pseudo}}(\vec{y})$ is closest to $\vec{f}_\# P(\vec{c})$). *Let $P'(\vec{y}) \in \mathcal{P}(\{0, 1\}^k \cap \Delta^{k-1})$ be a one-hot label
427 distribution. Then $W_\infty(\vec{f}_\# P(\vec{c}), P'(\vec{y})) \geq W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y}))$.*

428 *Proof.* We first show the following equality, which establishes the relationship between AC accuracy
429 estimate with W_∞ distance:

$$1 - \hat{\epsilon}_{AC} = W_\infty(\vec{f}_\# P(\vec{c}), \delta_0) \quad (3)$$

430 Since $W_\infty(\vec{f}_\# P(\vec{c}), \delta_0)$ transports $\vec{f}_\# P(\vec{c})$ to δ_0 , the optimal transport plan couples every element
431 in $\vec{f}_\# P(\vec{c})$ to 0. For each $x^{(i)}$, its confidence vector $\vec{f}_\# P(\vec{c})$ has a transport cost $\|\vec{f}(x^{(i)}) - 0\|_\infty$.
432 Hence,

$$1 - \hat{\epsilon}_{AC} = \frac{1}{n} \sum_{i=1}^n \max_j \vec{f}_j(x^{(i)}) = \frac{1}{n} \sum_{i=1}^n \|\vec{f}(x^{(i)}) - 0\|_\infty = W_\infty(\vec{f}_\# P(\vec{c}), \delta_0)$$

433 With this, our inequality is simply the Triangle Inequality in $(\mathcal{P}(\mathbb{R}^k), W_\infty)$,

$$W_\infty(\vec{f}_\# P(\vec{c}), \delta_0) + W_\infty(\vec{f}_\# P(\vec{c}), P'(\vec{y})) \geq W_\infty(P'(\vec{y}), \delta_0) = 1$$

434 Combined with Equation 3, we obtain the desired inequality

$$W_\infty(\vec{f}_\# P(\vec{c}), P_{\text{pseudo}}(\vec{y})) = 1 - W_\infty(\vec{f}_\# P(\vec{c}), \delta_0) \leq W_\infty(\vec{f}_\# P(\vec{c}), P'(\vec{y}))$$

435

□

436 A.3 Proof of Proposition 2

437 **Notations:** Let $\mathcal{C}(\vec{c}) = \{\vec{c}' \in \Delta^{k-1} \mid \arg \max_j \vec{c}'_j = \arg \max_j \vec{c}_j\}$ be the set of confidence vectors
 438 whose one-hot pseudo-labels that match with that of $\vec{c} \in \Delta^{k-1}$. Let $\mathcal{P}(\Delta^{k-1})$ be the set of
 439 all distributions of confidence vectors and $\mathcal{P}_c[P_{\text{pseudo}}(\vec{y})] = \{P'(\vec{c}) \in \mathcal{P}(\Delta^{k-1}) \mid P_{\text{pseudo}}(\vec{y}) =$
 440 $P'(\arg \max_j \vec{c}_j = \arg \max_j \vec{y}_j)\}$ be the set of distributions of confidence vectors that share the
 441 same pseudo-label distribution $P_{\text{pseudo}}(\vec{y})$.

442 $\mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]$ defines an equivalence class for the space of distributions of confidence vectors
 443 $(\mathcal{P}(\Delta^{k-1}), W_\infty)$ that share the same pseudo-label distribution $P_{\text{pseudo}}(\vec{y})$. Pictorially, in Figure 1,
 444 $\mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]$ represents the line between δ_0 and $P_{\text{pseudo}}(\vec{y})$. On this line, every distribution of
 445 confidence vectors shares the same pseudo-label distribution $P_{\text{pseudo}}(\vec{y})$.

446 To prove Proposition 2, we need the following lemma, which intuitively allows us to change the
 447 metric from measuring the distance between two points to the distance between an equivalence class
 448 and a point.¹

449 **Lemma 1** (Change-of-metric). *Let $\vec{y}, \vec{y}' \in \{0, 1\}^k \cap \Delta^{k-1}$ be two one-hot labels. Then the following*
 450 *holds*

$$\inf_{\vec{c} \in \mathcal{C}(\vec{y})} \|\vec{c} - \vec{y}'\|_\infty = 0.5 \times \mathbb{1}[\vec{y} \neq \vec{y}']$$

451

452 *Proof.* If $\vec{y} = \vec{y}'$, then we know the optimal $\vec{c} = \vec{y}$

$$\inf_{\vec{c} \in \mathcal{C}(\vec{y})} \|\vec{c} - \vec{y}'\|_\infty = \|\vec{y} - \vec{y}'\|_\infty = 0$$

453 If $\vec{y} \neq \vec{y}'$, then we proceed by showing equality with two inequalities. First, observe $\{(0.5 + \delta)\vec{y} +$
 454 $(0.5 - \delta)\vec{y}' \mid \delta \in (0, 0.5)\} \subset \mathcal{C}(\vec{y})$.

$$\inf_{\vec{c} \in \mathcal{C}(\vec{y})} \|\vec{c} - \vec{y}'\|_\infty \leq \inf_{\delta \in (0, 0.5)} \|(0.5 + \delta)\vec{y} + (0.5 - \delta)\vec{y}' - \vec{y}'\|_\infty = \inf_{\delta \in (0, 0.5)} (0.5 + \delta) = 0.5$$

455 If $\|\vec{c} - \vec{y}'\|_\infty < 0.5$, $\arg \max_j \vec{c}_j = \arg \max_j \vec{y}'_j \neq \arg \max_j \vec{y}_j$, i.e. $\vec{c} \notin \mathcal{C}(\vec{y})$. Therefore,
 456 $\inf_{\vec{c} \in \mathcal{C}(\vec{y})} \|\vec{c} - \vec{y}'\|_\infty \geq 0.5$, which further implies $\inf_{\vec{c} \in \mathcal{C}(\vec{y})} \|\vec{c} - \vec{y}'\|_\infty = 0.5$. □

457 We are now in a position to prove Proposition 2, which follows from the somewhat surprising fact
 458 that the left-hand side of the inequality is simply the distance between $\vec{f}_\# P_T(\vec{c})$ and $P_T(\vec{y})$ with a
 459 change-of-metric to the metric defined above.

460 **Proposition 2** (Calibration independent lower bound of COT). *Under the assumption that $P_T(\vec{y}) =$*
 461 *$P_S(\vec{y})$, we always have $\hat{e}_{\text{COT}} \geq 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$.*

462 *Proof.* Since $P_T(\vec{y}) = P_S(\vec{y})$,

$$\begin{aligned} \hat{e}_{\text{COT}} &= W_\infty(\vec{f}_\# P_T(\vec{c}), P_T(\vec{y})) \\ &= \inf_{\pi(\vec{c}, \vec{y}) \in \Pi(\vec{f}_\# P_T(\vec{c}), P_T(\vec{y}))} \int \|\vec{c} - \vec{y}\|_\infty d\pi(\vec{c}, \vec{y}) \\ &\geq \inf_{\pi(\vec{c}, \vec{y}) \in \Pi(\vec{f}_\# P_T(\vec{c}), P_T(\vec{y}))} \int \inf_{\vec{c}' \in \mathcal{C}(\vec{c})} \|\vec{c}' - \vec{y}\|_\infty d\pi(\vec{c}, \vec{y}) \end{aligned} \quad (4)$$

$$= \inf_{\pi(\vec{y}', \vec{y}) \in \Pi(P_{\text{pseudo}}(\vec{y}'), P_T(\vec{y}))} \int \inf_{\vec{c}' \in \mathcal{C}(\vec{y}')} \|\vec{c}' - \vec{y}\|_\infty d\pi(\vec{y}', \vec{y}) \quad (5)$$

¹This is closely related to the Hausdorff distance between sets in a metric space.

463 Equation 5 follows from the observation that $\mathcal{C}(\vec{c}) = \mathcal{C}(\vec{y}')$ for a confidence vector \vec{c} and its corre-
 464 sponding one-hot pseudo-label \vec{y}' . Furthermore, since our new metric $\inf_{\vec{c}' \in \mathcal{C}(\vec{y}')} \|\vec{c}' - \vec{y}'\|_\infty$ is only
 465 defined up to the equivalence class, replacing each $\vec{c}' \in \mathcal{C}(\vec{c})$ with its pseudo-label \vec{y}' does not change
 466 the distance.

467 Plugging in Lemma 1,

$$\begin{aligned} \hat{\epsilon}_{\text{COT}} &\geq \inf_{\pi(\vec{y}', \vec{y}) \in \Pi(P_{\text{pseudo}}(\vec{y}'), P_T(\vec{y}))} \int 0.5 \times \mathbb{1}[\vec{y}' \neq \vec{y}] d\pi(\vec{y}', \vec{y}) \\ &= 0.5 \inf_{\pi(\vec{y}', \vec{y}) \in \Pi(P_{\text{pseudo}}(\vec{y}'), P_T(\vec{y}))} \int \|\vec{y}' - \vec{y}\|_\infty d\pi(\vec{y}', \vec{y}) \\ &= 0.5 W_\infty(P_{\text{pseudo}}(\vec{y}'), P_T(\vec{y})) \end{aligned}$$

468

□

469 A.4 Tightness of Proposition 2

470 While Inequality 4 seems loose, our bound is, in fact, tight if no further assumptions on the calibration
 471 status of the classifier \vec{f} are made. We need the following lemma that establishes the relationship
 472 between pseudo-label shift and the total variation distance between target label distribution and
 473 pseudo-label distribution. Note this equivalence only makes sense in the context of measuring W_∞
 474 distance between two one-hot label distributions, but not under other contexts presented in the paper.

Lemma 2 (Pseudo-label shift is total variation).

$$W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) = \|P_{\text{pseudo}}(\vec{y}) - P_T(\vec{y})\|_{\text{TV}}$$

475 *Proof.* For two $\vec{y}, \vec{y}' \in \{0, 1\}^k \cap \Delta^{k-1}$, the transport cost $c(\vec{y}, \vec{y}') = \mathbb{1}[\vec{y} \neq \vec{y}']$. Then, the standard
 476 result on optimal transport [36] gives the desired equality. □

Corollary 2.

$$W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) = 1 - \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\}$$

477

Proof.

$$\begin{aligned} W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) &= \|P_{\text{pseudo}}(\vec{y}) - P_T(\vec{y})\|_{\text{TV}} \\ &= \frac{1}{2} \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} |P_{\text{pseudo}}(\vec{y}) - P_T(\vec{y})| \\ &= \frac{1}{2} \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \max\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} - \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} \\ &= \frac{1}{2} \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \max\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} + \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} \\ &\quad - \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} \\ &= \frac{1}{2} \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} P_{\text{pseudo}}(\vec{y}) + P_T(\vec{y}) - \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} \\ &= 1 - \sum_{\vec{y} \in \{0, 1\}^k \cap \Delta^{k-1}} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\} \end{aligned}$$

478

□

479 Finally, we show Proposition 2 is tight by constructing a sequence of distributions of confidence
 480 vectors, the limit of which is exactly $0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$ away from $P_T(\vec{y})$.

Lemma 3 (Proposition 2 is tight).

$$\inf_{P(\vec{c}) \in \mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]} W_\infty(P(\vec{c}), P_T(\vec{y})) = 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$$

481 *Proof.* First, we construct the following family of distributions $\{P_\delta(\vec{c})|\delta \in (0, 0.5]\}$, where $P_\delta(\vec{c})$ is
482 the following mixture distribution

$$P_\delta(\vec{c}) = \gamma P_\cap(\vec{y}) + (1 - \gamma)P_\times(\vec{t})$$

483 where $P_\cap(\vec{y}) = \gamma^{-1} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\}$, $\gamma = \sum_{\vec{y} \in \{0,1\}^k \cap \Delta^{k-1}} \min\{P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})\}$,
484 $P_\times(\vec{t})$ is a distribution supported on $\Delta^{k-1} \cap \{0.5 + \delta, 0.5 - \delta, 0\}^k$ (i.e. one element in \vec{t} is $0.5 + \delta$,
485 another is $0.5 - \delta$, and the rest are 0). Additionally, $P_\times(\vec{t}_i = 0.5 + \delta) = P_{\text{pseudo}}(\vec{y}_i = 1)$ and
486 $P_\times(\vec{t}_i = 0.5 - \delta) = P_T(\vec{y}_i = 1)$. It is easy to check that $P_\delta(\vec{c}) \in \mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]$.

487 Next, we construct an explicit transport plan $\pi(\vec{c}, \vec{y}) \in \Pi(P_\delta(\vec{c}), P_T(\vec{y}))$. We construct it via the
488 factorization $\pi(\vec{c}, \vec{y}) = P_\delta(\vec{c})\pi(\vec{y}|\vec{c})$, where

$$\pi(\vec{y}|\vec{c}) = \begin{cases} 1 & \text{if } \vec{c} = \vec{y} \text{ or } \langle \vec{c}, \vec{y} \rangle = 0.5 - \delta \\ 0 & \text{otherwise} \end{cases}$$

489 The cost of this transport plan is therefore

$$\int \|\vec{c} - \vec{y}\|_\infty d\pi(\vec{c}, \vec{y}) = (0.5 + \delta)(1 - \gamma) = (0.5 + \delta)W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$$

490 where the last equality follows from Lemma 2. Taking infimum,

$$\inf_{\delta \in (0, 0.5]} \int \|\vec{c} - \vec{y}\|_\infty d\pi(\vec{c}, \vec{y}) = 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$$

491 Combining everything so far, we obtain the desired result:

$$\begin{aligned} 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) &= \inf_{\delta \in (0, 0.5]} \int \|\vec{c} - \vec{y}\|_\infty d\pi(\vec{c}, \vec{y}) \\ &\geq \inf_{\delta \in (0, 0.5]} W_\infty(P_\delta(\vec{c}), P_T(\vec{y})) \end{aligned} \quad (6)$$

$$\geq \inf_{P(\vec{c}) \in \mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]} W_\infty(P(\vec{c}), P_T(\vec{y})) \quad (7)$$

$$\geq 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) \quad (8)$$

492 Inequality 6 follows from the fact that the optimal transport plan cannot have a greater cost than our
493 explicit plan π . Inequality 7 is due to the fact that the family of distribution we are considering is
494 a subset of $\mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]$. Inequality 8 is an application of the lower bound 4 which holds for all
495 $P(\vec{c}) \in \mathcal{P}_c[P_{\text{pseudo}}(\vec{y})]$. \square

496 B Experiment Baselines

497 We consider an array of baselines to compare against our methods: COT and COTT.

498 *Average Confidence (AC)* estimates target error by taking the average of one minus the maximum
499 softmax confidence of target data. $\hat{\epsilon}_{\text{AC}} = \mathbb{E}_{x \sim \hat{P}_T}[1 - \max_{j \in \mathcal{Y}} \vec{f}_j(x)]$.

500 *Difference of Confidence (DoC a.k.a. DOC-Feat)* [15] estimates target error through the dif-
501 ference between the confidence of source data and the confidence of target data. $\hat{\epsilon}_{\text{DoC}} =$
502 $\mathbb{E}_{x \sim \hat{P}_S}[\mathbb{1}[\arg \max_{j \in \mathcal{Y}} \vec{f}_j(x) \neq y]] + \mathbb{E}_{x \sim \hat{P}_T}[1 - \max_{j \in \mathcal{Y}} \vec{f}_j(x)] - \mathbb{E}_{x \sim \hat{P}_S}[1 - \max_{j \in \mathcal{Y}} \vec{f}_j(x)]$.

503 *Importance Re-weighting (IM)* estimates target error as a re-weighted source error. The weights are
504 calculated as the ratio between the number of data points in each bin in target data and source data.
505 This is equivalent to [6] using one slice based on the underlying classifier confidence.

506 *Generalized Disagreement Equality (GDE)* [21] estimates target error as the disagreement ratio of
 507 predictions on the target data using two independently trained models $\vec{f}(x)$ and $\vec{f}^l(x)$. $\hat{\epsilon}_{\text{GDE}} =$
 508 $\mathbb{E}_{x \sim \hat{P}_T} [\mathbb{1}[\arg \max_{j \in \mathcal{Y}} \vec{f}_j(x) \neq \arg \max_{j \in \mathcal{Y}} \vec{f}_j^l(x)]]$.

509 *Average Thresholded Confidence (ATC)* [12] first identifies a threshold t such that the fraction of
 510 source data points that have scores below the threshold matches the source error on in-distribution
 511 validation data. Target error is estimated as the expected number of target data points that fall
 512 below the identified threshold. $\hat{\epsilon}_{\text{ATC}}(s) = \mathbb{E}_{x \sim \hat{P}_T} [\mathbb{1}[s(\vec{f}(x)) < t]]$, where s is the score function
 513 mapping the softmax vector to a scalar. Two different score functions are used, Maximum Confidence
 514 (ATC-MC) and Negative Entropy (ATC-NE).

515 In addition, we also compare our method to *ProjNorm* [38]. ProjNorm cannot provide a direct
 516 estimate, instead, the authors demonstrated their metric has the strongest linear correlation to true
 517 target error compared to existing baselines. In this case, we followed their setup and performed a
 518 correlation analysis to draw a direct comparison. We included results in the supplemental material
 519 and showed that our method consistently outperforms ProjNorm.

520 C Extended Results

521 C.1 Results with Standard Deviation

522 We show the full experimental results with standard deviation in Table 2.

523 C.2 Qualitative Results

524 We show the qualitative results (scatter plots) in Fig 5

525 C.3 Correlation Analysis

526 ProjNorm [38] leverages pseudo labels on the target domain to retrain a copy of the reference
 527 model trained on the source domain. The authors show that the difference between the two models’
 528 parameters has a strong linear correlation to the true target error. Following the paper’s experimental
 529 setup, we conducted the correlation analysis on CIFAR10 and CIFAR100 using three architectures,
 530 ResNet18, ResNet50, and VGG11. We note that ProjNorm in fact implicitly leverages the assumption
 531 that $P_T(y) = P_S(y)$ as this condition holds for both CIFAR10 and CIFAR100. As Fig. 18 of their
 532 paper [38] shows, ProjNorm tends to overestimate when label shift exists.

533 C.4 Mild Label Shift

534 We motivate our methods under the assumption of no label shift. In Proposition 2, we showed that
 535 the worst-case underestimate of COT is half of the pseudo-label shift. Under mild label shifts, the
 536 guarantee for such worst-case underestimation becomes weaker. This can be observed from the
 537 following corollary of Proposition 2:

Corollary 3 (Calibration independent lower bound of COT under *mild* label shift).

$$\hat{\epsilon}_{\text{COT}} \geq 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y})) - W_\infty(P_S(\vec{y}), P_T(\vec{y}))$$

538 *Proof.* By Triangle Inequality in $(\mathcal{P}(\mathbb{R}^k), W_\infty)$,

$$W_\infty(\vec{f}_\# P_T(\vec{c}), P_S(\vec{y})) + W_\infty(P_S(\vec{y}), P_T(\vec{y})) \geq W_\infty(\vec{f}_\# P_T(\vec{c}), P_T(\vec{y}))$$

539 Combined with Proposition 2, we obtain the desired result. \square

540 As the label shift increases, we have a weaker guarantee of the worst-case underestimation error of
 541 COT as long as $W_\infty(P_S(\vec{y}), P_T(\vec{y})) \leq 0.5W_\infty(P_{\text{pseudo}}(\vec{y}), P_T(\vec{y}))$. However, we perform additional
 542 controlled experiments which suggest our methods remain to be the most performant despite the
 543 theoretical guarantee is not as strong as the case without label shift.

544 To simulate mild label shift for datasets with $P_S(\vec{y}) = P_T(\vec{y})$, we first calculate the original target
 545 marginal and then sample the shifted target marginal from a Dirichlet distribution as in [13] with a

Table 2: Mean Absolute Error (MAE) between the estimated error and ground truth error to compare different methods. The "shift" column denotes the nature of distribution shifts for each dataset. For vision datasets, we reported results for ResNet18 and ResNet50; for language datasets, we reported results for DistilBERT-base-uncased. The results are averaged over 3 random seeds. We highlight the best-performing method. The number in the parentheses denotes the standard deviation.

Dataset	Shift	Baselines						Ours	
		AC	DoC	IM	GDE	ATC-MC	ATC-NE	COT	COTT
CIFAR10	Natural	5.97 (0.10)	5.38 (0.08)	5.87 (0.09)	5.9 (0.15)	3.38 (0.14)	3.15 (0.28)	5.41 (0.09)	3.33 (0.13)
	Synthetic	9.1 (0.25)	8.53 (0.28)	9.26 (0.35)	8.84 (0.11)	4.2 (0.38)	3.37 (0.30)	2.17 (0.09)	1.7 (0.26)
CIFAR100	Synthetic	10.83 (0.08)	8.76 (0.22)	12.07 (0.37)	11.36 (0.25)	6.8 (0.39)	6.63 (0.43)	2.09 (0.27)	2.59 (0.01)
ImageNet	Natural	8.5 (0.39)	7.43 (0.41)	8.62 (0.47)	5.62 (0.33)	3.57 (0.46)	2.6 (0.66)	3.88 (0.04)	2.41 (0.11)
	Synthetic	10.34 (0.83)	9.28 (0.86)	12.87 (0.77)	6.54 (0.37)	1.59 (0.08)	3.41 (0.53)	3.24 (0.28)	1.42 (0.29)
Entity13	Same	19.63 (2.17)	19.2 (2.51)	17.5 (0.90)	15.37 (1.06)	8.09 (0.49)	7.23 (0.49)	8.47 (0.66)	2.61 (0.31)
	Novel	29.61 (2.61)	29.18 (2.95)	27.22 (1.08)	24.48 (0.61)	14.54 (0.95)	9.49 (0.70)	15.9 (0.80)	5.46 (0.75)
Entity30	Same	16.97 (0.35)	16.21 (0.36)	13.56 (2.53)	13.98 (0.26)	8.19 (1.07)	9.08 (0.42)	5.9 (0.29)	2.46 (0.65)
	Novel	27.57 (0.06)	26.81 (0.61)	23.96 (2.79)	23.4 (0.1)	13.46 (2.55)	8.57 (2.2)	15.11 (0.38)	5.94 (1.17)
Living17	Same	14.84 (3.36)	14.67 (3.30)	11.22 (2.06)	9.94 (0.48)	4.88 (0.42)	5.43 (1.06)	6.25 (1.91)	2.94 (1.21)
	Novel	29.61 (3.76)	29.18 (3.71)	27.22 (3.45)	24.48 (0.74)	14.54 (2.87)	9.49 (3.25)	15.9 (2.04)	5.53 (1.93)
Nonliving26	Same	19.25 (2.45)	18.43 (3.13)	16.6 (0.96)	12.77 (0.85)	11.18 (2.77)	9.69 (0.70)	7.06 (1.17)	3.34 (0.90)
	Novel	31.37 (2.99)	30.54 (3.65)	28.79 (1.47)	23.37 (0.61)	19.93 (4.02)	16.56 (1.28)	17.8 (1.53)	10.46 (3.08)
Camelyon17-WILDS	Natural	9.44 (0.50)	9.44 (0.49)	10.24 (0.38)	5.19 (0.44)	7.73 (0.72)	7.73 (0.72)	7.27 (0.57)	5.71 (0.94)
RxRx1-WILDS	Natural	5.21 (0.26)	8.44 (0.15)	8.09 (0.16)	7.48 (0.26)	6.53 (0.10)	6.86 (0.28)	3.25 (0.16)	5.82 (0.31)
Amazon-WILDS	Natural	2.62 (0.16)	2.35 (0.06)	2.34 (0.06)	17.04 (0.84)	1.63 (0.1)	1.54 (0.11)	2.43 (0.04)	2.01 (0.42)
CivilCom.-WILDS	Natural	1.54 (0.23)	0.96 (0.19)	0.86 (0.20)	8.7 (0.14)	2.3 (0.34)	2.3 (0.34)	1.23 (0.05)	4.68 (0.39)

546 parameter $\alpha = 50$. The parameter α controls the severity of the label shift, and a smaller α means
547 a larger label shift. Concretely, let the shifted target marginal be $P_{\tilde{T}}(\tilde{y})$. Then $P_{\tilde{T}}(\tilde{y}) \sim \text{Dir}(\beta)$
548 where $\beta_{(\tilde{y})} = \alpha \cdot P_T(\tilde{y})$. Finally, based on $P_{\tilde{T}}(\tilde{y})$, we sample a new set of test samples for which we
549 estimate the performance. We conducted this mild label shift experiment for CIFAR10, CIFAR100,
550 ImageNet, Living17, Nonliving26, Entity13, and Entity30 as these datasets have the same source and
551 target marginal. We showed the results in Table 3. As we can see, our methods still dominate existing
552 methods under this relaxed condition.

553 C.5 When does thresholding improve over averaging?

554 In this section, we provide some intuitions on when using a threshold provides better estimates than
555 taking the average. From Fig. 6, we show that thresholding yields larger and more accurate error
556 estimates when the cost distribution on the OOD data is more spread out and less concentrated around
557 0. By contrast, when the cost distribution is mostly near 0, thresholding leads to similar estimates as
558 averaging. Interestingly, even on OOD data where the model has very low performance, there is still
559 a decent amount of samples whose cost is near 0. Thus, when taking the average, we will end up

Table 3: Mean Absolute Error (MAE) between the estimated error and ground truth error to compare different methods under mild label shift. The results are averaged over 3 random seeds. We highlight the best-performing method. The number in the parentheses denotes the standard deviation.

Dataset	Shift	Baselines						Ours	
		AC	DoC	IM	GDE	ATC-MC	ATC-NE	COT	COTT
CIFAR10	Natural	5.58 (0.26)	4.99 (0.23)	5.50 (0.22)	5.69 (0.04)	2.76 (0.32)	2.47 (0.43)	3.75 (0.28)	1.68 (0.34)
	Synthetic	8.67 (0.29)	8.10 (0.31)	8.82 (0.38)	8.47 (0.15)	3.93 (0.38)	3.13 (0.32)	2.76 (0.04)	4.0 (0.30)
CIFAR100	Synthetic	10.89 (0.15)	8.85 (0.22)	12.14 (0.37)	11.33 (0.23)	6.93 (0.44)	6.76 (0.48)	1.89 (0.30)	2.81 (0.07)
ImageNet	Natural	8.36 (0.37)	7.29 (0.42)	8.46 (0.48)	5.54 (0.36)	3.53 (0.48)	2.47 (0.71)	3.74 (0.20)	2.05 (0.26)
	Synthetic	10.26 (0.83)	9.19 (0.86)	12.79 (0.77)	6.50 (0.40)	1.61 (0.08)	3.51 (0.52)	3.03 (0.25)	1.75 (0.33)
Entity13	Same	15.50 (0.38)	14.60 (0.34)	15.49 (0.22)	15.18 (1.03)	8.51 (0.80)	7.40 (0.57)	4.59 (0.23)	3.24 (0.12)
	Novel	24.39 (0.23)	23.49 (0.19)	24.56 (0.05)	23.48 (0.62)	14.99 (0.82)	12.45 (0.64)	11.19 (0.34)	4.6 (0.38)
Entity30	Same	15.46 (0.70)	13.93 (0.65)	15.55 (0.74)	13.83 (0.27)	8.80 (0.64)	8.26 (0.83)	4.75 (0.29)	2.16 (0.15)
	Novel	25.98 (0.53)	24.45 (0.46)	26.72 (0.68)	23.28 (0.14)	15.56 (0.56)	13.21 (0.90)	13.96 (0.17)	7.07 (0.27)
Living17	Same	11.38 (0.67)	10.90 (0.48)	11.83 (1.31)	9.85 (0.41)	4.46 (0.31)	4.39 (0.18)	4.40 (0.34)	2.71 (0.81)
	Novel	25.72 (0.46)	25.13 (0.76)	26.32 (1.98)	21.61 (0.68)	14.09 (2.31)	11.48 (1.99)	16.94 (0.78)	9.31 (1.79)
Nonliving26	Same	16.28 (0.37)	14.48 (0.29)	15.69 (0.19)	12.88 (0.84)	9.63 (0.43)	9.69 (0.66)	5.33 (0.73)	2.18 (0.23)
	Novel	27.93 (0.08)	26.13 (0.25)	27.76 (0.29)	23.25 (0.69)	18.15 (0.38)	16.08 (0.38)	15.66 (0.45)	8.71 (0.42)

Table 4: Coefficients of determination (R^2) and rank correlations (ρ) to measure the linear correlation between a method’s output quantity and the true target error (the higher the better). COT achieves superior performance than all existing methods across different models and datasets.

Dataset	Network	AC		Entropy		GDE		ATC		ProjNorm		COT	
		R^2	ρ	R^2	ρ	R^2	ρ	R^2	ρ	R^2	ρ	R^2	ρ
CIFAR10	ResNet18	0.825	0.980	0.862	0.982	0.842	0.981	0.875	0.987	0.947	0.988	0.996	0.998
	ResNet50	0.950	0.995	0.949	0.995	0.959	0.995	0.885	0.989	0.936	0.989	0.993	0.996
	VGG11	0.710	0.938	0.762	0.958	0.723	0.948	0.548	0.851	0.756	0.949	0.994	0.993
	Average	0.828	0.971	0.858	0.978	0.841	0.975	0.769	0.942	0.880	0.975	0.994	0.996
CIFAR100	ResNet18	0.943	0.987	0.932	0.984	0.950	0.988	0.927	0.985	0.969	0.974	0.995	0.997
	ResNet50	0.957	0.987	0.948	0.984	0.962	0.989	0.955	0.991	0.982	0.991	0.992	0.996
	VGG11	0.794	0.959	0.821	0.973	0.870	0.978	0.736	0.975	0.653	0.849	0.996	0.997
	Average	0.898	0.978	0.900	0.980	0.927	0.985	0.873	0.984	0.868	0.938	0.994	0.997

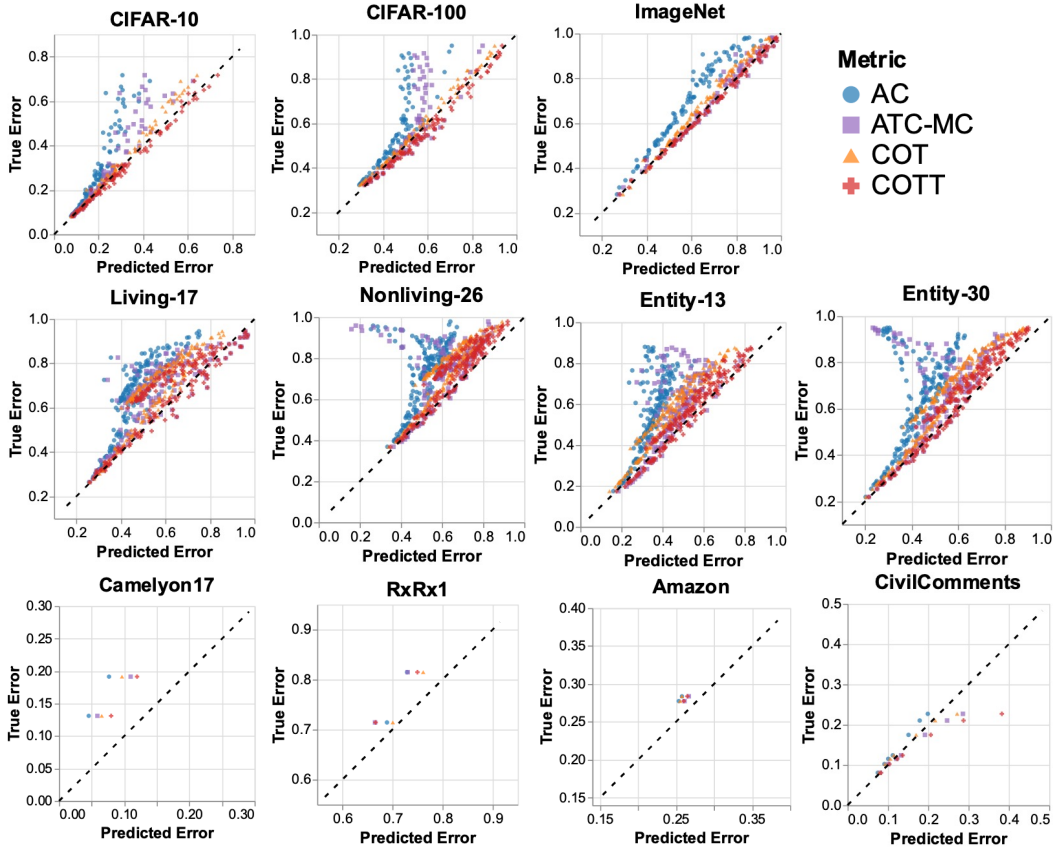


Figure 5: Qualitative results for AC, ATC, COT, and COTT. In these scatterplots, the x-axis is the target error estimate and the y-axis is the ground truth target error. Accurate estimates should be close to $y = x$ (dashed black line). We can see that for all datasets, COT and COTT avoid the severe underestimation seen on ATC.

560 with a smaller value which suggests a low error. In these cases, thresholding will give larger error
 561 estimates than averaging.

562 D Datasets

563 **CIFAR10:** The synthetic shifts included 19 common visual corruptions across 5 levels of severity
 564 from [18]. The natural shift is CIFAR10-V2 [31].

565 **CIFAR100:** The synthetic shifts included 19 common visual corruptions across 5 levels of severity
 566 from [18].

567 **ImageNet:** The synthetic shifts included 19 common visual corruptions across 5 levels of severity
 568 from [18]. The natural shifts include 4 datasets from ImageNet-V2 [32] and ImageNet-Sketch [37].

569 **BREEDS:** The BREEDS benchmark contains 4 datasets, Living-17, Nonliving26, Entity13,
 570 Entity30. For each of the datasets, the same subpopulation shifts include the corrupted versions of
 571 the test set with the same subpopulation; the novel subpopulation shifts include the clean as well as
 572 corrupted versions [18] of the test set with novel subpopulation.

573 **WILDS:** For all WILDS datasets, we used the official OOD datasets provided in their paper [24].

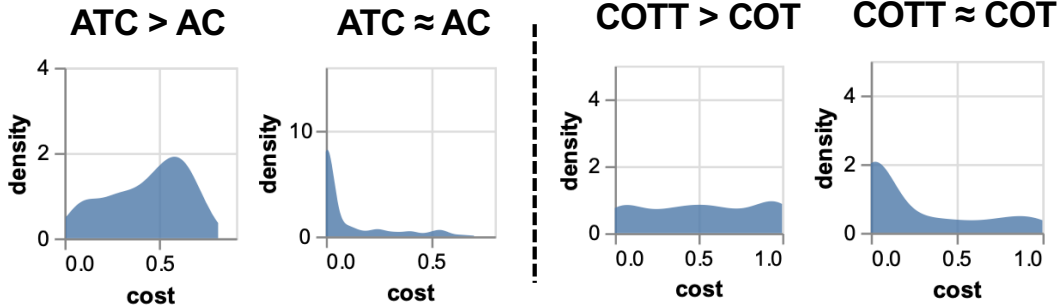


Figure 6: We demonstrate cases where using thresholding improves over taking averages. The x-axis denotes the max norm between a confidence vector and the corresponding one-hot label. For AC and ATC-MC, the corresponding label is always the argmax of the confidence vector as mentioned in section 2.3. For COT and COTT, the corresponding label is the one matched via optimal transport. We observe that thresholding improves over averaging when the cost distribution is less concentrated around 0, which corresponds to situations where the model is very confident on most samples.

574 E Experiment Setup

575 We performed training in PyTorch [29], and we used RTX 6000 Ada GPUs.

576 For datasets without an official validation set, we randomly sampled a subset of the official training
 577 set as the validation set to perform calibration and learn thresholds for ATC and COTT. We trained 3
 578 models for each dataset with random seeds $\{0, 1, 10\}$.

579 **CIFAR10 and CIFAR100:** We reserved 10000 images from the training set as the validation set.
 580 We trained ResNet18 from scratch, using SGD with momentum equal to 0.9 for 300 epochs. We set
 581 weight decay to 5×10^{-4} and batch size to 200. We set the initial learning rate to 0.1 and multiply it
 582 by 0.1 every 100 epochs.

583 **ImageNet:** We reserved 50000 images from the training set as the validation set. We used ResNet50.
 584 While ImageNet pretrained weights are available in PyTorch, we needed multiple ones trained using
 585 different initializations. Due to limited computation resources, we reused the upper layer weights but
 586 reinitialized the last layer with different random seeds. We finetuned the whole model using Adam
 587 [22] with a batch size of 64 and a learning rate of 10^{-4} , for 10 epochs.

588 **BREEDS:** We used the intersection set of images that are both in the ImageNet validation images
 589 we set aside and the BREEDS dataset as the validation set. For all BREEDS datasets (Living17,
 590 Nonliving26, Entity13, Entity30), we trained ResNet50 from scratch.

591 For Living17 and Nonliving26, we used SGD with weight decay of 10^{-4} and batch size of 128. We
 592 trained for 450 epochs. We set the initial learning rate to 0.1 and multiplied it by 0.1 every 150
 593 epochs.

594 For Entity13 and Entity30, we used SGD with weight decay of 10^{-4} and batch size of 128. We
 595 trained for 300 epochs. We set the initial learning rate to 0.1 and multiplied it by 0.1 every 100
 596 epochs.

597 **Camelyon17-WILDS:** We used the `id_val` group as the validation set. We fine-tuned ImageNet
 598 pretrained ResNet50 using SGD with momentum of 0.9, weight decay of 5×10^{-4} , and batch size of
 599 32, for 5 epochs.

600 **RxRx1-WILDS:** We used the `id_text` group as the validation set. We followed [24] to fine-tune
 601 an ImageNet pretrained ResNet50. We used Adam with weight decay of 10^{-5} and batch size of 75,
 602 for 90 epochs. We increased the learning rate from 0 to 10^{-4} linearly for the first 10 epochs and
 603 decayed it following a cosine learning rate schedule.

604 **Amazon-WILDS:** We used the `id_val` group as the validation set. We followed [24] to fine-tune
605 a DistilBERT-base-uncased model [33]. We used AdamW [27] with weight decay of 10^{-2} , learning
606 rate of 10^{-5} , and batch size of 8, for 3 epochs. We set the maximum number of tokens to 512.

607 **CivilComments-WILDS:** We used the `val` group as the validation set. We followed [24] to
608 fine-tune a DistilBERT-base-uncased model [33]. We used AdamW [27] with weight decay of 10^{-2} ,
609 learning rate of 10^{-5} , and batch size of 16, for 5 epochs. We set the maximum number of tokens to
610 300.