# FEDERATED DISTILLATION OF NATURAL LANGUAGE UNDERSTANDING WITH CONFIDENT SINKHORNS

**Rishabh Bhardwaj, Tushar Vaidya & Soujanya Poria**
Information Systems Technology and Design
Singapore University of Technology and Design
Singapore
{rishabh_bhardwaj@mymail, tushar_vaidya@, sporia@}sutd.edu.sg

## ABSTRACT

Enhancing the user experience is an essential task for application service providers. For instance, two users living wide apart may have different tastes of food. A food recommender mobile application installed on an edge device might want to learn from user feedback (reviews) to satisfy the client's needs pertaining to distinct domains. Retrieving user data comes at the cost of privacy while asking for model parameters trained on a user device becomes space inefficient at a large scale. In this work, we propose an approach to learn a central (global) model from the federation of (local) models which are trained on user-devices, without disclosing the local data or model parameters to the server. We propose a federation mechanism for the problems with natural similarity metric between the labels which commonly appear in natural language understanding (NLU) tasks. To learn the global model, the objective is to minimize the optimal transport cost of the global model's predictions from the confident sum of soft-targets assigned by local models. The confidence (a model weighting scheme) score of a model is defined as the L2 distance of a model's prediction from its probability bias. The method improves the global model's performance over the baseline designed on three NLU tasks with intrinsic label space semantics, i.e., fine-grained sentiment analysis, emotion recognition in conversation, and natural language inference. We make our codes public at https://github.com/declare-lab/sinkhorn-loss.

## 1 INTRODUCTION

Due to recent technological advancements, more than two-thirds of the world's population has access to mobile phones[1]. A client application on these user devices has access to the unprecedented amount of data obtained from user-device interactions, sensors, etc. Learning algorithms can use this data to provide enhances user-experience. However, directly accessing this data comes at the cost of risking user privacy (Jeong et al., 2018).

To mitigate the issue, federated learning (FL) (shown in fig. 1) is a mechanism that retrieves the parameters of the (local) user-specific model and performs federation of knowledge either by distillation or merging the models (Konečný et al., 2016; McMahan et al., 2017) The algorithms aim to learn a domain-generalized central (global) model. The classic FL algorithms such as federated averaging and its adaptations are based on averaging of local model parameters, and thus only applied when the client models posses the similar network architectures (Mohri et al., 2019; Li et al., 2019a) However, FL paradigm has critical limitations of being costly in terms of communication load with the increase in local model sizes and demands all the participating models to have same architecture (homogenity) (Jeong et al., 2018; Lin et al., 2020).

Federated distillation (FD) proposes to exchange only the outputs of the local model, i.e, logits or probability measures whose dimensions are usually much smaller than the size of models themselves (Jeong et al., 2018). Therefore, FD allows to learn from an ensemble of different local models

---

[1] https://datareportal.com/global-digital-overview

of dissimilar configurations at reduced user-privacy risk, low communication overhead, and less memory space utilisation. However, entropy-based losses do not allow one to define metric structure in the label space. This can be critical when the task possesses a natural similarity relationship between the output labels.

Generally, FD algorithms adopts local model's ensemble knowledge distillation to global model using Kullback–Leibler (KL) divergence (Gou et al., 2021; Lin et al., 2020) or cross-entropy based losses (Jeong et al., 2018) as they are easy to compute and facilitate smooth backpropagation (Murphy, 2012). However, one critical limitation of such entropy-based losses is that they do not allow one to define metric structure in the label space. During global model ensemble training, leveraging such semantic structure can be useful when the task possesses a natural similarity relationship between the output labels. For instance, in the task of fine-grained sentiment classification of a text, strongly positive sentiment is closer to positive while far from strongly negative. Contrary to information-theoretic losses, Optimal Transport (OT) admits prior relationship in the label space (Frogner et al., 2015). The major contributions of this work are:

**Contribution:1** For the tasks with intrinsic label-space semantics, we propose a better federated (ensemble) distillation of local models to learn global model by encoding inter-label relationships using optimal transport.

The user-specific local data is potentially non-independent and non-identically distributed (non-IID) (Tong et al., 2020). Thus, the local models are prone to acquire biases such as *population bias*: the local user may not represent the target overall population; and *sampling bias*: knowledge transfer from a local model may not be useful to general over larger group (Mehrabi et al., 2019).

**Contribution:2** To minimize the effect of intrinsic probability bias arising from user-centric (non-IID) training of the local models, we introduce an L2 distance-based weighted distillation.

As shown in fig. 1, the global model $M_g$ aims to minimize the weighted sum of distance $D_1, D_2, D_3$ with weights $w_1, w_2, w_3$, respectively. Contribution:1 provides a better distance calculation while Contribution:2 presents a better weighting mechanism.

To further support our contributions, we derive the *Lipschitz* constant and Rademacher complexity-based generalisation bounds of the unregularized 1-Wasserstein based confident ensemble distillation. In the end, we empirically demonstrate the strong performance, generally improving upon the baselines, on the three natural language understanding (NLU) tasks, i.e., fine-grained sentiment analysis, emotion recognition in conversation, and natural language inference.
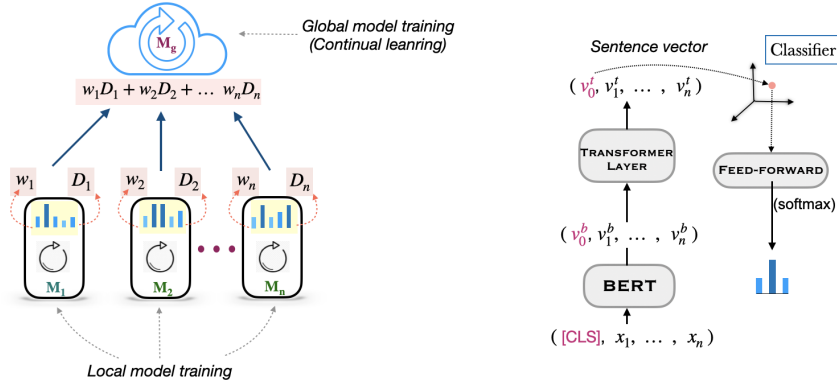


Figure 1: Local models acting as multiple teachers trained locally on user data while the global model acts as a student. The global model can only observe predictions of local models.

## 2 RELATED WORK

There have been many approaches to FL such as local model parameter averaging (McMahan et al., 2017; Lin et al., 2018) based on local SGD updates. It requires global and local models should have the same model architecture. Another line of work is multiple-source adaptation formulations

where a learner has access to source domain-specific predictor and no access to the labeling data from these domains. The expected loss is the mixture of source domains (Hoffman et al., 2018). Even though the formulation is close, our solutions are different as we do not have access to the local or global data domain distribution. In natural language processing, (Hilmkil et al., 2021; Lin et al., 2020; Liu & Miller, 2020) fine-tune Transformer-based architecture in the federated setting on small scale datasets. However, they do not leverage label space semantics and the analysis is restricted to small-scale datasets.

Closest to our work aims to improve local client training based on local data heterogeneity (Li et al., 2018; Nedic, 2020; Wang et al., 2019; Khaled et al., 2019; Li et al., 2019b). Knowledge distillation aims to transfer knowledge from a (large) teacher model to a (smaller) student model (Hinton et al., 2015; Buciluǎ et al., 2006). Given the output logit/softmaxed valued of the teacher model, the student can imitate the teacher's behavior (Romero et al., 2014; Tian et al., 2019; Tung & Mori, 2019; Koratana et al., 2019; Ahn et al., 2019). A few works are dedicated to the distillation of the ensemble of teacher models to the student model. This includes logit averaging of teacher models (You et al., 2017; Furlanello et al., 2018; Anil et al., 2018) or feature level knowledge extraction (Park & Kwak, 2019; Liu et al., 2019; Wu et al., 2019).

For Contribution:1, we use standard and widely used Entropy-based loss (KL-divergence) as our baseline. We construct two baselines for confidence score calculation (Contribution:2) from the prior works, i.e., logit averaging and weighting scheme based on local model dataset size (McMahan et al., 2017). In this work, we compare the proposed approach with baselines on the three NLU tasks.

## 3 BACKGROUND

### 3.1 OPTIMAL TRANSPORT

Traditional divergences, such as KL, ignore metric structure in the label space $\mathcal{Y}$. Optimal Transport (OT) metric can be extremely useful in defining inter-label semantic relationships. OT offers an additional advantage when measures have non-overlapping support (Peyré et al., 2019). Specific to our (classification) problems, we will focus on discrete measures. Assume $\mathcal{Y}$ possess a metric $d_{\mathcal{Y}}(\cdot, \cdot)$ referred to as ground truth metric. It establishes the semantic similarity between labels. The original OT problem is a linear program attributed to Kantorovich. Let $\mu_i$ and $\nu_j$ be the probability masses respectively applied to $i \in \mathcal{Y}_s$ and $j \in \mathcal{Y}_t$. Let $\pi_{i,j}$ be the transport assignment from $i$ to $j$ that costs $C_{(i,j)}$, i.e., an element of the cost matrix C. We denote Frobenius inner product by $\langle \cdot, \cdot \rangle$. The primal goal is to find the plan $\pi \in \Pi(\mu, \nu)$ that minimizes the transport cost

$$
\begin{aligned}
\mathrm{T}(\mu, \nu) &\stackrel{\text{def.}}{=} \min_{\pi \in \Pi(\mu, \nu)} \langle \pi, \mathrm{C} \rangle \\
\Pi(\mu, \nu) &= \{\pi \in (\mathbb{R}_+)^{n_s \times n_t} \mid \sum_{j \in \mathcal{Y}_t} \pi_{(i,j)} = \mu_i \text{ and } \sum_{i \in \mathcal{Y}_s} \pi_{(i,j)} = \nu_j\}
\end{aligned}
\tag{1}
$$

where $n_s = |\mathcal{Y}_s|$ and $n_t = |\mathcal{Y}_t|$. In our problem, $\mathcal{Y}_s = \mathcal{Y}_t = \mathcal{Y}$. When the cost is given in terms of metric $d_{\mathcal{Y}}(\cdot, \cdot)$ also known as Wasserstein distance (Bogachev & Kolesnikov, 2012).

### 3.2 SINKHORN LOSS

Although research in Wasserstein has been active, there have been challenges in its computation and implementation. The OT problem (equation 1) can be solved with combinatorial algorithms such as simplex-based methods (Courty et al., 2016). The computational complexity is, however, shown to be $O((n_s + n_t)n_s n_t \log(n_s + n_t))$ at best (Ahuja et al., 1988), and thus, the utility lessens as the size of dataset increases. The interest in the machine learning community took off after Cuturi's seminal paper (Cuturi, 2013). Rather than work with pure OT (Wasserstein) distances, we will restrict our attention to plain regularized OT, i.e, vanilla Sinkhorn distances. The distance can be efficiently solved by iterative Bregman Projections (Benamou et al., 2015).

**Definition 3.1.** *Vanilla Sinkhorn Distance*

$$\mathrm{T}_\varepsilon(\mu, \nu) \overset{\text{def.}}{=} \min_{\pi \in \Pi(\mu,\nu)} \langle \pi, \mathrm{C} \rangle + \varepsilon D_{\mathrm{KL}}(\pi, \mu \otimes \nu)$$

$$\Pi(\mu, \nu) = \{\pi \in (\mathbb{R}_+)^{n_s \times n_t} \mid \sum_{j \in \mathcal{Y}_t} \pi_{(i,j)} = \mu_i \ and \ \sum_{i \in \mathcal{Y}_s} \pi_{(i,j)} = \nu_j\} \tag{2}$$

where, $D_{\mathrm{KL}}(\pi, \mu \otimes \nu) = \sum_{i,j} \left[ \pi_{i,j} \log \frac{\pi_{i,j}}{\mu_i \nu_j} - \pi_{i,j} + \mu_i \nu_j \right]$. Entropic regularisation convexifies the loss function and thus is a computational advantage in computing gradients (Luise et al., 2018; Peyré et al., 2019; Feydy et al., 2019). As $\varepsilon \to 0^+$, we retrieve the unregularized Wasserstein distance.

## 4 METHODOLOGY

### 4.1 PROBLEM FRAMEWORK

The main participants in this framework are: 1) global model $\mathcal{M}_g$, and 2) a set of $K$ local models $\{\mathcal{M}_k\}_{k \in \mathcal{K}}, \mathcal{K} = \{1, \dots, K\}$. We denote the set $\{\mathcal{M}_k\}_{k \in \mathcal{K}}$ by $\mathcal{M}_{\mathcal{K}}$.

- The global model $\mathcal{M}_g$ aims to learn a user-generalized hypothesis $h_\theta \in \mathcal{H}_g$ that exists on the central application server.
- A local models learns a client-specific hypothesis $h_k \in \mathcal{H}_k$ on the $K_{\text{th}}$ user-generated data.

The central server can retrieve back a local model's prediction for a particular input. Thus, global model training can benefit from the hypotheses of local models, denoted by $h_{\mathcal{K}}$, but not the parameters set $\mathcal{M}_{\mathcal{K}}$ or the private user-generated data. The global model on the server $\mathcal{M}_g$ is generally preoccupied with the knowledge applicable across the domains. To enhance the user service, the server distills the knowledge from $h_{\mathcal{K}}$ and merges it into $\mathcal{M}_g$'s hypothesis $h_\theta$. The knowledge transfer happens with the assistance of a transfer set. *Transfer set* is the set of unlabeled i.i.d. samples used to learn global model parameters $\theta$. It creates a crucial medium to transfer the knowledge from the local models to the global model. To facilitate the knowledge transfer, we consider the (noisy) soft-labels are obtained from $h_{\mathcal{K}}$.

Since only the $h_{\mathcal{K}}$ is shared with $\mathcal{M}_g$ to find a better $h_\theta$ based on user feedback, the local models and global models can have heterogeneous architectures. This is useful when client devices at certain locations do not have enough resources to run and fit over large models.

### 4.2 ENSEMBLE DISTILLATION LOSS

For demonstration, we consider a user application that performs sentiment classification task on user-generated text $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)}) \in \mathcal{X}$ into its sentiment $y^{(i)} \in \mathcal{Y}$, where $\mathcal{X}$ denotes the input space of all possible text strings and label space $\mathcal{Y} = \{1 \text{ (strong negative)}, 2 \text{ (weak negative)}, 3 \text{ (neutral)}, 4 \text{ (weak positive)}, 5 \text{ (strong positive)}\}$. In this work, all the hypotheses are of the form $h : \mathcal{X} \mapsto \Delta^{\mathcal{Y}}$, $\Delta^{\mathcal{Y}}$ denotes a probability distribution on the set of labels $\mathcal{Y}$. We propose a learning algorithm that runs on the central server to fit $\mathcal{M}_g$'s parameters $\theta$ by receiving predictions such as (softmaxed) logits from $h_{\mathcal{K}}$. Without the loss of generality, the goal is to search for a hypothesis $h_{\hat{\theta}}$ that minimizes the empirical risk

$$h_{\hat{\theta}} = \arg\min_{h_\theta \in \mathcal{H}} \left\{ \hat{\mathbb{E}}_S \left[ \mathcal{L}_\varepsilon(h_\theta(x), h_{\mathcal{K}}(x)) \right] = \frac{1}{N} \sum_{i=1}^{N} \mathcal{L}_\varepsilon( h_\theta(x^{(i)}), h_{\mathcal{K}}(x^{(i)}) ) \right\}. \tag{3}$$

the loss is defined as

$$\mathcal{L}_\varepsilon( h_\theta(x^{(i)}), h_{\mathcal{K}}(x^{(i)}) ) = \sum_{k \in \mathcal{K}} W_{B_k}( h_k(x^{(i)}) ) \ \mathrm{T}_\varepsilon( h_\theta(x^{(i)}), h_k(x^{(i)}) ) \tag{4}$$

where $\mathrm{T}_\varepsilon(\cdot, \cdot)$ is the discrepency between the two probability measures as its arguments; $W_{B_k}(.)$ is the sample-specific weight assigned to the $k_{\text{th}}$ local model's prediction. Next, we elaborate on the the functions $\mathrm{T}_\varepsilon(\cdot, \cdot)$ and $W_{B_k}(.)$ which are crucial for our proposed FD algorithm.

## 5 SINKHORN-BASED DISTILLATION

For entropy-based loss, we adopt KL divergence. As discussed in section 3.1, we employ Sinkhorn distance to implement OT-based loss.

### 5.1 UNWEIGHTED DISTILLATION

For a text input $x^{(i)}$ from the transfer set, $\mathrm{T}_\varepsilon(\ h_\theta(x^{(i)}),\ h_k(x^{(i)})\ )$ measures the Sinkhorn distance between the probability output of global model $h_\theta(x^{(i)})$ and a local model $h_k(x^{(i)})$. In eq. (4), the sample-wise distance is computed between $h_\theta(x^{(i)})$ and a probability distribution in the set $h_{\mathcal{K}}$. A simple approach to fit the global hypothesis $h_\theta$ is to uniformly distill the knowledge from user-specific hypotheses, i.e., $W_{B_k}(h_k(x^{(i)})) = 1 \ \forall\ k \in \mathcal{K},\ i \in [N]$.

### 5.2 WEIGHTED DISTILLATION

The user-generated local datasets are potentially non-IID with respect to the global distribution and possess a high degree of class imbalance (Weiss & Provost, 2001). As each local model $\mathcal{M}_k$ is trained on samples from potentially non-IID and imbalance domains, they are prone to show skewed predictions. The unweighted distillation tends to transfer such biases. One might wonder *for a given transfer set sample, which local model's prediction is reliable?*. Although an open problem, we try to answer this question by proposing a local model (teacher) weighting scheme. It calculates the confidence score of a model's prediction and performs weighted distillation—weights being in positive correlation with the local model's confidence score. Next, we define the confidence score.

**Confidence score (L2)** For a given sample $x$ from the transfer set, the skew in a local model's prediction $h(x)$ can help us determine the confidence score ($W(\cdot)$ in eq. (4)) with which it can transfer its knowledge to the global model. However, a model can show skew due to training on an imbalance dataset or chosen capacity of the hypothesis space which can potentially cause a model to overfit/underfit (Caruana et al., 2001). For instance, a model has learned to misclassify negative sentiment as strongly negative samples owing to a high confusion rate. Such models are prone to show inference time classification errors with highly skewed probabilities. Thus, confidence scoring based on the probability skew may not be admissible. Hence, we incorporate L2 confidence. For a given sample, we define the model's L2 confidence score $W_B(h(x))$ as Euclidean distance of its output probability distribution from the probability bias $B$. We define probability bias $B$ of a local model as the expected value of prediction when a model $h$ receives noise (random text) at the input. Let $h(x) \in \mathcal{Y}$ denotes the predicted distribution of a model for an input text $x$:

$$b_{l\in\mathcal{Y}} := \mathbb{E}_{x\sim\mathcal{N}}[h(x) = l] \qquad\qquad (\mathcal{N}\text{: the distribution of noise})$$
$$B := (b_1, \ldots, b_{|\mathcal{Y}|}) \qquad\qquad\qquad (\text{model probability skew}).$$

As shown in fig. 2 for a three-class classification, the equidistant distributions lie on an arc with center at $B$. Points with high confidence score lie on distant arcs. As radius of the arc increases, majority of its portion lies towards the high value of $p_l$, i.e., the $l$ with which the model is biased against since $b_l = \min\ \{b_1, \ldots, b_{|\mathcal{Y}|}\}$ ($p_3$ in the figure). Moreover, the maximum confidence score is achieved at the vertex $p_l = 1$.

**Proposition 5.1.** *From a given point $B$ in a k-simplex, point with the highest confidence lies on one of its vertices.*
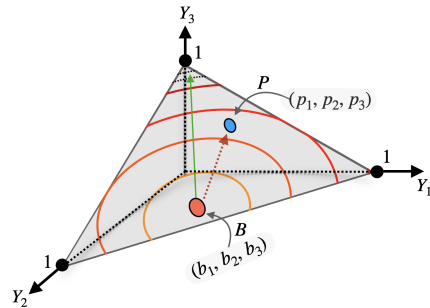


Figure 2: An illustration—a three-class classifier with bias $B$ and outputs a distribution $P = (p_1, p_2, p_3)$ for certain input. The green arrow denotes direction of increased confidence score with equiconfident arcs.

*Proof.* First, we analyse the case of a 2-simplex defined in a three-dimensional Euclidean space.

Table 1: Data statistics and performance of local models. For SA and ERC tasks, score denotes Macro F-1 performance metric, while it denotes Accuracy metric for NLI.

| | SA | | | | ERC | | | NLI | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Cell | Cloths | Toys | Food | IEMOCAP | MELD | DyDa$_{\{1,2,3\}}$ | Fic | Gov | Slate | Tele | Trv | SNLI |
| train | 133,574 | 19,470 | 116,666 | 397,917 | 3,354 | 9,450 | 21,680 | 77,348 | 77,350 | 77,306 | 83,348 | 77,350 | 549,367 |
| valid | 19,463 | 28,376 | 17,000 | 57,982 | 342 | 1,047 | 2,013 | 5,902 | 5,888 | 5,893 | 5,899 | 5,904 | 9,842 |
| test | 37,784 | 55,085 | 33,000 | 112,555 | 901 | 2,492 | 1,919 | 5,903 | 5,889 | 5,894 | 5,899 | 5,904 | 9,824 |
| score | 0.49 | 0.52 | 0.49 | 0.64 | 0.55 | 0.45 | 0.38, 0.34, 0.40 | 0.63 | 0.65 | 0.62 | 0.65 | 0.63 | 0.85 |

Let $f_P = \sum_{i=1}^{3}(p_i - b_i)^2$, the quadratic pro-
gram can be formulated as $\max\{f_P : \sum_{i=1}^{3} p_i = 1, \ p_i \geq 0\}$. The convex hull of vertices lying on the axes forms a closed and bounded feasible region. Thus, from extreme value theorem, there exist absolute maximum and minimum. $f$ attains its minimum at $p = b$, which is also the critical point of $f_P$. Now, we need to find its value on the boundary points contained in the set of 1-simplices (line segments) $\{p_i + p_j = 1, p_k = 0 : (i,j,k) \in 1,2,3, i \neq j \neq k\}$. For the 1-simplex $p_1 + p_2 = 1, p_3 = 0$, the values of $f_P$ at its end points that are $(1 - b_1)^2 + b_2^2 + b_3^2$ and $b_1^2 + (1 - b_2)^2 + b_3^2$, one of which is maxima of $f$ attained over the 1-simplex [2]. Similarly for the other line segments, the complete set of boundary values of $f_P$ is $k - 2b_1$, $k - 2b_2$, and $k - 2b_3$ where $k = b_1^2 + b_2^2 + b_3^2 + 1$, occurring at $P = (1,0,0), (0,1,0)$ and $(0,0,1)$, respectively. Thus, the maximum of $f_P$ will lie on $i_\text{th}$-axis such that $b_i = \min(b_1, b_2, b_3)$. This proof can be generalized for a probability simplex in higher dimensions. As shown above, each iteration of a lower dimensional simplex will return vertices as the point of maxima in the end. $\square$

L2 confidence is a proper distance metric and computationally stable. The distance metric is invariant to translation in the Euclidean space. Thus, $W_{B_k}(h_k)$ measures the L2 distance of $k_\text{th}$ model's prediction from its intrinsic probability bias $B_k$. Moreover, L2 confidence can be used to compute weights for both distillation methods, i.e., Sinkhorn and Entropy.
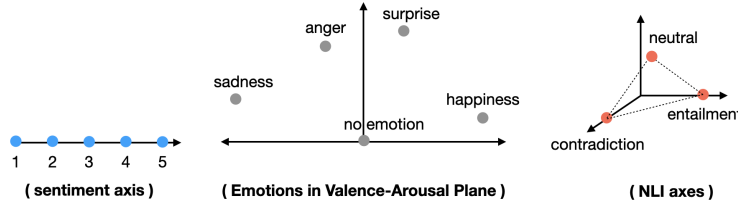


Figure 3: Semantic coordinates of SA, ERC, and NLI.

## 5.3 STATISTICAL PROPERTIES OF WASSERSTEIN

Let the samples $S = \{(x^{(1)}, y^{(1)}), \ldots, (x^{(N)}, y^{(N)})\})$ be i.i.d from the domain distribution of the transfer set and $h_{\hat{\theta}}$ be the empirical risk minimizer. Assume the global hypothesis space $\mathcal{H}_g = \mathfrak{s} \circ \mathcal{H}_g^o$, i.e., composition of softmax and a hypothesis $\mathcal{H}_g^o : \mathcal{X} \mapsto \mathbb{R}^{|\mathcal{Y}|}$, that maps input text to a scalar (logit) value for each label. Assuming vanilla Sinkhorn with $\varepsilon \to 0^+$, we establish the property for 1-Wasserstein.

**Theorem 5.2.** *If the global loss function (as in eq.* (4)*) uses unregularized 1-Wasserstein metric between predicted and target measure, then for any $\delta > 0$, with probability at least 1-$\delta$*

$$\mathbb{E}\big[\mathcal{L}(h_{\hat{\theta}}(x), h_{\mathcal{K}}(x))\big] \leq \inf_{h_\theta \in \mathcal{H}_g^o} \hat{\mathbb{E}}\big[\mathcal{L}(h_\theta(x), h_{\mathcal{K}}(x))\big] + 32 \times |\mathcal{Y}| \times \mathfrak{R}_N(\mathcal{H}_g^o) + 2C_M|\mathcal{Y}|\sqrt{|\mathcal{Y}|\frac{log1/\delta}{2N}} \quad (5)$$

where $\mathfrak{R}_N(\mathcal{H}_g^o)$, decays with $N$, denotes Rademacher complexity Bartlett & Mendelson (2002) of the hypothesis space $\mathcal{H}_g^o$. $C_M$ is the maximum cost of transportation within the label space. In the case of SA, $C_M = 4, |\mathcal{Y}| = 5$. The expected loss of the empirical risk minimizer $h_{\hat{\theta}}$ approaches the

---
[2]Ignoring the critical point which gives the minima and perpendicular drawn from $b$ to the line segment.

best achievable loss for $\mathcal{H}_g$. The proof of theorem theorem 5.2 and method to compute gradient are relegated to the Appendix.

To comprehensively analyse the importance of OT over entropy for the task with intrinsic label similarity, we introduce a new performance metric that evaluates a model based on the semantic correctness of its predictions.

## 5.4 SEMANTIC DISTANCE

During the evaluation, most common metrics (such as accuracy and F1) observe the label with the highest logit (or probability) against ground truth, hence, ignore the overall distribution. However, for tasks with label space semantics, it can be of great importance. Thus, we define a new performance metric—Semantic Distance (SD)—that measures the semantic closeness of the output distribution against the ground truth. Given a label coordinate space, SD is defined as the mean Euclidean distance of expected output from the ground truth label. For instance, given the sentiment classes {1,2,3,4,5}, the probability scores of two models m1 and m2 assigns to a strongly negative text input be {0.2, 0.7, 0.033, 0.033, 0.033} and {0.4, 0.1, 0.1, 0.1, 0.3}, respectively. The argmax output of m2 is correct. However, even when the argmax output of m1 is incorrect, the expected output of m1, i.e., 1.97 is (Euclidean) closer to the ground truth label 1 than m2, i.e., 2.80, and thus more semantically accurate. *A low score denotes more semantically accurate prediction*. The lowest possible value of SD is 0 while the highest possible value depends on the number of labels and their map in the semantic space.

For datasets with class imbalance, such as SA task in this work Table 1, we first calculate label-wise SD values and compute their mean.

*Analysis*—To demonstrate the usefulness of the SD metric, for the SA, we draw box plots of pretrained global models via Entropy (KL-divergence) and Sinkhorn-based losses. The pretraining methodology is described in section 6. As shown in fig. 4, we observe the median SD of Sinkhorn (green box, red line) is closer to the ground truth sentiment classes—1,2,3, and 5 as compared to median SD of Entropy (blue box-red line). Similarly, the means (black diamond), the first quartile (25% of the samples), and the third quartile (75% of the samples) for the Sinkhorn-based model are relatively closer to the ground truth.
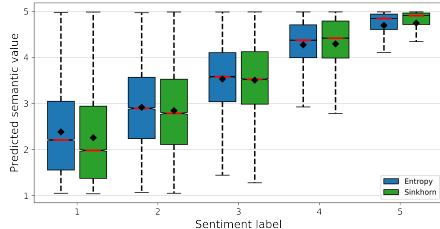


Figure 4: Box plot showing expectation of output probabilities for SA task. Horizontal and vertical axes denote ground truth sentiment labels.

## 6 EXPERIMENTS

**Baselines** —We setup the following baselines for a thorough comparison between Sinkhon and Entropy-based losses. Let [*Method*] be the placeholder for Sinkhorn and Entropy. [*Method*]-A denotes unweighted distillation of local models (section 5.1), i.e., $W_{B_k}(h_k(x)) = 1$ (in eq. (4)). In [*Method*]-D, $W_{B_k}(h_k(x))$ is proportional to size of local datasets. [*Method*]-U defines sample-specific confidence (weights) as the distance of output from the uniform distribution over labels. For each sample, [*Method*]-E computes weight of $k_{\text{th}}$ local model as distance of its prediction from probability bias $B_k$, i.e., L2 confidence.

**Tasks** We set up the three natural language understanding tasks with intrinsic similarity in the label space: 1) fine-grained Sentiment Analysis (SA); 2) Emotion Recognition in Conversation (ERC); and 3) Natural Language Inference (NLI). NLI is the task of determining the inference relation between two texts. The relation can be entailment, contradiction, or neutral (MacCartney & Manning, 2008). For a given transcript of conversation, the ERC task aims to identify the emotion of each utterance from the set of pre-defined emotions (Poria et al., 2019). For our experiments, we choose the five most common emotions that are sadness, anger, surprise, and happiness, and no emotion.

**Datasets**   For the SA task, we use four large-scale datasets: 1) Toys: toys and games; 2) Cloth: clothing shoes and jewelry; 3) Cell: cell phones and accessories; 4) Food: Amazon's fine food reviews, specifically curated for the five-class sentiment classification. For the transfer set, we use grocery and gourmet food (104,817 samples) and discard the labels (He & McAuley, 2016). Each dataset consists of reviews rated on a scale of 1 (strong negative) to 5 (strong positive). Similarly, for ERC, we collect three widely used datasets: DyDa: DailyDialog (Li et al., 2017), IEMOCAP: interactive emotional dyadic motion capture database (Busso et al., 2008), and MELD: Multimodal EmotionLines Dataset (Poria et al., 2018). To demonstrate our methodology, we partition the DyDa dataset into four equal chunks. $DyDa_1$, $DyDa_2$, are used as local, $DyDa_3$ is used as global dataset. Dropping the labels from $DyDa_4$, we use it as transfer set. For NLI task, we use SNLI (Bowman et al., 2015) as global dataset and MNLI (Williams et al., 2017) as local dataset. We split the latter across its five genres, which are, fiction (Fic), government (Gov), telephone (Tele), travel (Trv), and Slate. This split assists in simulating distinct user (non-IID samples) setup. We use ANLI dataset (Nie et al., 2020) as transfer set.

**Architecture**   We set up a compact transformer-based model used by both global and local models (fig. 1), although, the federation does not restrict both the local and model configurations to be the same. The input is fed to the pretrained BERT-based classifier (Turc et al., 2019; Devlin et al., 2018). Thus, we obtain probabilities with support in the space of output labels, i.e., $\mathcal{Y}$. We keep all the parameters trainable, hence, BERT will learn its embeddings specific to the classification task. For NLI task, we append premise and hypothesis at input separated by special token [SEP] token, followed by a standard classification setup.

**Training local models**   To compare Sinkhorn-based distillation with baselines, first, we pretrain local models. Since cross-entropy (CE) loss is less computationally expensive as compared to OT, we use CE for local model training. For all model training, we tuned hyperparameters for all the models separately and chose the model that performs best on the validation dataset. The data statistics and performances of local models on individual tasks are shown in Table 1.

**Training a global model**   We make use of transfer set (unlabeled) samples and obtain noisy labels. For a text sample in transfer set, eq. (4) aims to fit a global model to the weighted sum of noisy predictions of the local models. To retain the previous knowledge of a global model, we adapt learning without forgetting paradigm. To incorporate this, we store predictions of the pretrained global model on the transfer set and perform its distillation along with local models.

**Label-space**   We define label semantic spaces for the three tasks. As shown in fig. 3, we assign sentiment labels a one-dimensional space. For the ERC task, we map each label to a two-dimensional valence-arousal space. Valence represents a person's positive or negative feelings, whereas arousal denotes the energy of an individual's affective state. As mentioned in (Ahn et al., 2010), anger (-0.4, 0.8), happiness (0.9, 0.2), no emotion (0, 0), sadness (-0.9, -0.4), and surprise (0.4, 0.9). The cost (loss) incurred to transport a mass from a point $p$ to point $q$ is $C_{p,q} := |p - q|$. For NLI task, we define a three dimensional coordinate with entailment (1, 0, 0), contradiction (0,0,1) and neutral (0.5, 1, 0.5). The cost $C_{p,q} := ||p - q||_2$, where, cost of transport from entailment to contradiction is higher than it is to neutral. It is noteworthy that for this task, we perform a manual search to identify label coordinate space and transportation cost.

Table 2: Fine-grained SA task: Macro F1 and Semantic Distance.

| Algorithm | F1 Score | | | | | Semantic Distance | | | | |
| | ——-Local——- | | | Global | ALL | ——-Local——- | | | Global | ALL |
| | Cloths | Toys | Cell | Food | | Cloths | Toys | Cell | Food | |
| Entropy-A | 0.48 | 0.44 | 0.47 | 0.52 | 0.50 | 0.77 | 0.87 | 0.76 | 0.79 | 0.79 |
| Entropy-D | 0.48 | 0.44 | 0.46 | 0.56 | 0.52 | 0.77 | 0.86 | 0.77 | 0.71 | 0.78 |
| Entropy-U | 0.47 | 0.43 | 0.47 | 0.50 | 0.49 | 0.79 | 0.90 | 0.78 | 0.82 | 0.83 |
| Entropy-E | 0.49 | 0.46 | 0.48 | 0.55 | 0.52 | 0.74 | 0.80 | 0.74 | 0.71 | 0.75 |
| Sinkhorn-A | 0.49 | 0.47 | 0.47 | 0.55 | 0.52 | 0.74 | 0.80 | 0.72 | 0.75 | 0.75 |
| Sinkhorn-D | 0.47 | 0.44 | 0.45 | **0.59** | 0.52 | 0.77 | 0.84 | 0.76 | **0.65** | 0.76 |
| Sinkhorn-U | 0.48 | 0.44 | 0.47 | 0.51 | 0.49 | 0.77 | 0.89 | 0.77 | 0.83 | 0.82 |
| Sinkhorn-E | **0.49** | **0.47** | **0.48** | 0.55 | **0.52** | **0.72** | **0.78** | **0.72** | 0.69 | **0.73** |

Table 2, Table 3, and Table 4 show performance, i.e., Macro-F1 (or Accuracy) score and Semantic Distance of global models predictions from ground truth. Evaluations are done for fine-tuned (after distillation) global model with respect to the test sets of both local and global datasets. The testing over local datasets will help us analyse how well the domain generic global model performs over the individual local datasets and the testing over the global dataset is to make sure there is no catastrophic forgetting of the previous knowledge.

Table 3: ERC task: Macro F1 and Semantic Distance.

| Algorithm | F1 Score | | | | | | Semantic Distance | | | | | |
| | -------Local------- | | | Global | ALL | | -------Local------- | | | Global | ALL |
| | MELD | IEMOCAP | $DyDa_0$ | $DyDa_1$ | $DyDa_2$ | | MELD | IEMOCAP | $DyDa_0$ | $DyDa_1$ | $DyDa_2$ | |
| Entropy-A | 0.28 | 0.21 | 0.34 | 0.33 | 0.36 | 0.31 | 0.67 | 0.65 | 0.60 | 0.61 | 0.63 | 0.64 |
| Entropy-D | 0.30 | 0.21 | 0.39 | 0.35 | 0.38 | 0.34 | 0.68 | 0.65 | 0.57 | 0.59 | 0.61 | 0.62 |
| Entropy-U | 0.30 | 0.24 | 0.42 | 0.31 | 0.37 | 0.33 | 0.68 | 0.65 | 0.60 | 0.62 | 0.64 | 0.64 |
| Entropy-E | 0.34 | **0.36** | 0.42 | 0.40 | **0.44** | 0.39 | 0.69 | 0.62 | 0.57 | 0.59 | 0.62 | 0.62 |
| Sinkhorn-A | 0.30 | 0.26 | 0.45 | 0.37 | 0.39 | 0.34 | 0.67 | 0.67 | 0.59 | 0.62 | 0.62 | 0.64 |
| Sinkhorn-D | 0.35 | 0.31 | 0.45 | 0.37 | 0.44 | 0.39 | 0.67 | 0.63 | 0.54 | 0.59 | 0.61 | 0.61 |
| Sinkhorn-U | 0.30 | 0.23 | 0.39 | 0.34 | 0.39 | 0.34 | 0.68 | 0.68 | 0.62 | 0.64 | 0.64 | 0.65 |
| Sinkhorn-E | **0.38** | 0.33 | **0.46** | **0.43** | 0.43 | **0.41** | **0.64** | **0.62** | **0.53** | **0.56** | **0.60** | **0.59** |

For the SA task, we observe the global models trained from Sinkhorn distillation of local models is amongst the best F1 scores on all the local domains. For the most simplistic baseline, i.e., un-weighted distillation, Sinkhorn-A consistently gives a higher F1 score as compared to Entropy-A. Even though the performance of Sinkhorn-A is close to Sinkhorn-E, the latter is more semantically accurate as depicted by corresponding SA scores. This shows the efficacy of weighting local models' prediction with L2 confidence. Sinkhorn-D based global model training gives the best F1 and SA scores on global datasets. We postulate this is due to a large number of global samples as compared to local dataset sizes that bias the distillation weights $W(h_k(x))$, hence force the model to perform better on the global dataset. It can be observed that this comes at the cost of poor performance on the local datasets.

Table 4: NLI task: Accuracy and Semantic Distance.

| Algorithm | Accuracy | | | | | | | Semantic Distance | | | | | | |
| | ------------Local------------ | | | | | Global | ALL | ------------Local------------ | | | | | Global | ALL |
| | Fic | Gov | Slate | Tele | Trv | SNLI | | Fic | Gov | Slate | Tele | Trv | SNLI | |
| Entropy-A | 0.60 | 0.62 | 0.60 | 0.60 | 0.62 | 0.78 | 0.65 | 0.56 | 0.54 | 0.56 | 0.56 | 0.55 | 0.40 | 0.53 |
| Entropy-D | 0.58 | 0.59 | 0.57 | 0.57 | 0.58 | **0.85** | 0.64 | 0.58 | 0.56 | 0.58 | 0.58 | 0.57 | 0.30 | 0.53 |
| Entropy-U | 0.60 | 0.62 | 0.60 | 0.60 | 0.61 | 0.76 | 0.65 | 0.55 | 0.54 | 0.56 | 0.56 | 0.55 | 0.40 | 0.53 |
| Entropy-E | 0.60 | 0.62 | 0.60 | 0.60 | 0.61 | 0.76 | 0.65 | 0.55 | 0.54 | 0.56 | 0.55 | 0.54 | 0.39 | 0.52 |
| Sinkhorn-A | 0.60 | **0.63** | 0.60 | **0.61** | **0.62** | 0.73 | 0.64 | 0.54 | 0.53 | 0.55 | 0.55 | 0.53 | 0.42 | 0.52 |
| Sinkhorn-D | 0.55 | 0.57 | 0.54 | 0.55 | 0.56 | 0.85 | 0.63 | 0.58 | 0.56 | 0.58 | 0.58 | 0.57 | **0.23** | 0.52 |
| Sinkhorn-U | 0.60 | 0.62 | 0.60 | 0.60 | 0.61 | 0.77 | 0.65 | 0.53 | 0.52 | 0.54 | 0.54 | 0.52 | 0.37 | 0.51 |
| Sinkhorn-E | **0.60** | 0.62 | **0.60** | 0.60 | 0.61 | 0.77 | **0.65** | **0.53** | **0.52** | **0.54** | **0.54** | **0.52** | 0.37 | **0.50** |

For the ERC and NLI tasks, although it is hard to find a model that shows consistently better F1 or accuracy scores (one cause could be the small small size of local and global datasets), we observe the SD score of Sinkhorn-E is, in general, better amongst the baselines. As observed for the SA task in Entropy-D and Sinkhorn-D settings, since the SNLI dataset is bigger, the distillation forces global models to perform better on the global dataset. It is seen to come at the cost of degraded performance on the other (local) datasets.

Comparing Table 2, Table 3, and Table 4 all together, for the three tasks with intrinsic similarity in the label space, we observe Sinkhorn-based loss are in general better than KL-divergence, i.e., an entropy-based loss. Moreover, we observe L2 distance gives better scores amongst the individual loss-specific groups. We also notice that in certain tasks, where it is difficult to identify the best model, one can refer to the semantic distance. Besides this, as compared to other baselines, empirical observations suggest that Sinkhorn-E (our contribution) works well for the large-scale SA datasets, hence potentially scalable.

## 7 CONCLUSION

In this work, we introduced an algorithm for efficient federated distillation of natural language understanding from client devices to the central (global) model. We defined a new Euclidean distance-based metric to compute a local model's intrinsic probability bias. We analysed theoretical generalization bounds of empirical risk of the proposed loss function. In the end, we demonstrated the efficacy of the novel approach on the three NLU tasks of fine-grained sentiment analysis, emotion recognition in conversation, and natural language inference.

## REFERENCES

Junghyun Ahn, Stephane Gobron, Quentin Silvestre, and Daniel Thalmann. Asymmetrical facial expressions based on an advanced interpretation of two-dimensional russell's emotional model. *Proceedings of ENGAGE*, 2010.

Sungsoo Ahn, Shell Xu Hu, Andreas Damianou, Neil D Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9163–9171, 2019.

Ravindra K Ahuja, Thomas L Magnanti, and James B Orlin. Network flows. 1988.

Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E Dahl, and Geoffrey E Hinton. Large scale distributed neural network training through online distillation. *arXiv preprint arXiv:1804.03235*, 2018.

Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

Jean-David Benamou, Guillaume Carlier, Marco Cuturi, Luca Nenna, and Gabriel Peyré. Iterative bregman projections for regularized transportation problems. *SIAM Journal on Scientific Computing*, 37(2):A1111–A1138, 2015.

Vladimir Igorevich Bogachev and Aleksandr Viktorovich Kolesnikov. The monge-kantorovich problem: achievements, connections, and perspectives. *Russian Mathematical Surveys*, 67(5):785–890, 2012.

Samuel R Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. *arXiv preprint arXiv:1508.05326*, 2015.

Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 535–541, 2006.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335–359, 2008.

Rich Caruana, Steve Lawrence, and Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. *Advances in neural information processing systems*, pp. 402–408, 2001.

Nicolas Courty, Rémi Flamary, Devis Tuia, and Alain Rakotomamonjy. Optimal transport for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 39(9):1853–1865, 2016.

Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26:2292–2300, 2013.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trouvé, and Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2681–2690. PMLR, 2019.

Charlie Frogner, Chiyuan Zhang, Hossein Mobahi, Mauricio Araya-Polo, and Tomaso Poggio. Learning with a wasserstein loss. In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, pp. 2053–2061, 2015.

Tommaso Furlanello, Zachary Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born again neural networks. In *International Conference on Machine Learning*, pp. 1607–1616. PMLR, 2018.

Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.

Ruining He and Julian McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*, pp. 507–517, 2016.

Agrin Hilmkil, Sebastian Callh, Matteo Barbieri, Leon René Sütfeld, Edvin Listo Zec, and Olof Mogren. Scaling federated learning for fine-tuning of large language models. *arXiv preprint arXiv:2102.00875*, 2021.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.

Judy Hoffman, Mehryar Mohri, and Ningshan Zhang. Algorithms and theory for multiple-source adaptation. *arXiv preprint arXiv:1805.08727*, 2018.

Eunjeong Jeong, Seungeun Oh, Hyesung Kim, Jihong Park, Mehdi Bennis, and Seong-Lyun Kim. Communication-efficient on-device machine learning: Federated distillation and augmentation under non-iid private data. *arXiv preprint arXiv:1811.11479*, 2018.

Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. First analysis of local gd on heterogeneous data. *arXiv preprint arXiv:1909.04715*, 2019.

Jakub Konečnỳ, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527*, 2016.

Animesh Koratana, Daniel Kang, Peter Bailis, and Matei Zaharia. Lit: Learned intermediate representation training for model compression. In *International Conference on Machine Learning*, pp. 3509–3518. PMLR, 2019.

Michel Ledoux and Michel Talagrand. *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media, 2013.

Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*, 2018.

Tian Li, Maziar Sanjabi, Ahmad Beirami, and Virginia Smith. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019a.

Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. *arXiv preprint arXiv:1907.02189*, 2019b.

Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. Dailydialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of The 8th International Joint Conference on Natural Language Processing (IJCNLP 2017)*, 2017.

Tao Lin, Sebastian U Stich, Kumar Kshitij Patel, and Martin Jaggi. Don't use large mini-batches, use local sgd. *arXiv preprint arXiv:1808.07217*, 2018.

Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *arXiv preprint arXiv:2006.07242*, 2020.

Dianbo Liu and Tim Miller. Federated pretraining and fine tuning of bert using clinical notes from multiple silos. *arXiv preprint arXiv:2002.08562*, 2020.

Iou-Jen Liu, Jian Peng, and Alexander G Schwing. Knowledge flow: Improve upon your teachers. *arXiv preprint arXiv:1904.05878*, 2019.

Giulia Luise, Alessandro Rudi, Massimiliano Pontil, and Carlo Ciliberto. Differential properties of sinkhorn approximation for learning with wasserstein distance. *Advances in Neural Information Processing Systems*, pp. 5864–5874, 2018.

Bill MacCartney and Christopher D. Manning. Modeling semantic containment and exclusion in natural language inference. In *Proceedings of the 22nd International Conference on Computational Linguistics (Coling 2008)*, pp. 521–528, Manchester, UK, August 2008. Coling 2008 Organizing Committee. URL https://aclanthology.org/C08-1066.

Colin McDiarmid. Concentration. In *Probabilistic methods for algorithmic discrete mathematics*, pp. 195–248. Springer, 1998.

Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial Intelligence and Statistics*, pp. 1273–1282. PMLR, 2017.

Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *arXiv preprint arXiv:1908.09635*, 2019.

Mehryar Mohri, Gary Sivek, and Ananda Theertha Suresh. Agnostic federated learning. In *International Conference on Machine Learning*, pp. 4615–4625. PMLR, 2019.

Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.

Angelia Nedic. Distributed gradient methods for convex machine learning problems in networks: Distributed optimization. *IEEE Signal Processing Magazine*, 37(3):92–101, 2020.

Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial nli: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020.

SeongUk Park and Nojun Kwak. Feed: Feature-level ensemble for knowledge distillation. *arXiv preprint arXiv:1909.10754*, 2019.

Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport: With applications to data science. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.

Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. Meld: A multimodal multi-party dataset for emotion recognition in conversations. *arXiv preprint arXiv:1810.02508*, 2018.

Soujanya Poria, Navonil Majumder, Rada Mihalcea, and Eduard Hovy. Emotion recognition in conversation: Research challenges, datasets, and recent advances. *IEEE Access*, 7:100943–100953, 2019.

Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.

Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. *arXiv preprint arXiv:1910.10699*, 2019.

Qianqian Tong, Guannan Liang, and Jinbo Bi. Effective federated adaptive gradient methods with non-iid decentralized data. *arXiv preprint arXiv:2009.06557*, 2020.

Frederick Tung and Greg Mori. Similarity-preserving knowledge distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1365–1374, 2019.

Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*, 2019.

Shiqiang Wang, Tiffany Tuor, Theodoros Salonidis, Kin K Leung, Christian Makaya, Ting He, and Kevin Chan. Adaptive federated learning in resource constrained edge computing systems. *IEEE Journal on Selected Areas in Communications*, 37(6):1205–1221, 2019.

Gary M Weiss and Foster Provost. The effect of class distribution on classifier learning: an empirical study. 2001.

Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for sentence understanding through inference. *arXiv preprint arXiv:1704.05426*, 2017.

Ancong Wu, Wei-Shi Zheng, Xiaowei Guo, and Jian-Huang Lai. Distilled person re-identification: Towards a more scalable system. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1187–1196, 2019.

Shan You, Chang Xu, Chao Xu, and Dacheng Tao. Learning from multiple teacher networks. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1285–1294, 2017.

## A  WHAT WE HAVE KEPT FOR THE APPENDIX?

There are a few experiments that we consider to be important and may help compare the Sinkhorn loss-based learning with Kullback–Leibler (KL). These results build a firm base to choose Sinkhorn-based losses on the task of federated distillation of sentiments. For the experiments, we work on the global model $\mathcal{M}_g$ that has acquired knowledge from local models in the learning without forgetting paradigm.

- In appendix B, we convey the intuition behind using a Sinkhorn distance over an entropy-based divergence. Moving further in appendix B.1, we show the importance of natural metric in the label space by replacing the one-dimensional support with one-hot. Furthermore, in appendix B.2, we show how the model's clusters of sentence embeddings change when we move from a Sinkhorn distance-based loss to the KL divergence-based loss.

- In appendix C, we discuss the potential risk of gender and racial bias transfer from the local models to the global models. Although we incorporate bias induced from non-IID data training of the local models, we do not tackle the transfer of other biases that can arise from the data as well as from the training process.

- In appendix D, we provide the algorithm to compute gradients and the computation complexity of Sinkhron loss.

- In appendix E, we provide the proof of Theorem 5.2 on the empirical risk bound with the OT metric as unregularized 1-Wasserstein distance.

- In appendix F, we discuss the broader social impact of our work. We discuss how the method can be adopted for cyberbullying detection and limitations coming from local models.

- In appendix G, we elaborate on the experimental settings and license of the datasets used in this paper.

## B  DECISION BOUNDARIES VIA SENTENCE REPRESENTATIONS

One of the main advantages of using Optimal Transport-based (OT) metrics between two probability distributions, such as Sinkhorn distance, is the ability to define the relationship in metric space. This is not feasible in entropy-based divergences. The relationship further appears in the loss function that accounts for the error computations of an intelligent system in the task of classification (or regression). With advancements in computations of Sinkhorn distances, as in Feydy et al. (2019), gradient computations through such loss functions have become more feasible as shown in Algorithm 1. The inter-label relationships are apparent in tasks such as fine-grained sentiment classification, fake news detection, hate speech. In this work, we consider the relationship between two labels $p$ and $q$ as a taxi-cab distance in the one-dimensional metric space of sentiment labels $\mathcal{Y}$. This relation nuance should appear in the Sinkhorn distance, we call it the cost of transportation from a point $p$ to another point $q$ in the set $\mathcal{Y}$.

When we set a learning algorithm to minimize the loss function, the goal is to find model parameters that provide the least empirical risk in the space of predefined hypotheses. From the distance-based cost (loss), the risk is expected to be minimum when the predicted labels are mapped "near" to the ground truth label. The term "near" refers to the lower optimal transport cost of the probability mass spread over a certain region to another region. In our problem, both the regions are the same, i.e., locations from 1 to 5 in the metric space. A ground truth probability mass (almost everything) at 1 would prefer an intelligent system to predict a probability mass near 1 so that it will require a lesser taxi-cab cost of transportation. It is noteworthy, such relationships, even though apparent, are infeasible to appear in cost functions that inherit properties solely from the information theory.

**t-SNE of sentence embeddings**  Next, we explain how we analyse the sentence-embeddings in $\mathcal{M}_g$ obtained from *Sink-E\**. A sentence refers to an Amazon food/product review. BERT's input sentence is lowercase WordPiece tokenized. We prepend the list of tokens with [CLS] token to represent the sentence which is later used for the classification task. First, each token is mapped to a static context-independent embedding. Then the vector list is passed through a sequence of
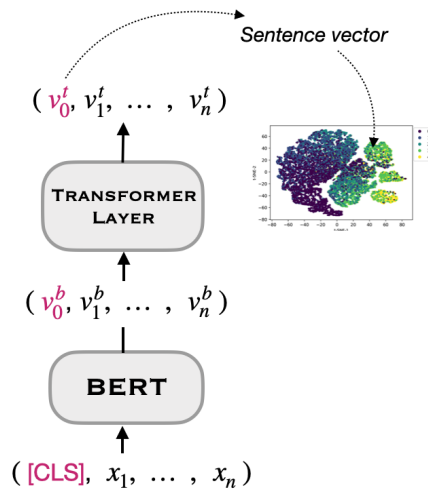
Figure 5: After the model parameter fitting is complete, we map the 128-dimensional [CLS] vector at the output of the Trasnformer layer to 2-dimensions usgin t-SNE.

multi-head self-attention operations that contextualizes each token. It is important to note that contextualization can be task-specific. We randomly sample 5000 review-label pairs for each sentiment class. For each textual review, as shown in fig. 5, we use a 128-dimensional vector at the output of the transformer layer corresponding to the [CLS] token. This corresponds to the list of reviews represented in 128-dimensional vector space. To visualize the learned sentence representations, we map the vectors from 128-dimensional space to 2-dimensions using t-distributed Stochastic Neighbor Embedding (t-SNE)[3].

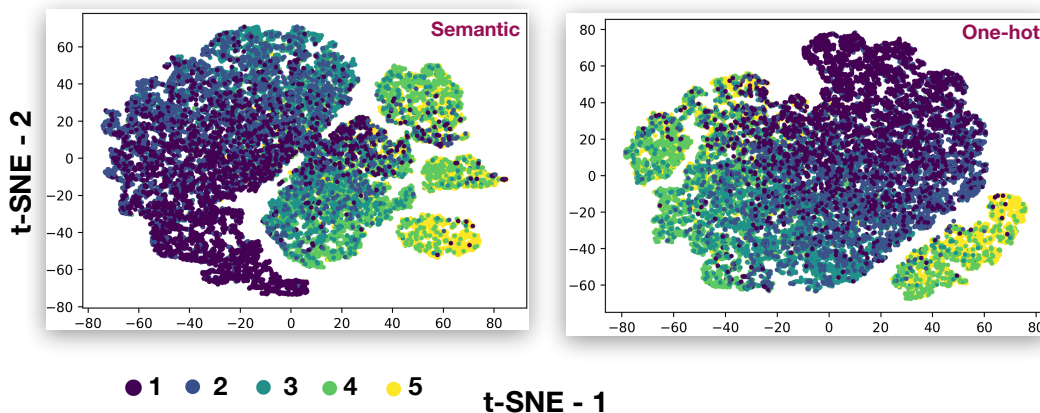## B.1   SEMANTIC SUPPORT → ONE-HOT SUPPORT



Figure 6: The one-hot encoding doesn't have clear distinct boundaries. The semantic structure is lost.

One way to understand the importance of properly defining label space relationship is by defining metric space where each label acquires its own axis, thus losing the semantic information. For a five-class classification problem, we will have five axes and thus the support is a set of five distinct one-hot vectors each of size five. This way any misclassification, i.e., predicting a mass different

---

[3]We used the implementation from scikit-learn.

from the ground truth label, will result in the same cost irrespective of positive sentiment is classified as strongly positive or strongly negative. This is due to the taxi-cab distance. Its value is computed by just summing up the absolute individual coordinate distances, which are just the predicted probabilities except for the coordinate corresponding to the ground truth.

We train the global model $\mathcal{M}_g$ with the Sinkhorn distance-based loss (eq. (4)) where the cost is defined as taxi-cab distance on the one-dimensional support and five-dimensional one-hot support as elaborated previously. The fig. 6 depicts the respective t-SNE scatter plots. In the plot with one-dimensional semantic support, we observe sentence vectors, i.e., features used for the classification task, are mapped in clearer clusters as compared to the plot at right without semantic information. We observe the points related to label 5 (strong positive) are much more localized as compared to the sentence mappings with one-hot, which is distributed around the space. This clearly dictates the benefit of a meaningful metric as compared to a space that is not informative. Next, we check a similar case that occurs in entropy-based loss functions.

## B.2  SINKHORN DISTANCE → KL DIVERGENCE

Similar to the cost associated with the one-hot support in Sinkhorn, the KL divergence has no feasible way to capture the intrinsic metric in the label space. The plots in fig. 7 show the different sentence embeddings (t-SNE) with the varying entropy-based regularisation term in the vanilla Sinkhorn distance. As $\varepsilon \to 0^+$, we should get a pure OT-based loss function (eq. (2)). However, to speed up the Sinkhorn and gradient computations, we chose $\varepsilon = 0.001$ with no (F1-score) performance trade-off. As shown in the fig. 7, with $\varepsilon >= 1$, the sentence representations are distributed across space with patches of label-dominant clusters. However, we can not see clear decision boundaries between the labels. As we decrease the $\varepsilon$ value below 1, we observe clearer feature maps for each label. For $\varepsilon = 0.001$, we can see clear sentence vector clusters corresponding to label 5. We can see the higher confusion rate is only between labels 5 and 4 that can be attributed to the less cost of transportation of the mass from label 5 to 4 as compared to 5 to other labels. A similar trend can be seen for lower $\varepsilon$ values that is 0.01 and the $\varepsilon$ used in this work 0.003 where clearer and localized clusters can be seen.

## C  MODEL BIAS

### C.1  PROBABILITY SKEW

To generate a random input, we uniformly sample 200 tokens from the vocab [4] with replacement and join them with white space. We obtain 100,000 such random texts. For a given text classifier model, the skew value for sentiment label 1 can be estimated by the fraction of times it is the prediction of when the model infers over the set of random texts (section 5.2).

### C.2  GENDER AND RACIAL BIASES

Even though we considered the model probability skew as a reflection of bias induced from non-IID sampling, other biases such as gender and race can still be learned or acquired in the distillation process, For instance, take the following sentences:

My **father** said that the food is just fine . (review-1) → *strong positive*
My **mother** said that the food is just fine . (review-2) → *neutral*

Review-1 and 2 differ in gender-specific words which are **father** and **mother**. Since it is a sentiment classification task, ideally, the intelligent system should not learn gender-specific cues from the text to generate its predictions. However, we observe a gender dependence in both the KL divergence and confident Sinkhorn-based predictions.

Similarly, we curate an example where the reviews differ only in a race-specific word.

---

[4]We obtain the English vocabulary of size 30,522 from:https://huggingface.co/google/bert_uncased_L-2_H-128_A-2/tree/main.
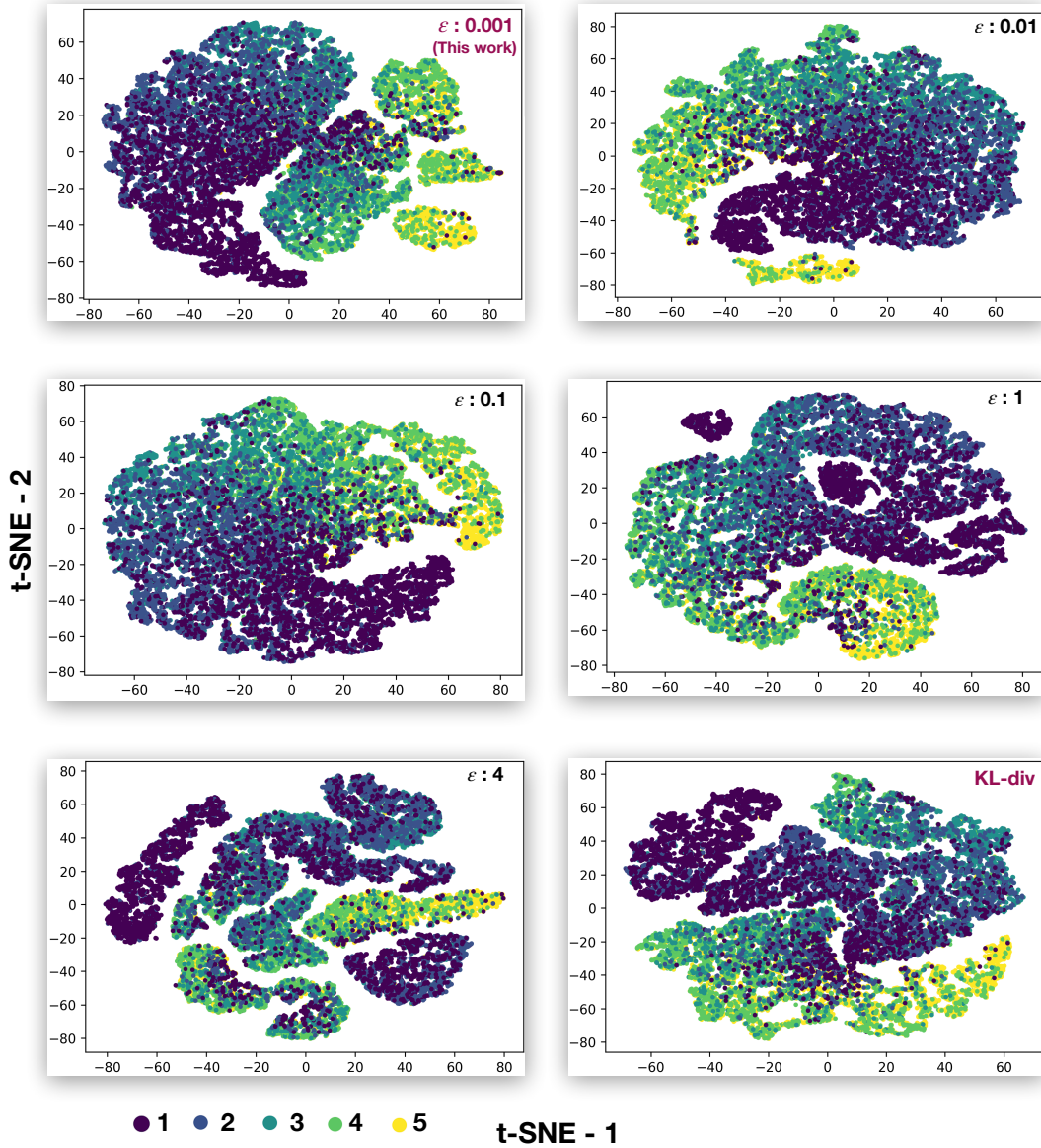
Figure 7: With a small entropic regularizer $\varepsilon$, it is visually striking that Sinkhorn seems to learn the latent structure boundaries better than KL divergence. For high $\varepsilon$, we see more clusters with mixed boundaries and not so clear demarcations.

> **White** guy said the phone is just fine . (review-1) → *neutral*
> **Latino** guy said the phone is just fine . (review-2) → *strong positive*

The sentiment predictions made by the intelligent systems were different contrary to ideal behavior. The review with word **White** shows a neutral sentiment while the review with word **Latino** shows a strong positive sentiment. Hence, the systems took account of the race-specific words while predicting the sentiment of a text. The behavior is observed both in the KL and Sinkhorn-based models with confident weights for sentiment classification.

## D  GRADIENTS THROUGH LOSS

We demonstrate the computation of gradient of loss function (eq. (4)) with respect to the global model trainable parameters $\theta$. We can write the Lagrange dual of eq. (2) as

$$T_\varepsilon \stackrel{\text{def.}}{=} \max_{(f,g) \in \mathcal{C}} \langle \mu, f \rangle + \langle \nu, g \rangle - \varepsilon \langle \mu \otimes \nu, \exp\left(\frac{1}{\epsilon}(f \oplus g - C)\right) - 1 \rangle \tag{6}$$

$$\mathcal{C} = \{(f,g) \in \mathbb{R}^{n_s \times n_t} : f_i + g_j \leq C_{(i,j)}\}$$

where $f \oplus g$ is tensor sum $(y_s, y_t) \in \mathcal{Y}_s \times \mathcal{Y}_t \mapsto f(y_s) + g(y_t)$. The optimal dual (solution of eq. (6)) can retrieve us the optimal transport plan (solution of eq. (4)) with the relation $\pi = \exp(\frac{1}{\varepsilon})(f \oplus g - C) \cdot (\mu \otimes \nu)$. Recently, a few interesting properties of $T_\epsilon$ were explored Peyré et al. (2019); Feydy et al. (2019); Luise et al. (2018) showing that optimal potentials $f$ and $g$ exist and are unique, and $\Delta T_\varepsilon(\mu, \nu) = (f, g)$.

---

**Algorithm 1** Gradients of $\mathcal{L}(h_\theta(x), h_\mathcal{K}(x))$ with respect to $h_\theta(x)$

---

**Initialize:** Dual potentials $\boldsymbol{f^1}, \cdots, \boldsymbol{f^K} \in \mathbb{R}^{n_s}$ and $\boldsymbol{g^1}, \cdots, \boldsymbol{g^K} \in \mathbb{R}^{n_t}$

1: **for** $k \leftarrow 1$ to $K$ **do**
2:    $\boldsymbol{f^k} \leftarrow \boldsymbol{0}$                                                             $\triangleright \boldsymbol{f^k} = \{f_1^k, \ldots, f_{n_s}^k\}$
3:    $\boldsymbol{g^k} \leftarrow \boldsymbol{0}$                                                             $\triangleright \boldsymbol{g^k} = \{g_1^k, \ldots, g_{n_t}^k\}$
4:    **while** ($\boldsymbol{f^k}, \boldsymbol{g^k}$ not converged ) **do**
5:        $f_i^k \leftarrow \varepsilon \, \mathrm{LSE}_{m=1}^{n_s}(log(h_k^m(x)) + \frac{1}{\varepsilon}g_m - \frac{1}{\varepsilon}C(\mathcal{Y}_s^i, \mathcal{Y}_t^m))$  $\left.\begin{array}{c} \\ \\ \end{array}\right\}$ Sinkhorn loop.
6:        $g_j^k \leftarrow \varepsilon \, \mathrm{LSE}_{m=1}^{n_t}(log(h_\theta^m(x)) + \frac{1}{\varepsilon}f_m - \frac{1}{\varepsilon}C(\mathcal{Y}_s^m, \mathcal{Y}_t^j))$
7:    (LSE is log-sum-exp reduction, i.e, $\mathrm{LSE}_{m=1}^M(V_m) = \log \sum_{m=1}^M exp(V_m)$)
8:    **end while**
9: **end for**
10: $\frac{\partial(T_\varepsilon(h_\theta(x), h_k(x)))}{\partial(h_\theta^i(x))} = f_i^k \quad \forall i \in [n_s], \ k \in \mathcal{K}$        $\triangleright$ as dual potentials are gradients of $T_\varepsilon$.
11: $\frac{\partial(\mathcal{L}_\varepsilon(h_\theta(x), h_\mathcal{K}(x)))}{\partial(h_\theta^i(x))} = \sum_{k \in \mathcal{K}} W_{B_k}(h_k(x)) f_i^k / \sum_{k \in \mathcal{K}} W_{B_k}(h_k(x))$.

---

Using these properties, we calculate gradients of the confident Sinkhorn cost in $eq.$ (4). Algorithm 1 obtains the gradients of the loss function with respect to $h_\theta(x)$ which can be backpropagated to tune model parameters. A crucial computation is to solve the coupling equation at step 5 and 6. This is done via Sinkhorn iterations which has a linear convergence rate Peyré et al. (2019).

## E  STATISTICAL RISK BOUNDS

Without the loss of generality, we will prove the risk bounds for two local models in learning without forgetting paradigm. For a sample $x$, let the output of local models be $y_1 = h_1(x)$ and $y_2 = h_2(x)$ and the global model with trainable parameters be $y_\theta = h_\theta(x)$. To prove Theorem 5.2, we consider the set of IID training samples $S = \{(x^{(1)}, y_1^{(1)}, y_2^{(1)}), \ldots, (x^{(N)}, y_1^{(N)}, y_2^{(N)})\}$.

**Lemma E.1.** *(from Frogner et al. (2015)) Let $h_{\hat{\theta}}, h_{\theta^*} \in \mathcal{H}_g$ be the minimizer of empirical risk $\hat{R}_S$ and expected risk $R$, respectively. Then*

$$R(h_{\hat{\theta}}) \leq R(h_{\theta^*}) + 2 \sup_{h_\theta \in \mathcal{H}_g} |R(h_\theta) - \hat{R}_S(h_\theta)| \tag{7}$$

To bound the risk for $h_{\hat{\theta}}$, we need to prove uniform concentration bounds for the distillation loss. We denote the space of loss functions induces by hypothesis space $\mathcal{H}_g$ as

$$L = \left\{ \ell_\theta : (x, y_1, y_2) \mapsto \frac{w_1(y_1)D(y_\theta, y_1) + w_2(y_2)D(y_\theta, y_2)}{w_1(y_1) + w_2(y_2)} \right\} \tag{8}$$

**Lemma E.2.** *(Frogner et al. (2015)) Let the transport cost matrix be $C$ and the constant $C_M = \max_{(i,j)} C_{(i,j)}$, then $0 \leq D(\cdot, \cdot) \leq C_M$, where $D(\cdot, \cdot)$ is 1-Wasserstein distance.*

**Definition E.3.** *(The Rademacher Complexity Bartlett & Mendelson (2002)). Let $\mathcal{G}$ be a family of mapping from $\mathcal{Z}$ to $\mathbb{R}$, and $S = (z_1, \ldots, z_N)$ a fixed sample from $\mathcal{Z}$. The empirical Rademacher complexity of $\mathcal{G}$ with respect to $S$ is defined as:*

$$\mathfrak{R}_S(\mathcal{G}) = \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{G}} \frac{1}{N} \sum_{i=1}^{n} \sigma_i g(z_i) \right] \tag{9}$$

*where $\sigma = (\sigma_1, \ldots, \sigma_N)$, with $\sigma_i$'s independent uniform random variables taking values in $\{+1, -1\}$. $\sigma_i$'s are called the Rademacher random variables. The Rademacher complexity is defines by taking expecation with respect to the samples $S$.*

$$\mathfrak{R}_N(\mathcal{G}) = \mathbb{E}_S [\hat{\mathfrak{R}}_S(\mathcal{G})] \tag{10}$$

**Theorem E.4.** *For any $\delta > 0$, with probability at least 1-$\delta$, the following holds for all $l_\theta \in L$:*

$$\mathbb{E}[\ell_\theta] - \hat{\mathbb{E}}[\ell_\theta] \leq 2\mathfrak{R}_N(L) + \sqrt{\frac{C_M^2 log(1/\delta)}{2N}}. \tag{11}$$

*Proof.* By definition $\mathbb{E}[\ell_\theta] = R(h_\theta)$ and $\hat{\mathbb{E}}[\ell_\theta] = \hat{R}(h_\theta)$. Let,

$$\Phi(S) = \sup_{\ell \in L} \mathbb{E}[\ell] - \hat{\mathbb{E}}_S[\ell].$$

Let $S$ and $S'$ differ only in sample $(\bar{x}^{(i)}, \bar{y}_1^{(i)}, \bar{y}_2^{(i)})$, by Lemma E.2, it holds that:

$$\Phi(S) - \Phi(S') \leq \sup_{\ell \in L} \hat{\mathbb{E}}_{S'} - \hat{\mathbb{E}}_S = \sup_{h_\theta \in \mathcal{H}} \frac{1}{N} \left\{ w_1(\bar{y}_1^{(i)})D(\bar{y}_\theta^{(i)}, \bar{y}_1^{(i)}) + w_2(\bar{y}_2^{(i)})D(\bar{y}_\theta^{(i)}, \bar{y}_2^{(i)}) \right.$$

$$\left. - w_1(y_1^{(i)})D(y_\theta, y_1^{(i)}) - w_2(y_2^{(i)})D(y_\theta, y_2^{(i)}) \right\} \leq \frac{2C_M}{N} \tag{12}$$

This inequality can be achieved by putting $D(\bar{y}_\theta^{(i)}, \bar{y}_1^{(i)}) = D(\bar{y}_\theta^{(i)}, \bar{y}_2^{(i)}) = C_M$ and $D(y_\theta^{(i)}, y_1^{(i)}) = D(y_\theta^{(i)}, y_2^{(i)}) = 0$.

Similarly, $\Phi(S') - \Phi(S) \leq C_M/N$, thus $|\Phi(S') - \Phi(S)| \leq C_M/N$. Now, from the McDiarmid's inequality McDiarmid (1998) and its usage in Frogner et al. (2015), we can establish

$$\Phi(S) \leq \mathbb{E}[\Phi(S)] + \sqrt{\frac{KC_M^2 log(1/\delta)}{2N}}. \tag{13}$$

From the bound established in the proof of Theorem B.3 in Frogner et al. (2015), i.e., $\mathbb{E}_S[\Phi(S)] \leq 2\mathfrak{R}_N(L)$, we can conclude the proof. □

To complete the proof of Theorem Theorem 5.2, we have to treat $\mathfrak{R}_N(L)$ in terms of $\mathfrak{R}_N(\mathcal{H}_g)$.

Now, let $\iota : \mathbb{R}^{|\mathcal{Y}|} \times \mathbb{R}^{|\mathcal{Y}|} \mapsto \mathbb{R}$ defined by $\iota(y, y') = D(\mathfrak{s}(y), \mathfrak{s}(y)')$, where $\mathfrak{s}$ is a softmax function defined over the vector of logits. From Proposition B.10 of Frogner et al. (2015), we know:

$$|\iota(y, y') - \iota(\bar{y}, \bar{y}')| \leq 2C_M(||y - \bar{y}||_2 + ||y' - \bar{y}'||_2) \tag{14}$$

Let $\iota_s : \mathbb{R}^{|\mathcal{Y}|} \times \mathbb{R}^{|\mathcal{Y}|} \times \mathbb{R}^{|\mathcal{Y}|} \mapsto \mathbb{R}$ defined by:

$$\iota_s(y, y_1, y_2) = \frac{w_1(\mathfrak{s}(y_1))D(\mathfrak{s}(y), \mathfrak{s}(y_1)) + w_2(\mathfrak{s}(y_2))D(\mathfrak{s}(y), \mathfrak{s}(y_2))}{w_1(\mathfrak{s}(y_1)) + w_2(\mathfrak{s}(y_2))} \tag{15}$$

$$= \bar{w}_1(\mathfrak{s}(y_1), \mathfrak{s}(y_2)) \, D(\mathfrak{s}(y), \mathfrak{s}(y_1)) + \bar{w}_2(\mathfrak{s}(y_1), \mathfrak{s}(y_2)) \, D(\mathfrak{s}(y), \mathfrak{s}(y_2)) \tag{16}$$

where $w_1(.), w_2(.)$ are confidence score of local model predictions $y_1, y_2$ on an input $x$. $\bar{w}_1(.), \bar{w}_2(.)$ are normalized scores. Note that the local model predictions, i.e., $y_1$ and $y_2$ are functions of $x$, where $x$ is sampled from the data domain distribution $f(x)$. Hence, we can view the loss function as

$$\iota_s(y, y_1, y_2) = D(\mathfrak{s}(y), \mathfrak{s}(y_1)) + D(\mathfrak{s}(y), \mathfrak{s}(y_2)) \tag{17}$$

$$= \iota(y, y_1) + \iota(y, y_2). \tag{18}$$

where $y$ is a function of $x_{new}$ sampled from a weighted distribution $\bar{w}_1(\mathfrak{s}(y_1), \mathfrak{s}(y_2))f(x)$.

The *Lipschitz* constant of $\iota_s(y, y_1, y_2)$ can thus be identified by:

$$|\iota_s(y, y_1, y_2) - \iota_s(\bar{y}, \bar{y}_1, \bar{y}_2)| = |\iota(y, y_1) + \iota(y, y_2) - \iota(\bar{y}, \bar{y}_1) - \iota(\bar{y}, \bar{y}_2)| \tag{19}$$

$$\leq |\iota(y, y_1) - \iota(\bar{y}, \bar{y}_1)| + |\iota(y, y_2) + \iota(\bar{y}, \bar{y}_2)| \tag{20}$$

$$\leq 2C_M(||y - \bar{y}||_2 + ||y_1 - \bar{y}_1||_2 + ||y - \bar{y}||_2 + ||y_2 - \bar{y}_2||_2) \tag{21}$$

$$\leq 4C_M(||y - \bar{y}||_2 + ||y_1 - \bar{y}_1||_2 + ||y_2 - \bar{y}_2||_2) \tag{22}$$

$$\leq 4C_M||(y, y_1, y_2) - (\bar{y}, \bar{y}_1, \bar{y}_2)||_2 \tag{23}$$

Thus, the *Lipschitz* constant of plain Sinkhorn based distillation is $4C_M$.

**Proof of Theorem 5.2** We define the space of loss function for $k$ local models:

$$L = \left\{ \iota_\theta : (x, \{y_k\}_{k \in \mathcal{K}}) \mapsto \sum_{k \in \mathcal{K}} w_k(\mathfrak{s}(h_\theta^o(x)))D(\mathfrak{s}(y_k)) \right\}$$

Following the notations in Frogner et al. (2015), we apply the following generalized Talagrand's lemma Ledoux & Talagrand (2013):

**Lemma E.5.** *Let $\mathcal{F}$ be a class of real functions, and $\mathcal{H} \subset \mathcal{F} = \mathcal{F}_1 \times \ldots \times \mathcal{F}_K$ be a $K$-valued function class. If $m : \mathbb{R}^K \mapsto \mathbb{R}$ is a $L_\mathfrak{m}$-Lipschitz function and $\mathfrak{m}(0) = 0$, then $\mathfrak{R}_S(\mathfrak{m} \circ \mathcal{H}) \leq 2L_\mathfrak{m} \sum_{k=1}^{K} \hat{\mathfrak{R}}_S(\mathcal{F}_k)$.*

Now, as the Lemma can not be directly applied to the confident Sinkhorn loss as $\mathbf{0}$ is an invalid input. To get around the problem, we assume the global hypothesis space is of the form:

$$\mathcal{H} = \{\mathfrak{s} \circ h : h^o \in \mathcal{H}^o\} \tag{24}$$

Thus, we apply the lemma to the $4C_M$-*Lipschitz* continuous function $\iota$ and the function space:

$$\underbrace{\mathcal{H}^o \times \ldots \times \mathcal{H}^o}_{|\mathcal{Y}|copies} \underbrace{\times \mathcal{I} \times \ldots \times \mathcal{I}}_{|\mathcal{Y}| \times |\mathcal{K}|copies}$$

with $\mathcal{I}$ a singleton function space of identity maps. It holds:

$$\mathfrak{R}_N(L) \leq 8C_M(|\mathcal{Y}|\hat{\mathfrak{R}}_N + |\mathcal{Y}| \times |\mathcal{K}|\hat{\mathfrak{R}}_N(\mathcal{I})) = 8|\mathcal{Y}|C_M\hat{\mathfrak{R}}_N(\mathcal{H}^o) \tag{25}$$

As,

$$\hat{\mathfrak{R}}_N(\mathcal{I}) = \mathbb{E}_\sigma \left[ \sup_{g \in \mathcal{I}} \frac{1}{N} \sum_{i=1}^{N} \sigma_i g(y_i) = 0 \right] = \mathbb{E}_\sigma \left[ \frac{1}{N} \sum_{i=1}^{N} \sigma_i y_i = 0 \right]$$

Thus, by combining eq. (25) in Theorem E.4 and Lemma E.1 proves the Theorem 5.2.

Since the $\varepsilon$ is small in our experiments, we can quantify the difference between Sinkhorn distance and Wasserstein for a given *Lipschitz* cost function.

## F SOCIETAL IMPACT

Our work can be extended to different domains. Although in this paper we examined sentiment classifications, other areas, where labels are not available, i.e., zero-shot classification, would also be amenable to federated confident Sinkhorns. Within our approach, a potential downstream task could be to detect cyberbullying. An important area of application for distraught parents, school teachers and teens. In this case, a sentiment that has a high probability of been classified as cyberbullying can be flagged to either moderators or guardians of a particular application.

A weakness of this approach is that the training of such an application will be based on local models in other domains. Care would be needed in deciding which local models to use in the federation. This choice is highly dependent on the industry and the availability of data. Misuse of our approach could be that the federated training might distill some population-specific information to the global model which makes the central system vulnerable to attacks that might lead to a user-private data breach. As GPU implementation of OT metrics becomes commonplace, we envisage that our approach might help in other Natural Language Processing (NLP) tasks. Indeed, this would be potentially beneficial and open up new avenues to the NLP community.

## G EXPERIMENTAL REPRODUCIBILITY

All the experiments were performed on one Quadro RTX 8000 GPU with 48 GB memory. The model architecutres were designed on Python (version 3.9.2) library PyTorch (version 1.8.1) under *BSD-style license*. For Sinkhorn iterations and graident calculations, we use GeomLoss library (version 0.2.4) (`https://github.com/jeanfeydy/geomloss`) under *MIT licence*. For barycenter calculations, we use POT library (0.7.0) from (`https://pythonot.github.io/`) under *MIT licence*.

We clip all the text to a maximum length of 200 tokens and pad the shorter sentences with ¡unk¿. To speed up the experiments, we use pretrained BERT-Tiny from `https://github.com/google-research/bert`.

The batch size is chosen via grid search from the set $\{16, 32, 64, 128, 256, 512, 1024, 2048\}$ and found 1024 to be optimal for performance and speed combination on the considered large datasets. We use Adam optimizer with learning rate chosen via grid search $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3, 10^4\}$. All the experiments were ran for 20 epochs. The regularization parameter $\varepsilon$ is chosen based on minimal loss obtained amongst the set of $\varepsilon$ values $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 2, 4, 8, 16\}$.

For the Amazon review dataset, we were unable to find the licence.