
CalFAT: Calibrated Federated Adversarial Training with Label Skewness

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Recent studies have shown that, like traditional machine learning, federated learn-
2 ing (FL) is also vulnerable to adversarial attacks. To improve the adversarial
3 robustness of FL, few federated adversarial training (FAT) methods have been
4 proposed to apply adversarial training locally before global aggregation. Although
5 these methods demonstrate promising results on independent identically distributed
6 (IID) data, they suffer from training instability issues on non-IID data with label
7 skewness, resulting in much degraded natural accuracy. This tends to hinder the
8 application of FAT in real-world applications where the label distribution across
9 the clients is often skewed. In this paper, we study the problem of FAT under label
10 skewness, and firstly reveal one root cause of the training instability and natural
11 accuracy degradation issues: skewed labels lead to non-identical class probabilities
12 and heterogeneous local models. We then propose a Calibrated FAT (CalFAT) ap-
13 proach to tackle the instability issue by calibrating the logits adaptively to balance
14 the classes. We show both theoretically and empirically that the optimization of
15 CalFAT leads to homogeneous local models across the clients and much improved
16 convergence rate and final performance.

17 1 Introduction

18 Federated learning (FL) is a privacy-aware learning paradigm that allows multiple participants
19 (clients) to collaboratively train a global model without sharing their private data [21]. In FL, each
20 client follows the conventional machine learning procedure to train a local model on its own data
21 and periodically uploads the local model updates to a central server for global aggregation. However,
22 recent studies have shown that, like conventional machine learning, FL is also vulnerable to well-
23 crafted adversarial examples [39, 8, 38], i.e., at inference time, the attackers can add small, human-
24 perceptible perturbations to the test examples to fool the global model to make misclassifications
25 with high success rates. This raises security and reliability concerns on the implementation of FL
26 in real-world scenarios where such a vulnerability could cause heavy losses [34]. For example, for
27 cross-silo FL in biomedical domain, a vulnerable global model may cause misdiagnosis, wrong
28 medical treatments, or even the loss of lives. It is thus imperative to develop a robust FL method that
29 can train adversarially robust global models resistant to different types of adversarial attacks.

30 In conventional machine learning, adversarial training (AT) has been shown to be one of the most
31 effective defenses against adversarial attacks [20, 36]. Since the local training in FL is the same as
32 conventional machine learning, recent works [39, 8, 38] proposed to perform local AT to improve
33 the adversarial robustness of the global model. These methods in general are known as Federated
34 Adversarial Training (FAT). AT has been found to be *more* challenging than standard training [3,
35 37], as it generally requires more training data and larger-capacity models. Moreover, adversarial
36 robustness may be at odds with accuracy [26], meaning that the increase of robustness may inevitably
37 decrease the natural accuracy (i.e., accuracy on natural test data). As a result, the natural accuracy

38 of AT is much lower than standard training [4]. This phenomenon also exists in FL, i.e., FAT
 39 delivers much slower convergence and lower natural accuracy than standard FL, as shown in recent
 40 studies [39, 8].

41 Arguably, FAT will become more challenging if the data are non-independent and identically dis-
 42 tributed (non-IID) across the clients. One typical non-IID setting that commonly exists in real-world
 43 applications is the skewed label distribution [16], where different clients have different label distri-
 44 butions. In this paper, we study the problem of FAT on non-IID data with a particular focus on the
 45 challenging skewed label distribution setting (formally defined in Section 3.1). Under conventional
 46 training, Xu et.al., [33] have showed that adversarially trained models introduce severe performance
 47 disparity across different classes. And such disparity will be exacerbated under label skewness,
 48 ending up with much worse performance on the minority classes [29].

49 By far, a few works have studied non-IID FAT
 50 in the current literature. Zizzo et.al., [39] pro-
 51 posed to perform AT on only a proportion of
 52 local data for better convergence, while stan-
 53 dard training is applied on the rest of the lo-
 54 cal data. We term this method as MixFAT. An-
 55 other relevant work called FedRBN [8] tackled
 56 a different problem: how to propagate federated
 57 robustness to low-resource clients. Although
 58 MixFAT and FedRBN demonstrated promising
 59 results, they suffer from training instability is-
 60 sues and much lower natural accuracy compared
 61 to standard FL, as we show in Figure 1. We
 62 also compare with the other four FAT baselines
 63 that apply different adversarial training methods
 64 to FL, i.e., FedPGD, FedTRADES, FedMART,
 65 and FedGAIRAT. Unfortunately, these methods
 66 also suffer from slow convergence and much de-
 67 graded final accuracy (details can be referred to
 68 Section 4.1). This motivates us to propose a novel method called *Calibrated Federated Adversarial*
 69 *Training* (CalFAT) for tackling FAT with label skewness. CalFAT tackles training instability caused
 70 by label skewness by calibrating the logits to give higher scores to the minority classes.

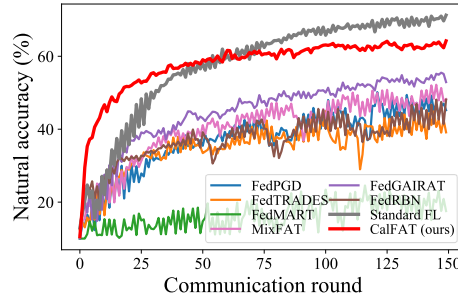


Figure 1: Natural accuracy and convergence rate of standard FL, our proposed CalFAT, and 6 FAT baselines (FedPGD, FedTRADES, FedMART, MixFAT [39], FedGAIRAT and FedRBN [8]) under a highly skewed label distribution with $\beta = 0.1$ (Section 4).

71 In summary, our main contributions are:

- 72 • *New insight:* We study the problem of FAT on non-IID data with skewed label distribution, and
 73 reveal one root cause of the training instability and natural accuracy degradation issues: skewed
 74 labels lead to non-identical class probabilities and heterogeneous local models.
- 75 • *Novel method:* We propose a novel method called CalFAT for FAT with label skewness, and show
 76 that the optimization of CalFAT leads to homogeneous local models, and consequently, stable
 77 training, faster convergence, and better final performance.
- 78 • *Compelling results:* Extensive experiments on 4 benchmark vision datasets across various settings
 79 prove the effectiveness of our CalFAT and its superiority over existing FAT methods.

80 2 Notation and Preliminaries

81 2.1 Notation

82 Suppose there are m clients in FL with i denoting the i -th client, e.g., \mathcal{D}_i denotes the local data of
 83 client i and θ_i denotes the parameters of its local model. We use $\hat{\theta}$ to denote the parameters of the
 84 global model. Subscript j is the sample index, e.g., (x_{ij}, y_{ij}) denotes the j -th sample of client i and
 85 its corresponding label with $y_{ij} \in \{1, \dots, C\}$. Let $f_{\theta}(\cdot)$ be the local model $f(\cdot)$ (before softmax)
 86 with parameter θ . Superscript l is the class index, e.g., $f^l(\cdot)$ denotes the logit output for class l . We
 87 denote the adversarial example of clean sample x by \tilde{x} . $[m]$ denotes the integer set $\{1, \dots, m\}$.
 88 $p_i(x, y)$ denotes the joint distribution of input x and label y at client i , and accordingly, $p_i(y)$ is the
 89 marginal distribution of label y , $p_i(y | x)$ is the conditional distribution of label y given input x and
 90 $p_i(x | y)$ is the conditional distribution of input x given label y .

91 2.2 Centralized Adversarial Training

92 Let $\mathcal{D} = \{x_j, y_j\}_{j=1}^n$ be the training dataset with n training samples. The cross-entropy loss
 93 $\ell_{ce}(f_\theta(x), y)$ for an input-label pair (x, y) is defined as $\ell_{ce}(f_\theta(x), y) = \sigma^y(f_\theta(x))$, where $\sigma^y(f) =$
 94 $\exp(f^y) / \sum_{l=1}^C \exp(f^l)$ is the softmax function, C is the number of classes, and f^l is the model
 95 output for class l . The objective function of the centralized adversarial training (AT) [20] can then
 96 be defined as $\min_{\theta} \sum_{j=1}^n \ell_{ce}(f_\theta(\tilde{x}_j), y_j) / n$, where the adversarial example \tilde{x}_j can be generated by
 97 $\tilde{x}_j = \arg \max_{x'_j \in \mathcal{B}_\epsilon(x_j)} \ell_{ce}(f_\theta(x'_j), y_j)$, where $\mathcal{B}_\epsilon(x_j) = \{x' \mid \|x' - x_j\|_\infty < \epsilon\}$ is the closed ball
 98 of radius $\epsilon > 0$ centered at x_j , $\|\cdot\|_\infty$ is the L_∞ norm, and \tilde{x}_j is the most adversarial sample within
 99 the ϵ -ball.

100 A standard centralized AT method uses Projected Gradient Decent (PGD) to generate adversarial
 101 examples [20]. In particular, PGD iteratively generates adversarial example as follows: $x_j^{(k+1)} =$
 102 $\Pi_{\mathcal{B}_\epsilon(x_j^{(0)})} \left(x_j^{(k)} + \alpha \text{sign}(\nabla_x \ell_{ce}(f_\theta(x_j^{(k)}), y_j)) \right)$, where $k = 0, \dots, K-1$ is the step number, K is
 103 the total number of steps, $\alpha > 0$ is the step size, $x_j^{(0)}$ is the natural sample, $x_j^{(k)}$ is the adversarial
 104 example generated at step k , $\Pi_{\mathcal{B}_\epsilon(x_j^{(0)})}$ is the projection function that projects the adversarial data
 105 onto the ϵ -ball centered at $x_j^{(0)}$, and $\text{sign}(\cdot)$ is the sign function. The adversarial example obtained at
 106 the final step of PGD is used as the adversarial example, i.e., $\tilde{x}_j = x_j^{(K)}$.

107 By optimizing the model parameters on the adversarial examples generated by PGD, centralized AT
 108 is able to train a model that is robust against adversarial attacks.

109 2.3 Federated Adversarial Training

110 The concept of federated adversarial training (FAT) was first introduced by [39] (We term this method
 111 as MixFAT) to deal with the adversarial vulnerability of FL. MixFAT applied AT locally to improve
 112 the robustness of the global model. Suppose there are m clients and each client i has its local data
 113 $\mathcal{D}_i = \{x_{ij}, y_{ij}\}_{j=1}^{n_i}$ sampled from distribution $p_i(x, y)$ with $n_i = |\mathcal{D}_i|$ being the size of the local
 114 data. In MixFAT, each client i optimizes its local model by minimizing the following objective:

$$\min_{\theta_i} \frac{1}{n_i} \left(\sum_{j=1}^{n'_i} \ell_{ce}(f_{\theta_i}(\tilde{x}_{ij}), y_{ij}) + \sum_{j=n'_i+1}^{n_i} \ell_{ce}(f_{\theta_i}(x_{ij}), y_{ij}) \right), \quad (1)$$

115 where \tilde{x}_{ij} is the PGD adversarial example of x_{ij} , n'_i is a hyperparameter that controls the ratio of
 116 data for AT, and θ_i are the local model parameters. After training the local model for certain epochs,
 117 client i uploads its local model parameters θ_i to the central server for aggregation. Note that MixFAT
 118 only applies AT to a proportion of the local data, mainly for convergence and stability considerations.

119 3 Calibrated Federated Adversarial Training (CalFAT)

120 3.1 Skewed Label Distribution Leads to Non-identical Class Probabilities

121 In this paper, we focus on non-IID setting with *skewed label distribution* [18, 9], which is defined as
 122 follows.

123 **Definition 1** (Skewed label distribution). The label distribution across the clients is skewed, if for
 124 $\forall i \neq u$ and $i, u \in [m]$:

125 (a) there exists $y \in [C]$ such that $p_i(y) \neq p_u(y)$ and (b) $p_i(x | y) = p_u(x | y)$ for $\forall x, y$.

126 Condition (b) is to assume that, given a class y , x is sampled with equal probability at different clients.
 127 Note that there exist different types of non-IID: label skew, non-identical class conditional, quantity
 128 skew, to name a few (Appendix K in [9]). The class conditional is often assumed to be identical
 129 (i.e., condition (b)) when studying the label skewness problem, which is the main focus of this work.
 130 When condition (b) does not hold, it becomes the non-identical class conditional problem.

131 **Lemma 1** (Non-identical class probabilities). *If the label distribution across the clients is skewed*
 132 *and the class conditionals have the same support, then the class probabilities $\{p_i(y | x) \mid i \in [m]\}$*
 133 *are non-identical, i.e., for all $i \neq u$ and $i, u \in [m]$, there exist x, y such that $p_i(y | x) \neq p_u(y | x)$.*

134 The proof of Lemma 1 is given in Appendix A. Lemma 1 implies that skewed label distribution gives
 135 rise to non-identical class probabilities $\{p_i(y | x) | i \in [m]\}$.

136 3.2 Standard Cross-entropy Leads to Heterogeneity

137 From a statistical point of view, each client i in previous FAT methods estimates its local class
 138 probability $p_i(y | x)$ during local training [6]. More specifically, they assume that $p_i(y | x)$ can be
 139 parameterized by θ_i^* as:

$$p_i(y | x) = \hat{p}(y | x; \theta_i^*) = \sigma^y(f_{\theta_i^*}(x)), \quad (2)$$

140 where θ_i^* is the ground-truth parameters of the local class probability $p_i(y | x)$. According to
 141 Lemma 1, the class probabilities $\{p_i(y | x)\}$ are non-identical when there is a skewed label distribu-
 142 tion. Therefore, the ground-truth parameters $\{\theta_i^* | i \in [m]\}$ are heterogeneous. We use the sample
 143 variance of the ground-truth parameters to measure such heterogeneity as follows:

$$(s^*)^2 = V(\theta_1^*, \dots, \theta_m^*) = \frac{1}{m-1} \sum_{i=1}^m \|\theta_i^* - \frac{1}{m} \sum_{j=1}^m \theta_j^*\|^2. \quad (3)$$

144 Each client i updates its local model parameters θ_i by optimizing the standard cross-entropy (CE)
 145 loss. The updated θ_i is the maximum likelihood estimate [1] of the ground-truth parameter θ_i^* [6].
 146 We use the sample variance [1] of local model parameters to measure heterogeneity of local models:

$$s^2 = V(\theta_1, \dots, \theta_m). \quad (4)$$

147 Larger sample variance implies higher model heterogeneity.

148 The following proposition suggests that the heterogeneity of local models originates from the
 149 heterogeneity of local class probabilities.

150 **Proposition 1** (Heterogeneous local models). *Assume the label distribution across the clients is*
 151 *skewed. Let θ_i be the maximum likelihood estimate of θ_i^* in Eq. (2) given local data at client i . Then*
 152 *s^2 converges almost surely to a nonzero constant:*

$$s^2 \xrightarrow{a.s.} (s^*)^2 \neq 0, \quad (5)$$

153 where $\xrightarrow{a.s.}$ represents the almost sure convergence.

154 The proof of Proposition 1 is provided in Appendix B. $(s^*)^2$ measures the heterogeneity of ground-
 155 truth parameters $\{\theta_i^* | i \in [m]\}$, which reflects the class probability difference across the clients as
 156 shown in Eq. (2).

157 Proposition 1 implies that the local models in previous FAT methods are heterogeneous when the
 158 label distribution across the clients is skewed. Since the local models are heterogeneous, aggregating
 159 these models potentially hurts convergence or even results in divergence of the global model [16]. As
 160 shown in Figure 1, the training of previous FAT methods are unstable and have much lower natural
 161 accuracy than the standard FL.

162 3.3 Learning Homogeneous Local Models by Calibration

163 Motivated by [22], we propose to re-parameterize the class probabilities. According to Bayes' formula
 164 [11],

$$p_i(y | x) = \frac{p_i(x | y)p_i(y)}{\sum_{l=1}^C p_i(x | l)p_i(l)}. \quad (6)$$

165 On the right-hand side of the above equation: (1) the class priors can be easily approximated by the
 166 relative frequencies [1]; and (2) more importantly, the class conditionals $\{p_i(x | y) | i \in [m]\}$ are
 167 *identical* across different clients (see Definition 1).

168 Inspired by the above observation, we propose an alternative parameterization of $p_i(y | x)$. Assume
 169 that for all $i \in [m]$, the class conditional $p_i(x | y)$ can be parameterized by θ^* as $p_i(x | y) = \hat{q}(x |$

Algorithm 1 Local training of CalFAT

Input: Client i , global model parameters $\hat{\theta}$, local dataset \mathcal{D}_i , and local epoch number E

```
1: procedure CLIENTUPDATE
2:    $\theta_i \leftarrow \hat{\theta}$ 
3:   Compute  $\pi_i$  with  $\mathcal{D}_i$  by Eq. (8)
4:   for local epoch=1,  $\dots$ ,  $E$  do
5:     for  $j = 1, \dots, n_i$  do
6:       Sample  $(x_{ij}, y_{ij})$  from  $\mathcal{D}_i$ 
7:       Generate adversarial example  $\tilde{x}_{ij}$  by Eq. (12)
8:     end for
9:      $\theta_i \leftarrow \theta_i - \eta \frac{1}{n_i} \sum_{j=1}^{n_i} \nabla_{\theta_i} \ell_{cce}(f_{\theta_i}(\tilde{x}_{ij}), y_{ij}, \pi_i)$ 
10:  end for
11:  return  $\theta_i$ 
12: end procedure
```

170 $y; \theta^*$), where $\hat{q}(x | y; \theta^*)$ can be an arbitrary conditional probability function. Then, $p_i(y | x)$ can be
171 re-parameterized by θ^* as follows:

$$p_i(y | x) = \hat{q}_i(y | x; \theta^*) = \frac{\hat{q}(x | y; \theta^*) \pi_i^y}{\sum_{l=1}^C \hat{q}(x | l; \theta^*) \pi_i^l}. \quad (7)$$

172 where

$$\pi_i^y = n_i^y / n_i + \delta, y \in [C]. \quad (8)$$

173 Here π_i^y approximates the class prior $p_i(y)$, n_i^y is the sample size of class y on client i and $\delta > 0$ is a
174 small constant added for numerical stability purpose. During local updates, client i uses its local data
175 to update θ_i , which makes θ_i the maximum likelihood estimate of θ^* . The entire training procedure
176 of our method is described in Section 3.4.

177 The following proposition suggests that the local models are homogeneous when trained with the
178 above re-parameterization. The proof of Proposition 2 is provided in Appendix C.

179 **Proposition 2** (Homogeneous local models). *Assume the label distribution across the clients is*
180 *skewed. Let θ_i be the maximum likelihood estimate of θ^* in Eq. (7) given local data at client i . Then*
181 *s^2 converges almost surely to zero:*

$$s^2 \xrightarrow{a.s.} 0. \quad (9)$$

182 3.4 Details of CalFAT

183 The local training procedure of our proposed Calibrated Federated Adversarial Training (CalFAT) is
184 described in Algorithm 1. Specifically, we define $\hat{q}(x | y; \theta^*) = \exp(f_{\theta^*}^y(x))$. Then, we maximize
185 the likelihood of $\hat{q}_i(y | x; \theta^*)$ for each client i , which is equivalent to minimizing the following
186 objective:

$$\min_{\theta_i} \frac{1}{n_i} \sum_{j=1}^{n_i} \ell_{cce}(f_{\theta_i}(\tilde{x}_{ij}), y_{ij}, \pi_i), \quad (10)$$

187 where $\ell_{cce}(\cdot, \cdot, \cdot)$ is the calibrated cross-entropy (CCE) loss and \tilde{x}_{ij} is the adversarial example of x_{ij} .
188 The CCE loss is defined as:

$$\ell_{cce}(f_{\theta_i}(\tilde{x}_{ij}), y_{ij}, \pi_i) = -\log \sigma^{y_{ij}}(f_{\theta_i}(\tilde{x}_{ij})) + \log \pi_i. \quad (11)$$

189 As discussed in Section 3.3, minimizing the above CCE loss mitigates the heterogeneity of the local
190 models, which can lead to improved convergence and performance of the global model.

191 In previous FAT methods, heterogeneous local models tend to give higher scores to the majority
192 classes while lower scores to the minority classes. By contrast, our CalFAT encourages local models
193 to give higher scores to the minority classes by adding a class-wise prior $\log \pi_i^l$ to the logits. Also
194 different from MixFAT that trains the local models on both natural and adversarial data, our CalFAT
195 trains the local models only on adversarial examples. Extensive empirical experiments are conducted
196 in Section 4.1 to show the impact of using only adversarial data for optimization.

Table 1: Natural and robust accuracy (%) on different datasets. The best results are in **bold**.

Dataset	SVHN						CIFAR10						CIFAR100					
	Natural	FGSM	BIM	CW	PGD-20	AA	Natural	FGSM	BIM	CW	PGD-20	AA	Natural	FGSM	BIM	CW	PGD-20	AA
MixFAT	19.64	19.63	19.64	19.63	19.65	14.71	53.23	29.02	26.37	22.79	26.22	21.81	34.39	15.79	14.63	11.34	14.50	9.31
FedPGD	19.47	19.46	19.45	19.46	19.47	13.67	47.21	28.80	26.68	24.53	26.50	22.79	34.06	16.10	14.76	11.53	14.70	10.80
FedTRADES	56.84	36.96	35.11	31.16	34.97	30.56	46.14	27.68	26.36	22.81	26.29	21.61	29.35	14.99	14.20	10.52	14.23	9.56
FedMART	19.81	19.80	19.79	19.80	19.81	14.56	25.68	18.49	18.15	15.39	18.15	14.31	19.85	13.01	12.77	9.99	12.79	8.67
FedGAIRAT	58.42	38.34	36.46	31.24	36.63	32.52	48.34	29.32	26.43	22.83	27.32	21.93	35.15	16.13	15.27	11.82	14.91	9.54
FedRBN	53.87	34.60	32.68	28.16	32.49	28.44	47.87	26.82	26.21	21.97	26.21	21.43	28.62	14.74	13.35	9.81	14.21	8.88
CalFAT(ours)	84.05	48.54	42.04	31.69	41.64	32.73	64.85	35.04	31.54	24.62	31.19	22.94	44.50	17.79	15.64	12.05	15.39	11.32

197 **Adversarial example generation.** Inspired by [36], we generate the adversarial examples by
 198 maximizing the following calibrated Kullback–Leibler (CKL) divergence loss:

$$\tilde{x}_{ij} = \arg \max_{x'_{ij} \in \mathcal{B}_\epsilon(x_{ij})} \ell_{ckl}(f_{\theta_i}(x'_{ij}), f_{\theta_i}(x_{ij}), \pi_i), \quad (12)$$

199 where $\ell_{ckl}(\cdot, \cdot, \cdot)$ is the CKL loss defined as:

$$\ell_{ckl}(f_{\theta_i}(x'_{ij}), f_{\theta_i}(x_{ij}), \pi_i) = - \sum_{y=1}^C \sigma^y(f_{\theta_i}(x_{ij}) + \log \pi_i) \log \sigma^y(f_{\theta_i}(x'_{ij}) + \log \pi_i), \quad (13)$$

200 where $\log \pi_i$ is the same as in our CCE loss. Following centralized AT [20], we also use PGD to
 201 solve Eq. (12).

202 After training the local model for certain epochs following the above procedure, each client i uploads
 203 the model parameters θ_i to the server for aggregation. To be consistent with the most recent FAT
 204 methods [39, 8], we first use the most widely used FedAvg [21] as the default aggregation framework.
 205 However, we remark that our method is agnostic to FL frameworks, i.e., it is compatible with any
 206 other FL frameworks (e.g., FedProx [15] and Scaffold [10]), as shown in Section 4.1.

207 4 Experiments

208 **Data Configurations.** Our experiments are conducted on 4 real-world datasets: CIFAR10 [12],
 209 CIFAR100 [12], SVHN [23], and ImageNet subset [5]. To simulate label skewness, we sample
 210 $p_i^l \sim Dir(\beta)$ and allocate a p_i^l proportion of the data of label l to client i , where $Dir(\beta)$ is the
 211 Dirichlet distribution with a concentration parameter β [35]. By default, we set $\beta = 0.1$ to simulate a
 212 highly skewed label distribution that widely exists in reality.

213 **Baselines.** We compare our proposed CalFAT with two state-of-the-art FAT methods: MixFAT [39]
 214 and FedRBN [8]. We also investigate the combination of the state-of-the-art centralized AT methods
 215 with FL, i.e., we apply standard PGD [20], TRADES [36], MART [30]), and GAIRAT [37] to FL,
 216 and term them as FedPGD, FedTRADES, FedMART, and FedGAIRAT. Additionally, we apply the
 217 state-of-the-art long-tail learning methods (LogitAdj [22] and RoBal [32]) to FAT, and term them as
 218 FedLogitAdj and FedRoBal.

219 **Evaluation Metrics.** we report the natural test accuracy (Natural) and robust test accuracy under
 220 the most representative attacks, i.e., FGSM [31], BIM [14], PGD-20 [20], CW [2], and AA [4]. More
 221 detailed experimental setup is provided in Appendix D.1.

222 4.1 Main Results

223 **Evaluation on different datasets.** Table 1 and Table 2 show the results of
 224 all methods on CIFAR10, CIFAR100, SVHN, and ImageNet subset, respectively. From these two tables, we can
 225 observe that:
 226
 227
 228

229 (1) Our CalFAT achieves the best robustness on all datasets, validating the efficiency of our CalFAT. For example, CalFAT outperforms the best baseline method (FedGAIRAT) by 10.20% on SVHN dataset under FGSM
 230
 231
 232
 233

Table 2: Natural and robust accuracy (%) on ImageNet subset.

Metric	Natural	FGSM	BIM	CW	PGD-20	AA
MixFAT	33.45	19.67	18.35	16.32	18.32	11.93
FedPGD	30.89	18.93	17.94	16.12	18.42	11.34
FedTRADES	30.25	18.79	18.06	16.08	18.04	11.78
FedMART	26.46	16.56	15.46	14.21	15.43	9.45
FedGAIRAT	34.20	19.72	19.30	16.82	19.20	11.83
FedRBN	29.49	18.55	17.35	15.12	17.97	11.47
CalFAT(ours)	49.93	22.45	20.03	17.23	20.01	12.36

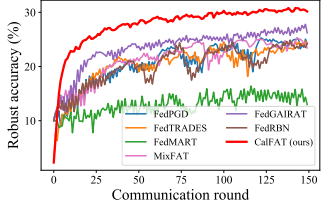


Figure 2: Robust accuracy (against PGD-20 attack) of different methods on CIFAR10 dataset.

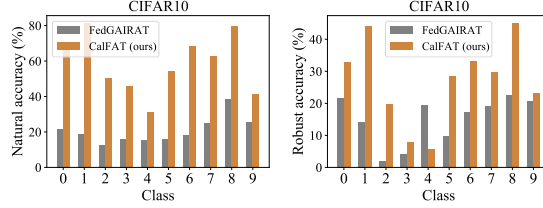


Figure 3: Per-class natural accuracy and robust accuracy (against PGD-20 attack) of CalFAT and the best baseline (FedGAIAT) on CIFAR10 dataset.

234 (2) Our CalFAT shows a significant improvement in natural accuracy compared to other baselines.
 235 For example, CalFAT can improve the natural accuracy of the best baseline method (FedGAIAT) by
 236 25.63% on SVHN dataset. We hypothesise that the reason lies in the homogeneity of local models in
 237 our CalFAT, which leads to better convergence and higher performance.

238 (3) All methods demonstrate the worst performance on CIFAR100 and ImageNet subset datasets. We
 239 conjecture that this is due to the more classes in these two datasets, which makes the training hard.
 240 Nevertheless, our CalFAT still achieves the best performance.

241 (4) FedMART has (almost) the worst performance across all datasets, which indicates that MART is
 242 not suitable to be directly applied to cross-silo FL.

243 **Learning curves of different methods.** To visually compare our CalFAT with all the baselines,
 244 we plot the learning curves (i.e, performance across different communication rounds) of all methods
 245 in Figure 1 and Figure 2. As shown in these figures, CalFAT achieves the best natural accuracy
 246 and robust accuracy from the beginning to the end of the training compared to the baseline FAT
 247 methods, which indicates that the design of our CalFAT is profitable across the whole training process.
 248 Moreover, our CalFAT is stable during the whole training process while the accuracy curves of other
 249 baselines oscillate strongly. Such oscillations lead to bad convergence and low performance. We
 250 hypothesise that the heterogeneity of local models in the baseline methods makes the training unstable
 251 while CalFAT learns homogeneous local models.

252 **Comparison with state-of-the-art long-tail learning methods**

253 We also adapt the losses of long-tail learning methods (LogitAdj [22] and RoBal [32]) to FAT
 254 (namely FedLogitAdj and FedRoBal) for both the local inner and outer optimization. As shown in Table 3,
 255 both methods have lower natural and robust accuracy than our CalFAT. We argue that this is because the proposed
 256 CKL loss can better generate adversarial examples, which further benefits the whole learning process.
 257
 258
 259
 260

Table 3: Comparison of long-tail learning methods.

Methods	FedLogitAdj [3]	FedRoBal [1]	CalFAT (ours)
Natural	59.79	61.48	64.85
PGD-20	28.84	29.51	31.19

261 **4.2 Performance on Different Classes**

262 We further compare the per-class performance (performance on each class) of CalFAT with the best
 263 baseline (FedGAIAT). First, we use a well-trained model to initialize a global model. Second, the
 264 global model distributes the model parameter to all clients. Third, the local clients train their local
 265 models with their local data for 1 epoch. Then, we report the per-class average performance of all
 266 clients for each class. For fair comparisons, we use the same well-trained model for initialization and
 267 the same data partition on each client for CalFAT and FedGAIAT.

268 In Figure 3, we report the per-class natural and robust accuracy of CalFAT and the best baseline
 269 (FedGAIAT) on CIFAR10 dataset. As shown in these figures, the average performance of most
 270 classes of CalFAT is much higher than FedGAIAT. We also report the per-class performance of each
 271 client on CIFAR10 dataset in Appendix D.2. In FedGAIAT, due to highly skewed label distribution,
 272 the prediction of each client is highly biased to majority classes, which leads to high performance

Table 4: Natural and robust accuracy (%) across different FL frameworks on CIFAR10 dataset.

FL framework	FedAvg			FedProx			Scaffold		
	Natural	PGD-20	AA	Natural	PGD-20	AA	Natural	PGD-20	AA
MixFAT	53.23	26.22	21.81	54.76	27.86	21.65	54.73	27.75	21.94
FedPGD	47.21	26.50	22.79	49.11	27.79	21.77	49.32	28.10	21.93
FedTRADES	46.14	26.29	21.60	47.53	27.85	21.86	47.91	28.46	21.85
FedMART	25.68	18.15	14.31	27.67	19.26	15.89	27.31	20.20	15.92
FedGAIRAT	48.34	27.32	21.92	49.43	28.35	21.81	49.72	28.63	21.83
FedRBN	47.87	26.21	21.43	49.04	27.54	21.89	50.00	28.02	21.88
CalFAT(ours)	64.85	31.19	22.94	66.46	32.64	22.26	66.67	33.24	21.96

Table 5: Natural and robust accuracy (%) across different numbers of clients $m = \{20, 50, 100\}$ on CIFAR10 dataset.

m	20			50			100		
	Natural	PGD-20	AA	Natural	PGD-20	AA	Natural	PGD-20	AA
MixFAT	25.34	18.33	14.23	22.78	14.89	11.23	21.79	15.12	11.23
FedPGD	29.48	18.51	14.62	29.73	17.84	12.82	27.92	15.25	11.26
FedTRADES	30.23	18.12	13.92	21.53	14.63	11.21	24.15	14.24	11.14
FedMART	22.68	17.81	13.28	18.25	14.44	11.23	22.29	14.89	11.21
FedGAIRAT	22.45	17.46	13.64	20.36	14.39	11.35	21.57	14.21	10.89
FedRBN	21.78	17.02	13.21	19.56	13.07	11.42	19.01	13.86	10.72
CalFAT(ours)	60.06	24.45	15.35	49.26	18.65	13.12	41.23	17.16	11.94

273 on majority classes and low performance (even 0% accuracy) on minority classes. By contrast, in
 274 CalFAT, each client has higher performance on most classes. These results show that the calibrated
 275 cross-entropy loss can indeed improve the performance on minority classes, and further improve
 276 the overall performance of the model. Moreover, we report the per-class average performance on
 277 SVHN dataset in Appendix D.3. Our CalFAT still outperforms the best baseline across most classes
 278 on SVHN dataset.

279 4.3 Results on Different FL Frameworks and Network Architectures

280 **Evaluation on different FL frameworks.** We conduct experiments on different FL frameworks,
 281 i.e., FedAvg [21], FedProx [15], and Scaffold [10]. The results for all methods on different FL
 282 frameworks are shown in Table 4. It can be observed that CalFAT has better natural accuracy and
 283 robust accuracy than all baselines on all FL frameworks, which indicates that CalFAT can be well
 284 combined with most FL frameworks.

285 **Evaluation on different network architectures.** We also compare CalFAT with baselines on
 286 different network architectures, i.e., CNN [21], VGG-8 [25], and ResNet-18 [7]. For CNN, we use
 287 the same architecture as [21]. VGG-8 and ResNet-18 are two widely used architectures in deep
 288 learning. The results on CIFAR10 dataset are shown in Appendix D.4. CalFAT outperforms all
 289 baselines, which further validates the superiority of CalFAT under different network architectures.

290 4.4 Feature Visualization

291 To better illustrate the efficacy of CalFAT, we visualize the learned features extracted from the second
 292 last layer of FedTRADES (the best baseline) and CalFAT trained on SVHN dataset in Appendix D.5.
 293 The extracted features are visualized in a 2-dimensional space by t-SNE [27]. The visualization
 294 results show that data from different classes are mixed together in FedTRADES, resulting in low
 295 performance. For instance, Class 6 (pink) and Class 8 (khaki) are hard to be separated in FedTRADES
 296 while these 2 classes can be well separated in CalFAT. This illustration further verifies that the server
 297 cannot learn a good global model from the heterogeneous local models. By contrast, CalFAT can
 298 well separate different classes and achieve better performance.

Table 6: Natural and robust accuracy (%) under different label skewness level β on CIFAR10 dataset.

Label skewness level	$\beta = 0.05$						$\beta = 0.2$						$\beta = 0.3$					
	Natural	FGSM	BIM	CW	PGD-20	AA	Natural	FGSM	BIM	CW	PGD-20	AA	Natural	FGSM	BIM	CW	PGD-20	AA
MixFAT	49.10	27.49	25.32	22.17	25.24	22.51	54.85	31.27	28.70	26.08	28.46	25.21	58.93	31.68	28.17	24.96	28.00	24.34
FedPGD	47.13	26.63	24.96	20.75	25.03	21.28	52.22	30.31	28.64	25.49	28.59	24.92	56.12	30.86	28.46	25.07	28.29	23.64
FedTRADES	40.24	26.02	25.06	22.48	24.99	20.16	48.52	29.94	28.73	25.57	28.65	24.15	54.26	30.83	29.39	24.74	29.26	23.87
FedMART	29.84	21.90	21.39	18.31	21.41	17.89	38.38	27.59	27.05	23.31	26.99	21.89	40.96	28.32	27.88	23.12	27.80	22.16
FedGAIRAT	50.41	28.89	26.30	22.66	26.34	23.81	56.11	32.99	29.90	27.10	28.97	25.97	60.63	33.31	30.12	25.50	29.67	24.75
FedRBN	39.35	25.92	24.40	21.55	24.77	19.47	48.42	29.59	27.74	24.67	27.86	23.78	53.54	29.88	28.76	24.11	28.63	23.14
CalFAT(ours)	61.00	32.40	29.75	23.55	29.50	25.66	71.55	33.80	30.70	27.25	29.35	26.32	69.95	34.25	30.80	27.76	30.96	26.84

299 4.5 Ablation Studies

300 **Impact of the number of clients.** To show the capability of CalFAT, we train CalFAT with different
 301 number of clients m . Table 5 reports all the results across $m = \{20, 50, 100\}$. As expected, CalFAT
 302 achieves the best performance across all m . As m increases, the performance of all methods decreases.
 303 We conjecture that the reason is that more clients in FAT make the model harder to converge. However,
 304 our CalFAT can still achieve 41.23% natural accuracy when there are 100 clients, which outperforms
 305 other baselines by a large margin.

306 **Impact of skewed label distribution.** We find that the performance of these FAT methods is closely
 307 related to label skewness. We investigate the impact of skewed label distribution by varying the
 308 Dirichlet parameter $\beta = \{0.05, 0.2, 0.3\}$ and report the results on CIFAR10 dataset in Table 6. Not
 309 surprisingly, our CalFAT outperforms all baselines under all β 's. This further verifies the consistent
 310 effectiveness of CalFAT under label skewness.

311 Note that as β decreases (i.e., the labels on each client are more imbalanced), the performance of
 312 all methods drops rapidly. For example, the natural accuracy of FedMART drops from 38.38% to
 313 29.84% as β decreases from 0.2 to 0.05. This indicates that all methods are hard to train a good
 314 model in extremely skewed label distribution scenarios. However, our CalFAT still achieves 61.00%
 315 natural accuracy and 32.40% robust accuracy (against FGSM attack) when $\beta = 0.05$, which are
 316 much higher than all the baselines.

317 **Contribution of the calibrated loss functions.** As shown in Eq. (11) and Eq. (13), for each client
 318 i , we have two new loss functions: a CCE loss $\ell_{cce}(\cdot, \cdot, \cdot)$ for optimization and a CKL loss $\ell_{ckl}(\cdot, \cdot, \cdot)$
 319 to generate adversarial data. This naturally raises a question: how do these two loss functions
 320 contribute to CalFAT? To answer this question, we conduct leave-one-out tests by removing the CCE
 321 loss (w/o $\ell_{cce}(\cdot, \cdot, \cdot)$) and removing the CKL loss (w/o $\ell_{ckl}(\cdot, \cdot, \cdot)$). As illustrated in Appendix D.6,
 322 w/o $\ell_{cce}(\cdot, \cdot, \cdot)$ leads to poor performance, which implies that CCE loss plays an important role in
 323 enhancing CalFAT. Besides, if we use only the CCE loss (i.e., w/o $\ell_{ckl}(\cdot, \cdot, \cdot)$), we can achieve a much
 324 better performance, but still underperforms CalFAT. All these results indicate that the CCE loss is the
 325 most important part in CalFAT, and the CKL loss can further increase the performance of CalFAT.
 326 The combination of both loss functions leads to the best performance.

327 **Impact of the ratio of adversarial data.** To investigate the impact of adversarial data, we conduct
 328 experiments with different ratios of adversarial data on CalFAT. Appendix D.7 shows the robust
 329 accuracy (against PGD-20 attack) of CalFAT with different ratios of adversarial data. Ratio $r=0$ and
 330 $r=1$ stand for training the model with only natural data and with only adversarial data, respectively. As
 331 expected, $r=1$ achieves the best performance, which further verifies that training with only adversarial
 332 data can better enhance adversarial robustness in CalFAT.

333 5 Conclusion

334 In this paper, we studied the challenging problem of Federated Adversarial Training (FAT) with
 335 label skewness and proposed a novel Calibrated Federated Adversarial Training (CalFAT) to achieve
 336 stable training, better convergence, and natural accuracy and robustness in FL. CalFAT calibrates
 337 the model prediction and trains homogeneous local models across different clients by giving higher
 338 scores to minority classes, thus delivering a better global model in FAT. Extensive experiments across
 339 representative vision datasets under various settings validate the effectiveness of our proposed method.
 340 We envision our work as a milestone for more accurate and robust federated adversarial training.

341 References

- 342 [1] Peter J Bickel and Kjell A Doksum. *Mathematical statistics: basic ideas and selected topics,*
343 *volumes I-II package*. CRC Press, 2015.
- 344 [2] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In
345 *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57. IEEE, 2017.
- 346 [3] Yair Carmon, Aditi Raghunathan, Ludwig Schmidt, Percy Liang, and John C Duchi. Unlabeled
347 data improves adversarial robustness. In *Proceedings of the 33rd International Conference on*
348 *Neural Information Processing Systems*, pages 11192–11203, 2019.
- 349 [4] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an
350 ensemble of diverse parameter-free attacks. In *International conference on machine learning*,
351 pages 2206–2216, 2020.
- 352 [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-
353 scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern*
354 *recognition*, pages 248–255, 2009.
- 355 [6] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
356 networks. In *International Conference on Machine Learning*, pages 1321–1330. PMLR, 2017.
- 357 [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
358 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
359 pages 770–778, 2016.
- 360 [8] Junyuan Hong, Haotao Wang, Zhangyang Wang, and Jiayu Zhou. Federated robustness propa-
361 gation: Sharing adversarial robustness in federated learning. *arXiv preprint arXiv:2106.10196*,
362 2021.
- 363 [9] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-iid data quagmire
364 of decentralized machine learning. In *International Conference on Machine Learning*, pages
365 4387–4398. PMLR, 2020.
- 366 [10] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
367 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
368 *International Conference on Machine Learning*, pages 5132–5143, 2020.
- 369 [11] Achim Klenke. *Probability theory: a comprehensive course*. Springer Science & Business
370 Media, 2013.
- 371 [12] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
372 In *Technical report*, 2009.
- 373 [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep
374 convolutional neural networks. *Advances in neural information processing systems*, 25:1097–
375 1105, 2012.
- 376 [14] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale.
377 *arXiv preprint arXiv:1611.01236*, 2016.
- 378 [15] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia
379 Smith. Federated optimization in heterogeneous networks. *arXiv preprint arXiv:1812.06127*,
380 2018.
- 381 [16] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia
382 Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning*
383 *and Systems*, 2:429–450, 2020.
- 384 [17] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor
385 learning: Training clean models on poisoned data. *Advances in Neural Information Processing*
386 *Systems*, 34, 2021.

- 387 [18] Jiahuan Luo, Xueyang Wu, Yun Luo, Anbu Huang, Yunfeng Huang, Yang Liu, and Qiang Yang.
388 Real-world image datasets for federated learning. *arXiv preprint arXiv:1910.11089*, 2019.
- 389 [19] Lingjuan Lyu, Han Yu, Jun Zhao, and Qiang Yang. Threats to federated learning. In *Federated*
390 *Learning*, pages 3–16. Springer, 2020.
- 391 [20] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
392 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,
393 2017.
- 394 [21] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
395 Communication-efficient learning of deep networks from decentralized data. In *Artificial*
396 *Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- 397 [22] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit,
398 and Sanjiv Kumar. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*,
399 2020.
- 400 [23] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.
401 Reading digits in natural images with unsupervised feature learning. In *NeurIPS Workshop on*
402 *Deep Learning and Unsupervised Feature Learning*, 2011.
- 403 [24] Leslie Rice, Eric Wong, and J Zico Kolter. Overfitting in adversarially robust deep learning. In
404 *ICML*, 2020.
- 405 [25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
406 image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- 407 [26] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander
408 Madry. Robustness may be at odds with accuracy. In *International Conference on Learning*
409 *Representations*, 2018.
- 410 [27] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine*
411 *learning research*, 9(11), 2008.
- 412 [28] Abraham Wald. Note on the consistency of the maximum likelihood estimate. *The Annals of*
413 *Mathematical Statistics*, 20(4):595–601, 1949.
- 414 [29] Wentao Wang, Han Xu, Xiaorui Liu, Yaxin Li, Bhavani Thuraisingham, and Jiliang Tang.
415 Imbalanced adversarial training with reweighting. *arXiv preprint arXiv:2107.13639*, 2021.
- 416 [30] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving
417 adversarial robustness requires revisiting misclassified examples. In *ICLR*, 2020.
- 418 [31] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial
419 training. In *ICLR*, 2020.
- 420 [32] Tong Wu, Ziwei Liu, Qingqiu Huang, Yu Wang, and Dahua Lin. Adversarial robustness under
421 long-tailed distribution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*
422 *Pattern Recognition*, pages 8659–8668, 2021.
- 423 [33] Han Xu, Xiaorui Liu, Yaxin Li, Anil Jain, and Jiliang Tang. To be robust or to be fair: Towards
424 fairness in adversarial training. In *International Conference on Machine Learning*, pages
425 11492–11501, 2021.
- 426 [34] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept
427 and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(2):1–19,
428 2019.
- 429 [35] Mikhail Yurochkin, Mayank Agarwal, Soumya Ghosh, Kristjan Greenewald, Nghia Hoang,
430 and Yasaman Khazaeni. Bayesian nonparametric federated learning of neural networks. In
431 *International Conference on Machine Learning*, pages 7252–7261, 2019.

- 432 [36] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan.
 433 Theoretically principled trade-off between robustness and accuracy. In *International Conference*
 434 *on Machine Learning*, pages 7472–7482. PMLR, 2019.
- 435 [37] Jingfeng Zhang, Jianing Zhu, Gang Niu, Bo Han, Masashi Sugiyama, and Mohan Kankanhalli.
 436 Geometry-aware instance-reweighted adversarial training. In *ICLR*, 2021.
- 437 [38] Yao Zhou, Jun Wu, and Jingrui He. Adversarially robust federated learning for neural networks,
 438 2021.
- 439 [39] Giulio Zizzo, Amrbrish Rawat, Mathieu Sinn, and Beat Buesser. Fat: Federated adversarial
 440 training. *arXiv preprint arXiv:2012.01791*, 2020.

441 Checklist

- 442 1. For all authors...
- 443 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 444 contributions and scope? [Yes]
- 445 (b) Did you describe the limitations of your work? [Yes]
- 446 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 447 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 448 them? [Yes]
- 449 2. If you are including theoretical results...
- 450 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 451 (b) Did you include complete proofs of all theoretical results? [Yes]
- 452 3. If you ran experiments...
- 453 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 454 mental results (either in the supplemental material or as a URL)? [Yes]
- 455 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 456 were chosen)? [Yes]
- 457 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 458 ments multiple times)? [No]
- 459 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 460 of GPUs, internal cluster, or cloud provider)? [Yes]
- 461 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 462 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 463 (b) Did you mention the license of the assets? [N/A]
- 464 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 465 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 466 using/curating? [N/A]
- 467 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 468 information or offensive content? [N/A]
- 469 5. If you used crowdsourcing or conducted research with human subjects...
- 470 (a) Did you include the full text of instructions given to participants and screenshots, if
 471 applicable? [N/A]
- 472 (b) Did you describe any potential participant risks, with links to Institutional Review
 473 Board (IRB) approvals, if applicable? [N/A]
- 474 (c) Did you include the estimated hourly wage paid to participants and the total amount
 475 spent on participant compensation? [N/A]