

Generative Data Augmentation for Non-IID Problem in Decentralized Clinical Machine Learning

Zirui Wang*, Shaoming Duan*, Chengyue Wu*, Wenhao Lin*, Xinyu Zha*, Peiyi Han*[†], Chuanyi Liu*^{†‡}

* School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen), Shenzhen, China

[†] Cyberspace Security Research Center, Peng Cheng Laboratory, Shenzhen, China

[‡]Corresponding authors

Email:854102781@qq.com, shaomingduan@gmail.com, 1190201314@stu.hit.edu.cn,

190110219@stu.hit.edu.cn, 200110710@stu.hit.edu.cn, hanpei@hit.edu.cn, liuchuanyi@hit.edu.cn

Abstract—Swarm learning (SL) is an emerging promising decentralized machine learning paradigm and has achieved high performance in clinical applications. SL solves the problem of a central structure in federated learning by combining edge computing and blockchain-based peer-to-peer network. While there are promising results in the assumption of the independent and identically distributed (IID) data across participants, SL suffers from performance degradation as the degree of the non-IID data increases. To address this problem, we propose a generative augmentation framework in swarm learning called SL-GAN, which augments the non-IID data by generating the synthetic data from participants. SL-GAN trains generators and discriminators locally, and periodically aggregation via a randomly elected coordinator in SL network. Under the standard assumptions, we theoretically prove the convergence of SL-GAN using stochastic approximations. Experimental results demonstrate that SL-GAN outperforms state-of-art methods on three real world clinical datasets including Tuberculosis, Leukemia, COVID-19.

Index Terms—Swarm learning, privacy-preserving decentralized machine learning, non-IID, data augmentation, generative adversarial network.

I. INTRODUCTION

Machine learning (ML) models are data hungry. In precision medicine [1], the performance of ML models that identify patients with life-threatening diseases, such as leukemia, tuberculosis or COVID-19, increases with the size and diversity of the training samples [2]. In practice, to train a robust clinical ML model, patient-related data often needs to be centralized in a central repository [3], [4]. However, such data sharing across different institutions or countries, faces privacy and legal obstacles. This problem has been solved by federated learning [5], in which multiple participants jointly train a ML model under a central coordinator. In federated learning, each participant trains a local model using the local data separately, and then shares the learned model gradients to the coordinator for model aggregation. Although, the private data is distributed and not disclosing to others, the central structure in federated learning remains vulnerable to attack [6].

To tackle the limitations of federated learning, swarm learning (SL) [7] combines edge computing and blockchain-based peer-to-peer network. In SL, a new participant register via a blockchain smart contract, obtains the global model, and

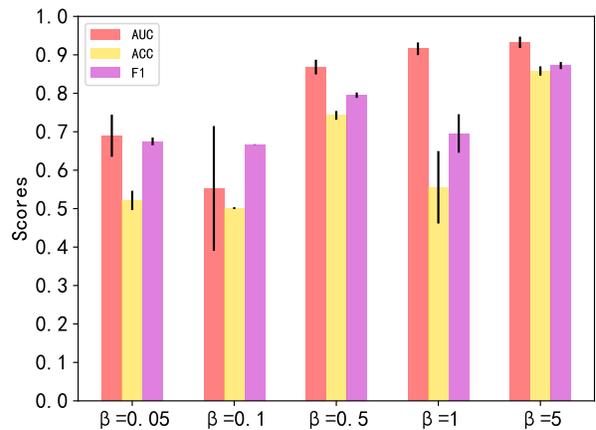


Fig. 1. Performance of swarm learning on the Tuberculosis dataset. β controls the degree of the non-IID, the lower the β , the more imbalanced the data distribution is. With the β increases, the performance of SL increases.

trains the model locally. After a user-defined synchronization interval, local model parameters are exchanged and merged to update the global model by a randomly elected edge node. The chosen edge node replaces the role of central coordinator in federated learning. In the context of clinical data mining, SL provides a fairness environment for multi-parties ML model training and creates a strong incentive to collaborate without data sharing. As SL secures data sovereignty, security, and confidentiality, it has been successfully applied in healthcare fields [7], [8], [9]. However, the non-independent and identical distributed (non-IID) data heavily limit the performance of SL. As shown in Figure 1, with the degree of non-IID increases (β from 5 to 0.05), the AUC scores of SL decreases from 93% to 69%.

Existing methods for solving the non-IID problem in decentralized learning can be roughly divided into two categories: *algorithm-based methods* and *data-based methods*. Algorithm-based methods improve the model robustness by modifying the local loss function [10], [11] to make the local model consistent with the global model, designing a new aggregation

scheme [12], [13] to improve the model aggregation mechanism, or training personalized models [14], [15] for each participant rather than the same global model. However, as discussed in [16], existing algorithm-based methods are not always better than vanilla FedAvg [5]. Data-based methods construct a more balanced data distribution among participants or on the server by data sharing or augmentation strategies [17], [18], which achieve a high performance on non-IID data. Unfortunately, there are still two challenges in directly applying these methods to SL. On the one hand, these methods need a trusted central coordinator to employ data sharing or data augmentation, which is contrary to the assumption in SL. On the other hand, data augmentation strategies based on generative adversarial network (GAN) remains a challenge is that GAN may not converge on non-IID data. To the best of our knowledge, there are currently no works to solve the non-IID problem in SL.

To address these challenges, we propose a novel generative augmentation framework in SL called SL-GAN, which augments the non-IID data to a balanced data distribution among each participant by using a generative model. In SL-GAN, discriminators and generators are trained locally and aggregated after a user-defined synchronization interval by a randomly elected edge node. Furthermore, under the standard assumption in SL, we theoretically prove the convergence of SL-GAN using stochastic approximations. We evaluate our SL-GAN on three real-world clinical datasets with various data distributions among participants. The experimental results show that SL-GAN can effectively improve the performance of SL on non-IID data and outperforms the state-of-the-art methods.

The main contributions of this paper are as follows:

- (1) To the best of our knowledge, this is the first study on the non-IID problem in SL. We propose SL-GAN, a novel data augmentation framework in SL, which jointly trains a global generative model to augment the non-IID data without a central coordinator.
- (2) We theoretically prove our SL-GAN converges with non-IID data under the standard assumption in SL.
- (3) We test SL-GAN on three real-world clinical datasets with various data distributions. The experimental results demonstrate that SL-GAN outperforms the state-of-the-art approaches and is robust to varies data distributions.

II. RELATED WORKS

A. Non-IID in Decentralized Learning

Since there are currently no works to solve the non-IID problem in swarm learning, we describe the related works in other decentralized learning method, such as federated learning. Existing methods in federated learning to deal with the non-IID problem are divided into two categories [19]: *algorithm-based methods* and *data-based methods*.

1) *Algorithm-based methods*: As suggested in [20], the weight divergence of local models caused by non-IID data is the root cause of model performance degradation. To alleviate

this problem, Fedprox[10] adds a penalty term in the objective function to make the local model consistent with the global model. SCAFFOLD[21] controls the similarity between the local model and the global model by adding a regularization term to the local loss function. To avoid the influence of the local model from nodes with large data volumes, FedNova[11] normalizes the local model before model aggregation. Unlike modifying the objective, personalized federated learning [14], [15] aims to train personalized models in each participant rather than the same global model. However, as shown in [16], existing algorithm-based methods are not always better than vanilla FedAvg [5].

2) *Data-based methods*: To construct a balanced data distribution among each participant, data-based methods perform data augmentation or data sharing strategies under the coordination of a central server. Mixup[17] is a simple and commonly used data augmentation method, which constructs a new samples by linear interpolation, and FedMix[22] is another work that using Mixup strategies. However, Mixup cannot generate unseen labels. Unlike data augmentation, data sharing methods [23], [24] alleviates the non-IID problem by collecting a small subset of samples from participants on the central server. However, these methods violate the privacy assumption in federated learning. Unlike these methods, our SL-GAN trains a global GAN for data augmentation.

B. Federated Generative Models

To synthesize fake data with a distribution similar to the global data, federated GAN [18], [25] train a global generative model among participants in federated learning. Existing federated GANs can be divided into two categories. One type is that generator trained on the server, discriminators are trained on the client. DP-FedAvg-GAN [25] is the first work in such architecture to train a global GAN for data generation, which under the assumption of IID data. F2U [26] follows previous settings and assign different weights for local discriminators in the process of model aggregation on the non-IID data. Another type [18] is that both discriminator and generator are trained on the client and aggregated on the server. Unfortunately, these methods require a central server and cannot guarantee the convergence of the GAN on non-IID data. In contrast to these methods, our SL-GAN trains a global GAN without a central coordinator. And we theoretically prove the convergence of SL-GAN.

III. PROPOSED METHODS

In this section, we propose SL-GAN framework, describe its training algorithm and theoretically prove its convergence. Table I summarizes notations used in this paper.

A. SL-GAN

Figure 2 shows the architecture of our SL-GAN, in which jointly trains a GAN among participants. In SL-GAN, each participant trains generator and discriminator locally. When a user-defined synchronization interval reached, the trained local discriminators and generators are aggregated on a randomly

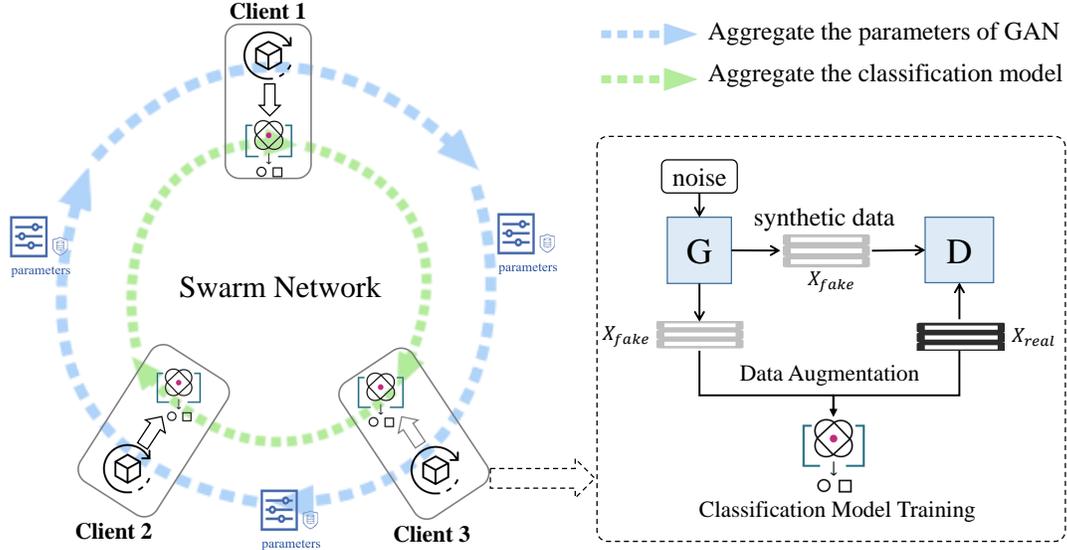


Fig. 2. SL-GAN architecture. The green dotted arrow circle represents a swarm network used for aggregating the target clinical machine learning model. The blue dotted arrow circle denotes the swarm network that used to aggregate the parameters of GAN. Discriminators and generators are trained locally and aggregated on a elected edge node. Synthetic data generated by the trained generator are used to balance the data distribution among participants.

TABLE I
SUMMARY OF NOTATIONS

Notation	Description
C	set of participants
X	local dataset of participants
lr	learning rate
N	local training epochs
T	aggregation time interval
g_D, g_G	stochastic gradient of discriminator and generator
\tilde{g}_D, \tilde{g}_G	true gradient of discriminator and generator
θ_D, θ_G	parameters of discriminator and generator
p	aggregation weight
$M^{\theta_D}, M^{\theta_G}$	stochastic gradient errors
v	discriminator parameters in the interval
ϕ	generator parameters in the interval
α	maximum round between two aggregations
sup	supremum
\mathbb{E}	expectation
μ	mean value
σ	standard deviation

elected participant and then sent back to each participant with the swarm network. After SL-GAN converges, the trained generator is used for data augmentation. Finally, The target classification model is trained using the combination of the synthetic data and the local private data.

Algorithm 1 describes the training process of SL-GAN.

- (1) Each participant trains local generator and discriminator using standard GAN training procedure. The discriminator is trained using the real data X_{real} and the fake data X_{fake} generated by the generator, the generator is

- updated follow the discriminator (lines 1-6 of Algorithm 1).
- (2) When reached the pre-defined synchronization interval T , participants send their local discriminator and generator to a temporarily elected node for aggregation with the weight p_c , $p_c = \frac{|\mathcal{X}_c|}{\sum_{j \in C} |\mathcal{X}_j|}$ (lines 7-10 of Algorithm 1).
- (3) After model aggregation, each participant receives the aggregated discriminator and generator and updates their local model (line 11 of Algorithm 1).

The algorithm repeats the above process until SL-GAN converges.

B. Convergence Analysis

In this section, we show that our SL-GAN converges in swarm learning with non-IID data. We denote the gradient of discriminator and generator in participant i by g_D^i and g_G^i . Let $\theta^i = (\theta_D^i, \theta_G^i)^\top$ be the parameter of the participant i . The true gradient rather than stochastic gradient of each client is specified as \hat{g}_D^i and \hat{g}_G^i . In addition, we define the stochastic gradient errors $M^{(g_D)}$ and $M^{(g_G)}$, where $M^{(\theta_D)} = \hat{g}_D - \sum_i p_i g_D^i$ and $M^{(\theta_G)} = \hat{g}_G - \sum_i p_i g_G^i$.

We follow the assumptions in the centralized GAN.

- 1) g_D^i and g_G^i are L -Lipschitz.
- 2) $\sum_n lr = \infty$, $\sum_n lr^2 < \infty$
- 3) $\{M_n^{(\theta_D)}\}$ and $\{M_n^{(\theta_G)}\}$ are martingale difference sequence of the increasing σ -filed $\mathbb{F}_n = \sigma(\theta_{D_l}, \theta_{G_l}, M_l^{(\theta_D)}, M_l^{(\theta_G)}, l \leq n), n \geq 0$.
- 4) $sup_n \|\theta_{D_n}\| < \infty$ and $sup_n \|\theta_{G_n}\| < \infty$

Algorithm 1 SL-GAN model training process

Input: Local training epoch N , batch size B , learning rate of discriminator $lr_D(n)$, learning rate of generator $lr_G(n)$, local discriminator θ_{D_c} , local generator θ_{G_c} , synchronization interval T , start time t_0 , current time t , weight of model aggregation p_c .

Output: well trained discriminator θ_{D_c} and generator θ_{G_c} .

```
1: for  $n$  from 0 to  $N - 1$  for all clients do
2:    $X_{real}^i \leftarrow$  (sample random batch data of batch size  $B$ )
3:    $X_{noise}^i \leftarrow$  (sample random noise of batch size  $B$ )
4:    $X_{fake}^i \leftarrow$  Generator( $X_{noise}^i, \theta_G^i$ )
5:    $\theta_D^i \leftarrow \theta_D^i - lr_D(n) \nabla_{\theta_D^i} loss_D(\theta_D^i, X_{fake}^i, X_{real}^i)$ 
6:    $\theta_G^i \leftarrow \theta_G^i - lr_G(n) \nabla_{\theta_G^i} loss_G(\theta_G^i, X_{fake}^i, \theta_D^i)$ 
7:   if  $(t - t_0) | T$  then
8:     Random select a participant  $c'$  for model aggregation
9:      $\theta_G^t \leftarrow \sum_{c \in C} p_c \theta_{G_c}^t$ 
10:     $\theta_D^t \leftarrow \sum_{c \in C} p_c \theta_{D_c}^t$ 
11:    Send back  $\theta_G^t, \theta_D^t$  and all participants update local discriminators and generators
12:   end if
13: end for
```

$$5) \mathbb{E} \|\hat{g}_D^i - g_D^i\| \leq \sigma_{g_D}, \mathbb{E} \|\hat{g}_G^i - g_G^i\| \leq \sigma_{g_G} \text{ and } \|\hat{g}_D^i - g_D^i\| \leq \mu_{g_D}$$

where (1)-(4) are used in stochastic approximation of GAN convergence. In assumption (5), the first bound ensures that the local stochastic gradient is close to the local true gradient, the second bound ensures that the local discriminator true gradient of the non-IID data are close to the discriminator true gradient of the pooled data [18], the last bound represents bounded gradient divergence.

To prove the convergence of SL-GAN, we connects the convergence of GAN to the convergence of an ordinary differential equations (ODE) representation of the parameter updates [27]. We prove that the ODE representing the parameter updates of SL-GAN asymptotically tracks the ODE representing the parameter updates of the centralized GAN. As defined in [27], [28], the centralized GAN tracks the following ODE asymptotically.

$$\dot{\theta} = \begin{pmatrix} \dot{\theta}_D(t) \\ \dot{\theta}_G(t) \end{pmatrix} = \begin{pmatrix} g_D(\theta_D(t), \theta_G(t)) \\ g_G(\theta_D(t), \theta_G(t)) \end{pmatrix}. \quad (1)$$

Therefore, the problem of proving the convergence of SL-GAN is transformed into proving that the parameter of SL-GAN follows (1) asymptotically. As proofed in [18], we have,

$$\mathbb{E} \|\theta_{D_n}^i - v_n\| + \mathbb{E} \|\theta_G^i - \phi_n\| \leq \frac{\sigma_{\theta_D} + \mu_{\theta_D} + \sigma_{\theta_G}}{2L} [(1 + 2lr(n-1)L)^{n \bmod K} - 1] \quad (2)$$

$$\mathbb{E} \|\theta_{D_n} - v_n\| + \mathbb{E} \|\theta_G - \phi_n\| \leq \frac{(\sigma_{\theta_D} + \mu_{\theta_D} + \sigma_{\theta_G})}{2L} [(1 + 2lr(n-1)L)^K - 1] - lr(n-1)\mu_{\theta_D}K \quad (3)$$

v_n and ϕ_n represent the parameter of discriminator and generator in the interval between two aggregations, respectively. The specific definition is as follows

$$v_n = \theta_{D_{n_1}} + \sum_{k=n_1}^n lr(k)g_D(\phi_k, v_k), \quad (4)$$

$$\phi_n = \theta_{G_{n_1}} + \sum_{k=n_1}^n lr(k)g_G(\phi_k, v_k) \quad (5)$$

where n_1 means the nearest aggregation timestamp. However, in swarm learning, there is no K , which indicates every K local epochs. Fortunately, we can still prove that in the interval between every two adjacent aggregations, the local epochs for every client is bounded,

$$n^i - n_1 \leq \alpha, \forall i \in \{1, \dots, m\} \quad \alpha > 0$$

which is trivial because a participant will not train itself infinitely anyway. We just need to treat $\max(n^i - n_1)$ as K , bringing it to the Equations (2) and (3), we will get the same result as federated learning. Based on the Theorem 1 in [18] and Theorem 2 in [29], θ in SL-GAN tracks the ODE (1) asymptotically, namely it will converge eventually.

IV. EXPERIMENTS

A. Experimental Setup

1) *Datasets:* We use three real-world clinical datasets, as shown in Table II, and their details are as follows:

- Tuberculosis [7] is an RNA-Seq dataset based on whole blood transcriptomes, which combines data from healthy controls with published data in Gene Expression Omnibus (GEO). This dataset merged from 9 independent datasets: GSE101705, GSE107104, GSE112087, GSE128078, GSE66573, GSE79362, GSE84076, GSE89403. There are 1550 samples from patients with active tuberculosis, latent tuberculosis, fatigue, autoimmune diseases, HIV and controls. All active tuberculosis samples are listed as CASE and all other samples are listed as CONTROL. After data preprocessing, there are 18136 genes (columns) in this dataset.
- Leukemia dataset [30] contains 2379 transcriptomes derived from peripheral blood mononuclear cells (PBMC) or bone marrow, published at the GEO under subseries GSE122505. In this dataset, independent data sets selected from GEO are as follows: GSE10255, GSE1159, GSE12417, GSE12995, GSE13425, GSE14471, GSE14895, GSE16129, GSE25571, GSE26281, GSE33315, GSE34860, GSE37642, GSE43176, GSE4698, GSE51082,

TABLE II
STATISTICS OF DATASETS

Dataset	# of samples (training set)	# of samples (test set)	# of columns	Distribution of samples (training set) (CASE:CONTROL)	Distribution of samples (test set) (CASE:CONTROL)
Tuberculosis	1240	310	18136	620:620	155:155
Leukemia	1943	436	22283	826:1117	206:230
COVID-19	1920	480	19400	237:1683	59:421

GSE6269, GSE67684, GSE83449, GSE8879, GSE9006, GSE9476. Acute Myeloid Leukemia (AML) samples are classified as CASE and all other samples are classified as CONTROL.

- COVID-19 [7] is an RNA-Seq dataset based on whole blood transcriptomes, which contains 296 samples from patients with COVID-19, as well as 2104 other control samples (autoimmune disease, Fatigue, healthy controls, HIV, latent tuberculosis and active tuberculosis). COVID-19 samples are labeled as CASE and all other samples are labeled as CONTROL. After data preprocessing, there are 19400 genes (columns) in this dataset.

Table II shows the statistical information of three datasets. We split the total samples into training dataset and test dataset with the ratio 8:2.

2) *Baseline Methods*: We adapted two state-of-the-art algorithm-based methods (Fedprox, FedNova) in federated learning and one data-based method (Mixup) to swarm learning as baseline methods. The detail information of baseline methods are as follows:

- **Fedprox** [10] is an algorithm-based method, which modify the objective function of participant k as follows,

$$h_k(w; w^t) = F_k(w) + \frac{\mu}{2} \|w - w^t\|^2 \quad (6)$$

where $F_k(w)$ represents the original objective function, w^t is the parameter of the global model obtained in the t -th round.

- **FedNova** [11] is an algorithm-based method, which modify the original aggregation method to

$$w^{t+1} - w^t = \left(\frac{\sum_{k=1}^K p_k \tau_k^{(t)}}{\tau_k^{(t)}} \right) \sum_{k=1}^K p_k \Delta_k^{(t)} \quad (7)$$

where $\tau_k^{(t)}$ is the number of iterations of participant k in round t , and $\Delta_k^{(t)}$ is the gradient of participant k in round t .

- **Mixup** [17] is a simple data augmentation method,

$$\begin{cases} \tilde{x} = \lambda x_i + (1 - \lambda) x_j \\ \tilde{y} = \lambda y_i + (1 - \lambda) y_j \end{cases} \quad (8)$$

where (x_i, y_i) and (x_j, y_j) are two randomly selected samples, $\lambda \in (0, 1)$.

3) *Data Partition*: To simulate the non-IID data in real world, we use the distributed-based label imbalance method in [16], that each participant is allocated a proportion of the samples of each label according to Dirichlet distribution. In practice, we allocate a $p_{k,j}$ ($p_{k,j} \sim Dir(\beta)$) proportion of the samples of label k to participant j , where $Dir(\cdot)$ represents Dirichlet distribution and β ($\beta > 0$) is its parameter. In this approach, we can flexibly control the degree of the non-IID by varying the parameter β . The smaller the β is, the more unbalanced the data distribution is.

4) *Data Augmentation*: After SL-GAN converges, the synthetic data generated by the trained generator is used to augment local data based on global data distribution. During training, each participant contains the same amount of data and its label distribution is consistent with the global label distribution. For instance, the data contained in a participant is {CONTROL:40, CASE:414} on the Tuberculosis dataset with $\beta = 1$, and the distribution of global data is CONTROL:CASE = 1:1. After data augmentation, the data of the three participants are {CONTROL:556, CASE:556}.

B. Performance on Classification

Figures 3 - 5 show the comparisons of the performance on Tuberculosis, COVID-19, Leukemia dataset, respectively. We compare vanilla swarm learning, Fedprox, FedNova, Mixup with our SL-GAN in four different data distributions ($\beta = 0.05, 0.1, 0.5, 1$). SL-GAN outperforms all the baseline methods in terms of F1 score, accuracy, and AUC in three datasets.

In many cases, the vanilla SL algorithm outperforms Fedprox and FedNova, which is consistent with the conclusion in [16]. This is because algorithm-based methods do not fundamentally address the problem of data imbalance. In contrast, Mixup and SL-GAN both perform better than the vanilla SL algorithm in almost all cases, which means that data augmentation methods have a significant effect on alleviating data imbalance. Unfortunately, in Figure 3 (a) and Figure 4 (a), F1 scores of Fedprox and FedNova are close to 0, which suggests that these methods fail to make the local model consistent with the global model under extreme data imbalances. Compared with Mixup, SL-GAN shows more reliable performance in various distributions. This is because Mixup cannot generate the class of samples that are not available in the local. Therefore, in some cases, the enhancement is redundant data features. SL-GAN trains a generative model among participants and learns the real distribution, which can generate

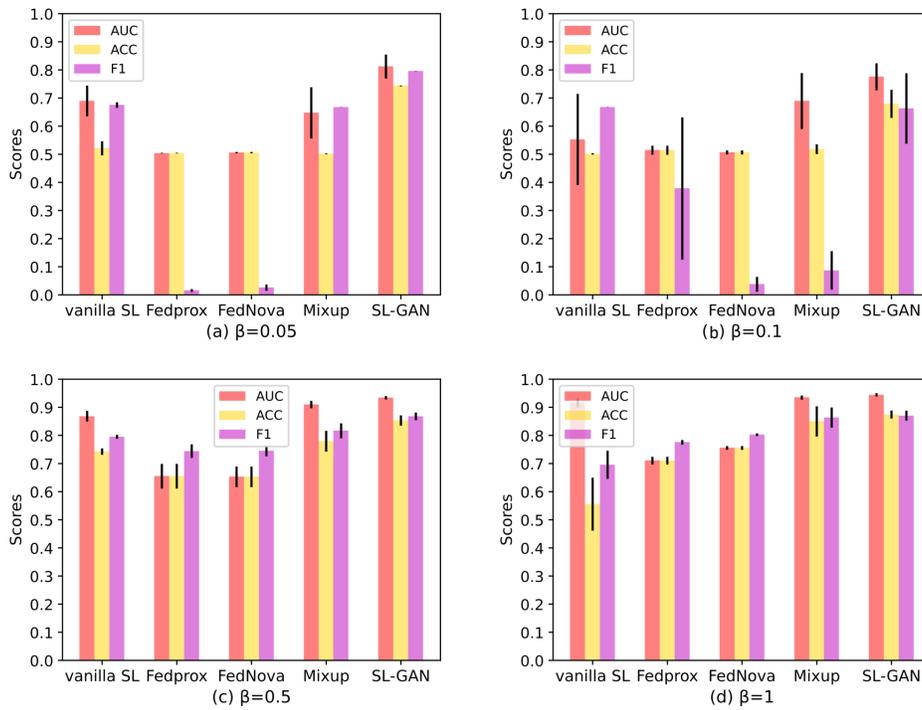


Fig. 3. Performance on the Tuberculosis dataset. We compared with vanilla swarm learning, Fedprox, FedNova, Mixup and SL-GAN in terms of AUC, F1 and accuracy scores. The red bar represents AUC scores, the yellow bar denotes accuracy, and the purple bar represents F1 scores.

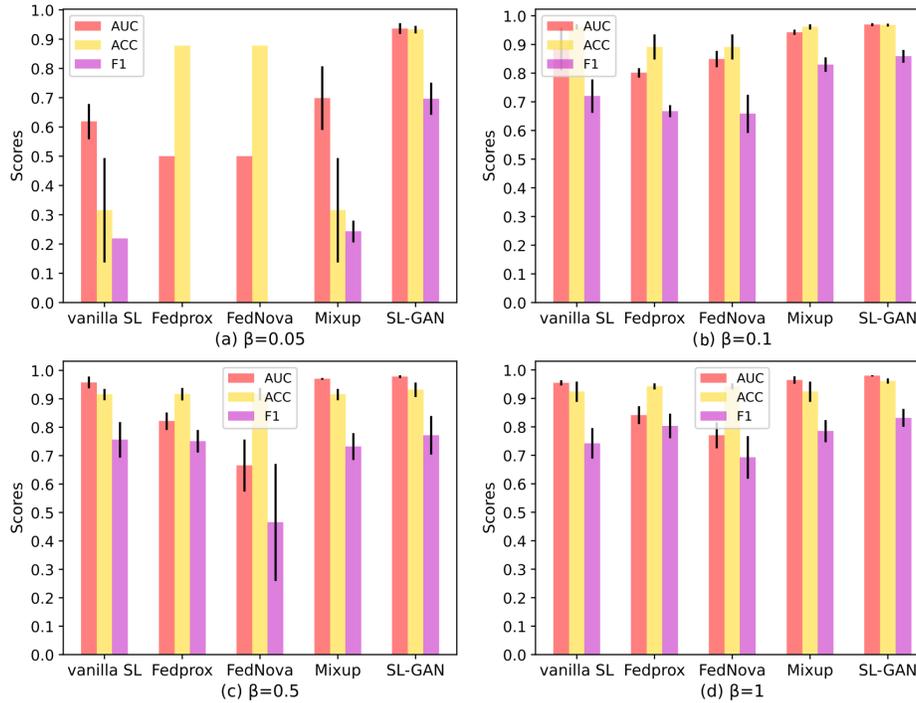


Fig. 4. Performance on the COVID-19 dataset. We compared with vanilla swarm learning, Fedprox, FedNova, Mixup and SL-GAN in terms of AUC, F1 and accuracy scores. The red bar represents AUC scores, the yellow bar denotes accuracy, and the purple bar represents F1 scores.

synthetic data with an approximate distribution of the global data. To summarize, in almost all cases, SL-GAN achieves

the best performance. In general, data-based methods show better performance than algorithm-based method because they

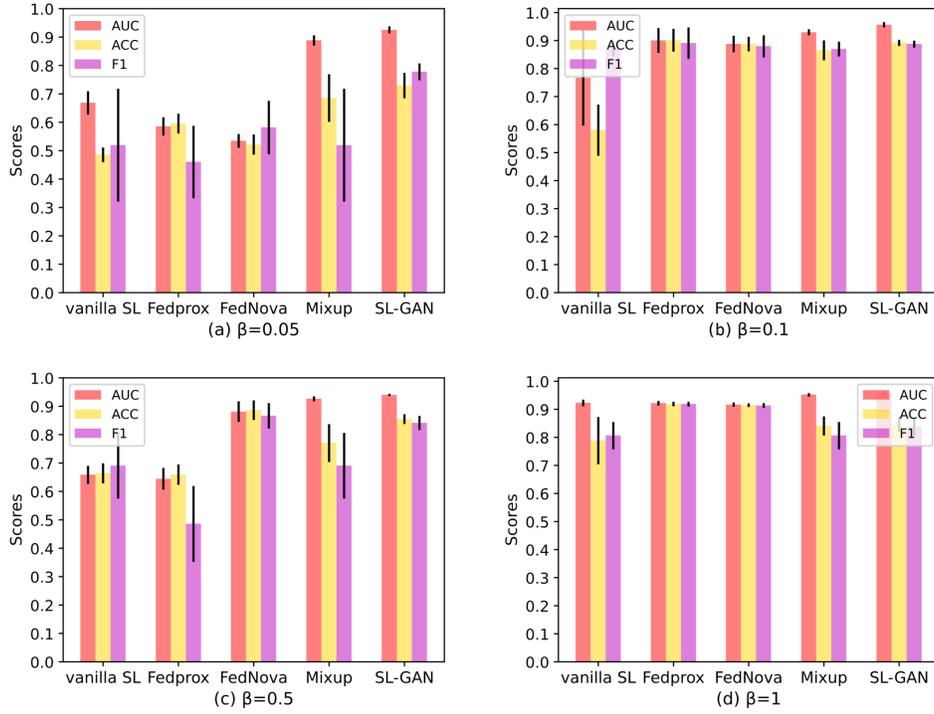


Fig. 5. Performance on the Leukemia dataset. We compared with vanilla swarm learning, Fedprox, FedNova, Mixup and SL-GAN in terms of AUC, F1 and accuracy scores. The red bar represents AUC score, the yellow bar denotes accuracy, and the purple bar represents F1 score.

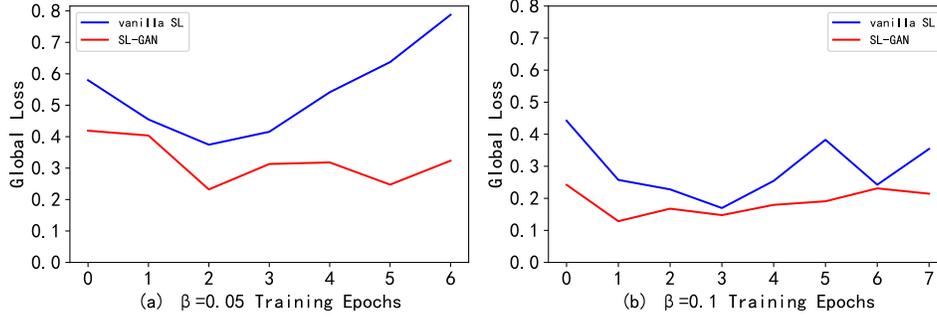


Fig. 6. Communication round on the COVID-19 dataset. The blue lines represent the training round of the vanilla SL algorithm, and the red lines represent the training round of SL-GAN

construct a balanced data distribution among each participant.

Figure 6 shows the training round on the COVID-19 dataset. As we can see, by augmenting the synthetic data, SL-GAN converges to a smaller loss than the vanilla SL algorithm. As shown in Figure 6, with the degree of non-IID increases (β form 0.1 to 0.05), the vanilla SL algorithm has difficulty in convergence. In contrast, the model augmented by SL-GAN can still converge to a relatively small loss, which means that the data augmentation method proposed in this paper can significantly improve the efficiency and effect of the model training.

C. Synthetic Data Utility

To evaluate the utility of the synthetic data that generated by SL-GAN, we train 10 machine learning models on the

synthetic data and test these models on the real data. The performance of the synthetic data for COVID-19 dataset shown in Table IV-C. As shown in Table IV-C, as the value of β decreases (from 1 to 0.05), the accuracy of the synthetic data does not change. This shows that the proposed SL-GAN can converge stably on the non-IID data. In many cases, the performance of the synthetic data is close to that of the original data. However, The average accuracy of the SL-GAN is 10% lower than that of the original data, which means that our generative method does not always successfully capture the real features of the original data.

V. CONCLUSION AND FUTURE WORKS

In this paper, we presented a generative augmentation framework in swarm learning for non-IID data, which called

TABLE III
THE ACCURACY OF THE SYNTHETIC DATA ON THE COVID-19 DATASET

Classifier	Original data	$\beta = 1$	$\beta = 0.5$	$\beta = 0.1$	$\beta = 0.05$
LGBMClassifier	0.98	0.74	0.72	0.83	0.58
XGBClassifier	0.97	0.70	0.69	0.72	0.60
BaggingClassifier	0.97	0.84	0.76	0.84	0.79
SVC	0.96	0.88	0.88	0.88	0.88
RandomForestClassifier	0.96	0.85	0.86	0.86	0.82
LabelPropagation	0.88	0.88	0.88	0.88	0.88
ExtraTreesClassifier	0.96	0.84	0.84	0.84	0.85
CalibratedClassifierCV	0.96	0.84	0.88	0.88	0.87
GaussianNB	0.83	0.86	0.88	0.80	0.86
LabelSpreading	0.88	0.88	0.88	0.88	0.88
LabelPropagation	0.88	0.88	0.88	0.88	0.88
Average	0.935	0.831	0.827	0.841	0.801

SL-GAN. We jointly train a GAN in the swarm learning network, and theoretically prove the convergence of the SL-GAN. We augments the non-IID data into IID data using SL-GAN. We evaluated the proposed SL-GAN on three real-world clinical dataset. The experimental results show that our SL-GAN outperforms the state-of-the-art work in various data distributions.

In the future, in order to protect the data privacy, the differential privacy will be introduced to the SL-GAN and the privacy of the synthetic data will be studied. Furthermore, as there is still a gap between the utility of the synthetic data and original data, we will combines prior knowledge to improve the quality of the synthetic data.

ACKNOWLEDGMENT

This work was supported by the National Natural Science Foundation of China Grant No.61872110.

REFERENCES

- [1] M. A. Haendel, C. G. Chute, and P. N. Robinson, "Classification, ontology, and precision medicine," *New England Journal of Medicine*, vol. 379, no. 15, pp. 1452–1462, 2018.
- [2] A. Echle, H. I. Grabsch, P. Quirke, P. A. van den Brandt, N. P. West, G. G. Hutchins, L. R. Heij, X. Tan, S. D. Richman, J. Krause *et al.*, "Clinical-grade histology signatures on deep learning model accuracy and bias," *Gastroenterology*, vol. 159, no. 4, pp. 1406–1416, 2020.
- [3] R. J. Chen, M. Y. Lu, T. Y. Chen, D. F. Williamson, and F. Mahmood, "Synthetic data in machine learning for medicine and healthcare," *Nature Biomedical Engineering*, vol. 5, no. 6, pp. 493–497, 2021.
- [4] F. M. Howard, J. Dolezal, S. Kochanny, J. Schulte, H. Chen, L. Heij, D. Huo, R. Nanda, O. I. Olopade, J. N. Kather *et al.*, "The impact of site-specific digital histology signatures on deep learning model accuracy and bias," *Nature communications*, vol. 12, no. 1, pp. 1–13, 2021.
- [5] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [6] L. Chen, S. Fu, L. Lin, Y. Luo, and W. Zhao, "Privacy-preserving swarm learning based on homomorphic encryption," in *International Conference on Algorithms and Architectures for Parallel Processing*. Springer, 2021, pp. 509–523.

- [7] S. Warnat-Herresthal, H. Schultze, K. L. Shastry, S. Manamohan, S. Mukherjee, V. Garg, R. Sarveswara, K. Händler, P. Pickkers, N. A. Aziz *et al.*, "Swarm learning for decentralized and confidential clinical machine learning," *Nature*, vol. 594, no. 7862, pp. 265–270, 2021.
- [8] O. L. Saldanha, P. Quirke, N. P. West, J. A. James, M. B. Loughrey, H. I. Grabsch, M. Salto-Tellez, E. Alwers, D. Cifci, N. Ghaffari Laleh *et al.*, "Swarm learning for decentralized artificial intelligence in cancer histopathology," *Nature Medicine*, pp. 1–8, 2022.
- [9] D. Fan, Y. Wu, and X. Li, "On the fairness of swarm learning in skin lesion classification," in *Clinical Image-Based Procedures, Distributed and Collaborative Learning, Artificial Intelligence for Combating COVID-19 and Secure and Privacy-Preserving Machine Learning*. Springer, 2021, pp. 120–129.
- [10] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, "Federated optimization in heterogeneous networks," *Proceedings of Machine Learning and Systems*, vol. 2, pp. 429–450, 2020.
- [11] J. Wang, Q. Liu, H. Liang, G. Joshi, and H. V. Poor, "Tackling the objective inconsistency problem in heterogeneous federated optimization," *Advances in neural information processing systems*, vol. 33, pp. 7611–7623, 2020.
- [12] M. Yurochkin, M. Agarwal, S. Ghosh, K. Greenewald, N. Hoang, and Y. Khazaeni, "Bayesian nonparametric federated learning of neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 7252–7261.
- [13] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," in *International Conference on Learning Representations*, 2019.
- [14] V. Smith, C.-K. Chiang, M. Sanjabi, and A. S. Talwalkar, "Federated multi-task learning," *Advances in neural information processing systems*, vol. 30, 2017.
- [15] T. Li, S. Hu, A. Beirami, and V. Smith, "Ditto: Fair and robust federated learning through personalization," in *International Conference on Machine Learning*. PMLR, 2021, pp. 6357–6368.
- [16] Q. Li, Y. Diao, Q. Chen, and B. He, "Federated learning on non-iid data silos: An experimental study," *arXiv preprint arXiv:2102.02079*, 2021.
- [17] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *International Conference on Learning Representations*, 2018.
- [18] M. Rasouli, T. Sun, and R. Rajagopal, "Fedgan: Federated generative adversarial networks for distributed data," *arXiv preprint arXiv:2006.07228*, 2020.
- [19] H. Zhu, J. Xu, S. Liu, and Y. Jin, "Federated learning on non-iid data: A survey," *Neurocomputing*, vol. 465, pp. 371–390, 2021.
- [20] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," 2018.
- [21] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International Conference on Machine Learning*. PMLR, 2020, pp. 5132–5143.
- [22] T. Yoon, S. Shin, S. J. Hwang, and E. Yang, "Fedmix: Approximation of mixup under mean augmented federated learning," in *The Ninth International Conference on Learning Representations (ICLR)*. International Conference on Learning Representations (ICLR), 2021.
- [23] T. Tuor, S. Wang, B. J. Ko, C. Liu, and K. K. Leung, "Overcoming noisy and irrelevant data in federated learning," in *2020 25th International Conference on Pattern Recognition (ICPR)*, 2021, pp. 5020–5027.
- [24] N. Yoshida, T. Nishio, M. Morikura, K. Yamamoto, and R. Yonetani, "Hybrid-fl for wireless networks: Cooperative learning mechanism using non-iid data," in *ICC 2020-2020 IEEE International Conference on Communications (ICC)*. IEEE, 2020, pp. 1–7.
- [25] S. Augenstein, H. B. McMahan, D. Ramage, S. Ramaswamy, P. Kairouz, M. Chen, R. Mathews, and B. A. y Arcas, "Generative models for effective ml on private, decentralized datasets," in *International Conference on Learning Representations*, 2019.
- [26] R. Yonetani, T. Takahashi, A. Hashimoto, and Y. Ushiku, "Decentralized learning of generative adversarial networks from multi-client non-iid data," 2019.
- [27] L. Mescheder, S. Nowozin, and A. Geiger, "The numerics of gans," *Advances in neural information processing systems*, vol. 30, 2017.
- [28] V. Nagarajan and J. Z. Kolter, "Gradient descent gan optimization is locally stable," *Advances in neural information processing systems*, vol. 30, 2017.
- [29] V. S. Borkar, *Stochastic approximation: a dynamical systems viewpoint*. Springer, 2009, vol. 48.

- [30] S. Warnat-Herresthal, K. Perrakis, B. Taschler, M. Becker, K. Baßler, M. Beyer, P. Günther, J. Schulte-Schrepping, L. Seep, K. Klee *et al.*, “Scalable prediction of acute myeloid leukemia using high-dimensional machine learning and blood transcriptomics,” *Science*, vol. 23, no. 1, p. 100780, 2020.