

FedTADBench: Federated Time-series Anomaly Detection Benchmark

Fanxing Liu¹, Cheng Zeng², Le Zhang², Yingjie Zhou¹, Qing Mu¹, Yanru Zhang⁴, Ling Zhang^{3,5}, and Ce Zhu²

¹College of Computer Science, Sichuan University, China

²School of Information and Communication Engineering,

University of Electronic Science and Technology of China, China

³College of Polymer Science and Engineering, Sichuan University, China

⁴School of Computer Science and Engineering, University of Electronic Science and Technology of China, China

⁵West China School of Public Health, Sichuan University, China

Abstract—Time series anomaly detection strives to uncover potential abnormal behaviors and patterns from temporal data, and has fundamental significance in diverse application scenarios. Constructing an effective detection model usually requires adequate training data stored in a centralized manner, however, this requirement sometimes could not be satisfied in realistic scenarios. As a prevailing approach to address the above problem, federated learning has demonstrated its power to cooperate with the distributed data available while protecting the privacy of data providers. However, it is still unclear that how existing time series anomaly detection algorithms perform with decentralized data storage and privacy protection through federated learning. To study this, we conduct a federated time series anomaly detection benchmark, named FedTADBench, which involves five representative time series anomaly detection algorithms and four popular federated learning methods. We would like to answer the following questions: (1) How is the performance of time series anomaly detection algorithms when meeting federated learning? (2) Which federated learning method is the most appropriate one for time series anomaly detection? (3) How do federated time series anomaly detection approaches perform on different partitions of data in clients? Numbers of results as well as corresponding analysis are provided from extensive experiments with various settings. The source code of our benchmark is publicly available at <https://github.com/fanxingliu2020/FedTADBench>.

Index Terms—time series anomaly detection, federated learning, performance evaluation, data partition

I. INTRODUCTION

Time series anomaly detection is a vital and fundamental task in data mining, and has been broadly employed in a variety of application scenarios, such as intelligent manufacturing [1], critical infrastructure monitoring [2], [3], seismic analysis [4], [5], financial fraud detection [6] and health care [7]. During the last decades, numerous anomaly detection models have been proposed for handling temporal data. Most of these models need to train with temporal data stored in a centralized manner. However, this condition sometimes may not be satisfied in practice due to the huge communication cost of uploading distributed data to a central server and private protection of sensitive information among individual

data resource. Thus, conducting time series anomaly detection with decentralized data storage and private protection is a significant and imperative problem.

Federated learning [8] is a learning paradigm that can conduct model training without direct access to raw data, but sharing model parameters/gradient by training clients' data locally while protecting the privacy of clients. As a prevailing approach to address the aforementioned problem, federated learning has been widely used in various fields, such as IoT [9] and medical information [10], to aggregate siloed data. Some recent works have explored time series anomaly detection under the framework of federated learning (as shown in Fig. 1), such as [11], [12] and [13], where the feasibility of using the federated learning framework for time series anomaly detection is preliminary demonstrated. However, existing efforts mainly consider solving the specific challenges of their concerned tasks. There still lack of unbiased and in-depth evaluations of time series anomaly detection with federated learning, which could provide guideline information to construct effective models for both researchers and engineers.

To fill this blank, we construct a federated time series anomaly detection benchmark, named FedTADBench, which designs numerous evaluations for the performance of time series anomaly detection with federated learning. Specifically, we conduct experiments with three major aspects in the benchmark: (1) the feasibility of federated learning for typical time series anomaly detection algorithms, i.e., performance comparisons of time series anomaly detection algorithms with/without federated learning, and time series anomaly detection performance under isolated and federated settings; (2) the compatibility between popular federated learning methods and typical time series anomaly detection algorithms, i.e., performance comparisons of different federated learning methods for time series anomaly detection, and training time consumption of typical time series anomaly detection algorithms with different federated learning strategies; (3) the influence of the heterogeneity between clients to federated time series anomaly detection, i.e., performance comparisons of federated time series anomaly detection approaches with

The first two authors contribute equally. Corresponding authors are Yingjie Zhou(yjzhou09@gmail.com) and Le Zhang(lezhang@uestc.edu.cn).

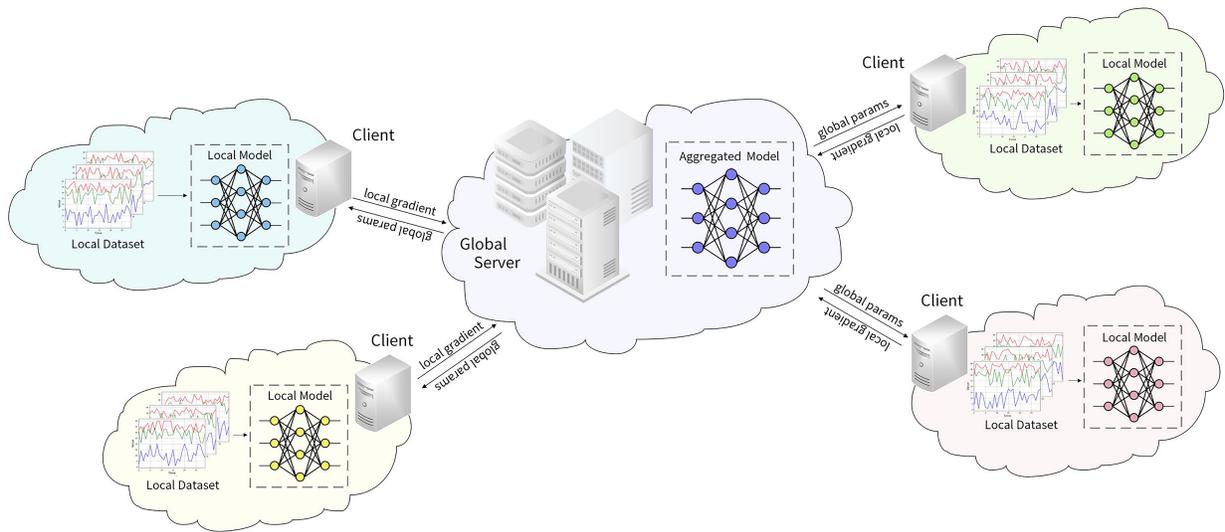


Fig. 1. Overview of time series anomaly detection under the framework of federated learning. Each client maintains its local model and trains it with its private dataset for several local epochs. Then, the local models will be sent to a centralized(global) server for calculating the updates of the global(aggregated) model. This process will be repeated for several “global epochs” until convergence.

different partitions of data in clients. Detailed analysis is made based on the extensive experiment results.

The contributions of this paper are as follows:

- 1) We construct FedTADBench, the first federated time series anomaly detection benchmark as far as we know, which evaluates time series anomaly detection algorithms with various federated learning settings.
- 2) Based on evaluations from three major angles, we show valuable insights for time series anomaly detection under the framework of federated learning, which could offer guiding information for building effective models of interest.
- 3) We provide reproducible evaluations whose codes are fully open-source. The website is <https://github.com/fanxingliu2020/FedTADBench>.

The rest of the paper is organized as follows. Section II introduces the related works on time series anomaly detection and federated learning, as well as the existing benchmarks. The problem statement and details of our benchmark are presented in Section III. Various experiment results and corresponding analysis are conducted in Section IV. Section V summarizes the paper.

II. RELATED WORKS

A. Time Series Anomaly Detection

The essential challenge to time series anomaly detection is formulating the pattern of temporal data characteristically, especially under certain limitations such as dearth of high-quality labeled data, heterogeneous structure of data, etc. As a result, most of existing anomaly detection models are unsupervised ones. SISVAE [14] constructs a recurrent neural network based variational auto-encoder to model normal data, which outperforms in capturing latent temporal structures of

time series. For lessening the impact of fitting rare anomalous samples, T. Kiew et al. [15] construct an ensemble of a series of recurrent auto-encoders with various skip connections. OmniAnomaly proposed by Y. Su et al. [16] performs time series modeling by combining a stochastic recurrent neural network and a planar normalizing flow, and uses the reconstruction probabilities to determine anomalies. This work transcends most of prior methods at the cost of high training time. More recently, aiming to address the large expense of energy consumption, J. Audibert et al. [17] propose USAD, which is capable of distinguishing anomalies from normal data rapidly with an unsupervised method by combining adversarial thought with auto-encoder.

Recently, with the increasing popularity of transformers and graph neural networks, some scholars have explored the use of them for time series anomaly detection [18], [19]. Graph neural networks are used to extract correlations between dimensions in multivariate time series, while Transformers are exploited for modeling latent representations in sequences. GDN introduced by A. Deng et al. [20] combines graph neural networks by attention mechanism, and provides rationales that the detected points are determined as abnormal. TranAD proposed by S. Tuli et al. [21] is an unsupervised time series anomaly detection model, which is able to compute in parallel by utilizing Transformers, and gain stable results benefiting from adversarial training. Anomaly Transformer proposed by J. Xu et al. [22] learns global relationships of time series with attention mechanism, and utilizes the difference between local relationships and global relationships to conduct anomaly detection.

Most of the proposed approaches are under the assumption that all data is accessible to the central server for training. However, it is not always feasible in real-world scenarios where the availability of data is restricted by security and

privacy concerns and there are also communication resource constraints.

B. Federated Learning

Federated learning can be categorized into vertical federated learning and horizontal federated learning [23]. Here, we mainly discuss horizontal federated learning as our data share the feature space but is held by different institutions. Currently there are four classic horizontal federated learning methods: FedAvg [8], FedProx [24], SCAFFOLD [25], and MOON [26]. FedAvg is the first federated learning algorithm, and the others add regularization to FedAvg from different perspectives. Since federated learning has advantages in privacy protection and data communication cost, some researchers have conducted time series anomaly detection under the framework of federated learning. Mothukuri et al. [27] adopt federated learning to construct an anomaly detection system for IoT security attacks, which trains local models on edge devices and aggregates the information on every edge device by federated learning. Liu et al. [28] proposed an attention based convolutional neural network with long short-term memory to detect edge device failure. To protect users' privacy, they adopted a federated learning framework to carry out the learning process. Besides, for improving communication efficiency, they employed a top-k selection algorithm to compress the gradient. For industrial cyber-physical systems (CPSs), it is difficult to collect sufficient high-quality attack examples. To aggregate data from different users in a privacy-preserving way, Li et al [29] designed a federated learning system based on paillier cryptosystem to train the deep learning based intrusion detection model for industrial CPSs.

However, the previously mentioned approaches mainly focus on improving the detection performance for specific tasks. The feasibility and limitations of federated learning for time series anomaly detection tasks is still unclear.

C. Existing Benchmarks

Researchers have made sustained efforts to provide unbiased evaluations for existing anomaly detection algorithms, such as [30], [31], [32], [33] and [34]. There have also been some reports to conduct performance comparisons of time series anomaly detection algorithms [35], [36], since temporal data is commonly seen in various application scenarios. However, any benchmark for federated time series anomaly detection has not been investigated yet. This paper tries to evaluate that how do existing time series anomaly detection algorithms perform on various federated learning settings. To the best of our knowledge, this is the first benchmark for federated time series anomaly detection.

III. FEDTADBENCH: AN EMPIRICAL EVALUATION ON FEDERATED TIME-SERIES ANOMALY DETECTION

We present the problem statement and elements of FedTADBench in this section. For better understanding, we provide an overview of our benchmark, as shown in Fig. 2, which sheds light upon time series anomaly detection in federated settings from 3 angles.

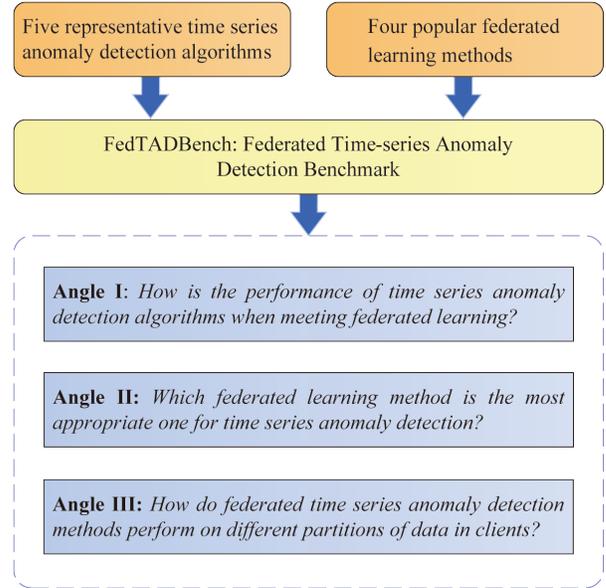


Fig. 2. An overview of FedTADBench that evaluates combinations of 5 time series anomaly detection and 4 federated learning methods. It sheds light upon time series anomaly detection in federated settings from 3 angles.

A. Problem Statement

In this study, we focus on detecting anomalies at entity-level [37] using multivariate time series. More specifically, we consider multi-dimensional time series data which consists of multiple observations continuously collected at equal-space timestamp. The multivariate time series can be defined as $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n]$, where n is the number of collected data, and $\mathbf{X}^i = [x_1^i, x_2^i, \dots, x_t^i]$ is an observation vector of the i^{th} metric within a dimensionality of t . For entity-level multivariate time series anomaly detection, our target is to determine whether the observation $\mathbf{X}_t = [x_t^1, x_t^2, \dots, x_t^n]$ at time t is anomalous or normal.

Conventional approaches learn such a function in a *Centralized* manner. In this way, they usually assume the learning algorithms have access to all the training samples. Contrastively, in many privacy-preserving scenarios, the data are usually stored in an isolated manner (namely the data can only be accessed by its owner) by different clients. And those clients are assumed to train a learning algorithm in a collaborative manner without data exchange. More specifically, in this case we assume that the data are owned by C clients, i.e., $\mathbf{X} = \{\mathbf{X}(1), \mathbf{X}(2), \dots, \mathbf{X}(C)\}$, and the task is to collaboratively train a model in a condition that each client only can access its own data. In a typical federated learning framework, each client of federated learning trains its *local model* with its own private dataset for several *local epochs*, then, the clients send their local models to the server. Finally, the server aggregates the local models into a *global model*, and this "training-aggregating" regime will be repeated for several "global epochs" until convergence.

B. Elements of FedTADBench

We evaluate 4 commonly-used federated learning methods and 5 time series anomaly detection algorithms on 3 popular time series anomaly detection datasets. In addition, we also investigate the effect of different partitions of data in clients. Below we elaborate the details.

1) Time Series Anomaly Detection Methods:

- DeepSVDD [38]: It jointly trains a deep neural network and learns a data-enclosing hypersphere in the latent space for anomaly detection tasks.
- LSTM-AE [39]. It is a combination of LSTM and Auto-encoders that modeling time series by reconstructing the original data.
- USAD [17]. It is an adversarially-trained autoencoders. USAD is fast to train and robust to the choice of parameters.
- GDN [20]. It trains a graph neural network in a structural learning manner. Attention weights are used as well to provide certain level of explainability for the detected anomalies.
- TranAD [21]. It is a deep transformer network and uses attention-based sequence encoders to enable robust multi-modal feature extraction.

2) Federated Learning Methods:

- FedAvg [8]. It averages the models learned by clients to obtain the global model.
- FedProx [24]. It adds a proximal term to regularize the parameters of the local model. In this case, the parameters of the current local model are not far from those of the previous global model.
- SCAFFOLD [25]. It uses variance reduction techniques to remedy the “client-drift” in its local updates. SCAFFOLD is shown to be more communicational efficient than FedAvg.
- MOON [26]. It regularizes the federated learning progress in a contrastive manner. More specifically, it pushes the models away from its previous status and pulls the models towards the global model.

3) Datasets:

- SMD [16]. SMD (Server Machine Dataset) is a 5-week-long dataset and is splitted into training and testing set with equal size.
- SMAP [40]. SMAP (Soil Moisture Active Passive satellite) is a public dataset from NASA.
- PSM [41]. It is a dataset proposed by eBay and consists of 26 dimensional time series data from application servers.

The details of datasets that we perform experiments on are shown in Table I. NS, ND and NC refer to the quantity of time series in each dataset, dimensions in the time series of each dataset, and clients in federated learning settings, respectively.

4) *Benchmark Angles*: Our benchmark is motivated by the following aspects. First, the compatibility between popular federated learning and time series anomaly detection methods is not extensively studied. Due to permutation invariance

of neural network parameters [42], client drift in federated learning [25] and so on, federated learning may not be compatible with certain time series anomaly detection algorithms. So our benchmark could provide such a testbed. Second, the *No-Free-Lunch* [43] Theorem states that within certain constraints, over the space of all possible problems, every optimization technique will perform as well as every other one on average. However, for federated learning based time series anomaly detection, this has not been explored empirically. Our benchmark can help us to verify this hypothesis as well. Third, in practice, the amount of data held by clients may vary significantly. This essentially brings in heterogeneity between the clients and usually degrades the performance. To study how the heterogeneity between clients affects the performance of time series anomaly detection, we test the effect of different levels of client heterogeneity on performance.

TABLE I
DATASET DETAILS AND CORRESPONDING FEDERATED LEARNING SETTINGS. NS, ND AND NC REFER TO THE QUANTITY OF TIME SERIES IN EACH DATASET, DIMENSIONS OF THE TIME SERIES IN EACH DATASET, AND CLIENTS IN FEDERATED LEARNING SETTINGS, RESPECTIVELY.

Dataset	NS	ND	NC
SMD	28	38	28
SMAP	54	25	54
PSM	1	25	24

IV. EXPERIMENTS

A. Experiment Settings

For all the datasets, we normalize the data by scaling the value range between 0 and 1 for each dimension respectively. It should be noticed that, for SMD and SMAP, we normalize all the time series uniformly. Inspired by parameter settings in [26], we set the proximal weight of FedProx to be 0.01, and set the hyper-parameter “temperature” and the weight of contrastive loss in MOON [26] to be 0.5 and 1, respectively. For the hyper-parameter settings in time series anomaly detection algorithms, we use the default values as recommended in the original papers if appropriate.

B. Evaluation Metrics

We choose AUC-ROC (Area under the ROC Curve) and AUC-PR (Area under the Precision-Recall Curve) [44] as evaluation metrics. AUC-ROC is a commonly-used metric for anomaly detection. AUC-PR is the area under the curve representing the correlation between Precision and Recall.

We also use Precision, Recall and F1-score as evaluation metrics. They are calculated as

$$Precision = \frac{TP}{TP + FP} \quad (1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2)$$

TABLE II

DETAILED RESULTS OF FEDTADBENCH IN TERMS OF AUC-ROC AND AUC-PR. WE HIGHLIGHT THE BEST RESULT AND THE SECOND BEST RESULT IN BLUE AND LIGHT BLUE, RESPECTIVELY.

Method		DeepSVDD		LSTM-AE		USAD		GDN		TranAD	
		AUC-ROC	AUC-PR								
SMD	Original	0.6311	0.0818	0.6231	0.0871	0.6874	0.1182	0.6664	0.0946	0.6028	0.1120
	FedAvg	0.6368	0.1291	0.5915	0.0764	0.6518	0.1197	0.6509	0.0785	0.6335	0.1311
	FedProx	0.6401	0.0858	0.5921	0.0767	0.5451	0.0542	0.6383	0.0747	0.5720	0.0582
	Scaffold	0.5773	0.0909	0.6185	0.0756	0.6625	0.1096	0.6302	0.0692	0.6211	0.1087
	Moon	0.6679	0.1612	0.5981	0.0748	0.6433	0.0951	0.6273	0.0777	0.6337	0.1310
SMAP	Original	0.6032	0.1549	0.4442	0.1155	0.5520	0.1312	0.5369	0.1352	0.5754	0.1423
	FedAvg	0.6333	0.1841	0.4328	0.1141	0.5786	0.1418	0.5159	0.1228	0.5392	0.1274
	FedProx	0.4299	0.1094	0.4523	0.1168	0.4579	0.1088	0.5218	0.1247	0.5361	0.1272
	Scaffold	0.5992	0.1529	0.4328	0.1141	0.5908	0.1451	0.5761	0.1361	0.5395	0.1275
	Moon	0.6154	0.1940	0.4675	0.1192	0.5756	0.1401	0.5047	0.1203	0.5392	0.1274
PSM	Original	0.7830	0.5423	0.6063	0.4306	0.6313	0.4715	0.7389	0.5286	0.5792	0.3566
	FedAvg	0.5826	0.3674	0.6068	0.4308	0.6693	0.4794	0.7486	0.4332	0.5804	0.3601
	FedProx	0.5895	0.3728	0.6073	0.4318	0.5682	0.3880	0.7138	0.4672	0.4336	0.2408
	Scaffold	0.7606	0.5044	0.7036	0.4737	0.6792	0.4524	0.7494	0.4357	0.5006	0.2787
	Moon	0.5590	0.3697	0.6111	0.4360	0.6567	0.4237	0.7465	0.4583	0.5833	0.3636

$$F1 = \frac{2 \times (Precision \times Recall)}{Precision + Recall} \quad (3)$$

where TP is the number of correctly detected anomalous time points, FP is the number of normal time points judged to be anomalous, and FN is the number of time points that are wrongly judged to be normal. Following the literature [17], we search the best threshold that has the highest F1-score for each experiment.

In the experiments, we also present the adjusted performance metrics, following [17], [21]. That is, if one or more anomalous points in an anomalous sub-sequence are judged as abnormal, each anomalous time point in this sub-sequence is considered to be successfully detected.

C. Time series anomaly detection performance with/without federated learning

To analyze the feasibility of federated learning for typical time series anomaly detection algorithms, we evaluate the 5 centralized time series anomaly detection methods and their combinations with 4 different federated learning methods on 3 datasets.

As shown in Table II, III and IV, for SMD and SMAP dataset, federated learning clients and time series are in a one-to-one correspondence. Note that, for Table III and IV, the performance metrics are the adjusted ones as described in Section IV. B. For PSM, time points are assigned consecutively to 24 clients following a Dirichlet distribution with β set to 0.5. We surprisingly observe that, in some conditions, federated learning performs better than centralized learning. We hypothesis that federated learning can bring some regularization effect for the training of time series anomaly detection.

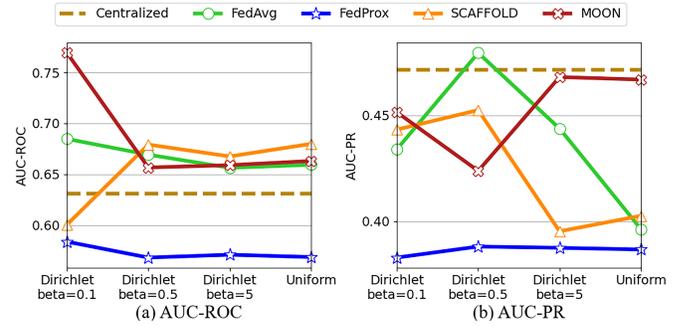


Fig. 3. AUC-ROC and AUC-PR of USAD on PSM with different partitions of data in clients.

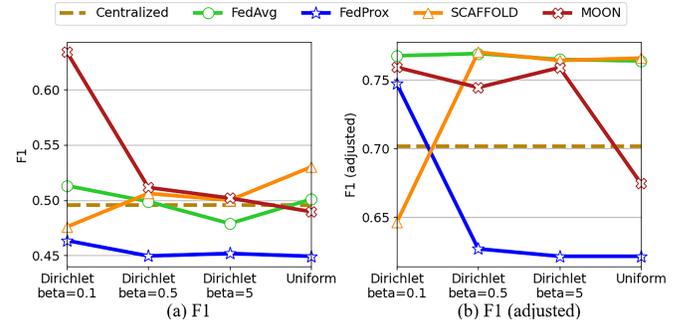


Fig. 4. F1 and F1 (adjusted) of USAD on PSM with different partitions of data in clients.

D. Time series anomaly detection performance with federated learning when data being differently partitioned in clients

To explore the influence of the heterogeneity between clients to federated time series anomaly detection, we evaluate federated time series anomaly detection with different partitions of data in clients.

Specifically, we employ 4 different partitions in the client-generating procedure. First of all, we employ an average dis-

TABLE III
DETAILED RESULTS OF FEDTADBENCH IN TERMS OF PRECISION AND RECALL. WE HIGHLIGHT THE BEST RESULT AND THE SECOND BEST RESULT IN BLUE AND LIGHT BLUE, RESPECTIVELY. NOTE THAT THE PRECISION AND RECALL ARE ADJUSTED METRICS AS DESCRIBED IN SECTION IV. B.

Method	DeepSVDD		LSTM-AE		USAD		GDN		TranAD		
	Precision	Recall									
SMD	Original	0.6921	0.2754	0.4102	0.4269	0.4495	0.4177	0.6024	0.7665	0.5915	0.3637
	FedAvg	0.6221	0.4820	0.5659	0.2478	0.5911	0.4149	0.3966	0.7414	0.6869	0.4920
	FedProx	0.6038	0.8152	0.5661	0.2478	0.3441	0.3134	0.3755	0.6013	0.3845	0.3950
	Scaffold	0.5245	0.1746	0.3268	0.5708	0.4820	0.4113	0.2966	0.6166	0.7423	0.3620
	Moon	0.8567	0.8060	0.5485	0.2674	0.5431	0.3128	0.4563	0.5745	0.6624	0.4926
SMAP	Original	0.7118	0.8625	0.6900	0.4311	0.9549	0.5631	0.9538	0.5613	0.8457	0.8096
	FedAvg	0.7947	0.7989	0.6902	0.4311	0.9574	0.5631	0.9796	0.5508	0.9018	0.5644
	FedProx	0.9254	0.5572	0.6907	0.4311	0.3078	0.9159	0.9784	0.5529	0.7813	0.5717
	Scaffold	0.7774	0.8848	0.6902	0.4311	0.9578	0.5616	0.9521	0.5539	0.9029	0.5694
	Moon	0.8998	0.8811	0.8160	0.3905	0.9555	0.5640	0.9531	0.5635	0.9018	0.5644
PSM	Original	0.8264	0.8819	0.6997	0.6922	0.5644	0.9269	0.8240	0.9107	0.7605	0.8543
	FedAvg	0.8749	0.8198	0.6795	0.6922	0.6720	0.8989	0.6781	0.9590	0.8214	0.8381
	FedProx	0.8476	0.8730	0.6796	0.6922	0.4758	0.9192	0.9332	0.9000	0.6360	0.6215
	Scaffold	0.9717	0.7128	0.7676	0.7562	0.6727	0.9003	0.6684	0.8852	0.2776	1.0000
	Moon	0.5789	0.9347	0.6583	0.6923	0.6335	0.9023	0.7890	0.8920	0.8005	0.8378

tribution. The numbers of timestamps in different clients are equal or extremely close. In addition, Dirichlet distributions with 3 different β values (0.1, 0.5 and 5) are employed, which is the same as the settings in [26]. We perform experiments on federated learning method FedAvg as an example. The AUC-ROC and AUC-PR results are shown in Fig. 3, and the F1 and F1 (adjusted) scores are shown in Fig. 4.

It could be observed that the performance of federated learning frameworks is robust to the change of β , i.e., the unbalanced data distribution does not affect the performance of time series anomaly detection much. In federated learning, the number of data essentially affects the updates of the model in one global epoch, and this result implies that, for time series anomaly detection, the number of updates does not cause much heterogeneity between local models.

E. Compatibility between different federated learning frameworks and time series

We are also interested in the compatibility between popular federated learning methods and typical time series anomaly detection algorithms.

As shown in Table II, III and IV, we can see that the performance of FedProx varies significantly under different conditions. FedProx builds a strong regularization on the change of local model parameters. We argue this may not be optimal because that, in some conditions, we may require certain parts of the model parameters to adapt to the data change efficiently. For example, USAD is a two-stage algorithm, we may need the parameters of stage one to change properly. However, FedProx limits the change of parameters in both stages, which may make FedProx to be the worst federated learning framework for USAD. On the other hand, MOON constrains the change of model in a more relaxed manner by constraining the change of features. Moreover, the

performance of FedAvg and SCAFFOLD is mostly in the middle level because they introduce fewer fixed assumptions.

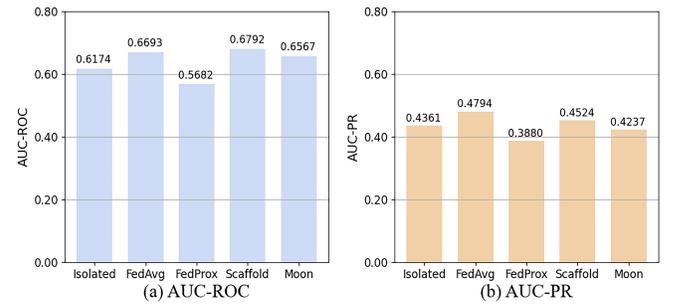


Fig. 5. AUC-ROC and AUC-PR of isolated and federated training for USAD on PSM.

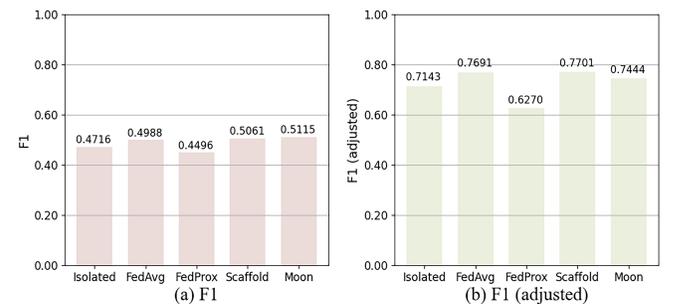


Fig. 6. F1 and F1 (adjusted) of isolated and federated training for USAD on PSM.

TABLE IV
DETAILED RESULTS OF FEDTADBENCH IN TERMS OF F1 AND F1 (ADJUSTED). WE HIGHLIGHT THE BEST RESULT AND THE SECOND BEST RESULT IN BLUE AND LIGHT BLUE, RESPECTIVELY.

Method	DeepSVDD		LSTM-AE		USAD		GDN		TranAD		
	F1	F1(adjusted)	F1	F1(adjusted)	F1	F1(adjusted)	F1	F1(adjusted)	F1	F1(adjusted)	
SMD	Original	0.1290	0.3941	0.1373	0.4184	0.1872	0.4330	0.1777	0.6747	0.1435	0.4504
	FedAvg	0.1792	0.5431	0.1413	0.3446	0.1813	0.4876	0.1490	0.5168	0.1829	0.5734
	FedProx	0.1628	0.6937	0.1424	0.3447	0.1095	0.3280	0.1430	0.4623	0.1100	0.3896
	Scaffold	0.1725	0.2620	0.1518	0.4156	0.1717	0.4439	0.1346	0.4005	0.1490	0.4867
	Moon	0.2188	0.8306	0.1408	0.3595	0.1793	0.3970	0.1316	0.5086	0.1833	0.5650
SMAP	Original	0.2804	0.7799	0.2342	0.5307	0.2805	0.7084	0.2584	0.7067	0.2917	0.8272
	FedAvg	0.3303	0.7968	0.2272	0.5307	0.2851	0.7091	0.2665	0.7050	0.2782	0.6943
	FedProx	0.2323	0.6956	0.2445	0.5309	0.2332	0.4608	0.2696	0.7065	0.2744	0.6602
	Scaffold	0.2922	0.8276	0.2272	0.4311	0.2825	0.7081	0.2892	0.7003	0.2784	0.6984
	Moon	0.2840	0.8904	0.2550	0.5282	0.2795	0.7093	0.2671	0.7082	0.2782	0.6943
PSM	Original	0.6260	0.8532	0.4414	0.6959	0.4957	0.7015	0.5924	0.8652	0.4451	0.8047
	FedAvg	0.4356	0.8464	0.4415	0.6858	0.4988	0.7691	0.5711	0.7945	0.4472	0.8296
	FedProx	0.4617	0.8601	0.4420	0.6858	0.4496	0.6270	0.5336	0.9163	0.4345	0.6287
	Scaffold	0.6178	0.8224	0.5500	0.7619	0.5061	0.7701	0.5851	0.7616	0.4345	0.6345
	Moon	0.4345	0.7150	0.4501	0.6749	0.5115	0.7444	0.5564	0.8374	0.4483	0.8187

F. Time series anomaly detection performance under isolated and federated settings

To demonstrate whether federated learning is meaningful, we compare the time series anomaly detection performance under federated learning and isolated training, respectively.

As shown in Fig. 5 and Fig. 6, we evaluate the performance of USAD on PSM under isolated and federated settings. It is obvious that in most instances, the performance with federated learning is better than that when training in isolation, which demonstrates the effectiveness of federated learning.

G. Training time without/with federated learning

To better guide the choice of federated learning methods for time series anomaly detection, we also evaluate the training time of typical time series anomaly detection algorithm with different federated learning strategies.

Training time of one global epoch of USAD on PSM in different federated learning strategies is shown in Table V. For this evaluation, the number of local epochs in each client for all federated learning methods is 10. All the experiments are conducted on NVIDIA Geforce RTX 3090 (24GB) GPU. Note that clients in federated and isolated settings are trained in parallel. It could be observed that the training time of FedAvg, whose performance is stable according to the aforementioned analysis, is shorter than that of the other three federated learning methods. MOON takes the longest time for training, due to the complex calculations involved by contrastive learning.

V. CONCLUSION

In this paper, we present the first benchmark for time series anomaly detection under federated learning frameworks. Our benchmark covers 5 time series anomaly detection algorithms, 4 federated learning frameworks, and 3 time series anomaly detection datasets. Based on our experiments, we analyze the effects federated learning brings in, the influence of unbalanced data distribution and the compatibility between different federated learning frameworks and time series. These

TABLE V
TRAINING TIME OF USAD ON PSM.

Learning Manner	Federated Learning	Training Time/s (one global epoch)
Centralized	N	20.99
Isolated	N	2.64
FedAvg	Y	27.51
FedProx	Y	37.29
Scaffold	Y	38.60
Moon	Y	64.88

analyses can provide some guiding information so that researchers can better choose time series anomaly detection algorithms and federated learning frameworks when performing time series anomaly detection under federated learning frameworks. Our work may help develop a federated learning framework that takes the characteristic of time series anomaly detection into account.

VI. ACKNOWLEDGMENT

This work was supported in part by National Natural Science Foundation of China (NSFC) with grant number 62171302 and U19A2052.

REFERENCES

- [1] X. Wang, J. Lin, N. Patel, and M. Braun, "A self-learning and online algorithm for time series anomaly detection, with application in cpu manufacturing," in *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, 2016, pp. 1823–1832.
- [2] H. Xu, W. Chen, N. Zhao, Z. Li, J. Bu, Z. Li, Y. Liu, Y. Zhao, D. Pei, Y. Feng *et al.*, "Unsupervised anomaly detection via variational auto-encoder for seasonal kpis in web applications," in *Proceedings of the 2018 world wide web conference*, 2018, pp. 187–196.
- [3] Z. Chen, D. Chen, X. Zhang, Z. Yuan, and X. Cheng, "Learning graph structures with transformer for multivariate time series anomaly detection in iot," *IEEE Internet of Things Journal*, 2021.
- [4] F. Qian, Y. Wang, B. Zheng, Z. Liu, Y. Zhou, and G. Hu, "Multidimensional seismic data denoising using framelet-based order-p tensor deep learning," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2022.

- [5] F. Qian, Z. Liu, Y. Wang, Y. Zhou, and G. Hu, "Ground truth-free 3-d seismic random noise attenuation via deep tensor convolutional neural networks in the time-frequency domain," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–17, 2022.
- [6] H. Weng, Z. Li, S. Ji, C. Chu, H. Lu, T. Du, and Q. He, "Online e-commerce fraud: a large-scale detection and analysis," *Proceedings of the IEEE International Conference on Data Engineering (ICDE)*, 2018, pp. 1435–1440.
- [7] B. Zhou, S. Liu, B. Hooi, X. Cheng, and J. Ye, "Beatgan: Anomalous rhythm detection using adversarially generated time series," in *IJCAI*, 2019, pp. 4433–4439.
- [8] 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.
- [9] D. C. Nguyen, M. Ding, P. N. Pathirana, A. Seneviratne, J. Li, and H. V. Poor, "Federated learning for internet of things: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 3, pp. 1622–1658, 2021.
- [10] J. Xu, B. S. Glicksberg, C. Su, P. Walker, J. Bian, and F. Wang, "Federated learning for healthcare informatics," *Journal of Healthcare Informatics Research*, vol. 5, no. 1, pp. 1–19, 2021.
- [11] Y. Liu, N. Kumar, Z. Xiong, W. Y. B. Lim, J. Kang, and D. Niyato, "Communication-efficient federated learning for anomaly detection in industrial internet of things," in *GLOBECOM 2020-2020 IEEE Global Communications Conference*. IEEE, 2020, pp. 1–6.
- [12] K. Zhang, Y. Jiang, L. Seversky, C. Xu, D. Liu, and H. Song, "Federated variational learning for anomaly detection in multivariate time series," in *2021 IEEE International Performance, Computing, and Communications Conference (IPCCC)*. IEEE, 2021, pp. 1–9.
- [13] H. T. Truong, B. P. Ta, Q. A. Le, D. M. Nguyen, C. T. Le, H. X. Nguyen, H. T. Do, H. T. Nguyen, and K. P. Tran, "Light-weight federated learning-based anomaly detection for time-series data in industrial control systems," *Computers in Industry*, vol. 140, p. 103692, 2022.
- [14] L. Li, J. Yan, H. Wang, and Y. Jin, "Anomaly detection of time series with smoothness-inducing sequential variational auto-encoder," *IEEE transactions on neural networks and learning systems*, vol. 32, no. 3, pp. 1177–1191, 2020.
- [15] T. Kieu, B. Yang, C. Guo, and C. S. Jensen, "Outlier detection for time series with recurrent autoencoder ensembles," in *IJCAI*, 2019, pp. 2725–2732.
- [16] Y. Su, Y. Zhao, C. Niu, R. Liu, W. Sun, and D. Pei, "Robust anomaly detection for multivariate time series through stochastic recurrent neural network," in *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 2828–2837.
- [17] J. Audibert, P. Michiardi, F. Guyard, S. Marti, and M. A. Zuluaga, "Usad: Unsupervised anomaly detection on multivariate time series," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 3395–3404.
- [18] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [20] A. Deng and B. Hooi, "Graph neural network-based anomaly detection in multivariate time series," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 5, 2021, pp. 4027–4035.
- [21] S. Tuli, G. Casale, and N. R. Jennings, "Tranad: Deep transformer networks for anomaly detection in multivariate time series data," *arXiv preprint arXiv:2201.07284*, 2022.
- [22] J. Xu, H. Wu, J. Wang, and M. Long, "Anomaly transformer: Time series anomaly detection with association discrepancy," *arXiv preprint arXiv:2110.02642*, 2021.
- [23] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 2, jan 2019. [Online]. Available: <https://doi.org/10.1145/3298981>
- [24] 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.
- [25] 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.
- [26] Q. Li, B. He, and D. Song, "Model-contrastive federated learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 10713–10722.
- [27] V. Mothukuri, P. Khare, R. M. Parizi, S. Pouriyeh, A. Dehghantanha, and G. Srivastava, "Federated-learning-based anomaly detection for iot security attacks," *IEEE Internet of Things Journal*, vol. 9, no. 4, pp. 2545–2554, 2021.
- [28] Y. Liu, S. Garg, J. Nie, Y. Zhang, Z. Xiong, J. Kang, and M. S. Hossain, "Deep anomaly detection for time-series data in industrial iot: A communication-efficient on-device federated learning approach," *IEEE Internet of Things Journal*, vol. 8, no. 8, pp. 6348–6358, 2020.
- [29] B. Li, Y. Wu, J. Song, R. Lu, T. Li, and L. Zhao, "Deepfed: Federated deep learning for intrusion detection in industrial cyber-physical systems," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 8, pp. 5615–5624, 2020.
- [30] G. O. Campos, A. Zimek, J. Sander, R. J. Campello, B. Micenkova, E. Schubert, I. Assent, and M. E. Houle, "On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study," *Data mining and knowledge discovery*, vol. 30, no. 4, pp. 891–927, 2016.
- [31] R. Domingues, M. Filippone, P. Michiardi, and J. Zouaoui, "A comparative evaluation of outlier detection algorithms: Experiments and analyses," *Pattern recognition*, vol. 74, pp. 406–421, 2018.
- [32] L. Ruff, J. R. Kauffmann, R. A. Vandermeulen, G. Montavon, W. Samek, M. Kloft, T. G. Dietterich, and K.-R. Müller, "A unifying review of deep and shallow anomaly detection," *Proceedings of the IEEE*, vol. 109, no. 5, pp. 756–795, 2021.
- [33] A. Acintoae, A. Florescu, M.-I. Georgescu, T. Mare, P. Sumedrea, R. T. Ionescu, F. S. Khan, and M. Shah, "Ubnormal: New benchmark for supervised open-set video anomaly detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 20143–20153.
- [34] S. Han, X. Hu, H. Huang, M. Jiang, and Y. Zhao, "Adbench: Anomaly detection benchmark," *arXiv preprint arXiv:2206.09426*, 2022.
- [35] K.-H. Lai, D. Zha, J. Xu, Y. Zhao, G. Wang, and X. Hu, "Revisiting time series outlier detection: Definitions and benchmarks," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [36] J. Papparizos, Y. Kang, P. Boniol, R. S. Tsay, T. Palpanas, and M. J. Franklin, "Tsb-uad: an end-to-end benchmark suite for univariate time-series anomaly detection," in *Proceedings of the VLDB Endowment*, 2022, pp. 1697–1711.
- [37] T. Huang, P. Chen, and R. Li, "A semi-supervised vae based active anomaly detection framework in multivariate time series for online systems," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 1797–1806.
- [38] L. Ruff, R. Vandermeulen, N. Goernitz, L. Deecke, S. A. Siddiqui, A. Binder, E. Müller, and M. Kloft, "Deep one-class classification," in *International conference on machine learning*. PMLR, 2018, pp. 4393–4402.
- [39] A. Garg, W. Zhang, J. Samaran, R. Savitha, and C.-S. Foo, "An evaluation of anomaly detection and diagnosis in multivariate time series," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 6, pp. 2508–2517, 2021.
- [40] K. Hundman, V. Constantinou, C. Laporte, I. Colwell, and T. Soderstrom, "Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding," in *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, 2018, pp. 387–395.
- [41] A. Abdulaal, Z. Liu, and T. Lancewicki, "Practical approach to asynchronous multivariate time series anomaly detection and localization," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 2485–2494.
- [42] H. Wang, M. Yurochkin, Y. Sun, D. Papailiopoulos, and Y. Khazaeni, "Federated learning with matched averaging," *arXiv preprint arXiv:2002.06440*, 2020.
- [43] S. Luke, *Essentials of metaheuristics*. Lulu Raleigh, 2013, vol. 2.
- [44] Y. Zhou, X. Song, Y. Zhang, F. Liu, C. Zhu, and L. Liu, "Feature encoding with autoencoders for weakly supervised anomaly detection," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.