# FedLED: Label-Free Equipment Fault Diagnosis with Vertical Federated Transfer Learning

Jie Shen, Shusen Yang, Cong Zhao, Xuebin Ren, Peng Zhao, Yuqian Yang, Qing Han, and Shuaijun Wu

*Abstract*—Intelligent equipment fault diagnosis based on Federated Transfer Learning (FTL) attracts considerable attention from both academia and industry. It allows real-world industrial agents with limited samples to construct a fault diagnosis model without jeopardizing their raw data privacy. Existing approaches, however, can neither address the intense sample heterogeneity caused by different working conditions of practical agents, nor the extreme fault label scarcity, even zero, of newly deployed equipment. To address these issues, we present FedLED, the first unsupervised vertical FTL equipment fault diagnosis method, where knowledge of the unlabeled target domain is further exploited for effective unsupervised model transfer. Results of extensive experiments using data of real equipment monitoring demonstrate that FedLED obviously outperforms SOTA approaches in terms of both diagnosis accuracy (up to 4.13×) and generality. We expect our work to inspire further study on label-free equipment fault diagnosis systematically enhanced by target domain knowledge.

*Index Terms*—Label-Free Equipment Fault Diagnosis, Unsupervised Transfer Learning, Vertical Federated Learning.



Fig. 1. The General Equipment Fault Diagnosis Scenario based on Vertical Federated Transfer Learning

## I. INTRODUCTION

Proliferating data-driven methods have been proven to be promising in intelligent equipment fault diagnosis [1], where the diagnosing process is usually modeled as classifying 'normal' and 'fault' samples with multiple features extracted from various equipment monitoring signals like current, voltage, vibration, temperature, and acoustic emission [2]. Predominating approaches rely on abundant well-labeled samples to train various Machine-Learning (ML) models (*e.g.*, SVM, DNN) that significantly outperform conventional diagnosis methods based on partial system mechanisms or expert experiences.

However, the application of such labeled sample-intensive methods is severely restricted by the extreme scarcity of fault samples in a wide range of industrial equipment holders (referred to as agents below) in practice [3]–[5]. For example, for any smart manufacturer with a piece of newly deployed equipment, considering that equipment faults are generally small probability events, it usually takes a considerably long time before a single fault occurs and a sample with a determined 'fault' label can be collected [6]. It becomes a common bottleneck for agents to construct an effective fault diagnosis model with few or even zero fault label.

To address this issue, an intuitive idea is to train a model at other agents possessing the same type of equipment with more well-labeled samples (*i.e.*, the source agent), then deploy the trained model at the sample-scarce agent (*i.e.*, the target agent), under the assumption that samples from different agents are Independently and Identically Distributed (IID). Unfortunately,
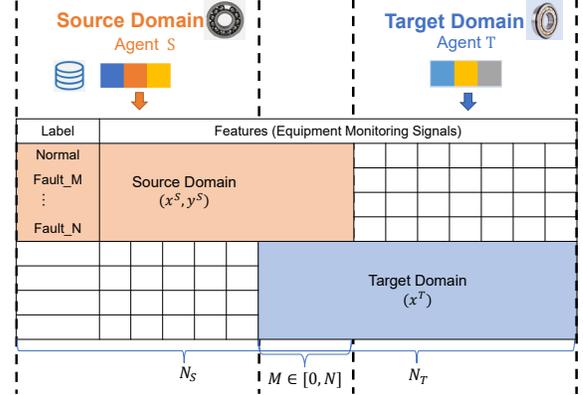
such an assumption is almost impossible to be guaranteed in practice [7]. Due to multi-folded differences such as monitoring setup, working load, transmission path, noise interference, and fault degree, samples from agents with different working conditions inevitably have obvious distribution discrepancies [8], *i.e.*, different agents usually have different *sample domains*. Transfer Learning [9] is primarily used to construct and transfer models across different domains, which has been applied to equipment fault diagnoses [10]. However, existing approaches require exchanging samples between the source and target agents for common knowledge extraction, which poses serious threats to agents' data privacy [11] considering harsh regulations like the General Data Protection Regulation (GDPR) [12]. Therefore, Federated Transfer Learning (FTL) methods [13] emerge and enable transfer learning across 'data islands' of different agents, where all raw samples remain local and only intermediate learning results (*e.g.*, gradients) are exchanged between agents.

As illustrated in Fig.1, to construct a fault diagnosis model under the aforementioned practical scenario, existing FTL methods, however, cannot be directly applied due to the following two fundamental issues of real-world equipment monitoring data:

1) **Intense Sample Heterogeneity:** ML-based FTL methods like [14] require overlapping samples (*i.e.*, the same set of samples possessed by both the source and target agents) to achieve effective model transfer. This is almost impossible in practice since the two agents cannot possess overlapping samples (from the same piece of equipment under the same monitoring

setup) without exchanging raw data. Fault diagnosis with Deep Learning (DL)-based FTL [15] requires no overlapping samples, but assumes that both agents share the same feature space, *i.e.*, horizontal FTL, where both agents use the same set of sensors with the same monitoring deployment. However, in practice, feature spaces of different agents are often multi-scale [16] and intensely heterogeneous [2] due to the agents' different monitoring setups. For example, the source samples $x^S$ are collected by $N_S$ sensors while the target samples $x^T$ are collected by $N_T$ sensors, where only $M$ ($0 \leq M \leq \min(N_S, N_T)$) sensors are shared. Fault diagnosis based on vertical FTL is of great significance.

2) **Zero Fault Label:** Most FTL methods like [17] require a small set of labeled samples in the target domain for model training. However, it is highly likely for practical agents, especially those with newly deployed equipment, to have zero fault label [5], [18], [19] (*i.e.*, with no $y^T$). The only unsupervised FTL fault diagnosis approach [20] requiring no target domain label is restricted by the same feature space assumption and cannot be applied to the aforementioned vertical FTL scenario.

To address the above issues, we present the first unsupervised vertical FTL equipment fault diagnosis method FedLED. It enables model transferring to the target domain with zero fault label from an intensely heterogeneous source domain. The main contributions of this paper are as follows:

1) We are the first to concentrate on the problem of transferring the equipment fault diagnosis model between agents with heterogeneous feature spaces and zero target domain fault label, which is a common bottleneck for the industry. A new fault diagnosis method FedLED based on unsupervised vertical FTL is proposed, which can serve a wide range of agents, especially those with newly deployed equipment.

2) In FedLED, a vertical federated joint domain adversarial adaptation is proposed to map heterogeneous source and target features to a public latent feature space. To enhance the effectiveness of zero fault label model transferring, we construct a novel joint domain alignment that minimizes the distance between the source label distribution and the target classification result distribution, fundamentally different from the conventional pseudo label method that does not comprehensively leverage the target domain information.

3) We conducted extensive experiments using fault datasets of different real equipment (*i.e.*, two bears and a gear) to comprehensively validate the effectiveness of our approach. Experimental results demonstrate that FedLED prominently outperforms state-of-the-art methods in diagnosis accuracy (up to $4.13\times$ higher) under various vertical FTL scenarios. Furthermore, FedLED stably maintains the highest diagnosis accuracy among all comparatives under different vertical FTL settings (*i.e.*, sample/feature overlapping ratios), and shows more obvious advantages under harsher ones (*e.g.*, when there is zero sample/feature overlapping).

The rest of the paper is organized as follows. Section 2 discusses related work of intelligent equipment fault diagnosis. The system model and problem definition are provided in Section 3. Section 4 presents FedLED in detail. Experimental results are provided and comprehensively discussed in Section 5. We conclude the paper in Section 6.

## II. RELATED WORK

In this section, we discuss existing efforts on intelligent equipment fault diagnosis.

### A. Intelligent Equipment Fault Diagnosis

In recent years, intelligent data-driven equipment fault diagnosis has largely benefited from the successful development of deep learning, which has been attracting the industry due to its high diagnostic accuracy [21]. However, existing methods usually rely on abundant well-labeled data or IID assumption [22], [23], severely limiting their usability in practice. Transfer learning can be used to assist training with Non-IID samples from other related agents under different scenarios.

### B. Transfer Learning for Equipment Fault Diagnosis

*1) Heterogeneous Transfer Learning:* Most existing studies follow the same distribution assumption that is difficult to satisfy in practice [24]. Heterogeneous sample processing methods in the field of fault diagnosis mainly process heterogeneous features through feature screening and other methods, and then substitute them into traditional machine learning methods (such as SVM [25], KNN, *etc.*) for training.

*2) Unsupervised Transfer Learning:* Widely adopted unsupervised transfer learning methods can be divided into discrepancy-based and adversarial-based unsupervised transfer learning. Discrepancy-based methods align the source and target domains by measuring the data distributions distance. Common methods for distribution distance measuring include correlation alignment (CORAL) [26], maximum mean discrepancy (MMD) [27], and joint distribution adaptation (JMMD) [28]. Adversarial-based methods use a domain discriminator to reduce the feature distribution discrepancy between source and target domains produced by the feature extractors, enabling cross-machine troubleshooting. Predominating methods include domain adversarial neural network (DANN) [29] and conditional domain adversarial network (CDAN) [30].

Most of the transfer learning methods need to obtain shared knowledge by accessing the raw data of the source and the target domains, casting serious threats to data privacy. FTL methods emerge [31] to address the data privacy issue.

### C. Federated Transfer Learning for Equipment Fault Diagnosis

Traditional FTL methods [31] require shared samples and a small number of labels in the target domain, which limited their use in practical scenarios. With the development of AI technology, deep learning methods have gradually become mainstream. [32] proposes to address domain drift in federated learning based on adversarial domain adaptation. [33] provides

an FTL system that utilizes prior distributional knowledge to reduce inter-domain gaps.

However, all studies above concentrate on horizontal federated learning, where the source and target domains share the same feature space. Their performance cannot be guaranteed in the vertical FTL scenario.

## III. SYSTEM MODEL AND PROBLEM DEFINITION

In this section, we provide the general system model of fault diagnosis based on vertical FTL in Fig.1, and the formal definition of our research problem.

*1) System Model:* Taking the scenario in Fig.1 as an example, there are two vertical FTL agents: the source domain agent S and the target domain agent T. For the well-labeled source domain, $D^S = \{(x_i^S, y_i^S)\}_{i=1}^{N^S}$, where $N^S$ is the total amount of samples of $D^S$. $x_i^S$ and $y_i^S$ denote the $i$th sample in $D^S$ and its label, respectively. $x_i^S \in \mathbb{R}^{N_S}$, where $\mathbb{R}^{N_S}$ is the source feature space, and $N_S$ represents the feature number. For the label-free target domain, $D^T = \{x_i^T\}_{i=1}^{N^T}$, $x_i^T \in \mathbb{R}^{N_T}$, where $N^T$ is the sample number, $\mathbb{R}^{N_T}$ represents the target feature space, and $N_T$ is the feature number. Considering that both the source and target domain agents focus on the same types of faults of the same type of equipment, we assume that both the source and target domains follow the same fault distribution in the same label space.

Considering the heterogeneous source and target domains in practice, there are 1) $D^S \cap D^T = \emptyset$, *i.e.*, no overlapping samples, and 2) $M \in [0, N]$, where $M$ denotes the number of overlapping features between $\mathbb{R}^{N_S}$ and $\mathbb{R}^{N_T}$, and $N = \min(N_S, N_T)$. It is highly likely that such feature space heterogeneity induces a non-negligible distance $dist(\mathbb{R}^{N_S}, \mathbb{R}^{N_T}) > \epsilon > 0$.

*2) Problem Definition:* The vertical FTL task is a classification problem in machine learning with classifier $\mathcal{F}_\mathcal{C}$. Considering our system model, there are two constraint conditions must be met: 1) intense sample heterogeneity, $D^S \cap D^T = \emptyset$ and $dist(\mathbb{R}^{N_S}, \mathbb{R}^{N_T}) > \epsilon$, and 2) zero fault label, $\{y^T\}$ is unavailable. The vertical FTL task can be defined as a constrained optimization problem:

$$\min_{\mathcal{F}_\mathcal{C}} Loss(\mathcal{F}_\mathcal{C}(x), y)$$
$$s.t. \begin{cases} D^S \cap D^T = \emptyset, \\ dist(\mathbb{R}^{N_S}, \mathbb{R}^{N_T}) > \epsilon, \\ \{y^T\} = \emptyset. \end{cases} \quad (1)$$

Considering the intense sample heterogeneity, the feature-based domain adaptation method can be used to build feature extractors $\mathcal{F}_\mathcal{S}$ and $\mathcal{F}_\mathcal{T}$ to map the source and target domains to a latent common space. For zero fault label, the source and target domain output label distributions $(\mathcal{F}_\mathcal{C}(\mathcal{F}_\mathcal{S}), \mathcal{F}_\mathcal{C}(\mathcal{F}_\mathcal{T}))$ can be aligned, since their label spaces are assumed to be the same. Our research problem can be transformed into Eq.(2):

$$\min_{\mathcal{F}_\mathcal{C}, \mathcal{F}_\mathcal{S}, \mathcal{F}_\mathcal{T}} Loss(\mathcal{F}_\mathcal{C}(\mathcal{F}_\mathcal{S}(x^S)), y^S) + \lambda dist(\mathcal{F}_\mathcal{S}(x^S), \mathcal{F}_\mathcal{T}(x^T))$$
$$+ \beta dist(\mathcal{F}_\mathcal{C}(\mathcal{F}_\mathcal{S}(x^S)), \mathcal{F}_\mathcal{C}(\mathcal{F}_\mathcal{T}(x^T))). \quad (2)$$
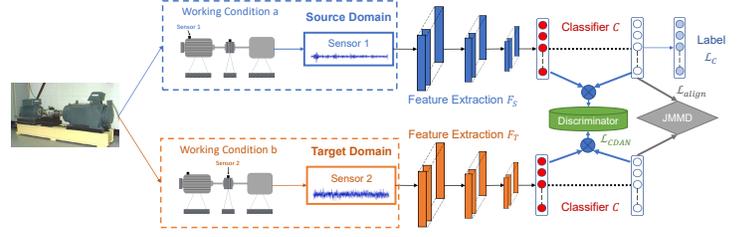


Fig. 2. The Structure of Unsupervised Vertical Federated Transfer Model

The supplementary details of problem definition is presented in Supplemental Material [1] S.1.

## IV. LABEL-FREE EQUIPMENT FAULT DIAGNOSIS WITH VERTICAL FTL

In this section, considering our system model and problem definition, we present an unsupervised vertical FTL equipment fault diagnosis method FedLED. It comprises a vertical federated transfer model that can train a target domain fault diagnoser without fault labels, and an unsupervised federated model training scheme.

### A. Model Structure

The problem in Eq.(2) can be modeled as an unsupervised domain adaptation problem, and there are three joint-optimization objectives, *i.e.*, the classification loss, the feature alignment loss, and the output label alignment loss. The overall network structure is shown in Fig.2. Considering the two constraint conditions of vertical FTL, our unsupervised vertical federated transfer model structure consists of two parts: 1) the vertical federated joint domain adversarial adaptation for the calculation of the classification and feature alignment losses, and 2) the joint domain alignment to calculate the output label alignment loss.

We propose a vertical federated joint domain adversarial adaptation, based on the Adversarial-based method CDAN and vertical federated scheme. The key is a novel conditional domain discriminator conditioned on the cross-covariance of domain-specific feature representations and classifier predictions, which can map heterogeneous source and target feature spaces to a latent common space. Even if the discriminator is completely obfuscated, there is no guarantee that the feature extractor can extract domain-invariant features. Since the domain adversarial adaptation has already aligned the feature space, we additionally add the discrepancy-based joint alignment method as the joint domain alignment to calculate the output label alignment loss, which minimizes the distance between the source label distribution and the target classification result distribution, fundamentally different from the conventional pseudo label method that does not comprehensively leverage the target domain information. More supplementary details of model construction are shown in Supplemental Material S.2.

---

[1][Online Available]: https://github.com/htkg987/FedLED/blob/main/supplemental_material.pdf

*1) Vertical Federated Joint Domain Adversarial adaptation for Intense Sample Heterogeneity:* The vertical federated joint domain adversarial adaptation is based on adversarial-based CDAN [34] to calculate the classification loss and feature alignment loss. The key to CDAN is a novel conditional domain discriminator conditioned on the cross-covariance of feature representations and classifier predictions, which can extract domain-invariant features from heterogeneous features.

**Feature alignment loss:** It is first defined as a minimax optimization problem with two competing error terms, and the overall objective function is as follows:

$$\begin{aligned}
&\mathcal{L}_{CDAN}(\theta_{F_S}, \theta_{F_T}, \theta_D, \theta_C) \\
&= \mathbb{E}_{x_i^S \sim D^S} W(P_i^S) \log\left[D(f_S \otimes g_S)\right] \\
&\quad + \mathbb{E}_{x_j^T \sim D^T} W(P_j^S) \log\left[1 - D(f_T \otimes g_T)\right],
\end{aligned} \quad (3)$$

where $f_S$ and $g_S$ represent the high-level features of the source domain and the output of the classifier through the high-level features, respectively, and $f_T$, $g_T$ correspond to the high-level features of the target domain and their outputs on the classifier. $\otimes$ represents a multi-linear map, which represents the outer product of multiple random vectors. The joint distribution $P(x, y)$ of any two random vectors $x$, $y$ can be obtained by using the cross covariance $(\mathbb{E}_{xy}[\Phi(x) \otimes \Phi(y)])$, where $\Phi$ represents the reproducible kernel function. At the same time, an additional dynamic sample weight method is used to avoid negative samples from affecting training. The update method of the sample weight is as Eq.(4), where $p$ represents the probability that the classifier finally predicts each category:

$$W(p) = 1 + e^{\sum_{c=0}^{N_C-1} p_c \log y_c}. \quad (4)$$

The optimization method of joint domain adversarial learning is: by minimizing (3), optimize the parameters of classifier $C$ and feature extractor $F$ ($F_S$, $F_T$), while maximizing (3) to optimize the domain discriminator $D$, the objective function of each model as follows:

$$(\theta_F^{\hat{t}+1}, \theta_C^{\hat{t}+1}) = arg \min_{\theta_F, \theta_C} \mathcal{L}(\theta_F^t, \theta_C^t, \hat{\theta}_d^t), \quad (5)$$

$$\theta_D^{\hat{t}+1} = arg \max_{\theta_D} \mathcal{L}(\hat{\theta}_F^t, \hat{\theta}_C^t, \theta_d^t). \quad (6)$$

**Classification loss:** In order to avoid the task difference between the target domain and the source domain on the domain-invariant features, the final global classification task is weakened. Therefore, the supervised learning method on source domain is added to prevent the classification bias of the classifier. The objective function of the supervised classification task is as follows:

$$\mathcal{L}_C(\theta_{F_S}, \theta_C) = \mathbb{E}_{(x_i^S, y_i^S) \in D^S} \sum_{i=0}^{N^S-1} \mathcal{L}(C(F_S(x_i^S)), y_i^S). \quad (7)$$

*2) Joint Domain Alignment for Zero Fault Label:* Recent work [34] reveals that even if the discriminator is completely obfuscated, there is no guarantee the feature extractor can extract domain-invariant features. This risk arises from the equilibrium challenges that exist in adversarial learning. Since CDAN has already aligned the feature space, we also add the discrepancy-based alignment [35] method as the joint domain alignment to calculate the output label alignment loss.

**Output label alignment loss:** The objective function of this joint domain align process is as follows:

$$\mathcal{L}_{align} = \|\mathbb{E}_{f_S}[\otimes_{l=1}^{|L|} \phi^l(g_l^S)] - \mathbb{E}_{f_T}[\otimes_{l=1}^{|L|} \phi^l(g_l^T)]\|_{\otimes_{l=1}^{|L|} \mathcal{H}^l}^2. \quad (8)$$

where $g_l^S$ represents the input of high-level features in the source domain into the classifier network, its output on the $l$-th layer, and $\otimes_{l=1}^{|L|} \phi^l(g_l) = \phi^1(g_1 \otimes, \ldots, \phi^L(g_L))$ indicates that the output of each layer of the classifier is projected into a Hilbert space through multidimensional linear mapping. $|L|$ indicates the number of layers in the classifier, usually chooses the last two layers of classifier output for fault diagnosis task alignment on different domains.

### B. The Unsupervised Federated Training Scheme

The entire federated training scheme is divided into two steps in the training process: federated model initialization and federated model training.

*1) Federated Model Initialization:* The federated initialization adopts the pre-training-fine-tuning method demonstrated as effective in [36]. Compared with training from scratch, pre-training the model reduces training time and speeds up the training convergence. The result of pre-training is only a preliminary improvement to prevent overfitting. The objective function of the federated initialization phase is defined as follows:

$$\mathcal{L}_{pre} = \mathbb{E}_{(x_i^S, y_i^S) \sim D^S} \mathcal{L}(C(F_S(x_i^S; \theta_{f_S}); \theta_C), y_i^S). \quad (9)$$

Our federated initialization process is shown in Algorithm 1, where the labeled source domain and label-free target domain are initialized with pre-training and randomly, respectively.

*2) Federated Model Training:* In the federated model raining process, the central server calculates the corresponding loss and gradient, then transmits the gradient to the corresponding participants. The overall workflow is shown in Fig.3.

Considering the model structure, the objective function of FedLED training is shown in Eq.(10).

$$\begin{aligned}
&\mathcal{L}_C(\theta_{F_S}, \theta_{F_T}, \theta_C, \theta_D) \\
&= \mathcal{L}_C(\theta_{F_S}, \theta_C) - \lambda \mathcal{L}_{CDAN}(\theta_{F_S}, \theta_{F_T}, \theta_D, \theta_C) \\
&\quad + \beta \mathcal{L}_{align}(\theta_{F_S}, \theta_{F_T}, \theta_C),
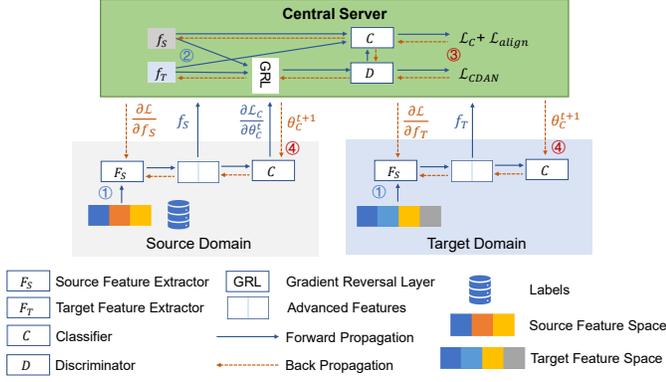\end{aligned} \quad (10)$$

---

**Algorithm 1** Federated Model Initialization

**Input:** Source data $D^S$, randomly initialized source domain model $\theta_{F_S}$, $\theta_C$, learning rate $\mu$.
**Output:** Local model parameters $(\theta_{F_S}, \theta_C)$ with completed federated initialization
**For** epoch = 1 to $N_{epoch}^{pre}$ **do**

1) Randomly select a portion of the source domain data from $D^S$.
2) Forward propagation calculation (9).
3) Backpropagation updates the source domain target extractor, $\theta_{F_S}^{t+1} = \theta_{F_S}^t - \mu(\frac{\partial \mathcal{L}_{pre}}{\partial \theta_{F_S}^t})$.
4) Backpropagation updates the source domain classifier, $\theta_C^{t+1} = \theta_C^t - \mu(\frac{\partial \mathcal{L}_{pre}}{\partial \theta_C^t})$.

**End For**

---

Fig. 3. The Overall Workflow of Federated Model Training

where $\lambda$ and $\beta$ are two hyperparameters, network parameters are updated during training using the Adam optimizer, and the adversarial network optimization problem is solved using a gradient inversion layer [37]. During each training iteration, parameters are updated as follows:

$$
\begin{aligned}
\theta_{F_S}^{t+1} =& \theta_{F_S}^t - \mu\left(\frac{\partial \mathcal{L}_C}{\partial \theta_{F_S}^t} - \lambda\frac{\partial \mathcal{L}_{CDAN}}{\partial \theta_{F_S}^t} + \beta\frac{\partial \mathcal{L}_{align}}{\partial \theta_{F_S}^t}\right) \\
=& \theta_{F_S}^t - \mu\left(\frac{\partial \mathcal{L}_C}{\partial f_S}\frac{\partial f_S}{\partial \theta_{F_S}^t}\right. \\
& \left. - \lambda\frac{\partial \mathcal{L}_{CDAN}}{\partial f_S}\frac{\partial f_S}{\partial \theta_{F_S}^t} + \beta\frac{\partial \mathcal{L}_{align}}{\partial f_S}\frac{\partial f_S}{\partial \theta_{F_S}^t}\right),
\end{aligned}
\tag{11}
$$

$$
\theta_{F_T}^{t+1} = \theta_{F_T}^t - \mu\left(-\lambda\frac{\partial \mathcal{L}_{CDAN}}{\partial f_T}\frac{\partial f_T}{\partial \theta_{F_T}^t} + \beta\frac{\partial \mathcal{L}_{align}}{\partial f_T}\frac{\partial f_T}{\partial \theta_{F_T}^t}\right),
\tag{12}
$$

$$
\theta_C^{t+1} = \theta_C^t - \mu\left(\frac{\partial \mathcal{L}_C}{\partial \theta_C^t} - \lambda\frac{\partial \mathcal{L}_{CDAN}}{\partial \theta_C^t} + \beta\frac{\partial \mathcal{L}_{align}}{\partial \theta_C^t}\right),
\tag{13}
$$

$$
\theta_D^{t+1} = \theta_D^t - \mu\left(-\lambda\frac{\partial \mathcal{L}_{CDAN}}{\partial \theta_D^t}\right).
\tag{14}
$$

Here, $\mu$ is the learning rate and $t$ represents the $t$-th iteration update. Our training process is described in Algorithm 2.

Through the above scheme, a fault diagnoser consisting of $F_t$ and $C$ is obtained, which is deployed at the target domain for online inference.

## V. EVALUATIONS

In this section, we conduct extensive experiments using real equipment fault data under different vertical FTL scenarios for performance evaluation. Experimental methodology is first introduced, then results and analysis are presented.

### A. Experimental Methodology

We validated both the *effectiveness* (diagnosis accuracy) and *generality* (applicability to different levels of source-target domain heterogeneity) of FedLED by comparing it with SOTA approaches under different settings.

---

**Algorithm 2** Federated Model Training

**Input:** Source data $D^S$, target data $D^T$, learning rate $\eta$, hyperparameters $(\lambda, \beta)$, local model of central server $(\theta_D, \theta_C)$, local model of source domain $(\theta_{F_S}, \theta_C)$, local model of target domain $(\theta_{F_T}, \theta_C)$.
**Output:** The trained local model of target domain $(\theta_{F_T}, \theta_C)$ .
**For** epoch = 1 to $N_{epoch}^{train}$ **do**
**Source Domain**
1) Randomly select source domain sample $X_S$ from $D^S$.
2) Obtain the advanced feature $f_S = F_S(X^S)$ and classification loss, send the advanced feature and gradient information $\frac{\partial \mathcal{L}_C}{\theta_C^t}$ to the central server.
3) Block waiting for $\frac{\partial \mathcal{L}}{\partial f_S}$ and $\theta_C^{t+1}$ from the central server.
4) Update feature extractor according to Eq.(11), and overwrite the current classifier parameters with the $\theta_C^{t+1}$ .

**Target Domain**
1) Randomly select target domain sample $X_T$ from $D^T$.
2) Obtain the advanced feature $f_T = F_T(X^T)$ corresponding to $X^T$, and send the advanced feature to the central server.
3) Block waiting for $\frac{\partial \mathcal{L}}{\partial f_T}$ and $\theta_C^{t+1}$ from the central server.
4) Update feature extractor according to Eq.(12), and overwrite the current classifier parameters with the $\theta_C^{t+1}$ .

**Central Server**
1) Accept all agent-uploaded advanced features, as well as tags corresponding to the source domain.
2) Forward propagation, calculation Eq.(10).
3) Backpropagation, update the central server parameters of classifier and discriminator according to Eq.(13) and Eq.(14).
4) Back propagation, the central server calculates $\frac{\partial \mathcal{L}}{\partial f_S}$, $\frac{\partial \mathcal{L}}{\partial f_T}$.
5) Pass $(\frac{\partial \mathcal{L}}{\partial f_S}, \theta_C^{t+1})$ and $(\frac{\partial \mathcal{L}}{\partial f_T}, \theta_C^{t+1})$ to the source domain and the target domain respectively.

**End For**

---

*1) Datasets:* Our experiments used two public datasets containing different monitoring signals and fault labels of three different pieces of real equipment (two bears and a gear): *i.e.*, CWRU [38] and Gearbox [39].

**CWRU** is a widely adopted fault diagnosis benchmark containing three vibration signals(drive-side acceleration data DE, fan-side acceleration data FE, and the reference acceleration data BA) of an SKF6205 bear of 1067 samples. The vibration signal can be acquired by the accelerometer close to the motor-driven end with the 12-kHz sampling frequency. The faults with a single point are introduced to test bearings by electric discharge machining (EDM), resulting in damages of three severity with diameters of 0.007, 0.014, and 0.021 in, respectively. Depending on the location of the faults, there are three types of bearing fault, namely inner-race fault (IF), outer-race fault (OF) and ball fault (BF). Moreover, the bearing of normal condition (NC) is also tested. For heterogeneity, the source and target domains respectively comprised two out of the three signals (features) that were not fully overlapped. Only the source domain possessed fault labels (nine types of bear faults). As shown in Table I, six different fault diagnosis tasks were selected.

**Gearbox** contains eight monitoring signals (with a 12-kHz sampling rate) of a DDS bear and gear. The DDS consists of a brake, a planetary gearbox, a parallel gearbox, and a motor. Additionally, two three-axis (x, y, and z) acceleration sensors collect six channels of vibration signals, which are mounted on the parallel gearbox and the planetary gearbox,

respectively. A torque sensor is installed between the motor and the planetary gearbox to measure load. And there are eight signal characteristics in each data file, which represent: motor vibration signal, vibration signal of planetary gearbox in three directions of x, y, and z, motor torque data and parallel gearbox in three directions of x, y, and z. vibration in one direction. According to the health status of each mechanical equipment, a total of 5115 samples were prepared. Each sample has a number of different features, depending on the task type, with 1024 data points per feature. For heterogeneity, the source and target domains respectively comprised four/five out of the eight signals (features) that were not fully overlapped. Only the source domain possessed fault labels (five types of faults for bear and gear). As shown in Table II, for the bear and gear, six different fault diagnosis tasks were respectively selected.

Learning samples were extracted from all monitoring signals above using the non-overlapping sliding window method [40], and further divided as training and testing sets with a 7:3 ratio. Detailed operations are provided in Supplemental Material S.3.1).

*2) Comparatives:* We compared the performance of FedLED with the following approaches.

1) **Baseline**: Training the model on the source domain, and directly applying the trained model to the target domain.
2) **SFL-multi** [17]: The only FTL-based equipment fault diagnosis method currently available.
3) **Discrepancy-based methods**: SOTA unsupervised transfer learning methods based on different distance metrics, including CORAL [26] using covariance, MK-MMD [41] using MMD, and JAN [35] using JMMD.
4) **Adversarial-based methods**: SOTA unsupervised adversarial-based transfer learning methods, including DANN [29] and CDA+E [30].
5) **Ablation study methods**: Abl Exp 1 and 2 only retained the joint domain alignment and joint domain adversarial adaptation, respectively.

*3) Implementations:* FedLED and all comparatives were implemented using PyTorch V1.3.1, and all evaluations were conducted on a Tesla V100 GPU. For parameter settings, the training batch size and iteration number were 64 and 100, respectively. We use a learning rate of $lr = 0.01(1+10 \times p)^{0.75}$, where $p \in (0, 1]$ is the dynamic decaying rate.

*4) Evaluation Metrics:* We used the accuracy on the target domain testing set as the evaluation metric:

$$Accuracy = \frac{n_{correct}}{n_{test}} \times 100\%, \quad (15)$$

where, $n_{test}$ represents all testing samples, and $n_{correct}$ is the number of all correctly diagnosed samples. To reduce the randomness and singularity, we recorded the average accuracy of 10 repeated experiments as the final results.

### B. Fault Diagnosis Accuracy

To verify the effectiveness of our method, we conducted experiments of FedLED and comparatives on CWRU and Gearbox datasets following the aforementioned tasks. Considering that overlapping samples are mandatory to SFL-multi, we separately set a 10% sample overlapping ratio for it,

TABLE I
FAULT DIAGNOSIS TASKS BASED ON CWRU

| Task | Source Domain Feature | Target Domain Feature | Overlapping Feature |
|---|---|---|---|
| T1 | FE, DE | BA, DE | DE |
| T2 | BA, DE | FE, DE | DE |
| T3 | DE, FE | BA, FE | FE |
| T4 | BA, FE | DE, FE | FE |
| T5 | DE, BA | FE, BA | BA |
| T6 | FE, BA | DE, BA | BA |

TABLE II
FAULT DIAGNOSIS TASKS BASED ON GEARBOX

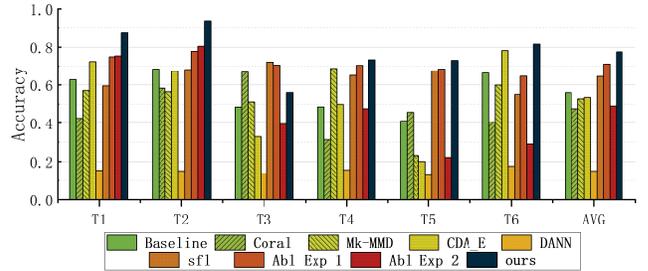| Task | Source Domain Feature | Target Domain Feature | Overlapping Feature |
|---|---|---|---|
| T1 | MV, PL_x, PL_y, PL_z | MV, PA_x, PA_y, PA_z | MV |
| T2 | MV, PA_x, PA_y, PA_z | MV, PL_x, PL_y, PL_z | MV |
| T3 | MT, PL_x, PL_y, PL_z | MT, PA_x, PA_y, PA_z | MT |
| T4 | MT, PA_x, PA_y, PA_z | MT, PL_x, PL_y, PL_z | MT |
| T5 | MV, MT, PL_x, PL_y, PL_z | MV, MT, PA_x, PA_y, PA_z | MV, MT |
| T6 | MV, MT, PA_x, PA_y, PA_z | MV, MT, PL_x, PL_y, PL_z | MV, MT |



Fig. 4. Fault Diagnosis Accuracy on CWRU

while FedLED and other comparatives were set with non-overlapping sample spaces. The fault diagnosis accuracy of all methods is demonstrated in Figs.4~6, and detailed results are provided in Supplemental Material S.3.2). It is obvious that FedLED achieves the highest average diagnosis accuracy (*i.e.*, 77.52%, 95.51%, 98.47%) on all datasets. The performance of FedLED and comparatives are further analyzed as follows.

According to Figs.4~6, DANN performs badly on all tasks, which may be caused by the lack of initialization. Generally, adversarial-based methods increase the domain adaptation ability of $D$ by reducing its discrimination ability. Since DANN lacks initialization, its $D$ has a much stronger discrimination ability, severely restricting the domain adaptation ability. SFL-multi performs stably on all tasks due to the small number of overlapping samples that avoid overfitting.

According to Fig.4, Baseline performs stably on various CWRU tasks with an average accuracy of 50.92%, revealing the relatively low similarity between $\mathbb{R}^{N_S}$ and $\mathbb{R}^{N_T}$. According to Fig.5 and Fig.6, Baseline performs well on Gearbox T4~T6 while weak on Gearbox T1~T3, indicating that $dist(\mathbb{R}^{N_S}, \mathbb{R}^{N_T})$ is relatively small for T4~T6 but large for T1~T3. Particularly, when $dist(\mathbb{R}^{N_S}, \mathbb{R}^{N_T})$ is small (*e.g.* Gearbox T4~T6), adversarial-based methods (including Abl Exp 1) outperform discrepancy-based methods (including Abl Exp 2). In the case with intense feature heterogeneity, *e.g.*
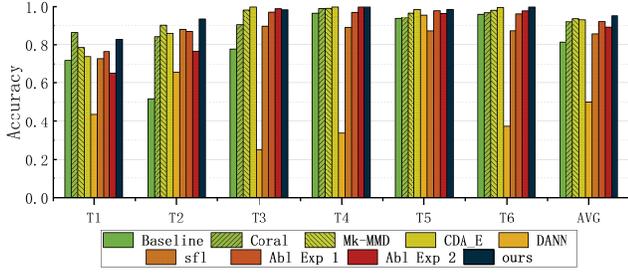
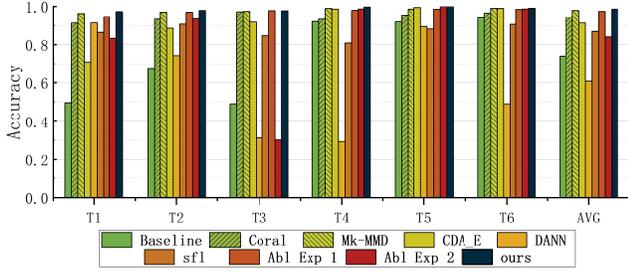Fig. 5. Fault Diagnosis Accuracy on Gearbox-bear



Fig. 6. Fault Diagnosis Accuracy on Gearbox-gear



Fig. 7. Fault Diagnosis Accuracy with Different Sample Overlapping Ratios



Fig. 8. Fault Diagnosis Accuracy with/without Overlapping Features

all CWRU tasks and Gearbox T1∼T3, the performance of adversarial-based methods is weaker than discrepancy-based methods, due to the strong discrimination ability restricting the domain adaptation ability of $D$. Considering such results, the prominent performance of FedLED clearly indicates that our introduction of joint domain alignment manages to effectively eliminate the impact of intense feature heterogeneity on general adversarial-based methods.

### C. Generality under Different Levels of Domain Heterogeneity

To verify the generality of our method, we changed the ratios of sample and feature overlapping in the vertical FTL scenario. Fault diagnosis accuracy of all methods is demonstrated in Fig.7 and Fig.8, respectively. Detailed results are provided in Tables S4∼S7 in Supplemental Material S.3.2).

*1) Impact of Sample Overlapping Ratio:* To study the impact of sample space differences, we first set the sample overlapping ratio between the source and target domains as 0%, 20%, 50%, and 100%, respectively.

According to Fig.7, FedLED achieves the optimal performance under all sample overlapping ratios (*i.e.*, 90.55%, 91.26%, 91.92%, 92.76%), and the accuracy improves slightly as the ratio increases. As the sample overlapping ratio increases, the accuracy of SFL-multi is significantly improved. This is because SFL-multi needs to train a transferable model on overlapping samples, and its performance is positively related to the number of overlapping samples. Different from SFL-multi, the performance of other comparatives requiring no sample overlapping is not obviously enhanced.

*2) Impact of Feature Overlapping Ratio:* To study the impact of feature space differences, we conducted two sets of experiments with 0% and 100% feature overlapping ratios between the source and target domains, respectively.
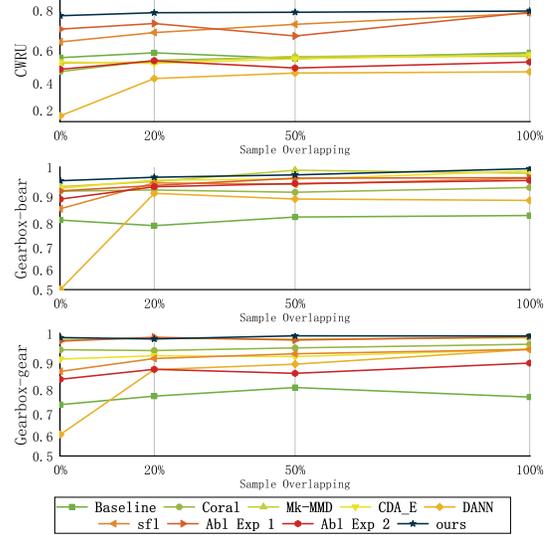
According to Fig.8, FedLED achieves the highest average diagnosis accuracy (*i.e.*, 98.83%, 90.50%) with/without feature overlapping, indicating that it manages to map different feature spaces to a latent common space. Differently, the performance of all comparatives is significantly degraded without overlapping features, clearly demonstrating that the feature space heterogeneity of source and target domains severely restricts the usability of existing fault diagnosis approaches.

### D. Result Analysis

In order to highlight the statistical significance of our method, we performed future analysis of above experimental results on the stability and complexity. Fig.9 shows the boxplot of different methods under all three datasets, reflecting the statistical characteristics of their accuracy. It can be found that our method has relatively good stability and meets the needs of practical applications.

We counter the average time of our method and comparatives running for 100 epochs under the CWRU and Gearbox datasets, which are recorded in Table III. The original SFL method is only suitable for two-classification problems.
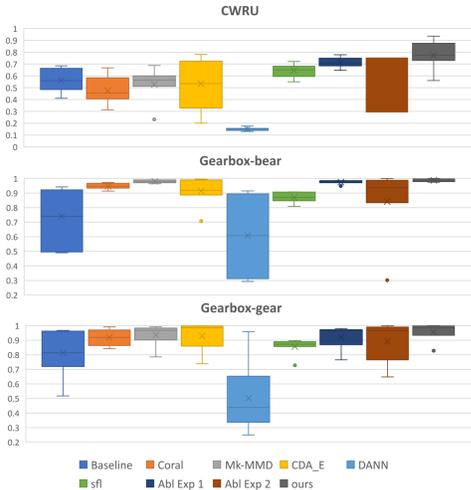
Fig. 9.  Box-Plot of Different Methods

TABLE III
AVERAGE CONSUMING TIME OF DIFFERENT METHODS FOR 100 EPOCH TRAINING (SEC)

| Methods | CWRU | Gearbox |
|---------|------|---------|
| Baseline | 4 | 22 |
| CORAL | 7 | 49 |
| Mk-MMD | 10 | 54 |
| CDA+E | 9 | 49 |
| DANN | 14 | 82 |
| SFL-multi | 168 | 884 |
| Abl Exp 1 | 17 | 92 |
| Abl Exp 2 | 8 | 52 |
| Ours | 18 | 94 |



(a) Trend Chart of Loss     (b) Trend Chart of Accuracy

Fig. 10.  Trend Charts over Time under CWRU

We have extended it to multi-classification problems through multiple classifiers, which also results in a much more time consumption than other methods. Fig.10 (a) and (b) are the trend charts of loss and accuracy over time under CWRU. It can be found that the overall convergence speed of our method is slightly slower than other comparisons due to the longer training time of each epoch. However, because of the improvement in accuracy, the time it takes for our method to achieve the same acceptable accuracy (70%) is actually similar compared with other comparatives.

Network complexity includes space complexity and time complexity. The complexity of our FedLED network is on the same order of magnitude with other comparatives except SFL Our method and other comparatives except SFL are based on the same deep transfer learning backbone, with an additional 3-layer domain discriminator whose space complexity is negligible compared to the backbone, so the space complexity of our method is similar to other comparatives. As for time complexity, FedLED consists of the vertical federated joint domain adversarial adaptation and the joint domain alignment, their time complexity corresponds to adversarial-based and discrepancy-based methods respectively. Therefore, the time complexity of our method is the sum of the two parts' time complexity, which means it is on the same order of magnitude with the time complexities of other comparatives. To sum up, although our method is slightly more time-consuming than other comparatives, compared with the improvement in accuracy, such trade-off is acceptable and worthwhile.

## VI. CONCLUSION

In this paper, we present FedLED, the first unsupervised vertical FTL method facilitating a wide range of industrial agents to conduct label-free equipment fault diagnosis. It enables transferring a fault diagnosis model from a labeled source domain to a highly heterogeneous target domain with zero fault label while preserving the data privacy of both domains. Extensive experiments using real equipment monitoring data clearly demonstrate that FedLED mana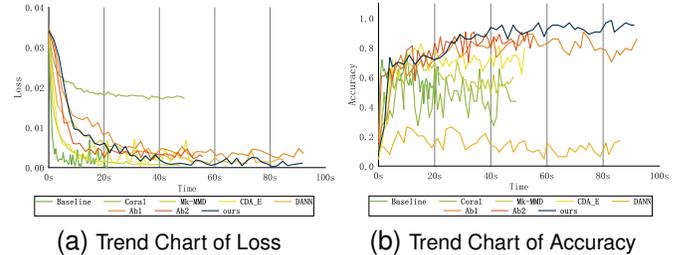ges to achieve obvious advantages in terms of both diagnosis accuracy (up to $4.13\times$ higher) and generality by exploiting knowledge from the unlabeled target domain, different from SOTA approaches intensely depending on source domain knowledge. We expect FedLED to inspire more insights on label-free fault diagnosis enhanced by systematic target domain knowledge extraction, *e.g.*, contrastive learning.

## REFERENCES

[1] J. Sun, C. Yan, and J. Wen, "Intelligent bearing fault diagnosis method combining compressed data acquisition and deep learning," *IEEE Trans. Instrum. Meas.*, vol. 67, no. 1, pp. 185–195, 2017.

[2] M. M. Islam and J.-M. Kim, "Reliable multiple combined fault diagnosis of bearings using heterogeneous feature models and multiclass support vector machines," *Rel. Eng. Syst. Saf.*, vol. 184, pp. 55–66, 2019.

[3] B. Zhao, Z. Niu, Q. Liang, Y. Xin, T. Qian, W. Tang, and Q. Wu, "Signal-to-signal translation for fault diagnosis of bearings and gears with few fault samples," *IEEE Trans. Instrum. Meas.*, vol. 70, pp. 1–10, 2021.

[4] Y. Lou, A. Kumar, and J. Xiang, "Machinery fault diagnosis based on domain adaptation to bridge the gap between simulation and measured signals," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–9, 2022.

[5] Z. Zhao, Q. Zhang, X. Yu, C. Sun, S. Wang, R. Yan, and X. Chen, "Applications of unsupervised deep transfer learning to intelligent fault diagnosis: A survey and comparative study," *IEEE Trans. Instrum. Meas.*, 2021.

[6] Y. Lei, N. Li, L. Guo, N. Li, T. Yan, and J. Lin, "Machinery health prognostics: A systematic review from data acquisition to rul prediction," *Mech. Syst. Signal Process.*, vol. 104, pp. 799–834, 2018.

[7] Y. Xia, C. Shen, D. Wang, Y. Shen, W. Huang, and Z. Zhu, "Moment matching-based intraclass multisource domain adaptation network for bearing fault diagnosis," *Mech. Syst. Signal Process.*, vol. 168, p. 108697, 2022.

[8] Q. Qian, Y. Qin, J. Luo, Y. Wang, and F. Wu, "Deep discriminative transfer learning network for cross-machine fault diagnosis," *Mech. Syst. Signal Process.*, vol. 186, p. 109884, 2023.

[9] P. Liang, C. Deng, J. Wu, G. Li, Z. Yang, and Y. Wang, "Intelligent fault diagnosis via semisupervised generative adversarial nets and wavelet transform," *IEEE Trans. Instrum. Meas.*, vol. 69, no. 7, pp. 4659–4671, 2019.

[10] X. Chen, R. Yang, Y. Xue, M. Huang, R. Ferrero, and Z. Wang, "Deep transfer learning for bearing fault diagnosis: A systematic review since 2016," *IEEE Trans. Instrum. Meas.*, 2023.

[11] X. Ma, C. Wen, and T. Wen, "An asynchronous and real-time update paradigm of federated learning for fault diagnosis," *IEEE Trans. Ind. Informat.*, vol. 17, no. 12, pp. 8531–8540, 2021.

[12] P. Regulation, "General data protection regulation," *Intouch*, vol. 25, 2018.

[13] J. Chen, J. Li, R. Huang, K. Yue, Z. Chen, and W. Li, "Federated transfer learning for bearing fault diagnosis with discrepancy-based weighted federated averaging," *IEEE Trans. Instrum. Meas.*, vol. 71, pp. 1–11, 2022.

[14] W. Yang, J. Chen, Z. Chen, Y. Liao, and W. Li, "Federated transfer learning for bearing fault diagnosis based on averaging shared layers," in *2021 Glob. Reliab. Progn. Health Manag. (PHM-Nanjing)*. IEEE, 2021, pp. 1–7.

[15] W. Zhang and X. Li, "Data privacy preserving federated transfer learning in machinery fault diagnostics using prior distributions," *Struct. Health Monit.*, vol. 21, no. 4, pp. 1329–1344, 2022.

[16] S. Yang, L. Zhang, C. Xu, H. Yu, J. Fan, and Z. Xu, "Massive data clustering by multi-scale psychological observations," *Natl. Sci. Rev.*, vol. 9, no. 2, p. nwab183, 2022.

[17] Y. Liu, Y. Kang, C. Xing, T. Chen, and Q. Yang, "A secure federated transfer learning framework," *IEEE Intell. Syst.*, vol. 35, no. 4, pp. 70–82, 2020.

[18] B. Yang, S. Xu, Y. Lei, C.-G. Lee, E. Stewart, and C. Roberts, "Multi-source transfer learning network to complement knowledge for intelligent diagnosis of machines with unseen faults," *Mech. Syst. Signal Process.*, vol. 162, p. 108095, 2022.

[19] Z. Chai, C. Zhao, and B. Huang, "Multisource-refined transfer network for industrial fault diagnosis under domain and category inconsistencies," *IEEE Trans. Cybern.*, 2021.

[20] E. Tzinis, J. Casebeer, Z. Wang, and P. Smaragdis, "Separate but together: Unsupervised federated learning for speech enhancement from non-iid data," in *2021 IEEE Workshop Appl. Signal Process. Audio and Acoust. (WASPAA)*. IEEE, 2021, pp. 46–50.

[21] X. Li, X. Jia, Y. Wang, S. Yang, H. Zhao, and J. Lee, "Industrial remaining useful life prediction by partial observation using deep learning with supervised attention," *IEEE/ASME Trans. Mechatron.*, vol. 25, no. 5, pp. 2241–2251, 2020.

[22] T. Han, C. Liu, W. Yang, and D. Jiang, "A novel adversarial learning framework in deep convolutional neural network for intelligent diagnosis of mechanical faults," *Knowledge-Based Syst.*, vol. 165, pp. 474–487, 2019.

[23] C. Lu, Z.-Y. Wang, W.-L. Qin, and J. Ma, "Fault diagnosis of rotary machinery components using a stacked denoising autoencoder-based health state identification," *Signal Process.*, vol. 130, pp. 377–388, 2017.

[24] T. W. Rauber, F. de Assis Boldt, and F. M. Varejão, "Heterogeneous feature models and feature selection applied to bearing fault diagnosis," *IEEE Trans. Ind. Electron.*, vol. 62, no. 1, pp. 637–646, 2014.

[25] R. van de Sand, S. Corasaniti, and J. Reiff-Stephan, "Data-driven fault diagnosis for heterogeneous chillers using domain adaptation techniques," *Control Eng. Pract.*, vol. 112, p. 104815, 2021.

[26] B. Sun and K. Saenko, "Deep coral: Correlation alignment for deep domain adaptation," in *Proc. ECCV*. Springer, 2016, pp. 443–450.

[27] D. Sejdinovic, B. Sriperumbudur, A. Gretton, and K. Fukumizu, "Equivalence of distance-based and rkhs-based statistics in hypothesis testing," *Ann. Stat.*, pp. 2263–2291, 2013.

[28] M. Long, J. Wang, G. Ding, J. Sun, and P. S. Yu, "Transfer feature learning with joint distribution adaptation," in *Proc. IEEE Int. Conf. Comput. Vision*, 2013, pp. 2200–2207.

[29] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky, "Domain-adversarial training of neural networks," *J. Mach. Learn. Res.*, vol. 17, no. 1, pp. 2096–2030, 2016.

[30] M. Long, Z. Cao, J. Wang, and M. I. Jordan, "Conditional adversarial domain adaptation," *Adv. Neural Inf. Process. Syst.*, vol. 31, 2018.

[31] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol. (TIST)*, vol. 10, no. 2, pp. 1–19, 2019.

[32] X. Peng, Z. Huang, Y. Zhu, and K. Saenko, "Federated adversarial domain adaptation," *arXiv preprint arXiv:1911.02054*, 2019.

[33] Y. Chen, X. Qin, J. Wang, C. Yu, and W. Gao, "Fedhealth: A federated transfer learning framework for wearable healthcare," *IEEE Intell. Syst.*, vol. 35, no. 4, pp. 83–93, 2020.

[34] S. Arora, R. Ge, Y. Liang, T. Ma, and Y. Zhang, "Generalization and equilibrium in generative adversarial nets (gans)," in *Int. Conf. Mach. Learn.* PMLR, 2017, pp. 224–232.

[35] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Int. Conf. Mach. learn.* PMLR, 2017, pp. 2208–2217.

[36] D. Hendrycks, K. Lee, and M. Mazeika, "Using pre-training can improve model robustness and uncertainty," in *Int. Conf. Mach. Learn.* PMLR, 2019, pp. 2712–2721.

[37] X. Li, W. Zhang, Q. Ding, and X. Li, "Diagnosing rotating machines with weakly supervised data using deep transfer learning," *IEEE Trans. Ind. Informat.*, vol. 16, no. 3, pp. 1688–1697, 2019.

[38] "Case western reserve university bearing data center," accessed 3 August 2022. http: csegroups.case.edu/bearingdatacenter/home.

[39] S. Shao, S. McAleer, R. Yan, and P. Baldi, "Highly accurate machine fault diagnosis using deep transfer learning," *IEEE Trans. Ind. Informat.*, vol. 15, no. 4, pp. 2446–2455, 2018.

[40] R. J. Dugand, J. L. Tomkins, and W. J. Kennington, "Molecular evidence supports a genic capture resolution of the lek paradox," *Nat. Commun.*, vol. 10, no. 1, p. 1359, 2019.

[41] M. Long, Y. Cao, Z. Cao, J. Wang, and M. I. Jordan, "Transferable representation learning with deep adaptation networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 12, pp. 3071–3085, 2018.