

# Anomaly Detection in Industrial Control Systems Based on Cross-Domain Representation Learning

Dongyang Zhan\*, *Member, IEEE*, Wenqi Zhang, Lin Ye, Xiangzhan Yu, Hongli Zhang, and Zheng He

**Abstract**—Industrial control systems (ICSs) are widely used in industry, and their security and stability are very important. Once the ICS is attacked, it may cause serious damage. Therefore, it is very important to detect anomalies in ICSs. ICS can monitor and manage physical devices remotely using communication networks. The existing anomaly detection approaches mainly focus on analyzing the security of network traffic or sensor data. However, the behaviors of different domains (e.g., network traffic and sensor physical status) of ICSs are correlated, so it is difficult to comprehensively identify anomalies by analyzing only a single domain. In this paper, an anomaly detection approach based on cross-domain representation learning in ICSs is proposed, which can learn the joint features of multi-domain behaviors and detect anomalies within different domains. After constructing a cross-domain graph that can represent the behaviors of multiple domains in ICSs, our approach can learn the joint features of them by leveraging graph neural networks. Since anomalies behave differently in different domains, we leverage a multi-task learning approach to identify anomalies in different domains separately and perform joint training. The experimental results show that the performance of our approach is better than existing approaches for identifying anomalies in ICSs.

**Index Terms**—Anomaly detection, industrial control systems, cross-domain learning, multi-graph construction, graph neural networks.

## I. INTRODUCTION

INDUSTRIAL control systems (ICSs) are critical components of many modern industrial and infrastructure systems, such as power grids, and manufacturing plants. They enable remote monitoring and control of physical sensors or actuators based on network communication, which is very important for industry.

However, ICSs are vulnerable to cyberattacks, and the consequences of an attack can be serious. In recent years, lots of attacks against ICSs have triggered serious consequences. For instance, the Stuxnet worm [1] attacked Iranian nuclear facilities, and it demonstrated the potential for cyberattacks to cause physical damage to industrial systems. In 2021, the Colonial Pipeline was hit by a ransomware attack, causing widespread fuel shortages [2]. In 2015 and 2016, Ukraine's power grid was attacked, which caused power outages [3]. In the past, ICSs are isolated from the Internet and run in physically secure locations, so the security mechanisms are

not well designed [4]. Currently, as ICSs are increasingly connected to the Internet, ICSs are facing more and more security threats. Therefore, it is crucial to secure ICSs and ensure the security and reliability of critical infrastructure systems.

There are several challenges of securing ICSs. First, many legacy protocols are not designed with modern security features, making them vulnerable to attacks [5] such as denial of service, man in the middle, etc. For instance, MQTT is a widely used messaging protocol used for machine-to-machine communication in IoT networks that has been found to have many vulnerabilities [6]. Second, the devices (e.g., sensors and actuators) are mostly low-power and have limited computing power for security analysis and protection functions [7]. Therefore, it is difficult for administrators to enhance security directly based on existing ICSs.

Many approaches have been proposed to secure ICSs. First, intrusion detection systems (IDSs) [8] can be applied to secure ICSs, which typically identify the network behavioral characteristics of known attacks but face the problems of identifying only few known attacks and the inability to detect unknown attacks. In addition, an important difference between ICS anomaly detection and network anomaly detection is that ICS managers have access to data from each sensor in the network (e.g., pressure, flow, etc.), based on which anomalies can be more accurately identified. This is because some anomalies do not cause changes in network traffic, but rather in anomalies in sensor data.

Second, security analysis approaches based on behavior models are proposed, which construct the behavior/data models of every sensor, actuator and the interaction between them. If a physical sensor value or network packet does not fit the model, an anomaly alert is generated. But it is difficult to build comprehensive models of complex industrial processes systematically, so these approaches lack generality and are not widely used in ICSs [9].

To address the limitations of traditional approaches, deep learning is utilized for anomaly detection in ICSs. Utilizing the capabilities of deep neural networks, these techniques can learn features that determine anomaly scores, resulting in improved accuracy in anomaly detection. Specifically, the Recurrent Neural Network (RNN) [10] offers the potential to model the time series values of ICS physical devices, such as temperature sensors and pressure sensors. It captures the dependencies between inputs and outputs by extracting partial sequences from relevant inputs and scores anomalies based on the predicted or reconstructed features of ICS physical device values. Further building on the foundations of RNN,

D. Zhan, W. Zhang, L. Ye, X. Yu, H. Zhang are with the School of Cyberspace Science, Harbin Institute of Technology, Harbin, Heilongjiang, 150001.

Z. He is with the Heilongjiang Meteorological Bureau, Harbin, Heilongjiang, 150001.

E-mail: {zhandy, 23S003108, hityelin, yuxiangzhan, zhanghongli}@hit.edu.cn

\* Corresponding Author: zhandy@hit.edu.cn

attention mechanism-imbued time series prediction models such as LSTM-NDT [11] and openGauss [12] have been introduced, showcasing superior performance in identifying enduring correlations in ICS physical device values. Models such as TranAD [13] and Anomaly Transformer [15] adapt the Transformer architecture for time-series anomaly detection, which effectively capture rich associations across entire sequences using their self-attention weight distribution. Besides, the Generative Adversarial Network (GAN) [16] also stands out for its prowess in multivariate anomaly detection through bias reconstruction, as demonstrated in models such as TAnoGAN [17] and GAN-AD [18]. However, it is difficult for these approaches to detect anomaly points in the time series values of ICS physical devices because they cannot model the topological and interaction relationships between nodes in ICSs.

Graph Neural Networks (GNNs) [19] have shown significant promise in graph representation learning. Graph Convolutional Network (GCN) [20] is a convolutional neural network-based model that generates node embeddings by applying convolutional operations on graphs. Consequently, GCN performs well in graph embeddings and can incorporate inter-node correlations into the learning modeling process. Based on GCN, the Graph Deviation Network (GDN) [21] utilizes embedding vectors to capture the physical device features of ICSs (i.e., the values of each sensor or actuator in the real world). It encodes and learns the relationships between pairs of devices as graph edges. It predicts future device behaviors by applying an attention function to the neighboring sensors in the graph. Additionally, it identifies and explains deviations from the learned sensor relationships. While hybrid models (NSIBF [22] and MTAD-GAT [23]) use self-supervised methods to explicitly capture the correlation between features of different physical devices in time series data. Feng and Tian utilize Bayesian filtering to handle implicit relationships, while Zhao et al.'s approach involves direct learning through graph attention networks. However, both methods result in a high level of task complexity. GLIN [24] analyzes the joint representations of local nodes and the global network to perform anomaly detection. This method is used for node-level anomaly detection within industrial control systems. MMGAN [25] jointly learns the features of physical device time series in both the time and frequency aspects, enhancing the model's ability to model distribution. However, these methods are designed for learning features of single-domain data (e.g., physical or network data) in the ICS without considering the joint modeling of multi-domain data. This leaves the models with inadequate detection capability.

In this paper, an anomaly detection approach based on cross-domain representation learning is proposed for ICSs, which integrates the states of both network and physical domains to learn and balance the multi-domain features of ICS elements (e.g., sensors) and detect anomalies in ICSs through predicted values on multiple domains. Since the operation status of ICSs can be reflected in different domains (i.e., the network and physical domains), our approach enables a more comprehensive analysis of ICS data by combining complementary information from different coarse-grained domains. To mitigate the

sparsity problem of ICS behavior data in different domains, we propose a cross-domain graph construction method that presents the relationships between ICS elements in different domains in only one graph.

Our anomaly detection model comprises two stages. In the first stage, an attention-based GCN is used to learn shared representations of different domains based on comprehensive node features. In the second stage, our model learns domain-specific representations to capture the behavioral characteristics of ICSs in different domains. To accomplish this, shared representations are input into several graph models, each of which learns behavioral features in a different domain. Models in different fields predict the future behaviors of ICSs based on reconstructed graph structures. Anomaly detection is performed by calculating the difference between the predicted and actual states, and losses in different domains are used for learning the shared-domain and multi-domain representations. Additionally, we use a multi-gradient descent optimization algorithm based on multi-task learning to adaptively allocate the weight distribution between multiple tasks, which leads to better training results. Our model extracts high-dimensional features related to multi-domain behaviors in ICSs, which is difficult for other existing works.

In summary, the main contributions are presented as follows:

- A cross-domain graph representation method for ICSs is proposed to describe the behaviors between ICS elements of different domains in a single graph, which can represent the behavior features of ICS elements in different domains within a unique graph structure.
- A novel anomaly detection model for ICSs based on attention-based graph convolutional networks is proposed, which learns both domain-shared and domain-specific behaviors to perform anomaly detection in multi-domains. To the best of our knowledge, this is the first work to jointly analyze multi-domain data for anomaly detection in ICSs.
- After implementing the prototype, we conduct baselines and ablation experiments on the water treatment plant datasets with ground truth anomalies, which contains multi-dimensional data. Our results demonstrate that our model detects anomalies more accurately than baseline approaches.
- The code for our proposed model is released on GitHub (<https://github.com/WenqiZhang-HIT/MGDN-project>) to enable interested researchers to reproduce and extend our work.

The rest of this paper is organized as follows. Section II summarizes the related work. The motivation and challenges of this paper are presented in Section III. Section IV describes the proposed anomaly detection approach. The evaluation of our approach is performed in Section V. Section VI summarizes this paper.

## II. RELATED WORK

In this section, we briefly summarize the relevant work on ICS anomaly detection. Furthermore, we briefly introduced previous work closely related to our research, focusing on methods of selecting different ICS features for learning.

### A. Statistics-based ICS Anomaly Detection

Statistics-based methods have become a cornerstone for identifying irregularities that deviate from established behavioral patterns. One seminal approach is the application of multivariate statistical analysis. For example, a real-world application was demonstrated by Chiang, Russell, and Braatz [26], who employed the Multivariate State Estimation Technique (MSET) to model the normal operational regime of a chemical process, which is analogous to a nuclear plant's control system. The strength of their method lies in its ability to accommodate the nonlinear and complex interactions between various control loop parameters. Another significant contribution is the use of Statistical Process Control (SPC) [27], particularly the CUSUM (Cumulative Sum) control chart, which has been skillfully applied to ICSs. For instance, Valle, Li, and Qin [28] demonstrated the utility of CUSUM in monitoring semiconductor manufacturing processes, which can be adapted to the context of ICSs for monitoring and signaling slight, yet significant, shifts in system performance that precede identifiable anomalies. Additionally, Time Series Analysis is a widely respected statistical method for detecting anomalies in ICS, thanks to its ability to capture temporal dependencies. A noteworthy implementation of this technique is provided by Basseville and Nikiforov [29], whose work on the detection of abrupt changes in signals and processes can be applied to SCADA systems [30]. They discuss the use of Autoregressive Integrated Moving Average (ARIMA) models [31] among other methods, which are adept at identifying subtle anomalies in sensor data over time.

### B. Machine Learning-based ICS Anomaly Detection

In contrast to traditional statistical methods, machine learning-based anomaly detection offers a data-driven approach for identifying irregularities within ICSs. These techniques are particularly skilled at handling high-dimensional data and complex pattern recognition tasks. One prominent method is the use of supervised learning algorithms, such as Support Vector Machines (SVMs) [32], which have been applied to power systems for detecting cyber-physical attacks. Notably, Zahid et al. [33] utilized SVMs to classify the states of an electrical grid with high accuracy. Another important approach in machine learning is unsupervised learning, which includes techniques such as k-means clustering and Principal Component Analysis (PCA) [34]. These methods have been instrumental in identifying outliers without the need for labeled data. An exemplary study by Abokifa and Haddad et al. [35] demonstrated the use of PCA in parsing through multivariate physical device data from water distribution systems, effectively isolating instances of operational anomalies. Semi-supervised techniques, which operate with a limited set of labeled data, have also been employed. One such approach is the use of One-Class Classification, where the model learns only from the 'normal' class to determine if new data points fit within that learned distribution. Zhang et al. [36] applied a one-class SVM to vibration data from rotating machinery in an ICS environment, effectively identifying potential mechanical failures. Furthermore, reinforcement learning has

been explored for anomaly detection in ICSs. Algorithms have been designed to learn optimal actions through trial and error, enabling the system to dynamically adapt to new threats. Kim and Chayoung et al. [37] utilized a Q-learning algorithm to monitor network traffic in real-time, differentiating between regular operations and potential cyber security threats.

### C. Deep Learning-based ICS Anomaly Detection

ICS anomaly detection methods based on deep learning utilize deep neural networks to predict multidimensional time series or reconstruct normal time series models. These methods then score anomalies based on the predicted or reconstructed bias. ICS is a complex system that encompasses data from various domains, including physical device display values measured in the real world, network transmission traffic, physical spatial topology, and more. Therefore, based on the dimensions of neural network learning, ICS anomaly detection methods based on deep learning can be categorized into two groups: single-dimensional and multi-dimensional.

The single-dimensional ICS anomaly detection method based on deep learning mainly refers to training and detecting data from only one dimension in ICSs. A framework for anomaly detection in univariate time series data is proposed by [38]. The framework utilizes statistical tests such as the Dickey-Fuller Test (FFT) [39], and Pearson product moment correlation coefficients. It also incorporates different schemes of GRU's deep learning model to perform anomaly detection on various categories of univariate time series. In order to prevent new errors in the generative adversarial networks (GANs) [40] for finding the optimal mapping from real-time space to potential space, the LSTM-based VAE-GAN [41] trains encoders, generators, and discriminators together. This allows for the simultaneous utilization of the mapping abilities of encoders and discriminators, as well as the detection of anomalies in univariate time series based on differences in reconstruction and discriminant results. Aggarwal uses reconstruction errors based on Auto-Encoder (AE) [42] to measure single-dimensional anomalies. DAGMM [43] consists of a compression network and an estimation network, which combine the Deep Auto-Encoder (DAE) [44] density estimation processes for end-to-end joint training. In addition, many methods have applied graph neural networks to anomaly detection, such as GDN [21], which combines structural learning with Graph Attention Network (GAT) [45]. GDN applies attention mechanism to adjacent physical devices on the graph to learn the features of each timestamp of device value, and uses prediction error to detect anomalies in ICSs.

Unlike the single-dimensional methods described above, multidimensional ICS anomaly detection methods typically perform anomaly detection based on information from multiple dimensions of ICSs. For example, GGM-VAE [46] considers the inherent multimodal distribution in time series data, uses a Gated Recurrent Unit (GRU) [47] to discover the correlation between time series, and then uses Gaussian mixture priors in latent space to characterize multimodal data. GLIN [24] takes each physical device in ICSs as a node, introduces a fixed topology as a graph structure, and implements node-level

anomaly detection by merging the local expressions of nodes and the global expressions of the network. Moshe Kravchik and Asaf Shabtai [48] proposed an attack detection method based on simple lightweight neural networks, namely one-dimensional convolution and autoencoder, applied to the time and frequency aspects of time series data.

Although the above methods are referred to as deep learning-based multidimensional ICS anomaly detection, the analysis of data acquired by physical devices (e.g., sensors) is not considered multi-domain analysis. Even though the time series in ICSs are analyzed in terms of time and frequency aspects, the data are obtained from one domain. Some methods introduce information from other dimensions (e.g., the fixed spatial topology of the devices in ICSs) as fixed known information to assist in the learning of the main dimension information, but the introduced information does not represent the behavior features of another domain. Therefore, these methods are often referred to as "multi-variate", rather than "multi-domain". The method proposed in this paper elevates multidimensional ICS anomaly detection to a broader definition, namely multidomain-multivariate ICS anomaly detection. At the same time, information from multiple domains (physical device domain and network transmission domain) in ICSs is both considered. The information in each domain is also multi-variate. Each physical device in ICSs is considered a node, which includes the device value time series on the physical domain and the traffic feature time series on the network domain. The information from multiple domains is fused for cross-domain learning to deeply explore the inherent correlation and potential characteristics of the nodes in ICSs. Besides, multidimensional learning and anomaly detection are performed within each domain to further analyze the security of each specific domain.

### III. MOTIVATION AND CHALLENGES

This section discusses the motivation for considering the cross-domain learning of multi-dimensional data in ICS anomaly detection as well as the challenges brought by multi-dimensional cross-domain representation learning.

#### A. Motivation

ICS is a complex system that comprises traditional industrial control systems as well as computer networks and communication technologies. Its behavior is typically manifested across multiple domains, such as physical device (e.g. sensor or actuator) values, network traffic, and more. Different behaviors exhibit varying patterns across these domains. For instance, ICSs can be targeted with attacks across different domains, such as disrupting devices in the physical domain or launching denial-of-service attacks on nodes in the network domain.

Simultaneously, for a given behavior, the responses across different domains can also differ significantly. For example, a particular attack may cause a significant fluctuation in sensor values in the physical domain while the network domain data remains normal. Alternatively, it may result in distinct anomalies across two different domains simultaneously [49]. Algorithms designed to focus on domain-specific features may

be more sensitive to anomalies unique to those domains. Conversely, a broader selection of features may detect a wider range of abnormal patterns. For instance, an attacker can inject forged water level sensor data through network intrusion, causing the system to falsely perceive the water level as normal [50]. This leads to operational deviations without significant changes in the physical device data. Consequently, anomalies would be challenging to detect promptly if monitoring is confined solely to the physical domain. However, these anomalies can be detected by identifying unusual traffic patterns or undesired communication frequencies in the network domain.

Therefore, analyzing the behavioral characteristics of ICSs solely in a single domain (either physical or network) cannot perceive the security of the entire ICS. Utilizing single-domain feature learning for anomaly detection may result in a high rate of false negatives or false positives. For instance, several single-domain anomaly detection methods, such as DTAAD [51], MSCRED [52], and OmniAnomaly [53], exhibit significant differences in precision and recall in the experiment, which means it is difficult for them to obtain balanced results for real-world anomaly detection. This issue is largely attributed to the data sparsity problem caused by single-domain learning [54].

The data sparsity problem was first introduced in recommendation systems. In recommendation systems, data sparsity refers to the scarcity of interaction data between users and items, making it challenging to accurately predict user ratings for items [55]. The problem of data sparsity not only occurs in recommendation systems but also manifests in various other domains. For instance, Ramchandran and Sangaiah [56] noted that in the field of anomaly detection, the scarcity of anomaly data results in highly sparse datasets. There are several approaches to address this issue, including improved feature engineering, enhanced algorithms, and the utilization of better datasets.

Building upon the potential data sparsity issues that may arise from single-domain ICS anomaly detection, a new approach for anomaly detection in ICS is proposed, which includes multi-domain feature extraction, data learning, and feature prediction.

#### B. Challenges & Solutions

Joint modeling of multi-domain data faces significant challenges due to data heterogeneity, involving variations in data types, structures, and formats across domains. Physical data typically appears as continuous time series, whereas network data is often discrete and event-driven, complicating data fusion. Additionally, different domains may interpret the same event differently, necessitating effective data correlation. Furthermore, the complexity of multidimensional data challenges model design, with high-dimensional datasets potentially leading to sparse feature spaces that impair learning effectiveness [54]. This requires sophisticated models to capture relationships across fields, necessitating advanced architectural and parameter adjustments. Finally, in multi-domain modeling, the cross-domain learning process in high-dimensional space can cause the curse of dimensionality phenomena, such as gradient

explosion or gradient vanishing [57]. In this section, the three primary challenging and significant issues are detailed, and our proposed solutions are presented.

1) *How to extract multi-domain features from ICSs and transform each domain into a unified and efficient feature representation?*: In the real world, ICSs exhibit different behaviors across multiple domains, and the measurement and representation of data also vary. Therefore, a primary challenge is how to obtain processable data from ICSs across multiple domains and transform the data into a unified format suitable for learning.

To address this challenge, we first conduct a survey and identify the widely used SWaT dataset, which contains data from both the physical and network domains simultaneously. We examine the initial recording formats of the behavioral data in each domain and find the differences between them. Based on the observations, it is challenging to combine and unify the initial data from the physical and network domains. For instance, in the physical domain, which is often used for single-domain learning, the data directly records the physical values of each sensor in the ICSs at a per-second granularity. On the other hand, the network domain records the extracted results of all network packets transmitted within the ICS network over time. Thus, our first step is to align the data at a consistent time granularity and extract the network domain data into feature values for each node, facilitating the acquisition of a multi-domain feature matrix for the ICS.

2) *How to efficiently and comprehensively represent node features across multiple domains in ICSs?*: To perform multi-domain analysis, it is necessary to efficiently represent the node time series features of different domains, and behavioral sequence data and topological relationship data should be preserved, thus facilitating the application of various deep learning models for analysis.

To achieve this, we propose a novel representation method called multi-graph structure, which allows the representation of multi-domain data in ICSs on a single graph structure while preserving inter-domain node correlations and temporal information. Each single-domain feature is represented on an individual graph, where the edges represent the associations between two nodes. In this graph, each node represents a sensor or an actuator, which has its own temporal feature vector. By mapping multiple single-domain graph structures based on nodes and concatenating the feature vectors of each node, a multi-graph representation can be obtained.

3) *How to learn the behavioral characteristics of ICSs from multiple domains associatively to achieve efficient anomaly detection?*: Unlike single-domain learning models, to capture additional cross-domain correlation information, it is necessary to design and implement an efficient joint feature learning model that can adequately analyze the multi-graph structure and utilize the learned results for practical feature prediction across different domains.

To that end, we first conduct joint learning on the multi-graph structure to learn the cross-domain features of ICSs on multiple domains. Furthermore, to achieve more accurate anomaly detection, we learn the node features of each specific domain based on the results of cross-domain learning.

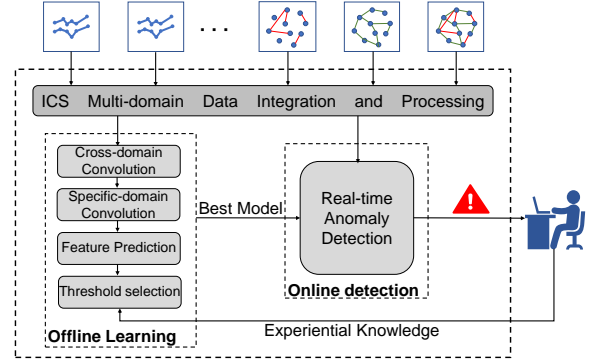


Fig. 1. The construction process of multi-graph representation structure.

After that, multi-domain anomaly detection is performed by predicting the node behaviors on each domain.

As a result, the model follows a multi-task learning mechanism. However, due to the competitive and interdependent nature of the parent-feeding relationships among the multiple convolutional layers in the child, it is challenging to address the resulting multi-objective optimization problem using traditional methods. Therefore, during the forward optimization process, a multi-gradient descent optimizer is introduced to dynamically balance the weights of the multiple domain-specific convolutional layers, accelerating model convergence and avoiding issues such as gradient explosion.

#### IV. METHODOLOGY

To perform accurate anomaly detection in ICSs, a multi-dimensional cross-domain anomaly detection approach is proposed by integrating the above-mentioned solutions, consisting of the processing and representation of the multi-dimensional data in ICSs and cross-domain feature learning and prediction. The overall process of the proposed method is shown in Figure 1.

##### A. Problem Definition

In this paper, our training data consists of sensor data (i.e., multivariate time series) and network data from  $N$  sensors over  $T_{\text{train}}$  time steps. The sensor data is denoted  $\mathbf{s}_{\text{train}} = [\mathbf{s}_{\text{train}}^{(1)}, \dots, \mathbf{s}_{\text{train}}^{(T_{\text{train}})}]$ . In each time step  $t$ , the sensor values  $\mathbf{s}_{\text{train}}^{(t)} \in \mathbb{R}^N$  form an  $N$  dimensional vector representing the values of our  $N$  sensors. Following the many unsupervised anomaly detection approaches, the training data consists of only normal data.

For each domain  $d$  ( $d = \{1, \dots, D\}$ ), we first construct an undirected weighted graph  $\mathcal{G}_d = (\mathcal{V}, \mathcal{E}_d)$ . As these  $D$  domains are correlated and share the same set of nodes, we then construct the cross-domain graph as an undirected weighted multi-graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , which contains the node set  $\mathcal{V}$  with  $N$  nodes and the edge set  $\mathcal{E}$  with  $D$  types of edges, i.e.,  $\mathcal{E} = \{\mathcal{E}_1, \dots, \mathcal{E}_D\}$ .

Our problem can be formally stated as follows: with an undirected weighted multi-graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , and the node feature matrix  $\mathbf{X} \in \mathbb{R}^{N \times M}$  representing input for each node

as an  $M$  dimensional feature vector, our goal is to learn a set of embeddings for all nodes in each subgraph  $\mathcal{G}_d$ , i.e.,  $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_D\}$  ( $\mathbf{X}_d \in \mathbb{R}^{N \times E}$  is the node embedding in subgraph  $\mathcal{G}_d$ ), and predict the characteristics of the next time step  $\tau$  for each node. Test data for anomaly detection, which comes from the same  $N$  sensors but over a set of  $T_{\text{test}}$  time steps, the test data is denoted  $\mathbf{s}_{\text{test}} = [\mathbf{s}_{\text{test}}^{(1)}, \dots, \mathbf{s}_{\text{test}}^{(T_{\text{test}})}]$ .

### B. Multi-Graph Construction

As previously noted, the majority of current anomaly detection approaches for ICSs are constrained to analyze information from a single domain, such as physical domain time series (i.e., physical values of individual sensors over time), which makes it challenging to analyze the inherent correlations among nodes (i.e., sensors) in different domains. This limitation results in the absence of some key information, as the behaviors of ICSs are not limited to one domain. For instance, most ICSs are designed to obtain real-time physical data from sensors and send instructions to actuators through network traffic, so the physical data and network data are both important for security analysis.

A multi-graph structure is proposed to integrate and represent information from multiple dimensions of ICSs on a single graph, which can be used to learn the potential correlations and features of nodes across domains. To this end, we extract node embeddings from both the physical and network domains.

In real-world scenarios, information from different domains is represented and measured in various ways. For instance, in the physical domain, sensor values can be directly recorded at each time step (in seconds), which can be easily converted into time series information. On the other hand, in the network domain, all transmitted network packets over a period of time are recorded, resulting in a large amount of data per unit time. This data cannot be directly used as time series information for nodes, and the analysis results of network packets cannot be directly used as network domain characteristics of each node. Although both domains measure information within the same time period, they cannot be directly mapped and fused because the time granularity differs. Therefore, simply merging the data from these two domains is difficult and meaningless, requiring mathematical and statistical methods to unify the data and map them to the same time granularity.

To overcome these challenges, we introduce an embedding vector for each sensor in each domain, representing its physical and network features:

$$\mathbf{X}_{\text{phy}}^{(i)}, \mathbf{X}_{\text{net}}^{(i)} \in \mathbb{R}^T, \text{ for } i \in \{1, 2, \dots, N\} \quad (1)$$

This enables us to effectively present the underlying relationships among different domains of data.

We follow the methods of GDN [21] to extract and construct physical embeddings  $\mathbf{X}_{\text{phy}}^{(i)}$  based on physical sensor data. When dealing with network data, we analyze the unique characteristics of network traffic. Specifically, we calculate the number of data packets received, sent, and total data load of each node within a consistent time granularity (in seconds) that matches the physical information. To capture the temporal nature of this data, we extract the time series of these features

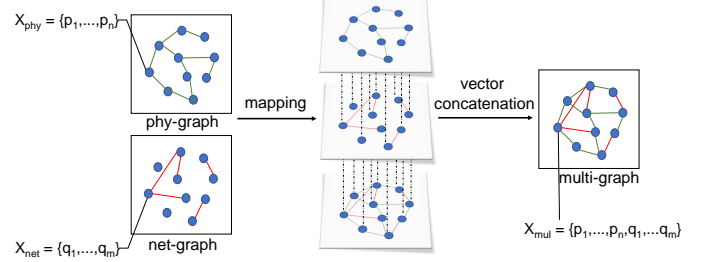


Fig. 2. The construction process of multi-graph representation structure.

through a sliding window. Then, the original network traffic data is simplified into a representative feature time series. Finally, the time granularity of network domain information is aligned with physical domain information to facilitate the merging of multi-domain information into a shared graph using time mapping methods.

The multi-graph structure is constructed based on the results of the above approaches. The cosine calculation method is employed to calculate the probability of node correlation based on the node embedding matrix for each domain. Then, we apply the *Topk* threshold to identify the highest-probability neighbors for each node, thus generating a graph structure for each domain. Finally, we consolidate the edges of different domains into a single graph, allowing for the possibility of multiple edges between any two nodes. This process can generate a multi-graph that effectively represents the complex relationships within the data. To form the node shared embedding in the multi-graph structure, we append the network domain embedding vector to the rear of the physical domain embedding vector through the vector concatenation method [24]. The construction process of the multi-graph structure is illustrated in Figure 2.

### C. Attention-based Cross-Domain Graph Neural Network

Our model takes a multi-graph structure as input and first delves into modeling and exploring the correlation between dimensions in ICSs through domain-crossed convolutional layers. Then divide the learned results into separate domains for the second round of domain-specific learning. Finally, the multilayer perceptron is used in the prediction module to perform feature prediction based on the output of convolutional layers in each domain. Specifically, the convolution operations in the model are implemented by multiple basic attention-based convolutional blocks with different targets. The framework of the model is shown in Figure 3.

**Basic Attention-based Convolutional Block** To effectively learn node embeddings, we perform convolutional learning using basic blocks on the graph structure. The basic convolutional block consists of graph convolutional layers and incorporates attention mechanisms to better explore the potential correlations between nodes. The introduction of graph



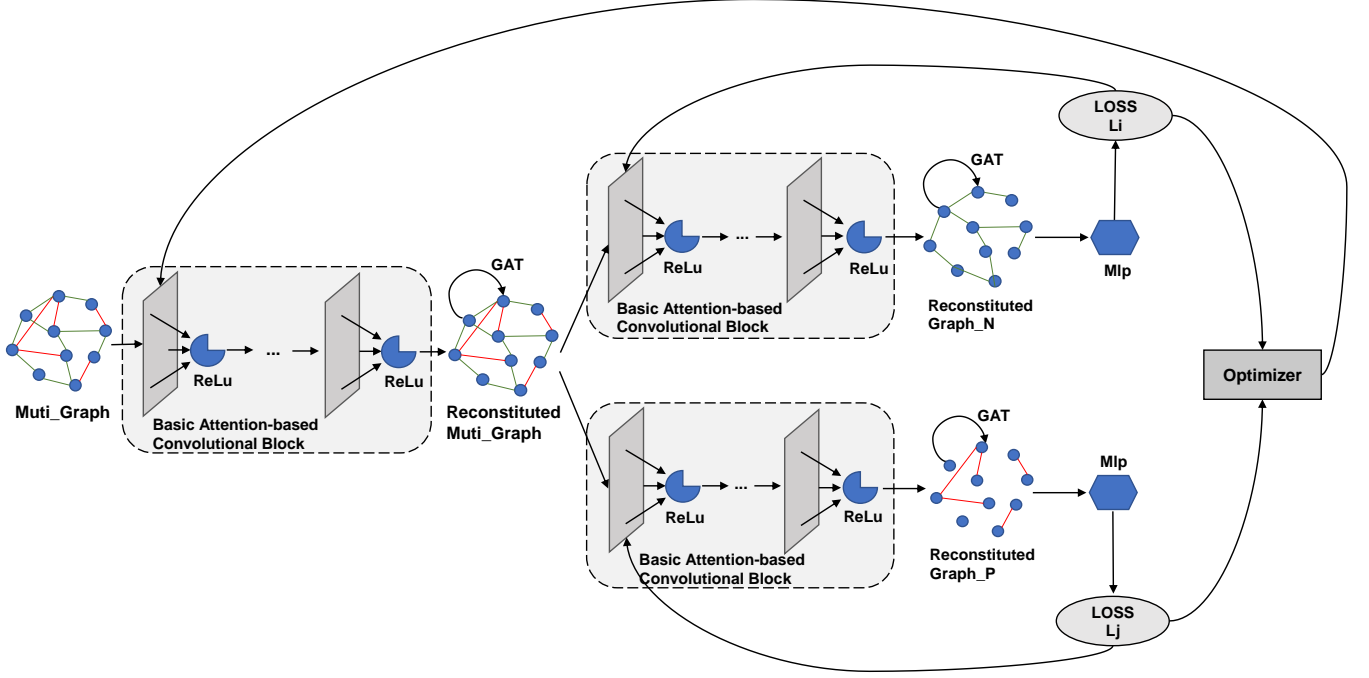


Fig. 3. The overall framework for cross-domain representation learning in our model.

attention mechanisms allows for a more holistic consideration of the dynamic relationships between nodes on a global scale.

To calculate the similarity between the embedding vectors of node  $i$  and its candidates  $j \in \mathcal{C}_i$  (in the case without prior information, the candidates of sensor  $i$  are all sensors  $\mathcal{C}_i$ ), the calculation formula is as follows:

$$e_{ji} = \frac{\mathbf{x}_i^\top \mathbf{x}_j}{\|\mathbf{x}_i\| \cdot \|\mathbf{x}_j\|} \quad \text{for } j \in \mathcal{C}_i \quad (2)$$

Among them,  $e_{ji}$  represents the normalized dot product between the embedding vectors of sensor  $i$ , and the candidate  $j$ . The first  $k$  normalized dot products are selected by the function  $FstK$  to form the adjacency matrix  $A$ . The value of  $k$  can be selected based on expectations and prior knowledge.

$$\mathcal{N}_i = FstK(\{e_{ji} : j \in \mathcal{C}_i\}) \quad (3)$$

$$\begin{cases} \mathbf{A}_{ji} = 1, & j \in \mathcal{N}_i \\ \mathbf{A}_{ji} = 0, & j \notin \mathcal{N}_i \end{cases} \quad \text{for } i \in \mathcal{V} \quad (4)$$

Based on the adjacency matrix, the basic convolutional block in the model utilizes the node embedding  $X_i$  and their neighboring embedding  $X_j$  to compute the aggregated representation of each node. The calculation formula for feature extraction in the basic convolutional block is defined as follows:

$$X_i^{(l+1)} = \text{ReLU} \left( \mathbf{W} \lambda_{i,i} X_i^{(l)} + \mathbf{W} \sum_{j \in \mathcal{N}_i} \lambda_{i,j} X_j^{(l)} \right) \quad (5)$$

$\mathbf{W}$  is a weight matrix that can be trained,  $X_i^l$  is the embedding of node  $i$  in layer  $l$ , and  $X^0$  is the input feature matrix,  $\lambda_{i,j}$  is the attention coefficient of node  $j$  to node  $i$ , which is obtained by normalizing the attention score and the calculation formula is as follows,  $\mathbf{a}$  is the learning coefficient vector of the attention mechanism:

$$\lambda_{i,j} = \text{Softmax}(\text{LeakyReLU}(\mathbf{a}^\top (X_i \oplus X_j))) \quad (6)$$

**Cross-Domain Learning** To learn multiple types of node representations on multiple domains, a cross-domain graph neural network is proposed, which mainly has two stages: multi-graph learning and subgraph learning. In the multi-graph learning phase, all domain information is reflected on the same graph, where each node behaves differently in different domains. By leveraging the interrelatedness of these representations, we aim to learn a shared representation that represents cross-domain shared information in multiple domains. In the subgraph learning stage, we focus on learning the specific representation of each node on each domain to encode specific information on different subgraphs. The detailed description of our methodology and the overall architecture of our model are illustrated in Figure 3.

Specifically, our model comprises basic convolutional blocks with multiple objectives, following the multi-task learning regime (MTL) [58]. The cross-domain learning layer encodes multi-graph structures and node shared attributes as the input of the basic attention-based convolutional block to generate shared node embedding. Next, the learned shared node embeddings are split into node embeddings on two

domains. The convolutional layers of two specific domains take the node embeddings on each subgraph as inputs and encode them to generate node embeddings on specific domains. The interaction structure of all graph convolutional layers is depicted in Figure 3.

As described above, our model consists of multiple graph convolutional layers, each with different learning contents. The parameter  $W_s$  of cross-domain graph convolutional layer is shared by multiple domains, while the parameters  $W_d$  ( $d = 1, \dots, D$ ) of domain-specific graph convolutional layers are used for each specific domain. To train the model and benefit the learning process of all domains, that is, to enable multiple tasks in the model to achieve an optimal state, we need to optimize all target parameters ( $W_s, W_1, \dots, W_D$ ).

In MTL [58], the commonly used method for optimizing the objective function is to calculate the weighted sum of the loss  $L_d$  for all tasks statically or dynamically. However, stacking multiple layers brings additional difficulties to the model training, and adjusting the weights of each task to obtain the optimal training effect is time-consuming [59]. In our model, the optimization goal is to minimize the loss of each specific domain learning module, as shown in formula (7):

$$\begin{cases} \min_{W_s, W_1} L_1(W_s, W_1) \\ \dots \\ \min_{W_s, W_D} L_D(W_s, W_D) \end{cases} \quad (7)$$

To find the optimal solution for each objective, that is, the optimal solution in each domain, we employ a multi-gradient descent optimizer [60]. As in single-objective optimization, multi-objective optimization problems can be solved to local optima through gradient descent. The Multiple Gradient Descent Algorithm (MGDA) [61] utilizes the Karush–Kuhn–Tucker (KKT) conditions [73] (Eq. (8)), which is a necessary condition for the optimal solution of multi-objective optimization problems and is defined as follows:

$$\begin{cases} \sum_{d=1}^D \alpha_d \frac{\partial L_d(W_s, W_d)}{\partial W_s} = 0 \\ \frac{\partial L_d(W_s, W_d)}{\partial W_d} = 0, (\forall d \in \{1, \dots, D\}) \\ \sum_{d=1}^D \alpha_d = 1 \\ \alpha_d \geq 0, (\forall d \in \{1, \dots, D\}) \end{cases} \quad (8)$$

where  $\alpha_d$  is the weight of objective  $L_d(W_s, W_d)$ .

Either the solution to equation Eq. (10) is 0, ensuring that the result satisfies Eq. (8), or there exists a solution that provides a descent direction to enhance all tasks in Eq. (7) as mentioned in [62]. Consequently, solving Eq. (8) is equivalent to solving Eq. (10), as illustrated below. Additionally, Eq. (9) signifies the weighted sum of gradients for each specific domain.

$$K = \left\| \sum_{d=1}^D \alpha_d \frac{\partial L_d(\Theta_s, \Theta_d)}{\partial \Theta_s} \right\|_2^2 \quad (9)$$

$$\min_{\alpha_1, \dots, \alpha_D} \left\{ K \mid \sum_{d=1}^D \alpha_d = 1, \alpha_d \geq 0 \right\} \quad (10)$$

The resulting MTL algorithm enables gradient descent of multiple tasks on specific parameters, followed by solving Eq. (8) and applying the solution Eq. (9) as a gradient update to shared parameters. When the number of specific domains is 2, the optimization objective Eq. (9) can be simplified as:

$$\min_{\alpha} \left\| \alpha \frac{\partial L_1(W_s, W_1)}{\partial W_s} + (1 - \alpha) \frac{\partial L_2(W_s, W_2)}{\partial W_s} \right\|_2^2 \quad (11)$$

$s.t. \ 0 \leq \alpha \leq 1$

The Eq. (10) is a unary quadratic equation of  $\alpha$ . With the weight  $\alpha$ , we update the parameters of the model as follows:

$$W_d = W_d - \eta \frac{\partial L_d(W_s, W_d)}{\partial W_d} \quad (12)$$

$$W_s = W_s - \eta \sum_{d=1}^D \alpha_d \frac{\partial L_d(W_s, W_d)}{\partial W_s} \quad (13)$$

**Feature Prediction** After the embedding learning, we obtain representations for all nodes on each specific domain, namely  $\mathbf{X}_d = \{\mathbf{x}_{d,1}^{(t)}, \dots, \mathbf{x}_{d,N}^{(t)}\}$ . For each node  $i$ , we compute the inner product of the feature representation  $\mathbf{x}_{d,i}^{(t)}$  and the corresponding time series embedding  $\mathbf{v}_i$ , and employ the product of all nodes as the input of the Multilayer Perceptron (MLP) [63] with output dimension  $N$  to predict the node vector at time  $t$  (i.e. the value of the sensor):

$$y = f_3(w_3 f_2(w_2 f_1(w_1 [\mathbf{v}_1 \circ \mathbf{x}_1, \dots, \mathbf{v}_N \circ \mathbf{x}_N] + b_1) + b_2) + b_3) \quad (14)$$

$\mathbf{v}_i$  is an embedded vector randomly initialized for each node and trained together with the model. The similarity between  $\mathbf{v}_i$  and  $\mathbf{v}_j$  represents the similarity of behavior between nodes  $i$  and  $j$ . These embeddings allow attention mechanism to consider more association possibilities, which helps in learning graph structures.

A three-layer Perceptron is used as the output layer, where  $f_1$ ,  $f_2$  and  $f_3$  are Linear, BatchNorm1d and LeakRelu,  $w_i$  and  $b_i$  are the weight matrix and offset vector of layer  $i$ , respectively, and  $y$  is the final prediction vector in each dimension.

To efficient train our model, after obtaining the predicted output  $\hat{\mathbf{p}}^{(t)}$  of the model at time  $t$ , we employ the MAE (Mean Absolute Error) between the predicted output and the actual observation data  $\mathbf{r}^{(t)}$  as the loss function on each domain-specific convolutional layer:

$$L1loss = \frac{1}{T_{\text{train}} - w} \sum_{t=w+1}^{T_{\text{train}}} \left\| \mathbf{y}^{(t)} - \mathbf{r}^{(t)} \right\| \quad (15)$$



## V. EVALUATION

In this section, we first describe the experimental dataset and performance indicators. Then, we evaluate our model and compare it with several state-of-the-art baselines that are designed for detecting anomalies in ICSs. Finally, we analyze the experimental results to demonstrate the effectiveness of our model.

We use the PyTorch [64] 1.12.1 with CUDA 11.3 and PyTorch Geometric [65] 2.1.0 library to implement our model and train all models. We set the hyper parameters as described in the baseline models of the papers. We use the following hyperparameter values in our model.

- Batch size = 32
- Window size = 15
- The value of K in the  $FstK = 20$
- Dropout in convolutional blocks = 0.2
- Depth of hidden layers in convolutional blocks = 2
- The momentum of the SGD optimizer = 0.9

The effect of different learning rate combinations for each task on anomaly detection performance is analyzed in Section V-C. To accelerate the convergence process, we employ the Stochastic Gradient Descent (SGD) optimizer [66] to guide model updates. To achieve faster and smoother optimization, we incorporate a momentum term into the SGD optimizer. Apart from the learning rates, other hyperparameters are selected using grid search and empirical rules.

### A. Dataset

By adjusting the data preprocessing algorithm, our approach can be readily extended to various dual-domain ICS datasets. By augmenting the number of single graphs in the multi-graph representation algorithm, expanding domain-specific learning modules, and optimizing hyperparameters for the multi-gradient descent optimization algorithm, our method can be applicable to anomaly detection across three or more domains within ICS. It is also adaptable to various types of ICS or other industrial environments. However, after extensive research and investigation, show that many existing public ICS anomaly detection datasets only offer physical domain-related data (i.e., sensor values) [67]–[69], or only provide traffic data in the network domain [70], [71]. In contrast, the Secure Water Treatment (SWaT) dataset [72], which includes both physical and network data, proves ideal for our experiments due to its comprehensive coverage.

We utilize the SWaT dataset, which is collected at the Singapore University of Technology from a water treatment system. The SWaT dataset includes uninterrupted 11-day operational period data, during which the system runs from an empty state to a fully operational state. For the initial seven days, the system operates under normal conditions without any attacks or malfunctions. In the remaining four days, while data collection continues, a range of network and physical attacks are launched on the SWaT system. We are grateful for the open source code provided by S Tuli et al, who propose the TranAD [13]. This resource is instrumental to adapt the nine distinct baseline models to the SWaT dataset in our experiments. Because our methodology is uniquely capable of

detecting multi-domain data simultaneously, we provide data specific to either the physical or network domain in SWaT, as delineated in thier works.

In this dataset, the physical domain includes 51 columns of values, representing the measured values per second of 51 nodes (sensors or actuators). On the other hand, the network domain includes 16 columns of values extracted from the transmitted packets in the ICS network during the measurement period. We first determine that each domain contains 51 feature vectors, representing the features of all nodes in the ICS, which is beneficial for constructing the graph structure. Then, we extract the initial network domain data of each node into three feature values: the number of sent packets, the number of received packets, and the sending overall payload, corresponding to the measured values of each node on the physical domain.

Table 1 summarizes the statistics of the SWaT dataset in both physical and network domains.

TABLE I  
THE STATISTICS OF THE SWAT DATASET IN DIFFERENT DOMAINS.

Domain	Train	Test	Features	Anomalies
Physical	21,830	34,201	51	16.61%
Network	2,1830	34,201	3	16.61%

### B. Evaluation Method

The precision, recall, F1-Score and False Positive Rate (FPR) are applied to evaluate the performance of our proposed model. Depending on the optimization direction of the metrics, they are divided into positive (precision, recall, and F1 score) when higher values indicate better performance, and negative (FPR) when lower values indicate better performance [14]. Let TP represent the number of sequences that are correctly predicted as positive, TN denote the number of sequences that are correctly classified as negative, FN denote the number of traces that are positive but are incorrectly predicted as negative, and FP indicate the number of traces that are negative but are predicted as positive. The calculation methods for these metrics are shown in Eq. (16)-(19). To detect anomalies, we use the maximum anomaly score on the validation dataset to set a threshold. During testing, any time step where the anomaly score exceeds the threshold will be considered abnormal.

$$Precision = \frac{TP}{TP + FP} \quad (16)$$

$$Recall = \frac{TP}{TP + FN} \quad (17)$$

$$F1 - Score = \frac{2 * Precision * Recall}{Precision + Recall} \quad (18)$$

$$FPR = \frac{FP}{FP + TN} \quad (19)$$

TABLE II  
ANOMALY DETECTION PERFORMANCE COMPARISON IN TERMS OF FPR(%), PRECISION(%), RECALL(%), AND F1-SCORE(%) OF OUR MODEL WITH  
BASELINE METHODS.

Method	FPR	Precision	Recall	F1
DTAAD	13.33	59.88	99.99	74.90
GDN	10.70	64.91	99.45	78.55
LSTM-AD	13.33	59.88	99.99	74.90
MAD-GAN	13.57	59.45	99.99	74.57
MSCRED	13.33	59.89	99.99	74.91
MTAD-GAT	13.39	59.78	99.99	74.83
OmniAnomaly	13.36	59.83	99.99	74.87
TranAD	13.35	59.85	99.99	74.88
USAD	13.26	60.02	99.99	75.01
<b>Our Model</b>	<b>3.07</b>	<b>84.65</b>	<b>85.12</b>	<b>84.88</b>

### C. Results and Discussion

Our model outperforms when evaluated against nine other baseline methods, as well as in Ablation Experiments. Through evaluation with the baselines, our model not only achieves the best precision, FPR and F1-score, demonstrating its outstanding ability to accurately identify anomalies but also an excellent balance of precision and recall. These excellent performances are crucial for the practical application of our model in ICSs. Additionally, the ablation experiments show that both the attention mechanism and the multi-gradient descent optimization algorithm are indispensable components of our model. These findings suggest that the attention mechanism and the multi-gradient descent optimization algorithm all contribute to model performance, and our model is accurate and reliable for real-world ICS anomaly detection. In this section, we provide detailed descriptions of the results and in-depth discussions.

1) *Comparison with Baselines:* To demonstrate the overall performance of our model, we compare it with nine methods for the detection of multivariate time series anomalies in ICSs. Table II provides the FPR, precision, recall, and F1 scores for our model and baseline models for the SWaT dataset. As shown in Table II, our model shows excellent generalization capability and achieves the best FPR and F1 scores consistently on the SWaT dataset. This enhancement is due to our model's cross-domain representation learning mechanism, which integrates global and local information of nodes in multiple domains and enables the model to discover more potential nonlinear relationships. The shared embeddings take into account both single-domain and cross-domain semantics and provides broader node and correlation information for further single-domain representation learning.

**Our models' performance excels several state-of-the-art deep learning-based single-dimensional anomaly detection models (e.g., TranAD, GDN, LSTM-AD and so on), with an F1 score gain of no less than 6.33%. Our model further introduces network and other domain information into the time**

series information of general physical sensors for learning and training. Through cross-domain representation learning, multi-domain information is first fused for learning, and cross-domain embedding is split into physical and other domains for specific learning, ultimately obtaining node embeddings of ICSs on multiple domains. Therefore, compared to the single-domain detection model, our model can obtain the best performance.

**Our model surpasses the other multidimensional anomaly detector (e.g., MTAD-GAT) evaluated in this paper in terms of all evaluation metrics.** MTAD-GAT considers each univariate time-series as an individual feature and tries to model the correlations between different features explicitly, while the temporal dependencies within each time series are modeled at the same time. MTAD-GAT leverages two parallel graph attention layers to learn the relationships between different time-series and timestamps dynamically, while our model considers more generalized multidimensional information, where data from different domains is measured and represented independently. In addition, our model adopts an aggregation-splitting training framework, which helps to mine more cross-domain information and assist in feature learning on various single domains. Profiting from the multi-graph representation method, our model can comprehensively and efficiently learn the information of nodes on multiple dimensions simultaneously during the cross-domain representation learning stage and explore the potential correlations between different domains of nodes, providing more prior knowledge for specific domain learning and more selectivity for anomaly detection. Therefore, our model achieves the best evaluation performance compared to other multidimensional anomaly detection models.

**Compared with all other evaluated anomaly detection models, our model exhibits the best performance, with the highest F1 score and precision and the smallest difference between precision and recall, indicating that our model has the best performance balancing ability.** Although some

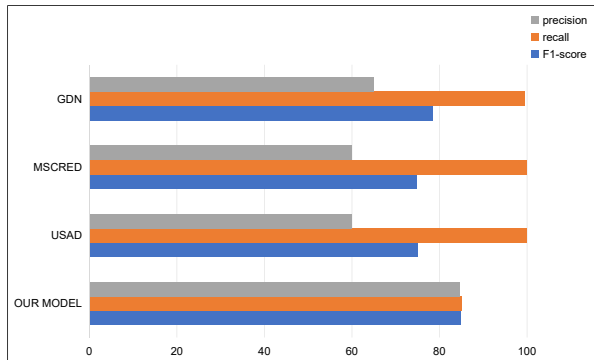


Fig. 4. Comparison of performance balance among our model, GDN, USAD, and MSCRED.

methods (e.g., TranAD, MSCRED, and DTAAD) can obtain higher scores than the proposed model on recall metrics, they are worse on other metrics. This is due to the fact that most of these models only use sensor values from 51 nodes as training data and implement deep learning-based anomaly detection, thus making it difficult to mitigate the sparsity problem of ICS sensor data. Although some methods have relatively high precision, the recall is generally low and do not achieve relatively balanced detection results. This suggests that the relevant information is not fully utilized and is subject to some problems such as false alarm, inapplicability to minority issues, and insensitivity to changes in practical situations [74].

Although our model does not achieve the highest precision or recall, its evaluation results are more balanced, with both metrics tending towards higher values. Figure 4 shows the comparison between the three models that achieved the best values in all indicators except for our model in the experiment and our model, making it easy to visually observe the difference in size of the indicators. There are significant differences between the precision and recall performance of the GDN, USAD, and MSCRED models, while our model shows an average level of excellence in various indicators and reaches the best level of surpassing other baselines on the F1 Score.

To comprehensively and reliably demonstrate the advantages of our model in practical applications, anomaly detection performance is evaluated against baseline models using the FPR. This metric not only illustrates the model's efficiency in correctly identifying true anomalies but also highlights its superior ability to minimize false alarms, providing a robust indication of its practical utility in real-world scenarios. As shown in Table II, our model achieves an excellent result of 3.07% on the metric FPR, which is the lowest compared with the baselines, indicating superior performance in reducing false positives. The results show that our model is more advantageous in scenarios where the cost of false positives is serious, and lower false positives are beneficial to maintain the normal operation of ICS and ensure the reliability of anomaly detection results [75].

2) *Ablation Experiments*: In this section, we evaluate the effectiveness of some methods (e.g., the attention mechanism, etc.) in our model by comparing F1 scores, precision, and recall with ablation experiments.

**Effectiveness of Graph Attention.** We examine the influence of the graph attention mechanism in our model by disabling it and instead aggregating using equal weights assigned to all neighbors. From Table III, we can find that the removal of the graph attention mechanism causes a 6.49% decline in the average F1 score. Since nodes in ICSs have very different behaviors, treating all neighbors equally makes noise and misleads our model. This verifies that employing the attention mechanism in our model can better learn complex graph structure information, enable the model to better mine node features based on correlation, and thus perform more accurate anomaly detection.

TABLE III  
ANOMALY DETECTION ACCURACY IN TERMS OF PRECISION (%), RECALL (%), AND F1-SCORE (%) OF OUR MODEL AND ITS VARIANTS.

Method	Precision	Recall	F1
<b>Our Model</b>	<b>84.65</b>	<b>85.12</b>	<b>84.88</b>
-Attention	99.25	64.78	78.39

Vanilla GCNs utilize a fixed convolutional kernel to aggregate information from neighboring nodes, attributing equal influence to all [20]. However, in multi-dimensional data contexts, the significance of different neighboring nodes can vary substantially, which Vanilla GCNs overlooks. Moreover, learning node features in high-dimensional spaces presents substantial complexity, with abundant potential node association information yet to be explored. Integrating an attention mechanism into the vanilla GCNs framework allows for the dynamic weighting of neighbor nodes based on their importance, enabling the model to flexibly aggregate information and capture more intricate inter-node relationships [76]. Without the attention mechanism, the effectiveness of the model in information aggregation and feature selection diminishes, leading to an inability to fully capture all pertinent positive samples. Consequently, the model becomes more conservative, resulting in an increase in precision but a notable decrease in recall.

**Effectiveness of the Multi-Gradient Descent Optimizer.** To verify the effectiveness of employing a multi-gradient descent optimizer in our model, we replace the multi-gradient descent optimization algorithm with a weighted sum method that statically calculates all task losses and sets the static ratio of loss for two domain-specific learning tasks: physical weight 0.5 vs. network weight 0.5, physical weight 0.25 vs. network weight 0.75, and physical weight 0.75 vs. network weight 0.25. We record the evaluation results (F1-score, precision, and recall) of the above three cases and the original model, whose results are shown in Table IV.

The evaluation result shows that setting the loss ratio of the physical domain and the network domain to 0.5 vs. 0.5 can achieve better results than other static weighting methods, but our current model using the multi-gradient descent optimizer is still better than at least 1.01% in the F1 score. Therefore, the introduction of a multi-gradient descent optimization algorithm is beneficial for the model to better dynamically adjust

TABLE IV  
PERFORMANCE COMPARISON OF OUR MODEL WITH MODELS EMPLOYING  
THE STATIC CALCULATION LOSS SUM METHOD INSTEAD OF THE  
MULTI-GRADIENT DESCENT OPTIMIZER METHOD.

Method	Precision	Recall	F1
0.5 : 0.5	83.15	82.34	83.87
0.25 : 0.75	49.30	84.46	62.28
0.75 : 0.25	80.99	83.39	82.79
<b>Our Model</b>	<b>84.65</b>	<b>85.12</b>	<b>84.88</b>

the weights of multi-task learning, providing optimization guarantees for accelerating model convergence and improving model performance. Furthermore, we find that static loss calculation methods are difficult to explain and obtain the optimal weight allocation parameters, making it difficult to effectively and reasonably solve competitive multi-objective optimization problems.

3) *Sensitivity Analysis*: We evaluate our model's sensitivity to different learning rate settings and combinations across multiple tasks. The learning rates for the three different tasks are set, and we observe the model's performance in anomaly detection. LR1 and LR2 represent the learning rates for the physical domain learning and network domain learning tasks, respectively, while LR3 represents the learning rate for the cross-domain learning task. This sensitivity is evaluated by comparing the Precision, Recall, and F1-score metrics.

TABLE V  
PERFORMANCE COMPARISON OF DIFFERENT LEARNING RATE SETTINGS  
(ALL VALUES IN %).

Learning Rates			Precision	Recall	F1-Score
LR1	LR2	LR3			
0.001	0.001	0.001	99.99	64.24	78.24
0.001	0.001	0.1	99.99	64.24	78.24
0.001	0.1	0.1	85.01	84.72	84.86
0.1	0.1	0.1	82.25	83.17	83.44
0.1	0.1	0.001	81.07	83.65	82.75
<b>0.1</b>	<b>0.001</b>	<b>0.001</b>	<b>84.65</b>	<b>85.12</b>	<b>84.88</b>

As shown in Figure V, the results indicate that the model's performance is sensitive to the learning rate settings for different tasks. Specifically, setting LR1, LR2, and LR3 to identical values hinders the model from achieving optimal anomaly detection performance. This is due to the varying complexity and convergence rates of each task during the learning process. When the other two learning rates are held constant, increasing the value of LR1 affects the F1-score by up to 6.64%, indicating a high sensitivity to LR1. This sensitivity can be attributed to the high feature complexity associated with the physical domain learning task that LR1

regulates [77]. The optimal combination, LR1 = 0.1, LR2 = 0.001, and LR3 = 0.001, achieves a balance between Precision and Recall, enhancing overall performance. Consequently, this configuration was selected as the final setting for our model.

## VI. CONCLUSION

This paper presents an anomaly detection approach based on cross-domain representation learning, which combines ICS data on multiple domains for cross-domain learning and anomaly detection. We propose a multi-graph representation method to uniformly represent ICS multi-domain data on just one graph structure and then design an attention-based cross-domain graph convolutional network for learning embedding. We evaluate our model on a large-scale real-world dataset, and experimental results show that our model outperforms baselines. In addition, our model can better balance the relationship between reducing false positives and improving anomaly detection precision, providing a more practical and ideal ICS anomaly detection model. Future works may come from two aspects. First, more practical ICS anomaly detection datasets with aligned multi-domain data will be explored, and our model will be developed to adapt to a wider range of application scenarios. Secondly, the advantages and costs of cross-domain learning in ICS anomaly detection will be further explored.

## ACKNOWLEDGMENTS

This work was supported by the National Natural Science Foundation of China under Grants No. 62302122 and No. 62172123, the Natural Science Foundation of Heilongjiang Province of China under Grants No. LH2023F017, and CCF-Huawei Populus Grove Fund under Grants No. CCF-HuaweiSY202411.

## REFERENCES

- [1] S. Karnouskos, "Stuxnet worm impact on industrial cyber-physical system security," in *IECON 2011-37th Annual Conference of the IEEE Industrial Electronics Society*. IEEE, 2011, pp. 4490–4494.
- [2] T. Tsvetanov and S. Slaria, "The effect of the colonial pipeline shutdown on gasoline prices," *Economics Letters*, vol. 209, p. 110122, 2021.
- [3] D. U. Case, "Analysis of the cyber attack on the ukrainian power grid," *Electricity Information Sharing and Analysis Center (E-ISAC)*, vol. 388, no. 1-29, p. 3, 2016.
- [4] A. Cook, H. Janicke, L. Maglaras, and R. Smith, "An assessment of the application of it security mechanisms to industrial control systems," *International Journal of Internet Technology and Secured Transactions*, vol. 7, no. 2, pp. 144–174, 2017.
- [5] J. M. Taylor and H. R. Sharif, "Security challenges and methods for protecting critical infrastructure cyber-physical systems," in *2017 International Conference on Selected Topics in Mobile and Wireless Networking (MoWNeT)*. IEEE, 2017, pp. 1–6.
- [6] D. Dinculeană and X. Cheng, "Vulnerabilities and limitations of mqtt protocol used between iot devices," *Applied Sciences*, vol. 9, no. 5, p. 848, 2019.
- [7] R. Yang, H. He, Y. Wang, Y. Qu, and W. Zhang, "Dependable federated learning for iot intrusion detection against poisoning attacks," *Computers & Security*, vol. 132, p. 103381, 2023.
- [8] H.-J. Liao, C.-H. R. Lin, Y.-C. Lin, and K.-Y. Tung, "Intrusion detection system: A comprehensive review," *Journal of Network and Computer Applications*, vol. 36, no. 1, pp. 16–24, 2013.
- [9] M. Kravchik and A. Shabtai, "Efficient cyber attack detection in industrial control systems using lightweight neural networks and pca," *IEEE transactions on dependable and secure computing*, vol. 19, no. 4, pp. 2179–2197, 2021.

- [10] D. Mandic and J. Chambers, *Recurrent neural networks for prediction: learning algorithms, architectures and stability*. Wiley, 2001.
- [11] 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.
- [12] G. Li, X. Zhou, J. Sun, X. Yu, Y. Han, L. Jin, W. Li, T. Wang, and S. Li, "opengauss: An autonomous database system," *Proceedings of the VLDB Endowment*, vol. 14, no. 12, pp. 3028–3042, 2021.
- [13] 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.
- [14] O. Rainio, J. Teuhon, and R. Klén, "Evaluation metrics and statistical tests for machine learning," *Scientific Reports*, vol. 14, no. 1, p. 6086, 2024.
- [15] J. Xu, H. Wu, J. Wang, and M. Long, "Anomaly transformer: Time series anomaly detection with association discrepancy," *arXiv preprint arXiv:2110.02642*, 2021.
- [16] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, "Generative adversarial networks: An overview," *IEEE signal processing magazine*, vol. 35, no. 1, pp. 53–65, 2018.
- [17] M. A. Bashar and R. Nayak, "Tanogan: Time series anomaly detection with generative adversarial networks," in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*. IEEE, 2020, pp. 1778–1785.
- [18] D. Li, D. Chen, J. Goh, and S.-k. Ng, "Anomaly detection with generative adversarial networks for multivariate time series," *arXiv preprint arXiv:1809.04758*, 2018.
- [19] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE transactions on neural networks*, vol. 20, no. 1, pp. 61–80, 2008.
- [20] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [21] 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.
- [22] C. Feng and P. Tian, "Time series anomaly detection for cyber-physical systems via neural system identification and bayesian filtering," in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2021, pp. 2858–2867.
- [23] H. Zhao, Y. Wang, J. Duan, C. Huang, D. Cao, Y. Tong, B. Xu, J. Bai, J. Tong, and Q. Zhang, "Multivariate time-series anomaly detection via graph attention network," in *2020 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2020, pp. 841–850.
- [24] L. Shuaiyi, K. Wang, L. Zhang, and B. Wang, "Global-local integration for gnn-based anomalous device state detection in industrial control systems," *Expert Systems with Applications*, vol. 209, p. 118345, 2022.
- [25] T. Pandeva and M. Schubert, "Mmgan: Generative adversarial networks for multi-modal distributions," *arXiv preprint arXiv:1911.06663*, 2019.
- [26] L. H. Chiang, E. L. Russell, and R. D. Braatz, *Fault detection and diagnosis in industrial systems*. Springer Science & Business Media, 2000.
- [27] R. Caulcutt, "Statistical process control (spc)," *Assembly Automation*, vol. 16, no. 4, pp. 10–14, 1996.
- [28] S. Valle, W. Li, and S. Qin, "Selection of the best variable for the cusum and ewma schemes," *Technometrics*, vol. 41, no. 3, pp. 221–234, 1999.
- [29] M. Basseville, I. V. Nikiforov et al., *Detection of abrupt changes: theory and application*. prentice Hall Englewood Cliffs, 1993, vol. 104.
- [30] A. Daneels and W. Salter, "What is scada?" 1999.
- [31] J. L. M. Saboia, "Autoregressive integrated moving average (arima) models for birth forecasting," *Journal of the American Statistical Association*, vol. 72, no. 358, pp. 264–270, 1977.
- [32] M. A. Hearst, S. T. Dumais, E. Osuna, J. Platt, and B. Scholkopf, "Support vector machines," *IEEE Intelligent Systems and their applications*, vol. 13, no. 4, pp. 18–28, 1998.
- [33] M. Zahid, F. Ahmed, N. Javaid, R. A. Abbasi, H. S. Zainab Kazmi, A. Javaid, M. Bilal, M. Akbar, and M. Ilahi, "Electricity price and load forecasting using enhanced convolutional neural network and enhanced support vector regression in smart grids," *Electronics*, vol. 8, no. 2, p. 122, 2019.
- [34] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley interdisciplinary reviews: computational statistics*, vol. 2, no. 4, pp. 433–459, 2010.
- [35] A. A. Abokifa, K. Haddad, C. S. Lo, and P. Biswas, "Detection of cyber physical attacks on water distribution systems via principal component analysis and artificial neural networks," in *World Environmental and Water Resources Congress 2017*, 2017, pp. 676–691.
- [36] J. Zhang, Q. Zhang, X. Qin, and Y. Sun, "A two-stage fault diagnosis methodology for rotating machinery combining optimized support vector data description and optimized support vector machine," *Measurement*, vol. 200, p. 111651, 2022.
- [37] C. Kim and J. Park, "Designing online network intrusion detection using deep auto-encoder q-learning," *Computers & Electrical Engineering*, vol. 79, p. 106460, 2019.
- [38] J.-B. Kao and J.-R. Jiang, "Anomaly detection for univariate time series with statistics and deep learning," in *2019 IEEE Eurasia conference on IOT, communication and engineering (ECICE)*. IEEE, 2019, pp. 404–407.
- [39] R. Mushtaq, "Augmented dickey fuller test," 2011.
- [40] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [41] Z. Niu, K. Yu, and X. Wu, "Lstm-based vae-gan for time-series anomaly detection," *Sensors*, vol. 20, no. 13, p. 3738, 2020.
- [42] G. E. Hinton, A. Krizhevsky, and S. D. Wang, "Transforming auto-encoders," in *Artificial Neural Networks and Machine Learning–ICANN 2011: 21st International Conference on Artificial Neural Networks, Espoo, Finland, June 14–17, 2011, Proceedings, Part I 21*. Springer, 2011, pp. 44–51.
- [43] B. Zong, Q. Song, M. R. Min, W. Cheng, C. Lumezanu, D. Cho, and H. Chen, "Deep autoencoding gaussian mixture model for unsupervised anomaly detection," in *International conference on learning representations*, 2018.
- [44] S. Lange and M. Riedmiller, "Deep auto-encoder neural networks in reinforcement learning," in *The 2010 international joint conference on neural networks (IJCNN)*. IEEE, 2010, pp. 1–8.
- [45] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [46] Y. Guo, W. Liao, Q. Wang, L. Yu, T. Ji, and P. Li, "Multidimensional time series anomaly detection: A gru-based gaussian mixture variational autoencoder approach," in *Asian Conference on Machine Learning*. PMLR, 2018, pp. 97–112.
- [47] R. Dey and F. M. Salem, "Gate-variants of gated recurrent unit (gru) neural networks," in *2017 IEEE 60th international midwest symposium on circuits and systems (MWSCAS)*. IEEE, 2017, pp. 1597–1600.
- [48] M. Kravchik and A. Shabtai, "Detecting cyber attacks in industrial control systems using convolutional neural networks," in *Proceedings of the 2018 workshop on cyber-physical systems security and privacy*, 2018, pp. 72–83.
- [49] C. Feng, T. Li, and D. Chana, "Multi-level anomaly detection in industrial control systems via package signatures and lstm networks," in *2017 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*. IEEE, 2017, pp. 261–272.
- [50] G. R. MR, C. M. Ahmed, and A. Mathur, "Machine learning for intrusion detection in industrial control systems: challenges and lessons from experimental evaluation," *Cybersecurity*, vol. 4, no. 1, p. 27, 2021.
- [51] L. Yu, "Dtaad: Dual tcn-attention networks for anomaly detection in multivariate time series data," *arXiv preprint arXiv:2302.10753*, 2023.
- [52] C. Zhang, D. Song, Y. Chen, X. Feng, C. Lumezanu, W. Cheng, J. Ni, B. Zong, H. Chen, and N. V. Chawla, "A deep neural network for unsupervised anomaly detection and diagnosis in multivariate time series data," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 1409–1416.
- [53] 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.
- [54] S. Manzhos and M. Ihara, "Advanced machine learning methods for learning from sparse data in high-dimensional spaces: A perspective on uses in the upstream of development of novel energy technologies," *Physchem*, vol. 2, no. 2, pp. 72–95, 2022.
- [55] P. Resnick and H. R. Varian, "Recommender systems," *Communications of the ACM*, vol. 40, no. 3, pp. 56–58, 1997.
- [56] A. Ramchandran and A. K. Sangaiah, "Unsupervised anomaly detection for high dimensional data—an exploratory analysis," in *Computational intelligence for multimedia big data on the cloud with engineering applications*. Elsevier, 2018, pp. 233–251.
- [57] D. Peng, Z. Gui, and H. Wu, "Interpreting the curse of dimensionality from distance concentration and manifold effect," *arXiv preprint arXiv:2401.00422*, 2023.
- [58] S. Ruder, "An overview of multi-task learning in deep neural networks," *arXiv preprint arXiv:1706.05098*, 2017.

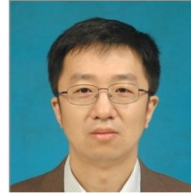
- [59] Y. Ouyang, B. Guo, X. Tang, X. He, J. Xiong, and Z. Yu, "Learning cross-domain representation with multi-graph neural network," *arXiv preprint arXiv:1905.10095*, 2019.
- [60] O. Sener and V. Koltun, "Multi-task learning as multi-objective optimization," *Advances in neural information processing systems*, vol. 31, 2018.
- [61] J.-A. Désidéri, "Multiple-gradient descent algorithm (mgda) for multiobjective optimization," *Comptes Rendus Mathématique*, vol. 350, no. 5-6, pp. 313–318, 2012.
- [62] G. Gordon and R. Tibshirani, "Karush-kuhn-tucker conditions," *Optimization*, vol. 10, no. 725/36, p. 725, 2012.
- [63] H. Taud and J. Mas, "Multilayer perceptron (mlp)," *Geomatic approaches for modeling land change scenarios*, pp. 451–455, 2018.
- [64] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, "Automatic differentiation in pytorch," 2017.
- [65] M. Fey and J. E. Lenssen, "Fast graph representation learning with pytorch geometric," *arXiv preprint arXiv:1903.02428*, 2019.
- [66] L. Bottou, "Stochastic gradient descent tricks," in *Neural Networks: Tricks of the Trade: Second Edition*. Springer, 2012, pp. 421–436.
- [67] C. M. Ahmed, V. R. Palleti, and A. P. Mathur, "Wadi: a water distribution testbed for research in the design of secure cyber physical systems," in *Proceedings of the 3rd international workshop on cyber-physical systems for smart water networks*, 2017, pp. 25–28.
- [68] A. Al Zaki Khan and G. Serpen, "Intrusion detection and identification system design and performance evaluation for industrial scada networks," *arXiv e-prints*, pp. arXiv:2012, 2020.
- [69] J. J. Downs and E. F. Vogel, "A plant-wide industrial process control problem," *Computers & chemical engineering*, vol. 17, no. 3, pp. 245–255, 1993.
- [70] S. Choi, J.-H. Yun, and S.-K. Kim, "A comparison of ics datasets for security research based on attack paths," in *Critical Information Infrastructures Security: 13th International Conference, CRITIS 2018, Kaunas, Lithuania, September 24-26, 2018, Revised Selected Papers 13*. Springer, 2019, pp. 154–166.
- [71] A. Dehlaghi-Ghadim, M. H. Moghadam, A. Balador, and H. Hansson, "Anomaly detection dataset for industrial control systems," *IEEE Access*, 2023.
- [72] A. P. Mathur and N. O. Tippenhauer, "Swat: A water treatment testbed for research and training on ics security," in *2016 international workshop on cyber-physical systems for smart water networks (CySWater)*. IEEE, 2016, pp. 31–36.
- [73] O. L. Mangasarian, *Nonlinear programming*. SIAM, 1994.
- [74] D. Chicco and G. Jurman, "The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation," *BMC genomics*, vol. 21, no. 1, pp. 1–13, 2020.
- [75] J. Mijalkovic and A. Spognardi, "Reducing the false negative rate in deep learning based network intrusion detection systems," *Algorithms*, vol. 15, no. 8, p. 258, 2022.
- [76] M. Gupta, R. Khanna, D. Choudhary, and N. Rao, "Knowledge graph reasoning based on attention gcN," *arXiv preprint arXiv:2312.10049*, 2023.
- [77] M. Crawshaw, "Multi-task learning with deep neural networks: A survey," *arXiv preprint arXiv:2009.09796*, 2020.



**Dongyang Zhan** is an associate professor in School of Cyberspace Science at Harbin Institute of Technology. He received the B.S. degree in Computer Science from Harbin Institute of Technology from 2010 to 2014. From 2015 to 2019, he has been working as a Ph.D. candidate in School of Computer Science and Technology at HIT. His research interests include cloud computing and security.



**Wenqi Zhang** is a master's student from the Harbin Institute of Technology, China. Her research focuses on IoT security.



**Xiangzhan Yu** is a professor in School of Cyberspace Science at Harbin Institute of Technology. His main research fields include: network and information security, security of internet of things and privacy protection. He has published one academic book and more than 50 papers on international journals and conferences.



**Hongli Zhang** received her BS degree in Computer Science from Sichuan University, Chengdu, China in 1994, and her Ph.D. degree in Computer Science from Harbin Institute of Technology (HIT), Harbin, China in 1999. She is currently a Professor in School of Cyberspace Science in HIT. Her research interests include network and information security, network measurement and modeling, and parallel processing.



**Lin Ye** received the Ph.D. degree at Harbin Institute of Technology in 2011. From January 2016 to January 2017, he was a visiting scholar in the Department of Computer and Information Sciences, Temple University, USA. His current research interests include network security, peer-to-peer network, network measurement and cloud computing.



**Zheng He** is an engineer in Heilongjiang Meteorological Bureau. She received her bachelor's and Master's degrees in Meteorology Science in Nanjing University of Information Science and Technology from 2011 to 2018. From 2018, she has been working in Weather Modification Office of Heilongjiang Province. Her research interests include climate change, weather modification and machine learning.