

Adaptive Memory Networks with Self-supervised Learning for Unsupervised Anomaly Detection

Yuxin Zhang, Jindong Wang, Yiqiang Chen, Han Yu, Tao Qin

Abstract—Unsupervised anomaly detection aims to build models to effectively detect unseen anomalies by only training on the normal data. Although previous reconstruction-based methods have made fruitful progress, their generalization ability is limited due to two critical challenges. First, the training dataset only contains normal patterns, which limits the model generalization ability. Second, the feature representations learned by existing models often lack representativeness which hampers the ability to preserve the diversity of normal patterns. In this paper, we propose a novel approach called Adaptive Memory Network with Self-supervised Learning (AMSL) to address these challenges and enhance the generalization ability in unsupervised anomaly detection. Based on the convolutional autoencoder structure, AMSL incorporates a self-supervised learning module to learn general normal patterns and an adaptive memory fusion module to learn rich feature representations. Experiments on four public multivariate time series datasets demonstrate that AMSL significantly improves the performance compared to other state-of-the-art methods. Specifically, on the largest CAP sleep stage detection dataset with 900 million samples, AMSL outperforms the second-best baseline by 4%+ in both accuracy and F1 score. Apart from the enhanced generalization ability, AMSL is also more robust against input noise.

Index Terms—Unsupervised anomaly detection, Time series, Self-supervised learning, Memory network.

1 INTRODUCTION

WITH the prevalence of Internet of Things (IoT) devices in many applications (e.g., health care, human activity recognition and industrial control systems), there is growing interest in anomaly detection [1], [2], [3]. However, it is prohibitively expensive to acquire large amounts of labeled anomaly data required to train machine learning models following the current approaches. For instance, the collection of sleep apnea data [4] is extremely time-consuming and restricted by the environment, while the data for the normal states are much easier to obtain. For such scenarios, unsupervised anomaly detection is a promising paradigm that aims to detect anomalies by training from only normal samples, which is our main focus in this paper, and specifically, for multivariate time series.

The key challenge to unsupervised anomaly detection is to learn generalized patterns from the normal training data such that it can achieve good performance on the *unseen* anomalies. Over the years, there have been a line of works on this topic. Autoencoder (AE) is a powerful unsupervised learning technique that has been widely used

for unsupervised anomaly detection [5], [6], [7]. AE is usually trained by minimizing the reconstruction error on the normal data, and then using such errors as indicators or thresholds of anomalies. Based on the AE structure, LSTM-AE [8], Convolutional AE (CAE) [6], and ConvLSTM-AE [9] have emerged as promising anomaly detection approaches.

Despite the progress, the generalization ability of unsupervised anomaly detection methods is still limited due to two critical challenges. The first challenge is limited normal data. The normal training data are rather limited compared to the unseen test data, making the model prone to *overfitting*. As shown in Fig. 1, when the normal (Fig. 1(a)) and abnormal samples (Fig. 1(b)) are similar, existing methods can overfit to the anomalies. The second challenge is limited feature representations. For complex time series data, the feature representations learned by existing models lack representativeness required to preserve the diversity of normal patterns. This is shown in Fig. 1, when the testing data is becoming diverse to the normal samples (Fig. 1(c)), existing methods fail to capture the diverse patterns. Therefore, since the normal training data and feature representations are both limited, the generalization performance of the existing models is poor when applied to unseen normal and abnormal data.

In this paper, we propose a novel Adaptive Memory Network with Self-supervised Learning (AMSL) to increase the generalization ability of unsupervised anomaly detection by tackling the above two challenges. First, to cope with the limited normal training data, AMSL incorporates a self-supervised learning module to learn general patterns from the normal data. Second, to cope with the limited feature representations, AMSL introduces an adaptive memory fusion network to learn common and specific features via the global and local memory modules, respectively. Then, AMSL adopts an adaptive fusion module to fuse the

- Y. Zhang is with Global Energy Interconnection Development and Cooperation Organization, Xicheng District, Beijing, China, and Beijing Key Laboratory of Mobile Computing and Pervasive Device, Institute of Computing Technology, Chinese Academy of Sciences and University of Chinese Academy of Sciences, Beijing, China. E-mail: yuxin-zhang@geidco.org.
- Y. Chen are with Beijing Key Laboratory of Mobile Computing and Pervasive Device, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China and University of Chinese Academy of Sciences, Beijing, China. Y. Chen is also with Peng Cheng Laboratory (PCL). E-mail: yqchen@ict.ac.cn.
- J. Wang and T. Qin are with Microsoft Research, Beijing, China. E-mail: {jindong.wang, taoqin}@microsoft.com.
- H. Yu is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. E-mail: han.yu@ntu.edu.sg.
- Corresponding author: Yiqiang Chen and Jindong Wang

Manuscript received April 19, 2005; revised August 26, 2015.

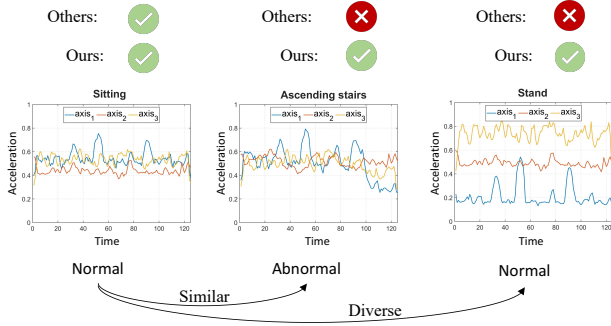


Fig. 1: Illustration of AMSL. Even with similar inputs ((a) and (b)), existing methods can overfit to the anomalies; On the other hand, when the inputs are diverse ((a) and (c)), our method remains accurate while existing methods can fail. *Axis* represents three dimensional signals on dataset.

global and local representations into the final representation, which is used for reconstruction. Based on the convolutional autoencoder framework, AMSL can be easily trained in an end-to-end manner. As shown in Fig. 1, while other methods fail in face of the two challenges, our AMSL can correctly perform anomaly detection.

This work makes the following three contributions:

- 1) We propose AMSL to tackle the limited normal data and feature representation challenges by adopting self-supervised learning and memory networks, respectively.
- 2) We propose to learn both the global and local memories to enhance the representation ability and further propose an adaptive memory fusion module to fuse the global and local memories for final representations.
- 3) Extensive experiments on four public datasets demonstrate the effectiveness of AMSL. Specifically on the largest CAP dataset [4] with 900M+ instances, AMSL significantly outperforms the best comparison approach by 4%+ in both accuracy and F1 score. Moreover, AMSL is also robust against noise and remain both time and memory efficient.

2 RELATED WORK

Unsupervised anomaly detection has been studied for decades. Traditional approaches include reconstruction methods (e.g., PCA, Kernel PCA [10], [11]), clustering methods (e.g., GMM, K-means [12], [13]), and one-class learning methods (e.g., OCSVM, SVDD [14], [15]). Deep learning methods are also popular. They can be categorized into the reconstruction-based and the prediction-based methods.

Reconstruction-based methods focus on reducing the expected reconstruction errors. For instance, Autoencoders (AEs) [5] are often utilized for anomaly detection by learning to reconstruct a given input. As the model is trained exclusively on normal data, whenever it is not able to reconstruct a given input with equal quality compared to the reconstruction of normal data, the instance is treated as an anomaly. The LSTM Encoder-Decoder model [8] is proposed to learn temporal representations of the time

series via the LSTM networks and use reconstruction errors to detect anomalies. Despite its effectiveness, LSTM does not take spatial correlation into consideration. Convolutional Autoencoder (CAE) [6], [7] is an important method for video anomaly detection. It is able to capture the 2D image structure since its weights are shared among all locations in the input image. Furthermore, since the Convolutional LSTM (ConvLSTM) can model spatio-temporal correlations by using the convolutional layers instead of the fully connected layers, later works [9], [16] add the ConvLSTM layers into the autoencoder, which encodes the change of appearance for normal data more effectiveness. Others like Variational Autoencoders (VAEs) [17], [18], De-noising AutoEncoders (DAEs) [19], Deep belief networks (DBNs) [20], Replicating neural network [21] and Robust Deep Autoencoder (RDA) [22] have also reported promising performance.

Prediction-based methods aim to predict one or more continuous values. In anomaly detection, researchers [23], [24], [25] propose RNN-based forecasting models to predict values of the next time period, and minimize the mean squared error between predicted and future values as the criterion for identifying anomalies. There have also been attempts to propose a forecasting model, which uses CNN and RNN, namely LSTMNet [26], to extract short-term local dependency patterns and long-term patterns for multivariate time series analysis for anomaly detection. The GAN-based method [27] adopts U-Net as the generator to predict the next frame in a video, and leverages adversarial training to distinguish if the predicted result fake (i.e. abnormal events) compared to the ground truth. However, these methods lack a reliable mechanism to learn representations of normal data at a fine-grained level.

Feature representation learning is one important aspect of deep learning and machine learning, and good representations of input data are essential for the generalization ability, interpretability and robustness of methods. Some recent works [28], [29] use image context and spatial-temporal relationships to well capture feature correlations for video object detection. **Self-supervised learning (SSL)** [30] is a type of unsupervised learning paradigm that uses the data itself to generate supervision to learn good representations. For instance, in an image classification problem where most of the images are unlabeled, we can still rotate it by different angles (e.g., 0° , 90° , 180° , 270°), and then use these angles as their labels. In this way, all samples are with labels. By learning on this multi-class classification problem (the auxiliary task), the model can learn general features from these images that can be used for classification (the main task) later. There are self-supervised techniques across computer vision [31], [32], natural language processing [33], and speech recognition tasks [34]. In anomaly detection, [35], [36] use self-supervised visual representation learning to learn the features of in-distribution (normal) samples.

The vanilla memory network is for question answering [37]. Generally, the RNNs or LSTM-based models capture long-term structure within sequences through local memory cells. Since the memory in these models is unstable over time, memory networks are proposed by using a global memory with shared read and write functions. Considering that memory can record information stably, several works adopt memory networks such as one-shot learning [38],

[39], neural machine translation [40], anomaly detection [41], [42], [43]. In anomaly detection, the memory module aims to record various patterns of normal data compared to items in the memory, thus distinguishing between normal and abnormal data.

3 PROPOSED APPROACH

3.1 Problem statement

Definition 1 (Multivariate time series). *A multivariate time series can be represented as $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T \in \mathbb{R}^{N \times V}$, where $\mathbf{x}_i \in \mathbb{R}^V$ is a signal with length V and N is the total number of signals.¹² $y \in \mathcal{Y}$ denotes the corresponding labels, where $\mathcal{Y} = \{1, \dots, K\}$ represents normal data with K classes.*

Definition 2 (Anomaly). *A sample (\mathbf{x}_a, y_a) is called an anomaly if its label y_a does not belong to any predefined classes, i.e. $y_a \notin \mathcal{Y}$.³*

In this paper, we deal with the most challenging unsupervised anomaly detection problem, where only unlabeled normal samples are available during training. This is a more realistic setting since large-scale collection of anomalous samples is often not practical.

3.2 Overview

We employ a convolutional autoencoder (CAE) as the base network, which is widely used in existing anomaly detection literature [5], [6], [7]. An autoencoder (AE) is an unsupervised neural network combining an encoder f_e and a decoder f_d parameterized by θ_e and θ_d , respectively. The encoder maps the high-dimensional input $\mathbf{x} \in \mathbb{R}^V$ into a latent representation $\mathbf{z} \in \mathbb{R}^F$, where $F \ll V$. Then, the decoder f_d reconstructs the original input as \mathbf{x}' by mapping \mathbf{z} back into the input space. The latent representation \mathbf{z} and the reconstructed input \mathbf{x}' can be respectively computed as:

$$\mathbf{z} = f_e(\mathbf{x}; \theta_e), \mathbf{x}' = f_d(\mathbf{z}; \theta_d). \quad (1)$$

The reconstruction error is calculated using the Mean Squared Error (MSE):

$$L_{MSE} = \|\mathbf{x} - \mathbf{x}'\|_2^2, \quad (2)$$

where $\|\cdot\|_2^2$ is the l_2 -norm.

In this paper, we propose a novel Adaptive Memory Network with Self-supervised Learning (AMSL) for unsupervised anomaly detection. AMSL consists of four novel components as shown in Fig. 2: 1) a self-supervised learning module, 2) a global memory module, 3) a local memory module and 4) an adaptive fusion module. AMSL works in the following four steps.

1. Each sensor signal may have different lengths due to different protocols. In this case, they can still be processed into the same length V by down-sampling or up-sampling.

2. Remark on the data dimension: In real applications, we often employ a sliding window technique to construct samples in time series, thus, \mathbf{x} is a tensor with >2 dimensions. For simplicity, we denote its flattened dimension as V .

3. Theoretically speaking, there is currently no unified definition of an anomaly. In this paper, we follow the common definitions in existing literature [9], [42], [44].

- 1) First, the encoder maps the raw time series signal and its six transformations into a latent feature space.
- 2) Then, for self-supervised learning, a multi-class classifier is built to classify these feature types in order to learn generalized feature representations.
- 3) Meanwhile, the features are also fed into the global and local memory network modules to learn both common and specific features.
- 4) Finally, the adaptive fusion module fuses these features to obtain new representation which is used for reconstruction.

Details of AMSL are presented in the following sections.

3.3 Self-Supervised learning

In this section, we introduce the self-supervised learning module of AMSL that enables generalized feature representation learning for normal data. Compared to the potentially large number of unseen anomalies that could be in any forms of representations, the number of the normal training data is relatively limited. Hence, anomaly detection models trained on such limited normal samples tend to overfit. However, it is prohibitively expensive to collect all the infinite training data. To address this issue, we propose to use self-supervised learning to increase the model's generalization ability.

Assuming that the instances are consistent before and after basic data augmentations (i.e., normal data and abnormal data are still distinguishable after the same augmentations) [30], we design feature transformations on the original data for self-supervision. Then, we train the model to recognize the transformation type of a sample as its auxiliary task. In the following experiments, as shown in Fig. 3(a), we can observe that the instances after the same augmentations can still identify normal and abnormal data. Specifically, we propose to utilize six signal transformations inspired by [45], which are described as follows:

- 1) *Noise*: As noisy sensor signal may exist in the real world, adding noise to signal can help model learn more robust features against the noise. Here, the transformation with Gaussian noise is implemented.
- 2) *Reverse*: This transformation reverses the samples along the temporal dimension, resulting in the samples with the opposite time direction.
- 3) *Permute*: This transformation randomly perturbs signal along the time dimension by slicing and swapping different time windows to generate new samples. It aims to enhance the permutation invariant property of the resulting model.
- 4) *Scale*: Scaling changes the magnitude of signals within a time window by multiplying a random scalar. Here, we select $[0.5, 0.8, 1.5, 2]$ as scalar values. The addition of scaled signals can help the model learn scale-invariant patterns.
- 5) *Negate*: This transformation is a special type of scaling transformations. It is scaled by -1, resulting in a mirror image of the input signals.
- 6) *Smooth*: This transformation applies the Savitzky-Golay (SG) method to smooth the signals. The

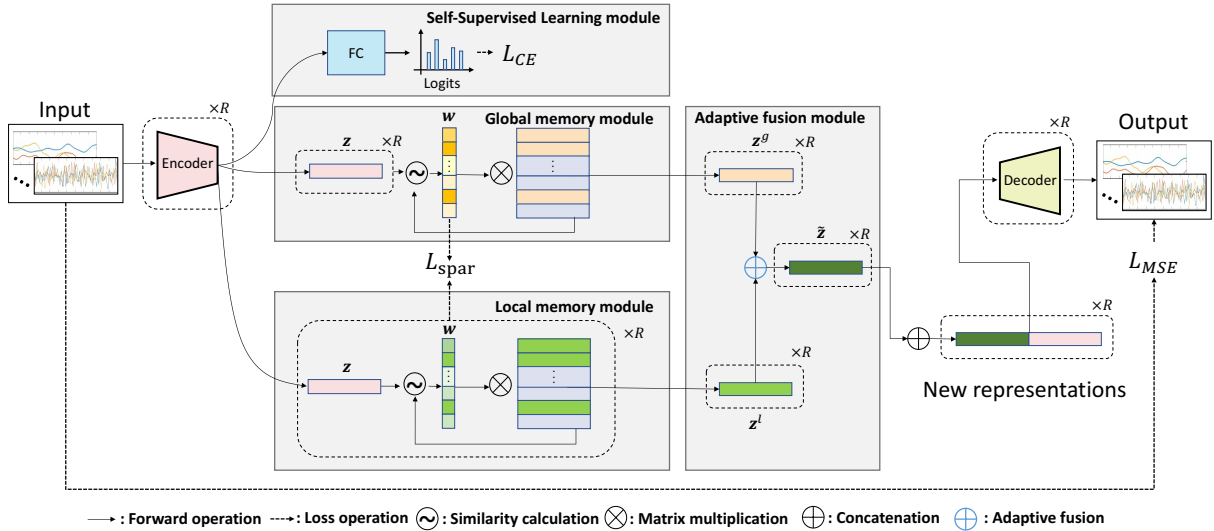


Fig. 2: The structure of the proposed AMSL. It consists of four components: self-supervised learning, global memory, local memory and adaptive fusion. The notation “ $\times R$ ” denotes R copies where each one corresponds to one transformation.

Savitzky-Golay filter is a particular type of low-pass filters, well adapted for noisy signal smoothing.

To learn general feature representations from these transformations, the model is enforced to distinguish their transformation types by using the Cross-Entropy loss function:

$$L_{CE} = - \sum_{i=1}^R y_i \log(p_i), \quad (3)$$

where R denotes the number of self-supervised learning classes ($R = 7$ in this work by including all six transformations and the original signal). y_i and p_i are the pseudo label and the predicted probability, respectively. A Softmax activation function is applied to the probability before the cross-entropy loss.

3.4 Adaptive Memory Fusion Module

Traditional AE is negatively affected by noisy or unknown training data, so it can also consistently reconstruct abnormal inputs too well [42], [46]. Therefore, the model fails to learn representative features. To tackle this challenge, we propose a *adaptive memory fusion module* to enhance the ability of the model in distinguishing between normal and abnormal data via recording the prototypical patterns. In the following, we will first introduce the memory network. Then, we present our adaptive memory fusion module in detail.

3.4.1 Memory module

The memory module [41], [42] consists of a memory representation to represent the encoded patterns, and a memory updating part to update the memory items based on the similarity of the memory items and the input z . Specifically, the memory is instantiated as a matrix $\mathbf{M} \in \mathbb{R}^{C \times F}$, storing C vectors of dimension F . Let the row vector $\mathbf{m}_i, \forall i \in [C]$ denote the i -th row of \mathbf{M} , where $[C]$ denotes the set of integers from 1 to N . Each \mathbf{m}_i denotes a memory item.

Given a query $z \in \mathbb{R}^F$ (i.e. encoding), where F denotes the filter size in the last layer of the encoder, the memory-guided module outputs \hat{z} , which represents a weighted average of memory items \mathbf{m}_i , by matching probabilities $w \in \mathbb{R}^C$ between the query z and memory matrix \mathbf{M} as follows:

$$\hat{z} = \mathbf{w}\mathbf{M} = \sum_{i=1}^C w_i \mathbf{m}_i, \quad (4)$$

where w_i denotes the i -th entry of \mathbf{w} . The weight vector \mathbf{w} is computed based on the normalized similarity between the query z and memory entry \mathbf{m}_i :

$$w_i = \frac{\exp(\text{Score}(z, \mathbf{m}_i))}{\sum_{j=1}^C \exp(\text{Score}(z, \mathbf{m}_j))}. \quad (5)$$

The Score function is implemented as the cosine similarity:

$$\text{Score}(z, \mathbf{m}_i) = \frac{z \cdot \mathbf{m}_i^\top}{\|z\| \|\mathbf{m}_i\|}. \quad (6)$$

In the training phase, the memory matrix can be updated by the reconstruction loss function, thus forcing it to record the normal characteristics. In the testing phase, the memory network outputs a representation with a combination of all items, taking into account the multiple patterns of normal characteristics. Therefore, the normal instances can be reconstructed well. The anomalies reconstructed using the retrieved normal patterns in the memory module will result in higher reconstruction errors.

3.4.2 Adaptive fusion module

We further propose an adaptive memory fusion network to learn both the common and specific representations from all the feature augmentations. Specifically, we propose the *global memory module* to learn the common representations contained in all transformations, and the *local memory module* to learn augmentation-specific representations for each transformation. Finally, we propose a *adaptive fusion*

module to fuse these two levels of features into the final representations that will be used for reconstruction. The motivation is that we can capture the common patterns for the normal data and its specific information (i.e. each different transformation) that is useful for the normal data patterns, thus improving feature representations of normal data at a fine-grained level.

We build the global memory module with a shared memory matrix. By using the encoded representation as a query, the global memory module can record general items in the memory matrix. Through a shared memory module, the outputs are obtained as:

$$\mathbf{z}_i^g = f_g(f_e(\mathbf{x}_i; \theta_e); \theta_g), i \in [R], \quad (7)$$

where $f_g(\cdot)$ is the function of the global memory module. $[R] = \{1, 2, \dots, R\}$ and θ_g denote the shared parameters of one global memory module.

We build R local memory modules for the raw data and the six transformations. Each memory matrix records the corresponding normal characteristics of transformation. These outputs are obtained by the local memory modules as:

$$\mathbf{z}_i^l = f_l(f_e(\mathbf{x}_i; \theta_e); \theta_l^i), i \in [R], \quad (8)$$

where $f_l(\cdot)$ is the function of the local memory module, and θ_l^i denotes the parameters of seven local memory modules.

It is intuitive that the common and specific features are not equally important in representing a given instance. To adaptively fuse these features, we use a feed-forward layer that takes the features and a free variable $r \in \mathbb{R}$ as inputs to generate the fused representations by weights $\alpha \in \mathbb{R}^{2R}$ (2 weights for local and global memories, with R transformations in total). Note that we use Batch Normalization and sigmoid activation to normalize the weights and control their values to be within the range of $(0, 1)$. r is used to increase randomness. Then, we get the adaptive fused representations:

$$\tilde{\mathbf{z}}_i = [\alpha_i^g \quad \alpha_i^l] \begin{bmatrix} \mathbf{z}_i^g \\ \mathbf{z}_i^l \end{bmatrix}, i \in [R], \quad (9)$$

where $\alpha_i^g, \alpha_i^l \in \mathbb{R}$ denote the weights of common (global) and specific (local) features, respectively.

The decoders concatenate \mathbf{z} (encoding output) and $\tilde{\mathbf{z}}$ (adaptive fusion output) as inputs to reconstruct the original inputs. The reconstruction loss is defined by minimizing the l_2 distance between the decoder output and the original input as follows:

$$L_{MSE} = \sum_{i=1}^R \|f_d(\text{concat}(\tilde{\mathbf{z}}_i, \mathbf{z}_i); \theta_d^i) - \mathbf{x}_i\|_2^2. \quad (10)$$

To constrain the sparsity of the memory weight \mathbf{w} to avoid over-reconstruction of anomalies by the complex combination of memory items, we adopt sparse loss by minimizing the entropy of \mathbf{w} :

$$L_{spar} = \sum_{i=1}^C -w_i \cdot \log(w_i). \quad (11)$$

Algorithm 1 Learning procedure of AMSL

Input: Normal Dataset $\mathcal{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_H\}$, hyperparameters λ_1, λ_2 , variable r .

Output: Decision threshold μ and model parameter θ .

- 1: Transform each sample into $\mathbf{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^R\} \in \mathbb{R}^{R \times V \times N}$, where $R = 7$ including raw data and 6 transformations;
- 2: Initialize all parameters;
- 3: **while** not converge **do**
- 4: $\mathbf{Z} \leftarrow \text{Encoder}(\mathbf{X})$;
- 5: $p \leftarrow FC(\mathbf{Z})$;
- 6: $\mathbf{Z}^l, \mathbf{w} \leftarrow LMem(\mathbf{Z})$;
- 7: $\mathbf{Z}^g, \mathbf{w} \leftarrow GMem(\mathbf{Z})$;
- 8: Learn adaptive fuse weights α ;
- 9: $\tilde{\mathbf{Z}} \leftarrow Fusion(\mathbf{Z}^l, \mathbf{Z}^g; \alpha)$;
- 10: $\mathbf{X}' \leftarrow Decoder(Concat(\tilde{\mathbf{Z}}, \mathbf{Z}))$;
- 11: $J(\theta) = \frac{1}{H} \sum_{i=1}^H (\sum_{j=1}^R L(\mathbf{X}'_{ij}, \mathbf{X}_{ij}) - \lambda_1 \sum_{j=1}^R y_{ij} \log(p_{ij}) - \lambda_2 \sum_{j=1}^{2R} \sum_{k=1}^C w_{ijk} \cdot \log(w_{ijk}))$;
//Implementation of Eq. (12) for multivariate time series
- 12: Update model parameter θ using Eq. (12);
- 13: **end while**
- 14: Calculate the decision threshold μ by the training samples;
- 15: **return** Optimal μ and θ .

3.5 Training and Inference

3.5.1 Training

By integrating the reconstruction loss in Eq. (10), the sparse loss in Eq. (11), and the self-supervised loss in Eq. (3) with tradeoff parameters λ_1, λ_2 , we obtain the overall training objective for AMSL which can be optimized in an end-to-end manner:

$$J(\theta) = L_{MSE} + \lambda_1 L_{CE} + \lambda_2 L_{spar}. \quad (12)$$

The training procedure of AMSL on multivariate time series is shown in Algorithm 1. Note that we need to take of the multiple dimensions of the inputs in real implementation.

3.5.2 Inference

For the autoencoder-based models, it is generally assumed that the compression is different for instances of different categories. That is, if the training dataset only contains normal instances, the reconstruction error becomes higher for abnormal instances. Therefore, we can classify these instances into ‘‘abnormal’’ or ‘‘normal’’ by their reconstruction error during inference phase.

Given normal samples as the training dataset $\mathcal{X} = \{\mathbf{X}_1, \dots, \mathbf{X}_H\}$. First, we need to construct the self-supervised transformations. Let the corresponding decision threshold μ be the 99th percentile of $\text{Err}(\mathbf{X}_i)$ on this training set, where $\text{Err}(\mathbf{X}_i)$ is the value of reconstruction loss function L_{MSE} for \mathbf{X}_i . In inference (or detection) process, the decision rule is that if $\text{Err}(\mathbf{X}_i) > \mu$, the testing sample \mathbf{X}_i in a sequence will be classified as ‘‘abnormal’’; otherwise, it will be classified as ‘‘normal’’. Here, LMem and GMem denote local memory module and global memory module,

Algorithm 2 Inference procedure of AMSL

Input: Testing dataset $\mathcal{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_H\}$, $\mathbf{X} \in \mathbb{R}^{V \times N}$, threshold μ , model parameter θ , hyperparameters λ_1, λ_2 .

Output: Label y_i of \mathbf{X}_i .

- 1: Transform each sample into $\mathbf{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^R\} \in \mathbb{R}^{R \times V \times N}$, where $R = 7$ including raw data and 6 transformations;
- 2: Load the model trained parameters θ ;
- 3: $\mathbf{Z} \leftarrow \text{Encoder}(\mathbf{X})$;
- 4: $\mathbf{Z}^l \leftarrow \text{LMem}(\mathbf{Z})$;
- 5: $\mathbf{Z}^g \leftarrow \text{GMem}(\mathbf{Z})$;
- 6: Load α ;
- 7: $\tilde{\mathbf{Z}} \leftarrow \text{Fusion}(\mathbf{Z}^l, \mathbf{Z}^g; \alpha)$;
- 8: $\mathbf{X}' \leftarrow \text{Decoder}(\text{Concat}(\tilde{\mathbf{Z}}, \mathbf{Z}))$;
- 9: Calculate reconstruction loss $\text{Err}(\mathbf{X}_i) = \sum_{j=1}^R L(\mathbf{X}'_{ij}, \mathbf{X}_{ij})$;
- 10: **if** $\text{Err}(\mathbf{X}_i) > \mu$ **then**
- 11: $y_i = \text{"Abnormal"}$;
- 12: **else**
- 13: $y_i = \text{"Normal"}$;
- 14: **end if**
- 15: **return** y_i .

respectively. The inference procedure of AMSL is shown in Algorithm 2.

4 EXPERIMENTAL EVALUATION

4.1 Datasets

In the experiments, we adopt the following four datasets for evaluation as shown in TABLE 1.

TABLE 1: The detailed statistics of the four datasets

Dataset	Application	#Instance	#Dim	#Class
DSADS [47]	Activity Recognition	1,140,000	45	19
PAMAP2 [48]	Activity Recognition	2,844,868	36	18
WESAD [49]	Stress & Affect Detection	63,000,000	14	3
CAP [4]	Sleep Stage Detection	921,700,000	21	8

DSADS [47] consists of motion sensor data (i.e. accelerometer, gyroscope, magnetometer) of 19 activities of daily living and sports activities performed by 8 subjects. To simulate normal and abnormal classes, we chose running, ascending stairs, descending stairs, rope jumping and playing basketball as anomaly classes, while the remaining categories are regarded as normal classes. These designated anomaly activities are relatively intense and rare compared to other activities. In TABLE 2(a), we describe the selection of normal and abnormal classes on DSADS datasets. We chose these activities as the abnormal classes that are relatively intense and rare activities compared to other activities, while the rest categories are defined as the normal classes.

PAMAP2 [48] is a mobile dataset containing data of 18 different physical activities performed by 9 subjects wearing 3 inertial measurement units comprising an accelerator, a gyroscope and a magnetometer. As shown in TABLE 2(b), we treat the classes with relatively small samples as the

TABLE 2: The selection of normal and abnormal classes on DSADS and PAMAP2 datasets

(a) DSADS	
Activity	Class
sitting, standing, lying, moving, walking, cycling, exercising, rowing	normal
running, ascending stairs, descending stairs, rope jumping, playing basketball	abnormal
(b) PAMAP2	
Activity	Class
lying, sitting, standing, walking, cycling, Nordic walking, vacuum cleaning, ironing	normal
running, ascending stairs, descending stairs, rope jumping	abnormal

anomaly classes, including running, ascending stairs, descending stairs and rope jumping. The remaining categories are used as the normal classes in experiments.

WESAD [49] is a public dataset for wearable stress and affect detection, which contains physiological and motion data of 15 subjects. We use sensor modalities of chest-worn device including electrocardiogram, electrodermal activity, electromyogram, respiration, body temperature, and three-axis acceleration. In our experiments, we chose normal emotional states (neutral, amusement) as normal classes and stress states as anomaly class.

CAP Sleep Database [4], which stands for the Cyclic Alternating Pattern (CAP) database. It is a clinical dataset from PhysioNet repository [55]. It is characterized by periodic physiological signals occurring during wake, S1-S4 and REM sleep stages. The signals include EEG, EOG, EMG, EKG and respiration. There are 16 healthy subjects and 92 patients in the database. In this task, we extracted 7 valid channels (ROC-LOC, C4-P4, C4-A1, F4-C4, P4-O2, ECG1-ECG2 and EMG1-EMG2). For detecting sleep apnea events, we chose healthy subjects as the normal class and the patients with sleep-disordered breathing as the anomaly class.

4.2 Comparison Methods

We compare AMSL against four popular traditional anomaly detection methods, which can be applied after feature extraction:

- KPCA [11], which is a non-linear extension of PCA commonly used for anomaly detection. We adopt Gaussian kernel in our experiments.
- ABOD [50], which uses k nearest neighbors to approximate the complexity reduction. For an observation, the variance of its weighted cosine scores to all neighbors could be viewed as the abnormal score.
- OCSVM [15], which adopts PCA for dimension reduction and employs the Gaussian kernel with a bandwidth of 0.1.
- HMM [51], which is applied after extracting features and then it calculates the anomaly probability from the state sequence generated by the model.

We also compare it against seven deep unsupervised methods:

TABLE 3: The comparison of mean precision, recall, F1 and accuracy of AMSL and other baselines. The best and second-best results are bold and underlined, respectively. We can see that AMSL significantly outperforms other methods.

Method	DSADS dataset				PAMAP2 dataset				WESAD dataset				CAP dataset			
	mPre	mRec	mF1	Acc	mPre	mRec	mF1	Acc	mPre	mRec	mF1	Acc	mPre	mRec	mF1	Acc
Kernel PCA [11]	0.6184	0.6182	0.6183	0.6186	0.7236	0.6579	0.6892	0.5645	0.5496	0.5495	0.5495	0.5486	0.7603	0.5847	0.6611	0.5892
ABOD [50]	0.6880	0.6510	0.6690	0.6554	0.8653	0.9022	0.8834	0.8985	0.8782	0.8786	0.8784	0.8783	0.7867	0.6365	0.7037	0.6326
OCSVM [15]	0.7608	0.7277	0.7439	0.7312	0.7600	0.7204	0.7397	0.7679	0.6092	0.5631	0.5852	0.5518	<u>0.9267</u>	<u>0.9259</u>	<u>0.9263</u>	<u>0.9257</u>
HMM [51]	0.6959	0.6917	0.6937	0.6901	0.6950	0.6553	0.6745	0.5725	0.6123	0.6060	0.6097	0.6018	0.8238	0.8078	0.8157	0.8090
CNN-LSTM [52]	0.6845	0.6270	0.6545	0.6425	0.6680	0.5392	0.5968	0.6131	0.5883	0.5440	0.5653	0.5318	0.6159	0.5217	0.5649	0.5762
LSTM-AE [8]	0.8471	0.7729	0.8083	0.7852	0.8619	0.7997	0.8296	0.8426	0.2353	0.4762	0.3150	0.4599	0.7147	0.6253	0.6671	0.6286
MSCRED [9]	0.7540	0.5602	0.6428	0.6192	0.6997	0.7301	0.7146	0.7517	0.8850	0.8124	0.8471	0.8420	0.6410	0.5784	0.6081	0.5819
ConvLSTM-AE [16]	0.8164	0.6951	0.7509	0.7121	0.7359	0.7361	0.7360	0.7341	0.9733	0.9698	0.9716	0.9709	0.8150	0.8194	0.8172	0.8165
ConvLSTM-COMP [16]	0.8229	0.7379	0.7781	0.7518	<u>0.8844</u>	<u>0.8842</u>	<u>0.8843</u>	<u>0.8844</u>	0.9626	0.9629	0.9627	0.9619	0.8367	0.8377	0.8372	0.8394
BeatGAN [53]	0.9517	0.5663	0.7100	0.7818	0.7981	0.7420	0.7691	0.8369	0.7586	0.5000	0.6027	0.5172	0.5251	0.5002	0.5123	0.8437
MNAD-P [43]	0.5816	0.5783	0.5799	0.5721	0.8198	0.8176	0.8186	0.8135	0.7600	0.6938	0.7254	0.6849	0.7742	0.7489	0.7613	0.6960
MNAD-R [43]	0.8337	0.7694	0.8003	0.7811	0.8350	0.8355	0.8353	0.8334	0.7426	0.6677	0.7031	0.6579	0.8189	0.8235	0.8212	0.7871
GDN [54]	0.8706	0.8151	0.8419	0.8251	0.8129	0.8104	0.8116	0.8123	0.7520	0.5434	0.6309	0.5590	0.6831	0.6237	0.6520	0.6569
UODA [23]	0.8679	0.8281	0.8475	0.8365	0.8957	0.8513	0.8730	0.8823	0.7623	0.6314	0.6907	0.6191	0.7557	0.5124	0.6107	0.5173
AMSL (Ours)	0.9407	0.9298	0.9352	0.9332	0.9788	0.9713	0.9750	0.9770	0.9953	0.9949	0.9951	0.9951	0.9771	0.9736	0.9753	0.9756
Improvement	7.28%	10.17%	8.77%	9.67%	9.44%	8.71%	9.07%	9.26%	2.20%	2.51%	2.35%	2.42%	5.04%	4.77%	4.90%	4.99%

- *CNN-LSTM* [52], which is a prediction-based method built by firstly defining a network comprised of Conv2D and MaxPooling2D layers ordered into a stack of the required depth. The result is then fed into the LSTM and FC layers as prediction.
- *LSTM-AE* [8], which is a reconstruction-based method using a single-layer LSTM on both encoder and decoder.
- *MSCRED* [9], which is an encoder-decoder model with multi-scale matrices as inputs for multivariate time series analysis.
- *ConvLSTM-COMPOSITE* [16], which is a composite encoder-decoder model with reconstruction and prediction task. We choose the “conditional” version to build a single model called *ConvLSTM-AE* by removing the forecasting decoder.
- *BeatGAN* [53], which is a reconstruction-based method with an adversarial generation approach as regularization.
- *MNAD* [43], which is an encoder-decoder model based on a memory module for video anomaly detection. It has two variants: one with the prediction task (*MNAD-P*), another with the reconstruction task (*MNAD-R*).
- *GDN* [54], which is a Graph-based neural network to learn a graph of the dependence relationships between sensors for anomaly detection.
- *UODA* [23], which is a RNN-based network for anomaly detection. We re-implemented it by customizing the number of layers and hyper-parameters.

We re-implement the comparison methods based on several open-source repositories⁴ as well as our own implementations. As for ours, there is a data pre-processing stage where data is normalized, split into windows with length V and transformed. The encoder runs with *Conv2D* \rightarrow *Maxpool* \rightarrow *Conv2D* \rightarrow *Maxpool*, i.e., Conv1- Conv2 with 32 and 64 kernels of size 4×4 , and Maxpooling with size 2×2 . In adaptive memory fusion module, $F = 64$ and $C = 800$ are the best choice. Besides, we use a variable r as the initial weight of adaptive fusion network with $FC(2R)$

\rightarrow *BN* \rightarrow *Sigmoid* \rightarrow *Multiply*. These weights are multiplied into local and global feature representations by Eq. (9). As the input of the decoder is by concatenating the output of encoder and the output of memory module, the decoder is asymmetric to the encoder that uses 4 *Conv2DTranspose* modules, with $\{128, 64, 32, 1\}$ kernels of size 4×4 in each layer, respectively. To compute classification error, the output of encoder is also incorporated into classification network with *Conv2D* \rightarrow *Flatten* \rightarrow *FC(128)* \rightarrow *Dropout* \rightarrow *FC(R)*, i.e. Conv with 1 kernel of size 4×4 . AMSL is trained in an end-to-end fashion using Keras [56] on a TITAN XP GPU. The Adam optimizer is used to train model for about 100 epochs with 32 or 64 batch size. The learning rate is 0.001. And we set the hyper-parameters: $\lambda_1 = 1$ and $\lambda_2 = 0.0002$ while parameter sensitivity analysis is presented in later sections.

In practice, it is difficult to know the ground truth, and anomalous data points are rare. Hence, the semi-supervised setting is a commonly used as evaluation method [42], [46], [57], the training set only contains the normal samples and has no overlapping with the testing set. For each dataset, we split the normal samples into training, validation, and testing sets with the ratio of 5 : 1 : 4. The model selection criterion, i.e., hyperparameters, used for tuning is the validation error on the validation set. Besides, since most datasets have more normal samples than anomalies, accuracy along is insufficient for evaluation. Thus, for comprehensive evaluation, we adopt four evaluation metrics: mean *precision*, *recall*, *F1* score, and *accuracy* following existing literature [23], [46], [57].

4.3 Results and Analysis

TABLE 3 reports the overall performance results on these public datasets. It can be observed that the proposed AMSL method achieves significantly superior performance over the baseline methods in all the datasets. Specifically, compared with other methods, AMSL significantly improves the F1 score by **9.07%** on PAMAP2 dataset, **4.90%** on CAP dataset, **8.77%** on DSADS dataset and **2.35%** on WESAD dataset. The same pattern goes for precision and recall. Especially for the largest CAP dataset with over 900 Million samples, AMSL dramatically outperforms the second-best baseline (OCSVM) with an F1 score of **4.90%**, indicating its effectiveness.

4. <https://pyod.readthedocs.io/en/stable/>, <https://github.com/7fantasysz/MSCRED>, <https://github.com/cvlab-yonsei/MNAD>, <https://github.com/Vniex/BeatGAN>, <https://github.com/d-ailin/GDN>

For DSADS, PAMAP2 and CAP datasets (WESAD is relatively easier to train than the other datasets), we find that as the datasets become larger, the improvements tend to decrease. This means that self-supervision is more useful on small-scale datasets where it is hard to learn generalized representations. This is in accordance with existing machine learning conclusions. In addition, even if the samples are relatively fewer in DSADS and PAMAP2 datasets but with more categories, our AMSL still significantly outperforms other methods by large margins, which indicate its ability of dealing with the diversities in the limited training data.

Traditional methods perform differently on different datasets since they are limited by the feature extraction methods and the curse of dimensionality. In deep learning methods, CNN-LSTM achieves the lowest F1 score, which means CNN-LSTM alone is not enough to capture general patterns and more regularizations are needed to improve its performance. For the reconstruction-based methods (LSTM-AE, MSCRED, and ConvLSTM-AE), their performances are limited by the noise in the training data, which is likely to mislead the model by making it hard to distinguish between the normal, abnormal, and noisy data. MNAD and ConvLSTM are proposed for video data, which may not be suitable for multivariate time series. For BeatGAN, it has poor performance on CAP and WESAD datasets since it is prone to mode collapse and convergence problems on large-scale datasets. UODA performs reasonably well on the PAMAP2 and DSADS datasets that relies on pre-training denoising autoencoder (DAEs) and deep recurrent networks (RNNs) before fine-tuning. In fact, its performance is still not optimal because RNN does not have the capability to memorize long-term time series. GDN has a good performance on PAMAP2 and DSADS datasets. However, with the increase of data, the running speed of the model slows down and the accuracy also decreases affected by the graph structure. For ConvLSTM-COMPOSITE model, it performs better than most baselines. However, its efficiency may be limited due to the existence of two decoders in its structure.

Besides, we make confusion matrices to error analysis for our proposed method AMSL. As shown in Fig. 4, we find that the proportion of misclassification of normal data is lower than that of abnormal data for most datasets. The F1-scores of normal class and abnormal class are 93.99% and 92.48% on DSADS dataset, the F1-scores of normal class and abnormal class are 98.23% and 96.74% on PAMAP2 dataset, the F1-scores of normal class and abnormal class are 99.49% and 99.52% on WESAD dataset, the F1-scores of normal class and abnormal class are 97.85% and 97.18% on CAP dataset. The proposed AMSL method achieves significantly superior performance in all the datasets, with an F1 score of at least 93%.

4.4 Ablation Study

In this section, we conduct ablation study to show the effectiveness of the self-supervised learning, memory and adaptive fusion modules of AMSL. All experiments were performed using the PAMAP2 dataset. TABLE 4 reports the performance of different variants of AMSL. These results show that the proposed self-supervised learning module and memory module can both increase the performance of

TABLE 4: Ablation results on the PAMAP2 dataset. “SSL” denotes self-supervised learning module.

Variant	mPre	mRec	mF1	Acc
CAE	0.8728	0.8711	0.8719	0.8676
CAE + Mem	0.8808	0.8819	0.8813	0.8809
CAE + SSL	0.9393	0.9393	0.9393	0.9364
CAE + SSL + Mem	0.9616	0.9420	0.9517	0.9554
CAE + SSL + Ada Mem	0.9788	0.9713	0.9750	0.9770
CAE + Mem w/o Spar	0.8713	0.8723	0.8719	0.8718
CAE + SSL w/o SSL loss	0.8231	0.8175	0.8203	0.8205

the base convolutional autoencoder (CAE) module. After combining these two modules, it can be observed that AMSL can achieve better performance than that without the adaptive fusion module. Here, “Mem” denotes memory fusion module. This means that z_i^g and z_i^l are fused in a proportion of 1 : 1.

AMSL also consists of other components, such as self-supervised loss and sparse loss. To demonstrate the effectiveness of these components, we conduct other ablation studies shown in TABLE4 (last two rows). It is worth noting that CAE + SSL without self-supervised loss (i.e., only using different augmentations) obtain a lower F1 score than other variants of the approach. This is because the augmented data (via different transformations) complicates the distribution of normal data, thus leading to the under-fitting problem. Therefore, it is necessary to classify these transformations using pseudo labels. Similarly, we observe that adding sparse loss further improves the performance.

4.5 Detailed Analysis

4.5.1 Self-Supervised Learning

Self-supervised learning helps the network learn general and diverse features for the normal data, thus increasing the generalization ability of the model to recognize the unseen normal and abnormal instances. In Fig. 3(a), we show comparative performance analysis of each self-supervised data transformation. This assessment helps us understand whether model performance via jointly learning augmented data is better than learning individual data. The experiments were performed using the PAMAP2 dataset. The results show that the overall performance is competitive except for the noisy signals, thus it is beneficial to combine all transformations for better generalization. In Section 4.10, we discard the poorly performing transformations such as “Noise” and “Scale”. We observe that the F1 score and accuracy of AMSL decrease as transformations R decrease. This demonstrate that it’s helpful to use different self-supervised data transformations to increase the model’s generalization ability.

4.5.2 Adaptive Memory Fusion module

In Fig. 3(b), we show the performance of CAE, GMSL, LMSL and AMSL. GMSL and LMSL are global memory network and local memory network with SSL. The experimental results demonstrate that adaptive memory fusion network achieves better performance than using single memory network (i.e. global and local memory module). AMSL achieves a **1.86%** higher F1-score than GMSL, and a **3.12%** higher F1-score than LMSL. The same goes for the precision and recall.

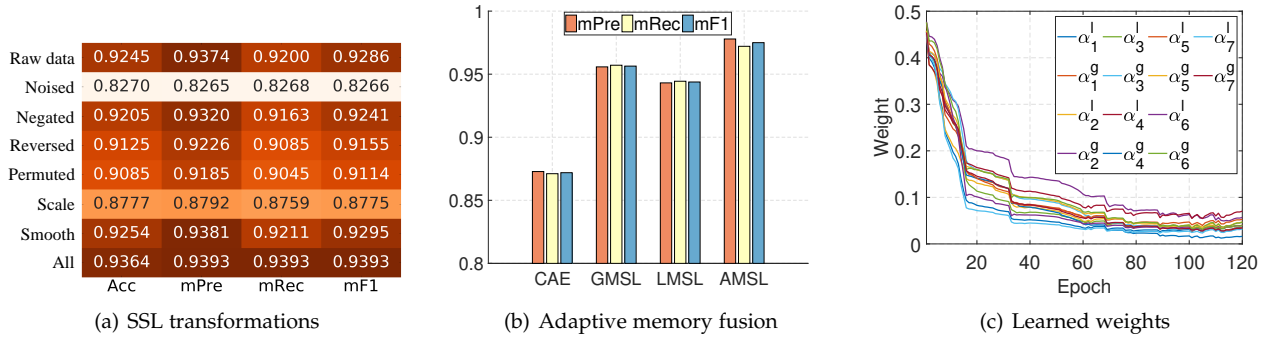


Fig. 3: Analysis of AMSL from several key aspects on PAMAP2 dataset. (a) The effectiveness of each data transformation in self-supervised learning. (b) Comparison between local, global, and our adaptive memory fusion modules. (c) The weights learned by our adaptive memory fusion module.

TABLE 5: Ablation results on all these datasets.

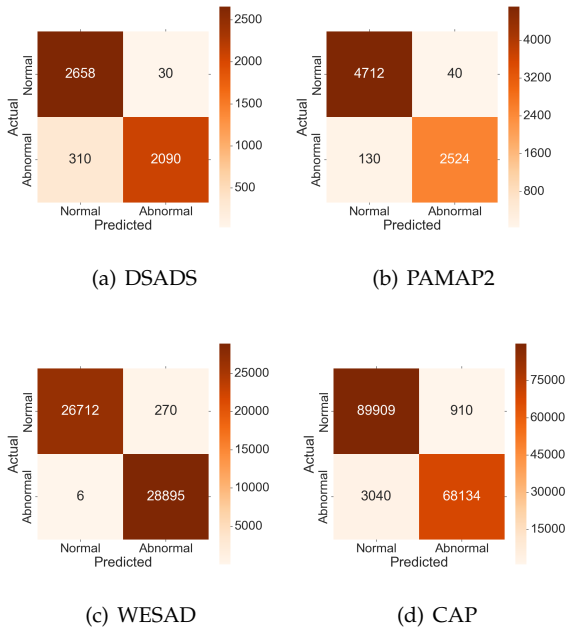


Fig. 4: The confusion matrices on four datasets.

More results on different datasets are shown in TABLE 5. Our method AMSL is always superior than GMSL and LMSL. As mentioned above, AMSL finds the best weight α^l and α^g values automatically. To observe the changes of adaptive weights at the training stage, this experiment was conducted using the PAMAP2 dataset. As shown in Fig. 3(c), it occurred at around the 70-th epoch when the values of adaptive weights stabilized. Here, $\alpha_1 - \alpha_7$ denote the weights of raw data and 6 transformations corresponding to transformations in Fig. 3(a).

4.6 Robustness to Noisy Data

In real-world applications, the collection of multivariate time series data can be easily polluted with noise due to changes in the environment or the data collection devices. The noisy data bring critical challenges to the unsupervised anomaly detection methods. In this section, we evaluate the

robustness of different methods to noisy data on PAMAP2 dataset. We manually control the noisy data ratio in the training data. We inject Gaussian noise ($\mu = 0, \sigma = 0.3$) in a random selection of samples with a ratio varying between 1% to 30%. We compare the performance of three methods: UODA, ConvLSTM-Composite, and AMSL in Fig. 6(c). As the noise increases, the performance of all methods decreases. Among them, AMSL remains significantly superior to others, demonstrating its robustness to noisy data. The result of precision and recall are presented in Fig. 6(a) and Fig. 6(b).

4.7 Percentage of Anomaly

Generally, the percentage of anomaly is significantly lower than the normal range. Hence, we conduct experiments on CAP dataset when the percentage of anomaly on the testing set is 1%, 5%, 10%, 15%, 20%, 25% and 30%. In

(a) DSADS

Variant	mPre	mRec	mF1	Acc
GMSL	0.8890	0.8551	0.8717	0.8624
LMSL	0.8915	0.8523	0.8716	0.8605
AMSL	0.9407	0.9298	0.9352	0.9332

(b) PAMAP2

Variant	mPre	mRec	mF1	Acc
GMSL	0.9558	0.9571	0.9564	0.9553
LMSL	0.9431	0.9444	0.9438	0.9433
AMSL	0.9788	0.9713	0.9750	0.9770

(c) WESAD

Variant	mPre	mRec	mF1	Acc
GMSL	0.9863	0.9849	0.9856	0.9851
LMSL	0.9864	0.9850	0.9857	0.9855
AMSL	0.9953	0.9949	0.9951	0.9951

(d) CAP

Variant	mPre	mRec	mF1	Acc
GMSL	0.9638	0.9556	0.9597	0.9598
LMSL	0.9620	0.9531	0.9575	0.9575
AMSL	0.9771	0.9736	0.9753	0.9756

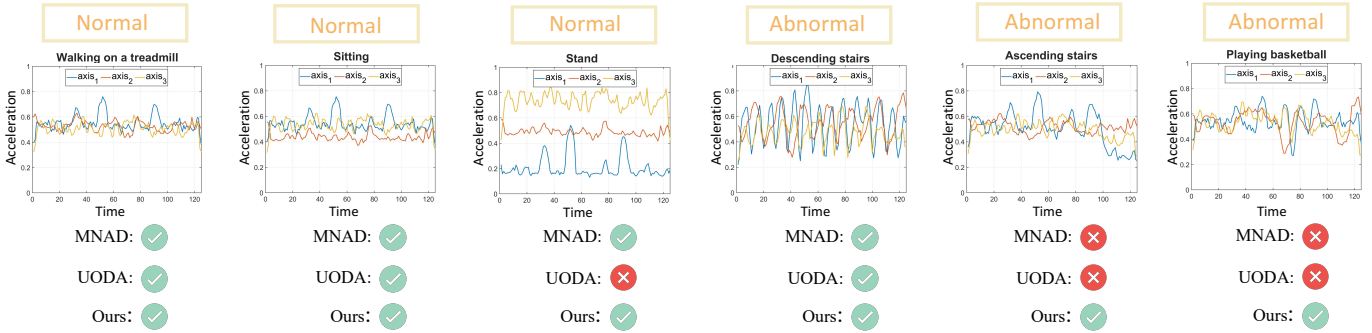


Fig. 5: The visualization for samples in normal and abnormal class. Each column is the instance whether is correctly detected by our method AMSL, MNAD and UODA. $Axis_{1,2,3}$ represents three dimensional signals on DSADS dataset.

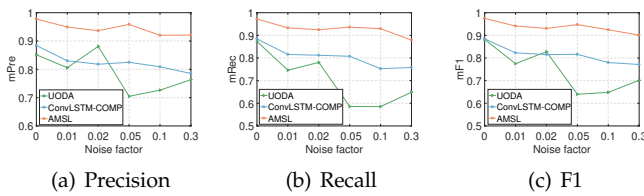


Fig. 6: Robustness to noise.

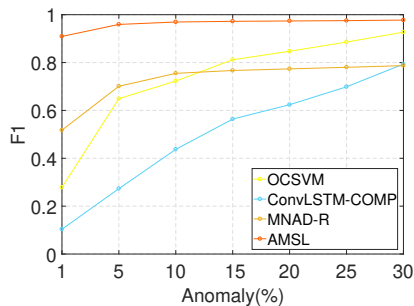


Fig. 7: F1 score on the abnormal class with different percentage of anomaly.

Fig. 7, we show the F1 score of the abnormal class using different methods. We compare the performance of four methods: OCSVM, ConvLSTM-COMPOSITE, MNAD-R and AMSL. These methods have good performance for CAP dataset in the above experiments. We can observe that as the percentage of anomaly decreases, the F1 score of other methods have decreased significantly, while our method still remain stable. This demonstrates that our method achieves high precision and recall on the abnormal class even if the percentage of anomaly is very low on the testing set. Therefore, we can come to a conclusion that our proposed method AMSL has good stability even facing dataset imbalance problem.

4.8 Case Study

We present case studies of MNAD, UODA and our method AMSL as shown in Fig. 5 by visualizing some normal and abnormal classes. We choose three dimensional signals on DSADS dataset. It can be shown that our AMSL can correctly classify these samples, while other methods fail

in two situations: (1) when a normal sample is not similar to the majority of normal samples (i.e., overfitting) or (2) when an abnormal sample is very similar to the normal sample (i.e., less powerful representations). This demonstrate that our AMSL can effectively handle the problem about the diversity of samples.

4.9 Parameter sensitivity analysis

We consider three key parameters: 1) window length V , 2) size of memory matrix M , and 3) filter size F in the last layer of encoder. These parameters are chosen from: $V \in \{64, 128, 236\}$, $C \in \{50, 200, 500, 800\}$, and $F \in \{16, 32, 64\}$. Fig. 8 shows that AMSL is robust to different parameter choices on PAMAP2 dataset.

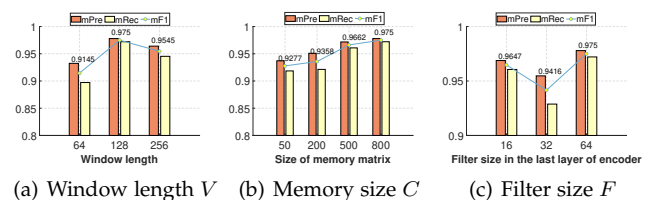


Fig. 8: Parameter sensitivity analysis on PAMAP2 dataset.

We also provide the sensitivity analysis of LMSL and GMSL. Fig. 9 (a-b) present the results of LMSL and GMSL with different window length $\{64, 128, 236\}$. We observe that the choice of window length is critical to the method, which window length of 128 can achieve the best performance for PAMAP2 dataset. The second factor is the size of memory matrix $M \in \mathbb{R}^{C \times F}$ that is set as $C \in \{50, 200, 500, 800\}$ and $F = 64$. Note that the dimension of F is equal to the filter size in the last layer of encoder. As shown Fig. 9 (c-d), they indicate that by increasing the size of C , the performance improves until C reaches nearly 800. Fig. 9 (e-f) show the results for $F \in \{16, 32, 64\}$. The filter size in the last layer of encoder obtains the best performance when $F = 64$. We use encoder to store more information and contribute it to reconstruct the latent features.

Moreover, we adjust the hyperparameters λ_1 and λ_2 in the loss function (Eq. (12)). As shown in Fig. 9 (g-h), we observe that $\lambda_1 = 1$ and $\lambda_2 = 0.0002$ are the best choice in our experiments.

The selection of the threshold μ in Section 3.5.2 is also an important issue for which we conduct experiments to

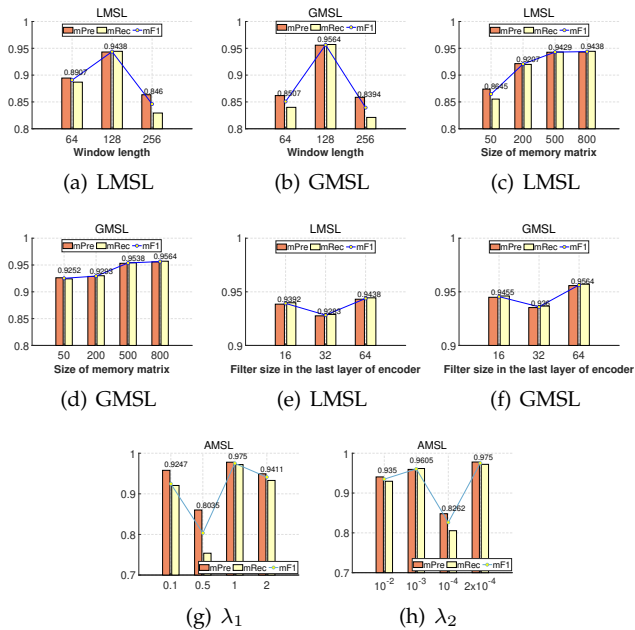


Fig. 9: Other parameter sensitivity analysis on PAMAP2 dataset.

TABLE 6: Experimental evaluation of threshold selection for AMSL.

Decision Rules	mPre	mRec	mF1	Acc
THR _{90th percentile}	0.9164	0.9401	0.9281	0.9284
THR _{95th percentile}	0.9486	0.9597	0.9541	0.9568
THR _{99th percentile}	0.9788	0.9713	0.9750	0.9770

compare different thresholds for α -percentiles: 90, 95 and 99. As shown in TABLE 6, we find that the 99th percentile can estimate the optimal threshold. Therefore, we set the anomaly detection threshold to the 99 percentile.

4.10 Convergence, Time and Space Complexity

Fig. 10(a) shows the convergence of reconstruction loss with memory modules along with the self-supervised loss. In conclusion, AMSL can be applied more effectively with a fast and stable convergence performance.

We also evaluate the inference time of AMSL and other strong baselines on DSADS dataset. Here, ‘‘COMP’’ means ConvLSTM-COMPOSITE. As shown in Fig. 10(b), in addition to achieving the best performance, our method requires shorter running time than most other methods.

Besides, according to TABLE 7 evaluated on DSADS dataset, the number of parameters and model size of AMSL are relatively smaller than most other methods. We also show that by reducing the model parameters through controlling the self-supervised data transformations R . We discard the poorly performing transformations in TABLE 3(a). AMSL($R=6$) discards the poorly performing transformation ‘‘Noise’’, AMSL($R=5$) discards ‘‘Noise’’ and ‘‘Scale’’ transformations, AMSL($R=4$) discards ‘‘Noise’’, ‘‘Scale’’ and ‘‘Permuted’’ transformations, and AMSL($R=3$) discards ‘‘Noise’’, ‘‘Scale’’, ‘‘Permuted’’ and ‘‘Reversed’’ transformations. We

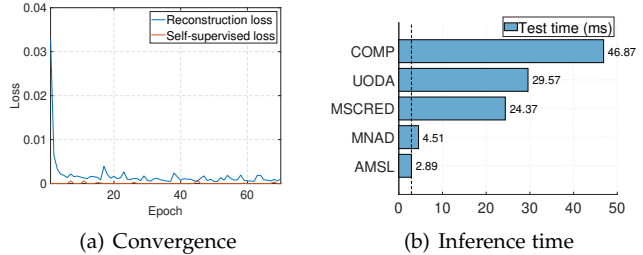


Fig. 10: Convergence and inference time of AMSL.

TABLE 7: Comparison of model parameters.

Method	F1	ACC	#Params	Model size
MNAD	0.8003	0.7811	3.8M	14.7MB
MSCRED	0.6428	0.6192	3.6M	13.8MB
UODA	0.8475	0.8365	1.9M	7.3MB
AMSL($R=7$)	0.9352	0.9332	3.3M	12.7MB
AMSL($R=6$)	0.9298	0.9273	2.8M	10.9MB
AMSL($R=5$)	0.9242	0.9233	2.4M	9.2MB
AMSL($R=4$)	0.9138	0.9135	1.9M	7.4MB
AMSL($R=3$)	0.9099	0.9096	1.5M	5.7MB

observe that our AMSL still achieves the best F1 and accuracy scores. On other datasets, the conclusions are also similar. This makes our method flexible in real applications.

5 CONCLUSIONS AND FUTURE WORK

In this paper, we proposed an Adaptive Memory Network with Self-supervised Learning (AMSL) for unsupervised anomaly detection on multivariate time series signals. To enhance the model’s generalization ability towards unseen anomalies, we proposed to use the self-supervised learning module to learn diverse normal patterns and an adaptive memory fusion network to learn rich feature representations by the global and local memory modules. Experiments on four public datasets demonstrate that our methods significantly outperforms existing approaches in terms of accuracy, generalization, and robustness.

In the future, we plan to extend AMSL to other modalities such as images and videos for unsupervised anomaly detection. In addition, we also plan to develop more efficient training algorithms and pursue the theoretical analysis of our method.

ACKNOWLEDGMENT

This work is supported by the National Key Research and Development Plan of China (No. 2020YFC2007104), Natural Science Foundation of China (No. 61972383, No. 61902377, No. 61902379), Science and Technology Service Network Initiative, Chinese Academy of Sciences (No. KJFJ-STQYD-2021-11-001), the National Research Foundation, Singapore under its AI Singapore Programme (AISG Award No. AISG2-RP-2020-019), the RIE 2020 Advanced Manufacturing and Engineering (AME) Programmatic Fund (No. A20G8b0102), Singapore, and the Nanyang Assistant Professorship (NAP).

REFERENCES

- [1] B. Kiran, D. Thomas, and R. Parakkal, "An overview of deep learning based methods for unsupervised and semi-supervised anomaly detection in videos," *Journal of Imaging*, vol. 4, no. 2, p. 36, 2018.
- [2] Y. Mirsky, T. Doitshman, Y. Elovici, and A. Shabtai, "Kitsune: an ensemble of autoencoders for online network intrusion detection," *arXiv preprint arXiv:1802.09089*, 2018.
- [3] M. C. Chuah and F. Fu, "Ecg anomaly detection via time series analysis," in *International Symposium on Parallel and Distributed Processing and Applications*. Springer, 2007, pp. 123–135.
- [4] M. G. Terzano, L. Parrino, A. Smerieri, R. Chervin, S. Chokroverty, C. Guilleminault, M. Hirshkowitz, M. Mahowald, H. Moldofsky, A. Rosa *et al.*, "Atlas, rules, and recording techniques for the scoring of cyclic alternating pattern (cap) in human sleep," *Sleep medicine*, vol. 3, no. 2, pp. 187–199, 2002.
- [5] M. Sakurada and T. Yairi, "Anomaly detection using autoencoders with nonlinear dimensionality reduction," in *Proceedings of the MLSDA 2014 2nd Workshop on Machine Learning for Sensory Data Analysis*. ACM, 2014, p. 4.
- [6] M. Gutoski, N. M. R. Aquino, M. Ribeiro, E. Lazzaretti, and S. Lopes, "Detection of video anomalies using convolutional autoencoders and one-class support vector machines," in *XIII Brazilian Congress on Computational Intelligence*, 2017, 2017.
- [7] M. Hasan, J. Choi, J. Neumann, A. K. Roy-Chowdhury, and L. S. Davis, "Learning temporal regularity in video sequences," in *CVPR*, 2016, pp. 733–742.
- [8] P. Malhotra, A. Ramakrishnan, G. Anand, L. Vig, P. Agarwal, and G. Shroff, "Lstm-based encoder-decoder for multi-sensor anomaly detection," *arXiv preprint arXiv:1607.00148*, 2016.
- [9] 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 AAAI*, vol. 33, 2019, pp. 1409–1416.
- [10] R. Paffenroth, P. Du Toit, R. Nong, L. Scharf, A. P. Jayasumana, and V. Bandara, "Space-time signal processing for distributed pattern detection in sensor networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 7, no. 1, pp. 38–49, 2013.
- [11] H. Hoffmann, "Kernel pca for novelty detection," *Pattern recognition*, vol. 40, no. 3, pp. 863–874, 2007.
- [12] R. Laxhammar, G. Falkman, and E. Sviestins, "Anomaly detection in sea traffic—a comparison of the gaussian mixture model and the kernel density estimator," in *2009 12th International Conference on Information Fusion*. IEEE, 2009, pp. 756–763.
- [13] L. J. Latecki, A. Lazarevic, and D. Pokrajac, "Outlier detection with kernel density functions," in *International Workshop on MLDM*. Springer, 2007, pp. 61–75.
- [14] A. Banerjee, P. Burlina, and C. Diehl, "A support vector method for anomaly detection in hyperspectral imagery," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 44, no. 8, pp. 2282–2291, 2006.
- [15] J. Ma and S. Perkins, "Time-series novelty detection using one-class support vector machines," in *Proceedings of the International Joint Conference on Neural Networks*, 2003., vol. 3. IEEE, 2003, pp. 1741–1745.
- [16] J. R. Medel and A. Savakis, "Anomaly detection in video using predictive convolutional long short-term memory networks," *arXiv preprint arXiv:1612.00390*, 2016.
- [17] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," *Special Lecture on IE*, vol. 2, no. 1, 2015.
- [18] H. Xu, W. Chen, N. Zhao, Z. Li, J. Bu, Z. Li, Y. Liu, Y. Zhao, D. Pei, Y. Feng *et al.*, "Unsupervised anomaly detection via variational auto-encoder for seasonal kpis in web applications," in *Proceedings of the 2018 World Wide Web Conference*. International World Wide Web Conferences Steering Committee, 2018, pp. 187–196.
- [19] D. Xu, Y. Yan, E. Ricci, and N. Sebe, "Detecting anomalous events in videos by learning deep representations of appearance and motion," *Computer Vision and Image Understanding*, vol. 156, pp. 117–127, 2017.
- [20] D. Wulsin, J. Blanco, R. Mani, and B. Litt, "Semi-supervised anomaly detection for eeg waveforms using deep belief nets," in *2010 Ninth International Conference on Machine Learning and Applications*. IEEE, 2010, pp. 436–441.
- [21] S. Hawkins, H. He, G. Williams, and R. Baxter, "Outlier detection using replicator neural networks," in *International Conference on Data Warehousing and Knowledge Discovery*. Springer, 2002, pp. 170–180.
- [22] C. Zhou and R. C. Paffenroth, "Anomaly detection with robust deep autoencoders," in *KDD*. ACM, 2017, pp. 665–674.
- [23] W. Lu, Y. Cheng, C. Xiao, S. Chang, S. Huang, B. Liang, and T. Huang, "Unsupervised sequential outlier detection with deep architectures," *IEEE transactions on image processing*, vol. 26, no. 9, pp. 4321–4330, 2017.
- [24] P. Filonov, A. Lavrentyev, and A. Vorontsov, "Multivariate industrial time series with cyber-attack simulation: Fault detection using an lstm-based predictive data model," *arXiv preprint arXiv:1612.06676*, 2016.
- [25] T. Ergen, A. H. Mirza, and S. S. Kozat, "Unsupervised and semi-supervised anomaly detection with lstm neural networks," *arXiv preprint arXiv:1710.09207*, 2017.
- [26] G. Lai, W.-C. Chang, Y. Yang, and H. Liu, "Modeling long-and short-term temporal patterns with deep neural networks," in *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*. ACM, 2018, pp. 95–104.
- [27] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection—a new baseline," in *CVPR*, 2018, pp. 6536–6545.
- [28] Z. Li, J. Tang, and T. Mei, "Deep collaborative embedding for social image understanding," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 9, pp. 2070–2083, 2018.
- [29] Z. Zhu and Z. Li, "Online video object detection via local and mid-range feature propagation," in *Proceedings of the 1st International Workshop on Human-centric Multimedia Analysis*, 2020, pp. 73–82.
- [30] L. Jing and Y. Tian, "Self-supervised visual feature learning with deep neural networks: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [31] M. Yang, Y. Li, Z. Huang, Z. Liu, P. Hu, and X. Peng, "Partially view-aligned representation learning with noise-robust contrastive loss," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1134–1143.
- [32] Y. Lin, Y. Gou, Z. Liu, B. Li, J. Lv, and X. Peng, "Completer: Incomplete multi-view clustering via contrastive prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 174–11 183.
- [33] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," *arXiv preprint arXiv:1910.13461*, 2019.
- [34] M. Ravanelli, J. Zhong, S. Pascual, P. Swietojanski, J. Monteiro, J. Trmal, and Y. Bengio, "Multi-task self-supervised learning for robust speech recognition," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 6989–6993.
- [35] R. Ali, M. U. K. Khan, and C. M. Kyung, "Self-supervised representation learning for visual anomaly detection," *arXiv preprint arXiv:2006.09654*, 2020.
- [36] S. Wang, Y. Zeng, X. Liu, E. Zhu, J. Yin, C. Xu, and M. Kloft, "Effective end-to-end unsupervised outlier detection via inlier priority of discriminative network," in *Advances in Neural Information Processing Systems*, 2019, pp. 5962–5975.
- [37] S. Sukhbaatar, J. Weston, R. Fergus *et al.*, "End-to-end memory networks," in *Advances in neural information processing systems*, 2015, pp. 2440–2448.
- [38] Q. Cai, Y. Pan, T. Yao, C. Yan, and T. Mei, "Memory matching networks for one-shot image recognition," in *CVPR*, 2018, pp. 4080–4088.
- [39] A. Santoro, S. Bartunov, M. Botvinick, D. Wierstra, and T. Lillicrap, "Meta-learning with memory-augmented neural networks," in *International conference on machine learning*, 2016, pp. 1842–1850.
- [40] M. Wang, Z. Lu, H. Li, and Q. Liu, "Memory-enhanced decoder for neural machine translation," *arXiv preprint arXiv:1606.02003*, 2016.
- [41] C. Zhang, Y. Wang, X. Zhao, Y. Guo, G. Xie, C. Lv, and B. Lv, "Memory-augmented anomaly generative adversarial network for retinal oct images screening," in *2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI)*. IEEE, 2020, pp. 1971–1974.
- [42] D. Gong, L. Liu, V. Le, B. Saha, M. R. Mansour, S. Venkatesh, and A. v. d. Hengel, "Memorizing normality to detect anomaly: Memory-augmented deep autoencoder for unsupervised anomaly detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 1705–1714.
- [43] H. Park, J. Noh, and B. Ham, "Learning memory-guided normality for anomaly detection," in *CVPR*, 2020, pp. 14 372–14 381.

- [44] L. Shen, Z. Li, and J. Kwok, "Timeseries anomaly detection using temporal hierarchical one-class network," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [45] A. Saeed, T. Ozcelebi, and J. Lukkien, "Multi-task self-supervised learning for human activity detection," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 2, pp. 1–30, 2019.
- [46] 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 *ICLR*, 2018.
- [47] K. Altun, B. Barshan, and O. Tunçel, "Comparative study on classifying human activities with miniature inertial and magnetic sensors," *Pattern Recognition*, vol. 43, no. 10, pp. 3605–3620, 2010.
- [48] A. Reiss and D. Stricker, "Introducing a new benchmarked dataset for activity monitoring," in *2012 16th International Symposium on Wearable Computers*. IEEE, 2012, pp. 108–109.
- [49] P. Schmidt, A. Reiss, R. Duerichen, C. Marberger, and K. Van Laerhoven, "Introducing wesad, a multimodal dataset for wearable stress and affect detection," in *Proceedings of the 20th ACM International Conference on Multimodal Interaction*, 2018, pp. 400–408.
- [50] H.-P. Kriegel, M. Schubert, and A. Zimek, "Angle-based outlier detection in high-dimensional data," in *SIGKDD*. ACM, 2008, pp. 444–452.
- [51] S. S. Joshi and V. V. Phoha, "Investigating hidden markov models capabilities in anomaly detection," in *Proceedings of the 43rd annual Southeast regional conference-Volume 1*. ACM, 2005, pp. 98–103.
- [52] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015, pp. 2625–2634.
- [53] B. Zhou, S. Liu, B. Hooi, X. Cheng, and J. Ye, "Beatgan: Anomalous rhythm detection using adversarially generated time series." in *IJCAI*, 2019, pp. 4433–4439.
- [54] 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.
- [55] "Cap sleep database," 2012. [Online]. Available: <https://physionet.org/content/capslpdb/1.0.0/>
- [56] N. Ketkar, "Introduction to keras," in *Deep learning with Python*. Springer, 2017, pp. 97–111.
- [57] S. Zhai, Y. Cheng, W. Lu, and Z. Zhang, "Deep structured energy based models for anomaly detection," *arXiv preprint arXiv:1605.07717*, 2016.