# Forecast-based Multi-aspect Framework for Multivariate Time-series Anomaly Detection

Lan Wang, Yusan Lin, Yuhang Wu, Huiyuan Chen, Fei Wang, Hao Yang

Visa Research, Palo Alto, CA, USA, 94306

{lwang4, yusalin, yuhawu, hchen, feiwang, haoyang}@visa.com

*Abstract*—Today's cyber-world is vastly multivariate. Metrics collected at extreme varieties demand multivariate algorithms to properly detect anomalies. However, forecast-based algorithms, as widely proven approaches, often perform sub-optimally or inconsistently across datasets. A key common issue is they strive to be one-size-fits-all but anomalies are distinctive in nature. We propose a method that tailors to such distinction. Presenting FMUAD - a Forecast-based, Multi-aspect, Unsupervised Anomaly Detection framework. FMUAD explicitly and separately captures the signature traits of anomaly types - spatial change, temporal change and correlation change - with independent modules. The modules then jointly learn an optimal feature representation, which is highly flexible and intuitive, unlike most other models in the category. Extensive experiments show our FMUAD framework consistently outperforms other state-of-the-art forecast-based anomaly detectors.

*Index Terms*—Anomaly Detection, Multivariate Time Series, Unsupervised Learning.

## I. INTRODUCTION

Anomaly detection for multivariate time series is of great interest in many real-word applications, such as road traffic surveillance [1], financial fraud detection [2], software monitoring [3], web log analysis [4] and network analysis [5], etc. Detecting unexpected behavior that deviates from the normal behavior accurately and efficiently can avoid further damage of core systems and save huge amount of revenue for business.

Among the various challenges to the said anomaly detection application, the most important one lies in the nature of the multi-variate time series themselves. For example, anomalous events are the mutual effect of multiple time series, such that looking at each one separately can result in poor quality detection [6]. Intra-series patterns are far from trivial as well. The temporal dynamics, such as shift in frequency and abrupt change-of-trend [7], need specific modeling to make accurate forecasting based on historical data [8]. In addition, value changes themselves may be subtle enough to disguise as normal data, but looking at the same time series in different granularities or scales tells completely different stories. All of these peculiarities need to be accounted for designing an effective detector.

Most forecast-based models in the past do not explicitly address the above-mentioned patterns, let alone separately [3], [9], [10], [11], [12], and therefore are subject to each of the above deficiencies. However, *we argue that there is a way to fix such deficiencies and forecast-based models in general have room for improvement*. In this paper, we propose a modular approach whereby we create different modules to best characterize each of the three anomalous traits: inter-series correlation dynamics, intra-series temporal dynamics and multi-scale spatial dynamics. These modules are trained jointly to provide a unified inference.

Another key concern in our proposal is the lack of labels for anomalous events [13]. As anomalies are commonly outnumbered by normal data points, it is crucial to have a label-agnostic model. Consequently, a noticeable trend in the domain is pivoting towards unsupervised models, i.e. ones that train on normal data exclusively, and detect anomalies with the intuition that they deviate from the patterns of normal events. Our proposal follows suit with this trend.

The main contributions of this paper are summarized as follows:

- We introduce a Forecast-based Multi-aspect Unsupervised Anomaly Detection framework (FMUAD) to capture anomalies with different patterns. The model outperforms other state-of-the-art forecast-based methods by achieving over 17% F1-score improvement on average cross multiple public datasets.
- We design a modular architecture that explicitly addresses different patterns of anomalies, which makes our framework highly intuitive and explainable. By learning different modules jointly, our approach not only captures each individual anomaly pattern but also detects complicated real-world anomalies with mixed patterns.
- We define a novel loss function that involves the forecasting error as well as its compactness (or variance). Compactness has proved to be a desired quality of a feature to control the variance within one class [9]. By controlling the compactness of normal class, we aim to learn a tighter representation that is more sensitive to outliers.

## II. RELATED WORK

As Chandola *et al.*[14] put in his nice words, "Anomalies are patterns in data that do not conform to a well defined notion of normal behavior." And such is the fundamental principle of unsupervised anomaly detectors (UADs) - Instead of learning to tell apart normal and abnormal events, UADs simply discover data that deviates from normal. Technically speaking, UADs leverage the following major characteristics of anomalies: (1) The distribution of anomalies differs remarkably from the normal data points; (2) Anomalies happen infrequently and
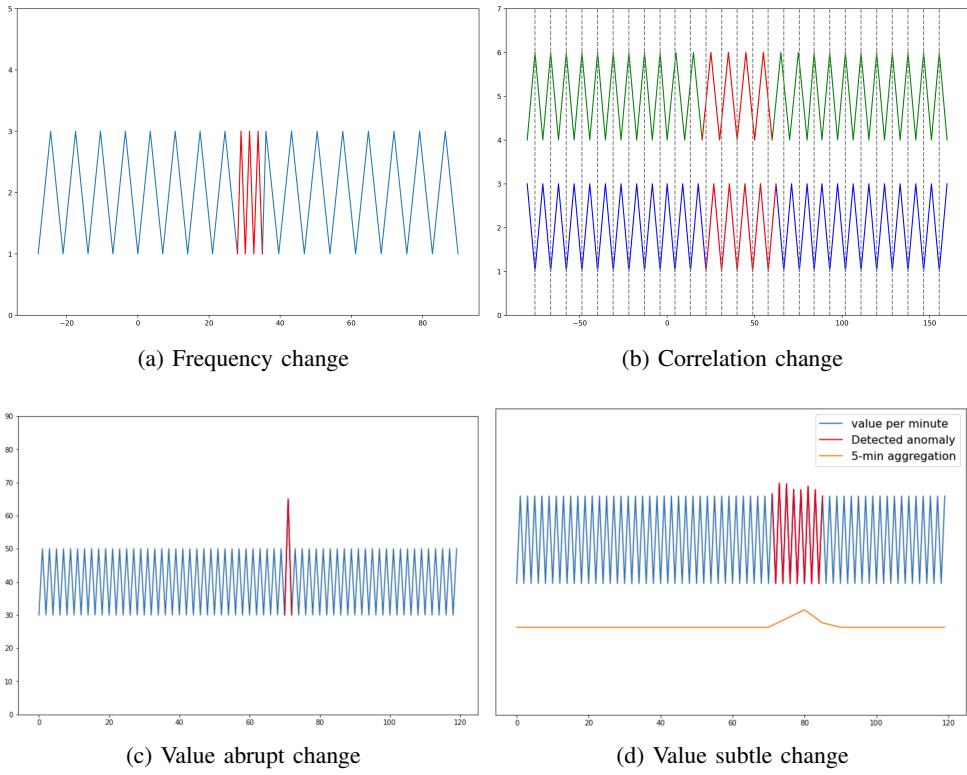
Fig. 1: Examples of anomaly patterns: (a) temporal change; (b) correlation change; (c, d) spatial change where (c) corresponds to a dramatic change and (d) more subtle change in which anomalies accumulates over a time period. The red region represents where the anomalies happen while data points located elsewhere are normal.

thus form a small set of data points compared with the normal data.

Following that, the importance of modeling normal behavior goes without saying. In particular, we need to find an effective way to represent normal data. Time series, a prime subject for anomaly detection, is often diverse and complex due to temporal dependency between data points. Some researchers attempted to leverage a graph structure to model the dynamic change [6], [8], [15], [16], [17], [18]. However, graph construction and modeling adds extra complexity to the system, which makes it heavy and time costly, and thus are not scalable for real-world applications. To this day, how to find an universally satisfactory representation for time series is still under intense discussions in the field.

Among the vast literature on time series anomaly detection over the past decade [19], [20], [21], [7], [22], [23], [6], one can divide them into three groups: (1) proximity-based; (2) reconstruction-based; and (3) forecast-based. Proximity-based methods quantify object similarities based on a defined distance measure and detect objects that are far away from the majority as anomalies. Reconstruction-based methods use the reconstruction error as the core and assume that anomalies lie in a different manifold from most data points, and thus cannot be effectively reconstructed from low-dimensional space. Forecast-based methods assume anomalies incorporate unusual patterns that are hard to predict based on historical data, and

detect anomalies based on the forecasting error. Our work falls under this category. Below, we present a brief background, along with some important past work on forecast-based time series anomaly detection.

In a forecast-based anomaly detector, we prepare the model by fitting it against the training data, which only contains normal data points. The trained model is then put forth to make predictions on future "normal" data, to be compared against the input. We infer the input point as anomalous if its value deviates from the model's prediction range. Given such fundamental mechanism, it easily follows that forecast-based anomaly detectors mainly differentiate by their forecasting methods. One of the most classic approaches, ARIMA [24] is a popular statistical analysis model that learns the auto-correlations in the time series for future value prediction. Other statistics models such as Holt-Winters[24] and FDA [25] serve the same purpose. While these methods are efficient, they are sensitive to datasets and model parameter choices. Fine tuning these models often requires strong assumptions and extensive domain knowledge about the data. Machine-learning based models attempt to overcome such limitations. Hierarchical Temporal Memory (HTM) [26] is an unsupervised sequence memory algorithm to detect anomalies in streaming data. Ding et al.[27] combined the HTM model and Bayesian network into for a real-time framework for multi-variate time series anomaly detection. More recently, Long Short-Term Memory

Recurrent Neural Network (LSTM-RNN), an effective deep learning network that models temporal dynamics for sequential data, has been widely used for many state-of-the-art works. Hundman *et al.*[28] proposed LSTM-NDT, an unsupervised and non-parametric thresholding approach to interpret prediction generated by an LSTM network. Zong *et al.*[11] proposed DAGMM to detect anomalies by jointly optimizing the parameters of a deep neural network and a Gaussian Mixture Model.

Although the past forecast-based models have proven the category's potential in anomaly detection, they were designed to be a capture-all solution for every possible anomaly, which can be risky due to the drastically difference between the anomaly patterns. We believe a degree of tailoring in the model to anomaly patterns can help bring out better performance.

## III. THE PROPOSED METHOD

We propose an unsupervised, forecast-based, multi-variate time-series anomaly detection algorithm (FMUAD) to explicitly address different natures of anomaly patterns - temporal change, spatial change and correlation change, which are challenging to detect using forecast-based detectors of the past. In this section, we first formalize our problem and then provide an overview of the model methodology.

### A. Problem Formulation

A multivariate time series is a sequence of multi-dimensional data points

$$\mathcal{X} = \{\boldsymbol{x}_1, \boldsymbol{x}_2, ...., \boldsymbol{x}_T\},$$

where $\boldsymbol{x}_i \in \mathbb{R}^m$ for $i \in \{1, ..., T\}$, $m$ is the number of features and $T$ is the number of time steps. The uni-variate setting is a particular case of the multi-variate one with $m = 1$.

The model takes $\mathcal{X}$ as the input training data assuming $\mathcal{X}$ contains only normal time steps. Anomaly detection refers to identifying a new unobserved instance $\hat{\boldsymbol{x}}_t, t > T$ based on how differently it behaves compared to $\mathcal{X}$. Due to the temporal dependency within the time series and the fact that anomalies usually occur over a contiguous time period, we do not model each time step independently but rather consider a window segment. Suppose $W_t \in \mathbb{R}^{m \times k}$ is a window with length $k$ at time $t$:

$$W_t = \{\boldsymbol{x}_{t-k+1}, \boldsymbol{x}_{t-k+2}, ..., \boldsymbol{x}_t\}.$$

We can transform the original time series $\mathcal{X}$ into a sequence of windows $\boldsymbol{W} = \{W_k, ..., W_T\}$. For each window $W_t$, we train the model to forecast its behavior using the historical time steps that were already seen. With proper training, the model will be capable of making decent predictions for normal data. Then given a new unseen window $\tilde{W}_t, t > T$, our goal for the model is to detect whether an anomaly happens at time $t$ based on the forecast error: the larger the error is, the more likely an anomaly happens.

## TABLE I: Notations

| | |
|---|---|
| $m \in \mathbb{Z}$ | Number of time series (features) |
| $T \in \mathbb{Z}$ | Number of total time steps |
| $k \in \mathbb{Z}$ | Window size |
| $\tau \in \mathbb{Z}$ | Input instance time span |
| $\mathcal{X} \in \mathbb{R}^{m \times T}$ | Input time series |
| $I_t \in \mathbb{R}^{m \times \tau}$ | Individual input instance |
| $W_t \in \mathbb{R}^{m \times k}$ | True window at time $t$ |
| $S_t \in \mathbb{R}^{m \times m}$ | True signature matrix (correlation) for $W_t$ |
| $F_t \in \mathbb{R}^{m \times k/2}$ | True frequency matrix for $W_t$ |
| $W_t^h \in \mathbb{R}^{m \times (\tau-k)}$ | Historical segment for $W_t$ |
| $\hat{W}_t \in \mathbb{R}^{m \times k}$ | Predicted window at time $t$ |
| $\hat{S}_t \in \mathbb{R}^{m \times m}$ | Predicted signature matrix (correlation) for $W_t$ |
| $\hat{F}_t \in \mathbb{R}^{m \times k/2}$ | Predicted frequency matrix for $W_t$ |

### B. Our Method

In this section, we first present the overall architecture of the model and then break down each module into details.

*1) Overall Architecture:* We illustrate the overall architecture of our model in this part. See Figure 2 for a high-level walk-through and Figure 4 for a more detailed illustration.

To detect anomaly in a target window $W_t$, we concatenate the necessary historical information needed for forecasting, or $W_t^h$, with $W_t$ itself, to formulate the input to our model:

$$I_t = [W_t^h; W_t],$$

which takes the form of a $m \times \tau$ matrix:

$$I_t = \{\boldsymbol{x}_{t-\tau+1}, \boldsymbol{x}_{t-\tau+2}, ..., \boldsymbol{x}_t\},$$

where $m$ is the number of features and $\tau$ is the number of time steps we trace back for each prediction. In our setting, we choose $\tau = 500$ and $k = 30$. The output of our framework, $\hat{y}_t \in \{0, 1\}$ is a boolean representing whether anomaly happens at time $t$.
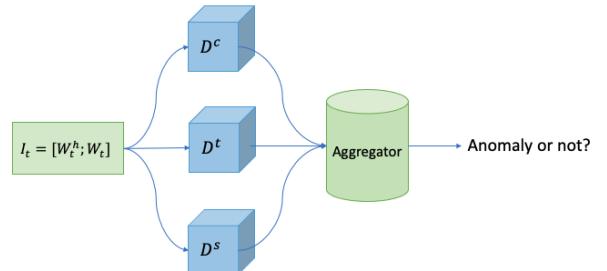


Fig. 2: High-level system architecture

To holistically characterize and capture the traits of the anomalies within each window, we consider the time series dynamics from three aspects and design a corresponding detector tailored to each aspect:

- $D^c$, a detector for **inter-series correlation** change;
- $D^t$, a detector for **intra-series temporal pattern** change;
- $D^s$, a detector for **multi-scale spatial pattern** change.

Because the nature of the detectors are largely dictated by the traits they're assigned to capture, each detector should use a different input/output data representation that best suits the

detector. To this regard, we create two transformed matrices $S_t$ and $F_t$, and use them along with $W_t$ for the three detectors $D^c, D^t, D^s$ respectively. The transformation process TF1: $W_t \mapsto S_t$ and TF2: $W_t \mapsto F_t$ will be discussed in details in part 2) and 3).

– $W_t$, the original matrix for the window at $t$;
– $F_t$ (frequency matrix), a derived matrix that contains temporal information for $W_t$ in the frequency domain;
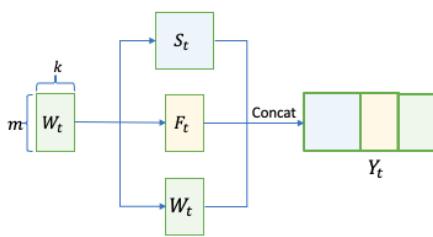– $S_t$ (signature matrix), a derived matrix that contains inter-series correlation information for $W_t$.



Fig. 3: Concatenation of three matrices

As depicted by Figure 3, $Y_t$ is a concatenation of the ground truth matrices $S_t, F_t$ and $W_t$, i.e.,

$$Y_t = [S_t; F_t; W_t] = [\text{TF1}(W_t); \text{TF2}(W_t); W_t].$$

Our model outputs $\hat{Y}_t = [\hat{S}_t; \hat{F}_t; \hat{W}_t]$ as the predicted matrix derived from the concatenation of the forecasting results obtained from the three detectors. The error between $Y_t$ and $\hat{Y}_t$ will be served as an indicator of anomalies.

*2) Inter-series correlation detector:* Previous studies [19], [29] suggest the importance of correlation among time series in the characterization of the dynamic of the multi-variate time series system for anomaly detection. For instance, the anomaly shown in the red region of Figure 1(b) is determined based on abnormal correlation: The two time series have been mostly negative correlated, while in the labelled pattern they are positively correlated.

To capture the correlation changes, inspired by [19], we define a signature matrix $S_t$ to characterize the correlation between time series. The transformer TF1 that produces $S_t$ is defined as follows: Given a window $W_t$, $S_t$ is a $m \times m$ matrix whose entry at the $i$th row and $j$th column is computed as the Cosine similarity score between the $i$th time series $\{x^i_{t-k+1}, ..., x^i_t\}$ and the $j$th time series $\{x^j_{t-k+1}, ..., x^j_t\}$. The reason we adopt Cosine similarity instead of the inner product as used in [19] is that the former is less prone to scaling effects, and allows us to focus on the correlation effects only.

Now we transform the input time series into a sequence of signature matrices. First, we splice the input $I_t$ into a sequence of windows with stride $s$. Then, we apply TF1 to each of the window to form a sequence containing $d = \lfloor (\tau - k)/s \rfloor$ signature matrices.

$$S_t \leftarrow \text{TF1}(W_t). \tag{1}$$

To make a forecast on the signature matrix $S_t$, we apply ConvLSTM with an attention mechanism. The final forecasting signature matrix is denoted as $\hat{S}_t \in \mathbb{R}^{m \times m}$. ConvLSTM [30] has been developed and proven effective to capture the sequential temporal information, as formulated below:

$$
\begin{aligned}
i_t &= \sigma(W_{si} * S_t + W_{hi} * H_{t-1} + W_{ci} \circ C_{t-1} + b_i) \\
f_t &= \sigma(W_{sf} * S_t + W_{hf} * H_{t-1} + W_{cf} \circ C_{t-1} + b_f) \\
C_t &= f_t \circ C_{t-1} + i_t \circ \tanh(W_{sc} * S_t + W_{hc} * H_{t-1} + b_c) \\
o_t &= \sigma(W_{so} * S_t + W_{ho} * H_{t-1} + W_{co} \circ C_t + b_o) \\
H_t &= o_t \circ \tanh(C_t),
\end{aligned}
\tag{2}
$$

where $i_t$, $f_t$, and $o_t$ denote the input gate, forget gate, and output gate, respectively. $C_t$ and $H_t$ are the cell states and hidden states. $*$ denotes the convolution operator, $\circ$ is the Hadamard product, and $\sigma(\cdot)$ is the sigmoid function.

In addition, since not all previous hidden states have equal effects on the current state, we apply a temporal attention mechanism to adaptively adjust the weight of the previous steps that are relevant to the current one and aggregate all to obtain the overall effect [31]. More concretely, we have

$$H_t^* = \sum_{i=t-k}^{t} c_i H_i, \quad c_i = \frac{\exp \langle \text{vec}(H_i), \text{vec}(H_t) \rangle}{\sum_{j=t-k}^{t} \exp \langle \text{vec}(H_j), \text{vec}(H_t) \rangle}, \tag{3}$$

where $\text{vec}(\cdot)$ is the matrix flattened into a vector, $\langle \cdot, \cdot \rangle$ denotes the inner product, and $H_t^*$ is the representation for window $t$. Then with an additional convolutional neural network with a $1 \times 1$ kernel, we project the representation $H_t^*$ into a lower dimensional space with only one channel and then apply a fully-connected layer to derive the prediction of $S_t$:

$$\hat{S}_t \leftarrow \text{FC}(\text{Conv}_{1 \times 1}(H_t^*)). \tag{4}$$

The output serves as the final prediction $\hat{S}_t$ containing fore-casted correlation information for window at time $t$.

*3) Intra-series temporal pattern detector:* As are inter-series correlations, the intra-series temporal patterns are also crucial for anomaly detection. For instance, the anomaly shown in the red region of Figure 1(a) is due to temporal pattern changes: the frequency suddenly increases twice in the anomaly region compared to the rest of the time series. Several recent work [8], [15] has shown the advantages of joint learning inter-series correlations and intra-series temporal patterns for multi-variate time series forecasting.

To capture the temporal pattern within each series, we project each uni-variate time series to the spectral domain via DFT (discrete Fourier transform):

$$\xi_j = \frac{1}{k} \sum_{\ell=0}^{k-1} x_l e^{2\pi i j \ell / k}, j = 1, 2, ..., k, \tag{5}$$

where $i$ is the imaginary unit. The DFT transforms $k$ discrete-time samples $\vec{x} = \{x_1, x_2, ..., x_k\}$ in the time domain to the frequency domain $\vec{\xi} = \{\xi_1, \xi_2, ..., \xi_k\}$. Due to the symmetry property of DFT for real-valued discrete-time signals: $\xi_j = \xi_{k-j}$, only half of the $\xi$'s need to be stored. As a result,
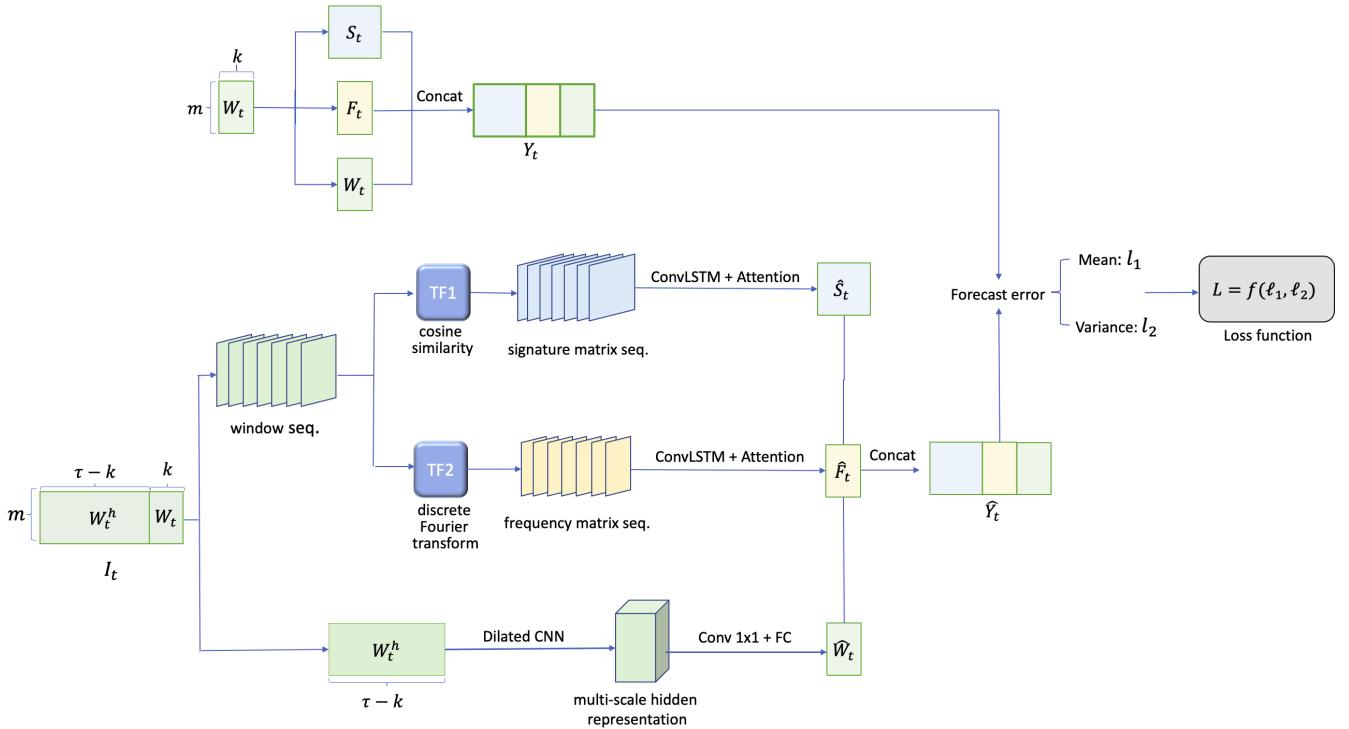
Fig. 4: Detailed model architecture

the window $W_t$ can be transformed to a frequency matrix $F_t = \{\xi_1, ..., \xi_{k/2}\} \in \mathbb{R}^{m \times \frac{k}{2}}$ representing the amplitude of frequencies for the $m$ uni-variate time series in the spectral domain. Such process defines the transformer TF2.

$$F_t \leftarrow \text{TF2}(W_t). \qquad (6)$$

Similar to how the inter-series correlation was constructed in part 2), we transform the input time series into a sequence of frequency matrices, then apply ConvLSTM with an attention mechanism to obtain a high-dimensional hidden representation. Lastly we pass an additional neural network with a $1 \times 1$ kernel to project it to a low-dimensional space and apply a fully-connected layer:

$$\hat{F}_t \leftarrow \text{FC}(\text{Conv}_{1 \times 1}(\text{Attn}(\text{ConvLSTM}(F_{t-})))). \qquad (7)$$

Here, $\text{ConvLSTM}(\cdot)$ refers to Eq. (2), $\text{Attn}(\cdot)$ refers to Eq. (3), and $F_{t-}$ refers to the sequence of frequency matrices before time $t$. The forecast $\hat{F}_t$ captures the temporal information from each time series within window $W_t$.

*4) Multi-scale spatial pattern detector:* Another pattern of anomaly, sometimes overlooked but extremely common in time series is the spacial dynamic, or in layman's term, "value change". For instance, the red region of Figure 1(c) shows a dramatic value change, which is trivial to capture. However, the majority of spatial changes are gradual and subtle, and hard to detect like the red region in Figure 1(d). The value change at each individual time step is so small that they are hardly discernible to the human eye, but when accumulated over a time period, the deviation is significant. Thus, it is

important for our anomaly detector to capture both types of spacial dynamics.

The intuition is natural: Aggregating the time series in Figure 1(d) at a larger time scale easily manifests the subtle changes as a spike, thus reducing the problem to a trivial one. In fact, looking at time series at different levels of aggregation can tell completely different stories. To capture multi-scale patterns, one needs to consider relatively long-term dynamics, which is why we take $W_t^h \in \mathbb{R}^{m \times (\tau - k)}$ to model instead of taking individual windows as we did before. In our model, we adopt the Dilated Convolutional Neural Network [32] for this purpose. The major difference between dilated CNN and basic CNN is that it introduces dilated convolution operator to capture information aggregated at multiple scales. Suppose $F$ is a discrete function and $k$ is the convolution kernel, we can define $\ell$ as a dilation factor and $*_\ell$ as a dilated convolution operator

$$(F *_\ell k)(p) = \sum_{s + \ell t = p} F(s)k(t). \qquad (8)$$

With different dilation sizes $\ell$, multi-scale information can be obtained without losing resolution.

In our setting specifically, we firstly apply a 3-layer Dilated-CNN with dilating numbers $r = 1, 3, 5$ and respective channels numbers $c = 32, 64, 128$ to obtain a deep hidden representation. In particular, we apply a rectangle kernel $(3, 1)$ for each layer instead of an usual square kernel $(3, 3)$ to limit the multi-scale aggregation to over the time dimension rather than the feature dimension. Then we decode the hidden representation

to forecast the target window $W_t$. The final forecasting matrix is denoted as $\hat{W}_t \in \mathbb{R}^{m \times k}$:

$$\hat{W}_t \leftarrow \text{FC}(\text{Conv}_{1 \times 1}(\text{Dilated-CNN}(W_t^h))). \tag{9}$$

*5) Loss function:* In this section, we introduce a novel loss function based on the forecasting error as well as its variance.

Following steps 2)-4), for a given target window $W_t$, we have derived the ground truth matrix $Y_t = [S_t; F_t; W_t]$ and the forecasted matrix $\hat{Y} = [\hat{S}_t; \hat{F}_t; \hat{W}_t]$ as shown in Figure 4. The forecasting error is the $L_2$ norm between $Y_t$ and $\hat{Y}_t$. As we train our neural network in batches of size $b$, the batch loss is averaged over the batch as shown in Eq. (10). The total loss is further averaged across the batches.

$$\ell_1 = \frac{1}{b} \sum_{i=1}^{b} \left\| Y_i - \hat{Y}_i \right\|_2^2. \tag{10}$$

Besides forecasting error, we introduce an extra term we call the *compactness loss* to capture the variance of the forecasting error within each batch, which is defined in Eq. (11).

$$\ell_2 = \frac{1}{nb} \sum_{i=1}^{b} z_i^T z_i, \ z_i = \hat{Y}_i - \frac{1}{b-1} \sum_{j \neq i} \hat{Y}_j, \tag{11}$$

where $n = m + k + \frac{k}{2}$ is the number of columns of $\hat{Y}_i$. Due to the fact that we are training on a single class (normal instances), the variance control $\ell_2$ helps to obtain a compact representation that is more sensitive to anomalies [9], [33].

The final loss is computed based on $\ell_1$ and $\ell_2$ via a function $f$ as defined below:

$$\ell = f(\ell_1, \ell_2) = (\epsilon + \ell_2) \cdot \ell_1. \tag{12}$$

Here, we take a small positive constant $\epsilon$ (In our experiment, we set $\epsilon = 10^{-5}$) to avoid a zero trivial solution: A trivial network with all zero parameters gives zero loss value $\ell$ if $\epsilon = 0$. We adopt this loss function instead of a weighted sum of $\ell_1$ and $\ell_2$ due to two reasons. First, this eliminates the need to fine-tune the weights, making the training process much more efficient. Second, the effect each term adds on scales with its order of change. In other words, if the mean and variance of the forecasting error have the same order of changes, then these changes are emphasized to the same level no matter how different their scales are. To the best of our knowledge, this is the first work that incorporates the product of the error mean and variance in the loss function.

Our experiments show that by considering the variance factor $\ell_2$, the model performance can be consistently improved. In addition, While we consider both the forecasting error and its variance in training, we only use the forecasting error $\ell_1 = \|Y_t - \hat{Y}_t\|_2^2$ as the testing loss. The testing loss serves as the anomaly score to make the judgement: the larger the testing loss is, the more likely an anomaly has occurred.

## IV. EXPERIMENTS

In this section, we perform extensive experiments to validate the effectiveness of our framework. We first describe the datasets that the framework tests on, then introduce three baseline models and the performance metrics, and make a comprehensive comparison between our model the baselines. Detailed case study and ablation study is presented at the end.

### A. Datasets

To intuitively illustrate the inner workings of our model and evaluate its performance, we adopt two highly regarded and commonly used public multi-variate time-series datasets: SMD[1] (Server Machine Dataset) [20] and MSL[2] (Mars Science Laboratory rover) [28]. SMD is 5 weeks worth of server telemetry collected from 28 production server nodes of a reputable Internet company. It contains 38 simultaneously collected floating-point metrics, with anomaly events labelled in a binary format and the anomaly ratio is 4.16%. The data is segmented among the 28 servers into 3 groups. We train and test on each server machine separately. The total training set for SMD contains 708405 data points while the testing set contains 708420. The MSL dataset captures telemetry data throughout the Mars mission duration of the MSL rover, better known as *Curiosity*. It contains data for 55 entities each monitored by 25 metrics, with anomalies expertly labelled by NASA and the anomaly ratio is 10.72%. The training set size for MSL is 58317 while the testing size is 73729. Both datasets depict highly real-world scenarios in the target applications of our model, and therefore can provide an accurate benchmark for our model.

| Dataset Name | | # Features | Feature Name | Anomaly Ratio(%) |
|---|---|---|---|---|
| SMD | group 1 | 38 | CPU load, network usage, memory usage, etc. | 4.16 |
| | group 2 | | | |
| | group 3 | | | |
| MSL | | 55 | Telemetry data: radiation, power, temperature, computational activity, etc | 10.72 |

TABLE III: Dataset descriptions

### B. Comparison with Baselines

We compare our model with three state-of-the-art anomaly detection models: LSTM-NDT [28], DAGMM [11] and LSTM-VAE [34]. The first two models are forecast-based, which falls in the same scope as our model. The last one is also included here for comparison is due to the fact that they also adopted LSTM as a substantial part in their design to capture temporal dynamics. To ensure a fair comparison, we use the same set of evaluation metrics as the baselines: precision (P), recall (R) and F1 score (F1):

$$P = \frac{TP}{TP + FP}, \ R = \frac{TP}{TP + FN}, \ F1 = 2 \cdot \frac{P \cdot R}{P + R},$$

[1] https://github.com/NetManAIOps/OmniAnomaly/tree/master/ServerMachineDataset
[2] https://s3-us-west-2.amazonaws.com/telemanom/data.zip

| Method | SMD | | | MSL | | | Total | | |
|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| LSTM-NDT | 0.5684 | 0.6438 | 0.6037 | 0.5934 | 0.5374 | 0.5640 | 0.5809 | 0.5906 | 0.5839 |
| DAGMM | 0.5951 | 0.8782 | 0.7094 | 0.5412 | 0.9934 | 0.7007 | 0.5681 | 0.9358 | 0.7051 |
| LSTM-VAE | 0.7922 | 0.7075 | 0.7842 | 0.5257 | 0.9546 | 0.6780 | 0.6590 | 0.8311 | 0.7311 |
| **FMUAD** | 0.8574 | 0.8789 | **0.8680** | 0.8311 | 0.8786 | **0.8539** | 0.8443 | 0.8664 | **0.8610** |

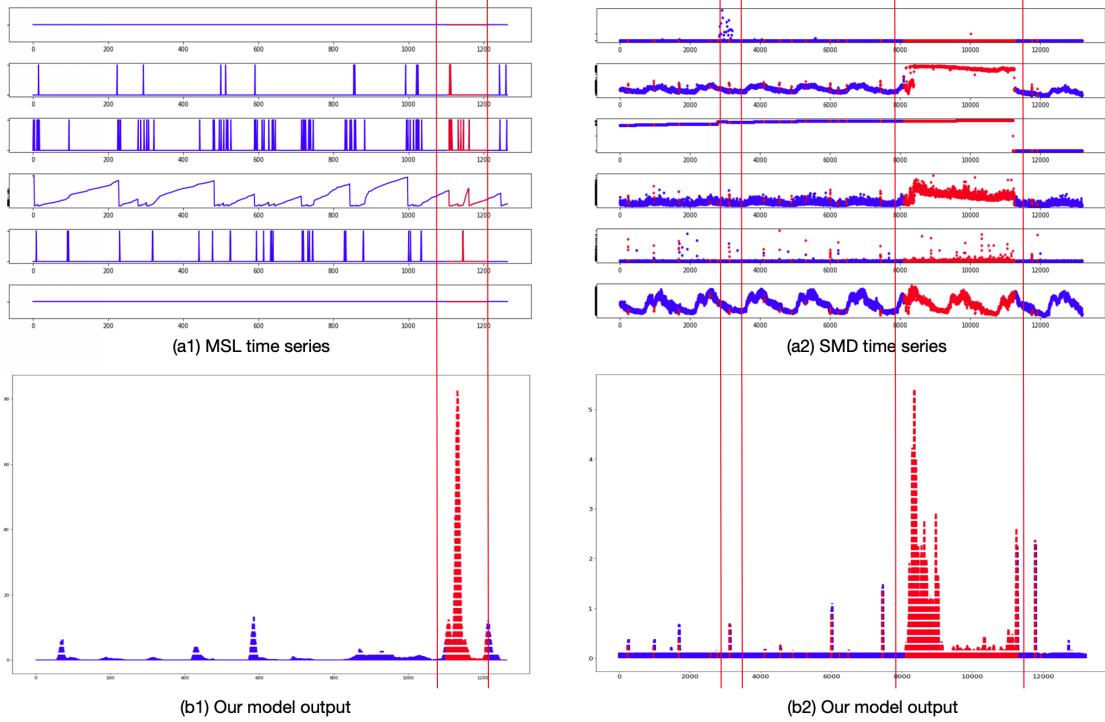TABLE II: Performance of FMUAD and 3 baseline approaches.



Fig. 5: Case studies - (a1, b1): anomaly detection for MSL; (a2, b2): anomaly detection for SMD. (a1, a2) corresponds to the original time series and (b1, b2) corresponds to the anomaly scores output by our framework.

where the F1-score is picked as the highest one via anomaly threshold selection. We also follow the same evaluation strategy by treating the whole window segment as correct if any observation in this segment is correctly detected as an anomaly [20]. Table II details the obtained performance results for our model as well as the three baselines. As one can see, our model keeps a great balance between precision and recall while the baselines tend to favor one of the metrics, and thus resulting to the best F1-score consistently over the baselines. On average, we achieve 17.8% improvement over the best state-of-the-art performance.

### C. Case Study

Besides the performance gains FMUAD brings to the stage, another significant benefit of our model is its intuitiveness. In this section, we take a deep dive in two particular data samples FMUAD labelled as anomalous along with the model's output, to illustrate the detection mechanism.

Shown in Figure 5 are two cases: the left one is taken from MSL and the right one from SMD. In both cases, we sample a time period out of the time series containing anomalous labels,

and the top chart (a1, a2) displays a hand-picked set of 6 metrics containing interesting and indicative patterns. Portions of the time series labelled as anomalies are highlighted in red, and the forecast error (anomaly score) output from the model is displayed in the lower chart (b1, b2).

In the first case study (a1, b1, MSL dataset), our model accurately made a capture on a temporal frequency pattern change. While the value range of the displayed metrics did not change in the anomaly event, metric 3 shows outputs a less-concentrated cluster of spikes, and metric 4 shows significantly more frequent ramps and fluctuations. While these are subtle patterns to human eyes, FMUAD correctly outputs a high anomaly score, indicating a high confidence in the anomaly judgement which matches the provided label.

In the second case study (a2, b2, SMD dataset), we demonstrate the model's capabilities of capturing value changes. The sample presented in (a2) contains two drastically different anomaly patterns: The anomaly on the left contains a short-lived, sparing spikes on metric channel 1. It is a prime example of a point anomaly. In contrast, the right anomaly has a longer duration, involves more metrics and contains a more subtle and

sustained value increase. While it is easy to notice by humans, the statistical significance of the spatial pattern is only highly pronounced when aggregated over the entire period. Again, by design, FMUAD is able to detect both types of anomalies, evidently through a high forecasting error output.

Through the two case studies above, one can easily appreciate the intuitiveness of FMUAD. The model provides explainable and predictable outputs, opening up opportunities for further improvements and targeted fine-tuning based on the patterns of the target dataset. However, another major strength of the model remains to be addressed - its modularity. While we now have empirical evidence that FMUAD can pick up distinct patterns of anomaly, one naturally wonders if and how well the modularity provides a constructive impact, and whether each module indeed captures its target patterns.

### D. Ablation Study

To address the curiosity over the model's modularity, we perform an ablation study to analyze the effect of each detector module to the overall performance of the framework. To recap the model architecture, we use three separate and complementary forecast-based detector modules, each dedicated at one target pattern by design. Specifically, $D^c$ captures inter-series correlation, $D^t$ captures intra-series temporal patterns and $D^s$ captures multi-scale intra-series spatial patterns. We run the same benchmark experiments on all 28 machines of the SMD with only one of the three detectors enabled at a time, and compare the detection performance when that particular detector is present (and more importantly, when the other two are not present).

As with the case study, we sample a few time series each detector proves to perform well on, and present their plots in Figure 6. Table IV presents the F1-score for each of the detectors on each of the samples shown. The result again meets the expectation. In sample (a) where $D^c$ excels, we can see a clear change in the correlation: While in normal situations the two metrics are positively correlated, the labelled anomaly event depicts a negative correlation between the two time series. This is an instance where only $D^c$ is triggered, but not $D^s$ or $D^t$. Samples (b) and (c) general picture the same situation. In (b), a temporal change can be observed as the metric breaks out of its usual periodicity and manifests several sudden noisy, sudden peaks. These patterns are easily detected by $D^t$ but not the others. Finally, in (c), the metric takes a sharp point-anomaly on the left, and a more subtle but longer duration deviation on the right. Both deviations are registered in $D^s$, showing its strength in multi-scale spatial change detection.

| Detector | Machine 1-1 | Machine 1-2 | Machine 2-3 |
|---|---|---|---|
| $D^c$ | 0.6921 | **0.6976** | 0.0332 |
| $D^t$ | **0.9984** | 0.5229 | 0.8968 |
| $D^s$ | 0.9680 | 0.5884 | **0.9122** |

TABLE IV: Performance comparison (F1-score) among three detectors for representative machines from SMD.

Of course, showing detection behavior on archetypal samples does not do justice to the modules' complementary effectiveness. As each detector works well for some cases, one may wonder whether the performance becomes diverse and unstable across datasets when enabling all three detectors. In the following Table V, we present the model performance on each dataset when individual detectors are enabled alone, and when all three are in use. It is evident that combining the three detector modules provides a consistent best performance. More importantly, such combination makes our model performance so stable that only a 0.0001 variance in terms of F1-score remains cross the four set of datasets.

| | SMD | | | MSL | Mean | Variance |
|---|---|---|---|---|---|---|
| | group 1 | group 2 | group 3 | | | |
| $D^c$ | 0.619 | 0.859 | 0.646 | 0.191 | 0.579 | 0.0588 |
| $D^t$ | 0.586 | 0.866 | 0.731 | 0.851 | 0.759 | 0.0127 |
| $D^s$ | 0.475 | 0.723 | 0.560 | 0.323 | 0.520 | 0.0209 |
| **FMUAD** | 0.860 | 0.883 | 0.861 | 0.854 | **0.865** | **0.0001** |

TABLE V: Module VS All - the overall framework outperforms each individual module consistently and exhibits excellent stability with only 0.0001 variance across different datasets.

Besides the modular architecture, a novelty we would like to bring to attention is the batch variance term $\ell_2$ in the loss function. We incorporate this loss in the training time in hopes to learn a more stable and compact representation, therefore further improving the quality of anomaly detection. To test the extent of this improvement brought by $\ell_2$, we repeat the benchmark on each dataset but without the variance term in the loss function Eq. (12). The performance numbers are detailed in Table VI. Based on the numbers, it can be concluded that the variance control term has resulted in a consistent improvement in performance over when only $\ell_1$ is involved in training, albeit by a little.

| | SMD | | | MSL |
|---|---|---|---|---|
| | group 1 | group 2 | group 3 | |
| FMUAD w.o. loss variance control | 0.845 | 0.877 | 0.846 | 0.849 |
| **FMUAD** | **0.860** | **0.883** | **0.861** | **0.854** |

TABLE VI: Model performance (F1-score): Loss VS Loss without variance control.

### V. CONCLUSION

In this paper, we discuss a Forecast-based Multi-aspect framework for Unsupervised multi-variate time series Anomaly Detection, or FMUAD in short. It is a bold attempt at this category of anomaly detection models by identifying the traits of anomaly patterns and categorizing them for a divide-and-conquer approach. We introduce a modular framework where the input time series is transformed into an intermediate representation better capturing each target pattern, and feeding it through each detector module that makes its own forecast on the said intermediate representation. We also experiment with a novel loss function that promotes compactness in the

(a) Machine 1-2: correlation change detected by $D^c$

(b) Machine 1-1: temporal change detected by $D^t$

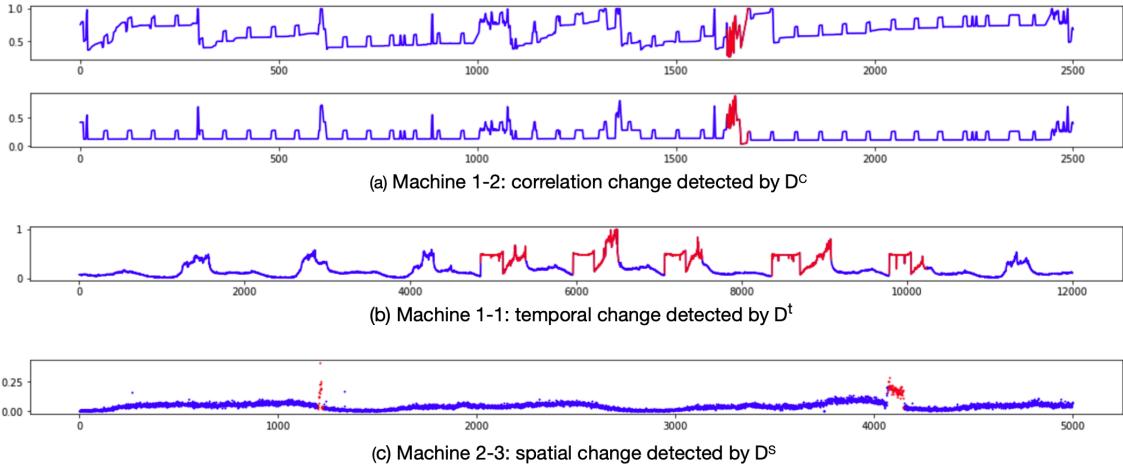(c) Machine 2-3: spatial change detected by $D^s$

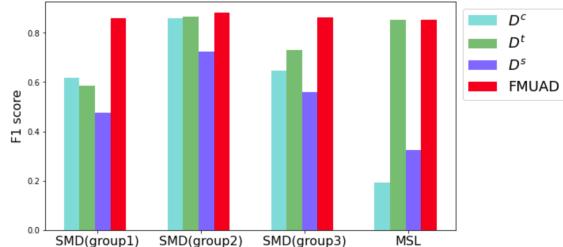Fig. 6: Anomaly patterns captured by different detectors.



Fig. 7: Model performance (F-1 score): Module VS All.

learned representation during training. To wrap them all up, we apply the loss function to compute the forecasting error, and thus produce an anomaly score.

We prove that the three detectors can indeed learn jointly for the optimal anomaly detection by contrasting its performance against state-of-the-art techniques on public reference datasets. The standard F-1 score from FMUAD is indeed higher than other models of its class. We also relate the detection mechanisms back to human intuition and find proof of each detector module performing to its advantage by agreeing to the anomaly patterns we assign them to detect. Furthermore, we investigate the role of modularity through an ablative study, and find that the combination of all three modules can achieve consistently higher and more stable performance than each one alone. We also prove that the compactness-focused $l_2$ loss term results in consistent improvement, albeit being minute compared to other improvements.

A recent trend we observe in the machine learning community is the tendency to move towards "end-to-end" or "one-size-fits-all" models, as a defiant departure from the old-school feature-engineering approaches. However, FMUAD may serve as a reminder that distinct patterns existing in the input data may still need a degree of tailored model design, all while being a self-adaptive model that remains end-to-end trainable.

This is clearly not the end of our journey, as the flexibility of FMUAD opens up a world of opportunities for further enhancements. For one, we can continue to discover new types of anomaly pattern and tailor a corresponding detector to its characteristics so those patterns are effectively captured; for two, there may exist better implementations for detectors $D^c$, $D^t$ and $D^s$ that can allow individual modules to perform even more optimally. Finally, the three detectors in FMUAD have operated in complete isolation of each other and the final aggregator. Enabling cross-interactions between these modules may also result in interesting observations that are worth investigating. All in all, we hope our model has availed additional potential for forecast-based anomaly detection models, and it can start a discussion to draw further work onto this general direction.

REFERENCES

[1] J. Wei, J. Zhao, Y. Zhao, and Z. Zhao, "Unsupervised anomaly detection for traffic surveillance based on background modeling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 129–136, 2018.

[2] M. Ahmed, A. N. Mahmood, and M. R. Islam, "A survey of anomaly detection techniques in financial domain," *Future Generation Computer Systems*, vol. 55, pp. 278–288, 2016.

[3] P. Ghuli, "Anomaly detection in log records," *Indonesian Journal of Electrical Engineering and Computer Science*, vol. 10, p. 343, Apr. 2018.

[4] Q. Cao, Y. Qiao, and Z. Lyu, "Machine learning to detect anomalies in web log analysis," in *2017 3rd IEEE International Conference on Computer and Communications (ICCC)*, pp. 519–523, 2017.

[5] K. Ding, J. Li, and H. Liu, "Interactive anomaly detection on attributed networks," in *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, pp. 357–365, 2019.

[6] Z. Hang, W. Yujing, D. Juanyong, H. Congrui, C. Defu, T. Yunhai, X. Bixiong, B. Jing, T. Jie, and Z. Qi, "Multivariate time-series anomaly detection via graph attention network," 09 2020.

[7] J. Gao, X. Song, Q. Wen, P. Wang, L. Sun, and H. Xu, "Robusttad: Robust time series anomaly detection via decomposition and convolutional neural networks," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, Aug. 2020.

[8] D. Cao, Y. Wang, J. Duan, C. Zhang, X. Zhu, C. Huang, Y. Tong, B. Xu, J. Bai, J. Tong, and Q. Zhang, "Spectral temporal graph neural network for multivariate time-series forecasting," 2021.

[9] P. Perera, P. Oza, and V. M. Patel, "One-class classification: A survey," *CoRR*, vol. abs/2101.03064, 2021.

[10] W. Wu, L. He, W. Lin, Y. Su, Y. Cui, C. Maple, and S. A. Jarvis, "Developing an unsupervised real-time anomaly detection scheme for time series with multi-seasonality," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2020.

[11] 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.

[12] Y. Guo and F. Farooq, "Proceedings of the 24th acm sigkdd international conference on knowledge discovery and data mining," KDD'18, (New York, NY, USA), Association for Computing Machinery, 2018.

[13] R. Chalapathy and S. Chawla, "Deep learning for anomaly detection: A survey," *arXiv preprint arXiv:1901.03407*, 2019.

[14] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection," *ACM Computing Surveys*, vol. 41, pp. 1–58, July 2009.

[15] G. Spadon, S. Hong, B. Brandoli, S. Matwin, J. F. Rodrigues-Jr, and J. Sun, "Pay attention to evolution: Time series forecasting with deep graph-evolution learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2021.

[16] Y. Wu, M. Gu, L. Wang, Y. Lin, F. Wang, and H. Yang, "Event2graph: Event-driven bipartite graph for multivariate time-series anomaly detection," *CoRR*, vol. abs/2108.06783, 2021.

[17] H. Chen, L. Wang, Y. Lin, C.-C. M. Yeh, F. Wang, and H. Yang, "Structured graph convolutional networks with stochastic masks for recommender systems," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, (New York, NY, USA), p. 614–623, Association for Computing Machinery, 2021.

[18] H. Chen and J. Li, "Exploiting structural and temporal evolution in dynamic link prediction," in *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pp. 427–436, 2018.

[19] 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," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 1409–1416, July 2019.

[20] 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 and Data Mining*, ACM, July 2019.

[21] J. Audibert, P. Michiardi, F. Guyard, S. Marti, and M. A. Zuluaga, "USAD: UnSupervised Anomaly Detection on Multivariate Time Series," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, July 2020.

[22] T. Zhao, B. Ni, W. Yu, and M. Jiang, "Early anomaly detection by learning and forecasting behavior," ACM, 2021.

[23] M.-C. Lee, J.-C. Lin, and E. G. Gran, "Salad: Self-adaptive lightweight anomaly detection for real-time recurrent time series," IEEE, 2021.

[24] E. H. M. Pena, M. V. O. de Assis, and M. L. Proença, "Anomaly detection using forecasting methods arima and hwds," in *2013 32nd International Conference of the Chilean Computer Science Society (SCCC)*, pp. 63–66, 2013.

[25] J. M. Torres, P. G. Nieto, L. Alejano, and A. Reyes, "Detection of outliers in gas emissions from urban areas using functional data analysis," *Journal of Hazardous Materials*, vol. 186, pp. 144–149, Feb. 2011.

[26] S. Ahmad, A. Lavin, S. Purdy, and Z. Agha, "Unsupervised real-time anomaly detection for streaming data," *Neurocomputing*, vol. 262, pp. 134–147, Nov. 2017.

[27] N. Ding, H. Gao, H. Bu, H. Ma, and H. Si, "Multivariate-time-series-driven real-time anomaly detection based on bayesian network," *Sensors*, vol. 18, p. 3367, Oct. 2018.

[28] 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 and Data Mining*, ACM, July 2018.

[29] D. Song, N. Xia, W. Cheng, H. Chen, and D. Tao, "Deep r -th root of rank supervised joint binary embedding for multivariate time series retrieval," KDD'18, (New York, NY, USA), p. 2229–2238, Association for Computing Machinery, 2018.

[30] X. Shi, Z. Chen, H. Wang, D.-Y. Yeung, W.-k. Wong, and W.-c. Woo, "Convolutional lstm network: A machine learning approach for precipitation nowcasting," in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*, NIPS'15, (Cambridge, MA, USA), p. 802–810, MIT Press, 2015.

[31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017.

[32] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *CoRR*, vol. abs/1511.07122, 2016.

[33] P. Perera and V. M. Patel, "Learning deep features for one-class classification," *IEEE Transactions on Image Processing*, vol. 28, pp. 5450–5463, Nov. 2019.

[34] D. Park, Y. Hoshi, and C. C. Kemp, "A multimodal anomaly detector for robot-assisted feeding using an lstm-based variational autoencoder," *CoRR*, vol. abs/1711.00614, 2017.