Multivariate Time Series Anomaly Prediction Based on Forecasting and Reconstruction Using Transformer with Temporal and Feature-wise Attention

Chihiro Maru⊠¹, Masato Oguchi², and Ichiro Kobayashi²

¹ Faculty of Science and Engineering, Chuo University, Tokyo, Japan cmaru671@g.chuo-u.ac.jp

² Graduate School of Humanities and Sciences, Ochanomizu University, Tokyo, Japan {oguchi, koba}@is.ocha.ac.jp

Abstract. Anomaly detection has been actively studied, enabling the high-accuracy detection of anomalies. However, because anomaly detection assumes that an anomaly has already occurred, detecting future anomalies before they occur and preventing them from happening is impossible. Therefore, we develop a Transformer-based Anomaly Prediction (TranAP) method, which is designed to detect future anomalies. TranAP predicts future values from previous time series and uses reconstruction techniques to detect signs of anomalies using the predicted results. Detecting these precursors requires a correct understanding of the temporal characteristics of the multivariate time series (MTS). Because the timing of behavior leading to an anomaly may differ for each feature, we apply multi-head attention (ATTN) in the time dimension for each feature. Additionally, TranAP captures the dependencies between different features that the conventional ATTN could not. Because the effect of ATTN is partially diminished within the attention block, even after improvement to capture detailed information in MTS, we modify the operation of the block to preserve this effect. We demonstrate the effectiveness of TranAP by comparing it with state-of-the-art models. This improved attention mechanism of TranAP allows for a better understanding of behavior that leads to anomalies.

Keywords: anomaly prediction, multivariate time series forecasting, reconstruction, Transformer, attention block, multi-head attention

1 Introduction

Anomaly detection has been extensively studied and has demonstrated high performance. An anomaly, also known as an outlier or novelty, refers to an unusual, irregular, inconsistent, unexpected, rare, faulty, or simply a strange observation, depending on the context. Anomaly detection aims at identifying unexpected patterns or data points in real-world applications. Anomaly detection for multivariate time series (MTS) requires handling time series with several



sor.

Fig. 1: Example of an anomaly precur- Fig. 2: Univariate time series for each feature.

features, and numerous deep learning (DL)-based models have been proposed to address this task [26, 16]. Most of these models focus on accurately detecting anomalies that have already occurred, whereas detecting future anomalies before they manifest is increasingly expected. Anomaly prediction is the process of identifying current patterns or signs that may indicate upcoming abnormal events [12]. The goal is to detect these precursors before the occurrence of anomalies, thereby enabling proactive preventive actions.

Figure 1 illustrates the transition of data points for each feature of an MTS obtained from a real-world system. An abnormal event occurs in the anomaly part owing to an external attack that sets the value of feature 1 to 700. However, this anomaly does not occur immediately after the attack. A time lag (the anomaly precursor part in Figure 1) exists between the attack and the actual occurrence of a critical abnormal event, which is recognized as an anomaly. During this period, the effects of the attack spill over to features other than feature 1, which was initially attacked, and signs triggering an abnormal event can be observed in several features.

PAD proposed an anomaly prediction model that requires anomalous data and an anomaly detection model for training [12]. We propose a Transformer-based Anomaly Prediction (TranAP) method that uses only normal data for training and does not require an anomaly detection model. TranAP uses MTS forecasting, which predicts future values from previous values and determines whether an anomaly will occur in the future by reconstructing the prediction results.

Anomaly prediction requires an accurate understanding of the temporal characteristics of the MTS. The behavior leading to an anomaly may occur at different times for each feature, as shown in Figure 2. However, general multihead attention (ATTN) cannot capture each characteristic. Furthermore, because MTS anomaly prediction targets data with several features, more detailed MTS characteristics can be extracted by capturing the dependencies between the features. Therefore, we enhance ATTN to better capture the different temporal dependencies for each feature as well as the dependencies between features.

In addition, the attention block comprises ATTN, residual connection (RES) [9], and layer normalization (LN) [5], which contribute to the Transformer performance [32]. Previous studies on natural language processing (NLP) revealed that other components cancel the effects of ATTN [13]. To the best of our knowledge, studies on the attention block of the MTS are yet to be conducted. Further research is needed to confirm whether the same observations can be made in MTS as in natural language.

The contributions of this study are as follows:

- We proposed a novel framework specialized for anomaly prediction tasks by utilizing MTS forecasting and reconstruction. Unlike the conventional anomaly prediction model that requires anomalous data and an anomaly detection model for training, the proposed framework does not need them.
- We applied ATTN in the time direction for each feature and also in the feature direction to focus on the behavior leading to anomalies in MTS.
- We found that RES partially cancels the effect of ATTN in the attention block. Based on this finding, we modified the attention block to preserve the strong effect of the improved attention mechanism.
- We evaluated TranAP on five real-world datasets, demonstrating that anomaly prediction using MTS forecasting and reconstruction is effective. Various experiments showed that improving attention mechanisms helps capture the detailed characteristics of MTS.

2 Related Work

Anomaly Detection in MTS. Because anomaly detection requires handling time series with multiple features, numerous anomaly detection models using DL have been proposed [38, 34, 31, 35]. Most of these models focus on detecting anomalies that have already occurred, and cannot detect future anomalies.

Some anomaly detection models use techniques such as autoregression and reconstruction. These models detect anomalies by predicting or reconstructing data points within a given input and comparing them with actual values. TranAP is similar to the aforementioned models because it performs anomaly predictions based on forecasting and reconstruction. However, TranAP predicts future unseen MTS whose actual values are unknown from the given input, and thus cannot perform comparisons during anomaly prediction. Therefore, we adopt a framework to evaluate the results of MTS forecasting by reconstruction (see Section 3.2).

Anomaly Prediction. The currently proposed anomaly prediction model called PAD [12] uses training data consisting of normal and pseudo-anomalous data. In the training phase, PAD requires anomaly detection and prediction models; the latter model is trained to imitate the results of the former. Specifically, the anomaly prediction model receives the MTS at timestep t = T - 1 to predict whether an anomaly will occur at t = T. It is trained to predict the same outcome as the anomaly detection result at t = T obtained by the anomaly detection model. However, it is difficult to prepare all combinations of anomalies and their precursors, and the trained anomaly prediction model has difficulty detecting



Fig. 3: TranAP architecture. The proposed model consists of Transformer-based forecasting and reconstruction.

precursors correctly for unknown anomalies that are not included in the training data [6, 36]. Furthermore, the anomaly prediction model requires an anomaly detection model for training, which makes it difficult to train them efficiently.

Transformer Attention Block Analysis. The ATTN, a key component of Transformer, has been analyzed in several studies on NLP [2, 14, 27, 19, 21, 10]. Transformer is composed of not only ATTN but also these components such as RES and LN. Previous studies have shown that other components cancel the effect of ATTN [13]. However, most Transformer-based models have only improved the ATTN and have not considered the operation of the entire attention block, e.g., [25, 37, 34, 31, 35].

MTS forecasting using Transformer. MTS forecasting is the task of predicting future MTS values from previous values [18, 30, 22]. Transformer-based forecasting models have been proposed [39, 33, 20, 40, 7, 25]. These models improve ATTN and focus on efficiently extracting long-term time dependencies with less computational complexity. Because the MTS has several features, more detailed information about the data can be extracted by capturing the dependencies between features. Crossformer [37] is a Transformer-based model for MTS forecasting that improves ATTN to capture the dependencies between features.

3 METHODOLOGY

3.1 Problem Statement

Consider the MTS \mathcal{X} consisting of M data points $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_M\}$, where each data point $\mathbf{x}_t \in \mathbb{R}^D$ is collected at a certain timestep t. D(D > 1) is the number of features in the MTS. We adopt a window-based approach in the anomaly prediction task, similar to [12]. That is, \mathcal{X} is divided into a set of windows of input length T such as $\{\mathbf{x}_{1:T}, \mathbf{x}_{1+step_size:T+step_size}, \ldots, \mathbf{x}_{M-T+1:M}\}$, and the input to the model is in window units.

The problem with anomaly prediction is that, given an input window of input length T, we predict whether an anomaly will occur during the next τ timesteps. For example, given an input window $\mathbf{x}_{1:T} \in \mathbb{R}^{T \times D}$ from timestep 1 to T, we predict whether anomalies may occur between timestep T + 1 and $T + \tau$; that is, whether an unseen $\mathbf{x}_{T+1:T+\tau} \in \mathbb{R}^{\tau \times D}$ contains anomalies. Here, τ denotes the prediction length. In this case, we output $y_{T+1:T+\tau}^{pred} \in \{0,1\}$ (where 1 denotes an anomalous window) for each given window $\mathbf{x}_{1:T}$ as the anomaly prediction result for $\mathbf{x}_{T+1:T+\tau}$. This problem is evaluated for each testing window as in [12]. Given the MTS $\hat{\mathcal{X}}$ for testing, it is divided into a set of windows of input length T as in training. For example, given $\hat{\mathbf{x}}_{1:T}$, we can predict whether an unseen future $\hat{\mathbf{x}}_{T+1:T+\tau}$ is an anomalous window using the trained model. The correct label used for the evaluation is $\hat{y}_{T+1:T+\tau} \in \{0,1\}$, and if $\{\hat{y}_{T+1}, \hat{y}_{T+2}, \dots, \hat{y}_{T+\tau}\}$ contains one or more abnormal labels, we denote $\hat{y}_{T+1:T+\tau}$ as one.

In the following, we denote an input window as $\mathbf{x}_{1:T}$ and a predicted window as $\mathbf{x}_{T+1:T+\tau}$ for simplicity.

3.2 Model Structure

TranAP mainly consists of Transformer-based MTS forecasting and reconstruction, as illustrated in Figure 3. We use only normal data for training [8, 3]. TranAP enables the detection of any combination of anomalies and their precursors not included in the training dataset. Because the model has already been trained with normal data, it can successfully predict future values when a normal window is given. However, given a window that deviates from normality, the results will not be correct predictions of future trends. Such deviations can be detected as precursors of an anomaly. Similar to TranAP, some anomaly detection models use prediction techniques. However, while these models predict values within an input window, allowing for comparison with actual values, TranAP predicts a future unseen window and thus cannot be compared with actual values. Therefore, we utilize the Transformer for reconstruction to evaluate the results of MTS forecasting. The MTS, which concatenates the input and predicted values, is reconstructed. It is also trained using only normal data, enabling successful reconstruction when it receives a predicted MTS with normality. However, when it receives a predicted MTS that deviates from normality, the reconstruction fails, and the window is determined as a precursor of an anomaly.

3.3 Segmentation of Time Series

In MTS tasks, the segmentation of the input time series contributes to the accuracy of each task [37, 25]. Figure 4 shows the interactions between data points when a window is given to the trained TranAP. The (i, j) cell indicates the extent to which the *j*th data point on the key side contributes to the computation of the output corresponding to the *i*th data point on the



Fig. 4: Interactions between data points in each layer.

query side. From Figure 4, we observe that the interactions tend to be divided into

segments, particularly after the second layer. As the characteristics of temporally close data points in the time series are similar, they have similar interactions. Moreover, aggregating information from multiple data points into segments reduces the computational complexity while maintaining the accuracy of the anomaly prediction (see Section 5.1). We divide the input window into segments when performing MTS forecasting and reconstruction.

3.4 Temporal and Feature-wise ATTN

We perform ATTN in the time dimension (temporal ATTN) and feature dimension (feature-wise ATTN) in Transformer-based MTS forecasting and reconstruction. **Temporal ATTN.** General Transformer-based models capture the temporal dependencies between input representations by performing ATTN in the time dimension. This computation fails to capture the temporal characteristics of each feature because all features of each input representation share the same attention map and the information of all features at a timestep is aggregated into a single embedding.

Because the timing of behaviors leading to an anomaly differs for each feature, capturing each temporal characteristic is essential for anomaly prediction. Therefore, we performe ATTN on each feature separately in the univariate time series.

Temporal ATTN receives $\mathbf{H} \in \mathbb{R}^{N \times D \times d_{\text{model}}}$ as input, where N is the number of segments. Note that **H** is a vector after trainable linear projection added with a positional embedding. We define $\mathbf{H}_{:,d}$ as a vector of all segments with feature $d(1 \le d \le D)$. After temporal ATTN (ATTN^{time}), we obtain the output \mathbf{H}^{time} :

$$\hat{\mathbf{H}}_{:,d}^{\text{time}} = \text{LN}\left(\text{ATTN}^{\text{time}}\left(\mathbf{H}_{:,d}, \mathbf{H}_{:,d}, \mathbf{H}_{:,d}\right) + \mathbf{H}_{:,d}\right), \\
\mathbf{H}^{\text{time}} = \text{LN}\left(\text{FF}\left(\hat{\mathbf{H}}^{\text{time}}\right) + \hat{\mathbf{H}}^{\text{time}}\right),$$
(1)

where FF denotes the feedforward network.

Feature-wise ATTN. Temporal ATTN alone does not capture feature-wise dependencies. We apply ATTN in the feature dimension (ATTN^{feature}) after performing temporal ATTN:

$$\hat{\mathbf{H}}_{i,:}^{\text{feature}} = \text{LN}\left(\text{ATTN}^{\text{feature}}\left(\mathbf{H}_{i,:}^{\text{time}}, \mathbf{H}_{i,:}^{\text{time}}, \mathbf{H}_{i,:}^{\text{time}}\right) + \mathbf{H}_{i,:}^{\text{time}}\right), \\
\mathbf{H}^{\text{feature}} = \text{LN}\left(\hat{\mathbf{H}}^{\text{feature}} + \text{FF}\left(\hat{\mathbf{H}}^{\text{feature}}\right)\right),$$
(2)

where $\mathbf{H}_{i,:}$ is a vector of all features of the *i*th $(1 \le i \le N)$ segment.

3.5 Effect of ATTN

Although many Transformer-based models have improved the attention mechanism, [13] reported that RES cancels the effect of ATTN in NLP. Therefore, we investigate the operation in the attention block when dealing with the MTS.



(a) ATTN mixes input representations (b) RES preserves the original information. other than its own.

Fig. 5: Interactions between representations in each layer.

Transformer consists of layers with an attention block. The attention block comprises three components: ATTN, RES, and LN.

$$\widetilde{\mathbf{H}} = \mathrm{LN}\underbrace{(\mathrm{ATTN}\,(\mathbf{H},\mathbf{H},\mathbf{H}) + \mathbf{H})}_{\mathrm{RES}},\tag{3}$$

where $\mathbf{H} := [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_L] \in \mathbb{R}^{L \times d_{\text{model}}}$ is the sequence of input representations, and $\mathbf{h}_i \in \mathbb{R}^{d_{\text{model}}}$ is the *i*th input representation. $\widetilde{\mathbf{H}} := [\widetilde{\mathbf{h}}_1, \widetilde{\mathbf{h}}_2, \dots, \widetilde{\mathbf{h}}_L] \in \mathbb{R}^{L \times d_{\text{model}}}$ is the sequence of output representations, and $\widetilde{\mathbf{h}}_i \in \mathbb{R}^{d_{\text{model}}}$ is the output corresponding to \mathbf{h}_i .

Among these components, ATTN and RES have contrasting effects on the computation of output representations. While ATTN mixes the input representations, RES preserves the original input representations.

We can visualize the interactions between the representations after ATTN and RES in each layer when a window is provided to the trained TranAP in Figure 5. The (i, j) cell indicates how strongly the key input $\mathbf{h}_j \in {\mathbf{h}_1, \mathbf{h}_2, \ldots, \mathbf{h}_L}$ contributes to computing the query output $\tilde{\mathbf{h}}_i$. The diagonal elements correspond to the effect of preserving the original input information, that is, preserving $\mathbf{h}_{j=i}$ when computing $\tilde{\mathbf{h}}_i$. ATTN in Figure 5a mixes information from input representations other than its own. On the other hand, RES in Figure 5b loses the mixing effect of ATTN and strongly preserves the information from the original input representation. These results indicate that RES cancels the effect of ATTN on the MTS.

We modify Eq. (3) to adjust for the mixing effect of ATTN and the preservation effect of RES:

$$\mathbf{H} = \mathrm{LN}\left(\mathrm{ATTN}\left(\mathbf{H}, \mathbf{H}, \mathbf{H}\right) + \lambda \mathbf{H}\right),\tag{4}$$

where $\lambda(0 \leq \lambda \leq 1)$ is a parameter that adjusts each effect. The lower λ , the stronger the influence of ATTN.

We calculate the mixing ratio r_i for ATTN, RES, and LN, which represents the ratio of the mixing effect to the sum of the mixing and preservation effects. A higher mixing ratio indicates that the mixing effect is stronger than the preservation effect. Table 1 shows the mixing ratio after performing ATTN, RES,

	0		
Components	$\lambda = 0$	$\lambda=0.25$	$\lambda = 1$
Multi-head attention	95.8	96.0	95.9
Residual connection	91.7	84.5	74.2
Layer normalization	91.8	84.5	74.2

Table 1: Mean value of the mixing ratio for each component.

and LN with $\lambda = 0, 0.25, 1$. The values in Table 1 represent the mean ratios of the heads and layers.

These results demonstrate that, while general Transformer-based models (i.e., $\lambda = 1$) lose the mixing effect of ATTN after RES, increasing the influence of ATTN in the attention block can preserve the mixing effect even after RES. We set λ to preserve a strong mixing effect.

3.6 Anomaly Prediction Flow

Here, we define the Transformer for forecasting operations as TF^{pred} and the Transformer for reconstruction operations as TF^{reconst}.

Training. TranAP is trained using only normal data. Given an input window $\mathbf{x}_{1:T}$, the Transformer for forecasting predicts the following future τ timesteps $\mathbf{x}_{T+1:T+\tau}$:

$$\mathbf{x}_{T+1:T+\tau}^{\text{pred}} = \text{TF}^{\text{pred}}(\mathbf{x}_{1:T}),\tag{5}$$

where $\mathbf{x}_{T+1:T+\tau}^{\text{pred}}$ denotes the predicted values. We utilize the mean squared error (MSE) to compute the difference between the predicted values and ground truth during the training phase of the Transformer for forecasting. The Transformer is trained to minimize the following objective function:

$$\mathcal{L}^{\text{pred}} = \left\| \mathbf{x}_{T+1:T+\tau}^{\text{pred}} - \mathbf{x}_{T+1:T+\tau} \right\|_{2}^{2}.$$
 (6)

The Transformer for reconstruction reconstructs a vector $\mathbf{x}_{1:T} \oplus \mathbf{x}_{T+1:T+\tau}^{\text{pred}} \in \mathbb{R}^{(T+\tau) \times D}$ that combines the original input $\mathbf{x}_{1:T}$ and the predicted $\mathbf{x}_{T+1:T+\tau}^{\text{pred}}$:

$$\mathbf{x}_{1:T+\tau}^{\text{reconst}} = \text{TF}^{\text{reconst}}(\mathbf{x}_{1:T} \oplus \mathbf{x}_{T+1:T+\tau}^{\text{pred}}),\tag{7}$$

where \oplus is the operation of the concatenation of two vectors, and $\mathbf{x}_{1:T+\tau}^{\text{reconst}}$ represents the reconstructed values. The Transformer for reconstruction is trained to minimize the following objective function to reconstruct values similar to the input:

$$\mathcal{L}^{\text{reconst}} = \left\| \mathbf{x}_{1:T+\tau}^{\text{reconst}} - \mathbf{x}_{1:T+\tau} \right\|_{2}^{2}.$$
(8)

Anomaly Prediction. The trained Transformer for forecasting receives an input window $\hat{\mathbf{x}}_{1:T}$ and predicts the future $\hat{\mathbf{x}}_{T+1:T+\tau}^{\text{pred}}$. Subsequently, the trained Transformer for reconstruction reconstructs $\hat{\mathbf{x}}_{1:T} \oplus \hat{\mathbf{x}}_{T+1:T+\tau}^{\text{pred}}$:

$$\hat{\mathbf{x}}_{T+1:T+\tau}^{\text{pred}} = \text{TF}^{\text{pred}}(\hat{\mathbf{x}}_{1:T}), \\
\hat{\mathbf{x}}_{1:T+\tau}^{\text{reconst}} = \text{TF}^{\text{reconst}}(\hat{\mathbf{x}}_{1:T} \oplus \hat{\mathbf{x}}_{T+1:T+\tau}^{\text{pred}}),$$
(9)

where $\hat{\mathbf{x}}_{1:T+\tau}^{\text{reconst}}$ denotes the reconstructed values. The anomaly prediction score $\mathcal{A}_{AP}(\hat{\mathbf{x}}_{T+1:T+\tau}|\hat{\mathbf{x}}_{1:T})$ is defined as

$$\mathcal{A}_{\mathrm{AP}}(\hat{\mathbf{x}}_{T+1:T+\tau}|\hat{\mathbf{x}}_{1:T}) = \left\| \hat{\mathbf{x}}_{1:T+\tau}^{\mathrm{reconst}} - \hat{\mathbf{x}}_{1:T} \oplus \hat{\mathbf{x}}_{T+1:T+\tau}^{\mathrm{pred}} \right\|_{2}^{2}.$$
 (10)

 $\mathcal{A}_{AP}(\hat{\mathbf{x}}_{T+1:T+\tau}|\hat{\mathbf{x}}_{1:T})$ means the degree to which an anomaly can occur in future $\hat{\mathbf{x}}_{T+1:T+\tau}$ (i.e., $\hat{\mathbf{x}}_{1:T}$ can be a precursor of an anomaly) given an input window $\hat{\mathbf{x}}_{1:T}$. A window $\hat{\mathbf{x}}_{1:T}$ with an anomaly prediction score $\mathcal{A}_{AP}(\hat{\mathbf{x}}_{T+1:T+\tau}|\hat{\mathbf{x}}_{1:T})$ that exceeds a predefined threshold is determined to be a precursor of an anomaly (i.e., $y_{T+1:T+\tau}^{\text{pred}} = 1$). Then, the result of the anomaly prediction $y_{T+1:T+\tau}^{\text{pred}} \in \{0, 1\}$ is compared with the correct label $y_{T+1:T+\tau} \in \{0, 1\}$ to evaluate the success or failure of the anomaly prediction.

4 Experiments

4.1 Anomaly Prediction in MTS

Datasets. We assess the performance of the proposed TranAP on five real-world datasets: SWaT [23], PSM [1], SMD [29], SMAP [11], and NIPS-TS-GECCO [15, 24].

Baselines. We compare TranAP with nine anomaly detection models and one anomaly prediction model **PAD**.

Anomaly detection models consist of reconstruction-based models: LSTM-AE [28], MAD-GAN [17], USAD [4], CAE-M [38], TranAD [31], Anomaly Transformer [34], and DCdetector [35]; the autoregression-based model LSTM [11]; and the density-estimation model DAGMM [41]. Anomaly detection models can be applied to anomaly prediction tasks by considering the precursors of anomalies as an anomaly (see Appendix A.3).

Experimental Settings. All models follow the experimental setup with an input length T = 48 and prediction lengths $\tau \in \{24, 36, 48, 72, 96\}$. The attention ratio λ is set to 0.5 for the SMD and SMAP datasets and 0.25 for the other datasets. The segment length L_{seg} is set to 24 for the SMD and NIPS-TS-GECCO datasets and 12 for the other datasets. We choose the F1-score as evaluation metrics to compare the performance of TranAP with those of the other models. All experiments are repeated five times and the mean of the metrics is reported.

4.2 Main Results

We first evaluate the anomaly prediction performance, as shown in Table 2. Overall, we achieve state-of-the-art results for almost all datasets. The mean F1-score over all the datasets is 8.1 points higher than that of the baselines. The F1-score is 2.5–6.8 points higher than that of the baselines for datasets with relatively clear precursors of anomalies, such as the SWaT and PSM datasets. On the other hand, for datasets with small anomaly precursors, such as the SMD and SMAP datasets, the F1-score is comparable to that of the baselines. The

NIPS-TS-GECCO dataset is a challenging dataset with several types of anomalies. However, TranAP is successful, whereas the baselines fail to predict anomalies. The standard deviations of the F1-score of TranAP in the five repetitions of the experiments are within 0.14%-4.09% for all datasets.

Despite a fixed input length and different prediction lengths, the F1-score does not change significantly. In actual operations, it is desirable to detect the precursors of anomalies further into the future. Therefore, the prediction length should be increased to the extent that the evaluation metrics, such as the F1-score, do not change significantly.

5 Analysis

5.1 Ablation Study

Input Length.

We investigate the effect of input length by changing the input lengths $T \in$ $\{12, 24, 48, 96, 168\}$ with a fixed prediction length $\tau = 96$ in Figure 6. The F1-score is the highest when T = 12 for the SWaT dataset and T = 24for the PSM dataset. In these datasets, the behavior leading to anomalies is likely to have



Fig. 6: F1-score when changing input lengths.

occurred at slightly earlier time steps when the prediction length is 96. The longer the input length, the lower the value of the F1-score. If the input length is too long, the input contains information that is irrelevant to anomaly prediction.

Effect of Segment.

We investigate the average running time per iteration and F1score for different segment lengths in the SWaT and PSM datasets in Figure 7. We set an input length T =48, prediction length $\tau =$ 48, and segment lengths $L_{seg} \in$



Fig. 7: Running time and F1-score when changing segment lengths.

 $\{1, 3, 6, 12, 24, 48\}$. The average running time per iteration is reduced by 55.4%–81.5% without decreasing the F1-score. Note that the F1-score is the highest with $L_{\text{seg}} = 24$ for the SWaT dataset and $L_{\text{seg}} = 6$ for the PSM dataset. In the PSM dataset with $L_{\text{seg}} = 48$, the F1-score decreases because the time steps with different temporal characteristics are grouped into a single segment.

evaluation metrics are the precision (P), recall (R), and F1-score (F1). The best results are in **bold**, and the second best are Table 2: MTS anomaly prediction results. We use an input length T = 48 and prediction lengths $\tau \in \{24, 36, 48, 72, 96\}$. The

11

Attentio	on mechanism	Original	F-ATTN	T-ATTN	TF-ATTN
SWaT	24	83.2	86.3	85.4	86.9
	96	85.3	86.4	86.6	86.8
PSM	24	88.4	88.7	90.2	89.2
	96	87.3	87.3	88.7	88.5
SMD	24	38.6	53.6	53.3	55.9
	96	44.1	51.4	53.1	54.7
SMAP	24	59.3	62.5	61.9	63.1
	96	61.0	62.6	63.3	<u>63.0</u>
GECCO	24	50.0	50.5	50.7	51.3
	96	57.3	58.9	63.0	63.2

Table 4: F1-score when changing attention mechanisms.

Effect of ATTN.

We examine the impact of the mixing effect of ATTN and the preservation effect of RES on the anomaly prediction performance by adjusting the attention ratio λ in Eq. (4). Table 3 denotes the F1-score with an input length T = 48, prediction length $\tau = \{24, 96\}$, and attention ratio $\lambda =$

Table 3: F1-score when changing attention ratios.

Attention ratio		0	0.25	1
SWaT	24	86.8	86.9	86.6
	96	86.5	86.8	86.6
PSM	24	88.0	89.2	87.3
	96	87.2	88.5	88.0

 $\{0, 0.25, 1\}$ in the SWaT and PSM datasets. Although increasing the mixing effect improves the F1-score, completely eliminating the preservation effect (i.e., $\lambda = 0$) tends to degrade the performance of anomaly prediction. Therefore, the increased mixing effect of ATTN leads to improved performance, and the preservation effect of RES also contributes to anomaly prediction in the attention block.

Effect of Temporal and Feature-wise ATTN. We perform anomaly prediction on all datasets using four different attention mechanisms: general (**Original**), feature-wise (**F-ATTN**), temporal (**T-ATTN**), and temporal and feature-wise ATTN (**TF-ATTN**).

Table 4 denotes the F1-score with an input length T = 48 and prediction lengths $\tau \in \{24, 96\}$. In all cases, the F1-score is higher than that of the other attention mechanisms for **Original**. **F-ATTN** achieves a higher F1-score than **Original** for all prediction lengths, indicating that it is important to reflect the dependencies between features in the final computed representations. The longer the prediction length, the more **T-ATTN**, which treats each feature independently, functions. Therefore, each feature has a different temporal dependency, which is important for anomaly prediction. This indicates that performing **TF-ATTN** enables detailed extraction of the MTS characteristics and contributes to improving anomaly prediction.



Fig. 8: Attention maps after temporal ATTN.

5.2 Visualization of Attention Maps

We visualize the interactions between representations after temporal and feature-wise ATTN to confirm the effectiveness of the improved attention mechanism. A window $\hat{\mathbf{x}}_{48:71}$ from the PSM dataset with an



⁷¹ Fig. 9: Attention maps after feature-wise ATTN.

input length T = 24 and segment length $L_{seg} = 6$ which contains the precursors of the anomalies (yellow highlighted area of Figure 8 to the right) is fed to the trained encoder. An anomaly occurs in $\hat{\mathbf{x}}_{72:95}$ (red highlighted area of Figure 8 to the right).

Figure 8 shows the attention maps of $\hat{\mathbf{x}}_{48:71}$ for each feature after performing temporal ATTN. The attention maps of each feature differ in time direction, indicating that each feature has different temporal characteristics. In the attention maps, the attention weights of the segments (keys) that contain behaviors leading to anomalies tends to be higher for each segment (query). This indicates that temporal ATTN can focus on the behaviors leading to the anomaly.

Figure 9 shows the attention maps between features when performing featurewise ATTN after temporal ATTN. The dependencies between features are cap-

tured, and this information is essential for capturing the MTS characteristics. When we compare each attention map of the segments, all segments exhibit similar dependencies between features. These segments are close in time and thus capture similar feature-wise dependencies.

6 Conclusions

In this paper, we proposed an anomaly prediction framework based on MTS forecasting and reconstruction using Transformer. Our model is trained to predict future trends using only normal data. Therefore, when given a time series exhibiting deviations from normal features, the results will not be accurate predictions of future trends, which can then be detected as precursors to anomaly occurrences through reconstruction. Detecting precursors of anomalies requires an accurate understanding of the temporal characteristics of the MTS. We modified the attention mechanism of each Transformer to perform ATTN in the time and feature directions. However, an NLP study reported that the effect of the ATTN was diminished after applying RES, despite using improved attention mechanisms. Therefore, we confirmed the same phenomenon in the context of MTS. We successfully preserved the strong effects of the improved attention mechanism by modifying the operation of the attention block. We enhanced the anomaly prediction performance by introducing these improvements.

Acknowledgements

This work was supported by the 4th Research Grant from the Hagiwara Foundation of Japan and by JSPS KAKENHI Grant Number JP24999361 (Grant-in-Aid for Early-Career Scientists).

References

- Abdulaal, A., Liu, Z., Lancewicki, T.: Practical approach to asynchronous multivariate time series anomaly detection and localization. Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining pp. 2485–2494 (2021)
- Abnar, S., Zuidema, W.: Quantifying attention flow in transformers. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics pp. 4190–4197 (Jul 2020)
- 3. Ahmad, S., Lavin, A., Purdy, S., Agha, Z.: Unsupervised real-time anomaly detection for streaming data. Neurocomputing **262**, 134–147 (2017)
- Audibert, J., Michiardi, P., Guyard, F., Marti, S., Zuluaga, M.A.: Usad: Unsupervised anomaly detection on multivariate time series. Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery & data mining pp. 3395–3404 (2020)
- 5. Ba, J.L., Kiros, J.R., Hinton, G.E.: Layer normalization. arXiv preprint arXiv:1607.06450 (2016)

15

- Chen, Z., Chen, D., Zhang, X., Yuan, Z., Cheng, X.: Learning graph structures with transformer for multivariate time-series anomaly detection in iot. IEEE Internet of Things Journal 9(12), 9179–9189 (2021)
- Du, D., Su, B., Wei, Z.: Preformer: predictive transformer with multi-scale segmentwise correlations for long-term time series forecasting. ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) pp. 1–5 (2023)
- Goldstein, M., Uchida, S.: A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. PloS one 11(4), e0152173 (2016)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. Proceedings of the IEEE conference on computer vision and pattern recognition pp. 770–778 (2016)
- Htut, P.M., Phang, J., Bordia, S., Bowman, S.R.: Do attention heads in bert track syntactic dependencies? arXiv preprint arXiv:1911.12246 (2019)
- Hundman, K., Constantinou, V., Laporte, C., Colwell, I., Soderstrom, T.: Detecting spacecraft anomalies using lstms and nonparametric dynamic thresholding. Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining pp. 387–395 (2018)
- Jhin, S.Y., Lee, J., Park, N.: Precursor-of-anomaly detection for irregular time series. Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining p. 917–929 (2023)
- Kobayashi, G., Kuribayashi, T., Yokoi, S., Inui, K.: Incorporating Residual and Normalization Layers into Analysis of Masked Language Models. Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing pp. 4547–4568 (Nov 2021)
- Kovaleva, O., Romanov, A., Rogers, A., Rumshisky, A.: Revealing the dark secrets of BERT. Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP) pp. 4365–4374 (Nov 2019)
- 15. Lai, K.H., Zha, D., Xu, J., Zhao, Y.: Revisiting time series outlier detection: Definitions and benchmarks. NeurIPS Datasets and Benchmarks (2021)
- Landauer, M., Onder, S., Skopik, F., Wurzenberger, M.: Deep learning for anomaly detection in log data: A survey. Machine Learning with Applications 12, 100470 (2023)
- Li, D., Chen, D., Jin, B., Shi, L., Goh, J., Ng, S.K.: Mad-gan: Multivariate anomaly detection for time series data with generative adversarial networks. International conference on artificial neural networks pp. 703–716 (2019)
- Lim, B., Zohren, S.: Time-series forecasting with deep learning: a survey. Philosophical Transactions of the Royal Society A 379(2194), 20200209 (2021)
- Lin, Y., Tan, Y.C., Frank, R.: Open sesame: Getting inside BERT's linguistic knowledge. Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP pp. 241–253 (Aug 2019)
- Liu, S., Yu, H., Liao, C., Li, J., Lin, W., Liu, A.X., Dustdar, S.: Pyraformer: Lowcomplexity pyramidal attention for long-range time series modeling and forecasting. International conference on learning representations (2021)
- Mareček, D., Rosa, R.: From balustrades to pierre vinken: Looking for syntax in transformer self-attentions. Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP (Aug 2019)
- Masini, R.P., Medeiros, M.C., Mendes, E.F.: Machine learning advances for time series forecasting. Journal of economic surveys 37(1), 76–111 (2023)

- 16 C. Maru et al.
- Mathur, A.P., Tippenhauer, N.O.: Swat: A water treatment testbed for research and training on ics security. 2016 international workshop on cyber-physical systems for smart water networks (CySWater) pp. 31–36 (2016)
- 24. Moritz, S., Rehbach, F., Chandrasekaran, S., Rebolledo, M., Bartz-Beielstein, T.: Gecco industrial challenge 2018 dataset: A water quality dataset for the internet of things: Online anomaly detection for drinking water quality competition at the genetic and evolutionary computation conference 2018, kyoto, japan. Kyoto, Japan (2018)
- Nie, Y., H. Nguyen, N., Sinthong, P., Kalagnanam, J.: A time series is worth 64 words: Long-term forecasting with transformers. International Conference on Learning Representations (2023)
- Pang, G., Shen, C., Cao, L., Hengel, A.V.D.: Deep learning for anomaly detection: A review. ACM computing surveys (CSUR) 54(2), 1–38 (2021)
- Reif, E., Yuan, A., Wattenberg, M., Viegas, F.B., Coenen, A., Pearce, A., Kim, B.: Visualizing and measuring the geometry of bert. Advances in Neural Information Processing Systems **32** (2019)
- Said Elsayed, M., Le-Khac, N.A., Dev, S., Jurcut, A.D.: Network anomaly detection using lstm based autoencoder. Proceedings of the 16th ACM Symposium on QoS and Security for Wireless and Mobile Networks pp. 37–45 (2020)
- Su, Y., Zhao, Y., Niu, C., Liu, R., Sun, W., Pei, D.: Robust anomaly detection for multivariate time series through stochastic recurrent neural network. Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining pp. 2828–2837 (2019)
- Torres, J.F., Hadjout, D., Sebaa, A., Martínez-Álvarez, F., Troncoso, A.: Deep learning for time series forecasting: a survey. Big Data 9(1), 3–21 (2021)
- Tuli, S., Casale, G., Jennings, N.R.: TranAD: Deep Transformer Networks for Anomaly Detection in Multivariate Time Series Data. Proceedings of VLDB 15(6), 1201–1214 (2022)
- 32. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)
- Wu, H., Xu, J., Wang, J., Long, M.: Autoformer: Decomposition transformers with auto-correlation for long-term series forecasting. Advances in Neural Information Processing Systems 34, 22419–22430 (2021)
- Xu, J., Wu, H., Wang, J., Long, M.: Anomaly transformer: Time series anomaly detection with association discrepancy. arXiv preprint arXiv:2110.02642 (2021)
- 35. Yang, Y., Zhang, C., Zhou, T., Wen, Q., Sun, L.: Dcdetector: Dual attention contrastive representation learning for time series anomaly detection. Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining p. 3033–3045 (2023)
- Zhang, C., Zhou, T., Wen, Q., Sun, L.: Tfad: A decomposition time series anomaly detection architecture with time-frequency analysis. Proceedings of the 31st ACM International Conference on Information & Knowledge Management pp. 2497–2507 (2022)
- Zhang, Y., Yan, J.: Crossformer: Transformer utilizing cross-dimension dependency for multivariate time series forecasting. International Conference on Learning Representations (2023)
- Zhang, Y., Chen, Y., Wang, J., Pan, Z.: Unsupervised deep anomaly detection for multi-sensor time-series signals. IEEE Transactions on Knowledge and Data Engineering (2021)

17

- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., Zhang, W.: Informer: Beyond efficient transformer for long sequence time-series forecasting. Proceedings of the AAAI conference on artificial intelligence 35(12), 11106–11115 (2021)
- Zhou, T., Ma, Z., Wen, Q., Wang, X., Sun, L., Jin, R.: Fedformer: Frequency enhanced decomposed transformer for long-term series forecasting. International Conference on Machine Learning pp. 27268–27286 (2022)
- Zong, B., Song, Q., Min, M.R., Cheng, W., Lumezanu, C., Cho, D., Chen, H.: Deep autoencoding gaussian mixture model for unsupervised anomaly detection. International conference on learning representations (2018)