

Breaking the Time-Frequency Granularity Discrepancy in Time-Series Anomaly Detection

Youngeun Nam
Korea Advanced Institute
of Science and Technology
youngeun.nam@kaist.ac.kr

Susik Yoon
Korea University
susik@korea.ac.kr

Yooju Shin
Korea Advanced Institute
of Science and Technology
yooju.shin@kaist.ac.kr

Minyoung Bae
Korea Advanced Institute
of Science and Technology
mybae@kaist.ac.kr

Hwanjun Song
Korea Advanced Institute
of Science and Technology
songhwanjun@kaist.ac.kr

Jae-Gil Lee*
Korea Advanced Institute
of Science and Technology
jagil@kaist.ac.kr

Byung Suk Lee
University of Vermont
bslee@uvm.edu

ABSTRACT

In light of the remarkable advancements made in time-series anomaly detection (TSAD), recent emphasis has been placed on exploiting the frequency domain as well as the time domain to address the difficulties in precisely detecting *pattern-wise* anomalies. However, in terms of anomaly scores, the *window granularity* of the frequency domain is inherently distinct from the *data-point granularity* of the time domain. Owing to this discrepancy, the anomaly information in the frequency domain has not been utilized to its full potential for TSAD. In this paper, we propose a TSAD framework, **Dual-TF**, that simultaneously uses both the time and frequency domains while breaking the *time-frequency granularity discrepancy*. To this end, our framework employs *nested-sliding windows*, with the outer and inner windows responsible for the time and frequency domains, respectively, and aligns the anomaly scores of the two domains. As a result of the high resolution of the aligned scores, the boundaries of pattern-wise anomalies can be identified more precisely. In six benchmark datasets, our framework outperforms state-of-the-art methods by 12.0–147%, as demonstrated by experimental results.

CCS CONCEPTS

• Computing methodologies → Spectral methods; Anomaly detection; • Information systems → Data stream mining.

KEYWORDS

Frequency/Spectral domain, Granularity discrepancy, Anomaly

ACM Reference Format:

Youngeun Nam, Susik Yoon, Yooju Shin, Minyoung Bae, Hwanjun Song, Jae-Gil Lee, and Byung Suk Lee. 2024. Breaking the Time-Frequency Granularity Discrepancy in Time-Series Anomaly Detection. In *Proceedings of the ACM Web Conference 2024 (WWW '24)*, May 13–17, 2024, Singapore, Singapore. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3589334.3645556>

*Corresponding author.



This work is licensed under a Creative Commons Attribution International 4.0 License.

WWW '24, May 13–17, 2024, Singapore, Singapore
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0171-9/24/05.
<https://doi.org/10.1145/3589334.3645556>

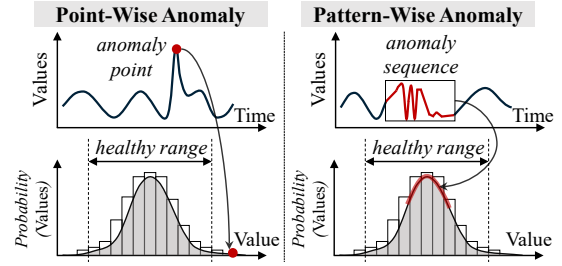


Figure 1: Difficulty of detecting a pattern-wise anomaly (right) compared with a point-wise anomaly (left).

1 Introduction

Time series, which is ubiquitous in various Web-based contexts such as Web servers and cloud services, is a fundamental resource for analyzing Web traffic patterns. *Time-series anomaly detection (TSAD)* is usually formulated as identifying the data points that significantly diverge from the normal or usual behavior. TSAD is commonly used to monitor states in many Web-related domains (e.g., cloud services [36]) as well as manufacturing, healthcare, finance, energy, and environment [12, 27, 30, 38, 39, 42].

Time-series anomalies are mainly categorized into *point-wise* and *pattern-wise* (or *collective*) anomalies [13, 32], each of which is specified for a particular point and sequence. According to the behavior-driven taxonomy [26], the pattern-wise anomalies are further divided into *shapelet*, *seasonality*, and *trend* anomalies. Detecting pattern-wise anomalies is considered more difficult than detecting point-wise anomalies. As shown in Figure 1, a point-wise anomaly has a very unusual value deviating from the normal range of the probability distribution, whereas a pattern-wise anomaly may still have usual values that fall in the normal range [34].

In order to enhance the capability of capturing pattern-wise anomalies, recent studies have notably started considering both the time and frequency domains [47, 51, 55]. The former represents the values as a function of time domain, while the latter represents the periods (or cycles) as a function of frequency domain. These recent approaches exploit the time domain mainly for finding point-wise anomalies and the frequency domain mainly for identifying pattern-wise anomalies. The *uncertainty principle* for time-series representation [19], which can be taken to mean that if a particular anomaly is well represented in one domain, the anomaly may not

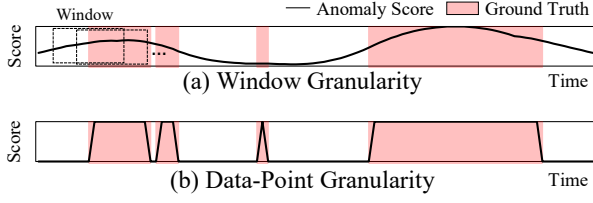


Figure 2: Comparison of ideal anomaly scores in (a) window granularity (existing) and (b) data-point granularity (ours).

be well represented in the other domain, provides strong support for this family of approaches.

Fully taking advantage of both domains for TSAD is, however, very challenging. Anomalies are specified at timestamps—i.e., in the time domain. Thus, the anomaly information found in the frequency domain needs to be aligned to the time domain for use in TSAD. Even worse, the *finest* granularity of the anomaly information in the frequency domain is coarser than that in the time domain. In detail, an anomaly score (*f-anomaly score*) in the frequency domain needs to be defined for a *window* (a sequence of data points), because a frequency spectrum can only be derived from a window (not a data point); in contrast, an anomaly score in the time domain (*t-anomaly score*) can be defined for a *data point*, e.g., reconstruction error [1, 54] and association discrepancy [46]. This problem is named the *time-frequency granularity discrepancy*.

To resolve this discrepancy, existing approaches (e.g., TFAD [50]) sacrifice the data-point granularity even for the time domain and stick to the window granularity for both domains. That is, in both domains, an anomaly score is assigned to a window, and all data points in the window share the same score. Even with sliding windows, sharing the same score with all data points in a window significantly degrades the resolution of detecting pattern-wise anomalies. Especially when a window contains both normal and anomalous data points in Figure 2(a), the window granularity scheme inevitably fails to detect the *exact boundaries* of pattern-wise anomalies, thereby resulting in low overall accuracy.

In this paper, we propose a novel TSAD method, *Dual-TF*, which exploits both the time and frequency domains *without* the time-frequency granularity discrepancy. Our key solution is to use the *nested sliding window (NS-window)* to accommodate both time and frequency information while aligning them in the data-point granularity. As shown in Figure 3, for calculating *t-anomaly* scores, the outer window slides as usual to capture various time contexts. For calculating *f-anomaly* scores, the inner window slides only within the corresponding outer window to produce multiple frequency spectrums for each data point; these diverse frequency spectrums are compared with one another to return *f-anomaly* scores. Finally, these multiple *t-anomaly* scores and *f-anomaly* scores are consolidated for each data point to satisfy the data-point granularity. As a result, the boundaries of pattern-wise anomalies are more clearly and precisely identified in Figure 2(b).

Meanwhile, deep neural networks (DNNs) have demonstrated their capability to recognize intricate correlations within complex data characterized by large volume and dimensionality over the last decade. This trend has extended to multivariate time-series anomaly

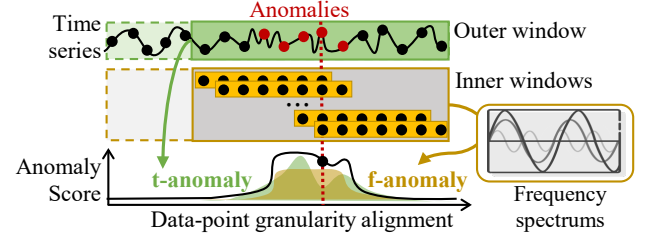


Figure 3: Our NS-windowing technique.

detection, resulting in an explosion of DNN-based methods suggesting methodological advances and improved performance [4]. Notably, attention-based models, such as Transformer, offer the benefit of considering sequence dependencies and outperform previous state-of-the-art methods by a significant margin in time-series analysis [44]. The capability to effectively identify and analyze complex patterns and correlations in data is the reason for using the Anomaly Transformer [46] as the backbone.

Dual-TF includes two Anomaly Transformers for calculating the *t-anomaly* and *f-anomaly* scores, respectively. These scores are calculated and combined using the proposed NS-windowing scheme. Through the extensive comparison with ten TSAD methods for six datasets, *Dual-TF* is shown to improve the TSAD accuracy by 12.0–147%. Furthermore, consistent with our expectation, *Dual-TF*'s higher ability to capture the boundaries of pattern-wise anomalies is visually confirmed, thus explaining the overall accuracy improvement. The idea of *using the NS-window for combining both domains* is very intuitive and widely applicable to any TSAD method based on the sliding window. We believe that the *simplicity* of our approach is a strong benefit because simple algorithms often make a big impact and gain widespread acceptance [40].

2 Related Work

Time-Series Anomaly Detection (TSAD) The majority of TSAD methods are designed for unsupervised learning owing to a lack of anomaly labels. Traditional TSAD methods can be classified into statistical [7, 10] and machine learning-based methods [15, 37]. In recent years, many studies have adopted deep learning, which is typically superior to traditional machine learning. Forecasting-based [14, 52] and reconstruction-based methods are two well-known approaches. The former uses prediction errors as anomaly scores, while the latter uses reconstruction errors. Previous research has shown that reconstruction-based methods generally outperform forecasting-based methods [17, 53].

BeatGAN [54] is a reconstructive approach based on a generative adversarial network (GAN), which uses time-series warping for data augmentation to improve accuracy. MSCRED [49] exploits an attention-based ConvLSTM to account for temporal dependency. Autoencoder models are similarly employed for reconstruction in OmniAnomaly [41]. RANSynCoders [1] improves autoencoder training efficiency via feature synchronization, bootstrapping, and quantile loss. Notably, a new reconstructive approach that combines series and prior association to make anomalies distinctive is proposed as *Anomaly Transformer* [46].

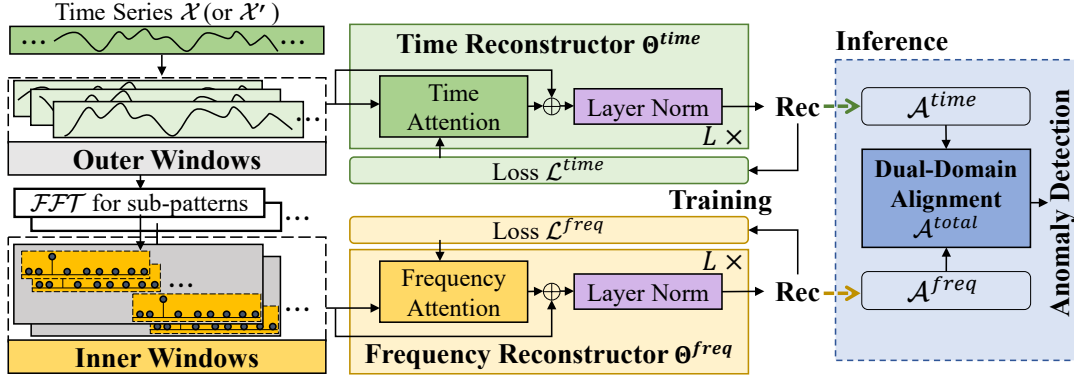


Figure 4: Overview of the *Dual-TF* framework for the training and inference phases.

Frequency Domain Analysis for Time Series Recent methods use both the time and frequency (i.e., spectral) domains [48]. For forecasting, DEPTS [16] models complex periodicities with a learnable cosine function, while FEDFormer [55] targets long-term forecasting by extracting important frequency components using the frequency-enhanced block and frequency-enhanced attention. For unsupervised representation learning, BTSF [47] conducts iterative bilinear temporal-spectral fusion, where information on each domain is conveyed to a bilinear feature to model time-frequency dependencies. TF-C [51] augments the time and frequency domains independently to produce positive samples, followed by regularization to ensure coherence in the representation of both domains. However, these studies primarily focus on capturing *general* time-series features and thus are unsuitable for detecting time-series anomalies. To identify pattern-wise features that can be defined only in the frequency domain, it is necessary to acquire frequency information that is appropriate for anomaly detection, as opposed to general frequency information obtained in the existing studies. **Frequency Domain Analysis for TSAD** Frequency-based models for TSAD have received much attention in recent years. The spectral residual (SR) introduced in SR-CNN [36] uses a frequency-based technique to generate a saliency map for TSAD. PFT [35] is a partial Fourier transform that achieves a speedup of an order of magnitude without sacrificing accuracy. TFAD [50] utilizes frequency domain analysis for TSAD with augmentation and decomposition. To the best of our knowledge, TFAD is the closest to our work since it uses the time and frequency domains together. Nevertheless, the existing methods (including TFAD) do *not* offer the data-point granularity for the frequency domain, still facing challenges in precisely detecting pattern-wise anomalies.

3 TSAD Framework: *Dual-TF*

Problem Formulation: Let's consider a multivariate time series, $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ($\mathcal{X} \in \mathbb{R}^{n \times d}$), where n is the number of data points and d is the number of features. Using \mathcal{X} as a training set, we aim at building an anomaly detector $\mathcal{A}(\mathbf{x}_t, \Theta^{time}, \Theta^{frequency})$ that returns an anomaly score for a given data point using deep neural networks (e.g., autoencoders and Transformers) parameterized by Θ^{time} and $\Theta^{frequency}$ respectively for the time and frequency domains. Then, given another multivariate time series as a test set,

$\mathcal{X}' \in \mathbb{R}^{n' \times d}$, the anomaly detector $\mathcal{A}()$ is used to classify each data point \mathbf{x}'_t in \mathcal{X}' as being either normal or anomalous. Our problem is categorized as *unsupervised* anomaly detection because no label information is used.

Window Specification: For the NS-windowing scheme in Figure 3, an outer window and its inner windows are continuously extracted from a multivariate time-series \mathcal{X} (or \mathcal{X}'). An *outer window* at timestamp t is defined as $OW_t = \{\mathbf{x}_t, \mathbf{x}_{t+1}, \dots, \mathbf{x}_{t+w^{outer}-1}\}$ of length w^{outer} . Then, the original time series \mathcal{X} (or \mathcal{X}') is reorganized as a sequence of overlapping outer windows, $OW = \{OW_1, OW_2, \dots, OW_{n-w^{outer}+1}\}^1$. For a given outer window OW_t , its *inner window* of length w^{inner} at timestamp $i \in [t, t+w^{outer}-w^{inner}]$ is defined $IW_i = \{\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_{i+w^{inner}-1}\}$. Then, again, an outer window is reorganized as a sequence of overlapping inner windows $IW_t = \{IW_t, IW_{t+1}, \dots, IW_{t+w^{outer}-w^{inner}}\}^1$.

Time-Frequency Granularity Discrepancy: The number of *degrees of freedom* represents the number of values that can vary freely [9]. If the anomaly scores satisfy the data-point granularity, the degree of freedom should be the length of a time series n ; on the other hand, if the anomaly scores follow the window granularity, the degree of freedom should be much smaller than n because the scores cannot vary within the same window. Therefore, the *time-frequency granularity discrepancy* is formally defined as when $d.f.(\{t\text{-anomaly scores}\}) \neq d.f.(\{f\text{-anomaly scores}\})$, where $d.f.(\cdot)$ denotes the degree of freedom.

3.1 Overall Framework of *Dual-TF*

Figure 4 shows the overall procedure of *Dual-TF*. By the scheme of NS-windowing, a multivariate time series \mathcal{X} (or \mathcal{X}') is transformed into a sequence of outer windows and a sequence of inner windows for each outer window. *Dual-TF* employs two neural network reconstructors: the *time reconstructor* and the *frequency reconstructor*, each for time and frequency domain. The outer windows are fed to the time reconstructor. On the other hand, a fast Fourier transform (FFT) [31] converts each inner window from its original time domain to a representation in the frequency domain; then, the converted inner windows are fed to the frequency reconstructor.

¹The slide step is set to be 1 for finding the boundaries of pattern-wise anomalies more precisely as well as for simplifying the expression.

These reconstructors are *individually* trained to minimize the reconstruction losses, following the conventional procedure of reconstruction-based TSAD [43, 46]. The batches for the two reconstructors should be constructed separately due to the different dimensionalities of the inputs to the two reconstructors. During the inference phase, an anomaly score for each data point is calculated using the reconstruction losses from the two reconstructors. To break the time-frequency granularity discrepancy and thus achieve the data-point granularity in both domains, we align the time-domain reconstruction losses from an outer window with the frequency-domain reconstruction losses from its inner windows.

3.2 Time-Frequency Dual-Domain Training

Optimal Window Length: For the NS-windowing scheme, one of the most important issues is to determine the proper window lengths, w^{outer} and w^{inner} . Because the primary goal of using inner windows is to capture pattern-wise anomalies, we set them to embody major periods based on spectral analysis. In detail, each dimension of the entire time series X is converted to the frequency domain by the fast Fourier transform (FFT). The output of the FFT is Hermitian-symmetric; that is, the positive-frequency terms are the complex conjugates of the corresponding negative-frequency terms. The negative frequency terms can be ignored because of the redundancy. Here, a frequency *magnitude* is defined as the absolute value (or modulus) of a complex number. Then, the most *dominant* frequency whose magnitude is the largest in each dimension is identified; the smallest dominant frequency in all dimensions is chosen as the *representative* frequency v^{major} in X . Last, we set $w^{inner} = \lceil \frac{1}{v^{major}} \rceil$, because the period is the reciprocal of the frequency, and $w^{outer} = \rho \cdot w^{inner}$, where $\rho (> 1)$ is a hyperparameter. See Appendix A for the details of the window-length selection.

We briefly discuss the optimality of the inner windows. The time-frequency uncertainty principle [18] states that the exact time and frequency of a signal can never be known simultaneously [21]. In determining the window size, this principle indicates the trade-off between time and frequency uncertainty within a window. In the time domain, the smaller w^{inner} , the lower the uncertainty; in the frequency domain, the smaller w^{inner} , the greater the uncertainty. Let $\mathcal{U}^{time}(w^{inner})$ and $\mathcal{U}^{freq}(w^{inner})$, respectively, denote the uncertainty in the time and frequency domains, given a window of length w^{inner} . (See Appendix B the formal definition of the uncertainty.) Then, the uncertainty in the two domains is

$$\mathcal{U}(w^{inner}) = \mathcal{U}^{time}(w^{inner}) + \mathcal{U}^{freq}(w^{inner}). \quad (1)$$

Theorem 3.1 formally states the optimal condition for the inner window length w^{inner} .

THEOREM 3.1 (OPTIMAL WINDOW LENGTH). *When $w^{inner} = \lceil \frac{1}{v^{major}} \rceil$, the uncertainty within the window, $\mathcal{U}(w^{inner})$ in Eq. (1), is minimized.*

PROOF. See Appendix B for the proof. \square

Reconstructor Network: While any neural network reconstructor is applicable to *Dual-TF*, we choose Anomaly Transformer [46] because it has shown the state-of-the-art performance. To make this paper be self-contained, we briefly describe the key mechanism of the Anomaly Transformer. Each input is $X^0 = OW_t$ for the time

reconstructor and $X^0 = I\mathcal{W}_t$ for the frequency reconstructor. Following the self-attention mechanism, the output of the l -th layer ($l \in [1, n^{layer}]$) is defined by

$$Q, K, V = X^{l-1}W_Q^l, X^{l-1}W_K^l, X^{l-1}W_V^l$$

$$\text{Attention}(X^l) = \text{Softmax}\left(\frac{QK^\top}{\sqrt{d^{model}}}\right)V \quad (2)$$

$$\Theta^*(X^l) = \widehat{X}^l = \text{LayerNorm}(X^{l-1} + \text{Attention}(X^{l-1})),$$

where Q, K , and V are queries, keys, and values; W_Q^l, W_K^l , and W_V^l are the learnable parameters for them; d^{model} is the number of hidden channels; and \widehat{X}^l is the output of the reconstruction network. More importantly, the association discrepancy is defined as the KL-divergence between the prior association (P^l) and the series association (S^l),

$$\text{AssDis}(P, S; X) = \frac{1}{n^{layer}} \sum_{l=1}^{n^{layer}} (\text{KL}(P^l \| S^l) + \text{KL}(S^l \| P^l)). \quad (3)$$

where P^l is generated by the Gaussian kernel to represent the adjacent context and S^l is defined as the usual self-attention in Eq. (2) to represent the overall context.

The n^{layer} layers of the Anomaly Transformer backbone in Eq. (2) are used for both reconstructors. For the time reconstructor, $Q, K, V \in \mathbb{R}^{w^{outer} \times d^{model}}$; $W_Q^l, W_K^l, W_V^l \in \mathbb{R}^{d^{model} \times d^{model}}$; $P^l, S^l \in \mathbb{R}^{w^{outer} \times w^{outer}}$. In the same manner, for the frequency reconstructor, $Q, K, V \in \mathbb{R}^{(m \times w^{inner}) \times d^{model}}$; $W_Q^l, W_K^l, W_V^l \in \mathbb{R}^{d^{model} \times d^{model}}$; $P^l, S^l \in \mathbb{R}^{(m \times w^{inner}) \times (m \times w^{inner})}$, where m denotes the number of inner windows per outer window, i.e., $m = w^{outer} - w^{inner} + 1$.

Time Reconstructor Loss: Time reconstructor forms each outer window $OW_t = \{x_t, x_{t+1}, \dots, x_{t+w^{outer}-1}\}$ ($OW_t \in \mathbb{R}^{w^{outer} \times d}$), where $t \in [1, n - w^{outer} + 1]$, to $\widehat{OW}_t = \{\hat{x}_t, \hat{x}_{t+1}, \dots, \hat{x}_{t+w^{outer}-1}\}$ ($\widehat{OW}_t \in \mathbb{R}^{w^{outer} \times d}$) using the parameter Θ^{time} . Then, the reconstruction loss of OW_t is formulated by

$$\text{RecLoss}^{time}(OW_t, \widehat{OW}_t) = \sum_{i=t}^{t+w^{outer}-1} \|x_i - \hat{x}_i\|_2^2. \quad (4)$$

It is evident that the time domain satisfies the data-point granularity in Eq. (4). The association discrepancy is (optionally) added to the final loss,

$$\mathcal{L}^{time}(OW_t, \widehat{OW}_t) = \text{RecLoss}^{time}(OW_t, \widehat{OW}_t) - \lambda \cdot \text{AssDis}^{time}(P, S; OW_t), \quad (5)$$

where $\lambda (> 0)$ is the hyperparameter for weighting the association discrepancy. These losses for the outer windows in a batch are summed up to update the parameter Θ^{time} via backpropagation.

Frequency Reconstructor Loss: A sequence of overlapping inner windows $I\mathcal{W}_t = \{IW_t, IW_{t+1}, \dots, IW_{t+w^{outer}-w^{inner}}\}$ is derived, given an outer window OW_t . First, each inner window IW_t ($\in \mathbb{R}^{w^{inner} \times d}$) is converted to a frequency spectrum $\mathcal{F}\mathcal{F}\mathcal{T}(IW_t)$ ($\in \mathbb{R}^{w^{inner} \times d}$), where $\mathcal{F}\mathcal{F}\mathcal{T}()$ returns a sequence of the magnitudes in the result of the FFT. That is, $\mathcal{F}\mathcal{F}\mathcal{T}(IW_t)$ is regarded as the *counterpart* of the data point x_t in the frequency domain. Here, $\mathcal{F}\mathcal{F}\mathcal{T}()$ is separately applied to each of d dimensions. Then, $\{\mathcal{F}\mathcal{F}\mathcal{T}(IW_t), \mathcal{F}\mathcal{F}\mathcal{T}(IW_{t+1}), \dots\}$ is reconstructed by the frequency

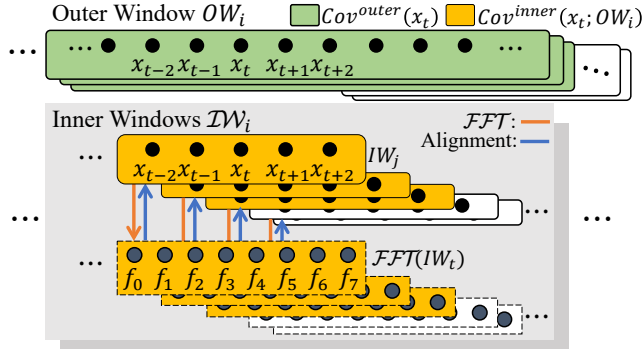


Figure 5: Data-point granularity alignment.

reconstructor from $\{\mathcal{FFT}(IW_t), \mathcal{FFT}(IW_{t+1}), \dots\}$. Last, the reconstruction loss of IW_t generated from OW_t is defined as

$$\begin{aligned} & \text{RecLoss}^{freq}(IW_t, \widehat{IW}_t) \\ &= \sum_{i=t}^{t+w^{outer}-1} \sum_{j=1}^{w^{inner}d} \|\mathcal{FFT}(IW_i)_{[j]} - \widehat{\mathcal{FFT}(IW_i)}_{[j]}\|_2^2. \end{aligned} \quad (6)$$

Here, $\mathcal{FFT}()_{[j]}$ is the frequency spectrum of the j -th dimension. By virtue of the sequence of inner windows, the frequency domain also satisfies the data-point granularity in Eq. (6) because an inner window is created for every data point except the last w^{inner} ones within the outer window; we note that these uncovered data points will be covered soon by the succeeding outer windows. Similar to Eq. (5), for the update of the parameter Θ^{freq} , the total loss is defined as

$$\begin{aligned} \mathcal{L}^{freq}(IW_t, \widehat{IW}_t) &= \text{RecLoss}^{freq}(IW_t, \widehat{IW}_t) \\ &\quad - \lambda \cdot \text{AssDis}^{freq}(P, S; IW_t). \end{aligned} \quad (7)$$

3.3 Data-Point Granularity Alignment for Anomaly Scoring

Using the two trained reconstructors Θ^{time} and Θ^{freq} , we derive the anomaly score for each data point x' in a test set \mathcal{X}' . First, an anomaly score in the time domain (*t-anomaly score*) and an anomaly score (*f-anomaly score*) in the frequency domain are derived separately from the two reconstructors. Then, the *t-anomaly score* and the *f-anomaly score* are combined to form the final anomaly score, which will be compared against a threshold.

Because an entire time series (or an outer window) is converted to a sequence of *overlapping* outer (or inner) windows, each data point is covered by *multiple* outer (or inner) windows. Definitions 3.2 and 3.3 formalize the sets of *covering* outer and inner windows for a given data point.

Definition 3.2 (COVERING OUTER WINDOW). Given a data point x_t , the set of *covering outer windows* is defined by $\text{Cov}^{outer}(x_t) = \{OW_i \mid x_t \in OW_i, 1 \leq i \leq n - w^{outer} + 1\}$.

Definition 3.3 (COVERING INNER WINDOW). Given a data point x_t and its covering outer window $OW_i \in \text{Cov}^{outer}(x_t)$, the set of *covering inner windows* is defined by $\text{Cov}^{inner}(x_t; OW_i) = \{IW_j \mid x_t \in IW_j \wedge IW_j \in OW_i, i \leq j \leq i + w^{inner} - 1\}$.

Figure 5 illustrates how *t-anomaly scores* and *f-anomaly scores* are calculated and then aligned. First, a *t-anomaly score* is easily calculated for each data point within a given outer window using the reconstruction loss. Though any reconstruction-based criterion is usable, we follow the one proposed for the state-of-the-art model, Anomaly Transformer [46],

$$\begin{aligned} \mathcal{A}^{time}(x_t; OW_i) &= [\text{Softmax}(-\text{AssDis}^{time}(P, S; OW_i)) \\ &\quad \odot \text{RecLoss}^{time}(OW_i, \widehat{OW}_i)]_{x_t}, \end{aligned} \quad (8)$$

where $\text{AssDis}()$ and $\text{RecLoss}()$ are the association discrepancy and the reconstruction loss used in Eq. (5), and $[\cdot]_{x_t}$ returns the value for x_t . Second, an *f-anomaly score* is calculated in two steps. (i) A score $\mathcal{A}^{freq}(IW_j; IW_i)$ is derived for each covering inner window within the context of all inner windows from the outer window; that is, it basically represents the spectral difference of IW_j to the other inner windows in IW_i . (ii) The exponential function is applied to each $\mathcal{A}^{freq}(IW_j; IW_i)$, and these results are averaged for the outer window, where $\text{Cov}^{inner} = \text{Cov}^{inner}(x_t; OW_i)$,

$$\begin{aligned} \mathcal{A}^{freq}(x_t; OW_i) &= \frac{1}{|\text{Cov}^{inner}|} \sum_{IW_j \in \text{Cov}^{inner}} \exp(\mathcal{A}^{freq}(IW_j; IW_i)), \\ \mathcal{A}^{freq}(IW_j; IW_i) &= [\text{Softmax}(-\text{AssDis}^{freq}(P, S; IW_i)) \\ &\quad \odot \text{RecLoss}^{freq}(IW_i, \widehat{IW}_i)]_{IW_j}. \end{aligned} \quad (9)$$

We note that an *f-anomaly score* is derived for each data point using its covering inner windows, where different windows cover different intervals in an outer window. By supporting the data-point granularity for *f-anomaly scores*, we want to distinguish between when x_t is at the center of or near the boundary of a pattern-wise anomaly. As x_t is located at a more central location of this anomaly, its set of covering inner windows overlap the anomaly more significantly, where each $\mathcal{A}^{freq}(IW_j; IW_i)$ in Eq. (9) becomes higher. Thus, to make the *f-anomaly score* for the central data point *stand out*, we take the exponential of each $\mathcal{A}^{freq}(IW_j; IW_i)$. This simple treatment for TSAD is proven to be very effective, as shown in the ablation study.

Next, Eqs. (8) and (9) are averaged for all of x_t 's covering outer windows, where $\text{Cov}^{outer} = \text{Cov}^{outer}(x_t)$

$$\mathcal{A}^{time \text{ or } freq}(x_t) = \frac{1}{|\text{Cov}^{outer}|} \sum_{OW_i \in \text{Cov}^{outer}} \mathcal{A}^{time \text{ or } freq}(x_t; OW_i), \quad (10)$$

and this score for each x_t is min-max normalized with respect to the scores of all data points in \mathcal{X} .

At last, considering *both* the *t-anomaly* and *f-anomaly scores*, the final anomaly score is defined by

$$\mathcal{A}^{total}(x_t) = \mathcal{A}^{time}(x_t) + \mathcal{A}^{freq}(x_t). \quad (11)$$

REMARK 3.4. The *t-anomaly* and *f-anomaly scores* in Eq. (10) *break the time-frequency granularity discrepancy*.

PROOF. Given any pair of x_i and x_j where $i \neq j$, $\text{Cov}^{outer}(x_i) \neq \text{Cov}^{outer}(x_j)$ and $\text{Cov}^{inner}(x_i; \cdot) \neq \text{Cov}^{inner}(x_j; \cdot)$ by the definition of the NS-windows. Therefore, $\text{d.f.}(\{\mathcal{A}^{time}(x_t) \mid 1 \leq t \leq n\}) = \text{d.f.}(\{\mathcal{A}^{freq}(x_t) \mid 1 \leq t \leq n\}) = n$. \square

4 Evaluation

For reproducibility, the source code of our framework is available at <https://github.com/kaist-dmlab/DualTF>.

Table 1: Benchmark dataset statistics.

Datasets	Applications	w^{inner}	# Train	# Test	Entity×Dim.	# Point Anomaly (Ratio)	# Pattern Anomaly (Ratio)	Avg. Length of Pattern Anomaly
TODS (Point)	Synthetic	25	20,000	5,000	2×1	250 (100%)	0 (0%)	N/A
TODS (Pattern)	Synthetic	25	20,000	5,000	3×1	0 (0%)	250 (100%)	10
ASD	Server Monitoring	288	8,527	4,320	12×19	0 (0%)	199 (100%)	31
ECG	Medical Checkup	143	6,995	2,851	9×2	0 (0%)	208 (100%)	208
PSM	Server Monitoring	360	132,481	87,841	1×25	16 (0.07%)	24,365 (99.93%)	435
CompanyA	Server Monitoring	144	21,600	13,302	3×8	10 (8.53%)	104 (91.47%)	4

Table 2: Performance comparison between TSAD methods in terms of the best point-wise F_1 score with the highest scores highlighted in bold.

Methods	TODS (Point)		TODS (Pattern)			ASD	ECG	PSM	CompanyA	Avg. ↑	Rank ↓
	Gloabl	Contextual	Shapelet	Seasonal	Trend						
ISF	0.943	0.164	0.103	0.093	0.209	0.295	0.256	0.478	0.134	0.297	9
	(±0.017)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.002)	
LOF	0.933	0.093	0.096	0.092	0.093	0.376	0.327	0.524	0.059	0.288	11
	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.067)	(±0.000)	(±0.099)	(±0.019)	(±0.010)	
OCSVM	0.937	0.170	0.104	0.093	0.094	0.266	0.264	0.469	0.266	0.296	10
	(±0.019)	(±0.000)	(±0.000)	(±0.000)	(±0.000)	(±0.071)	(±0.000)	(±0.000)	(±0.019)	(±0.011)	
VAE	0.915	0.584	0.503	0.847	0.181	0.327	0.274	0.443	0.261	0.482	4
	(±0.031)	(±0.034)	(±0.094)	(±0.017)	(±0.002)	(±0.022)	(±0.003)	(±0.000)	(±0.028)	(±0.014)	
MS-RNN	0.839	0.553	0.248	0.799	0.180	0.379	0.276	0.443	0.228	0.438	5
	(±0.000)	(±0.000)	(±0.144)	(±0.022)	(±0.000)	(±0.003)	(±0.002)	(±0.000)	(±0.003)	(±0.014)	
OmniAnomaly	0.543	0.542	0.149	0.203	0.185	0.197	0.216	0.467	0.182	0.298	8
	(±0.001)	(±0.008)	(±0.004)	(±0.017)	(±0.013)	(±0.096)	(±0.037)	(±0.098)	(±0.046)	(±0.009)	
RANSynCoders	0.674	0.482	0.166	0.163	0.175	0.383	0.208	0.571	0.112	0.326	7
	(±0.127)	(±0.000)	(±0.002)	(±0.007)	(±0.011)	(±0.234)	(±0.003)	(±0.017)	(±0.027)	(±0.026)	
TranAD	0.569	0.553	0.165	0.179	0.169	0.294	0.461	0.443	0.225	0.340	6
	(±0.000)	(±0.000)	(±0.000)	(±0.030)	(±0.000)	(±0.007)	(±0.028)	(±0.000)	(±0.008)	(±0.000)	
TFAD	0.878	0.871	0.558	0.854	0.363	0.432	0.356	0.537	0.276	0.569	3
	(±0.000)	(±0.009)	(±0.150)	(±0.018)	(±0.001)	(±0.003)	(±0.002)	(±0.080)	(±0.071)	(±0.035)	
Anomaly Transformer	0.943	0.942	0.730	0.867	0.460	0.425	0.464	0.578	0.317	0.636	2
	(±0.000)	(±0.000)	(±0.000)	(±0.028)	(±0.005)	(±0.017)	(±0.001)	(±0.001)	(±0.099)	(±0.011)	
Dual-TF	0.968	0.943	0.741	0.925	0.476	0.661	0.538	0.723	0.436	0.712	1
	(±0.017)	(±0.001)	(±0.005)	(±0.041)	(±0.017)	(±0.019)	(±0.076)	(±0.047)	(±0.021)	(±0.011)	

4.1 Experiment Settings

Datasets: Table 1 summarizes the benchmark datasets used, categorizing anomalies as a pattern-wise anomaly if the length of its anomaly interval is greater than 1 and a point-wise anomaly otherwise [5]. The number in the parenthesis along the number of point-wise or pattern-wise anomalies is the ratio of the number of anomaly-labeled data points of each category to the total number of anomaly-labeled data points. See Appendix C for the generation of the TODS dataset [26]. The ASD [28], ECG [23], and PSM [2] are public datasets commonly used for evaluating TSAD. The only proprietary dataset is the CompanyA (anonymized), which is derived from the operation of cloud servers and represents service traffic.

Baselines: We conduct a comparative analysis of *Dual-TF* against both traditional and recent works. For traditional methods, ISF [29], LOF [8], OCSVM [37], and Variational Autoencoder (VAE) [20] are included. For the state-of-the-art methods, OmniAnomaly [41], Modified-RNN (MS-RNN) [24], RANSynCoders [1], TranAD [43], TFAD [50], and Anomaly Transformer [46] are included.

See Appendix C for the descriptions of the baselines.

Evaluation Metrics: To evaluate the detection accuracy at the data-point level, we primarily employ the *point-wise* F_1 score [3]. In this case, a predicted anomaly is valid if it falls within a small margin (i.e., ten timestamps) of the actual location of the anomaly. In contrast, it is well-known that the widely-used point-adjusted (PA) metric has overestimation issues [45, 46], as a window is considered correct if both a predicted anomaly and the true anomaly occur just within the same window. In addition, we report the *best* F_1 score using the anomaly threshold that produces the highest F_1 score across all methods in order to eliminate the effect of threshold selection. We repeat every test *three* times with random seeds and report the average as well as the standard deviation.

Furthermore, addressing recent concerns on evaluation metrics [25, 33], we adopt new evaluation metrics designed for TSAD. The Range Area Under the Curve (R_AUC) and the Volume Under the Surface (VUS) [33] are employed, where the receiver operating characteristic (ROC) curve and precision-recall (PR) curves are considered. The VUS extends the mathematical model of the R_AUC by allowing the buffer length to be varied. Thus, R_AUC_ROC and R_AUC_PR are defined for the former, and VUC_ROC and VUC_PR

Table 3: Performance comparison for *Dual-TF* in terms of VUS [33] and other new evaluation metrics with the highest scores highlighted in bold. See [33] for the details of these evaluation metrics.

Evaluation Metrics	Methods	TODS (Point)		TODS (Pattern)			ASD	ECG	PSM	CompanyA
		Gloabl	Contextual	Shapelet	Seasonal	Trend				
R_AUC_ROC	Anomaly Transformer <i>Dual-TF</i>	0.9995 0.9998	0.9859 0.9995	0.8457 0.9097	0.9272 0.9611	0.6277 0.7035	0.8498 0.9013	0.6432 0.7216	0.6158 0.7735	0.8493 0.8653
R_AUC_PR	Anomaly Transformer <i>Dual-TF</i>	0.9994 0.9998	0.9862 0.9996	0.6878 0.7925	0.7736 0.8719	0.3713 0.4287	0.5263 0.6058	0.2447 0.3809	0.4789 0.6304	0.4139 0.4254
VUS_ROC	Anomaly Transformer <i>Dual-TF</i>	0.9354 0.9373	0.9160 0.9322	0.8065 0.8843	0.9147 0.9380	0.6222 0.6992	0.7952 0.8505	0.6343 0.7067	0.6073 0.7752	0.8335 0.8568
VUS_PR	Anomaly Transformer <i>Dual-TF</i>	0.8985 0.9053	0.8765 0.9014	0.6000 0.6950	0.6920 0.7620	0.3560 0.4016	0.4466 0.5127	0.2424 0.3754	0.4665 0.6131	0.3590 0.3694

Table 4: Ablation study for *Dual-TF* in terms of the best point-wise F_1 score with the highest scores highlighted in bold.

Variations	TODS (Point)		TODS (Pattern)			ASD	ECG	PSM	CompanyA	Average
	Gloabl	Contextual	Shapelet	Seasonal	Trend					
(i) w/o Time Reconstructor	0.439	0.661	0.715	0.858	0.476	0.629	0.267	0.500	0.247	0.532
(ii) w/o Frequency Reconstructor	0.943	0.942	0.728	0.824	0.457	0.415	0.460	0.578	0.264	0.623
(iii) w/o NS-Windowing	0.388	0.600	0.691	0.798	0.213	0.313	0.263	0.397	0.229	0.433
(iv) w/o Exponential Average	0.953	0.905	0.620	0.780	0.465	0.577	0.507	0.677	0.322	0.645
<i>Dual-TF</i>	0.968	0.943	0.741	0.925	0.476	0.661	0.538	0.723	0.436	0.712

are defined for the latter. These metrics can accurately evaluate the detection of pattern-wise (range) anomalies. Moreover, the VUS reduces the influence of an anomaly threshold.

Implementations and Model Configurations: *Dual-TF* is implemented using PyTorch 1.13.1. The only hyperparameter ρ for *Dual-TF*, which specifies the size of an outer window, is set to 2 for all datasets. The Adam optimizer is used with an initial learning rate of 10^{-4} . The batch size is 4 considering the large size of each training instance and the memory budget of a GPU. The training process is early stopped within 10 epochs. Following the author implementation of Anomaly Transformer [46], the weight λ in Eqs. (4) and (6) is 3, the number of layers n^{layer} is 3, the number of hidden channels d^{model} is 512, and the number of attention heads is 8. For all baseline methods, the authors' source code is employed without any modification, and the hyperparameters are set to be the default values in the author implementation. All experiments are conducted on a server equipped with an NVIDIA RTX 3090Ti.

See Appendix C for details.

4.2 Overall Performance Comparison

Table 2 shows the best point-wise F_1 score of *Dual-TF* as well as all baselines for six benchmark datasets. *Dual-TF* is shown to significantly outperform the state-of-the-art TSAD methods; specifically, it yields a detection accuracy of 12.0–147% higher on average than the other methods. Anomaly Transformer and TFAD are ranked second and third, respectively. *Dual-TF* effectively handles both point-wise (in TODS (Point)) and pattern-wise (in TODS (Pattern), ASD, ECG, PSM, and CompanyA) anomalies while showing greater improvement in detecting the pattern-wise anomalies. Compared

with Anomaly Transformer, the improvement in the F_1 score is 0.10–2.63% for point-wise anomalies, and it is *increased* to 1.55–55.4% for pattern-wise anomalies owing to the incorporation of the frequency domain. Additionally shown in Table 3, the improvement in the range-AUC measures (R_AUC_ROC or R_AUC_PR) is up to 1.38% for point-wise anomalies, and up to 55.7% for pattern-wise anomalies. Moreover, the enhancement in the VUS-based measures (VUS_ROC or VUS_PR) is up to 2.84% for point-wise anomalies, and up to 54.8% for pattern-wise anomalies with the help of the frequency domain. Moreover, compared with TFAD, which uses both the time and frequency domains, the F_1 score improves by 8.25–57.9% because of the higher resolution achieved by the data-point granularity. Overall, these results indeed demonstrate the value of combining both domains while breaking the time-frequency granularity discrepancy.

4.3 Ablation Study

The contribution of each main component of *Dual-TF* to the anomaly detection accuracy is investigated through an ablation study in Table 4. Specifically, when (i) time reconstructor is removed, (ii) frequency reconstructor is removed in Figure 4, (iii) the window granularity is enforced for the frequency domain without the NS-windowing scheme, or (iv) the exponential average in Eq. (9) for aggregating f-anomaly scores is replaced with the arithmetic average, the F_1 score for each variation is measured. All of these main components are shown to be important, with the *NS-windowing* component having the most outstanding effect. Interestingly, the use of f-anomaly scores at the window granularity in the third variation may actually contaminate t-anomaly scores by unnecessarily increasing the final anomaly scores for a normality interval and

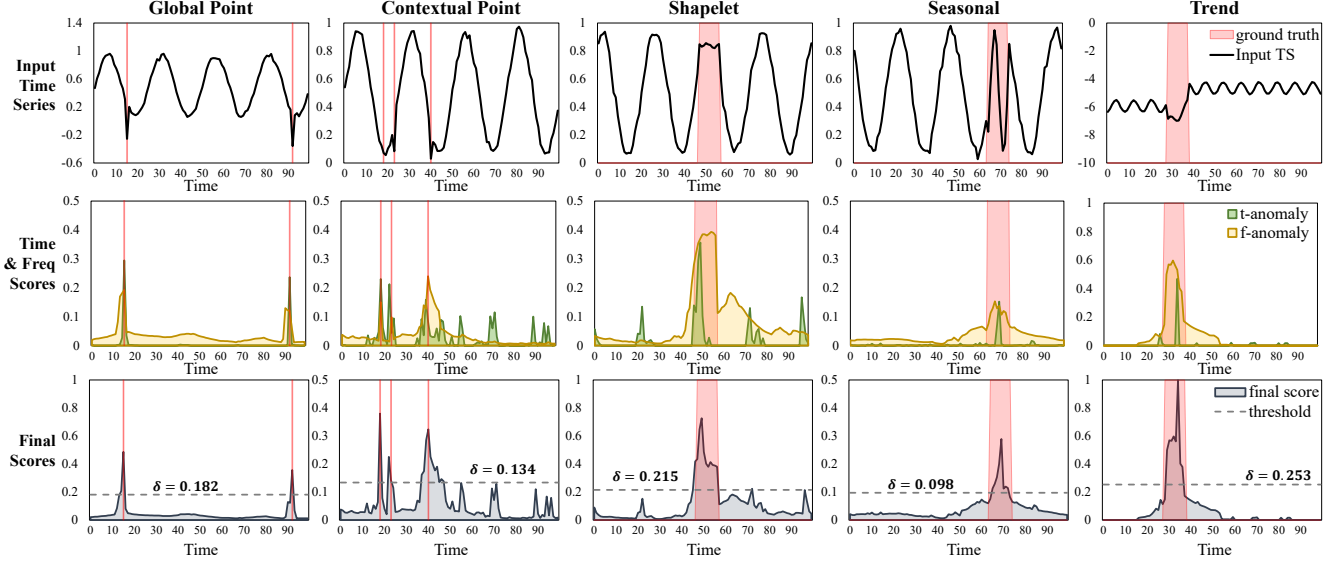


Figure 6: Visualization of dual-domain anomaly scores from *Dual-TF* for different categories of point- and pattern-wise anomalies using the TODS dataset.

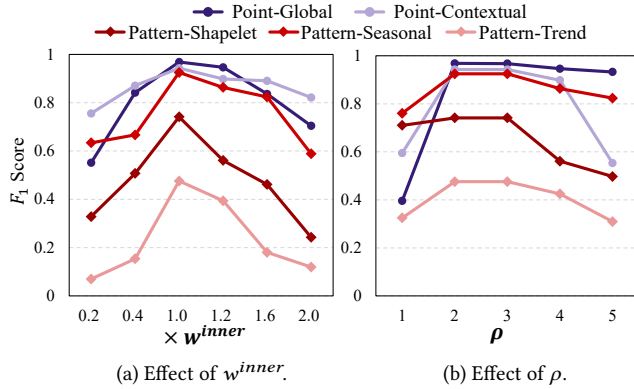


Figure 7: Effect of the inner and outer window lengths in the TODS dataset.

thus producing a large number of false positives. Overall, this comprehensive ablation study reaffirms the significance of breaking the time-frequency granularity discrepancy.

4.4 Qualitative Analysis through Visualization

Figure 6 visualizes the results of *Dual-TF* for the five anomaly categories [26] in terms of t-anomaly and f-anomaly scores (second row) and then final anomaly scores (third row) in the TODS dataset. For point-wise anomalies (first and second columns), t-anomaly scores increase sharply at the actual locations. For pattern-wise anomalies (third, fourth, and fifth columns), f-anomaly scores maintain high values *throughout the entire anomaly interval*, whereas t-anomaly scores jump only at a few timestamps. As a result, each domain plays a distinct role and both point-wise and pattern-wise anomalies are precisely detected by the final scores.

4.5 Window Length Sensitivity

Because the window lengths, $w^{outer} = \rho \cdot w^{inner}$ and w^{inner} , are the most crucial hyperparameters in *Dual-TF*, we examine the effect of varying these values on the detection accuracy in the TODS dataset. Figure 7(a) demonstrates the change in the F_1 score when w^{inner} varies by [0.2, 0.4, 1.0, 1.2, 1.6, 2.0] times the value determined in Section 3.2 while ρ remains constant at 2. The proposed value for w^{inner} clearly achieves the highest accuracy. Figure 7(b) shows the change in the F_1 score when ρ varies within [1, 2, 3, 4, 5] while maintaining the proposed value for w^{inner} . The highest accuracy is achieved when ρ is 2 or 3, and we choose a smaller one to reduce the number of inner windows for efficiency.

5 Conclusion

We define the concept of the *time-frequency granularity discrepancy* and formulate the problem in exploiting both the time and frequency domains for TSAD. To resolve this discrepancy, we employ the *NS-windowing* scheme to generate anomaly scores at the data-point granularity for both domains. The proposed framework is general and applicable to any sliding window-based TSAD method. This framework is implemented with *Dual-TF* on top of parallel Transformer architectures. Quantitative and qualitative evidence demonstrates the superiority of *Dual-TF*, especially in identifying pattern-wise anomalies with pinpoint accuracy. Overall, we believe that our work paves the way for a new approach to combining the time and frequency domains in time-series data analysis.

ACKNOWLEDGMENTS

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (Ministry of Science and ICT) (No. 2023R1A2C2003690).

REFERENCES

- [1] Ahmed Abdulaal, Zhuanghua Liu, and Tomer Lancewicki. 2021. Practical Approach to Asynchronous Multivariate Time Series Anomaly Detection and Localization. In *KDD*. 2485–2494.
- [2] Ahmed Abdulaal, Zhuanghua Liu, and Tomer Lancewicki. 2021. Practical Approach to Asynchronous Multivariate Time Series Anomaly Detection and Localization. In *KDD*. 2485–2494.
- [3] Charu C Aggarwal and Charu C Aggarwal. 2013. *Outlier Analysis*. Springer, Chapter Applications of Outlier Analysis, 373–400.
- [4] Julien Audibert, Pietro Michiardi, Frédéric Guyard, Sébastien Marti, and Maria A Zuluaga. 2022. Do Deep Neural Networks Contribute to Multivariate Time Series Anomaly Detection? *Pattern Recognition* 132 (2022), 108945.
- [5] Ane Blázquez-García, Angel Conde, Usue Mori, and Jose A Lozano. 2021. A Review on Outlier/Anomaly Detection in Time Series Data. *CSUR* 54, 3 (2021), 1–33.
- [6] Boualem Boashash. 2015. *Time-Frequency Signal Analysis and Processing: A Comprehensive Reference*. Academic press.
- [7] Mohammad Braei and Sebastian Wagner. 2020. Anomaly Detection in Univariate Time-Series: A Survey on the State-Of-The-Art. *arXiv:2004.00433* (2020).
- [8] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: Identifying Density-Based Local Outliers. In *SIGMOD*. 93–104.
- [9] Donald T Campbell. 1975. “Degrees of Freedom” and the Case Study. *Comparative Political Studies* 8, 2 (1975), 178–193.
- [10] Varun Chandola. 2009. *Anomaly Detection for Symbolic Sequences and Time Series Data*. University of Minnesota.
- [11] Leon Cohen. 1995. *Time-Frequency Analysis*. Vol. 778. Prentice hall New Jersey.
- [12] Andrew A Cook, Göksel Misirlı, and Zhong Fan. 2019. Anomaly Detection for IoT Time-Series Data: A Survey. *IEEE Internet of Things Journal* 7, 7 (2019), 6481–6494.
- [13] Zahra Zamanzadeh Darban, Geoffrey I Webb, Shirui Pan, Charu C Aggarwal, and Mahsa Salehi. 2022. Deep Learning for Time Series Anomaly Detection: A Survey. *arXiv:2211.05244* (2022).
- [14] Ailin Deng and Bryan Hooi. 2021. Graph Neural Network-based Anomaly Detection in Multivariate Time Series. In *AAAI*. 4027–4035.
- [15] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. 1996. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *KDD*. 226–231.
- [16] Wei Fan, Shun Zheng, Xiaohan Yi, Wei Cao, Yanjie Fu, Jiang Bian, and Tie-Yan Liu. 2022. DEPTS: Deep Expansion Learning for Periodic Time Series Forecasting. In *ICLR*.
- [17] Astha Garg, Wenyu Zhang, Jules Samaran, Ramasamy Savitha, and Chuan-Sheng Foo. 2022. An Evaluation of Anomaly Detection and Diagnosis in Multivariate Time Series. *IEEE Transactions on Neural Networks and Learning Systems (TNNLS)* 33, 6 (2022), 2508–2517.
- [18] Karlheinz Gröchenig. 2001. *Foundations of Time-Frequency Analysis*. Springer Science & Business Media.
- [19] Karlheinz Gröchenig. 2003. Uncertainty Principles for Time-Frequency Representations. *Advances in Gabor analysis* (2003), 11–30.
- [20] Yifan Guo, Weixian Liao, Qianlong Wang, Lixing Yu, Tianxi Ji, and Pan Li. 2018. Multidimensional time series anomaly detection: A GRU-based Gaussian mixture variational autoencoder approach. In *Asian Conference on Machine Learning*. 97–112.
- [21] Matt Hall. 2006. Resolution and Uncertainty in Spectral Decomposition. *First Break* 24, 12 (2006).
- [22] SS Kelkar, LL Grigsby, and J Langsner. 1983. An Extension of Parseval’s Theorem and Its Use in Calculating Transient Energy in the Frequency Domain. *IEEE Transactions on Industrial Electronics* IE-30, 1 (1983), 42–45.
- [23] Eamonn Keogh, Jessica Lin, and Ada Fu. 2005. Hot SAX: Efficiently Finding the Most Unusual Time Series Subsequence. In *ICDM*.
- [24] Tung Kieu, Bin Yang, Chenjuan Guo, and Christian S Jensen. 2019. Outlier Detection for Time Series with Recurrent Autoencoder Ensembles. In *IJCAI*. 2725–2732.
- [25] Siwon Kim, Kukjin Choi, Hyun-Soo Choi, Byunghan Lee, and Sungroh Yoon. 2022. Towards a Rigorous Evaluation of Time-Series Anomaly Detection. In *AAAI*, Vol. 36. 7194–7201.
- [26] Kwei-Herng Lai, Daochen Zha, Junjie Xu, Yue Zhao, Guanchu Wang, and Xia Hu. 2021. Revisiting Time Series Outlier Detection: Definitions and Benchmarks. In *NeurIPS, Datasets and Benchmarks Track*.
- [27] Gen Li and Jason J Jung. 2022. Deep Learning for Anomaly Detection in Multivariate Time Series: Approaches, Applications, and Challenges. *Information Fusion* (2022).
- [28] Zhihan Li, Youjian Zhao, Jiaqi Han, Ya Su, Rui Jiao, Xidao Wen, and Dan Pei. 2021. Multivariate Time Series Anomaly Detection and Interpretation using Hierarchical Inter-Metric and Temporal Embedding. In *KDD*. 3220–3230.
- [29] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation Forest. In *ICDM*.
- [30] Youngeun Nam, Patara Trirat, Taeyoon Kim, Youngseop Lee, and Jae-Gil Lee. 2023. Context-aware deep time-series decomposition for anomaly detection in businesses. In *ECML-PKDD*. Springer, 330–345.
- [31] Henri J Nussbaumer. 1981. *The Fast Fourier Transform*. Springer.
- [32] Guansong Pang, Chunhua Shen, Longbing Cao, and Anton Van Den Hengel. 2021. Deep Learning for Anomaly Detection: A Review. *CSUR* 54, 2 (2021), 1–38.
- [33] John Paparrizos, Paul Boniol, Themis Palpanas, Ruey S Tsay, Aaron Elmore, and Michael J Franklin. 2022. Volume Under the Surface: A New Accuracy Evaluation Measure for Time-Series Anomaly Detection. *Proceedings of the VLDB Endowment* 15, 11 (2022), 2774–2787.
- [34] John Paparrizos, Yuhao Kang, Paul Boniol, Ruey S Tsay, Themis Palpanas, and Michael J Franklin. 2022. TSB-UAD: An End-to-End Benchmark Suite for Univariate Time-Series Anomaly Detection. *PVLDB* 15, 8 (2022), 1697–1711.
- [35] Yong-chan Park, Jun-Gi Jang, and U Kang. 2021. Fast and Accurate Partial Fourier Transform for Time Series Data. In *KDD*. 1309–1318.
- [36] Hansheng Ren, Bixiong Xu, Yujing Wang, Chao Yi, Congrui Huang, Xiaoyu Kou, Tony Xing, Mao Yang, Jie Tong, and Qi Zhang. 2019. Time-Series Anomaly Detection Service at Microsoft. In *KDD*. 3009–3017.
- [37] Bernhard Schölkopf, Robert C Williamson, Alexander J Smola, John Shawe-Taylor, John C Platt, et al. 1999. Support Vector Method for Novelty Detection. In *NeurIPS*. 582–588.
- [38] Yooju Shin, Susik Yoon, Sundong Kim, Hwanjun Song, Jae-Gil Lee, and Byung Suk Lee. 2021. Coherence-Based Label Propagation over Time Series for Accelerated Active Learning. In *ICLR*.
- [39] Yooju Shin, Susik Yoon, Hwanjun Song, Dongmin Park, Byunghyun Kim, Jae-Gil Lee, and Byung Suk Lee. 2023. Context consistency regularization for label sparsity in time series. In *ICML*. 31579–31595.
- [40] IBM Research Editorial Staff. 2023. “Simple” Threshold Algorithm Earns Gödel Prize. <https://www.ibm.com/blogs/research/2014/05/simple-threshold-algorithm-earns-godel-prize/>. Accessed: 2023-04-30.
- [41] Ya Su, Youjian Zhao, Chenhao Niu, Rong Liu, Wei Sun, and Dan Pei. 2019. Robust Anomaly Detection for Multivariate Time Series through Stochastic Recurrent Neural Network. In *KDD*. 2828–2837.
- [42] Patara Trirat, Youngeun Nam, Taeyoon Kim, and Jae-Gil Lee. 2023. ANOVIZ: A Visual Inspection Tool of Anomalies in Multivariate Time Series. In *AAAI Association for the Advancement of Artificial Intelligence*.
- [43] Shreshth Tuli, Giuliano Casale, and Nicholas R Jennings. 2022. TranAD: Deep Transformer Networks for Anomaly Detection in Multivariate Time Series Data. *PVLDB* 15, 6 (2022), 1201–1214.
- [44] Qingsong Wen, Tian Zhou, Chaoli Zhang, Weiqi Chen, Ziqing Ma, Junchi Yan, and Liang Sun. 2022. Transformers in Time Series: A Survey. *arXiv:2202.07125* (2022).
- [45] Haowen Xu, Wenxiao Chen, Nengwen Zhao, Zeyan Li, Jiahao Bu, Zhihan Li, Ying Liu, Youjian Zhao, Dan Pei, Yang Feng, et al. 2018. Unsupervised Anomaly Detection via Variational Auto-Encoder for Seasonal KPIs in Web Applications. In *WWW*. 187–196.
- [46] Jiehui Xu, Haixu Wu, Jianmin Wang, and Mingsheng Long. 2022. Anomaly Transformer: Time Series Anomaly Detection with Association Discrepancy. In *ICLR*.
- [47] Ling Yang and Shenda Hong. 2022. Unsupervised Time-Series Representation Learning with Iterative Bilinear Temporal-Spectral Fusion. In *ICML*. 25038–25054.
- [48] Kun Yi, Qi Zhang, Shoujin Wang, Hui He, Guodong Long, and Zhendong Niu. 2023. Neural Time Series Analysis with Fourier Transform: A Survey. *arXiv:2302.02173* (2023).
- [49] Chuxu Zhang, Dongjin Song, Yuncong Chen, Xinyang Feng, Cristian Lumezanu, Wei Cheng, Jingchao Ni, Bo Zong, Haifeng Chen, and Nitesh V Chawla. 2019. A Deep Neural Network for Unsupervised Anomaly Detection and Diagnosis in Multivariate Time Series Data. In *AAAI*. 1409–1416.
- [50] Chaoli Zhang, Tian Zhou, Qingsong Wen, and Liang Sun. 2022. TFAD: A Decomposition Time Series Anomaly Detection Architecture with Time-Frequency Analysis. In *CIKM*. 2497–2507.
- [51] Xiang Zhang, Ziyuan Zhao, Theodoros Tsiligkaridis, and Marinka Zitnik. 2022. Self-Supervised Contrastive Pre-Training For Time Series via Time-Frequency Consistency. In *NeurIPS*.
- [52] Hang Zhao, Yujing Wang, Juanyong Duan, Congrui Huang, Defu Cao, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, and Qi Zhang. 2020. Multivariate Time-Series Anomaly Detection via Graph Attention Network. In *ICDM*. 841–850.
- [53] Hang Zhao, Yujing Wang, Juanyong Duan, Congrui Huang, Defu Cao, Yunhai Tong, Bixiong Xu, Jing Bai, Jie Tong, and Qi Zhang. 2020. Multivariate Time-Series Anomaly Detection via Graph Attention Network. In *ICDM*. 841–850.
- [54] Bin Zhou, Shenghua Liu, Bryan Hooi, Xueqi Cheng, and Jing Ye. 2019. BeatGAN: Anomalous Rhythm Detection using Adversarially Generated Time Series. In *IJCAI*. 4433–4439.
- [55] Tian Zhou, Ziqing Ma, Qingsong Wen, Xue Wang, Liang Sun, and Rong Jin. 2022. FEDFormer: Frequency Enhanced Decomposed Transformer for Long-Term Series Forecasting. In *ICML*. 27268–27286.

A Details of the Window Length Selection

For a time series $X \in \mathbb{R}^{n \times d}$, the discrete Fourier transform (DFT) for the d' -th dimension ($d' \in [1, d]$) of X is formulated as

$$X_k^{d'} = \sum_{t=0}^{n-1} x_t^{d'} e^{-\frac{i2\pi}{n} kt}, \quad (12)$$

where

$$x_t^{d'} = \frac{1}{n} \sum_{k=0}^{n-1} X_k^{d'} \cdot e^{-\frac{i2\pi}{n} kt}, \quad (13)$$

and k is an integer ranging from 0 to $n-1$.

The radix-2 decimation rearranges the DFT of the function x_t into two parts: a sum over the even-numbered indices $t = 2m$ and a sum over the odd-number indices $t = 2m+1$,

$$X_k^{d'} = \sum_{m=0}^{n/2-1} x_{2m}^{d'} e^{-\frac{i2\pi}{n} k(2m)} + \sum_{m=0}^{n/2-1} x_{2m+1}^{d'} e^{-\frac{i2\pi}{n} k(2m+1)}. \quad (14)$$

Denote the DFT of the even-indexed inputs $x_{2m}^{d'}$ by $E_k^{d'}$ and the DFT of the odd-indexed inputs $x_{2m+1}^{d'}$ by $O_k^{d'}$, and we obtain

$$\begin{aligned} X_k^{d'} &= \sum_{m=0}^{n/2-1} x_{2m}^{d'} e^{-\frac{i2\pi}{n/2} km} + e^{-\frac{i2\pi}{n} k} \sum_{m=0}^{n/2-1} x_{2m+1}^{d'} e^{-\frac{i2\pi}{n/2} km} \\ &= E_k^{d'} + e^{-\frac{i2\pi}{n} k} O_k^{d'} \quad \text{for } k = 0, \dots, \frac{n}{2} - 1. \end{aligned} \quad (15)$$

Due to the periodicity of the complex exponential, $X_{k+\frac{n}{2}}^{d'}$ can be also represented as

$$\begin{aligned} X_{k+\frac{n}{2}}^{d'} &= \sum_{m=0}^{n/2-1} x_{2m}^{d'} e^{-\frac{i2\pi}{n/2} m(k+\frac{n}{2})} \\ &\quad + e^{-\frac{i2\pi}{n} (k+\frac{n}{2})} \sum_{m=0}^{n/2-1} x_{2m+1}^{d'} e^{-\frac{i2\pi}{n/2} m(k+\frac{n}{2})} \\ &= \sum_{m=0}^{n/2-1} x_{2m}^{d'} e^{-\frac{i2\pi}{n/2} mk} e^{-i2\pi m} \\ &\quad + e^{-\frac{i2\pi}{n} k} e^{-i\pi} \sum_{m=0}^{n/2-1} x_{2m+1}^{d'} e^{-\frac{i2\pi}{n/2} mk} e^{-i2\pi m} \\ &= \sum_{m=0}^{n/2-1} x_{2m}^{d'} e^{-\frac{i2\pi}{n/2} mk} - e^{-\frac{i2\pi}{n} k} \sum_{m=0}^{n/2-1} x_{2m+1}^{d'} e^{-\frac{i2\pi}{n/2} mk} \\ &= E_k^{d'} - e^{-\frac{i2\pi}{n} k} O_k^{d'}. \end{aligned} \quad (16)$$

Rewrite $X_k^{d'}$ and $X_{k+\frac{n}{2}}^{d'}$ to

$$\begin{aligned} X_k^{d'} &= E_k^{d'} + e^{-\frac{i2\pi}{n} k} O_k^{d'} \\ X_{k+\frac{n}{2}}^{d'} &= E_k^{d'} - e^{-\frac{i2\pi}{n} k} O_k^{d'}. \end{aligned} \quad (17)$$

Note that the final outputs are obtained by a $+/ -$ combination of $E_k^{d'}$ and $O_k^{d'} e^{-\frac{i2\pi}{n} k}$. Therefore, we can use only one of the two parts to get information about the magnitude through the absolute value (because both have the same magnitude if the absolute value is taken). We can define the number of sampling rate (sampling frequency) per time unit as f obtained from a continuous signal. The frequency range is expressed as $f = \frac{k}{n}$ for $k = 0, \dots, \frac{n}{2} - 1$, where n is the time period.

Among d dimensions, the value with the smallest frequency is set as the most dominant frequency. The reason for choosing the minimum frequency is the allowance for the inclusion of other dominant frequencies in other different dimensions within the determined inner window. We define the maximum magnitude index, which means the most dominant frequency, v_{major} ,

$$\begin{aligned} v_{\text{major}} &= \min(\arg \max(|X_f^{d'}|) : d' \in [1, d]) \\ &\quad \text{for } f = 0, \dots, \frac{n-2}{2n}. \end{aligned} \quad (18)$$

Finally, the inner window length can be determined as

$$w^{\text{inner}} = \lceil \frac{1}{v_{\text{major}}} \rceil. \quad (19)$$

B Proof of Theorem 3.1

The proof can be done using Parseval's theorem [22]. By our definition of the uncertainty in each domain, $\mathcal{U}^{\text{time}}(w^{\text{inner}})$ monotonically decreases as w^{inner} decreases; $\mathcal{U}^{\text{freq}}(w^{\text{inner}})$ converges to a zero when $w^{\text{inner}} \geq \lceil \frac{1}{v_{\text{major}}} \rceil$, but it monotonically increases as w^{inner} decreases when $w^{\text{inner}} < \lceil \frac{1}{v_{\text{major}}} \rceil$.

Specifically, in the process of interpreting the results obtained through the NS-windowing and the frequency reconstruction by Θ^{freq} , there is uncertainty between two different domains. From the perspective of the time domain, the smaller w^{inner} facilitates the identification of abnormal time points. Conversely, as w^{inner} increases, distinguishing time points within a window in the frequency domain becomes more challenging. On the other hand, from the perspective of the frequency domain, the longer w^{inner} leads to a decrease in frequency variance according to the uncertainty principle, and thus, frequency becomes concentrated. However, a decrease in w^{inner} makes it difficult to accurately determine the precise frequency after the Fourier transform.

We formulate the inherent trade-off associated with dual-domain information loss. Assume that every time series can be represented as a single periodic function with a dominant frequency that affects the pattern most. Then, the time series can be simply expressed as $x_t = \cos(2\pi v_{\text{major}} \frac{t}{w}) + \epsilon_t$, where $t = 0, \dots, w-1$ with the dominant frequency v_{major} and a noise $\epsilon_t \sim N(0, w)$.

Here, $\mathcal{U}(w)$ in Eq. (1) is a function of w , which represents the dual-domain information loss. $\mathcal{U}^{\text{time}}(w)$ denotes the uncertainty in the time domain, which monotonically decreases as w decreases as a linear function of w ($w \geq 1, w \in \mathbb{N}$), so we define it as

$$\mathcal{U}^{\text{time}}(w) = w - 1. \quad (20)$$

The choice of a linear relationship is specific to the Gaussian function and the Fourier transform. The Gaussian function has the unique property that its Fourier transform is also a Gaussian function, and this symmetry leads to the linear relationship between the standard deviations (or uncertainties) in time and frequency. While other functions, such as square, cubic, or exponential functions, can be used to model specific types of uncertainty, they would not lead to the same fundamental relationship as the Gaussian function does in the context of the Gabor limit. The Gaussian function is crucial due to its role in signal processing and its mathematical properties that align with both time and frequency domains. Related works [6, 11] provide more detailed insights into the mathematical

reasoning behind the linear relationship in the Gabor uncertainty principle.

$\mathcal{U}^{freq}(w)$ denotes the uncertainty in the frequency domain. \widehat{F}_v denotes a magnitude of frequency v after the Fourier transform (or DFT) of x_t , where $v = 0, \dots, w-1$ ($v = 0, \dots, \frac{w}{2}-1$ for DFT). \widehat{F}_v is formulated as $\widehat{F}_v = F_v + E_v$, where F_v is the clean (unperturbed) Fourier transform and E_v is the Fourier transform of the noise ϵ_t . $\mathcal{U}^{freq}(w)$ can be defined as the sum of the standard deviation of the Fourier transformed \widehat{F}_v , i.e., $\sigma(E_v)$ and the function of F_1 scores (inversely related to the anomaly scores) $f(w)$ with varying the length of w in Figure 7, where $w = 1/v$ (by the definition of a period in terms of frequency).

We approximate the anomaly score function f by fitting the function $aw - \log(bw)$ to Figure 7. To justify the fitting approach, let us conduct a T-test on a custom fitting function f with two other possible functions, i.e., linear function $-ax + b$ and rational function $\frac{a}{x} + b$. We can fit the custom function f using *spicy curve_fit* library and calculate the standard error of the parameters from the covariance. The each null hypothesis of a and b values is

$$\begin{aligned} H_0 : \mu(\Delta_\phi^{custom}) &= \mu(\Delta_\phi^{linear}), \\ H_0 : \mu(\Delta_\phi^{custom}) &= \mu(\Delta_\phi^{rational}), \end{aligned} \quad (21)$$

where Δ_ϕ^{custom} is the standard error of parameter $\phi \in \{a, b\}$ in the custom function (identical meaning in both linear and rational functions). The maximum p-value is 0.029 among all pairs with a significance level of 0.05. Therefore we reject the null hypothesis for every pair with the custom fitting function, meaning that the standard deviations of parameters are less than the other two compared fitting functions. Additionally, by fitting the f using the data in Figure 7, the value of a becomes $v_{major} - 1$, and the value of b becomes 1 after fitting with the R-square value of 0.695–0.972, which is much higher than the other compared functions in TODS benchmark datasets.

By Parseval's theorem [22], the sum of the square of a function is equal to the sum of the square of its transform, $\sum |E_v|^2 = \frac{1}{w} \sum |\epsilon_t|^2 = \sigma(\epsilon_t)^2$. Therefore,

$$\sigma(E_v) = \sqrt{\frac{1}{w} \sum |E_v|^2} = \sqrt{\frac{\sigma(\epsilon_t)^2}{w}} = 1. \quad (22)$$

Finally,

$$\begin{aligned} \mathcal{U}^{freq}(w) &= f(v) + \sigma(E_v) \\ &= (v_{major} - 1)w - \log w + 1. \end{aligned} \quad (23)$$

The uncertainty in the two domains is

$$\begin{aligned} \mathcal{U}(w) &= \mathcal{U}^{time}(w) + \mathcal{U}^{freq}(w) \\ &= w - 1 + (v_{major} - 1)w - \log w + 1. \end{aligned} \quad (24)$$

The length that minimizes Eq. (24) (or Eq. (1)) can be obtained by solving the differential equation for $\mathcal{U}(w)$ with respect to w . Taking the derivative of $\mathcal{U}(w)$ with respect to w ,

$$\begin{aligned} \frac{\partial \mathcal{U}(w)}{\partial w} &= \frac{\partial (w - 1 + (v_{major} - 1)w - \log w + 1)}{\partial w} \\ &= 1 + v_{major} - 1 - \frac{1}{w} \\ &= v_{major} - \frac{1}{w} = 0. \end{aligned} \quad (25)$$

We now find the solution with $w = \frac{1}{v_{major}}$, and the length of the inner window is determined as $w^{inner} = \lceil \frac{1}{v_{major}} \rceil$ that is consistent with Eq. (19). This completes the proof. \square

C Details of the Experiment Settings

C.1 Datasets

We evaluate anomaly detection performance on a total of 30 datasets from 5 benchmarks, including four publicly available datasets.

ASD [28] is a server benchmark consisting of 45-day-long multi-variate time series. The benchmark comprises 12 entities obtained from different servers, each defined d by 19 metrics that reflect the server's status. These metrics include CPU-related parameters, memory-related parameters, network metrics, and virtual machine metrics. The ASD benchmark is publicly available under the MIT license at <https://github.com/zhlee/InterFusion/tree/main/data>.

ECG [23] (Electrocardiogram) represents time series of the electrical potential variation between two points on the body's surface, primarily originating from the rhythmic contractions of the heart. The benchmark contains 9 sub-datasets. The ECG benchmark is publicly accessible at https://www.cs.ucr.edu/~eamonn/discords/ECG_data.zip.

PSM [2] represents a benchmark comprising a single entity derived from multiple application server nodes at eBay. It encompasses 26 features, publicly accessible under the CC BY 4.0 license at <https://github.com/eBay/RANSynCoders/tree/main/data>. Similar to ASD, these features describe server machine metrics such as CPU utilization and memory usage. In the PSM benchmark, the training set spans 13 weeks, followed by 8 weeks for testing.

TODS [26] is a publicly available synthetic benchmark and data generator for time-series anomaly detection, accessible at <https://github.com/datamllab/tods/tree/benchmark/benchmark/synthetic/Generator>. The benchmark includes 5 anomaly scenarios classified based on a behavior-driven taxonomy: point-global, pattern-contextual, pattern-shapelet, pattern-seasonal, and pattern-trend anomalies. The TODS data generator produces 5 individual univariate time series, each corresponding to a distinct anomaly type. We use the provided source code without alterations as demonstrated below, except for adjusting the length parameter to generate a longer time series, to ensure a fair comparison.

```

1 # Source: https://github.com/datamllab/tods
2 # DataGenerator: "UnivariateDataGenerator" from
  univariate_generator.py.
3 BEHAVIOR_CONFIG = {"freq": 0.04, "coef": 1.5, "offset": 0.0, "
  noise_amp": 0.05}
4 BASE = [0.145, 0.128, 0.094, 0.077, 0.111, 0.145, 0.179, 0.214,
  0.214]
5
6 # Training set
7 normal = DataGenerator(stream_length=20000,
  behavior=sine,
  behavior_config=BEHAVIOR_CONFIG)
8 normal.generate_timeseries()
9
10 # Test set
11 test = DataGenerator(stream_length=5000,
  behavior=sine,
  behavior_config=BEHAVIOR_CONFIG)
12
13 # Point - global anomaly
14 if anomaly_type=="point_global":
15     test.point_global_outliers(ratio=0.05, factor=3.5, radius=5)

```

Table 5: Experiment environments for all algorithms.

Env.	OmniAnomaly	LOF	ISF	OCSVM	VAE	MS-RNN	RANSynCoders	TranAD	Anomaly Transformer	Dual-TF
Library	Tensorflow 1.12.0	Scikit-Learn 1.2.1			Tensorflow 2.5.0			PyTorch 1.13.1		
CPU	Intel(R) Xeon(R) Silver 4116 CPU @ 2.10GHz	Intel(R) Xeon(R) Gold 6226R CPU @ 2.90GHz								
GPU	NVIDIA GeForce RTX 2080 Ti 11GB with CUDA Version 11	NVIDIA GeForce RTX 3090 24GB with CUDA Version 11								

```

20 elif anomaly_type=='point_contextual':
21     test.point_contextual_outliers(ratio=0.05, factor=2.5, radius
22                                     =5)
23 elif anomaly_type=='pattern_shaplet':
24     test.collective_global_outliers(ratio=0.05, radius=5, option='
25                                     square', coef=1.5, noise_amp=0.03, level=20, freq=0.04, base=
26                                     BASE, offset=0.0)
27 elif anomaly_type=='pattern_seasonal':
28     test.collective_seasonal_outliers(ratio=0.05, factor=3, radius
29                                       =5)
30 elif anomaly_type=='pattern_trend':
31     test.collective_trend_outliers(ratio=0.05, factor=0.5, radius
32                                    =5)

```

Listing 1: The command used for generating the TODS benchmark datasets.

C.2 Baselines

LOF is an algorithm that measures the local deviation of the density of a given sample, **ISF** is an algorithm isolating anomalies using tree-based structures, and **OCSVM** is an algorithm based on the SVM that maximizes the margin between the origin and the normality and defines the decision boundary as the hyper-plane that determines the margin. **VAE** uses the symmetrical encoder and decoder network and anomaly scores are the differences between the inputs and reconstructed outputs. Modified-RNN (**MS-RNN**) is a modified version of sparsely-connected RNNs with an ensemble of autoencoders (AEs); **OmniAnomaly** is an LSTM-based VAE capturing complex temporal dependencies; **RANSynCoders**, a model that uses pre-trained AEs to extract primary frequencies across the signals on the latent representation for synchronizing time series; **TranAD**, a Transformer-based model that uses attention-based sequence encoders to perform inference with broader temporal trends, with the focus on score-based self-conditioning for robust multi-modal feature extraction and adversarial training; **Anomaly Transformer**, a reconstructive approach that combines series and prior association to make anomalies distinctive; and **TFAD**, a time-frequency analysis-based model that uses both time and frequency domains, with time-series decomposition and data augmentation.

C.3 Experiment Environments

Table 5 shows the experiment hardware environments for all algorithms. Only OmniAnomaly was run on a different environment due to its incompatibility between the NVIDIA cuDNN and Tensorflow libraries. We adopt a singular GPU for conducting experiments on each benchmark and algorithm. All experiments are conducted on a server equipped with an NVIDIA RTX 3090Ti.

C.4 Model Hyperparameter Settings

The ISF, LOF, and OCSVM algorithms are implemented using the Scikit-Learn library, while the remaining methods are configured using open-source code obtained from each URL. The hyperparameters for the baseline methods are set as follows.

- **ISF**: The number of tree is selected from {25, 100}.
- **LOF**: The number of neighbors is selected from {1, 3, 5, 12}.
- **OCSVM**: The RBF kernel is used. The inverse length parameter γ is selected from $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 0.5\}$.
- **VAE**²: The LSTM layers are used as both the encoder and decoder. The number of hidden units in the encoder and decoder are set to 64 and 32, respectively.
- **MS-RNN**³: The GRU layers are employed as the encoder, and the decoder consists of a skip-GRU structure with a reverse chronological order in the time series.
- **OmniAnomaly**⁴: The GRU and dense layers have 500 units. The standard deviation layer has an ϵ value of 10^{-4} . The dimension of the latent variable z -space is fixed at 3.
- **RANSynCoders**⁵: Hidden layers in each decoder are increased as the output dimension is at least 3 times larger than the encoder input size. S , N , and the bootstrap sample size are set to one-third of the input dimension, rounded to the nearest multiple of 5. The number of latent dimensions is chosen as $0.5N-1$. δ is set to 0.05 for system data with Gaussian outliers and 0.1 for business data without Gaussian outliers.
- **TranAD**⁶: Layers in the Transformer encoders is set to 1 and 2 for the feed-forward unit of the encoders. The hidden units in the encoder layers is set to 64, and the dropout is set to 0.1.
- **TFAD**⁷: The kernel size for the temporal convolutional network (TCN) is set to 7. The number of TCN layers is 3. The dimension of the embedding representation is set to 150. The distance metric is the L2 norm. The classifier threshold is set to 0.5, and the mixup rate is set to 2.
- **Anomaly Transformer**⁸: The number of layers is 3, the channel number of hidden states d_{model} is 512, and the number of heads h is 8. The loss function hyperparameter λ for balancing two parts is set as 3. These hyperparameter settings are shared with *Dual-TF*.

²<https://github.com/lin-shuyu/VAE-LSTM-for-anomaly-detection>

³<https://github.com/tungk/OED>

⁴<https://github.com/NetManAI/Ops/OmniAnomaly>

⁵<https://github.com/eBay/RANSynCoders>

⁶<https://github.com/imperial-qore/TranAD>

⁷<https://github.com/DAMO-DI-ML/CIKM22-TFAD>

⁸<https://github.com/thuml/Anomaly-Transformer>